UNIVERSIDAD NACIONAL DE HUANCAVELICA

(Creado por Ley 25265)

FACULTAD DE CIENCIAS DE INGENIERÍA

ESCUELA ACADÉMICA PROFESIONAL DE CIVIL (HUANCAVELICA)



TESIS

"TÉCNICAS NUMÉRICAS SOBRE PROCESADOR GRÁFICO, PARA LA SOLUCIÓN DE SISTEMAS DE ECUACIONES LINEALES EN ANÁLISIS ESTRUCTURAL"

LÍNEA DE INVESTIGACIÓN

ANÁLISIS ESTRUCTURAL

PRESENTADO POR:

Bach. CASTELLARES PAUCAR, William

PARA OPTAR EL TÍTULO PROFESIONAL DE:

INGENIERO CIVIL

HUANCAVELICA, PERÚ 2019



UNIVERSIDAD NACIONAL DE HUANCAVELICA

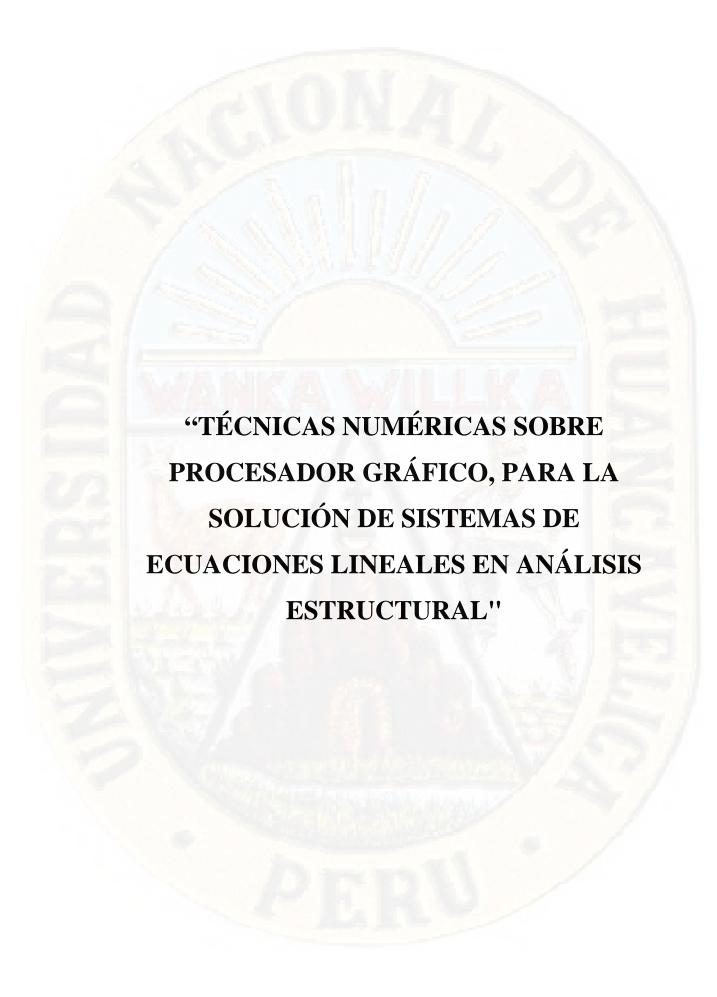


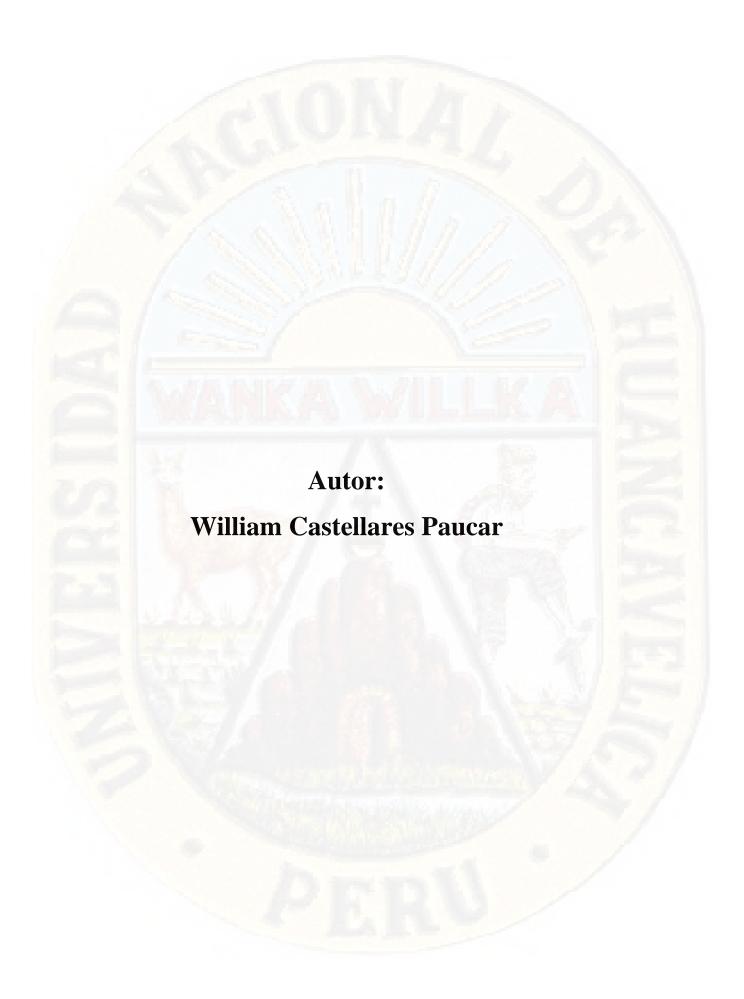
FACULTAD DE CIENCIAS DE INGENIERÍA

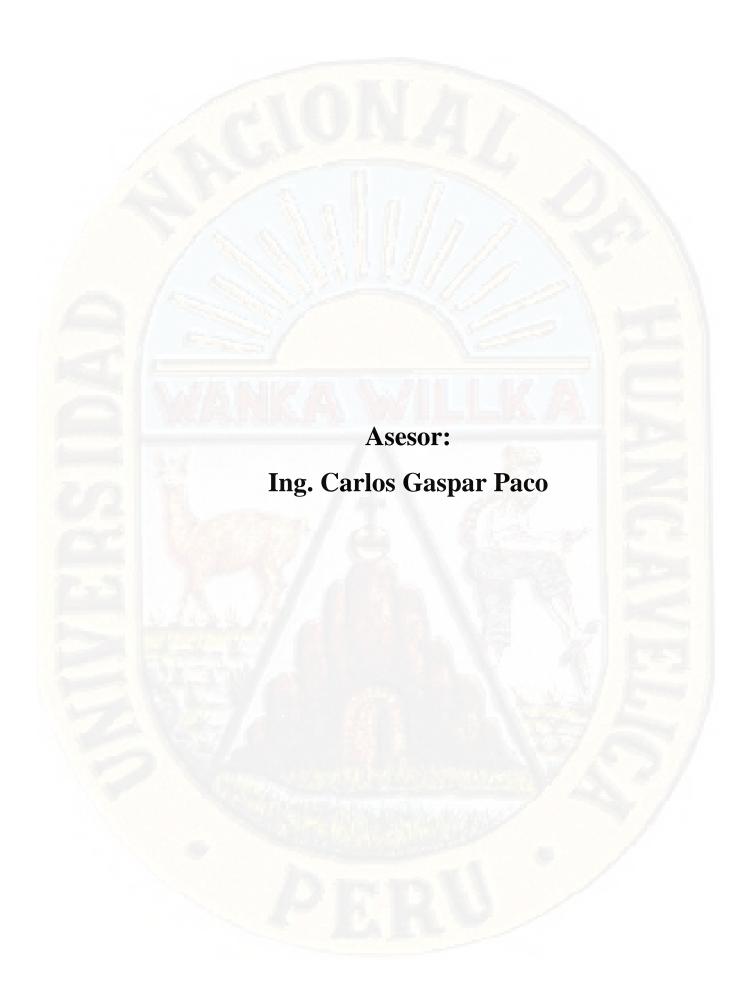
ACTA DE SUSTENTACIÓN DE TESIS

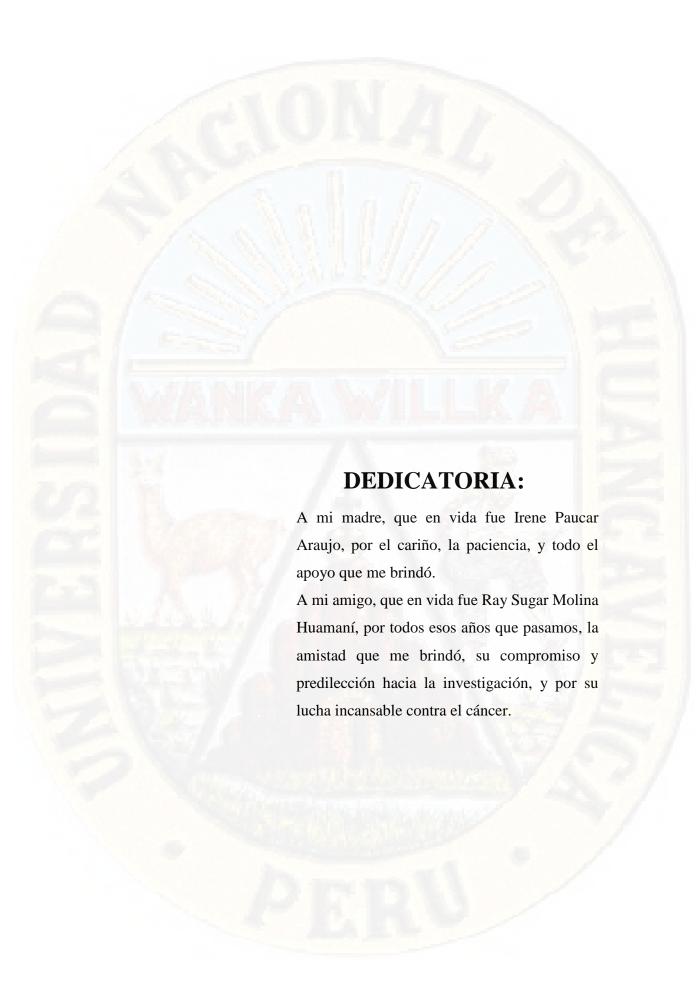
En el Auditórium de la Facultad de Ciencias de Ingeniería, a los 12 días del mes de diciembre del año 2019, a horas 12:00 m., se reunieron los miembros del Jurado Calificador conformado por los siguientes: M.Sc. Iván Arturo AYALA BIZARRO (PRESIDENTE), Arq. Abdón Dante OLIVERA QUINTANILLA (SECRETARIO), M.Sc. Hugo Rubén LUJAN JERI (VOCAL), designados con Resolución de Consejo de Facultad N° 338-2017-FCI-UNH, de fecha 03 de agosto del 2017, así mismo se cambió al Asesor y han sido reestructurado los miembros de Jurados Evaluadores con Resolución de Decano N° 196-2019-FCI-UNH, de fecha 30 de octubre del 2019 y con Resolución de Decano N° 281-2019-FCI-UNH de fecha 06 de diciembre del 2019 se ratifica al Asesor y Jurados evaluadores, a fin de proceder con la calificación de la sustentación del informe final de tesis titulado: "TÉCNICAS NUMÉRICAS SOBRE PROCESADOR GRÁFICO, PARA LA SOLUCIÓN DE SISTEMAS DE ECUACIONES LINEALES EN ANÁLISIS ESTRUCTURAL", presentado por el Bachiller William CASTELLARES PAUCAR, para optar el Título Profesional de Ingeniero Civil; en presencia del Ing. Carlos GASPAR PACO como Asesor del presente trabajo de tesis. Finalizado la evaluación a horas. L. 20. p.m.; se invitó al público presente y al sustentante abandonar el recinto. Luego de una amplia deliberación por parte de los Jurados, se llegó al siguiente resultado:

APROBADO		POR UNANIMIDAD	
DESAPROBADO			
En señal de conformi Presidente	dad, firmam	Secretario	Vocal



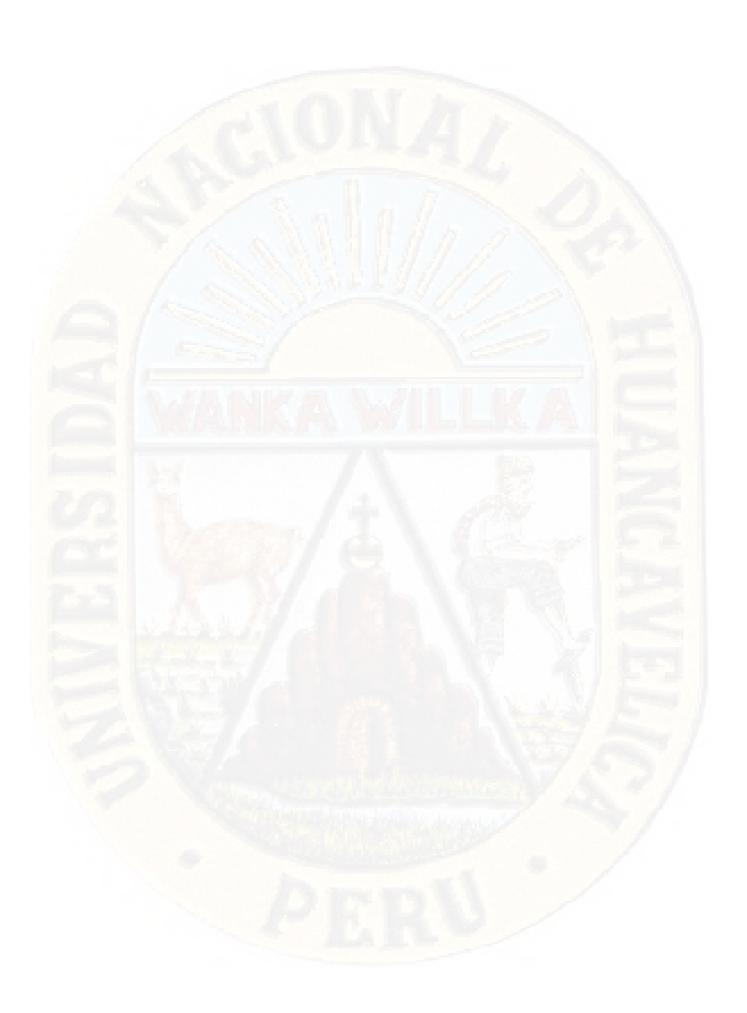






AGRADECIMIENTOS

- A la Escuela Profesional de Ingeniería Civil- Huancavelica; a los docentes que intentan dar lo mejor de sí mismos en la formación de buenos profesionales en esta difícil carrera, y a los amigos que hice durante mi permanencia en esta escuela.
- A mi familia, en especial a mi hermana y sobrino por su apoyo y paciencia.



Contenido

		MEN	
	ABST	RACT	xviii
	INTRO	ODUCCIÓN	xix
	CAPÍT	TULO I	22
	PROB	LEMA	22
1.1.	Descr	ripción del problema	22
1.2.	Form	ulación del problema	24
	1.2.1.	Problema general	24
	1.2.2.	Problemas específicos	24
1.3.		tivos	
	1.3.1.	Objetivo general	24
	1.3.2.	Objetivos específicos	25
1.4.	Justif	icación	25
1.5.	Limit	taciones	26
	CAPÍT	TULO II	27
		CO TEÓRICO	
2.1.	Ante	cedentes	27
2.2.	Bases	s teóricas sobre el tema de investigación	31
	2.2.1.	Sistema de ecuaciones lineales (SEL)	
	2.2.2.	Uso de procesadores gráficos (GPU) en propósitos generales	37
2.3.	Bases	s conceptuales	42
	2.3.1.	Programación de la GPU en OpenCL	42
	2.3.2.	Formato de almacenamiento CSC	
	2.3.3.	Formato de almacenamiento SKS	50

2.4.		ición de términos	
2.5.	Varia	bles	
	2.5.1.	Variable independiente	54
	2.5.2.	Variable dependiente	55
2.6.	Opera	acionalización de variables	55
	CAPÍT	'ULO III	56
		DOLOGÍA DE INVESTIGACIÓN	
3.1.	Tipo o	de investigación	56
3.2.	Nivel	de investigación	56
3.3.	Pobla	ción, muestra y muestreo	56
3.4.	Técni	cas e instrumentos de obtención de datos	57
3.5.		dimiento de obtención de datos	
3.6.		cas de procesamiento y análisis de datos	
		TULO IV	
	PRESE	ENTACIÓN DE RESULTADOS	59
4.1.	Análi	sis de información	59
	4.1.1.	Etapas	59
	4.1.2.	Compatibilidad en la indexación	63
	4.1.3.	Orden de almacenamiento preferencial	63
	4.1.4.	Alistando los kernels	65
	4.1.5.	Kernels sobre matrices densas	65
	4.1.6.	Kernels sobre matrices en formato CSC	83
	4.1.7.	Kernels sobre matrices en formato SKS	109
	4.1.8.	Medición del tiempo y memoria utilizada	
	4.1.9.	Características del dispositivo	131

	4.1.11.	Modelo 1: Viga en voladizo	132
	4.1.12.	Modelo 2: Pilar de puente	136
	4.1.13.	Balance entre memoria y tiempo de ejecución	141
	4.1.14.	Pruebas adicionales	144
4.1.	Discusion	ón de resultados	147
	CONCL	USIONES	151
	RECOM	ENDACIONES	154
	Bibliogra	ıfía	156
	ANEXOS	S	158

Índice de tablas:

Tabla 1. Correspondencia de tipo entre objetos en ArrayFire y en Julia	46
Tabla 2. Operacionalización de variables	55
Tabla 3. Primera etapa de la investigación	59
Tabla 4. Segunda etapa de la investigación	60
Tabla 5. Tercera etapa de la investigación	
Tabla 6. Cuarta etapa de la investigación	61
Tabla 7. Quinta etapa de la investigación	61
Tabla 8. Sexta etapa de la investigación	62
Tabla 9. Séptima etapa de la investigación	62
Tabla 10. Modelo 1: memoria utilizada en MB	
Tabla 11. Modelo 1: tiempo de ejecución en segundos	135
Tabla 12. Modelo 2: memoria utilizada en MB	
Tabla 13. Modelo 2: tiempo de ejecución en segundos	
Tabla 14. Modelo 2: tiempo de ejecución en segundos	140
Tabla 15. Modelo 2: memoria utilizada en MB para el formato SKS	142
Tabla 16. Modelo 2: tiempo de ejecución en matrices en formato SKS	143
Tabla 17. Diferencias en las características de los dispositivos	145
Tabla 18. Modelo 2: tiempo de ejecución en diferentes dispositivos	145

Índice de gráficos

Gráfico 1: Modelo 1: memoria utilizada en MB	134
Gráfico 2: Modelo 1: tiempo de ejecución en segundos	135
Gráfico 3: Modelo 1: tiempo de ejecución en segundos	136
Gráfico 4: Modelo 2: memoria utilizada en MB	139
Gráfico 5: Modelo 2: memoria utilizada en MB	139
Gráfico 6: Modelo 2: tiempo de ejecución en segundos	141
Gráfico 7: índice de dispersión para el formato SKS	143
Gráfico 8: tiempo de ejecución en matrices en formato SKS	144
Gráfico 9: tiempo de ejecución en diferentes dispositivos	146
Gráfico 10: tiempo de ejecución GT 630M y solver de OpenFEM	147
Gráfico 11: tiempo de ejecución GTX 1070 y solver de OpenFEM	147

Índice de figuras:

Figura 1. Ejemplo numérico de la deformación de una estructura (Félix, 2015) 29
Figura 2. Comparación de los tiempos de solución de los diferentes solvers (Félix,
2015)
Figura 3. Comparación de la capacidad de memoria usada por los diferentes solvers
paralelos (Félix, 2015)30
Figura 4. Número de supercomputadores que emplean aceleradores, extraída de
(Juliacomputing.com, s.f.) 39
Figura 5. Comparación de la aceleración del gradiente conjugado en diferentes
tamaños, extraído de (Juliacomputing.com, s.f.)
Figura 6. Escalabilidad automática, extraída de (NVIDIA Corporation)
Figura 7. Ejemplo de un Kernel que suma dos vectores
Figura 8. Ejemplo de la ejecución de un kernel usando la biblioteca de ArrayFire
(ArrayFire company, s.f.)45
Figura 9. Declaración de la función af_add de ArrayFire
Figura 10. Función <i>add</i> de ArrayFire.jl
Figura 11. Función con argumentos af_array que pueden ser llamados desde Julia. 48
Figura 12. Código que llama en Julia a la función de la Figura 1149
Figura 13. Contenido de Common.h
Figura 14. Gauss Jordan: Kernel para trabajar sobre una matriz densa
Figura 15. Función de dispositivo llamada por le kernel de la Figura 14
Figura 16. Cholesky: función de dispositivo <i>Chol_1_c</i>
Figura 17. Cholesky: función de dispositivo Chol_2_c70
Figura 18. Cholesky: función de dispositivo <i>Chol_3_c</i> 71
Figura 19. Cholesky: función de dispositivo <i>Chol_4_c</i>
Figura 20. Cholesky: función de dispositivo <i>Chol_5_c</i>
Figura 21. Cholesky: función de dispositivo <i>Chol_6_c</i>
Figura 22. Cholesky: Kernel para trabajar sobre una matriz densa
Figura 23. LDL ^T : función de dispositivo <i>ldlt_1_c</i>
Figura 24. LDL ^T : función de dispositivo <i>ldlt_2_c</i>
Figura 25. LDL ^T : función de dispositivo <i>ldlt 3 c</i>

Figura 26. LDL ^T : función de dispositivo <i>ldlt_4_c</i>	78
Figura 27. LDL ^T : función de dispositivo <i>ldlt_5_c</i>	79
Figura 28. LDL ^T : Kernel para trabajar sobre una matriz densa	
Figura 29. Gradientes conjugados: función principal a exportar	83
Figura 30. Estructura para objetos que almacenan la matriz en formato CSC	83
Figura 31. Estructura para objetos que almacenan la matriz en formato CSC co	n la
memoria del dispositivo	84
Figura 32. Conversión de una matriz densa a una matriz dispersa en formato CSC	. 84
Figura 33. Función para convertir un objeto del tipo SparseMS a af_SparseMS	85
Figura 34. Función que realiza la factorización de Cholesky simbólica	86
Figura 35. Cholesky: función de dispositivo sparse_fill.	87
Figura 36. Cholesky: función de dispositivo <i>Chol_sparse_1</i>	88
Figura 37. Cholesky: función de dispositivo <i>Chol_sparse_2</i>	90
Figura 38. Cholesky: función de dispositivo <i>Chol_sparse_3</i>	91
Figura 39. Cholesky: función de dispositivo <i>Chol_sparse_4</i>	92
Figura 40. Cholesky: función de dispositivo <i>Chol_sparse_5</i>	92
Figura 41. Cholesky: función de dispositivo <i>Chol_sparse_6</i>	93
Figura 42. Cholesky: Kernel para trabajar sobre una matriz en formato CSC	94
Figura 43. LDL ^T : función de dispositivo <i>sparse_fill_c</i>	
Figura 44. LDL ^T : función de dispositivo <i>ldlt_sparse_1_c</i>	98
Figura 45. LDL ^T : función de dispositivo <i>ldlt_sparse_2_c</i>	99
Figura 46. LDL ^T : función de dispositivo <i>ldlt_sparse_3_c</i>	
Figura 47. LDL ^T : función de dispositivo <i>ldlt_sparse_4_c</i>	
Figura 48. LDL ^T : función de dispositivo <i>ldlt_sparse_5_c</i>	103
Figura 49. LDL ^T : Kernel para trabajar sobre una matriz en formato CSC	104
Figura 50. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul1.	106
Figura 51. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul2.	107
Figura 52. Gradientes conjugados: función de dispositivo atom_add_double	108
Figura 53. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul3.	108
Figura 54. Gradientes conjugados: Kernel para multiplicar matriz en formato CSC	por
vector	109
Figura 55 Estructura para una matriz en formato SKS	110

Figura 56. Estructura para una matriz en formato SKS con memoria en el dispositivo
Figura 57. Conversión de una matriz densa a una matriz en formato SKS
Figura 58. Función para convertir un objeto del tipo sparse_sks_ms a
af_sparse_sks_ms111
Figura 59. Cholesky: función de dispositivo <i>chol_sparse_1_sks</i>
Figura 60. Cholesky: función de dispositivo chol_sparse_2_sks
Figura 61. Cholesky: función de dispositivo <i>chol_sparse_3_sks</i>
Figura 62. Cholesky: función de dispositivo <i>chol_sparse_4_sks</i>
Figura 63. Cholesky: función de dispositivo <i>chol_sparse_5_sks</i>
Figura 64. Cholesky: función de dispositivo <i>chol_sparse_6_sks</i>
Figura 65. Cholesky: Kernel para trabajar sobre una matriz en formato SKS 118
Figura 66. LDL ^T : función de dispositivo <i>ldlt_sparse_1_sks</i> 119
Figura 67. LDL ^T : función de dispositivo <i>ldlt_sparse_2_sks</i>
Figura 68. LDL ^T : función de dispositivo <i>ldlt_sparse_3_sks</i>
Figura 69. LDL ^T : función de dispositivo <i>ldlt_sparse_4_sks</i>
Figura 70. LDL ^T : función de dispositivo <i>ldlt_sparse_5_sks</i>
Figura 71. LDL ^T : kernel para trabajar sobre una matriz en formato SKS125
Figura 72. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul1_sks
Figura 73. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul2_sks
Figura 74. Gradientes conjugados: función de dispositivo atom_add_double 129
Figura 75. Gradientes conjugados: kernel para multiplicar matriz en formato SKS por
vector
Figura 76. malla de elementos finitos: viga en voladizo
Figura 77. Malla de elementos finitos: soporte central de puente

RESUMEN

El presente trabajo de investigación se centra principalmente en el uso del procesador gráfico (GPU) para paralelizar instrucciones en la solución de sistemas de ecuaciones lineales obtenidos de problemas de análisis estructural, es decir, dicho esto, nos limitamos a la solución de sistemas que tienen como matriz de coeficientes (matriz de rigidez): una matriz dispersa, cuadrada, simétrica y definida positiva. Los métodos considerados son: la factorización de Cholesky, la factorización LDL^T, el método de los gradientes conjugados y el método de Gauss Jordan.

Se han usado los formatos de almacenamiento CSC (*Compressed Sparse Column*) y SKS (*Skyline Storage*) para comprimir la matriz, adaptándolos para almacenar matrices simétricas (en donde solo la parte triangular inferior es almacenada), ya que la matriz, en su forma normal, almacena una gran cantidad de ceros, que en general, son innecesarios, con esto se busca mejorar tanto el rendimiento como el uso de memoria.

Se ha usado el OpenFEM, un conjunto de herramientas de elementos finitos, para analizar un problema determinado y obtener así el sistema de ecuaciones, es decir, la matriz de rigidez, luego, se ha comprimido la matriz usando uno de los formatos de almacenamiento, para luego resolver el sistema usando el kernel correspondiente. Los kernels (funciones que se ejecutan sobre la GPU) se han escrito usando la interfaz de OpenCL en C++, pero ejecutados por comodidad en Julia, mediante el uso de bibliotecas compartidas. Para poder trabajar con la memoria directamente en la GPU se usó la biblioteca de ArrayFire en C++, y el paquete ArrayFire.jl en Julia. Se ha medido la memoria utilizada y el tiempo de ejecución para diferentes órdenes de la matriz, y finalmente, se han preparado tablas para comparar resultados.

Palabras clave: Técnicas numéricas, procesador gráfico, GPU, Sistemas de ecuaciones lineales, análisis estructural.

ABSTRACT

The present research work focuses mainly on the use of the graphic processor (GPU) to parallelize instructions in the solution of systems of linear equations obtained from structural analysis problems, i.e. with that said, we limit ourselves to the solution of systems that have as a coefficient matrix (stiffness matrix): a dispersed, square, symmetric and positive-definite matrix. The methods considered are: Cholesky factorization, LDL^T factorization, the conjugate gradient method and the Gauss Jordan method.

The CSC (Compressed Sparse Column) and SKS (Skyline Storage) storage formats have been used to compress the matrix, adapting them to store symmetric matrices (where only the lower triangular part is stored), since the matrix, in its normal form, It stores a large number of zeros, which are generally unnecessary, this seeks to improve both performance and memory usage.

The OpenFEM, a set of finite element tools, has been used to analyze a given problem and thus obtain the system of equations, i.e. the stiffness matrix, then the matrix has been compressed using one of the storage formats, then solve the system using the corresponding kernel. The kernels (functions that run on the GPU) have been written using the OpenCL interface in C ++, but executed for convenience in Julia, through the use of shared libraries. In order to work with the memory directly on the GPU, the ArrayFire library in C ++ was used, and the ArrayFire.jl package in Julia. The memory used and the execution time for different matrix orders have been measured, and finally, tables have been prepared to compare results.

Keywords: numerical techniques, graphic processor, GPU, system of linear equations, structural analysis.

INTRODUCCIÓN

Los programas de análisis estructural más usados en la actualidad, usan el método de elementos finitos (MEF), método que consiste en discretizar sistemas continuos en elementos más simples, que se conectan mediante nodos (puntos en común), al analizar un problema con este método, se genera un **sistema de ecuaciones lineales** (**SEL**) que tiene como incógnitas los desplazamientos de los nodos, resolver este sistema de ecuaciones podría volverse un problema cuando el número de ecuaciones o incógnitas es considerablemente grande.

Se han propuesto infinidad de métodos para resolver sistemas de ecuaciones lineales de la forma Ax = b, donde A es la matriz de coeficientes del sistema, b es el vector de elementos constantes y x es el vector incógnita, el más general consiste en obtener la inversa de la matriz A, después del cual podemos determinar el vector incógnita mediante la expresión $x = A^{-1}b$, sin embargo, a medida que el orden de A crece, calcular su inversa se vuelve complicado. Es por esto, que el uso de métodos alternativos resulta más adecuado, la principal razón es que la gran mayoría de estos métodos utilizan mucho menos recursos computacionales que el que se utilizaría para invertir la matriz, sin embargo, no todos los métodos son aplicables a un sistema de ecuaciones en particular, ya que el tipo de sistema de ecuaciones varía según el tipo de problema que se resuelva, para problemas de análisis estructural que nos concierne, la matriz de coeficientes (matriz de rigidez) es dispersa, cuadrada, simétrica y definida positiva; estas tres últimas características la hacen adecuada a la aplicación de la factorización de Cholesky y el método de los gradientes conjugados, métodos que hemos considerado en este trabajo de investigación, además, está el método LDL^T, método que no exige que la matriz sea definida positiva, y adicionalmente, se ha usado el método de Gauss Jordan, método más flexible en cuanto a las exigencias para la matriz de coeficientes aplicada a la matriz en formato denso (sin comprimir).

Otra característica de la matriz de rigidez que se puede explotar es su forma dispersa, es decir, la gran cantidad de ceros que puede llegar a tener, esto, desde el punto de vista computacional, puede aprovecharse almacenando en memoria solo los elementos distintos de cero o por lo menos almacenando ceros en la menor medida posible, los formatos de almacenamiento considerados en este trabajo de investigación son: el formato denso (formato por defecto en donde todos los n² elementos de la matriz son almacenados), el formato CSC (Compressed Sparse Column) y el SKS (Skyline Storage); adaptando estos dos últimos para almacenar una matriz simétrica, es decir, solo la parte triangular inferior de la matriz es almacenada.

El desarrollo de la tecnología ha logrado simplificar significativamente el trabajo de resolver SEL's, sin embargo, paralelo a esto, los problemas también son cada vez más complejos, uno de los principales enfoques que se desarrollan en la actualidad es el uso de **procesadores gráficos** (GPU's) para propósitos generales, nuestro propósito en este trabajo de investigación, fue el de resolver SEL's generados de un análisis estructural, buscando principalmente aprovechar su paralelismo. Durante alguna etapa en la aplicación de los métodos para resolver SEL's citados anteriormente, nos topamos con ciertas operaciones que pueden realizarse independientemente, es decir, no dependen una de la otra, esto permite la posibilidad de ejecutar instrucciones en paralelo, para este fin, la programación del procesador gráfico resulta de gran ayuda, sin embargo, no todas las instrucciones pueden hacerse en paralelo, dependiendo del método, hay instrucciones que imperativamente deben ejecutarse una después de otra, es aquí en donde se incluye el concepto de heterogeneidad, que consiste básicamente en el uso combinado de la **GPU** y la **CPU**, la GPU principalmente cuando se requiera o se pueda ejecutar instrucciones en paralelo y la CPU cuando las instrucciones deban ejecutarse en serie.

En resumen, lo que se ha buscado en el presente trabajo de investigación, es mediante la programación del procesador gráfico (GPU), proponer kernels (funciones que se ejecutan usando la GPU) que resuelvan un sistema de ecuaciones lineales generado de un análisis estructural, usando los métodos y los formatos de almacenamiento citados anteriormente, y probar su funcionamiento. Para determinar cuál técnica es más efectiva, se ha medido la memoria utilizada en MB (Mega Bytes) y el tiempo de



CAPÍTULO I

PROBLEMA

1.1. Descripción del problema

Dentro del análisis de estructuras, casi siempre nos encontrarnos con problemas que consisten en encontrar el valor de ciertas variables, un ejemplo típico es el sistema de ecuaciones lineales $\mathbf{K}\mathbf{u} = \mathbf{F}$ derivados de un análisis por rigideces o elementos finitos, donde \mathbf{K} es la matriz de rigidez del sistema, \mathbf{F} el vector de fuerzas y \mathbf{u} el vector de desplazamientos que en este caso es la incógnita que tenemos que encontrar.

La solución de un **Sistema de ecuaciones lineales** (**SEL**) depende en gran medida de la capacidad de los ordenadores, como también del método que se emplee. El método más directo sería invertir la matriz, con la que se obtendría la solución directamente con $\boldsymbol{u} = \boldsymbol{K^{-1}F}$, para ordenes pequeñas de la matriz, calcular su inversa no es un problema y puede obtenerse relativamente fácil, el problema comienza cuando el orden de la matriz crece, ya que a la par con esto, los recursos computacionales necesarios también aumentan y los recursos disponibles podrían no ser suficientes.

Considerando inviable la opción de invertir la matriz **K** por el gran tamaño que puede llegar a tener, se usan métodos alternativos para obtener el vector **u** que satisface el sistema; básicamente estos métodos se dividen en: métodos directos (Gauss-Jordan, factorización LU, LDL^T, Cholesky, etc.) y métodos iterativos (Jacobi, Gauss-Seidel, relajación, minimización, etc.); algunos de estos métodos exigen que la matriz tenga características particulares, como ser simétrica, ser diagonalmente dominante, ser definida positiva, etc. Es decir, no todos se pueden aplicar a un problema de análisis estructural. La matriz de rigidez es cuadrada, dispersa, simétrica y definida positiva;

lo que la hace adecuada a la aplicación de la factorización de Cholesky y el método de los gradientes conjugados, el método LDL^T solo exige que la matriz sea simétrica y no necesariamente definida positiva, por tal también es adecuado para usarse en este tipo de problemas, de los métodos generales que no tienen limitaciones en las características de la matriz de rigidez, el método de Gauss Jordan podría ser el más idóneo para utilizarse. Los recursos computacionales necesarios que demanda la utilización de los métodos citados anteriormente son mucho menores a los que se necesitarían si quisiéramos invertir la matriz.

La matriz de rigidez también es dispersa, es decir contiene una gran cantidad de ceros; al almacenar y resolver la matriz sin tomar en cuenta esto también resulta en un problema, ya que mientras más sean los ceros, el uso de la memoria resulta deficiente y a la par con esto las operaciones realizadas sobre estos ceros también son un desperdicio; existen formatos de almacenamiento que hacen posible almacenar solo los elementos distintos de cero u otros en donde el almacenamiento de ceros es mínimo, formatos como el CSR (*Compressed Sparse Row*), CSC (*Compressed Sparse Column*), SKS (*Skyline Storage*), son ideales para este fin, además pueden adaptarse para almacenar matrices simétricas en donde solo una de las partes triangulares de la matriz es almacenada, así, el almacenamiento se reduciría considerablemente y las operaciones realizadas sobre la matriz serían más eficientes, sin embargo, el uso eficiente de memoria no necesariamente conlleva a una mejora en el rendimiento, esto nos sugiere que debería tenerse cuidado en almacenar la matriz en un formato que nos permita un acceso a memoria eficiente.

Al programar de forma tradicional, también nos encontramos con la imposibilidad de paralelizar instrucciones, o por lo menos, la imposibilidad de hacerlo a grandes escalas. Los métodos mencionados anteriormente, en su procedimiento, contienen una gran cantidad de operaciones que pueden realizarse independientemente, un ejemplo simple sería la suma de vectores, al sumar elemento a elemento, cada suma es independiente de la otra y de ser posible podrían hacerse en paralelo, un código estándar realizaría esta suma por medio de un bucle en donde las sumas se realizarían en serie. El desarrollo de la tecnología ha hecho posible paralelizar instrucciones, el

uso de procesadores gráficos nos ofrece esa posibilidad, aunque estos no hayan sido creados inicialmente para este fin.

En resumen, el problema de resolver sistemas de ecuaciones lineales se centra en, la búsqueda del método más idóneo, la utilización de un formato de almacenamiento que optimice mejor el uso de memoria y finalmente, la necesidad de paralelizar instrucciones; es a estos tres problemas a los que intentamos buscar solución en este trabajo de investigación.

1.2. Formulación del problema

1.2.1. Problema general

¿Cuáles son las características de las técnicas numéricas sobre procesador gráfico, al resolver sistemas de ecuaciones lineales en problemas de análisis estructural?

1.2.2. Problemas específicos

- ¿Cuál es el tiempo de ejecución de las técnicas numéricas sobre procesador gráfico, al resolver sistemas de ecuaciones lineales en problemas de análisis estructural?
- ¿Cuál es la memoria utilizada de las técnicas numéricas sobre procesador gráfico, al resolver sistemas de ecuaciones lineales en problemas de análisis estructural?

1.3. Objetivos

1.3.1. Objetivo general

❖ Determinar las características de las técnicas numéricas sobre procesador gráfico, al resolver sistemas de ecuaciones lineales en problemas de análisis estructural.

1.3.2. Objetivos específicos

- Determinar el tiempo de ejecución de las técnicas numéricas sobre procesador gráfico, al resolver sistemas de ecuaciones lineales en problemas de análisis estructural.
- ❖ Determinar la memoria utilizada de las técnicas numéricas sobre procesador gráfico, al resolver sistemas de ecuaciones lineales en problemas de análisis estructural.

1.4. Justificación

Un ingeniero civil se enfrenta en su día a día a la solución de problemas complejos que lo obligan a apoyarse en la tecnología, esto es necesario en casi todas las áreas de la ingeniería, un ejemplo típico es el cálculo de estructuras, que incluyen, análisis sísmicos, diseño de elementos, cálculo de cimentaciones, etc. Para todo esto siempre existe algún software especializado, sin embargo, la mayoría de programas informáticos que se utiliza a nivel local es extranjero, las desventajas, posiblemente la barrera de idioma, la complejidad en el uso y en muchos casos los precios poco accesibles. A pesar de lo dicho anteriormente, que casi para todo lo que quiera hacer un ingeniero existe algún software especializado, no debemos descartar el hecho de toparnos con problemas aislados, que no puedan ser abarcados o resueltos por los programas que se tenga a la mano, para estos casos, sería ideal tener a nuestro alcance herramientas que nos permitan desarrollar software especializado, anulando así, la necesidad de buscar ayuda externa. Cabe resaltar, que no sería descabellada la idea de desarrollar software a nivel local para resolver problemas típicos, las ventajas serían obvias, precios más cómodos, la eliminación de la barrera de idioma, el fácil manejo, el rápido acceso a soporte, entre otros.

En análisis estructural y en muchas otras ramas de la ingeniería, resolver sistemas de ecuaciones lineales es un paso obligado y probablemente el que más tiempo consume, este trabajo de investigación abarca a nivel computacional este tema, intentado explotar el paralelismo de una GPU, optimizar la memoria utilizada y mejorar el

rendimiento. A nivel local, las investigaciones que abarcan este tipo de temas son muy pocos, principalmente por la falta de interés, ya que de alguna manera se ven como algo secundario, algo un poco contradictorio con el hecho de que el uso de software es casi imperativo en todas las áreas de la ingeniería.

1.5. Limitaciones

- No se encontró bibliografía en nuestro idioma, de manuales o guías de programación en OpenCL, las guías de programación que vienen junto con el instalador están en inglés y son casi únicos, por lo menos durante la ejecución de la tesis, no se encontraron manuales o guías adicionales. Casos parecidos ocurrieron con la biblioteca de ArrayFire y la herramienta de elementos finitos OpenFEM.
- La incompatibilidad con la plataforma CUDA del dispositivo (GPU) usada para la tesis, que en principio era el entorno en donde se pensaba realizar la programación del procesador gráfico, obligó a la búsqueda de otras alternativas de programación, es así, que se elige la API de OpenCL para este propósito.
- La dificultad de programar el procesador gráfico y obtener resultados consistentes fue complicado en un principio, la obtención de errores en tiempo de ejecución, de compilación, de sintaxis y otros, fue algo recurrente al principio.
- Luego de superar las dificultades anteriores, pasamos a otro un poco más complicado de manejar, nos referimos al rendimiento, hablando en términos de memoria y tiempo de ejecución; a pesar de que se han mejorado los códigos paulatinamente, no se ha llegado a alcanzar el nivel de rendimiento deseado.
- Para órdenes de matrices considerablemente grandes, la escritura de archivos resultó un tanto lenta.

CAPÍTULO II MARCO TEÓRICO

2.1. Antecedentes

✓ (Sánchez Meza, 2000) Sánchez Meza, Roque Alberto. Procedimientos de gradiente conjugada en el análisis estructural. [ed.] Universidad Nacional de Ingeniería. Programa Cybertesis PERÚ. 2000.

Conclusiones: (Las más relevantes, propias de la referencia).

- En el análisis de pórticos planos, los valores de tiempo para el programa PT (método de Eliminación Gaussiana), son menores que aquellos que utilizan cualquier variante de gradiente conjugada, pero conforme aumenta el número de grados de libertad, la diferencia, en tiempo, se hace notable a favor de los métodos de gradiente conjugada.
- En el análisis tridimensional, las comparaciones de tiempo entre los programas SFE (análisis tridimensional de pórticos utilizando el método de Eliminación Gaussiana) y los métodos: gradiente conjugada con energía potencial (G.C. E.P.) y gradiente conjugada GC, son sustancialmente diferentes lo que nos ratifica que los procedimientos de gradiente conjugada produce los mismos resultados y además en un tiempo mucho menor. En todos los casos analizados los porcentajes de tiempo utilizados por los algoritmos de gradiente están muy por debajo del usado por el programa SFE.

- Los ejemplos tridimensionales no hicieron más que confirmar lo que ya se había avisorado en los ejemplos bidimensionales, es decir que el procedimiento empleado lleva a tiempos de cómputo más efectivos y rápidos, pero sin descuidar la precisión de los cálculos.
 - ✓ (Juan Carlos Camacho Puello, 2012) Juan Carlos Camacho Puello, Marlos de Jesús Romero Torres. Análisis estructural con el método de elementos finitos asistido por computadora. Cartagena de Indias: Universidad Tecnológica de Bolivar, 2012.
 - ✓ (Félix, 2015) **Félix, Miguel Vargas.** Factorización Cholesky simbólica. [En línea] 2015.

Conclusiones: La referencia (Juan Carlos Camacho Puello, 2012) se basa en los resultados de (Félix, 2015) aunque se detallan más en (Juan Carlos Camacho Puello, 2012) el cual no tiene conclusiones, estas son transcritas directamente de la referencia (Juan Carlos Camacho Puello, 2012).

- La descomposición de Cholesky es muy aplicada en programas estructurales que aplican el método de los elementos finitos, dado que la matriz de coeficientes es una matriz simétrica y definida positiva. BABEL 2.0¹ utiliza esta descomposición por ser fácil de implementar y por consumir menos recursos que algunos de los métodos tradicionales (pág. 58).
- En una prueba que llevo a cabo por (Félix, 2015) se utilizaron varios métodos de solución de ecuaciones para resolver una estructura compleja. La estructura (Figura 1) fue dividida en 264,250 elementos y 326,228 nodos. El tamaño de la matriz de rigidez es 978,684. Se resolvió el problema variando el número de procesadores usados. Como comparación, los resultados obtenidos usando el método del gradiente conjugado en paralelo

-

¹ Programa de análisis estructural producto de la tesis de referencia

para resolver el problema. El método del gradiente conjugado fue usado con y sin pre condicionamiento, el precondicionador usado es Jacobi. La tolerancia en la norma del gradiente usado para solvers iterativos es 1 x 10⁻⁵. En la Figura 2 y Figura 3, los valores entre paréntesis representan el número de procesadores usados en paralelo (pág. 64).

• En las imágenes se puede ver que la descomposición de Cholesky tarda menos tiempo en solucionar el problema que los otros métodos, pero es menos eficiente en términos de consumo de memoria. Cabe aclarar que el problema presenta una gran cantidad de variables que pueden generar que la descomposición de Cholesky parezca un método poco eficaz, pero debido a que BABEL 2.0 trabaja con un número de variables más reducido, el tiempo de resolución del sistema será mucho menor y habrá un consumo de memoria que no alterara el rendimiento normal del equipo (pág. 66).

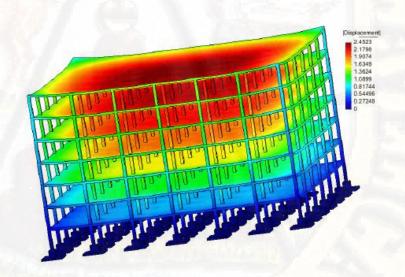


Figura 1. Ejemplo numérico de la deformación de una estructura (Félix, 2015)



Figura 2. Comparación de los tiempos de solución de los diferentes solvers (Félix, 2015)

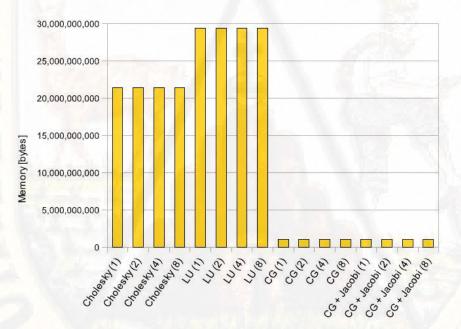


Figura 3. Comparación de la capacidad de memoria usada por los diferentes solvers paralelos (Félix, 2015)

2.2. Bases teóricas sobre el tema de investigación

2.2.1. Sistema de ecuaciones lineales (SEL)

Un sistema de ecuaciones lineales (SEL) es un conjunto de expresiones (ecuaciones) de primer grado que deben ser resueltas para un determinado número de variables.

El problema del SEL es uno de los más antiguos de la matemática y tiene una infinidad de aplicaciones, como en procesamiento digital de señales, análisis estructural, estimación, predicción y más generalmente en programación lineal así como en la aproximación de problemas no lineales de análisis numérico.²

En forma general un SEL tiene la siguiente forma:

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

 $a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$
...
 $a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m$

Donde x₁, x₂, x₃,...x_n, son incógnitas y a_{ij} junto con b_i {i=1, 2,3,..., m}, {j=1, 2,3,..., n} son constantes. Este conjunto de ecuaciones puede escribirse en forma matricial de la siguiente manera:

O de forma más compacta: $\mathbf{A}\mathbf{x}=\mathbf{b}$, en donde A_{mxn} es la matriz de coeficientes del sistema, \mathbf{x} un vector incógnita y \mathbf{b} un vector con términos constantes.

-

² Extraído de https://es.wikipedia.org/wiki/Sistema_de_ecuaciones_lineales el 06/03/2017

2.2.1.1.Métodos de solución de sistemas de ecuaciones lineales

Existen diversidad de métodos para la solución de SEL's estos se pueden dividir básicamente en dos grupos, los métodos directos y los métodos iterativos; los primeros suelen consistir en una cantidad finita de operaciones ordenadas que deben seguirse hasta determinar los valores de las incógnitas, los segundos no tienen una cantidad determinada de operaciones ya que dependen del número de iteraciones seguidas hasta llegar a una solución exacta o por lo menos aceptable, de aquí en adelante se considerará que la matriz **A** es cuadrada de orden "n" y por tanto el número de ecuaciones es igual al número de incógnitas.

2.2.1.1.1. Método de Gauss-Jordan

En matemática, la eliminación de Gauss-Jordan, llamada así debido a Carl Friedrich Gauss y Wilhelm Jordan, es un algoritmo del álgebra lineal para determinar las soluciones de un sistema de ecuaciones lineales, encontrar matrices e inversas. Un sistema de ecuaciones se resuelve por el método de Gauss cuando se obtienen sus soluciones mediante la reducción del sistema dado a otro equivalente en el que cada ecuación tiene una incógnita menos que la anterior. El método de Gauss transforma la matriz de coeficientes en una matriz triangular superior. El método de Gauss-Jordan continúa el proceso de transformación hasta obtener una matriz diagonal³.

El proceso generalizado consiste en n pasos (n es el orden de la matriz) y es el siguiente: considerar que en cada proceso el elemento de la diagonal principal es diferente de cero, si no es así podría intercambiarse las filas para cumplir esta condición.

Las operaciones se realizan sobre la matriz aumentada, es decir la matriz **A** unida al vector de elementos constantes, o sea:

_

³ Extraído de <u>https://es.wikipedia.org/wiki/Eliminación_de_Gauss-Jordan</u> el 06/03/17

Para un paso i, la intención es convertir toda la columna i en ceros, excepto el elemento correspondiente a la diagonal, para lograr esto, a cada fila (a excepción de la fila i), se le suma la fila i multiplicada por un factor, tal que el i-ésimo elemento de la fila se vuelva cero.

Luego de realizar los n pasos, debe obtenerse una matriz diagonal, es decir la matriz aumentada tendría la siguiente forma:

Luego de esto, encontrar cada x_i se vuelve sencillo mediante las divisiones: $x_i = b'_i/a'_{ii}$ con i = 1,2,3,...,n

2.2.1.1.2. Factorización LDLT

Es un método de factorización aplicable a matrices simétricas que consiste en descomponer la matriz A en tres submatrices, dos triangulares unitarias L y L^T , tal que una es la transpuesta de la otra y otra diagonal D tal que:

$$A = LDL^{T}$$

Con:

$$L = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ l_{21} & 1 & 0 & \dots & 0 \\ l_{31} & l_{32} & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & l_{n3} & \dots & 1 \end{bmatrix}$$

$$D = \begin{bmatrix} d_1 & 0 & 0 & \dots & 0 \\ 0 & d_2 & 0 & \dots & 0 \\ 0 & 0 & d_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & d_n \end{bmatrix}$$

Realizando la multiplicación obtenemos las siguientes fórmulas

$$d_{1} = a_{11}$$

$$l_{i1} = a_{i1}/d_{1}$$

$$d_{i} = a_{ii} - \sum_{k=1}^{i-1} l_{ik}^{2} d_{k}$$

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} d_k \quad d_j$$

La última expresión solo es válida cuando i>j ya que por ser una matriz triangular inferior si j>i entonces $l_{ij} = 0$.

Luego de conseguir cada elemento de L y D usando las expresiones anteriores, se procede a resolver los sistemas equivalentes: Lz=b, Dy=z y finalmente $L^Tx=y$, que es simple ya que L es triangular y D es diagonal.

2.2.1.1.3. Factorización de Cholesky

En matemática, la **factorización o descomposición de Cholesky** toma su nombre del matemático **André-Louis Cholesky**, quien encontró que una **matriz simétrica definida positiva** puede ser descompuesta como el producto de una matriz triangular inferior y la traspuesta de la matriz triangular inferior. La matriz triangular inferior es el triángulo de Cholesky de la matriz original positiva definida. El resultado de Cholesky ha sido extendido a matrices con entradas complejas. Es una manera de resolver sistemas de ecuaciones matriciales y se deriva de la factorización LU con una pequeña variación⁴.

Dado esto, si en el sistema **Ax=b**, la matriz A es simétrica y definida positiva puede descomponerse en la forma:

⁴ Extraída de <u>https://es.wikipedia.org/wiki/Factorización_de_Cholesky</u> el 06/03/2017

$$A = LL^T$$

Donde la matriz L es triangular inferior, luego al igual como se hace con la factorización LU la solución del sistema Ax=b, se reduce a resolver los sistemas equivalentes Ly=b, y $L^Tx=y$.

La forma de la matriz L es la siguiente:

$$L = \begin{bmatrix} l_{11} & 0 & 0 & \dots & 0 \\ l_{21} & l_{22} & 0 & \dots & 0 \\ l_{31} & l_{32} & l_{33} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & l_{n3} & \dots & l_{44} \end{bmatrix}$$

Al reemplazar en la expresión A=LL^T obtenemos las siguientes fórmulas:

$$l_{11=\sqrt{a_{11}}}$$

$$l_{i1} = a_{i1}/l_{11}$$

$$l_{ii} = a_{ii} - \sum_{k=1}^{i-1} l_{ik}^{2}$$

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} \quad l_{jj}$$

La última expresión solo es válida cuando i>j ya que por ser una matriz triangular inferior si j>i entonces $l_{ij} = 0$.

2.2.1.1.4. Método de los gradientes conjugados

En matemática, el método de **los gradientes conjugados** es un algoritmo para resolver numéricamente los sistemas de ecuaciones lineales cuyas matrices son simétricas y definidas positivas. Es un método iterativo, así que se puede aplicar a los sistemas dispersos que son demasiado grandes para ser tratados por métodos directos como la descomposición de

Cholesky. Tales sistemas surgen frecuentemente cuando se resuelve numéricamente las ecuaciones en derivadas parciales⁵.

Definimos la función $f: \mathbb{R}^n \to \mathbb{R}$ mediante:

$$f(x) = \frac{1}{2}x^T A x - b^T x$$

Ambos términos a la derecha representan un producto interno, ya que el rango de la función son números reales.

Si A es simétrica se cumple que:

$$\nabla f(x) = Ax - b$$

Si f(x) tiene un mínimo entonces este resultará de igualar su gradiente $\nabla f(x)$ a cero, es decir cuando Ax = b, por tanto resolver el sistema Ax = b es equivalente a buscar el mínimo de la función f.

Puede demostrarse que para que el mínimo exista entonces A debe ser definida positiva.

Tomando como primera aproximación el vector x_0 la primera dirección a usar será el negativo de la gradiente en éste punto es decir $p_0 = r_0 = b - Ax_0$, a partir de estos valores puede procederse con la iteración:

Entonces, teniendo una aproximación anterior x_k y una nueva dirección de avance p_k el valor de α_k que hace mínima la función en ésta dirección es:

$$\alpha_k = \frac{r_k^T p_k}{p_k^T A p_k}$$

La nueva aproximación se determina entonces con:

$$x_{k+1} = x_k + \alpha_k p_k$$

36

⁵ Extraído de <u>https://es.wikipedia.org/wiki/Método_del_gradiente_conjugado_</u>el 07/03/17

Multiplicando por A y considerando que $r_k = b - Ax_k$ se llega a:

$$r_{k+1} = r_k - \alpha_k A p_k$$

La siguiente dirección se determinará con:

$$p_{k+1} = r_{k+1} + \beta_k p_k$$

El valor de β_k que hace que las direcciones sean conjugadas es decir que el producto $p_k^T A p_{k+1}$ sea cero es:

$$\beta_k = -\frac{p_k^T A r_{k+1}}{p_k^T A p_k}$$

Se vuelve a determinar α_{k+1} para la nueva dirección y se continua el proceso hasta acercarnos todo que queramos a la solución del sistema.

A partir de una matriz A de rango n solo se pueden definir n vectores A-conjugados, por lo tanto el algoritmo de los gradientes conjugados garantiza la obtención de una solución en un máximo de n iteraciones (Vargas-Félix, 2015).

2.2.2. Uso de procesadores gráficos (GPU) en propósitos generales 2.2.2.1.GPGPU (Wikipedia, 2013)

GPGPU o General-Purpose Computing on Graphics Processing Units es un concepto reciente dentro de informática que trata de estudiar y aprovechar las capacidades de cómputo de una GPU.

Una **GPU** es un procesador diseñado para los cómputos implicados en la generación de gráficos 3D interactivos. Algunas de sus características (bajo precio en relación a su potencia de cálculo, gran paralelismo, optimización para cálculos en coma flotante), se consideran atractivas para su uso en aplicaciones fuera de los gráficos por computadora, especialmente en el ámbito científico y de simulación. Así, se han desarrollado técnicas para la implementación de simulaciones de fluidos, bases de datos, algoritmos de clustering, etc.

2.2.2.1.1. Modelo de programación GPU

Debido a las diferencias fundamentales entre las arquitecturas de la GPU y la CPU, no cualquier problema se puede beneficiar de una implementación en la GPU. En concreto, el acceso a memoria plantea las mayores dificultades. Las CPU están diseñadas para el acceso aleatorio a memoria. Esto favorece la creación de estructuras de datos complejas, con punteros a posiciones arbitrarias en memoria. En cambio, en una GPU, el acceso a memoria está mucho más restringido. Por ejemplo, en un procesador de vértices (la parte de una GPU diseñada para transformar vértice en aplicaciones 3D), se favorece el modelo scatter, en el que el programa lee en una posición predeterminada de la memoria, pero escribe en una o varias posiciones arbitrarias. En cambio, un procesador de píxeles, o fragmentos, favorece el modelo gather, pudiendo el programa leer de varias posiciones arbitrarias, pero escribir en solo una posición predeterminada.

La tarea del diseñador de algoritmos GPGPU consiste principalmente en adaptar los accesos a memoria y las estructuras de datos a las características de la GPU. Generalmente, la forma de almacenar datos es en un buffer 2D, en lugar de lo que normalmente sería una textura. El acceso a esas estructuras de datos es el equivalente a una lectura o escritura de una posición en la textura. Puesto que generalmente no se puede leer y escribir en la misma textura, si esta operación es imprescindible para el desarrollo del algoritmo, este se debe dividir en varias pasadas.

Pese a que cualquier algoritmo que sea implementable en una CPU lo es también en una GPU, esas implementaciones no serán igual de eficientes en las dos arquitecturas. En concreto, los algoritmos con un alto grado de paralelismo, sin necesidad de estructuras de datos complejas y con una alta intensidad aritmética son los que mayores beneficios obtienen de su implementación en la GPU.

2.2.2.Programación de GPU en Julia (Juliacomputing.com, s.f.)⁶

Los problemas científicos se han resuelto tradicionalmente usando poderosos grupos de CPU's homogéneos, conectados en una variedad de topologías de red. Sin embargo, el número de superordenadores que emplean aceleradores ha ido en constante aumento, La última lista Top 500 publicada en SC'15 muestra que el número de supercomputadores que emplean aceleradores ha aumentado a 109.

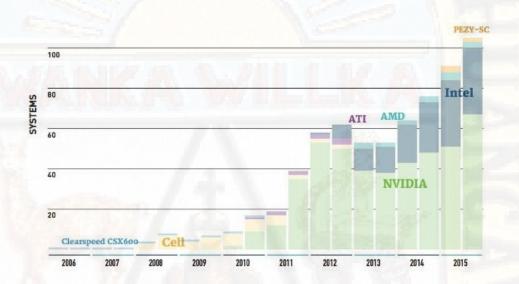


Figura 4. Número de supercomputadores que emplean aceleradores, extraída de (Juliacomputing.com, s.f.)

Los aceleradores que se emplean en la práctica son en su mayoría unidades de procesamiento gráfico (**GPU** por sus siglas en ingles), Xeon Phis y FPGAs. Estos aceleradores se aprovechan de las arquitecturas de núcleos múltiples que pueden utilizarse para explotar y cristalizar el paralelismo. Sin embargo el problema tradicional al usar las GPU's y otros aceleradores ha sido la facilidad (o falta de ello) a la hora de programarlos. Para este fin la corporación **NVIDIA** diseñó la actual Arquitectura Unificada de Dispositivos de Cómputo (CUDA por sus siglas en inglés) que incorpora en C una interface para la programación

39

⁶ La referencia originalmente en inglés, ha sido traducida para poder ser añadida en el marco teórico

científica y de propósito general. Esta fue una considerable mejora sobre sistemas previos como DirectX u OpenGL que requerían habilidades avanzadas de programación gráfica. Sin embargo, CUDA todavía se caracteriza por una baja curva de productividad, para programadores que tienen que afinar sus aplicaciones para diferentes dispositivos y algoritmos. En este contexto, la programación interactiva sobre GPU's proporcionaría grandes beneficios a científicos y programadores que no solo desearían poner a prueba sus aplicaciones, sino implementarlas con pocos o ningún cambio en su código.

2.2.2.2.1. Julia sobre GPU's

Julia ofrece a los programadores la capacidad de codificar interactivamente en la GPU. Hay varias librerías incluidas en Julia, dando a los usuarios acceso a acelerados BLAS, FFT's, rutinas para matrices dispersas y solucionadores, y aprendizaje profundo. Con una combinación de estos paquetes, los programadores pueden desarrollar interactivamente Kernels (funciones) GPU personalizadas. Un ejemplo es el de los gradientes conjugados, el cual fue comparado a continuación.

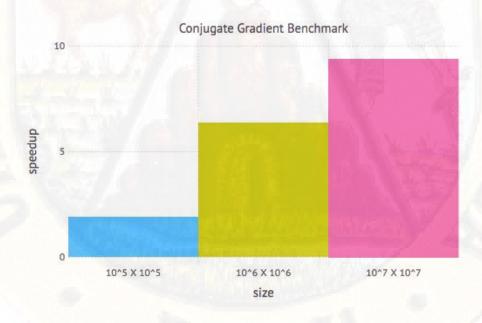


Figura 5. Comparación de la aceleración del gradiente conjugado en diferentes tamaños, extraído de (Juliacomputing.com, s.f.)

Sin embargo, se podría argumentar que las librerías de nivel inferior no aumentan en modo alguno la productividad del programador, ya que implican trabajar con interfaces de funciones oscuras. En tal caso, sería conveniente tener una interfaz limpia para matrices en la GPU con una librería estándar conveniente que pueda trabajar con matrices. Cada operación se sintonizaría con el dispositivo en cuestión para lograr un gran rendimiento. La gente de **ArrayFire** ha reunido una biblioteca de código abierto de alta calidad para trabajar en problemas científicos con GPU's.

2.2.2.2.2. ArrayFire.jl

ArrayFire.jl es un conjunto de enlaces de Julia a la biblioteca de ArrayFire, está diseñado para imitar la librería estándar de Julia en su versatilidad y facilidad de uso, proporcionando una interfaz fácil pero de potente formación que apunta hacia ubicaciones en la memoria de la GPU.

La capacidad de envío múltiple y programación genérica de Julia, hace posible que los usuarios escriban código matemático natural y aprovechen de manera transparente las GPU's para obtener rendimiento. Esto se hace definiendo un tipo **AFArray** como un subtipo de **AbstractArray**. **AFArray** que ahora actúa como una interfaz hacía una matriz en la memoria del dispositivo. Un conjunto de funciones se importan desde Base Julia y se distribuyen a través del nuevo tipo AFArray. Así los usuarios pueden ser capaces de escribir códigos en Julia que se ejecutan en la CPU, pudiendo ejecutarlos en la GPU con muy pocos cambios de código. Adicionalmente a las funciones que imitan la librería estándar de Julia, **ArrayFire.jl** proporciona poderosas funciones en procesamiento de imágenes y la visión por computador, entre otros.

2.3. Bases conceptuales

2.3.1. Programación de la GPU en OpenCL

Una función escrita para ejecutarse usando un dispositivo como un procesador gráfico (GPU), es comúnmente denominada como kernel, **OpenCL** brinda una interfaz fácil de usar para escribir estos kernels, un kernel es lanzado para una cierta cantidad de grupos de trabajo (work groups), cada una con una cantidad determinada de hilos (threads), estos hilos son los encargados de ejecutar una simple instrucción simultáneamente, la cantidad de grupos de trabajo y número de hilos por cada uno de estos se especifican antes de lanzar el kernel.

Existe un límite en el número de grupos de trabajo y de hilos según el dispositivo que se esté usando, estos están directamente relacionados con el número de multiprocesadores del dispositivo, es decir, mientras más multiprocesadores tenga un dispositivo, más instrucciones en paralelo podrá ejecutar y por lo tanto el tiempo de ejecución global será menor.

En la Figura 6, se muestra que tan rápido puede ejecutarse una instrucción con una GPU de cuatro multiprocesadores en comparación a una con solo dos multiprocesadores, el término escalabilidad automática, se refiere a que solo se especifica el número de grupos de trabajo (cada grupo de trabajo es un block en la figura) y estos son repartidos automáticamente de acuerdo al número de multiprocesadores disponibles.

Un kernel es definido mediante el uso del especificador __kernel, cada hilo que ejecuta el kernel, tiene un único ID global y un único ID dentro del grupo al que pertenece, que son obtenidos usando las funciones get_global_id y get_local_id respectivamente. Del mismo modo cada grupo de trabajo tiene un único ID que es obtenido usando la función get_group_id.

Para especificar que un argumento de la función es un vector en la memoria del dispositivo, debe anteponerse el prefijo __global o __local antes del tipo dentro de la definición de la función, ya sea que esté en la memoria global o local

respectivamente, en la Figura 7 se muestra un ejemplo que suma dos vectores y guarda el resultado en un tercer vector dentro de la memoria del dispositivo

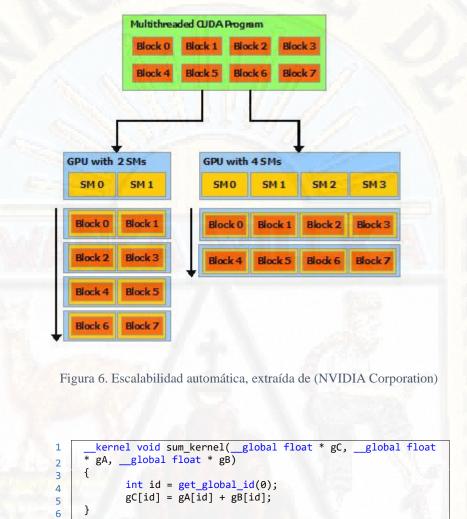


Figura 7. Ejemplo de un Kernel que suma dos vectores

2.3.1.1.Ejecutando un kernel usando la biblioteca de ArrayFire

ArrayFire nos permite ejecutar un kernel con la memoria directamente en la GPU, sin necesidad de trasferir datos entre la memoria central y el dispositivo y viceversa, el manual de **ArrayFire** (ArrayFire company, s.f.), en su sección *Interoperability with OpenCL* (interoperabilidad con OpenCL) nos muestra como:

Por defecto, ArrayFire administra su propio contexto, cola, memoria y crea ID's personalizados para los dispositivos, por tanto se deben seguir algunos pasos para poder integrar los kernels de OpenCL personalizados.

Si los kernels pueden operar en la misma cola que ArrayFire, deben seguirse los siguientes pasos:

- 1. Agregar un *include* para *af/opencl.h* en el proyecto.
- Obtener el contexto, el dispositivo y la cola de OpenCL usado por ArrayFire.
- 3. Obtener las referencias *cl_mem* a los objetos *af::array*.
- 4. Cargar, construir y usar los kernels
- 5. Retornar el control de la memoria af::array a ArrayFire.

Notar que ArrayFire usa una cola ordenada, por tanto cuando ArrayFire y los kernels operan en la misma cola, no es necesario realizar ningún tipo de sincronización.

En la Figura 8, se muestra un ejemplo completo, donde se muestra todo éste proceso, notar que el kernel es escrito directamente en el cuerpo de la función *main*, esto no será conveniente cuando el cuerpo del kernel sea demasiado grande, por este motivo es común escribir los kernels en un archivo de extensión cl, para luego ser cargados en un objeto del tipo *char**.

2.3.1.2.Escribir funciones usando kernels personalizados para usarse en Julia usando el paquete ArrayFire.jl

ArrayFire tiene una biblioteca compartida con funciones que pueden llamarse desde otro lenguaje de programación como Julia, el paquete de Julia **ArrayFire.jl** cumple este cometido y nos permite usar las funciones de ArrayFire en Julia, de esta forma es también posible escribir funciones personalizadas que usan kernels OpenCL personalizados para usarlos en Julia.

```
#include <arrayfire.h>
      // 1. Agregando el include para af/opencl.h en el proyecto
     #include <af/opencl.h>
3
4
     int main() {
5
             size_t length = 10;
6
             // creando vectores ArrayFire:
             af::array A = af::randu(length, f32);
             af::array B = af::constant(0, length, f32);
8
9
               ... operaciones ArrayFire adicionales aquí
             // 2. Obtener el contexto, el dispositivo y la cola de OpenCL
10
11
             // usado por ArrayFire
12
             static cl_context af_context = afcl::getContext();
             static cl_device_id af_device_id = afcl::getDeviceId();
13
14
             static cl_command_queue af_queue = afcl::getQueue();
15
             // 3. Obteniendo las referencias cl_mem a los objetos af::array
             cl mem * d A = A.device<cl mem>();
16
             cl_mem * d_B = B.device<cl_mem>();
17
18
             // 4. Cargando, construyendo y usando los kernels.
             //por legibilidad se ha omitido la comprobación de errores
19
20
             int status = CL_SUCCESS;
21
             // Un kernel de copia simple, usando la sintaxis de C++11 para
22
             // cadenas en múltiples líneas
             const char * kernel_name = "copy_kernel";
23
             const char * source = R"(
24
25
              void __kernel
              copy_kernel(__global float * gA, __global float * gB)
26
27
28
                  int id = get_global_id(0);
29
                  gB[id] = gA[id];
30
31
             // Creando el programa, construyendo el ejecutable y extrayendo
32
33
             //el punto de entrada para el kernel.
34
             cl_program program = clCreateProgramWithSource(af_context,
                     1, &source, NULL, &status);
35
36
             status = clBuildProgram(program, 1, &af_device_id, NULL,
                     NULL, NULL);
37
38
             cl_kernel kernel = clCreateKernel(program, kernel_name, &status);
39
             // Estableciendo los argumentos y ejecutando el kernel
             clSetKernelArg(kernel, 0, sizeof(cl_mem), d_A);
40
             clSetKernelArg(kernel, 1, sizeof(cl_mem), d_B);
41
42
             clEnqueueNDRangeKernel(af_queue, kernel, 1, NULL, &length, NULL,
43
                     0, NULL, NULL);
             // 5. Retornando el control de la memoria af::array a ArrayFire
44
             A.unlock();
45
46
             B.unlock();
             // ... continuar con las operaciones ArrayFire
47
48
             // Ya que los punteros del dispositivo, d_x y d_y, fueron
49
             // devueltas al control de ArrayFire mediante la función unlock,
50
             // no se necesita liberarlos usando clReleaseMemObject()
51
             return 0;
```

Figura 8. Ejemplo de la ejecución de un kernel usando la biblioteca de ArrayFire (ArrayFire company, s.f.)

En Julia una matriz ArrayFire es de un tipo *AFArray*, un objeto de este tipo tiene un campo .*arr* cuyo equivalente en ArrayFire es un objeto *af_array*, que es un puntero a una ubicación en la memoria del dispositivo.

Para llamar una función ArrayFire desde Julia deben conocerse los equivalentes de los tipos usados en la declaración de la función en ArrayFire con los tipos en Julia, estos tipos son principalmente dos:

Tipo en ArrayFire	Tipo en Julia (ArrayFire.jl)	
af_array	af_array	
af_array*	Ptr{af_array}	

Tabla 1. Correspondencia de tipo entre objetos en ArrayFire y en Julia

Otro tipo de correspondencias son los existentes entre C y Julia, los cuales pueden verse en la sección *Calling C and Fortran Code* de la documentación de Julia (JuliaLang, s.f.).

La llamada se hace de la misma forma en la que se llaman funciones C, por ejemplo, la función de ArrayFire que suma dos vectores *af_array* cuya declaración se muestra en la Figura 9, es llamada por ArrayFire.jl en Julia usando la función mostrada en la Figura 10.

```
AFAPI af_err af_add (af_array *out, const af_array lhs, const af_array rhs, const bool batch);
```

Figura 9. Declaración de la función af_add de ArrayFire

Por lo anterior dicho, si queremos usar una función ArrayFire en Julia, es necesario escribir las funciones con argumentos del tipo *af_array* o *af_array**, combinándolos con cualquier otro tipo C que tenga su equivalente en Julia.

Figura 10. Función add de ArrayFire.jl⁷

El objeto *af_array* de ArrayFire es similar al objeto *array*, puede obtenerse el objeto *af_array* a partir de un objeto *array* por medio de la función *get*, es decir, si A es una matriz del tipo *array*, obtendremos el tipo *af_array* con *A.get()*.

⁷ https://github.com/JuliaGPU/ArrayFire.jl/blob/master/src/wrap.jl

La Figura 11, muestra un ejemplo, usando el kernel de suma de vectores, de una función con argumentos del tipo af_array, viable para ser llamada desde Julia, esta función es llamada en Julia mediante la función de la Figura 12.

2.3.2. Formato de almacenamiento CSC

El formato de compresión CSC es el acrónimo de *Compressed Sparse Column* (columna dispersa comprimida), también denominada como CCS, acrónimo de *Compressed Column Storage* (Almacenamiento de columna comprimida), es un formato de compresión por columnas de una matriz dispersa, en donde solo los elementos distintitos de cero de la matriz son almacenados en memoria como un vector, adicionalmente, para poder identificar la posición de estos elementos en la matriz original, se almacenan otros dos vectores que contienen la extensión de las columnas y los índices de fila de cada elemento. Este formato es análogo al formato CSR (Compressed Sparse Row) o CRS (Compressed Row Storage) que almacena la matriz por filas.

Sea la matriz:

$$M = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 5 & 5 & 0 & 0 & 0 \\ 3 & 0 & 3 & 3 & 0 & 0 \\ 0 & 3 & 0 & 9 & 1 & 1 \\ 0 & 0 & 9 & 6 & 5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8 \end{bmatrix}$$

Esta matriz almacenada en formato CSC será representada por 3 vectores:

 Vector A de elementos diferentes de cero, se escriben solo los elementos distintos de cero columna a columna, es decir:

$$A = \begin{bmatrix} 1 & 2 & 3 & 5 & 3 & 5 & 3 & 9 & 3 & 9 & 6 & 1 & 5 & 1 & 8 \end{bmatrix}$$

```
1
       void sumaaf(af_array* out , af_array dA, af_array dB) {
3
                //para guardar el resultado
4
                af_array dC;
                af_copy_array(&dC, dA);
5
6
                // 2. Obteniendo el dispositivo, el contexto y la cola usada por ArrayFire
static cl_context af_context = afcl::getContext();
8
                static cl_device_id af_device_id = afcl::getDeviceId();
10
                static cl_command_queue af_queue = afcl::getQueue();
11
                dim_t _order[AF_MAX_DIMS];
12
                af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
13
14
                size_t order = _order[0];
15
                int status = CL_SUCCESS;
16
17
                // 3. Obteniendo las referencias cl_mem a los objetos af_array
18
                cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
CL_MEM_READ_ONLY, sizeof(float) * order,
19
20
                NULL, &status);
af_get_device_ptr((void**)d_A, dA);
21
23
24
                cl_mem *d_B = (cl_mem*)clCreateBuffer(af_context,
25
                          CL_MEM_READ_ONLY, sizeof(float) * order,
                NULL, &status);
af_get_device_ptr((void**)d_B, dB);
26
27
28
                29
30
                NULL, &status);
af_get_device_ptr((void**)d_C, dC);
31
32
33
                // 4. Cargando, construyendo y usando el kernel
                      por legibilidad omitiremos la comprobación de errores
                // Escribiendo el kernel.
                const char * kernel_name = "sum_kernel";
const char * source = R"(
37
38
39
                void <u>kernel</u>
                sum_kernel(__global float * gC, __global float * gA, __global float * gB)
40
41
                    int id = get_global_id(0);
gC[id] = gA[id]+gB[id];
42
43
44
45
46
                // Creando el programa, construyendo el ejecutable y extrayendo el punto de entrada
                // para el kernel.
                cl_program program = clCreateProgramWithSource(af_context, 1, &source,
                          NULL, &status);
50
                status = clBuildProgram(program, 1, &af_device_id, NULL, NULL);
51
                cl_kernel sumkernel = clCreateKernel(program, kernel_name, &status);
52
                // estableciendo los argumentos y ejecutando el kernel
                clSetKernelArg(sumkernel, 0, sizeof(cl_mem), d_C);
clSetKernelArg(sumkernel, 1, sizeof(cl_mem), d_A);
53
54
                clSetKernelArg(sumkernel, 2, sizeof(cl_mem), d_B);
clEnqueueNDRangeKernel(af_queue, sumkernel, 1, NULL, &order, NULL, 0,
55
56
                          NULL, NULL);
57
58
                // 5. Retornando el control de la memoria af_array a ArrayFire
59
60
                af_unlock_array(dA);
                af_unlock_array(dB);
61
62
                af_unlock_array(dC);
63
64
                //copiando el resultado al objeto af_array de salida
65
                af_copy_array(out, dC);
```

Figura 11. Función con argumentos af_array que pueden ser llamados desde Julia

```
function summaf(lhs::AFArray(Float32,1), rhs::AFArray(Float32,1))

out = ArrayFire.RefValue.af_array(0)

coall((:summaf,"ruta/al/dll"),

af_err,(Ftr{af_array}, af_array, af_array).out.lhs.arr, rhs.arr)

AFArray-Float32,1(out[])

end
```

Figura 12. Código que llama en Julia a la función de la Figura 11

 Vector colA donde cada elemento indicará la ubicación en A del primer elemento distinto de cero en cada columna de M, es decir:

```
colA = [1 \ 4 \ 6 \ 9 \ 12 \ 14 \ 16]
```

 Vector rowA donde se guardará los índices de fila de todos los elementos, es decir:

$$rowA = \begin{bmatrix} 1 & 2 & 3 & 2 & 4 & 2 & 3 & 5 & 3 & 4 & 5 & 4 & 5 \end{bmatrix}$$

Notar que el primer elemento de colA siempre será 1, y el último elemento representa la ubicación del primer elemento distinto de cero en A de la columna 7 de M (si existiera), este valor que a primera impresión no sirve, guarda la cantidad de elementos distintos de cero de M más uno, en algunos lenguajes de programación este valor necesita ser almacenado explícitamente. Si la indexación se basara en cero, este último elemento representaría exactamente la cantidad de elementos distintos de cero de M, puede obviarse este valor si n es necesario.

2.3.2.1. Formato de compresión CSC para matrices simétricas

Para matrices simétricas solo es necesario almacenar la parte triangular superior o inferior de la matriz, para este fin, puede usarse el formato de almacenamiento CSC asumiendo que una de las partes triangulares de la matriz tiene todos los elementos iguales a cero.

Si asumimos que los elementos de la diagonal de la matriz son todos diferentes de cero, puede obviarse el almacenamiento de sus índices en rowA por ser obvias.

Sea la matriz simétrica:

$$M = \begin{bmatrix} 1 & 2 & 3 & 0 & 0 & 0 \\ 2 & 5 & 0 & 3 & 0 & 0 \\ 3 & 0 & 3 & 0 & 9 & 0 \\ 0 & 3 & 0 & 9 & 6 & 0 \\ 0 & 0 & 9 & 6 & 5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8 \end{bmatrix}$$

Esta matriz almacenada en formato CSC será representada por 3 vectores, considerar que solo se guardará la parte triangular inferior, es decir, la parte triangular superior exceptuando la diagonal se asume que tiene todos los elementos iguales a cero:

 Vector A de elementos diferentes de cero, se escriben solo los elementos distintos de cero columna a columna de la parte triangular inferior, es decir:

$$A = [1 \ 2 \ 3 \ 5 \ 3 \ 3 \ 9 \ 9 \ 6 \ 5 \ 8]$$

 Vector colA donde cada elemento indicará la ubicación en A del primer elemento distinto de cero en cada columna de M correspondiente a la parte triangular inferior, es decir:

$$colA = [1 \ 4 \ 6 \ 8 \ 10 \ 11]$$

 Vector rowA donde se guardará los índices de fila de todos los elementos en la parte triangular inferior, exceptuando los índices de fila de los elementos que coincidan con la diagonal, es decir:

$$rowA = [2 \ 3 \ 4 \ 5 \ 5]$$

2.3.3. Formato de almacenamiento SKS

El formato de almacenamiento SKS, acrónimo de *Skyline Storage* (Almacenamiento horizontal), es un formato de compresión, ideal para matrices dispersas en donde los elementos distintos de cero estén concentrados cerca de la

diagonal, en una matriz simétrica solo los elementos de la parte triangular superior o inferior son almacenados. Una forma de utilizar este formato sobre matrices simétricas es análogo al formato CSC visto anteriormente, sin embargo en este caso no será necesario un tercer vector de índices de fila ya que estos pueden obtenerse usando el segundo vector.

Sea la matriz M:

$$M = \begin{bmatrix} 1 & 2 & 3 & 0 & 0 & 0 \\ 2 & 5 & 0 & 3 & 0 & 0 \\ 3 & 0 & 3 & 0 & 9 & 0 \\ 0 & 3 & 0 & 9 & 6 & 0 \\ 0 & 0 & 9 & 6 & 5 & 0 \\ 0 & 0 & 0 & 0 & 0 & 8 \end{bmatrix}$$

Esta matriz almacenada en formato SKS será representada por 2 vectores, considerar que solo se guardará la parte triangular inferior columna a columna y que los elementos de la diagonal siempre serán guardados aunque alguno de ellos sean ceros:

 Vector A de elementos diferentes de cero, se busca en la parte triangular inferior, el último elemento distinto de cero en cada columna y se escriben los elementos desde la diagonal hasta este, esto significa que si en este tramo existen elementos iguales a cero, también deben escribirse, es decir:

$$A = \begin{bmatrix} 1 & 2 & 3 & 5 & 0 & 3 & 3 & 0 & 9 & 9 & 6 & 5 & 8 \end{bmatrix}$$

Adicionalmente para cada columna, el último elemento distinto de cero debe tener índice de fila mayor o igual al índice de fila del elemento correspondiente a la columna anterior, esto para evitar realizar factorizaciones simbólicas en el caso de realizar una factorización de Cholesky o LDL^T, si esto no se cumple, debe almacenarse algunos ceros de esta columna para cumplir con lo mencionado.

• Vector idxA donde cada elemento indicará la ubicación en A del primer elemento distinto de cero en cada columna de M (es decir el elemento diagonal) correspondiente a la parte triangular inferior, es decir:

$$idxA = [1 \ 4 \ 7 \ 10 \ 12 \ 13]$$

Hay otras formas de almacenar en este formato: está el análogo almacenamiento por filas, en donde podría guardarse solo la parte triangular superior; o en donde se almacena los elementos de la diagonal en un vector diferente como puede verse en (stuff, s.f.) y en (Watkins, 2002).

Como se dijo anteriormente este formato es ideal para cuando los elementos distintos de cero de la matriz se concentran cerca de la diagonal, de no ser este el caso, este formato podría no ser adecuado, ya que como se vio, se almacenan algunos ceros extras; o puede usarse algún método para reenumerar filas y columnas (sin perder la simetría), con tal de lograr que los elementos distintos de cero se concentren más cerca de la diagonal.

2.4. Definición de términos

ArrayFire: Biblioteca de software de alto rendimiento para computación paralela con una API fácil de usar. Su conjunto de funciones basado en matrices hace que la programación paralela sea más accesible.

ArrayFire.jl: Paquete de Julia que adecua la biblioteca de ArrayFire para usarse en este lenguaje.

Cola: traducción de *queue* en OpenCL, directamente relacionada con la cola de comandos (*command-queue*), que es una estructura de datos para coordinar la ejecución de los kernels en los dispositivos.

Contexto: traducción de *context* en OpenCL, se refiere a una definición que hace el procesador central para la ejecución de kernels.

Dispositivo: Se refiere a cualquier dispositivo de cómputo que pueda programarse con OpenCL, principalmente procesadores gráficos.

En sitio: Se refiere a que las operaciones de escritura se realizan sobre el mismo objeto de lectura, perdiéndose los datos anteriores.

Factorización de Cholesky: Método de descomposición de una matriz cuadrada simétrica y definida positiva en dos submatrices triangulares tal que una es la traspuesta de la otra.

Factorización LDL^T: Método de descomposición de una matriz cuadrada simétrica $\bf A$ en tres factores: dos matrices triangulares $\bf L$ y $\bf L^T$ tal que una es la transpuesta de la otra y una matriz diagonal $\bf D$ tal que $A = LDL^T$

Factorización simbólica: proceso simbólico de factorización para obtener la cantidad de memoria máxima necesaria al realizar la factorización de Cholesky o LDL^T.

Función atómica: función en OpenCL usada para que un hilo pueda escribir en la misma ubicación en memoria que otros hilos ejecutando una instrucción en paralelo.

Función de dispositivo: Función secundaria escrita en OpenCL para ejecutar instrucciones en paralelo usando el procesador gráfico y que es llamada desde un kernel principal.

GPU: Acrónimo de Graphics Processing Unit (en español unidad de procesamiento gráfico)

Grupo de trabajo: Traducción de *work group*, denominación de un bloque de hilos en OpenCL

Hilos: Traducción de *threads* en OpenCL, forma parte de un grupo de trabajo (*work group*), se refiere básicamente a un hilo de ejecución, es decir ejecutan una instrucción en paralelo a la par con otros hilos dentro del grupo.

Julia: Lenguaje de programación de alto nivel.

Kernel: Función escrita en OpenCL para ejecutar instrucciones en paralelo usando el procesador gráfico.

Memoria global: Traducción de *global memory* en OpenCL, se refiere a la memoria que reside en la memoria del dispositivo, que permite leer y escribir a todos los hilos de todos los grupos.

Memoria local: Traducción de *local memory* en OpenCL, se refiere a la memoria que reside en la memoria del dispositivo, pero asignado solo a un grupo de trabajo.

Método de Gauss-Jordan: Llamada también eliminación de Gauss-Jordan, es un método de solución de sistemas de ecuaciones lineales, que consiste en convertir la matriz por medio de operaciones entre sus filas en una matriz diagonal.

Método de los gradientes conjugados: Método iterativo de solución de sistemas de ecuaciones lineales que consiste en buscar el mínimo de la función $f: R^n \to R$: $f(x) = \frac{1}{2}x^TAx - b^Tx$

OpenCL: Acrónimo de Open Computing Language (en español lenguaje de computación abierto), interfaz de programación del procesador gráfico que se usa con C++.

OpenFEM: Conjunto de herramientas de elementos finitos de distribución libre.

SEL: Acrónimo de sistema de ecuaciones lineales

Sistema de ecuaciones lineales: conjunto de expresiones (ecuaciones) de primer grado que deben ser resueltas para un determinado número de variables.

Sistema triangular: Un sistema de ecuaciones lineales cuya matriz de coeficientes es triangular superior o inferior.

2.5. Variables

2.5.1. Variable independiente

Técnicas numéricas implementadas sobre procesador gráfico (GPU): las técnicas numéricas consideradas en el presente estudio son implementadas sobre procesador gráfico, entiéndase técnica como la manera de usar un método en particular, en nuestro caso, los métodos serían: el método de Gauss Jordan, la factorización de Cholesky, la factorización LDL^T y el método de los gradientes conjugados; y la forma sería la aplicación de estos métodos sobre matrices densas, en formato CSC y en formato SKS.

2.5.2. Variable dependiente

Solución de sistemas de ecuaciones lineales en análisis estructural: las técnicas numéricas (variable independiente) se utilizan para resolver el sistema de ecuaciones, la solución está directamente relacionada con la memoria utilizada y el tiempo de procesamiento que toman las técnicas para resolver el sistema, son estas variables las que nos interesa medir.

2.6. Operacionalización de variables

La Operacionalización de variables se describe en la siguiente tabla:

Tabla 2. Operacionalización de variables

Variables	Dimensiones	Metodología
V. Independiente • Técnicas numéricas sobre procesador gráfico.	Método (formato de almacenamiento) Gauss Jordan (Matriz densa) Factorización de Cholesky (matriz densa) Factorización de Cholesky (matriz CSC) Factorización de Cholesky (matriz SKS) Factorización LDLT (Matriz densa). Factorización LDLT (Matriz CSC). Factorización LDLT (matriz SKS) Gradientes conjugados (Matriz densa). Gradientes conjugados (matriz CSC). Gradientes conjugados (matriz CSC).	 Revisión bibliográfica Escritura de funciones en Julia y ArrayFire. Escritura de kernels en OpenCL. Ejecución de funciones y/o kernels
V. Dependiente • Solución de sistemas de ecuaciones lineales en análisis estructural	 Tiempo de ejecución (s) Memoria utilizada (MB) 	Obtención del sistema de ecuaciones: OpenFEM, herramienta de elementos finitos Memoria: medición manual de acuerdo al número de elementos de la matriz y la precisión de los elementos. Tiempo: Macro @time de Julia Macro @elapsed de Julia

CAPÍTULO III

METODOLOGÍA DE INVESTIGACIÓN

3.1. Tipo de investigación

El tipo de investigación es **aplicada**, se busca aplicar conocimientos o procedimientos ya establecidos en la solución de problemas prácticos (Niño Rojas, 2011).

3.2. Nivel de investigación

El nivel de investigación es: **descriptivo**, se busca especificar propiedades, características y rasgos importantes del fenómeno que se analiza.

(Hernández Sampieri, Fernández Collado, & Baptista Lucio, 2010), cada una de las técnicas consideradas para resolver el sistema tendrá diferentes comportamientos computacionalmente hablando, las características principales son tiempo de ejecución y memoria utilizada.

El método que se usará en la realización del proyecto de investigación: **método deductivo**, de lo general a lo particular, lo general sería los conocimientos que se tienen sobre el método de solución de sistemas de ecuaciones lineales y lo particular sería la aplicación del método a problemas de análisis estructural.

3.3. Población, muestra y muestreo

3.1.1 Población:

Modelos estructurales

3.1.2 Muestra

- · Viga en voladizo
- Pilar de puente

3.4. Técnicas e instrumentos de obtención de datos

Cada sistema de ecuaciones lineales, fue generado al resolver problemas típicos de análisis estructural, se ha usado el **OpenFEM**, una herramienta de elementos finitos de libre distribución, al analizar una estructura por elementos finitos tenemos la ventaja de generar una matriz de rigidez del orden que queramos, basta con hacer los elementos cada vez más pequeños.

3.5. Procedimiento de obtención de datos

Este procedimiento comprenderá:

- Revisión de documentos: revisar la información existente y de artículos que involucren temas relacionados con la investigación, es decir en nuestro caso documentación, artículos, manuales o guías de programación de los lenguajes a usar.
- Analizar cada estructura usando el método adecuado.
- Cuando se use el método de elementos finitos, primero se decidirá en cuantos sub-elementos se dividirá cierto elemento estructural considerando una cantidad de nodos aceptable que nos genere una cantidad de incógnitas relativamente grande.

3.6. Técnicas de procesamiento y análisis de datos

Las técnicas para el procesamiento y análisis de los datos serán:

Cuantitativas: Los resultados (tiempo y memoria utilizada) son cuantificables.

La programación del procesador gráfico se ha realizado usando la interfaz de **OpenCL** en C++, para ejecutar los kernels escritos se ha usado la librería de **ArrayFire**, escribiendo funciones que serán incluidas en otra biblioteca compartida con la finalidad de llamarlas desde Julia.

Por comodidad estas funciones serán puestas a prueba en **Julia**, usando el paquete **ArrayFire.jl** que hará posible usar las funciones con la memoria directamente en el dispositivo.

Todo el software mencionado anteriormente es de libre distribución y se ha usado la **versión comunity** de **Visual Studio**.

La medición del tiempo de procesamiento se obtendrán usando las macros @time y/o @elapsed de Julia, la medición de la memoria será directa, usando el orden de la matriz en el caso de una matriz densa, y las longitudes de los vectores que representan la matriz en el caso de una matriz dispersa, estos serán multiplicados por la memoria en (MB) MegaBytes que utiliza cada elemento en las componentes de la matriz.

❖ Cualitativas: se empleará para la interpretación de resultados.

CAPÍTULO IV

PRESENTACIÓN DE RESULTADOS

4.1. Análisis de información

4.1.1. Etapas

La investigación se compuso principalmente de siete etapas, partiendo de la escritura de los kernels y finalizando en la prueba de estos para resolver los sistemas de ecuaciones lineales midiendo el tiempo y la memoria utilizada, estas etapas se resumen a continuación:

Tabla 3. Primera etapa de la investigación

ETAPA I		
Escritura de los kernels		
Instrumentos	 API de OpenCL en C++ Visual Studio comunity 	
Metodología	 Uso de manuales y guías de programación en C++ y OpenCL 	
Resultado	Kernels para resolver el sistema de ecuaciones lineales, con las siguientes combinaciones: método-formato de almacenamiento: Gauss Jordan-denso Cholesky-denso Cholesky-CSC Cholesky-SKS LDL ^T -denso LDL ^T -CSC LDL ^T -SKS Gradientes conjugados-denso Gradientes conjugados -CSC Gradientes conjugados -SKS	

Tabla 4. Segunda etapa de la investigación

ETAPA II Escritura de funciones que utilizan los kernels y que se exportan en una biblioteca dinámica (dll) Instrumentos Biblioteca de ArrayFire Visual Studio comunity Uso de manuales y guías de programación en ArrayFire y C++. Archivo de extensión dll que contiene las funciones que utilizan los kernels para ser llamados desde otro lenguaje de programación.

Tabla 5. Tercera etapa de la investigación

ETAPA III		
Escritura de funciones en Julia que llaman a las funciones exportadas		
Instrumentos	JuliaPaquete ArrayFire.jl de Julia	
Metodología	• Uso de manuales y guías de programación en Julia.	
Resultado	 Funciones <i>Julia</i> listas para recibir como argumentos la matriz de rigidez (en uno de los formatos de almacenamiento) y el vector de fuerzas, y devolver el resultado de resolver éste sistema de ecuaciones lineales. Adicionalmente se escribieron funciones <i>Julia</i> que devuelven el tiempo de ejecución de resolver el sistema de ecuaciones lineales. 	

Tabla 6. Cuarta etapa de la investigación

ETAPA IV Modelamiento de la estructura y obtención del sistema de ecuaciones lineales **OpenFEM Instrumentos** MatLab Uso de manuales y guías de programación en MatLab y OpenFEM. Metodología Afinamiento de la malla de elementos finitos para obtener diferentes órdenes de la matriz. Cinco archivos de texto por cada orden de la matriz de rigidez de extensiones .nz, .row, .col, .fzs y .def; los tres primeros representan la matriz en formato COO, el cuarto Resultado es el vector de fuerzas y el último es el vector de deformaciones o vector resultado arrojado por el solver del OpenFEM que se exporta exclusivamente con fines de validación.

Tabla 7. Quinta etapa de la investigación

	ETAPA V
	del sistema de ecuaciones lineales en Julia y construcción de iz con el formato de almacenamiento correspondiente
Instrumentos	JuliaPaquete ArrayFire.jl de Julia
Metodología	Uso de manuales y guías de programación de Julia
Resultado	• Matrices en formato denso, CSC o SKS, y vector de fuerzas, listas para ser resueltas usando las funciones escritas en la etapa 3.

Tabla 8. Sexta etapa de la investigación

14	ETAPA VI	
Resultados		
Instrumentos	JuliaPaquete IJulia de Julia	
Metodología	 Uso de manuales y guías de programación de Julia Uso de las funciones para medir el tiempo de ejecución de la etapa 3. Medición manual de la memoria utilizada en función al tamaño de los elementos. 	
Resultado	 Tiempo de ejecución en segundos que toma resolver el sistema de ecuaciones lineales. Memoria utilizada en MB. 	

Tabla 9. Séptima etapa de la investigación

ETAPA VII		
Conclusiones		
Instrumentos	• Excel	
Metodología	 Elaboración de gráficos y tablas con los resultados de la etapa 6. Comparación con los resultados arrojados en tiempo de ejecución y memoria por las mismas funciones resolviendo el mismo sistema de ecuaciones pero ejecutadas usando otro dispositivo (GPU). Comparación con los resultados arrojados en tiempo de ejecución y memoria por el solver del OpenFEM. 	
Resultado	 Entre las técnicas utilizadas encontrar la que tiene un mejor balance en tiempo de ejecución y memoria utilizada. Resultado numérico de las diferencias que se obtienen al usar las mismas funciones, resolviendo el mismo sistema de ecuaciones, en diferentes dispositivos (GPU's). Determinar la efectividad de estas técnicas en comparación al resultado arrojado por el solver de OpenFEM. 	

4.1.2. Compatibilidad en la indexación

La indexación en C++ y por lo tanto en OpenCL, está basada en cero, es decir el primer elemento en un vector tendrá índice cero. En Julia, la indexación está basada en uno, esta diferencia fue manejada mediante la disminución previa de los índices de los vectores en Julia antes de usarse en funciones C++, esto es solo necesario cuando se han almacenado índices o indicadores de posición en algún tipo de vector que será usado como argumento, por ejemplo, los vectores de índices creados al almacenar matrices dispersas en formato CSC o SKS, cada elemento de estos vectores se han reducido en uno antes de establecerlos como argumentos en alguna de las funciones compartidas de C++ que usan los kernels OpenCL.

4.1.3. Orden de almacenamiento preferencial

Los elementos de una matriz no son almacenadas de forma bidimensional como podría concebirse a priori, si no como un vector, en ArrayFire los elementos de una matriz (objetos del tipo *array* o *af_array*) se almacenan en un vector columna por columna, es decir, una matriz de *mxn* es almacenada en memoria como un vector de longitud *mxn*.

Por ejemplo sea la matriz de 6 x 6 siguiente:

Se almacena en memoria como:

En memoria solo existe el vector y no la matriz en sí, acceder a los elementos de este vector es simple y solo requiere un id unidimensional, estos id en este caso son valores entre 0 y 35 (recordar que la indexación es basada en cero) que cubrirá los 36 elementos de esta matriz, ahora dado un id bidimensional (row, col), podemos acceder al elemento correspondiente en el vector partiendo de estos subíndices usando la siguiente expresión:

$$x = row + h_A col$$

 h_A es la altura de la matriz, row y col son los índices de fila y columna respectivamente de un elemento.

Un método que usamos recurrentemente para acceder a un elemento en OpenCL, es el de declarar un apuntador hacia la memoria del dispositivo moviéndonos cierta cantidad de elementos, en OpenCL este movimiento es realizado mediante la siguiente declaración:

```
__global double* Colb = A + hA * col;
```

A es el apuntador normal a la memoria del dispositivo con los elementos de la matriz, hA es la altura de la matriz y col es el índice de columna, como resultado Colb apuntará a la memoria del dispositivo con un movimiento de h_Acol elementos hacia adelante, es decir, el elemento con id 0 de Colb equivaldrá al

elemento h_Acol de A. así para acceder al elemento (row, col) de A bastaría con acceder al elemento row de Colb es decir:

$$A_{row,col} = Colb_{row}$$

4.1.4. Alistando los kernels

Los kernels fueron escritos en archivos de extensión cl, para cargarlos se ha creado una variable de entorno CL_FILES con valor igual a la ruta a la carpeta que contiene los archivos cl, creamos una función que obtiene el valor de una variable de entorno dado su nombre y creamos otra función con dos argumentos del tipo *char** que contienen el nombre de un archivo de encabezado y el nombre del archivo cl, esta última función utiliza la función anterior para obtener el valor de la variable CL_FILES, es decir, la ruta a la carpeta con los archivos cl, luego busca los archivos con los nombres especificados por lo argumentos, une su contenido y lo guarda en un objeto del tipo *char** listo para ser cargado. El archivo de encabezado mencionado es una archivo .h ubicada en la misma carpeta, contiene definiciones y otras líneas comunes a todos los archivos cl, en nuestro caso el archivo es llamado *Common.h* cuyo contenido se muestra en la Figura 13. Los archivos cl contienen exclusivamente los kernels y las funciones de dispositivo.

```
//tamaño de bloque para ser usado en funciones normales y Kernels
#define BLOCK_SIZE 32

//activando doble precisión (double) para kernels
#pragma OPENCL EXTENSION cl_khr_fp64 : enable

//activando funciones atómicas para 32 y 64 bits
#pragma OPENCL EXTENSION cl_khr_global_int32_base_atomics : enable
#pragma OPENCL EXTENSION cl_khr_int64_base_atomics : enable
```

Figura 13. Contenido de Common.h

4.1.5. Kernels sobre matrices densas

Se presentan los kernels elaborados para trabajar sobre matrices densas. Ya que la intención es trabajar con matrices de orden relativamente grande, se trabajó con precisión doble (tipo *double*), ya que en matrices con precisión simple (tipo *float*), se encontró pérdida de precisión en la solución a medida que aumentaba el orden de la matriz, cabe mencionar que existen métodos para trabajar en simple

precisión conservando una precisión aceptable, sin embargo, la aplicación de dichos métodos quedan fuera del alcance de esta tesis.

4.1.5.1. Gauss Jordan

El kernel se muestra en la Figura 14, como puede observarse el kernel llama a la función de dispositivo *Gss_Jrd_c* que se muestra en la Figura 15, el argumento *A* de esta función, es el apuntador hacia la memoria del dispositivo que contiene los elementos de la matriz, *wA* es el ancho de la matriz y *hA* es el alto de la matriz, los demás argumentos son usados por la función de la siguiente manera:

Figura 14. Gauss Jordan: Kernel para trabajar sobre una matriz densa

La función resta a cada fila de *A* la fila *row* multiplicada por un factor (calculada en la línea 26), tal que el elemento ubicado en la posición *col* de cada fila sea cero luego de realizada esta operación, *fcol* indica la posición en cada fila desde donde se realizarán las operaciones, si es para la fila completa el valor correspondiente de *fcol* es cero, del mismo modo debe especificarse una fila (*frow*) para obviar las filas anteriores en las operaciones, para todas las filas este valor debe ser cero.

Como puede observarse esta función es generalizada para matrices no necesariamente cuadradas, en nuestro caso el argumento A deberá ser la matriz aumentada, es decir, con el vector de constantes añadido a la derecha de la matriz de coeficientes original, al final debe tenerse una matriz de $n \times n+1$, siendo n el orden de la matriz. Cada llamada al kernel realizará una de las n iteraciones que deben hacerse para terminar con la eliminación de Gauss Jordan, esta iteración será realizada en el host. El término host en este caso se refiere a la entorno de ArrayFire.

```
void Gss Jrd c( global double* A, int wA, int hA,
1
2
             int row, int col, int frow, int fcol)
3
             // fcol + indice local de grupo
4
5
             int bx = fcol + get_group_id(0);
6
             //indice del primer elemento en una columna
             int tx = frow + get_local_id(0);
8
9
10
             //factor de multiplicación
11
             double fm;
12
13
             //columna base, la que se vuelve cero, excepto el
14
             //valor en la posición frow
              global double* Colb = A + hA * col;
15
16
17
18
             //en este caso, cada grupo se encargará de las operaciones
19
             //sobre todos los elementos de una columna
20
             double keyEl = Colb[row];
21
22
             for (int x = tx; x < hA; x += get_local_size(0)) {</pre>
23
24
                     if (x != row && keyEl != 0) {
25
26
                              fm = Colb[x] / keyEl;
27
                              for (int y = bx; y < wA; y += get_num_groups(0)) {</pre>
28
29
                                         global double* ColEl = A + y
                                      ColEl[x] -= fm * ColEl[row];
30
31
32
33
34
```

Figura 15. Función de dispositivo llamada por le kernel de la Figura 14

Notar que el kernel solo tiene tres argumentos, los argumentos de *Gss_Jrd_c* son establecidos en función a estos. El número de iteración que se realice, será igual al índice de la columna que se hace cero en esa iteración, es decir si estamos en la quinta iteración, al finalizar esta, la quinta columna tendrá que ser cero (exceptuando obviamente el quinto elemento de esta columna).

El código de *Gss_Ird_c* está escrito para trabajar de la siguiente manera:

Cada grupo de trabajo y sus hilos se encargará de una columna, ya que la función está diseñada para trabajar partiendo de ciertos índices dados por *frow* y *fcol*, en las líneas 5 y 8 se suma a los índices de grupo y de hilo *fcol* y *frow* respectivamente, esto para hacer posible que cada grupo empiece a trabajar a partir del índice *fcol* y cada hilo a partir del índice *frow*.

La línea 11 declara una variable del tipo double para el factor de multiplicación que se calcula en la línea 26. La línea 15 declara un apuntador *Colb* hacía la memoria del dispositivo con inicio a una distancia *hA*col* elementos del inicio

de *A*, esto nos lleva a la columna base, la que se trabaja para convertir *n-1* de sus elementos en cero. En la línea 20 se obtiene el elemento *row* de *Colb* correspondiente a la fila base, la que multiplicaremos por el factor *fm* para restarla de las otras filas, en el kernel de Figura 14, equivaldría a obtener el elemento diagonal de A.

La iteración de la línea 22 a 33 se realiza con la finalidad de que el número de hilos en un grupo de trabajo cubra todos los elementos de la columna a trabajar, esto es necesario ya que el número de hilos por grupo, establecido al llamar el kernel, no depende del tamaño de la matriz. Si el método para llamar al kernel estableciera un número de hilos igual al número de filas que debe trabajarse, esta iteración no sería necesaria. De la misma manera ocurre con la iteración interior (líneas 27 a 31), que se hace con la finalidad de que los grupos disponibles cubran todas las columnas.

Un detalle más a tener en cuenta es el valor asignado a *fcol* en el kernel, este es igual a *row+1*, lo que significa que no se está trabajando la columna *row* (la que debe volverse cero), se puede omitir esta operación ya que no necesitamos el resultado final en esta columna por que no se utilizan en ninguna otra operación, estas serán implícitamente cero.

Al terminar con la iteración en el *host*, la solución del sistema estará dada por la división de los elementos de la última columna de la matriz aumentada (que corresponde a los elementos del vector de constantes del sistema) y los elementos de la diagonal principal.

4.1.5.2. Factorización de Cholesky

Para resolver un sistema de ecuaciones usando la factorización de Cholesky, son necesarias dos etapas, uno de la factorización misma y otro que consiste en resolver los sistemas triangulares generados después de la factorización.

El proceso de factorización fue dividido en tres partes, un paso en el que modificamos una columna *j* tal que:

$$A_{n,j} = \frac{A_{n,j}}{\sqrt{A_{j,j}}}, n > j$$

La función de dispositivo que realiza esta operación se muestra en la Figura 16, el trabajo de modificar los elementos de la columna se le asigna a todos los hilos de todos los grupos, es por eso que en vez de pedir el id local de hilo se pide el id global (línea 5), la columna a trabajar es la misma determinada por el argumento *step*.

```
void Chol_1_c(__global double* A, int order, int step)
2
3
4
             //step + 1 + indice global de hilo
            int tx = get_global_id(0);
             //número total de hilos
8
            int ntx = get_num_groups(0)*get_local_size(0);
9
10
             //columna base
             global double* colB = A + step * order;
11
12
13
             //raíz cuadrada del elemento A(step, step)
14
            double _r = sqrt(colB[step]);
15
16
             //cada hilo se encargará de modificar un elemento
17
             //de la columna
18
            for (int x = tx + step + 1; x < order; x += ntx) {
19
                     colB[x] = colB[x] / _r;
20
```

Figura 16. Cholesky: función de dispositivo *Chol_1_c*

Otro paso consiste en modificar las columnas siguientes a la columna *j* modificada en el paso anterior tal que:

$$A_{n,s} = A_{n,s} - A_{s,j}A_{n,j}, s > j, n \ge s$$

La función de dispositivo que realiza esta tarea se muestra en la Figura 17, en este caso, cada grupo y sus hilos se encargarán de modificar una columna.

Estos dos primeros pasos deben realizarse uno después de otro durante *n* iteraciones, un tercer paso es necesario para terminar la factorización, este modifica los elementos de la diagonal tal que:

$$A_{j,j} = \overline{A_{j,j}}$$

```
void Chol_2_c(__global double* A, int order, int step)
2
             //indice local de grupo
4
             int bx = get_group_id(0);
             //step + índice local de hilo
             int tx = get_local_id(0);
8
9
             //columna base
10
              _global double* colB = A + step * order;
11
12
             //parámetros de apoyo
             double _k;
13
14
             //en este caso, cada grupo se encargará de las
15
             //operaciones sobre todos los elementos de una
16
17
             //columna
             for (int y = bx + step + 1; y < order; y += get_num_groups(0)) {</pre>
19
                     _k = colB[y];
                     __global double* colI = A + y * order;
20
21
22
                     for (int x = tx + y; x < order;
23
                              x += get_local_size(0)) {
24
25
                              colI[x] -= _k * colB[x];
26
27
28
```

Figura 17. Cholesky: función de dispositivo Chol_2_c

La función de dispositivo que realiza esta tarea se muestra en la Figura 18, todos los hilos de todos los grupos se encargan cada uno de modificar un elemento de la diagonal, este paso solo se realiza una vez luego de realizada las *n* iteraciones con los dos pasos vistos previamente.

Terminada la factorización se procede a solucionar los sistemas equivalentes, es decir si Ax = b es el sistema original, luego de la factorización quedará:

$$LL^Tx = b$$

Haciendo $L^T x = y$ entonces L y = b. Primero se resuelve L y = b y finalmente $L^T x = y$, cuya solución en x es la solución del sistema A x = b.

Como puede observarse, la factorización se realiza en sitio, es decir los elementos de L son almacenados en la misma matriz A, luego de la factorización, A contendrá en su parte triangular inferior, los elementos de L, asumiendo implícitamente que los elementos de la parte triangular superior (sin contar la diagonal) son ceros.

```
void Chol_3_c(_global double* A, int order)
4
             //indice global de hilo
             int tx = get_global_id(0);
             //número total de hilos
             int ntx = get_num_groups(0)*get_local_size(0);
8
10
             //parámetros de apoyo
11
             int diag;
13
             for (int x = tx; x < order; x += ntx) {
                     //ubicación del elemento diagonal en memoria
14
                     diag = x * order + x;
15
16
                     A[diag] = sqrt(A[diag]);
17
```

Figura 18. Cholesky: función de dispositivo *Chol_3_c*

Para resolver el sistema Ly = b, usaremos la eliminación típica de Gauss Jordan aplicada a una matriz triangular inferior, en cada iteración los elementos de cada columna fuera de la diagonal van volviéndose cero progresivamente, hasta llegar a una matriz diagonal, en ningún momento debe hacerse uso de la parte triangular superior (son implícitamente ceros, aunque en memoria guardan la parte triangular superior de la matriz original A), las operaciones sobre L pueden obviarse y dejarla intacta para poder usarla al resolver el siguiente sistema $L^Tx = y$, es decir solo se modificará el vector de constantes b tal que para la iteración j:

$$b_n = b_n - b_j \frac{L_{n,j}}{L_{j,j}}, n > j$$

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 19, el número de iteración y por la tanto el índice de la columna implicada está dada por el valor del argumento step, esta función debe ser llamada n veces para terminar con la eliminación, luego será necesario dividir los elementos de b entre los elementos de la diagonal (función de dispositivo de la Figura 20), este vector será la solución del sistema Ly = b, que será usado como argumento al resolver el siguiente sistema triangular.

```
void Chol_4_c(__global double* L, __global double* b, int order,
2
             int step)
4
              //indice global de hilo
              int tx = get_global_id(0);
8
              //número total de hilos
9
              int ntx = get_num_groups(0)*get_local_size(0);
10
11
              //columna de L implicada
              __global double* colB = L + order * step;
13
              //elemento diagonal
14
             double diag = colB[step];
15
16
             double k = b[step] / diag;
17
18
             for (int x = tx + step + 1; x < order; x += ntx) {
    b[x] -= k * colB[x];</pre>
19
20
21
22
```

Figura 19. Cholesky: función de dispositivo Chol_4_c

Para resolver el sistema $L^T x = y$, se usará la misma matriz que almacena L, el artificio es asumir que está ordenada por filas y no por columnas, hecho esto la transposición de la matriz ya no será necesaria. L^T es una matriz triangular superior, se resolverá el sistema obteniendo el valor de las variables una a continuación de otra desde abajo, el valor de la variable en la posición i sería:

```
x_{i} = \frac{y_{i} - \sum_{j=i+1}^{n-1} x_{j} L_{i,j}}{L_{i,i}}
    void Chol_5_c(__global double* A, __global double* b, int order)
3
              //indice global de hilo
              int tx = get_global_id(0);
              //número total de hilos
              int ntx = get_num_groups(0)*get_local_size(0);
10
              //parámetros de apoyo
11
              int diag;
12
13
              for (int x = tx; x < order; x += ntx) {
14
                       //ubicación del elemento diagonal en memoria
                       diag = x * order + x;
15
                       b[x] = b[x] / A[diag];
16
17
18
           Figura 20. Cholesky: función de dispositivo Chol_5_c
```

Se usa el mismo vector y para almacenar la solución desde abajo. La función de dispositivo que resuelve el sistema $L^T x = y$ se muestra en la Figura 21, cada llamada a esta función calculará una variable. Ya que el sistema se resuelve

desde abajo, para un paso j se calculará el valor de la variable i tal que i=n-j-1, la fila de L usada para calcular esta variable tiene el mismo índice (se obtiene en la línea 14). Esta función de dispositivo usará memoria local para calcular la sumatoria $\sum_{j=i+1}^{n-1} x_j L_{i,j}$ por reducción en paralelo, el tamaño de esta memoria local será el mismo que el tamaño de un grupo de trabajo, por esta razón, en este caso solo se utilizará un grupo de trabajo al usar esta función (en la línea 26 se especifica que el grupo de trabajo utilizado es únicamente el con id cero).

```
1
3
4
            //indice de bloque
5
            int bx = get_group_id(0);
6
            //indice global de hilo
7
8
            int tx = get_local_id(0);
9
            //número total de hilos
10
11
            int ntx = get_local_size(0);
12
13
            //índice de la fila de L implicada
14
            int fL = order - step - 1;
15
16
            //fila de L implicada, data
17
             _global double* filB = L + fL * order;
18
19
            //elemento diagonal
20
            double diag = filB[fL];
21
            //valor donde se guardará una suma parcial
22
23
            double sum = 0.0;
24
25
            //solo se usará un bloque
26
            if (bx == 0) {
                    for (int x = tx + fL + 1; x < order; x += ntx) {
27
                            sum += filB[x] * y[x];
28
29
30
                    partialSum[tx] = sum;
31
                    //reducción en paralelo pare determinar la suma de
32
33
                    //todos los elementos de partialSum
                    for (int stride = ntx / 2; stride > 0; stride /= 2) {
34
35
                            barrier(CLK_LOCAL_MEM_FENCE);
36
                            if (tx < stride) {</pre>
37
                                    partialSum[tx] += partialSum[tx + stride];
38
39
                    if (tx == 0)
40
41
                            y[fL] = (y[fL] - partialSum[0]) / diag;
42
43
```

Figura 21. Cholesky: función de dispositivo *Chol_6_c*

En la línea 27 a 29 se calcula una suma parcial almacenándola en la variable *sum* (declarada en la línea 23), luego este valor es almacenado en memoria local (línea 30), los elementos de *partialsum* se suman en paralelo en el bloque

for de las líneas 34 a 39, finalmente se calcula la variable correspondiente en la línea 41, esta tarea solo la realiza un hilo (en la línea 40 se especifica el id del hilo que realiza la operación). Esta función de dispositivo debe ser llamada también n veces para calcular las n variables.

```
kernel void
    Cholesky_c(__global double* A,
                                     __global double* b,
2
3
              local double* partialSum, int order, int step, int sstep)
4
5
             //la iteración principal se realiza en el host
6
             //se modificará una sola columna
8
             if (sstep == 0) {
9
                     Chol_1_c(A, order, step);
10
11
             //se modificará las columnas siguientes a la columna
12
             //implicada
13
             else if (sstep == 1) {
14
                     Chol_2_c(A, order, step);
15
16
17
             //una vez usada las funciones anteriores todavía es
             //necesario modificar la diagonal D[i]=sqrt(D[i]) para
18
19
             //concluir la factorización
20
             else if (sstep == 2) {
                     Chol_3_c(A, order);
21
22
23
24
             ///hasta aquí la factorización está concluida, la parte
             //triangular inferior de A contendrá la factorización de
25
26
             //Cholesky es decir L
27
             //desde aquí se procederá a solucionar el sistema
28
             //L*transpose(L)x=b
29
30
             //se resolverá el sistema Ly=b, y=transpose(L)*x
31
             else if (sstep == 3) {
32
                     Chol_4_c(A, b, order, step);
33
34
             else if (sstep == 4) {
                     Chol_5_c(A, b, order);
35
36
37
38
             //se resolverá el sistema transpose(L)*x=y
             //"y" ya debe haberse obtenido usando la función anterior
39
40
             else {
41
                     Chol_6_c(A, b, order, step, partialSum);
42
43
```

Figura 22. Cholesky: Kernel para trabajar sobre una matriz densa

El kernel principal se muestra en la Figura 22, como puede observarse el kernel llama a las funciones de dispositivo mencionadas anteriormente, dependiendo del valor del argumento *sstep*, los valores 0, 1 y 2 indican que nos encontramos en el proceso de factorización y los valores 3, 4 y 5 que estamos en el proceso de solución de los sistemas triangulares.

El kernel es llamado n veces para sstep igual a 0 y 1 (con step de 0 a n-1), solo una vez para sstep=2, n veces para sstep=3 (con step de 0 a n-1), una vez para sstep=4 y finalmente n veces para sstep=5(con step de 0 a n-1).

4.1.5.3. Factorización LDL^T

El proceso para resolver un sistema usando la factorización LDL^T es similar al que utilizamos en la factorización de Cholesky.

Para resolver un sistema de ecuaciones usando la factorización LDL^T, son necesarias dos etapas, uno de la factorización misma y otro que consiste en resolver los sistemas triangulares generados después de la factorización.

El proceso de factorización fue dividido en dos partes, un paso en el que modificamos una columna *j* tal que:

$$A_{n,j} = \frac{A_{n,j}}{A_{j,j}}, n > j$$

```
void ldlt 1 c( global double* A, int order, int step)
2
3
             //indice global de hilo
5
             int tx = get_global_id(0);
6
             //número total de hilos
8
             int ntx = get_num_groups(0)*get_local_size(0);
9
             //columna base
             global double* colB = A + step * order;
11
12
13
             //elemento A(step, step)
14
             double _r = colB[step];
15
             //cada hilo se encargará de modificar un elemento
16
             //de la columna
17
18
             for (int x = tx + step + 1; x < order; x += ntx) {
                     colB[x] = colB[x] / _r;
19
20
```

Figura 23. LDL^T: función de dispositivo *ldlt 1 c*

La función de dispositivo que realiza esta operación se muestra en la Figura 23, el trabajo de modificar los elementos de la columna se le asigna a todos los hilos de todos los grupos, es por eso que en vez de pedir el id local de hilo se pide el id global (línea 5), la columna a trabajar es la misma determinada por el argumento *step*.

Otro paso consiste en modificar las columnas siguientes a la columna j modificada en el paso anterior tal que:

$$A_{n,s} = A_{n,s} - A_{s,j}A_{n,j}A_{j,j}, s > j, n \ge s$$

La función de dispositivo que realiza esta tarea se muestra en la Figura 24, en este caso, cada grupo y sus hilos se encargarán de modificar una columna.

```
void ldlt_2_c(__global double* A, int order, int step)
    {
             //indice de grupo
             int bx = get_group_id(0);
             //indice local de hilo
             int tx = get_local_id(0);
             //columna base
              _global double* colB = A + step * order;
10
11
12
             //elemento A(step, step)
13
             double _r = colB[step];
14
15
             //parámetros de apoyo
             double _k;
16
17
             //en este caso, cada grupo se encargará de las
18
             //operaciones sobre todos los elementos de una
19
20
21
             for (int y = bx + step + 1; y < order; y += get_num_groups(0)) {</pre>
                     _{k = _r * colB[y]};
22
                      _global double* colI = A + y * order;
23
24
                     for (int x = tx + y; x < order;
25
26
                              x += get local size(0)) {
27
                              coll[x] -= _k * colb[x];
28
29
30
```

Figura 24. LDL^T: función de dispositivo *ldlt_2_c*

Estos dos primeros pasos deben realizarse uno después de otro durante n iteraciones para terminar con la factorización. Terminada la factorización en sitio, A contendrá en su parte triangular inferior la factorización LDL^{T} , L en su parte estrictamente triangular inferior y D en su diagonal, debe recordarse que L por si sola tiene sus elementos diagonales iguales a 1, esto es implícito y debe tenerse en cuenta cuando llegue el momento de resolver el sistema.

Terminada la factorización se procede a solucionar los sistemas equivalentes, es decir si Ax = b es el sistema original, luego de la factorización quedará:

$$LDL^{T}x = b$$

Haciendo $DL^Tx = y$ entonces Ly = b, es el primer sistema a resolverse para y, luego haciendo $L^Tx = z$ entonces Dz = y es el segundo sistema a resolver en z, que es simple ya que D es una matriz diagonal, finalmente se resuelve $L^Tx = z$, cuya solución en x es la solución del sistema Ax = b.

Para resolver el sistema Ly = b, usaremos la eliminación típica de Gauss Jordan aplicada a una matriz triangular inferior, en cada iteración los elementos de cada columna fuera de la diagonal van volviéndose cero progresivamente, hasta llegar a una matriz diagonal, en ningún momento debe hacerse uso de la parte triangular superior (son implícitamente ceros, aunque en memoria guardan la parte triangular superior de la matriz original A), resolver un sistema triangular nos permite obviar las operaciones sobre L y dejarla intacta para poder usarla al resolver el sistema $L^Tx = z$, solo se modificará el vector de constantes b tal que para la iteración j:

$$b_n = b_n - b_j L_{n,j}, n > j$$

```
void ldlt_3_c(__global double* L, __global double* b, int order,
2
              int step)
3
4
5
              //indice global de hilo
              int tx = get_global_id(0);
6
8
              //número total de hilos
              int ntx = get_num_groups(0)*get_local_size(0);
9
10
              //columna de L implicada
11
12
              global double* colB = L + order * step;
13
              for (int x = tx + step + 1; x < order; x += ntx) {
    b[x] -= b[step] * colB[x];</pre>
14
15
16
17
```

Figura 25. LDL^T: función de dispositivo *ldlt_3_c*

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 25, el número de iteración y por la tanto el índice de la columna implicada está dada por el valor del argumento step, esta función debe ser llamada n veces para terminar con la eliminación, al final b contendrá la solución del sistema Ly = b, que será usado como argumento al resolver el siguiente sistema.

Luego se procede a resolver el sistema Dz = y, ya que D es diagonal, la solución se reduce a realizar la operación:

$$b_j = \frac{b_j}{D_{j,j}}$$

La función de dispositivo que realiza esta operación se muestra en la Figura 26, el kernel se ejecuta una sola vez.

```
void ldlt_4_c(__global double* D, __global double* b, int order)
             //indice global de hilo
             int tx = get_global_id(0);
             //número total de hilos
             int ntx = get_num_groups(0)*get_local_size(0);
9
10
             //parámetros de apoyo
11
             int diag;
12
             for (int x = tx; x < order; x += ntx) {
13
                      //ubicación del elemento diagonal en memoria
14
                      diag = x * order + x;
b[x] = b[x] / D[diag];
15
16
17
```

Figura 26. LDL^T: función de dispositivo *ldlt_4_c*

Para resolver el sistema $L^T x = z$, se usará la misma matriz que almacena L, el artificio es asumir que está ordenada por filas y no por columnas, con esto la transposición de la matriz ya no será necesaria. L^T es una matriz triangular superior, se resolverá el sistema obteniendo el valor de las variables una a continuación de otra desde abajo, recordar que los elementos de la diagonal de L son 1 aunque no estén almacenados, así, el valor de la variable en la posición i sería:

$$x_i = z_i - \sum_{j=i+1}^{n-1} z_j L_{i,j}$$

```
1
2
            //indice de bloque
4
5
            int bx = get_group_id(0);
6
            //indice global de hilo
            int tx = get_local_id(0);
8
9
10
            //número total de hilos
11
            int ntx = get_local_size(0);
12
13
            //indice de la fila de L implicada
14
            int fL = order - step - 1;
15
            //fila de L implicada, data
16
             _global double* filB = L + fL * order;
17
18
19
            //valor donde se guardará una suma parcial
20
            double sum = 0.0;
21
22
            //solo se usará un bloque
23
            if (bx == 0) {
24
                    for (int x = tx + fL + 1; x < order; x += ntx) {
25
                            sum += filB[x] * z[x];
26
27
                    partialSum[tx] = sum;
28
29
                    //reducción en paralelo pare determinar la suma de
30
                    //todos los elementos de partialSum
                    for (int stride = ntx / 2; stride > 0; stride /= 2) {
31
32
                            barrier(CLK_LOCAL_MEM_FENCE);
33
                            if (tx < stride) {</pre>
                                    partialSum[tx] += partialSum[tx + stride];
34
35
36
                    if (tx == 0)
37
                            z[fL] -= partialSum[0];
38
39
40
```

Figura 27. LDL^T: función de dispositivo *ldlt_5_c*

Se usa el mismo vector z para almacenar la solución desde abajo. La función de dispositivo que resuelve el sistema $L^Tx=z$ se muestra en la Figura 27, cada llamada a esta función calculará una variable. Ya que el sistema se resuelve desde abajo, para un paso j se calculará el valor de la variable i tal que i=n-j-1, la fila de L usada para calcular esta variable tiene el mismo índice (se obtiene en la línea 14). Esta función de dispositivo usará memoria local para calcular la sumatoria $\sum_{j=i+1}^{n-1} z_j L_{i,j}$ por reducción en paralelo, el tamaño de esta memoria local será el mismo que el tamaño de un grupo de trabajo, por esta razón, en este caso solo se utilizará un grupo de trabajo al usar esta función (en la línea 23 se especifica que el grupo de trabajo utilizado es únicamente el con id cero).

En la línea 24 a 26 se calcula una suma parcial almacenándola en la variable *sum* (declarada en la línea 20), luego este valor es almacenado en memoria

local (línea 27), los elementos de *partialsum* se suman en paralelo en el bloque *for* de las líneas 31 a 36, finalmente se calcula la variable correspondiente en la línea 38, esta tarea solo la realiza un hilo (en la línea 37 se especifica el id del hilo que realiza la operación). Esta función de dispositivo debe ser llamada también *n* veces para calcular las *n* variables.

```
kernel void
1
2
    ldlt_c(__global double* A, __global double* b,
              local double* partialSum, int order, int step, int sstep)
3
4
5
             //la iteración principal se realiza en el host
6
             //se modificará una sola columna
8
             if (sstep == 0) {
                     ldlt_1_c(A, order, step);
9
10
11
             //se modificará las columnas siguientes a la columna
12
             //implicada
13
             else if (sstep == 1) {
14
                     ldlt_2_c(A, order, step);
15
16
17
             ///hasta aquí la factorización está concluida, la parte
             //triangular inferior de A contendrá la factorización LDLT
18
19
             //es decir L y D
20
             //desde aquí se procederá a solucionar el sistema
            //L*D*transpose(L)x=b
21
22
23
             //se resolverá el sistema Ly=b, y=D*transpose(L)*x
24
            else if (sstep == 2) {
25
                     ldlt_3_c(A, b, order, step);
26
27
28
             //se resolverá el sistema Dz=y, z=transpose(L)*x
29
            else if (sstep == 3) {
30
                     ldlt_4_c(A, b, order);
31
32
33
             //se resolverá el sistema transpose(L)*x=z
34
            //"z" ya debe haberse obtenido usando la función anterior
35
            else {
                     ldlt_5_c(A, b, order, step, partialSum);
36
37
38
```

Figura 28. LDL^T: Kernel para trabajar sobre una matriz densa

El kernel principal se muestra en la Figura 28, como puede observarse el kernel llama a las funciones de dispositivo mencionadas anteriormente, dependiendo del valor del argumento *sstep*, los valores 0 y 1 indican que nos encontramos en el proceso de factorización y los valores 2, 3 y 4 que estamos en el proceso de solución de los sistemas triangulares.

El kernel es llamado n veces para sstep igual a 0 y 1 (con step de 0 a n-1), n veces para sstep=2 (con step de 0 a n-1), una vez para sstep=3 y finalmente n veces para sstep=4(con step de 0 a n-1).

4.1.5.4. Gradientes conjugados

La ventaja de este método es que solo requiere de operaciones sencillas en cada iteración, estas operaciones puede hacerse haciendo uso de las funciones incorporadas en ArrayFire, por lo que solo necesitamos escribir el código en esta interface, sin necesidad de escribir un kernel.

En la Figura 29 se muestra la función que se exporta para llamarla desde Julia, C es el objeto af_array* donde se guardará la solución del sistema, A es la matriz de coeficientes del sistema, b es el vector de constantes del sistema e Ierr es el error máximo para aceptar una solución. En las líneas 4 a 7 se extrae el orden de la matriz, las líneas 9 a 10 se extrae el tipo (float o double, cuyos equivalentes en ArrayFire son f32 y f64 respectivamente), en las líneas 12 a 26 se declaran objetos que se usarán durante las iteraciones, en las líneas 28 a 57 se establecen valores iniciales antes de empezar la iteración, y finalmente en las líneas 58 a 93 se procede con la iteración. Notar que la iteración termina una vez que la norma del vector r sea menor al valor de Ierr introducido como argumento (en las líneas 66 a 67 se hace este chequeo en cada iteración).

```
void AFire::SEL_gc(af_array* C, af_array A, af_array b,
2
              double Ierr) {
3
4
              dim_t _order[AF_MAX_DIMS];
5
              af_get_dims(&_order[0], &_order[1], &_order[2],
6
                      &_order[3], A);
7
              size t order = _order[0];
8
9
              af_dtype typef;
10
              af_get_type(&typef, A);
11
12
              double normr;
13
              //af array de ayuda
14
              af_array zero;
15
              dim_t d_order[] = { 1 };
16
              af_constant(&zero, 0, 1, d_order, typef);
17
              af_array Ax0;
18
              af_array rtxp;
19
              af array ptxz;
20
              af_array axp;
              af_array B;
21
22
              af_array axz;
              af_array copyr;
23
              af_array rtxz;
```

```
25
                af_array rsp;
26
                af_array Bxp;
27
28
                //x0=b
29
                af_array x0;
30
                af_copy_array(&x0, b);
31
32
                //r=b-A*x0
33
                af array r;
34
                af_matmul(&Ax0, A, x0, AF_MAT_NONE, AF_MAT_NONE);
35
                af_sub(&r, b, Ax0, false);
36
37
                //p = r
38
                af_array p;
                af_copy_array(&p, r);
39
40
41
                //z=A*p
42
                af_array z;
43
                af_matmul(&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
44
                //a = (r'*p)/(p'*z)
45
                af_array a;
46
               af_matmul(&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
af_matmul(&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
47
48
49
                af_div(&a, rtxp, ptxz, false);
50
51
                //x = x0 + a.*p
                af_array x;
52
53
                af_mul(&axp, a, p, true);
54
                af_add(&x, x0, axp, false);
55
56
                int contin;
57
                int I = 0;
                for (int i = 0; i < order; i++)</pre>
58
59
60
                         //r -= a.*z
61
                         af_copy_array(&copyr, r);
62
                         af_mul(&axz, a, z, true);
63
                         af_sub(&r, copyr, axz, false);
64
                         af_norm(&normr, r, AF_NORM_EUCLID, 1, 1);
if (normr <= Ierr)</pre>
65
66
67
                                  break;
68
69
                         //B = -(r'*z)/(p'*z)
                         af_matmul(&rtxz, r, z, AF_MAT_TRANS, AF_MAT_NONE);
70
                         af_matmul(&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
af_div(&rsp, rtxz, ptxz, false);
71
72
73
                         af_sub(&B, zero, rsp, false);
74
75
                         //p = r + B.*p
                         af_mul(&Bxp, B, p, true);
76
77
                         af_add(&p, r, Bxp, false);
78
79
                         //z = A*p
80
                         af_matmul(&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
81
                         //a = (r'*p)/(p'*z)
82
                         af_matmul(&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
af_matmul(&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
83
84
85
                         af_div(&a, rtxp, ptxz, false);
86
87
                         //x += a.*p
                         af_mul(&axp, a, p, true);
af_copy_array(&x0, x);
88
89
90
                         af_add(&x, x0, axp, false);
91
92
                         I++;
93
                }
94
```

```
//copiando el resultado en el argumento de salida
95
96
              af_copy_array(C, x);
97
98
              //liberando objetos af_array usados
99
              af_release_array(zero);
100
              af_release_array(Ax0);
              af_release_array(rtxp);
101
102
              af_release_array(ptxz);
203
              af release array(axp);
104
              af_release_array(B);
105
              af_release_array(axz);
106
              af_release_array(copyr);
107
              af_release_array(rtxz);
108
              af_release_array(rsp);
109
              af release array(Bxp);
110
              af_release_array(x0);
111
              af_release_array(r);
112
              af_release_array(p);
113
              af_release_array(z);
114
              af_release_array(a);
115
              af_release_array(x);
116
```

Figura 29. Gradientes conjugados: función principal a exportar

4.1.6. Kernels sobre matrices en formato CSC

Se presentan los kernels elaborados para trabajar sobre matrices en formato CSC. Se ha definido una estructura *SparseMS* para los objetos en este formato cuya definición en Julia se muestra en la Figura 30, los campos son: *elements* que guardan todos los elementos distintos de cero, *colIdx* donde cada elemento indicará la ubicación en *elements* del primer elemento distinto de cero en cada columna de la matriz original y *rowIdx* que guarda los índices de fila de todos los elementos. De la misma manera definimos la estructura *af_SparseMS* (Figura 31) análogo a *SparseMS* que guardará la matriz dispersa en formato CSC con los campos del tipo *AFArray*, que indica que son objetos en la memoria del dispositivo.

```
struct SparseMS{T,N}
elements::Array{T,1}
colIdx::Array{N,1}
rowIdx::Array{N,1}
end
```

Figura 30. Estructura para objetos que almacenan la matriz en formato CSC

```
struct af_SparseMS{T,N}
elements::AFArray{T,1}
colIdx::AFArray{N,1}
rowIdx::AFArray{N,1}
```

Figura 31. Estructura para objetos que almacenan la matriz en formato CSC con la memoria del dispositivo

```
function sparse ms (A::Array(T,2)) where {T<:Real}
        colIdx=Array(Int32,1)(UndefInitializer(),0);
        rowIdx=Array(Int32,1)(UndefInitializer(),0);
        elements=Array(T, 1) (UndefInitializer(), 0);
        cont=1;
        for j in 1:size(A)[1]
            push! (colIdx, cont);
 8
            for i in j:size(A)[1]
 9
                 if i--j
                     #el elemento diagonal siempre será guardado, pero su indice será
                     #implicito (no será almacenado en rowIdx)
                     push! (elements, A[i,j]);
13
                     cont+=1;
14
                 else
15
                     if A[i,j]!=0
16
                         push! (elements, A[i,j]);
                         push! (rowIdx, i);
                         cont+=1;
19
                     end
                 end
            end
21
        end
23
        SparseMS (elements, colIdx, rowIdx)
24
```

Figura 32. Conversión de una matriz densa a una matriz dispersa en formato CSC

La Figura 32 muestra el código en Julia para convertir una matriz simétrica en formato denso en una matriz en formato CSC, para poder usar las funciones exportadas todavía es necesario tener la memoria en el dispositivo, para esto, hacemos uso de la función incorporada de ArrayFire *AFArray*, que hace una copia en la memoria del dispositivo de un objeto del tipo *Array*, se ha sobrecargado la función para poder hacer una copia de un objeto del tipo *SparseMS* en la memoria del dispositivo, en otras palabras convertir del tipo *SparseMS* a *af_SparseMS*, esta función se muestra en la Figura 33.

```
function AFArray(A::SparseMS{T,N}) where {T<:Real,N<:Real}

Ael=AFArray(A.elements);

AcIdx=AFArray(A.colIdx);

ArIdx=AFArray(A.rowIdx);

af_SparseMS(Ael,AcIdx,ArIdx)
end</pre>
```

Figura 33. Función para convertir un objeto del tipo SparseMS a af_SparseMS

4.1.6.1. Factorización de Cholesky

Para resolver un sistema de ecuaciones usando la factorización de Cholesky, son necesarias dos etapas, uno de la factorización misma y otro que consiste en resolver los sistemas triangulares generados después de la factorización.

El proceso es similar al usado sobre matrices densas, sin embargo para empezar con la factorización es necesario realizar una factorización simbólica, esto con el fin de obtener la cantidad de memoria máxima requerida para almacenar los elementos de la factorización, en la Figura 34 se muestra el código en Julia propuesto, esta función obtiene un objeto del tipo *SparseMS* cuyos campos *colldx* y *rowldx* son los finales, sin embargo el campo *elements* solo contiene ceros, pero la cantidad de elementos es la necesaria para guardar la factorización.

Sea A la matriz original en formato CSC y L la matriz que contendrá la factorización de Cholesky también en formato CSC. Un paso previo es necesario antes de proceder con la factorización, este paso consiste en llenar los espacios del campo *elements* de L, con los correspondientes elementos del campo *elements* de A, la función de dispositivo que realiza esta tarea se muestra la Figura 35.

```
function symbolic chol fac(A::SparseMS(T,N)) where {T<:Real,N<:Real}
        Col=A.colIdx;
3
        Row-A.rowIdx;
4
        orden=length(Col);
5
 6
        #construyendo un vector de vectores, donde se separará por columna los
7
        #indices almacenados en Row
8
        Rows=Array(Array(N, 1), 1) (UndefInitializer(), 0);
9
10
        for i in 1:orden-1
            nEl=Col[i+1]-Col[i]-1;
11
12
            if nEl>0
13
                push! (Rows, Row[cont:cont+nEl-1]);
14
                cont += nEl;
15
16
                #significa que no hay elementos fuera de la diagonal
17
                push! (Rows, []);
18
            end
19
        end
21
      for i in 1:orden-1
            nEl=length(Rows[i]);
23
            #debe haber un minimo de 2 elementos
            if nEl>1
25
                for j in 1:nEl-1
26
                     union! (Rows[Rows[i][j]], Rows[i][j+1:end]);
                     sort!(Rows[Rows[i][j]]);
28
                end
29
            end
30
      end
        #generando los campos del objeto de salida
31
        colId=ones(N, orden);
33
        rowId=zeros(N,0);
34
        for i in 1:orden-1
35
            append! (rowId, Rows[1]);
36
            colId[i+1]=colId[i]+length(Rows[1])+1;
37
            popfirst! (Rows);
38
39
        Elm=zeros(T,colId[end]);
40
        SparseMS (Elm, colld, rowld)
41 end
```

Figura 34. Función que realiza la factorización de Cholesky simbólica

La función acepta seis argumentos, los campos de A y L en formato CSC (elmA, colA, rowA y elmL, colL, rowL respectivamente) y el tamaño del campo de elmA. En la línea 6 se obtiene el índice global del hilo, la línea 10 calcula el número total de hilos como la multiplicación del número de grupos y el tamaño de grupo, el bloque for de la línea 11 a 51 realiza la copia de elementos de elmA a elmL, cada hilo se encargará de realizar la copia de un elemento de elmA en el lugar correspondiente de elmL; el bloque while de las líneas 19 a 31 realiza la búsqueda de los índices de fila y columna del elemento actual en elmA, teniendo estos índices se calcula la ubicación del elemento en elmL, el bloque

if de las líneas 38 a 47 se encarga de esta tarea, finalmente en la línea 50 se ubica el elemento en el lugar correspondiente. Esta función es llamada una sola vez para rellenar los espacios de *elmL* con los elementos de *elmA* correspondientes, luego de esto la factorización se realizará en sitio sobre *elmL*.

```
2
3
4
5
            //indice global de hilo
6
            int tx = get_global_id(0);
7
8
            //número total de hilos
            int ntx = get_num_groups(0)*get_local_size(0);
9
10
11
            for (int x = tx; x < size_elmA; x += ntx) {</pre>
12
                    //variable donde se guardará el índice de columna y de fila
13
                    //respectivamente del elemento elmA[tx]
14
                    int col = 0;
15
                    int row = 0;
                    //variable que controlará la salida del bucle.
16
17
                    int key = 0;
18
19
                    while (key == 0) {
20
                            if (colA[col] > x) {
21
                                    col -= 1;
22
                                    row = rowA[x - col - 1];
23
                                    key = 1;
24
25
                            else if (colA[col] == x) {
26
                                    row = col;
27
                                    key = 1;
28
29
                            else
30
                                    col += 1;
31
32
                    //hasta aquí para un elemento de ElmA, ya se
33
                    //se cuenta con el índice de fila y de columna (row y col)
                    //ahora necesitamos ubicar este elemento en elmL
34
35
36
                    int posf;//posición final en elmL
37
38
                    if (col != row) {
39
                            int pos = colL[col] - col;
40
                            int cont = 0;
41
                            key = 0;
42
                            while (key == 0)
43
                                    rowL[pos + cont] == row ? key = 1 : cont++;
44
                            posf = colL[col] + cont + 1;
45
46
                    else
47
                            posf = coll[col];
48
                    //escribiendo elementos de elmA en el lugar
49
                    //correspondiente de elmL
50
                    elmL[posf] = elmA[x];
51
            }
```

Figura 35. Cholesky: función de dispositivo sparse_fill

Similar al proceso sobre matrices densas, la factorización fue dividida en tres partes, un paso en el que modificamos una columna *j* tal que:

$$A_{n,j} = \frac{A_{r,j}}{\sqrt{A_{j,j}}}, n > j$$

```
void Chol_sparse_1(__global double* elmL, __global int* colL,
1
2
             int order, int step)
3
             //es procedimiento solo deberá realizarse hasta
4
5
             //la columna order-2
             //Siendo order el orden de la matriz, la última
6
             //columna tendrá índice order-1, donde solo contiene
7
8
             //el elemento diagonal
9
             if (step < order - 1) {</pre>
                      //obteniendo la ubicación del primer elemento
10
11
                      //de la columna determinada por step
                     int pos = coll[step];
12
13
14
                      //posición siguiente
15
                     int pos_nxt = coll[step + 1];
16
17
                      //indice global de hilo
18
                     int tx = get_global_id(0);
19
20
                      //número total de hilos
                     int ntx = get_num_groups(0)*get_local_size(0);
21
22
23
                      //raíz cuadrada del elemento diagonal. El primer
24
                      //elemento guardado en elmA para cada columna es el
25
                      //elemento diagonal
26
                      double _r = sqrt(elmL[pos]);
27
                     //cada hilo se encargará de modificar un elemento
28
29
                      //de la columna
30
                     for (int x = tx + pos + 1; x < pos_nxt; x += ntx) {</pre>
31
                              elmL[x] = elmL[x] / _r;
32
33
34
```

Figura 36. Cholesky: función de dispositivo Chol_sparse_1

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 36, como puede observarse para realizar esta operación solo es necesario introducir como argumentos *elmL* y *colL*, la división se realiza sobre todos los elementos de una columna sin importar el índice de fila que tengan, cada hilo se encargará de realizar la modificación de uno de los elementos de la columna, en la línea 12 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna (la columna es determinada con el argumento *step*), de manera similar en la línea 15 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna siguiente, los elementos en *elmL* con índices entre estos dos valores son los que se tienen que modificar. En las líneas 18 y 21 se obtiene el índice global de hilo y el número total de hilos respectivamente, en la línea 26 se obtiene la raíz cuadrada del elemento diagonal y finalmente en la línea 31 se modifican los elementos de la columna.

Otro paso consiste en modificar las columnas siguientes a la columna j modificada en el paso anterior tal que:

$$A_{n,s} = A_{n,s} - A_{s,j}A_{n,j}, s > j, n \ge s$$

La función de dispositivo que realiza esta tarea se muestra en la Figura 37, en este caso, cada grupo y sus hilos se encargarán de modificar una columna. En las líneas 11 y 14 se obtiene el índice de grupo y el índice local de hilo respectivamente, la línea 17 obtiene la ubicación en *elmL* del primer elemento de la columna (determinada por *step*), en la línea 20 se calcula la ubicación en *rowL* de primer índice de fila almacenado para la columna, en la línea 23 se calcula el número de elementos en *rowL* que le corresponden a la columna, debe recordarse que no se almacena el índice de fila del elemento diagonal por ser obvio, el bloque *if* de las líneas 25 a 61 índica que el bloque *for* anidado (líneas 27 a 60) solo se utiliza cuando el número de elementos en *rowL* para la columna es por lo menos uno.

El bloque *for* de las líneas 27 a 60 realizará las modificaciones de todas las columnas siguientes a la columna actual (argumento *step*).

Estos dos primeros pasos deben realizarse uno después de otro *n-1* veces.

Un tercer paso es necesario para terminar la factorización, este modifica los elementos de la diagonal tal que:

$$A_{j,j} = A_{j,j}$$

La función de dispositivo que realiza esta tarea se muestra en la Figura 38, todos los hilos de todos los grupos se encargan cada uno de modificar un elemento de la diagonal, este paso solo se realiza una vez luego de realizada las *n-1* iteraciones con los dos pasos vistos previamente.

```
1
2
3
    {
4
            //es procedimiento solo deberá realizarse hasta
5
            //la columna order-2
            //Siendo order el orden de la matriz, la última
7
            //columna tendrá índice order-1, donde solo contiene
8
            //el elemento diagonal
            if (step < order - 1) {</pre>
9
10
                    //indice local de grupo
11
                    int bx = get_group_id(0);
12
                    //indice local de hilo
13
                    int tx = get_local_id(0);
14
15
                    //posición del primer elemento de la columna de elmL
16
                    int pos_elmL = colL[step];
17
18
                    //posición del primer elemento de la columna en rowL
19
20
                    int pos_rowL = pos_elmL - step;
21
                    //número de elementos en rowL para la columna
22
23
                    int nEl = colL[step + 1] - pos_elmL - 1;
24
25
                    if (nEl > 0) {
26
27
                             for (int b = bx; b < nEl; b += get_num_groups(0)) {</pre>
                                     //indice de la columna a modificar
28
29
                                     int colId = rowL[pos_rowL + b];
30
                                     //primer factor
31
                                     double k1 = elmL[pos_elmL + b + 1];
32
33
                                     for (int x = b + tx; x < nE1;
34
35
                                             x += get_local_size(0)) {
36
                                             //índice de la fila en la columna a
37
38
                                             //modificar
                                             int rowId = rowL[pos_rowL + x];
39
40
41
                                             //segundo factor
                                             double k2 = elmL[pos_elmL + x + 1];
42
43
                                             int posf;//posición final en elmL
44
45
                                             if (colId != rowId) {
46
47
                                                     int pos = colL[colId] - colId;
                                                     int cont = 0;
48
                                                     int key = 0;
49
50
                                                     while (key == 0)
                                                             rowL[pos + cont] ==
rowId ? key = 1 : cont++;
51
52
53
                                                     posf = colL[colId] + cont + 1;
54
55
                                             else
56
                                                     posf = coll[colId];
57
                                             elmL[posf] -= k1 * k2;
58
59
60
61
62
63
```

Figura 37. Cholesky: función de dispositivo Chol_sparse_2

Figura 38. Cholesky: función de dispositivo Chol_sparse_3

Terminada la factorización se procede a solucionar los sistemas equivalentes, es decir si Ax = b es el sistema original, luego de la factorización quedará:

$$LL^Tx = b$$

Haciendo $L^T x = y$ entonces Ly = b. Primero se resuelve Ly = b y finalmente $L^T x = y$, cuya solución en x es la solución del sistema Ax = b.

Para resolver el sistema Ly = b, usaremos la eliminación típica de Gauss Jordan aplicada a una matriz triangular inferior, en cada iteración los elementos de cada columna fuera de la diagonal van volviéndose cero progresivamente, hasta llegar a una matriz diagonal, las operaciones sobre L pueden obviarse y dejarla intacta para poder usarla al resolver el siguiente sistema $L^Tx = y$, es decir solo se modificará el vector de constantes b tal que para la iteración j:

$$b_n = b_n - b_j \frac{L_{n,j}}{L_{i,j}}, n > j$$

El kernel que realiza esta operación se muestra en la Figura 39, el número de iteración y por la tanto el índice de la columna implicada está dada por el valor del argumento step, esta función debe ser llamada n-l veces para terminar con la eliminación, luego será necesario dividir los elementos de b entre los elementos de la diagonal (función de dispositivo de la Figura 40), este vector será la solución del sistema Ly = b, que será usado como argumento al resolver el siguiente sistema triangular.

```
2
3
4
5
            if (step < order - 1) {
6
                    //indice global de hilo
                   int tx = get_global_id(0);
8
9
                   //número total de hilos
10
                   int ntx = get_num_groups(0)*get_local_size(0);
11
12
                   //posición del elemento diagonal
13
                   int diag = coll[step];
14
15
                   //número de elementos para una columna
                   int nEl = colL[step + 1] - diag;
16
17
18
                   //deben haber elementos aparte del elemento diagonal
                   if (nEl > 1) {
19
20
21
                           //elemento diagonal y elemento correspondiente
22
                           double diag_elm = elmL[diag];
23
24
                           double elmb = b[step];
25
                           //posición en rowL del primer elemento fuera de
26
27
                           //la diagonal (indice)
                           int pos_rowL = diag - step;
28
29
                           double k = elmb / diag_elm;
30
                           for (int x = tx; x < nEl - 1; x += ntx) {
31
32
                                   int rowId = rowL[pos_rowL + x];
                                   b[rowId] -= k * elmL[diag + x + 1];
33
34
35
36
37
```

Figura 39. Cholesky: función de dispositivo Chol_sparse_4

```
void Chol_sparse_5(__global double* elmL,
2
             __global int* colL, __global double* b, int order)
4
             //indice global de hilo
5
            int tx = get_global_id(0);
             //número total de hilos
             int ntx = get_num_groups(0)*get_local_size(0);
9
10
             for (int x = tx; x < order; x += ntx) {
11
                     //ubicación del elemento diagonal en elmL
12
13
                     int diag = colL[x];
14
                     b[x] = b[x] / elmL[diag];
15
     Figura 40. Cholesky: función de dispositivo Chol_sparse_5
```

Para resolver el sistema $L^T x = y$, se usará la misma matriz que almacena L, el artificio es asumir que está ordenada por filas y no por columnas, hecho esto la transposición de la matriz ya no será necesaria. L^T es una matriz triangular

superior, se resolverá el sistema obteniendo el valor de las variables una a continuación de otra desde abajo, el valor de la variable en la posición *i* sería:

$$x_{i} = \frac{y_{i} - \sum_{j=i+1}^{n-1} x_{j} L_{i,j}}{L_{i,i}}$$

```
3
4
5
             //indice de bloque
6
            int bx = get_group_id(0);
8
             //indice global de hilo
9
            int tx = get_local_id(0);
10
11
             //número total de hilos
            int ntx = get_local_size(0);
12
13
14
             //indice de la fila de L implicada
15
            int fL = order - step - 1;
16
17
             //ubicación de elemento diagonal de la fila en elmL
            int diagId = colL[fL];
18
19
20
             //número de elementos en la fila
21
            int nEl = (fL == order - 1 ? 1 : colL[fL + 1] - diagId);
22
23
             //ubicación del índice del primer elemento fuera de la
             //diagonal en rowL
24
25
            int pos_rowL = diagId - fL;
26
27
             //elemento diagonal
            double diag = elmL[diagId];
28
29
             //valor donde se guardará una suma parcial
30
31
             double sum = 0.0;
32
33
             //solo se usará un bloque
34
            if (bx == 0) {
                     for (int x = tx; x < nEl - 1; x += ntx) {
    int rowId = rowL[pos_rowL + x];</pre>
35
36
37
                             sum += elmL[diagId + x + 1] * y[rowId];
38
39
                     partialSum[tx] = sum;
40
41
                     //reducción en paralelo pare determinar la suma de
                     //todos los elementos de partialSum
42
                     for (int stride = ntx / 2; stride > 0; stride /= 2) {
43
44
                             barrier(CLK_LOCAL_MEM_FENCE);
45
                             if (tx < stride) {</pre>
46
                                     partialSum[tx] += partialSum[tx + stride];
47
48
                     if (tx == 0)
49
                             y[fL] = (y[fL] - partialSum[0]) / diag;
50
51
```

Figura 41. Cholesky: función de dispositivo *Chol_sparse_6*

Se usa el mismo vector y para almacenar la solución desde abajo. La función de dispositivo que resuelve el sistema $L^Tx=y$ se muestra en la Figura 41, cada llamada a esta función calculará una variable. Ya que el sistema se resuelve desde abajo, para un paso j se calculará el valor de la variable i tal que i=n-j-1, la fila de L usada para calcular esta variable tiene el mismo índice (se obtiene en la línea 15). Esta función de dispositivo usará memoria local para calcular la sumatoria $\sum_{j=i+1}^{n-1} x_j L_{i,j}$ por reducción en paralelo, el tamaño de esta memoria local será el mismo que el tamaño de un grupo de trabajo, por esta razón, en este caso solo se utilizará un grupo de trabajo al usar esta función (en la línea 34 se especifica que el grupo de trabajo utilizado es únicamente el con id cero).

```
kernel void
             _sparse_c(__global double* elmA, __global int* colA, __global int* rowA, __global double* elmL, __global int* colL,
     Cholesky_sparse_c(__global double* elmA,
4
               global int* rowL, int size_elmA, int size_colA,
               _global double* b, __local double* partialSum, int step,
             int sstep)
6
     {
8
9
             if (sstep == -1)
10
                      sparse_fill(elmA, colA, rowA, elmL, colL
11
                              rowL, size_elmA);
             //la iteración principal se realiza en el host
12
13
14
                      //se modificará una sola columna
             else if (sstep == 0)
15
16
                      Chol_sparse_1(elmL, colL, size_colA, step);
17
18
             //se modificará las columnas siguientes a la columna
             //implicada
19
20
             else if (sstep == 1)
                      Chol_sparse_2(elmL, colL, rowL, size_colA, step);
21
22
23
             //una vez usada las funciones anteriores todavía es
             //necesario modificar la diagonal D[i]=sqrt(D[i]) para
24
25
             //concluir la factorización
26
             else if (sstep == 2)
27
                      Chol_sparse_3(elmL, colL, size_colA);
28
             //hasta aquí la factorización está concluida, la parte
29
             //triangular inferior de A contendrá la factorización de
30
             //Cholesky es decir L
             //desde aquí se procederá a solucionar el sistema
31
32
             //L*transpose(L)x=b
33
34
             //se resolverá el sistema Ly=b, y=transpose(L)*x
35
             else if (sstep == 3)
                      Chol_sparse_4(elmL, colL, rowL, b, size_colA, step);
36
37
             else if (sstep == 4)
38
                      Chol_sparse_5(elmL, colL, b, size_colA);
39
40
             //se resolverá el sistema transpose(L)*x=y
             //"y" ya debe haberse obtenido usando la función anterior
41
42
             else
43
                      Chol_sparse_6(elmL, colL, rowL, b, size_colA, step,
44
                               partialSum):
```

Figura 42. Cholesky: Kernel para trabajar sobre una matriz en formato CSC

En la línea 35 a 37 se calcula una suma parcial almacenándola en la variable *sum* (declarada en la línea 31), luego este valor es almacenado en memoria local (línea 39), los elementos de *partialsum* se suman en paralelo en el bloque *for* de las líneas 43 a 48, finalmente se calcula la variable correspondiente en la línea 50, esta tarea solo la realiza un hilo (en la línea 49 se especifica el id del hilo que realiza la operación). Esta función de dispositivo debe ser llamada *n* veces para calcular las *n* variables.

El kernel principal se muestra en la Figura 42, como puede observarse el kernel llama a las funciones de dispositivo mencionadas anteriormente, dependiendo del valor del argumento *sstep*, un valor igual a -1 corresponde al proceso de copia de los elementos de *elmA* en el lugar correspondiente de *elmL*, los valores 0, 1 y 2 indican que nos encontramos en el proceso de factorización y los valores 3, 4 y 5 que estamos en el proceso de solución de los sistemas triangulares.

El kernel es llamado 1 vez para sstep=-1, n-1 veces para sstep igual a 0 y 1 (con step de 0 a n-2), solo una vez para sstep=2, n-1 veces para sstep=3 (con step de 0 a n-2), una vez para sstep=4 y finalmente n veces para sstep=5(con step de 0 a n-1).

4.1.6.2. Factorización LDL^T

El proceso para resolver un sistema usando la factorización LDL^T es similar al que utilizamos en la factorización de Cholesky.

Para resolver un sistema de ecuaciones usando la factorización LDL^T, son necesarias dos etapas, uno de la factorización misma y otro que consiste en resolver los sistemas triangulares generados después de la factorización.

El proceso es similar al usado sobre matrices densas, sin embargo para empezar con la factorización es necesario realizar una factorización simbólica, esto con el fin de obtener la cantidad de memoria máxima requerida para almacenar los elementos de la factorización, el espacio requerido para almacenar el campo *elements* de una factorización LDL^T es el mismo requerido para una

factorización de Cholesky, por tal la misma función usada para la factorización de Cholesky simbólica es usada para la factorización LDL^T simbólica (función Julia de la Figura 34 pág. 86), esta función obtiene un objeto del tipo *SparseMS* cuyos campos *colIdx* y *rowIdx* son los finales, sin embargo el campo *elements* solo contiene ceros, pero la cantidad de elementos es la necesaria para guardar la factorización.

```
__global int* rowL, int size_elmA)
3
4
            //indice global de hilo
6
            int tx = get_global_id(0);
            //número total de hilos
8
9
            int ntx = get_num_groups(0)*get_local_size(0);
10
            for (int x = tx; x < size_elmA; x += ntx) {</pre>
11
12
                    //variable donde se guardará el índice de columna y de fila
                    //respectivamente del elemento elmA[tx]
13
                    int col = 0;
14
15
                    int row = 0;
                    //variable que controlará la salida del bucle.
16
17
                    int key = 0;
18
                    while (key == 0) {
19
20
                            if (colA[col] > x) {
21
                                    col -= 1;
22
                                    row = rowA[x - col - 1];
23
                                    key = 1;
24
25
                            else if (colA[col] == x) {
26
                                    row = col;
27
                                    key = 1;
28
29
                            else
30
                                    col += 1;
31
                    //hasta aquí para un elemento de ElmA, ya se
32
33
                    //se cuenta con el índice de fila y de columna (row y col)
34
                    //ahora necesitamos ubicar este elemento en elmL
35
36
                    int posf;//posición final en elmL
37
                    if (col != row) {
38
39
                            int pos = coll[col] - col;
40
                            int cont = 0;
41
                            key = 0;
42
                            while (key == 0)
43
                                    rowL[pos + cont] == row ? key = 1 : cont++;
                            posf = colL[col] + cont + 1;
44
45
46
                    else
                            posf = colL[col];
47
48
                    //escribiendo elementos de elmA en el lugar
                    //correspondiente de elmL
49
                    elmL[posf] = elmA[x];
50
51
```

Figura 43. LDL^T: función de dispositivo *sparse_fill_c*

Sea A la matriz original en formato CSC y L la matriz que contendrá la factorización LDL^T también en formato CSC. Un paso previo es necesario antes de proceder con la factorización, este paso consiste en copiar los elementos del campo *elements* de A, en el lugar correspondiente del campo *elements* de L para poder realizar la factorización en sitio sobre este último, la función de dispositivo que realiza esta tarea se muestra la Figura 43.

La función acepta seis argumentos, los campos de *A* y *L* en formato CSC (*elmA*, *colA*, *rowA* y *elmL*, *colL*, *rowL* respectivamente) y el tamaño de *elmA*. En la línea 6 se obtiene el índice global del hilo, la línea 10 calcula el número total de hilos como la multiplicación del número de grupos y el tamaño de grupo, el bloque *for* de la línea 11 a 51 realiza la copia de elementos de *elmA* a *elmL*, cada hilo se encargará de realizar la copia de un elemento de *elmA* en el lugar correspondiente de *elmL*; el bloque *while* de las líneas 19 a 31 realiza la búsqueda de los índices de fila y columna del elemento actual en *elmA*, teniendo estos índices se calcula la ubicación del elemento en *elmL*, el bloque *if* de las líneas 38 a 47 se encarga de esta tarea, finalmente en la línea 50 se ubica el elemento en el lugar correspondiente. Esta función es llamada una sola vez para rellenar los espacios de *elmL* con los elementos de *elmA* correspondientes, luego de esto la factorización se realizará en sitio sobre *elmL*.

El proceso de factorización fue dividido en dos partes, un paso en el que modificamos una columna *j* tal que:

$$A_{n,j} = \frac{A_{n,j}}{A_{j,j}}, n > j$$

La función de dispositivo que realiza esta operación se muestra en la Figura 44, como puede observarse para realizar esta operación solo es necesario introducir como argumentos *elmL* y *colL*, la división se realiza sobre todos los elementos de una columna sin importar el índice de fila que tengan, cada hilo se encargará de realizar la modificación de uno de los elementos de la columna, en la línea 12 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna (la columna es determinada con el argumento *step*), de manera similar en la

línea 15 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna siguiente, los elementos en *elmL* con índices entre estos dos valores son los que se tienen que modificar. En las líneas 18 y 21 se obtienen el índice global de hilo y el número total de hilos respectivamente, en la línea 26 se obtiene el elemento diagonal y finalmente en la línea 31 se modifican los elementos de la columna.

```
3
4
            //es procedimiento solo deberá realizarse hasta
            //la columna order-2
6
            //Siendo order el orden de la matriz, la última
            //columna tendrá índice order-1, donde solo contiene
8
            //el elemento diagonal
9
            if (step < order - 1) {</pre>
10
                    //obteniendo la ubicación del primer elemento
11
                    //de la columna determinada por step
12
                    int pos = coll[step];
13
14
                    //posición siguiente
15
                    int pos_nxt = coll[step + 1];
16
17
                    //indice global de hilo
18
                    int tx = get_global_id(0);
19
20
                    //número total de hilos
21
                    int ntx = get_num_groups(0)*get_local_size(0);
22
23
                    //elemento diagonal. El primer
24
                    //elemento guardado en elmA para cada columna es el
25
                    //elemento diagonal
26
                    double _r = elmL[pos];
27
28
                    //cada hilo se encargará de modificar un elemento
29
                    //de la columna
30
                    for (int x = tx + pos + 1; x < pos_nxt; x += ntx) {</pre>
31
                            elmL[x] /= _r;
32
33
```

Figura 44. LDL^T: función de dispositivo *ldlt_sparse_1_c*

Otro paso consiste en modificar las columnas siguientes a la columna *j* modificada en el paso anterior tal que:

$$A_{n,s}=A_{n,s}-A_{s,j}A_{n,j}A_{j,j}, s>j, n\geq s$$

```
1
2
3
    {
4
            //es procedimiento solo deberá realizarse hasta
5
            //la columna order-2
            //Siendo order el orden de la matriz, la última
7
            //columna tendrá índice order-1, donde solo contiene
            //el elemento diagonal
8
9
            if (step < order - 1) {</pre>
10
                    //indice local de grupo
11
                    int bx = get_group_id(0);
12
13
                    //indice local de hilo
                    int tx = get_local_id(0);
14
15
16
                    //posición del primer elemento de la columna de elmL
                    int pos_elmL = colL[step];
17
18
19
                    //elemento diagonal
                    double diag = elmL[pos_elmL];
20
21
22
                    //posición del primer elemento de la columna en rowL
23
                    int pos_rowL = pos_elmL - step;
24
25
                    //número de elementos en rowL para la columna
26
                    int nEl = colL[step + 1] - pos_elmL - 1;
27
                    if (nEl > 0) {
28
29
30
                            for (int b = bx; b < nEl; b += get num groups(0)) {
                                     //índice de la columna a modificar
31
32
                                    int colId = rowL[pos_rowL + b];
33
34
                                    //primer factor
35
                                    double k1 = elmL[pos_elmL + b + 1];
36
37
                                    for (int x = b + tx; x < nE1;
38
                                            x += get_local_size(0)) {
39
40
                                            //índice de la fila en la columna
41
                                            //a modificar
42
                                            int rowId = rowL[pos_rowL + x];
43
44
                                            //segundo factor
45
                                            double k2 = elmL[pos_elmL + x + 1];
46
47
                                            int posf;//posición final en elmL
48
49
                                            if (colId != rowId) {
                                                    int pos = colL[colId] - colId;
50
                                                    int cont = 0;
51
52
                                                    int key = 0;
53
                                                    while (key == 0)
54
                                                            rowL[pos + cont] ==
                                                            rowId ? key = 1 : cont++;
55
56
                                                    posf = coll[colId] + cont + 1;
57
                                            else
58
59
                                                    posf = coll[colId];
60
61
                                            elmL[posf] -= diag * k1 * k2;
62
63
64
65
            }
66
```

Figura 45. LDL^T: función de dispositivo *ldlt_sparse_2_c*

La función de dispositivo propuesta que realiza esta tarea se muestra en la Figura 45, en este caso, cada grupo y sus hilos se encargarán de modificar una columna. En las líneas 11 y 14 se obtiene el índice de grupo y el índice local de hilo respectivamente, la línea 17 obtiene la ubicación en *elmL* del primer elemento de la columna actual (determinada por *step*), en la línea 20 se obtiene el elemento diagonal, en la línea 23 se calcula la ubicación en *rowL* de primer índice de fila almacenado para la columna, en la línea 26 se calcula el número de elementos en *rowL* que le corresponden a la columna, debe recordarse que no se almacena el índice de fila del elemento diagonal por ser obvio, el bloque *if* de las líneas 28 a 64 índica que el bloque *for* anidado (líneas 30 a 63) solo se utiliza cuando el número de elementos en *rowL* para la columna es por lo menos uno.

El bloque *for* de las líneas 30 a 63 realizará las modificaciones de todas las columnas siguientes a la columna actual (argumento *step*).

Estos dos primeros pasos deben realizarse uno después de otro *n-1* veces para terminar con la factorización. Terminada la factorización en sitio, el campo *elements* de *L* contendrá los elementos de la factorización LDL^T, este incluye los elementos de la diagonal *D*, almacenados como si fueran la diagonal de *L*, es decir, en formato CSC, *L* almacena una matriz triangular inferior, debe recordarse que L por si sola tiene sus elementos diagonales iguales a 1, esto es implícito y debe tenerse en cuenta cuando llegue el momento de resolver el sistema.

Terminada la factorización se procede a solucionar los sistemas equivalentes, es decir si Ax = b es el sistema original, luego de la factorización quedará:

$$LDL^Tx = b$$

Haciendo $DL^Tx = y$ entonces Ly = b, es el primer sistema a resolverse para y, luego haciendo $L^Tx = z$ entonces Dz = y es el segundo sistema a resolver en z, que es simple ya que D es una matriz diagonal, finalmente se resuelve $L^Tx = z$, cuya solución en x es la solución del sistema Ax = b.

Para resolver el sistema Ly = b, usaremos la eliminación típica de Gauss Jordan aplicada a una matriz triangular inferior, en cada iteración los elementos de cada columna fuera de la diagonal van volviéndose cero progresivamente, hasta llegar a una matriz diagonal, resolver un sistema triangular nos permite obviar las operaciones sobre L y dejarla intacta para poder usarla al resolver el sistema $L^Tx = z$, solo se modificará el vector de constantes b tal que para la iteración j:

$$b_n = b_n - b_i L_{n,i}, n > j$$

```
void ldlt_sparse_3_c(__global double* elmL,
               _global int* colL, __global int* rowL, __global double* b,
3
             int order, int step)
4
5
             if (step < order - 1) {</pre>
                      //indice global de hilo
6
                      int tx = get_global_id(0);
8
9
                      //número total de hilos
10
                      int ntx = get_num_groups(0)*get_local_size(0);
11
12
                      //posición del elemento diagonal
                      int diag = coll[step];
13
14
15
                      //número de elementos para una columna
                      int nEl = colL[step + 1] - diag;
16
17
18
                      //deben haber elementos aparte del elemento diagonal
                      if (nEl > 1) {
19
20
21
                               //elemento diagonal y elemento correspondiente en b
                               double diag_elm = elmL[diag];
22
23
                               double elmb = b[step];
24
25
                               //posición en rowL del primer elemento fuera de la
26
                               //diagonal (indice)
27
                               int pos_rowL = diag - step;
28
29
                               for (int x = tx; x < nEl - 1; x += ntx) {
                                       int rowId = rowL[pos_rowL + x];
b[rowId] -= elmb * elmL[diag + x + 1];
30
31
32
                               }
33
34
```

Figura 46. LDL^T: función de dispositivo *ldlt_sparse_3_c*

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 46, el índice de la columna actual está dada por el valor del argumento step, esta función debe ser llamada n-l veces para terminar con la eliminación, al final b contendrá la solución del sistema Ly = b, que será usado como argumento al resolver el siguiente sistema.

Luego se procede a resolver el sistema Dz = y, ya que D es diagonal, la solución se reduce a realizar la operación:

$$b_j = \frac{b_j}{D_{j,j}}$$

La función de dispositivo que realiza esta operación se muestra en la Figura 47, esta función se ejecuta una sola vez.

```
void ldlt_sparse_4_c(__global double* elmL,
              _global int* colL, __global double* b, int order)
4
             //indice global de hilo
            int tx = get_global_id(0);
6
8
             //número total de hilos
             int ntx = get_num_groups(0)*get_local_size(0);
10
             for (int x = tx; x < order; x += ntx) {
                     //ubicación del elemento diagonal en elmL
13
                     int diag = colL[x];
14
                     b[x] /= elmL[diag];
15
```

Figura 47. LDL^T: función de dispositivo *ldlt_sparse_4_c*

Para resolver el sistema $L^T x = z$, se usará la misma matriz que almacena L, el artificio es asumir que está ordenada por filas y no por columnas, con esto la transposición de la matriz ya no será necesaria. L^T es una matriz triangular superior, se resolverá el sistema obteniendo el valor de las variables una a continuación de otra desde abajo, recordar que los elementos de la diagonal de L son 1 aunque no estén almacenados, así, el valor de la variable en la posición i sería:

$$x_i = z_i - \sum_{j=i+1}^{n-1} z_j L_{i,j}$$

Se usa el mismo vector z para almacenar la solución desde abajo. La función de dispositivo propuesta que resuelve el sistema $L^Tx=z$ se muestra en la Figura 48, cada llamada a esta función calculará una variable. Ya que el sistema se resuelve desde abajo, para un paso j se calculará el valor de la variable i tal que i=n-j-1, la fila de L usada para calcular esta variable tiene el mismo índice (se obtiene en la línea 16). Esta función de dispositivo usará memoria local

para calcular la sumatoria $\sum_{j=i+1}^{n-1} z_j L_{i,j}$ por reducción en paralelo, el tamaño de esta memoria local será el mismo que el tamaño de un grupo de trabajo, por esta razón, en este caso solo se utilizará un grupo de trabajo al usar esta función (en la línea 32 se especifica que el grupo de trabajo utilizado es únicamente el con id cero).

```
3
              local double* partialSum)
4
6
             //indice de bloque
             int bx = get_group_id(0);
8
9
             //indice global de hilo
10
            int tx = get_local_id(0);
11
             //número total de hilos
12
13
            int ntx = get_local_size(0);
14
            //índice de la fila de L implicada
15
16
            int fL = order - step - 1;
17
             //ubicación de elemento diagonal de la fila en elmL
18
19
            int diagId = colL[fL];
20
            //número de elementos en la fila
21
22
            int nEl = (fL == order - 1 ? 1 : colL[fL + 1] - diagId);
23
            //ubicación del índice del primer elemento fuera de la
24
25
             //diagonal en rowL
            int pos rowL = diagId - fL;
26
27
28
             //valor donde se guardará una suma parcial
29
            double sum = 0.0;
30
31
             //solo se usará un bloque
32
             if (bx == 0) {
33
                    for (int x = tx; x < nEl - 1; x += ntx) {
34
                            int rowId = rowL[pos_rowL + x];
                            sum += elmL[diagId + x + 1] * y[rowId];
35
36
37
                    partialSum[tx] = sum;
38
39
                    //reducción en paralelo pare determinar la suma de
40
                    //todos los elementos de partialSum
                    for (int stride = ntx / 2; stride > 0; stride /= 2) {
41
                            barrier(CLK_LOCAL_MEM_FENCE);
42
43
                            if (tx < stride) {</pre>
44
                                     partialSum[tx] += partialSum[tx + stride];
45
46
47
                    if (tx == 0)
                            y[fL] -= partialSum[0];
48
49
```

Figura 48. LDL^T: función de dispositivo *ldlt_sparse_5_c*

En la línea 33 a 36 se calcula una suma parcial almacenándola en la variable *sum* (declarada en la línea 29), luego este valor es almacenado en memoria

local (línea 37), los elementos de *partialsum* se suman en paralelo en el bloque *for* de las líneas 41 a 46, finalmente se calcula la variable correspondiente en la línea 48, esta tarea solo la realiza un hilo (en la línea 47 se especifica el id del hilo que realiza la operación). Esta función de dispositivo debe ser llamada *n* veces para calcular las *n* variables.

```
1
       kernel void
               se_c(__global double* elmA, __global int*
_global int* rowA, __global double* elmL,
2
     ldlt_sparse_c(
                                               global int* colA,
3
               global int* colL, __global int* rowL, int size_elmA,
4
             int size_colA, __global double* b,
__local double* partialSum, int step, int sstep)
5
6
7
8
             if (sstep == -1)
                      sparse_fill_c(elmA, colA, rowA, elmL, colL,
9
10
                               rowL, size_elmA);
11
             //la iteración principal se realiza en el host
12
13
             //se modificará una sola columna
14
             else if (sstep == 0)
                      ldlt_sparse_1_c(elmL, colL, size_colA, step);
15
16
17
             //se modificará las columnas siguientes a la columna
             //implicada
18
19
             else if (sstep == 1)
20
                      ldlt_sparse_2_c(elmL, colL, rowL, size_colA, step);
21
22
             //hasta aquí la factorización está concluida, la parte
23
             //triangular inferior de A contendrá la factorización
             //LDLt es decir L y D
24
25
             //desde aquí se procederá a solucionar el sistema
26
             //L*D*transpose(L)x=b
27
28
             //se resolverá el sistema Ly=b, y=D*transpose(L)*x
29
             else if (sstep == 2)
                      ldlt_sparse_3_c(elmL, colL, rowL, b, size_colA, step);
30
31
32
             //se resolverá el sistema Dz=y, z=transpose(L)*x
33
             else if (sstep == 3)
                      ldlt_sparse_4_c(elmL, colL, b, size_colA);
34
35
36
             //se resolverá el sistema transpose(L)*x=z
37
             //"z" ya debe haberse obtenido usando la función anterior
38
             else
                      ldlt_sparse_5_c(elmL, colL, rowL, b, size_colA, step,
39
40
                               partialSum);
41
```

Figura 49. LDL^T: Kernel para trabajar sobre una matriz en formato CSC

El kernel principal se muestra en la Figura 49, como puede observarse el kernel llama a las funciones de dispositivo mencionadas anteriormente, dependiendo del valor del argumento *sstep*, para un valor de -1 copiará los elementos de *elmA* en el lugar correspondiente de *elmL*, *los* valores 0 y 1 indican que nos

encontramos en el proceso de factorización y los valores 2, 3 y 4 que estamos en el proceso de solución de los sistemas triangulares.

El kernel es llamado 1 vez para sstep=-1, n-1 veces para sstep igual a 0 y 1 (con step de 0 a n-2), n-1 veces para sstep=2 (con step de 0 a n-2), una vez para sstep=3 y finalmente n veces para sstep=4(con step de 0 a n-1).

4.1.6.3. Gradientes conjugados

La ventaja de este método es que solo requiere de operaciones sencillas en cada iteración, casi todas estas operaciones puede hacerse haciendo uso de las funciones incorporadas en ArrayFire, la única operación que necesita una función de dispositivo es la que se encargará de multiplicar matriz por vector, con la matriz almacenada en formato CSC.

La matriz en su forma densa puede descomponerse en tres términos: una matriz estrictamente triangular inferior, una matriz estrictamente triangular superior y una matriz diagonal, entonces se cumple que:

$$L + D + U = A$$

Siendo L la matriz estrictamente triangular inferior con los elementos de la parte estrictamente triangular inferior de A, U la matriz estrictamente triangular superior con los elementos de la parte estrictamente triangular superior de A y D la matriz diagonal con los elementos de la diagonal de A.

Ya que se trata de una matriz simétrica entonces $U = L^T$ por lo que:

$$L + D + L^T = A$$

Multiplicando por un vector *b*:

$$Ab = Lb + Db + L^Tb$$

La matriz en formato CSC almacena L y D, para L^T puede usarse la misma L asumiendo que está ordenada por filas como se hizo al resolver los sistemas triangulares de la factorización de Cholesky o LDL^T .

```
void sparse_mat_vec_mul1(__global double* elmA,
1
              __global int* colA, __global int* rowA, __global double* b
__global double* c, int order, __local double* partialSum)
2
                                                            _global double* b,
3
4
5
              //indice de bloque
              int bx = get_group_id(0);
6
7
              //indice local de hilo
8
9
              int tx = get_local_id(0);
10
              for (int row = bx; row < order; row += get_num_groups(0)) {</pre>
11
12
                       //ubicación de elemento diagonal de la fila en elmL
13
14
                       int dId = colA[row];
15
16
                       //número de elementos en la fila
17
                       int nEl = (row == order - 1 ? 1 : colA[row + 1] - dId);
18
19
                       //ubicación del índice del primer elemento fuera de la
20
                       //diagonal en rowL
                       int baseId = dId - row;
21
22
23
                       //valor donde se guardará una suma parcial
24
                       double sum = 0.0;
25
26
                       for (int x = tx; x < nEl - 1; x += get_local_size(0)) {</pre>
                               int cId = rowA[baseId + x];
27
                                sum += elmA[dId + x + 1] * b[cId];
28
29
30
                       partialSum[tx] = sum;
31
32
                       //reducción en paralelo pare determinar la suma de
33
                       //todos los elementos de partialSum
34
                       for (int stride = get_local_size(0) / 2;
35
                               stride > 0; stride /= 2) {
                               barrier(CLK LOCAL MEM FENCE);
36
37
                               if (tx < stride) {</pre>
38
                                         partialSum[tx] += partialSum[tx + stride];
39
40
                      if (tx == 0)
41
42
                               c[row] = partialSum[0];
43
```

Figura 50. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul1

La función de dispositivo propuesta para obtener L^Tb se muestra en la Figura 50, elmA, colA y rowA son los componentes de la matriz en formato CSC, b es el vector por el que se desea multiplicar, c es el vector donde se almacenará el resultado, order es el orden de la matriz y partialsum es la memoria local que se utilizará para acumular sumas parciales, en el bloque for de las líneas 26 a 29 se calcula una suma parcial y en la línea 30 se almacena está suma parcial en memoria local, posteriormente se suman todos los elementos de la memoria local por reducción en paralelo (bloque for de las líneas 34 a 40), esta suma total viene a ser un elemento del vector resultado que es almacenada en la línea 42 por solo uno de los hilos.

```
void sparse_mat_vec_mul2(__global double* elmA,
              global int* colA, __global int* rowA,
                                                         global double* b,
2
3
              _global double* c, int order)
4
5
             //indice de grupo
             int bx = get_group_id(0);
6
8
             //indice local de hilo
9
             int tx = get local id(0);
10
11
             for (int col = bx; col < order; col += get_num_groups(0)) {</pre>
                     //posición del elemento diagonal
12
                     int dId = colA[col];
13
14
15
                     //número de elementos para una columna
                     int nEl = (col == order - 1 ? 1 : colA[col + 1] - dId);
16
17
                     //posición en rowL del primer elemento fuera de la
18
19
                     //diagonal (indice)
20
                     int baseId = dId - col;
21
22
                     //elemento en b, por el que se multiplica
23
                     double mul = b[col];
24
25
                     for (int x = tx; x < nEl - 1; x += get_local_size(0)) {</pre>
26
                              int rId = rowA[baseId + x];
                              double sum = mul * elmA[dId + x + 1];
27
28
                              atom_add_double(&(c[rId]), sum);
29
                     barrier(CLK GLOBAL MEM FENCE);
30
31
32
```

Figura 51. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul2

La función de dispositivo propuesta para obtener *Lb* se muestra en la Figura 51, esta función tiene los mismos argumentos que la función anterior sin la memoria local, el resultado de la operación se acumula directamente en el vector resultado de la función anterior, cada grupo se encarga de una columna y cada hilo de este grupo acumulará en el vector resultado el producto del elemento en la columna correspondiente al hilo y el elemento correspondiente en *b*, ya que es posible que dos hilos intenten escribir simultáneamente en la misma ubicación, se hace uso de la función atómica personalizada *atom_add_double*, análoga a la función incorporada de OpenCL *atom_add*, adaptada para trabajar sobre valores del tipo *double*, está función de dispositivo se muestra en la Figura 52, esta fue escrita adaptando la función CUDA *atomicAdd* para tipos *double* para dispositivos con capacidad computacional menor a 6.0, presentado en la sección *Programming Guide* del CUDA *toolkit documentation* (NVIDIA Corporation).

```
double atom_add_double(__global double* address, double val) {
1
2
             __global long* address_as_ull =
                     (__global long*)address;
4
             long old = *address_as_ull;
5
             long assumed;
7
             do {
                     assumed = old;
8
9
                     old = atom_cmpxchg(address_as_ull, assumed,
10
                             as_long(val + as_double(assumed)));
                     // Note: uses integer comparison to avoid hang in case
11
                     //of NaN (since NaN != NaN)
12
             } while (assumed != old);
13
14
             return as double(old);
15
16
```

Figura 52. Gradientes conjugados: función de dispositivo atom_add double

Finalmente la función de dispositivo propuesta para obtener *Db* se muestra en la Figura 53, al igual que en la función anterior, el resultado se acumula directamente en el vector resultado anterior, quedando después de usar esta última función, la multiplicación final matriz por vector. Esta función es básica y directa ya que se trata de multiplicar cada elemento de *b* por el elemento diagonal correspondiente y acumular este valor en el vector resultado, no es necesario introducir como argumento el vector de índices de fila *rowA*.

```
void sparse_mat_vec_mul3(__global double* elmA,
2
                                                         global double* c,
              _global int* colA, __global double* b,
3
             int order)
4
    {
5
6
             //indice global de hilo
             int tx = get_global_id(0);
8
9
             //número total de hilos
             int ntx = get_num_groups(0)*get_local_size(0);
10
11
12
             for (int x = tx; x < order; x += ntx) {
                     //ubicación del elemento diagonal en elmL
13
                     int diag = colA[x];
14
15
                     c[x] += b[x] * elmA[diag];
16
17
```

Figura 53. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul3

El kernel encargado de realizar la tarea completa de multiplicación se muestra en Figura 54, como puede observarse se hace uso de las funciones de dispositivo anteriores dependiendo del valor del argumento step, para el valor de step=0 calcula L^Tb , para step=1 calcula Lb acumulando el resultado en el mismo vector resultado de la operación anterior, y finalmente para step=2

calcula *Db* acumulando de la misma manera el resultado en el vector resultado anterior.

```
1
      _kernel void
2
    sparse_mat_vec_mul(__global double* elmA, __global int* colA,
             __global int* rowA, __global double* b, __global double* c,
3
              _local double* partialSum, int order, int step)
5
             //multiplicación de matriz por vector
6
             //transpose(L)*b, los elementos de c serán reemplazados
             //c[i]=valor
9
10
             if (step == 0)
11
                     sparse_mat_vec_mul1(elmA, colA, rowA, b, c, order,
                             partialSum);
12
13
             //L*b, los elementos de c se acumularán a partir de su
14
             //valor actual
             //c[i]+=valor o valores
15
16
             else if (step == 1)
                     sparse_mat_vec_mul2(elmA, colA, rowA, b, c, order);
17
18
19
             //D*b, los elementos de c se acumularán a partir de su
20
             //valor actual
21
             //c[i]+=valor
22
             else if (step == 2)
23
                     sparse_mat_vec_mul3(elmA, colA, b, c, order);
24
             else:
25
```

Figura 54. Gradientes conjugados: Kernel para multiplicar matriz en formato CSC por vector

La función que se exporta es similar a la presentada para matrices densas (Figura 29 pág. 83), cambiando la parte donde se calcula la multiplicación matriz por vector, por la función personalizada que se escribe usando el kernel presentado en esta sección.

4.1.7. Kernels sobre matrices en formato SKS

Se presentan los kernels elaborados para trabajar sobre matrices en formato SKS. Se ha definido una estructura *sparse_sks_ms* para matrices en este formato cuya definición en Julia se muestra en la Figura 55, los campos son: *elements* que guardan todos los elementos distintos de cero e *index* que guarda índices que indican la ubicación en *elements* del primer elemento distinto de cero para cada columna de la matriz original. De la misma manera definimos la estructura *af_sparse_sks_ms* (Figura 56) análogo a *sparse_sks_ms* que guardará la matriz dispersa en formato SKS con los campos del tipo *AFArray*, que indica que son objetos con memoria en el dispositivo.

```
struct sparse_sks_ms{T,N}
elements::Array{T,1}
index::Array{N,1}
end
```

Figura 55. Estructura para una matriz en formato SKS

```
struct af_sparse_sks_ms{T,N}
elements::AFArray{T,1}
index::AFArray{N,1}
end
```

Figura 56. Estructura para una matriz en formato SKS con memoria en el dispositivo

La Figura 57 muestra el código en Julia para convertir una matriz simétrica en formato denso en una matriz en formato SKS, para poder usar las funciones exportadas todavía es necesario tener la memoria en el dispositivo, para esto, hacemos uso de la función incorporada de ArrayFire *AFArray*, que hace una copia en la memoria del dispositivo de un objeto del tipo *Array*, se ha sobrecargado la función para poder hacer una copia de un objeto del tipo *sparse_sks_ms* en la memoria del dispositivo, en otras palabras convertir del tipo *sparse_sks_ms* a *af_sparse_sks_ms*, esta función se muestra en la Figura 58.

```
function sparse sks ms(A::Array(T,2)) where {T<:Real}
        orden=size(A)[1];
 3
        index=ones(Int32, orden);
4
        elements=Array(T,1)(UndefInitializer(),0);
 5
        for j in 1:orden
7
            index[j] += length(elements);
8
            skl=key;
9
            key=orden;
10
            while A[key,j]==0 && key>j && key>skl
11
                key-=1;
12
13
            append! (elements, A[j:key, j])
14
        end
15
        sparse sks ms(elements, index)
16
```

Figura 57. Conversión de una matriz densa a una matriz en formato SKS

```
function AFArray(A::sparse_sks_ms(T,N)) where {T<:Real,N<:Real}
Ael=AFArray(A.elements);
AIdx=AFArray(A.index);
af_sparse_sks_ms(Ael,AIdx)
end</pre>
```

Figura 58. Función para convertir un objeto del tipo sparse_sks_ms a af_sparse_sks_ms

4.1.7.1. Factorización de Cholesky

Para resolver un sistema de ecuaciones usando la factorización de Cholesky, son necesarias dos etapas, uno de la factorización misma y otro que consiste en resolver los sistemas triangulares generados después de la factorización.

El proceso es similar al usado sobre matrices densas, la ventaja de almacenar en este formato es que no es necesaria una factorización simbólica, cosa que si fue necesaria al almacenar en formato CSC. En otras palabras, la memoria necesaria para almacenar la factorización es la misma que utiliza la matriz almacenada en formato SKS, siempre y cuando se siga el procedimiento descrito en la sección 2.3.3. Formato de almacenamiento SKS pág. 50.

Sea A la matriz original en formato SKS y L la matriz que contendrá la factorización de Cholesky también en formato SKS, ya que no hay diferencia en la cantidad de memoria usada para la factorización y para la matriz misma en formato SKS, la factorización puede ser realizada en sitio, es decir, si realizamos la factorización en sitio, A contendrá la final en su campo *elements* los elementos de L.

Similar al proceso sobre matrices densas, la factorización fue dividida en tres partes, un paso en el que modificamos una columna *j* tal que:

$$A_{n,j} = \frac{A_{n,j}}{\sqrt{A_{j,j}}}, n > j$$

```
void chol_sparse_1_sks(__global double* elmL,
1
             _global int* idxL, int order, int step)
2
3
    {
4
             //es procedimiento solo deberá realizarse hasta
             //la columna order-2
5
             //Siendo order el orden de la matriz, la última
7
             //columna tendrá índice order-1, donde solo contiene
8
             //el elemento diagonal
9
             if (step < order - 1) {</pre>
10
                     //obteniendo la ubicación del primer elemento
11
                     //de la columna determinada por step
                     int pos = idxL[step];
12
13
14
                     //posición siguiente
                     int pos_nxt = idxL[step + 1];
15
16
17
                     //indice global de hilo
18
                     int tx = get_global_id(0);
19
20
                     //número total de hilos
                     int ntx = get_num_groups(0)*get_local_size(0);
21
22
23
                     //elemento diagonal. El primer
24
                     //elemento guardado en elmA para cada columna es el
25
                     //elemento diagonal
26
                     double _r = sqrt(elmL[pos]);
27
                     //cada hilo se encargará de modificar un elemento
28
29
                     //de la columna
30
                     for (int x = tx + pos + 1; x < pos nxt; x += ntx) {
31
                              elmL[x] /= _r;
32
33
```

Figura 59. Cholesky: función de dispositivo chol_sparse_1_sks

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 59, *elmL* e *idxL* son los campos que representan a la matriz en formato SKS, cada hilo se encargará de realizar la modificación de uno de los elementos de la columna, en la línea 12 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna (la columna es determinada con el argumento *step*), de manera similar en la línea 15 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna siguiente, los elementos en *elmL* con índices entre estos dos valores son los que se tienen que modificar. En las líneas 18 y 21 se obtiene el índice global de hilo y el número total de hilos respectivamente, en la línea 26 se obtiene la raíz cuadrada del elemento diagonal y finalmente en la línea 31 se modifican los elementos de la columna.

Otro paso consiste en modificar las columnas siguientes a la columna *j* modificada en el paso anterior tal que:

$A_{n,s} = A_{n,s} - A_{s,j} A_{n,j}, s > j, n \ge s$

```
3
4
            //este procedimiento solo deberá realizarse hasta
            //la columna order-2
6
            //Siendo order el orden de la matriz, la última
            //columna tendrá índice order-1, donde solo contiene
8
            //el elemento diagonal
9
            if (step < order - 1) {</pre>
10
                    //indice local de grupo
11
                    int bx = get_group_id(0);
12
13
                    //indice local de hilo
14
                    int tx = get local id(0);
15
16
                    //posición del primer elemento de la columna en elmL
17
                    int pos_elmL = idxL[step];
18
19
                    //número de elementos en elmL para la columna
                    //sin contar la diagonal
20
                    int nEl = idxL[step + 1] - pos_elmL - 1;
21
22
23
                    if (nEl > 0) {
24
                            for (int b = bx; b < nEl; b += get_num_groups(0)) {</pre>
25
                                    //indice de la columna a modificar
26
                                    int colId = step + b + 1;
27
28
                                    //primer factor
                                    double k1 = elmL[pos_elmL + b + 1];
29
30
31
                                    for (int x = b + tx; x < nEl;
32
                                             x += get_local_size(0)) {
33
34
                                             //segundo factor
                                             double k2 = elmL[pos_elmL + x + 1];
35
36
37
                                             //posición final en elmL
38
                                             int
                                                    posf = idxL[colId] + x - b;
39
40
                                             elmL[posf] -= k1 * k2;
41
42
43
44
```

Figura 60. Cholesky: función de dispositivo chol_sparse_2_sks

La función de dispositivo propuesta que realiza esta tarea se muestra en la Figura 60, en este caso, cada grupo y sus hilos se encargarán de modificar una columna. En las líneas 11 y 14 se obtiene el índice de grupo y el índice local de hilo respectivamente, la línea 17 obtiene la ubicación en *elmL* del primer elemento de la columna (determinada por *step*), en la línea 21 se calcula el número de elementos en *elmL* para la columna descontando el elemento diagonal, el bloque *if* de las líneas 23 a 43 índica que el bloque *for* anidado

(líneas 24 a 42) solo se utiliza cuando el número de elementos para la columna fuera de la diagonal es por lo menos uno.

El bloque *for* de las líneas 24 a 42 realizará las modificaciones de todas las columnas siguientes a la columna actual (argumento *step*).

Estos dos primeros pasos deben realizarse uno después de otro *n-1* veces.

Un tercer paso es necesario para terminar la factorización, este modifica los elementos de la diagonal tal que:

$$A_{j,j} = A_{j,j}$$

La función de dispositivo propuesta que realiza esta tarea se muestra en la Figura 61, todos los hilos de todos los grupos se encargan cada uno de modificar un elemento de la diagonal, este paso se realiza una sola vez luego de realizada las *n-1* iteraciones con los dos pasos vistos previamente.

Terminada la factorización se procede a solucionar los sistemas equivalentes, es decir si Ax = b es el sistema original, luego de la factorización quedará:

$$LL^Tx = b$$

Figura 61. Cholesky: función de dispositivo chol_sparse 3 sks

Haciendo $L^T x = y$ entonces L y = b. Primero se resuelve L y = b y finalmente $L^T x = y$, cuya solución en x es la solución del sistema A x = b.

Para resolver el sistema Ly = b, usaremos la eliminación típica de Gauss Jordan aplicada a una matriz triangular inferior, en cada paso los elementos de

cada columna fuera de la diagonal van volviéndose cero progresivamente, hasta llegar a una matriz diagonal, sin embargo, al resolver un sistema triangular, las operaciones sobre L pueden hacerse implícitamente y dejarla intacta para poder usarla al resolver el siguiente sistema $L^Tx = y$, con esto solo sería necesario modificar el vector de constantes b tal que para el paso j:

$$b_n = b_n - b_j \frac{L_{n,j}}{L_{j,j}}, n > j$$

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 62, el número de iteración y por la tanto el índice de la columna actual está dada por el valor del argumento step, esta función debe ser llamada n-l veces para terminar con la eliminación, luego será necesario dividir los elementos de b entre los elementos de la diagonal (función de dispositivo de la Figura 63), este vector será la solución del sistema Ly = b, que será usado como argumento al resolver el siguiente sistema triangular.

```
void chol_sparse_4_sks(__global double* elmL,
2
                global int* idxL, __global double* b, int order,
3
              int step)
4
5
              if (step < order - 1) {</pre>
6
                       //indice global de hilo
                       int tx = get_global_id(0);
8
9
                       //número total de hilos
10
                       int ntx = get_num_groups(0)*get_local_size(0);
11
12
                       //posición del elemento diagonal
                       int diag = idxL[step];
13
14
15
                       //número de elementos para una columna
                       int nEl = idxL[step + 1] - diag;
16
17
18
                       //deben haber elementos aparte del elemento diagonal
19
                       if (nEl > 1) {
20
                               //elemento diagonal y elemento correspondiente en b
double k = b[step] / elmL[diag];
21
22
23
                                for (int x = tx; x < nEl - 1; x += ntx) {
24
25
                                        int rowId = step + x + 1;
                                        b[rowId] -= k * elmL[diag + x + 1];
27
28
29
```

Figura 62. Cholesky: función de dispositivo chol_sparse_4_sks

```
void chol_sparse_5_sks(__global double* elmL,
             __global int* idxL, __global double* b, int order)
2
             //indice global de hilo
             int tx = get_global_id(0);
             //número total de hilos
9
             int ntx = get_num_groups(0)*get_local_size(0);
10
11
             for (int x = tx; x < order; x += ntx) {
                     //ubicación del elemento diagonal en elmL
                     int diag = idxL[x];
13
                     b[x] /= elmL[diag];
14
```

Figura 63. Cholesky: función de dispositivo chol_sparse_5_sks

Para resolver el sistema $L^T x = y$, se usará la misma matriz que almacena L asumiendo que está ordenada por filas y que se trata de una matriz triangular superior, con esto ya no sería necesaria la transposición explícita de la matriz. Se resolverá el sistema obteniendo el valor de las variables una a continuación de otra desde abajo, el valor de la variable en la posición i sería:

$$x_{i} = \frac{y_{i} - \sum_{j=i+1}^{n-1} x_{j} L_{i,j}}{L_{i,i}}$$

Se usa el mismo vector y para almacenar la solución. La función de dispositivo propuesta que resuelve el sistema $L^Tx = y$ se muestra en la Figura 64, cada llamada a esta función calculará una variable. Ya que el sistema se resuelve desde abajo, para un paso j se calculará el valor de la variable i tal que i=n-j-1, la fila de L usada para calcular esta variable tiene el mismo índice (se obtiene en la línea 15). Esta función de dispositivo usará memoria local para calcular la sumatoria $\sum_{j=i+1}^{n-1} x_j L_{i,j}$ por reducción en paralelo, el tamaño de esta memoria local será el mismo que el tamaño de un grupo de trabajo, por esta razón, en este caso solo se utilizará un grupo de trabajo al usar esta función (en la línea 30 se especifica que el grupo de trabajo utilizado es únicamente el con id cero).

En la línea 31 a 34 se calcula una suma parcial almacenándola en la variable *sum* (declarada en la línea 27), luego este valor es almacenado en memoria local (línea 35), los elementos de *partialsum* se suman en paralelo en el bloque *for* de las líneas 39 a 44, finalmente se calcula la variable correspondiente en la línea 46, esta tarea solo la realiza un hilo (en la línea 45 se especifica el id

del hilo que realiza la operación). Esta función de dispositivo debe ser llamada n veces para calcular las n variables.

```
void chol_sparse_6_sks(__global double* elmL,
             __global int* idxL, __global double* y, int order, int step, __local double* partialSum)
2
3
4
             //indice de bloque
             int bx = get_group_id(0);
6
8
             //indice global de hilo
9
             int tx = get_local_id(0);
10
11
             //número total de hilos
12
             int ntx = get_local_size(0);
13
14
             //indice de la fila de L implicada
             int fL = order - step - 1;
15
16
17
             //ubicación del elemento diagonal de la fila en elmL
             int diagId = idxL[fL];
18
19
20
             //número de elementos en la fila
             int nEl = (fL == order - 1 ? 1 : idxL[fL + 1] - diagId);
21
22
23
             //elemento diagonal
24
             double diag = elmL[diagId];
25
26
             //valor donde se guardará una suma parcial
27
             double sum = 0.0;
28
29
             //solo se usará un bloque
30
             if (bx == 0) {
31
                      for (int x = tx; x < nEl - 1; x += ntx) {
32
                               int rowId = fL + x + 1;
                               sum += elmL[diagId + x + 1] * y[rowId];
33
34
35
                      partialSum[tx] = sum;
36
37
                      //reducción en paralelo pare determinar la suma de
38
                      //todos los elementos de partialSum
                      for (int stride = ntx / 2; stride > 0; stride /= 2) {
39
40
                               barrier(CLK_LOCAL_MEM_FENCE);
41
                               if (tx < stride) {</pre>
42
                                       partialSum[tx] += partialSum[tx + stride];
43
44
                      if (tx == 0)
45
                              y[fL] = (y[fL] - partialSum[0]) / diag;
46
47
48
```

Figura 64. Cholesky: función de dispositivo chol_sparse_6_sks

El kernel principal se muestra en la Figura 65, como puede observarse el kernel llama a las funciones de dispositivo mencionadas anteriormente, dependiendo del valor del argumento *sstep*, los valores 0, 1 y 2 indican que nos encontramos en el proceso de factorización y los valores 3, 4 y 5 que estamos en el proceso de solución de los sistemas triangulares.

```
kernel void
2
    chol_sparse_sks(__global double* elmL,
3
              _global int* idxL, int size_elmL, int size_idxL,
              _global double* b, __local double* partialSum, int step,
4
             int sstep)
5
    {
             //la iteración principal se realiza en el host
8
9
             //se modificará una sola columna
10
             if (sstep == 0)
                     chol_sparse_1_sks(elmL, idxL, size_idxL, step);
11
12
             //se modificará las columnas siguientes a la columna
13
14
             //implicada
15
             else if (sstep == 1)
16
                     chol_sparse_2_sks(elmL, idxL, size_idxL, step);
17
             //una vez usada las funciones anteriores todavía es
18
             //necesario modificar la diagonal D[i]=sqrt(D[i]) para
19
20
             //concluir la factorización
21
             else if (sstep == 2)
22
                     chol_sparse_3_sks(elmL, idxL, size_idxL);
23
             //hasta aquí la factorización está concluida, la parte
24
25
             //triangular inferior de A contendrá la factorización
26
             //de Cholesky
27
             //desde aquí se procederá a solucionar el sistema
28
             //L*transpose(L)x=b
29
30
             //se resolverá el sistema Ly=b, y=transpose(L)*x
31
             else if (sstep == 3)
                     chol_sparse_4_sks(elmL, idxL, b, size_idxL,
32
33
                              step);
34
             else if (sstep == 4)
35
                     chol_sparse_5_sks(elmL, idxL, b, size_idxL);
36
             //se resolverá el sistema transpose(L)*x=y
37
38
             <mark>//"y" ya deb</mark>e haberse obten<mark>ido</mark> usando la función anterior
39
             else
40
                     chol_sparse_6_sks(elmL, idxL, b, size_idxL,
41
                              step, partialSum);
```

Figura 65. Cholesky: Kernel para trabajar sobre una matriz en formato SKS

El kernel es llamado n-1 veces para sstep igual a 0 y 1 (con step de 0 a n-2), solo una vez para sstep=2, n-1 veces para sstep=3 (con step de 0 a n-2), una vez para sstep=4 y finalmente n veces para sstep=5(con step de 0 a n-1).

4.1.7.2. Factorización LDL^T

El proceso para resolver un sistema usando la factorización LDL^T es similar al que utilizamos en la factorización de Cholesky.

Para resolver un sistema de ecuaciones usando la factorización LDL^T, son necesarias dos etapas, uno de la factorización misma y otro que consiste en resolver los sistemas triangulares generados después de la factorización.

El proceso es similar al usado sobre matrices densas, la ventaja de usar este formato de almacenamiento es que ya no es necesaria una factorización simbólica, cosa que si fue necesaria al almacenar en formato CSC, si para almacenar la matriz en formato SKS se usó el procedimiento descrito en la sección 2.3.3. (Formato de almacenamiento SKS pág. 50), el tamaño de memoria usado para el campo *elements* de la matriz, tiene el mismo tamaño que la memoria necesaria para guardar el campo *elements* de su factorización.

Sea A la matriz original en formato SKS y L la matriz que contendrá la factorización LDL^T también en formato SKS, como se dijo, el campo *elements* de la matriz tiene el mismo tamaño que el campo *elements* de su factorización, esto nos permite realizar una factorización en sitio sobre el mismo campo *elements* de A, si hacemos la factorización en sitio, el campo *elements* de A contendrá al final los elementos correspondientes del campo *elements* de L, es decir A se convertirá en L.

```
void ldlt_sparse_1_sks(__global double* elmL,
             global int* idxL, int order, int step)
3
4
             //es procedimiento solo deberá realizarse hasta
5
             //la columna order-2
6
             //Siendo order el orden de la matriz, la última
7
             //columna tendrá índice order-1, donde solo contiene
8
             //el elemento diagonal
9
             if (step < order - 1) {</pre>
10
                      //obteniendo la ubicación del primer elemento
11
                     //de la columna determinada por step
                     int pos = idxL[step];
12
13
14
                     //posición siguiente
15
                     int pos_nxt = idxL[step + 1];
16
17
                     //indice global de hilo
                     int tx = get_global_id(0);
18
19
20
                     //número total de hilos
21
                     int ntx = get_num_groups(0)*get_local_size(0);
22
23
                     //elemento diagonal. El primer
24
                     //elemento guardado en elmA para cada columna es el
25
                     //elemento diagonal
26
                     double _r = elmL[pos];
27
28
                     //cada hilo se encargará de modificar un elemento
29
                     //de la columna
30
                     for (int x = tx + pos + 1; x < pos_nxt; x += ntx) {</pre>
31
                              elmL[x] /= _r;
32
33
```

Figura 66. LDL^T: función de dispositivo *ldlt_sparse_1_sks*

El proceso de factorización fue dividido en dos partes, un paso en el que modificamos una columna *j* tal que:

$$A_{n,j} = \frac{A_{n,j}}{A_{j,j}}, n > j$$

```
void ldlt_sparse_2_sks(__global double* elmL,
              global int* idxL, int order, int step)
3
4
             //este procedimiento solo deberá realizarse hasta
             //la columna order-2
             //Siendo order el orden de la matriz, la última
6
             //columna tendrá índice order-1, donde solo contiene
             //el elemento diagonal
8
9
             if (step < order - 1) {
10
                      //indice local de grupo
11
                     int bx = get_group_id(0);
12
13
                      //indice local de hilo
14
                     int tx = get_local_id(0);
15
16
                      //posición del primer elemento de la columna de elmL
                     int pos_elmL = idxL[step];
17
18
19
                      //elemento diagonal
20
                     double diag = elmL[pos_elmL];
21
22
                      //número de elementos en elmL para la columna
                      //sin contar la diagonal
23
24
                     int nEl = idxL[step + 1] - pos_elmL - 1;
25
                     if (nEl > 0) {
26
                              for (int b = bx; b < nEl; b += get_num_groups(0)) {</pre>
27
28
                                      //indice de la columna a modificar
                                      int colId = step + b + 1;
29
30
31
                                      //primer factor
32
                                      double k1 = elmL[pos_elmL + b + 1];
33
34
                                      for (int x = b + tx; x < nEl;
35
                                               x += get_local_size(0)) {
36
37
                                               //segundo factor
                                               double k2 = elmL[pos_elmL + x + 1];
38
39
40
                                               //posición final en elmL
41
                                                       posf = idxL[colId] + x - b;
42
43
                                               elmL[posf] -= diag * k1 * k2;
                                      }
44
45
46
47
48
```

Figura 67. LDL^T: función de dispositivo *ldlt_sparse_2_sks*

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 66, cada hilo se encargará de modificar un elemento de la columna, en la línea 12 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna (la columna es determinada con el argumento *step*), de manera similar

en la línea 15 se obtiene el índice de ubicación en *elmL* del primer elemento de la columna siguiente, los elementos en *elmL* con índices entre estos dos valores son los que se tienen que modificar. En las líneas 18 y 21 se obtienen el índice global de hilo y el número total de hilos respectivamente, en la línea 26 se obtiene el elemento diagonal y finalmente en la línea 31 se calcula y se escribe el valor final en *elmL*.

Otro paso consiste en modificar las columnas siguientes a la columna *j* modificada en el paso anterior tal que:

$$A_{n,s} = A_{n,s} - A_{s,j} A_{n,j} A_{j,j}, s > j, n \ge s$$

La función de dispositivo propuesta que realiza esta tarea se muestra en la Figura 67, en este caso, cada grupo y sus hilos se encargarán de modificar una columna. En las líneas 11 y 14 se obtiene el índice de grupo y el índice local de hilo respectivamente, la línea 17 obtiene la ubicación en *elmL* del primer elemento de la columna actual (determinada por *step*), en la línea 20 se obtiene el elemento diagonal, en la línea 24 se calcula el número de elementos en *elmL* fuera de la diagonal que le corresponden a la columna, el bloque *if* de las líneas 26 a 46 índica que el bloque *for* anidado (líneas 27 a 45) solo se utiliza cuando el valor calculado en la línea 24 es por lo menos uno.

El bloque *for* de las líneas 27 a 45 realizará las modificaciones de todas las columnas siguientes a la columna actual (argumento *step*).

Estos dos pasos deben realizarse uno después de otro n-l veces para terminar con la factorización. Terminada la factorización, el campo *elements* de L contendrá los elementos de la factorización LDL^T , este incluye los elementos de la diagonal D almacenados como si fueran la diagonal de L, debe recordarse que L por si sola tiene elementos diagonales iguales a 1, esto es implícito y debe tenerse en cuenta al momento de resolver el sistema.

Terminada la factorización se procede a solucionar los sistemas equivalentes, es decir si Ax = b es el sistema original, luego de la factorización quedará:

$$LDL^{T}x = b$$

Haciendo $DL^Tx = y$ entonces Ly = b, es el primer sistema a resolverse para y, luego haciendo $L^Tx = z$ entonces Dz = y es el segundo sistema a resolver en z, que es simple ya que D es una matriz diagonal, finalmente se resuelve $L^Tx = z$, cuya solución en x es la solución del sistema Ax = b.

Para resolver el sistema Ly = b, usaremos la eliminación típica de Gauss Jordan aplicada a una matriz triangular inferior, en cada paso, los elementos de cada columna fuera de la diagonal van volviéndose cero progresivamente, hasta llegar a una matriz diagonal. Al resolver un sistema triangular las operaciones sobre L pueden hacer implícitamente dejándola intacta para resolver los siguientes sistemas, con esto, solo sería necesario modificar el vector de constantes b tal que para el paso j:

$$b_n = b_n - b_j L_{n,j}, n > j$$

```
1
            int step)
4
            if (step < order - 1) {</pre>
6
                    //indice global de hilo
                    int tx = get_global_id(0);
8
9
                    //número total de hilos
10
                    int ntx = get_num_groups(0)*get_local_size(0);
11
12
                    //posición del elemento diagonal
13
                    int diag = idxL[step];
14
                    //número de elementos para una columna
15
16
                    int nEl = idxL[step + 1] - diag;
17
18
                    //deben haber elementos aparte del elemento diagonal
19
                    if (nEl > 1) {
20
21
                            //elemento diagonal y elemento correspondiente en b
22
                            double diag_elm = elmL[diag];
23
                            double elmb = b[step];
24
                            for (int x = tx; x < nEl - 1; x += ntx) {
25
                                    int rowId = step + x + 1;
b[rowId] -= elmb * elmL[diag + x + 1];
26
27
28
29
```

Figura 68. LDL^T: función de dispositivo *ldlt_sparse_3_sks*

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 68, el índice de la columna actual es el mismo determinado por el argumento step, esta función debe ser llamada n-1 veces para terminar con la eliminación, al final b contendrá la solución del sistema Ly = b, que será usado como argumento al resolver el siguiente sistema.

Luego se procede a resolver el sistema Dz = y, ya que D es diagonal, la solución se reduce a realizar la operación:

$$b_j = \frac{b_j}{D_{i,j}}$$

La función de dispositivo propuesta que realiza esta operación se muestra en la Figura 69, esta función se ejecuta una sola vez.

```
void ldlt sparse 4 sks( global double* elmL,
              _global int* idxL, __global double* b, int order)
2
3
             //indice global de hilo
             int tx = get_global_id(0);
             //número total de hilos
8
             int ntx = get_num_groups(0)*get_local_size(0);
9
11
             for (int x = tx; x < order; x += ntx) {
                     //ubicación del elemento diagonal en elmL
12
                     int diag = idxL[x];
13
                     b[x] /= elmL[diag];
14
15
```

Figura 69. LDL^T: función de dispositivo ldlt_sparse_4_sks

Para resolver el sistema $L^T x = z$, se usará la misma matriz que almacena L, asumiendo que está ordenada por filas y que se trata de una matriz triangular superior, con esto, no será necesaria la transposición explícita de la matriz, se resolverá el sistema obteniendo el valor de las variables una a continuación de otra desde abajo, recordar que los elementos de la diagonal de L son 1 aunque no estén almacenados, así, el valor de la variable en la posición i sería:

$$x_i = z_i - \sum_{j=i+1}^{n-1} z_j L_{i,j}$$

Se usa el mismo vector z para almacenar la solución. La función de dispositivo propuesta que resuelve el sistema $L^Tx=z$ se muestra en la Figura 70, cada llamada a esta función calculará una variable. Ya que el sistema se resuelve desde abajo, para un paso j se calculará el valor de la variable i tal que i=n-j-1, la fila de L usada para calcular esta variable tiene el mismo índice (se obtiene en la línea 15). Esta función de dispositivo usará memoria local para calcular la sumatoria $\sum_{j=i+1}^{n-1} z_j L_{i,j}$ por reducción en paralelo, el tamaño de esta memoria local será el mismo que el tamaño de un grupo de trabajo, por esta razón, en este caso solo se utilizará un grupo de trabajo al usar esta función (en la línea 27 se especifica que el grupo de trabajo utilizado es únicamente el con id cero).

```
void ldlt_sparse_5_sks(__global double* elmL,
             __global int* idxL, __global double* y, int order, int step, __local double* partialSum)
3
5
             //indice de bloque
             int bx = get_group_id(0);
8
             //indice global de hilo
9
             int tx = get_local_id(0);
10
11
             //número total de hilos
             int ntx = get_local_size(0);
12
13
14
             //indice de la fila de L implicada
15
             int fL = order - step - 1;
16
17
             //ubicación del elemento diagonal de la fila en elmL
             int diagId = idxL[fL];
18
19
20
             //número de elementos en la fila
             int nEl = (fL == order - 1 ? 1 : idxL[fL + 1] - diagId);
21
22
             //valor donde se guardará una suma parcial
23
             double sum = 0.0;
24
25
26
             //solo se usará un bloque
27
             if (bx == 0) {
28
                      for (int x = tx; x < nEl - 1; x += ntx) {
                               int rowId = fL + x + 1;
29
                               sum += elmL[diagId + x + 1] * y[rowId];
30
31
32
                      partialSum[tx] = sum;
33
                      //reducción en paralelo pare determinar la suma de
34
35
                      //todos los elementos de partialSum
                      for (int stride = ntx / 2; stride > 0; stride /= 2) {
36
                               barrier(CLK_LOCAL_MEM_FENCE);
37
38
                               if (tx < stride) {</pre>
                                        partialSum[tx] += partialSum[tx + stride];
39
40
41
                      }
if (tx == 0)
42
                              y[fL] -= partialSum[0];
43
44
```

Figura 70. LDL^T: función de dispositivo *ldlt_sparse_5_sks*

En la línea 28 a 31 se calcula una suma parcial almacenándola en la variable *sum* (declarada en la línea 24), luego este valor es almacenado en memoria local (línea 32), los elementos de *partialsum* se suman en paralelo en el bloque *for* de las líneas 36 a 41, finalmente se escribe el valor final de la variable correspondiente en la línea 43, esta tarea solo la realiza un hilo (en la línea 42 se especifica el id del hilo que realiza la operación). Esta función de dispositivo debe ser llamada *n* veces para calcular las *n* variables.

```
kernel void
1
2
    ldlt_sparse_sks(__global double* elmL,
             _global int* idxL, int size_elmL, int size_idxL,
_global double* b, _local double* partialSum, int step,
3
4
5
6
    {
             //la iteración principal se realiza en el host
7
8
9
              //se modificará una sola columna
             if (sstep == 0)
10
11
                      ldlt_sparse_1_sks(elmL, idxL, size_idxL, step);
12
13
             //se modificará las columnas siguientes a la columna
14
             //implicada
15
             else if (sstep == 1)
                      ldlt_sparse_2_sks(elmL, idxL, size_idxL, step);
16
17
18
             //hasta aquí la factorización está concluida, la parte
             //triangular inferior de A contendrá la factorización
19
20
             //LDLt es decir L y D
21
             //desde aquí se procederá a solucionar el sistema
22
             //L*D*transpose(L)x=b
23
24
             //se resolverá el sistema Ly=b, y=D*transpose(L)*x
25
             else if (sstep == 2)
26
                      ldlt_sparse_3_sks(elmL, idxL, b, size_idxL,
27
                              step);
28
29
             //se resolverá el sistema Dz=y, z=transpose(L)*x
             else if (sstep == 3)
30
                      ldlt_sparse_4_sks(elmL, idxL, b, size_idxL);
31
32
33
             //se resolverá el sistema transpose(L)*x=z
34
             //"z" ya debe haberse obtenido usando la función anterior
35
36
                      ldlt_sparse_5_sks(elmL, idxL, b, size_idxL,
37
                               step, partialSum);
```

Figura 71. LDL^T: kernel para trabajar sobre una matriz en formato SKS

El kernel principal se muestra en la Figura 71, como puede observarse el kernel llama a las funciones de dispositivo mencionadas anteriormente, dependiendo del valor del argumento *sstep*, los valores 0 y 1 indican que nos encontramos en el proceso de factorización y los valores 2, 3 y 4 que estamos en el proceso de solución de los sistemas triangulares.

El kernel es llamado n-1 veces para sstep igual a 0 y 1 (con step de 0 a n-2), n-1 veces para sstep=2 (con step de 0 a n-2), una vez para sstep=3 y finalmente n veces para sstep=4(con step de 0 a n-1).

4.1.7.3. Gradientes conjugados

La ventaja de este método es que solo requiere de operaciones sencillas en cada iteración, casi todas estas operaciones puede hacerse haciendo uso de las funciones incorporadas en ArrayFire, la única operación que necesita una función de dispositivo es la que se encargará de multiplicar matriz por vector, con la matriz almacenada en formato SKS.

La matriz en su forma densa puede descomponerse en dos términos: una matriz estrictamente triangular inferior y una matriz triangular superior, entonces se cumplirá:

$$L + U = A$$

Siendo L la matriz estrictamente triangular inferior con los elementos de la parte estrictamente triangular inferior de A y U la matriz triangular superior con los elementos de la parte triangular superior de A, tomar en cuenta que ya que se trata de una matriz simétrica entonces la parte estrictamente superior de U contiene la transpuesta de L.

Multiplicando por un vector *b*:

$$Ab = Lb + Ub$$

La matriz en formato SKS almacena L y la diagonal de A, es decir, la transpuesta de U. Para multiplicar por U puede usarse la misma matriz A en formato SKS, asumiendo que está ordenada por filas y que se trata de una matriz triangular superior, así, la trasposición explícita de la matriz ya no es necesaria.

```
void sparse_mat_vec_mul1_sks(__global double* elmA,
1
              _global int* idxA, __global double* b,
_global double* c, int order,
2
                local double* partialSum)
4
5
6
              //indice de bloque
             int bx = get_group_id(0);
8
9
             //indice local de hilo
10
             int tx = get_local_id(0);
11
              for (int row = bx; row < order; row += get_num_groups(0)) {</pre>
12
13
                       //ubicación del elemento diagonal de la fila en elmA
14
                      int dId = idxA[row];
15
16
17
                      //número de elementos en la fila
                      int nEl = (row == order - 1 ? 1 : idxA[row + 1] - dId);
18
19
20
                       //valor donde se guardará una suma parcial
21
                      double sum = 0.0;
22
23
                      for (int x = tx; x < nEl; x += get_local_size(0)) {</pre>
24
                               int cId = row + x;
25
                               sum += elmA[dId + x] * b[cId];
26
27
                      partialSum[tx] = sum;
28
29
                      //reducción en paralelo pare determinar la suma de
30
                      //todos los elementos de partialSum
31
                      for (int stride = get_local_size(0) / 2;
32
                               stride > 0; stride /= 2) {
                               barrier(CLK_LOCAL_MEM_FENCE);
33
                               if (tx < stride) {</pre>
34
35
                                        partialSum[tx] += partialSum[tx + stride];
36
37
                      if (tx == 0)
38
39
                               c[row] = partialSum[0];
40
```

Figura 72. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul1_sks

La función de dispositivo propuesta para obtener *Ub* se muestra en la Figura 72, *elmA* e *idxA* son los componentes de la matriz en formato SKS, *b* es el vector por el que se desea multiplicar, *c* es el vector donde se almacenará el resultado, *order* es el orden de la matriz y *partialsum* es la memoria local que se utilizará para acumular sumas parciales, en el bloque *for* de las líneas 23 a 26 se calcula una suma parcial y en la línea 27 se almacena esta suma parcial en memoria local, posteriormente se suman todos los elementos de esta por reducción en paralelo (bloque *for* de las líneas 31 a 37), esta suma total viene a ser un elemento del vector resultado que es almacenada en la línea 39 por solo uno de los hilos.

```
void sparse_mat_vec_mul2_sks(__global double* elmA,
1
              _global int* idxA, __global double* b,
_global double* c, int order)
2
3
4
5
              //indice de grupo
6
             int bx = get_group_id(0);
7
8
              //indice local de hilo
9
             int tx = get_local_id(0);
10
              for (int col = bx; col < order; col += get_num_groups(0)) {
11
12
                      //posición del elemento diagonal
                      int dId = idxA[col];
13
14
                      //número de elementos para una columna
15
                      int nEl = (col == order - 1 ? 1 : idxA[col + 1] - dId);
16
17
                      //elemento en b, por el que se multiplica
18
19
                      double mul = b[col];
20
21
                      for (int x = tx; x < nEl - 1; x += get_local_size(0)) {</pre>
22
                               int rId = col + x + 1;
                               double sum = mul * elmA[dId + x + 1];
23
                               atom_add_double(&(c[rId]), sum);
24
25
26
                      barrier(CLK_GLOBAL_MEM_FENCE);
28
```

Figura 73. Gradientes conjugados: función de dispositivo sparse_mat_vec_mul2_sks

La función de dispositivo propuesta para obtener *Lb* se muestra en la Figura 73, esta función tiene los mismos argumentos que la función anterior sin la memoria local, el resultado de la operación se acumula directamente en el vector resultado de la función anterior, cada grupo se encarga de una columna y cada hilo de este grupo acumulará en el vector resultado el producto del elemento en la columna correspondiente al hilo y el elemento correspondiente en *b*, ya que es posible que dos hilos intenten escribir simultáneamente en la misma ubicación, se hace uso de la función atómica personalizada *atom_add_double*, análoga a la función incorporada de OpenCL *atom_add*, adaptada para trabajar sobre valores del tipo *double*, está función de dispositivo se muestra en la Figura 74, esta fue escrita adaptando la función CUDA *atomicAdd* para tipos *double* para dispositivos con capacidad computacional menor a 6.0, presentado en la sección *Programming Guide* del CUDA toolkit documentation (NVIDIA Corporation).

```
double atom_add_double(__global double* address, double val) {
1
2
              _global long* address_as_ull =
                     (__global long*)address;
3
             long old = *address_as_ull;
4
5
             long assumed;
             do {
                     assumed = old;
8
                     old = atom_cmpxchg(address_as_ull, assumed,
10
                             as_long(val + as_double(assumed)));
11
                     // Note: uses integer comparison to avoid hang in case
12
                     //of NaN (since NaN != NaN)
13
             } while (assumed != old);
14
15
             return as double(old);
16
```

Figura 74. Gradientes conjugados: función de dispositivo atom_add_double

El kernel encargado de realizar la tarea completa de multiplicación se muestra en Figura 75, como puede observarse se hace uso de las funciones de dispositivo anteriores dependiendo del valor del argumento step, para el valor de step=0 calcula Ub y para step=1 calcula Lb acumulando el resultado en el mismo vector resultado de la operación anterior.

```
kernel void
     sparse_mat_vec_mul_sks(__global double* elmA,
3
              _global int* idxA, __global double* b,
               global double* c, _
4
                                   local double* partialSum,
             int order, int step)
5
6
    {
             //multiplicación de matriz por vector
8
9
             //U*b, los elementos de c serán reemplazados
10
             //c[i]=valor
11
             if (step == 0)
12
                     sparse_mat_vec_mul1_sks(elmA, idxA, b, c,
                              order, partialSum);
13
14
             //L*b, los elementos de c se acumularán a partir de su
15
             //valor actual
             //c[i]+=valor o valores
16
17
             else if (step == 1)
                     sparse_mat_vec_mul2_sks(elmA, idxA, b, c,
18
19
                              order);
20
             else;
```

Figura 75. Gradientes conjugados: kernel para multiplicar matriz en formato SKS por vector

La función que se exporta es similar a la presentada para matrices densas (Figura 29 pág. 83), cambiando la parte donde se calcula la multiplicación matriz por vector, por la función personalizada que usa el kernel presentado en esta sección.

4.1.8. Medición del tiempo y memoria utilizada

En general, el tiempo de ejecución de un kernel consiste de tres componentes: tiempo de preparación del kernel, tiempo de cómputo y tiempo de sincronización, el cuál es representado como:

$$T = \sum_{i=1}^{M} (t_0^{(i)} + t_C^{(i)} + t_S^{(i)})$$

Donde M es el número de lanzamientos del kernel, $t_O^{(i)}$ es el tiempo de preparación del kernel, $t_C^{(i)}$ es el tiempo de cómputo sobre una GPU y $t_S^{(i)}$ es el tiempo de sincronización para el i-ésimo lanzamiento del kernel, cada uno de estos componentes es influenciado por algunos factores, por ejemplo, el tiempo de preparación del kernel depende de la tasa de transferencia de datos del procesador central al dispositivo así como del tamaño del código y los parámetros del kernel. Para el tiempo de cómputo, este es afectado por los métodos de accesos a memoria, la organización de los hilos (número de hilos por grupo y número de grupos) en el kernel, etc. (Shucai Xiao, 2009).

El tiempo de lanzamiento del kernel es muy pequeño en comparación a los otros dos (Shucai Xiao, 2009), los resultados de tiempo de ejecución en el presente trabajo, incluyen estos tres componentes.

Respecto a la memoria utilizada, en matrices densas, ya que no se utiliza ningún formato de almacenamiento que ignore los ceros y considerando que en los métodos con factorización se realizan dichas factorizaciones en sitio, entonces, la cantidad de memoria en MB (Mega bytes) utilizada estará dada por (recordar que hacemos uso de precisión doble):

$$M_{MB} = \frac{8n(n+1)}{1x10^6}$$

Siendo n el orden de la matriz, la expresión anterior considera la matriz de rigidez (matriz de coeficientes del sistema) y el vector de fuerzas (vector de constantes).

Para matrices almacenadas en formato CSC se suma la memoria utilizada por los tres componentes de la matriz en este formato y la memoria utilizada por el vector de fuerzas, considerando que los vectores que almacenan los índices lo hacen como enteros de 32 bits (4 bytes), la memoria total utilizada sería:

$$M_{MB} = \frac{4(3l_e + 2n)}{1x10^6}$$

Siendo l_e el tamaño de los campos *elements* de la matriz en formato CSC y n es el orden de la matriz (o el tamaño del vector de fuerzas). También se está considerando que la matriz en formato CSC original es sustituida por la matriz L (también en formato CSC), en los métodos que usan factorización.

Finalmente, Para matrices almacenadas en formato SKS se suma la memoria utilizada por los dos componentes de la matriz en este formato y la memoria utilizada por el vector de fuerzas, considerando que el vector que almacenan los índices lo hace como enteros de 32 bits (4 bytes), la memoria total utilizada sería:

$$M_{MB} = \frac{4(2l_e + 3n)}{1x10^6}$$

Siendo l_e el tamaño del campo *elements* de la matriz en formato SKS y n el orden de la matriz (o el tamaño del vector de fuerzas). También se está considerando una factorización en sitio para los métodos que usan factorización.

4.1.9. Características del dispositivo

Las pruebas fueron realizadas en una **GeForce GT 630M**, cuyas características principales son: número de multiprocesadores 2, número de núcleos CUDA 96, con capacidad de cálculo (compute capability) 2.1, cantidad total de memoria global 2048 MB, cantidad máxima de hilos (o work ítems) en un grupo de trabajo 1024, cantidad total de memoria local por grupo de trabajo 48 KB. En general, al usarse otros dispositivos, el rendimiento mejorará según aumenten estos valores.

4.1.10. Validación de resultados

El vector de desplazamientos (solución del sistema de ecuaciones lineales) arrojado por las funciones, se validaron comparándolos con el vector de

desplazamientos arrojados por el solver incorporado en OpenFEM, el error máximo encontrado para la diferencia entre elementos correspondientes no excedió de $1x10^{-9}$, y de todo el vector (medido mediante la obtención de la norma de la diferencia de vectores) no excedió de $1x10^{-5}$.

4.1.11. Modelo 1: Viga en voladizo

Se ha comenzado analizando una viga simple de concreto en voladizo de 210 Kg/cm², con una sección de 0.30 x 0.50 m, longitud de 3.00 m y sometida a una carga distribuida de 2 tn/m (19613.3 N/m). Analizando la viga como un sólido, se hace el mallado correspondiente del elemento para analizarla usando el método de los elementos finitos, en la Figura 76 se muestra el mallado de la viga, se ha hecho el mallado dividiendo la sección y el largo de la viga en partes variables con tal de obtener una cantidad de elementos más densa, el gráfico es generado usando el complemento **Medit** de OpenFEM.

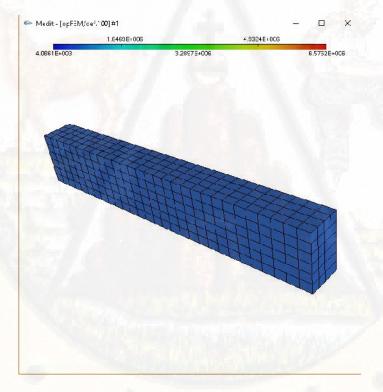


Figura 76. malla de elementos finitos: viga en voladizo

Para este modelo se han generado sistemas de ecuaciones con matrices que van del orden de 1080 a 10044, se ha considerado la matriz tal y como la devuelve la

herramienta de elementos finitos (OpenFEM), es decir sin ningún pivote previo, para analizar el comportamiento de las funciones y formatos de almacenamiento.

4.1.11.1. Memoria utilizada

Los resultados de memoria utilizada para los técnicas consideradas se muestran en la Tabla 10, puede notarse que para el formato denso (sin comprimir), la memoria es la misma para los cuatro métodos (resultado trivial), sin embargo, para el formato CSC existen diferencias entre la memoria utilizada por los métodos con factorización y los métodos sin factorización (método de los gradientes conjugados), algo que se esperaba, ya que los métodos con factorización requieren mayor cantidad de memoria que la original para almacenar su factorización, es por eso, que antes de realizar la factorización en este formato de almacenamiento, es necesaria una factorización simbólica previa. Lo anterior no pasa al almacenar en formato SKS ya que el almacenamiento en este formato se hace de tal manera que la cantidad de memoria necesaria para almacenar la matriz original es suficiente también para almacenar su factorización, es una ventaja del formato SKS sobre el CSC.

Tabla 10. Modelo 1: memoria utilizada en MB

Memoria utilizada (MB)												
		Denso			CSC				SKS			
orden	Gauss Jordan	('holesky LDL'		Cholesky	$\mathbf{LDL^{T}}$	Gradientes conjugados	Cholesky	LDLT	Gradientes conjugados			
1080	9.34	9.34	9.34	9.34	3.36	3.36	0.35	2.72	2.72	2.72		
2016	32.53	32.53	32.53	32.53	10.15	10.15	0.72	8.33	8.33	8.33		
3045	74.20	74.20	74.20	74.20	18.74	18.74	1.15	13.88	13.88	13.88		
4032	130.09	130.09	130.09	130.09	28.47	28.47	1.59	20.67	20.67	20.67		
5040	203.25	203.25	203.25	203.25	39.72	39.72	2.03	29.82	29.82	29.82		
6048	292.67	292.67	292.67	292.67	59.23	59.23	2.46	45.70	45.70	45.70		
7161	410.30	410.30	410.30	410.30	82.28	82.28	2.96	63.74	63.74	63.74		
8184	535.89	535.89	535.89	535.89	91.57	91.57	3.42	65.02	65.02	65.02		
9216	679.55	679.55	679.55	679.55	115.93	115.93	3.88	82.74	82.74	82.74		
10044	807.14	807.14	807.14	807.14	126.58	126.58	4.27	97.26	97.26	97.26		

Un aspecto importante a resaltar es la diferencia considerable en el uso de memoria entre el formato CSC sin factorizar y el mismo formato con factorización, en estas condiciones, para este modelo estructural, podemos apreciar que el método de los gradientes conjugados aplicada sobre la matriz en formato CSC es la que utiliza la memoria más eficientemente.

Esta diferencia puede notarse de forma más evidente en el siguiente gráfico donde *CSC_fact* es la memoria que utiliza los métodos con factorización (Cholesky y LDL^T). La memoria que utiliza el formato SKS y el formato CSC con factorización no se diferencian considerablemente.

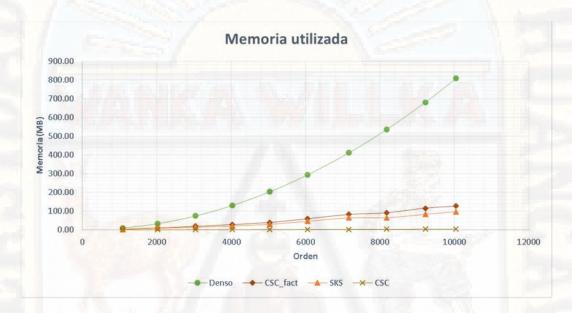


Gráfico 1: Modelo 1: memoria utilizada en MB

4.1.11.2. Tiempo de ejecución

La Tabla 11 muestra el tiempo de ejecución para los diferentes formatos de almacenamiento y métodos de solución. Puede observarse que ambos métodos de factorización aplicadas a la matriz en formato SKS arrojan los mejores resultados, los tiempos para los demás métodos y formatos de almacenamiento se alejan considerablemente de estos resultados.

El Gráfico 2 muestra los tiempos de ejecución en segundos para todas las técnicas, el Gráfico 3, muestra los tiempos de ejecución para las mejores cinco técnicas. Puede notarse las diferencias significativas de los resultados en tiempo de ejecución de los métodos con factorización sobre la matriz en

formato SKS con los demás tiempos arrojados, algo que se notó también analizando el otro modelo estructural.

Tabla 11. Modelo 1: tiempo de ejecución en segundos

	Tiempo de ejecución (s)									
NA	10	De	Denso CSC SKS				CSC			
orden	Gauss Jordan	Cholesky	LDLT	Gradientes conjugados	Cholesky	LDLT	Gradientes conjugados	Cholesky	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$	Gradientes conjugados
1080	0.61	0.39	0.88	0.36	1.90	1.89	3.72	0.30	0.25	4.77
2016	2.70	1.53	1.83	1.47	8.51	8.49	6.57	0.71	0.66	10.32
3045	9.32	4.33	4.33	4.25	16.94	17.04	11.05	1.15	1.07	18.80
4032	29.10	10.32	10.45	9.99	30.29	29.99	20.00	1.72	1.62	28.77
5040	42.67	17.48	17.48	11.89	45.87	45.77	24.89	2.53	2.52	44.58
6048	71.26	30.18	30.25	15.91	92.85	93.24	29.38	4.46	4.30	61.87
7161	134.42	55.91	56.01	25.90	159.43	160.01	34.12	6.41	6.28	83.83
8184	197.91	77.18	78.16	36.09	156.71	157.16	41.44	6.34	6.15	106.85
9216	290.66	115.33	115.46	50.10	238.10	238.93	49.01	9.18	8.99	123.64
10044	369.93	146.98	149.13	64.05	307.67	308.91	61.15	10.00	9.33	160.84



Gráfico 2: Modelo 1: tiempo de ejecución en segundos



Gráfico 3: Modelo 1: tiempo de ejecución en segundos

4.1.12. Modelo 2: Pilar de puente

Se ha modelado pilar de puente con columnas circulares, viga de sección variable y zapatas cuadradas (Ver Figura 77), las columnas tienen un diámetro de 1 m y altura de 10 m, la viga tiene un ancho de 1.04 m y un peralte variable de 1 m en el centro de la luz y final de las alas, y de 1.5 m en los extremos; la distancia entre ejes de columna es de 10.04 m, la longitud total de la viga es de 15.08 m, el material es de concreto de 280 Kg/cm² y la carga distribuida sobre la viga de 50 tn/m (490332.5 N/m), se hizo el mallado correspondiente del elemento para analizarla usando el método de los elementos finitos, en la Figura 77 se muestra el mallado de la estructura, se ha hecho el mallado dividiendo los elementos en partes cada vez más grandes con tal de obtener una cantidad de nodos más grande, el gráfico es generado usando el complemento **Medit** de OpenFEM.

Para este modelo se han generado sistemas de ecuaciones con matrices que van del orden de 1572 hasta 10449 para todos los métodos y hasta 205602 para pruebas adicionales con los más sobresalientes, se ha considerado la matriz con pivote previo de filas y columnas, usando el reordenamiento de Cuthill–McKee, con esto buscamos concentrar elementos distintos de cero más cerca de la diagonal y por tanto mejorar el rendimiento y optimizar aún más el uso de memoria.

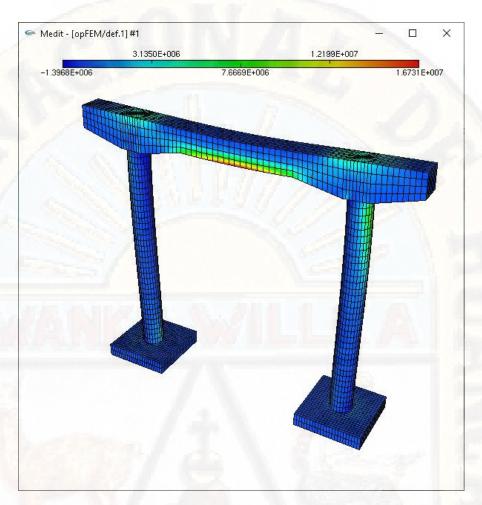


Figura 77. Malla de elementos finitos: soporte central de puente

4.1.12.1. Memoria utilizada

Los resultados de memoria utilizada para todos los métodos considerados se muestran en la Tabla 12, puede notarse que para el formato denso (sin comprimir), la memoria es la misma para los cuatro métodos, sin embargo, como ya se vio en los resultados de memoria para el modelo anterior, para el formato CSC existen diferencias entre la memoria utilizada por los métodos con factorización y los métodos sin factorización (método de los gradientes conjugados), algo que se esperaba, ya que los métodos con factorización requieren mayor cantidad de memoria que la original para almacenar la matriz factorizada, es por eso, que antes de realizar la factorización en este formato de almacenamiento, es necesaria una factorización simbólica previa. Lo anterior no pasa al almacenar en formato SKS ya que el almacenamiento en este formato

se hace de tal manera que la cantidad de memoria necesaria para almacenar la matriz original es suficiente también para almacenar su factorización, es una ventaja del formato SKS sobre el CSC.

Tabla 12. Modelo 2: memoria utilizada en MB

8	- 3	1/1	4	Me	moria utiliz (MB)	zada	11/11	111	M	WAT
		De	nso		CSC			SKS		
Orde n	Gauss Jordan	Cholesk y	LDL^{T}	Gradiente s conjugado s	Cholesk y	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$	Gradiente s conjugado s	Cholesk	$\mathbf{LDL}^{\mathrm{T}}$	Gradiente s conjugad os
1572	19.78	19.78	19.78	19.78	1.19	1.19	0.33	1.43	1.43	1.43
3534	99.94	99.94	99.94	99.94	3.21	3.21	1.01	3.26	3.26	3.26
5409	234.10	234.10	234.10	234.10	7.54	7.54	1.81	7.80	7.80	7.80
8331	555.31	555.31	555.31	555.31	13.29	13.29	2.51	12.94	12.94	12.94
10449	873.54	873.54	873.54	873.54	24.55	24.55	3.51	24.89	24.89	24.89
15129	1831.21	1831.21	1831.21	1831.21	30.86	30.86	5.16	30.48	30.48	30.48
21987	3867.60	3867.60	3867.60	3867.60	67.99	67.99	7.54	65.13	65.13	65.13
28773	6623.31	6623.31	6623.31	6623.31	125.73	125.73	10.44	118.88	118.88	118.88
39051	12200.16	12200.16	12200.16	12200.16	246.77	246.77	14.29	226.55	226.55	226.55
51156	20935.90	20935.90	20935.90	20935.90	322.65	322.65	19.78	290.52	290.52	290.52

Un aspecto importante a resaltar es la diferencia considerable en el uso de memoria entre el formato CSC sin factorizar y el mismo formato factorizado, en estas condiciones, para este modelo estructural, podemos apreciar que el método de los gradientes conjugados aplicada sobre la matriz en formato CSC es la que utiliza la memoria más eficientemente, algo que se notó también en el modelo estructural anterior.

En el Gráfico 4 puede notarse la diferencia significativa que existe al usar la matriz densa (sin comprimir) y los otros formatos de almacenamiento, en este gráfico, aparentemente, no existe diferencia significativa entre los formatos CSC (con factorización o sin factorización) y SKS, sin embargo, al quitar la línea de la matriz densa pueden apreciarse las diferencias (ver Gráfico 5). La memoria que utiliza el formato SKS y el formato CSC con factorización no se diferencian considerablemente, como se vio también en el modelo estructural anterior.



Gráfico 4: Modelo 2: memoria utilizada en MB

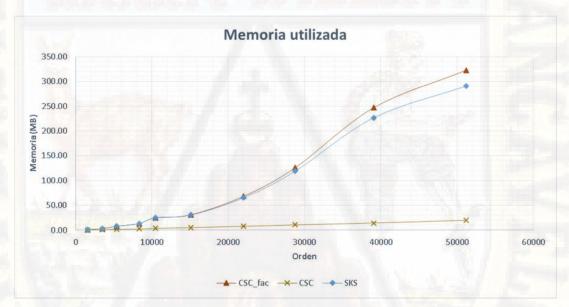


Gráfico 5: Modelo 2: memoria utilizada en MB

4.1.12.2. Tiempo de ejecución

La Tabla 13 muestra el tiempo de ejecución para los diferentes formatos de almacenamiento y métodos de solución. Puede observarse que ambos métodos de factorización aplicadas a la matriz en formato SKS arrojan los mejores resultados, los tiempos para los demás métodos y formatos de almacenamiento se alejan considerablemente de estos. Hay aspectos bastante importantes a

resaltar con el modelo anterior, y es que, aunque las técnicas que arrojan los mejores resultados se mantienen, para ordenes similares, el tiempo de ejecución para este modelo es mucho menor (10.00 s para el orden de 10044 para el modelo 1 y 2.61 s para el orden de 10449 para el modelo 2), esta diferencia es muy probablemente debido al reordenamiento previo que se hizo a la matriz para este modelo.

Otra observación importante es el deterioro del método de los gradientes conjugados con la matriz en formato CSC, en comparación a los resultados arrojados para el modelo estructural anterior y, por el contrario, la mejora significativa que tuvieron los métodos con factorización.

Tabla 13. Modelo 2: tiempo de ejecución en segundos

	Tiempo de ejecución (s)									
	Denso				CSC			SKS		
orden	Gauss Jordan	Cholesky	$\mathbf{LDL}^{\mathrm{T}}$	Gradientes conjugados	Cholesky	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$	Gradientes conjugados	Cholesky	LDL ^T	Gradientes conjugados
1572	1.47	0.88	0.81	4.13	0.47	0.42	13.68	0.25	0.20	14.45
3534	15.74	6.52	6.43	18.11	1.20	1.12	34.75	0.52	0.43	39.98
5409	57.01	21.64	21.56	38.77	3.28	3.14	61.69	1.00	0.96	73.68
8331	216.29	78.58	79.38	90.77	5.85	5.61	130.11	1.61	1.35	103.33
10449	300.88	125.99	127.66	139.14	14.61	14.44	281.13	2.61	2.27	175.44

Tabla 14. Modelo 2: tiempo de ejecución en segundos

	Tiempo de ejecución (s)						
	CS	SC	SKS				
Orden	Cholesky	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$	Cholesky	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$			
1572	0.47	0.42	0.25	0.20			
3534	1.20	1.12	0.52	0.43			
5409	3.28	3.14	1.00	0.96			
8331	5.85	5.61	1.61	1.35			
10449	14.61	14.44	2.61	2.27			
15129	17.37	16.90	3.39	2.94			
21987	56.19	55.72	6.69	6.06			
28773	176.26	176.27	12.47	11.73			

La razón por la que solo se probaron todas las técnicas hasta el orden de 10449 en el caso de los métodos sobre la matriz densa, es el límite de memoria del dispositivo (Ver sección 4.1.9. Características del dispositivo pág. 131), los resultados de memoria de la Tabla 12 (pág. 138) fueron calculadas manualmente, observar que el límite de dispositivo se alcanza apenas en el orden 15129, y que para el mayor orden en ésta tabla se necesitaría un dispositivo de más de 20 000 MB (20 GB) si se quisiera trabajar con la matriz sin comprimir.

Descartando las técnicas con los resultados más desfavorables, se hicieron pruebas adicionales (hasta el orden de 28773), para los cuatro técnicas con mejores tiempos, la Tabla 14 y el Gráfico 6 muestras los resultados. Puede observarse la similitud entre los tiempos que arrojan los métodos con factorización para un mismo formato de almacenamiento.



Gráfico 6: Modelo 2: tiempo de ejecución en segundos

4.1.13. Balance entre memoria y tiempo de ejecución

Para el formato CSC, la memoria utilizada es mucho menor al usar el método de los gradientes conjugados (ver Tabla 10 y Tabla 12, pág. 133 y 138 respectivamente) ya que la matriz original en formato CSC no requiere ningún tipo de factorización, sin embargo, la contraparte es que al usar el método de los gradientes conjugados con la matriz en formato CSC los tiempos de ejecución

resultantes son mucho mayores (ver Tabla 11 y Tabla 13, pág. 135 y 140 respectivamente).

Moviéndonos hacia los métodos que usan, en segundo lugar, más eficientemente la memoria, llegamos a los métodos que usan la matriz en formato SKS, podemos observar que para los métodos con factorización (Cholesky y LDL^T) en este formato, los tiempos son bastante mejores, no es difícil concluir entonces, que estos métodos en éste formato de almacenamiento, tienen el mejor balance de memoria y tiempo de ejecución. De aquí en adelante, para el modelo 2, se harán pruebas adicionales usando este formato de almacenamiento y estos métodos, ya que para órdenes más grandes, la tendencia en tiempo de ejecución para el formato CSC es obtener valores mucho más grandes, superior a los 170 segundos, lo cual es exagerado (ver Tabla 14, pág. 140).

Tabla 15. Modelo 2: memoria utilizada en MB para el formato SKS

Memoria utilizada (MB)						
Orden	Longitud del campo elements	Memoria utilizada <i>Mu</i>	Men_cia en formate denso	Mu/Ma		
1572	176955	1.43	19.78	0.073		
3534	402647	3.26	99.94	0.033		
5409	966955	7.80	234.10	0.033		
8331	1605480	12.94	555.31	0.023		
10449	3096149	24.89	873.54	0.028		
15129	3787050	30.48	1831.21	0.017		
21987	8107764	65.13	3867.60	0.017		
28773	14816912	118.88	6623.31	0.018		
39051	28260478	226.55	12200.16	0.019		
51156	36238341	290.52	20935.90	0.014		
64836	59114778	473.70	33630.17	0.014		
72570	89764866	718.99	42131.82	0.017		
89733	108759760	871.15	64416.81	0.014		
102630	121760423	975.31	84264.16	0.012		
117390	124626462	998.42	110244.24	0.009		
142422	185245727	1483.67	162273.35	0.009		

Se han realizado pruebas adicionales sobre órdenes de la matriz de rigidez más grandes, llegando al orden 142422, afinando más la malla de elementos finitos. Se ha utilizado el formato de almacenamiento SKS y los métodos de la factorización de Cholesky y LDL^T. La memoria utilizada por ambos métodos es la misma, los

resultados se muestran en la Tabla 15, donde puede observarse que a medida que el orden aumenta, el índice de dispersión (relación M_u/M_d) disminuye, para el orden máximo se alcanzó un valor de 0.009 lo que significa que la memoria utilizada es de un 0.9% de la memoria que se utilizaría con la matriz en formato denso (ver Gráfico 7).

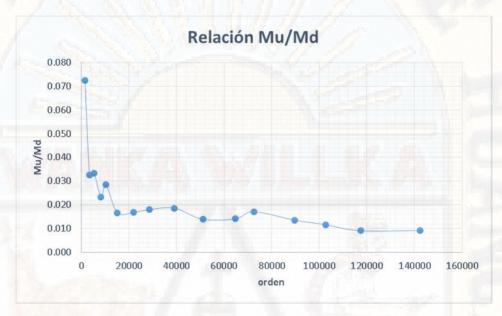


Gráfico 7: índice de dispersión para el formato SKS

Tabla 16. Modelo 2: tiempo de ejecución en matrices en formato SKS

Tiempo de ejecución (s)							
orden	Cholesky	LDL^{T}					
1572	0.25	0.20					
3534	0.52	0.43					
5409	1.00	0.96					
8331	1.61	1.35					
10449	2.61	2.27					
15129	3.39	2.94					
21987	6.69	6.06					
28773	12.47	11.73					
39051	24.28	23.30					
51156	29.57	28.47					
64836	52.98	50.48					
72570	94.96	93.66					
89733	121.05	119.09					
102630	132.26	129.67					
117390	135.51	132.64					
142422	225.44	223.96					

Los resultados de tiempo de ejecución se muestran en la Tabla 16 y Gráfico 8, puede observarse que los tiempos arrojados por ambos métodos son muy parecidos, el máximo tiempo obtenido para el orden 142422 es de 225.44 segundos para la factorización de Cholesky y de 223.96 segundos para la factorización LDL^T.



Gráfico 8: tiempo de ejecución en matrices en formato SKS

4.1.14. Pruebas adicionales

4.1.14.1. GT 630M vs GTX 1070

Para mostrar el impacto que tiene la utilización de los mismas técnicas al resolver el mismo sistema de ecuaciones lineales pero en otro dispositivo (otra GPU), se han realizado pruebas adicionales en otro dispositivo con mejores características, las características de ambos dispositivos (arrojados por la aplicación ocldevicequery.exe de OpenCL) se muestran en la Tabla 17, la GT 630M es el dispositivo que se usó en esta investigación; en esta tabla puede observarse las ventajas que tiene la GTX 1070 sobre la GT 630M, la diferencia más notable está en el número de núcleos CUDA para ambos dispositivos. Estas diferencias se reflejaron también significativamente al realizar las

pruebas, usando las mismas funciones (y por lo tanto los mismos kernels) y resolviendo los mismos sistemas de ecuaciones lineales para el modelo 2 (pilar de puente), estos resultados se muestran en la Tabla 18 y el Gráfico 9.

Tabla 17. Diferencias en las características de los dispositivos

GT 630M	GTX 1070
2	16
96	4294967280
2.1	6.1
2048 MB	8192 MB
1024	1024
48 KB	48 KB
	2 96 2.1 2048 MB

Tabla 18. Modelo 2: tiempo de ejecución en diferentes dispositivos

Tiempo de ejecución (s), para matrices en formato SKS usando diferentes dispositivos					
	GT 630M		GTX 1	1070	
orden	Cholesky	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$	Cholesky	$\mathbf{L}\mathbf{D}\mathbf{L}^{\mathrm{T}}$	
1572	0.25	0.20	0.06	0.06	
3534	0.52	0.43	0.12	0.12	
5409	1.00	0.96	0.2	0.19	
8331	1.61	1.35	0.29	0.28	
10449	2.61	2.27	0.37	0.39	
15129	3.39	2.94	0.52	0.51	
21987	6.69	6.06	0.91	0.9	
28773	12.47	11.73	1.5	1.48	
39051	24.28	23.30	2.86	2.88	
51156	29.57	28.47	3.41	3.4	
64836	52.98	50.48	5.95	6.02	
72570	94.96	93.66	10.92	11.11	
89733	121.05	119.09	12.34	12.58	
102630	132.26	129.67	13.34	13.57	
117390	135.51	132.64	13.85	14.23	
142422	225.44	223.96	23.11	23.86	

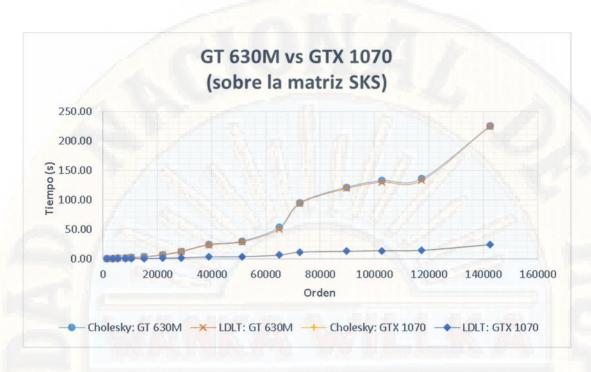


Gráfico 9: tiempo de ejecución en diferentes dispositivos

4.1.14.2. Comparación con el solver incorporado de OpenFEM

Como referencia, también se ha hecho la comparación con el solver incorporado de OpenFEM que usa la factorización LDL^T sobre la matriz en formato COO, los resultados no son tan alentadores para la GT 630M (ver Gráfico 10), sin embargo, a comparación con la GTX 1070 se observó una mejora significativa (ver Gráfico 11), para esta última prueba se llegó al orden 205609.

Estos últimos resultados son solo referenciales, ya que, incluso sabiendo el método y formato de almacenamiento que usa el solver de OpenFEM, se desconoce el tipo de dispositivo que usa para resolver el sistema.

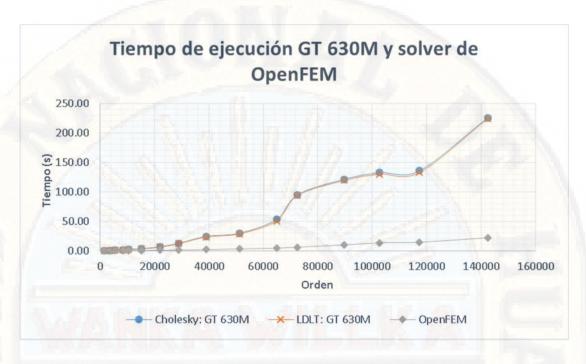


Gráfico 10: tiempo de ejecución GT 630M y solver de OpenFEM

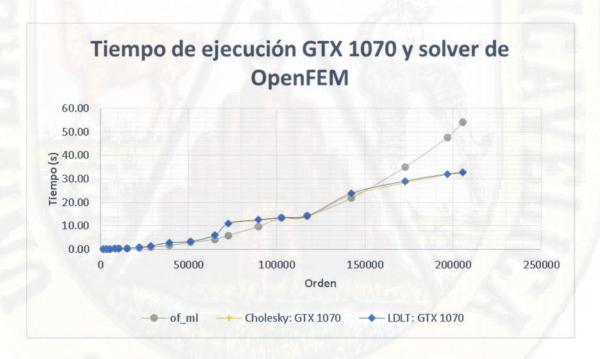


Gráfico 11: tiempo de ejecución GTX 1070 y solver de OpenFEM

4.1. Discusión de resultados

Según los visto en la sección anterior, entre las técnicas usadas en ésta tesis, al aplicar los métodos de la factorización de Cholesky y la factorización LDL^T en matrices en

formato SKS se obtienen los mejores balances de memoria y tiempo de ejecución; uno de los principales motivos por el que la aplicación de estos métodos sobre matrices en formato CSC tienen muy bajo rendimiento es por los accesos a memoria, los hilos en un grupo de trabajo podrían estar apuntado a ubicaciones en memoria que no son continuos, como dice el *OpenCL Programming for the CUDA Architecture*: "asegurándonos que los accesos a memoria global se fusionen tan seguido como sea posible es una de las más importantes optimizaciones para la arquitectura CUDA" (NVIDIA).

El problema con la memoria fusionada (o *memory coalesced*) al realizar la factorización de Cholesky en matrices en formato CSC se explica a continuación:

Sea una matriz simétrica donde los elementos de la parte triangular inferior se distribuyen como sigue:

$$A = \begin{bmatrix} 1 & * & * & * & * & * & * \\ 2 & 2 & * & * & * & * & * \\ 0 & 5 & 4 & * & * & * & * \\ 0 & 6 & 0 & 3 & * & * & * \\ 3 & 7 & 0 & 0 & 7 & * & * \\ 0 & 0 & 4 & 6 & 2 & 9 & * \\ 0 & 0 & 0 & 0 & 9 & 6 & 9 \end{bmatrix}$$

Luego de realizada la compresión a formato CSC y realizada la factorización simbólica, los elementos que compondrán el campo *elements* de la matriz en este formato son los que aparecen en negrita, recordar que los ceros que se agregan al campo dejarán de ser ceros en algún momento de la factorización, en memoria el campo *elements* sería:

Durante la factorización, al usar la función *Chol_sparse_2* (Ver función en la Figura 37 pág. 90), en el primer paso, se modifican las columnas siguientes a la primera columna, en este caso se modifican las columnas con índices del 1 hasta el 6 (recordar que la indexación es con base en cero, es decir, la primera columna tiene índice 0), cada grupo de trabajo con sus hilos se encarga de una columna, así, el grupo de trabajo

con índice 0 se encargará de modificar la segunda columna, es decir, la columna con índice 1, pero en base a los elementos de la primera columna, esto según la expresión:

$$A_{n,1} = A_{n,1} - A_{1,0}A_{n,0}, n$$
 1

En este simple ejemplo el grupo de trabajo encargado de modificar la segunda columna solo modificará los elementos correspondientes a los elementos 2 y 3 de la primera columna, es decir 2 y 7, estos elementos se muestran en negrita en la matriz siguiente en donde el símbolo de asterisco simboliza los elementos que no están almacenados en memoria.

En el campo elements sería:

Ya que solo son dos elementos a modificar entonces solo los hilos con ID's 0 y 1 realizarán la operación en la siguiente forma:

Hilo 0
$$A_{1,1} = A_{1,1} - A_{1,0}$$
 $A_{1,0} = 2 - 2$ $2 = -2$
Hilo 1 $A_{4,1} = A_{4,1} - A_{1,0}$ $A_{4,0} = 7 - 2$ $3 = 1$

La matriz modificada quedaría:

$$A = \begin{bmatrix} 1 & * & & & & \\ 2 & -2 & & & & \\ * & 5 & 4 & & & \\ * & 6 & 0 & 3 & & & \\ 3 & 1 & 0 & 0 & 7 & & \\ * & * & 4 & 6 & 2 & 9 & \\ & & & 9 & 6 & 9 \end{bmatrix}$$

Y el campo elements quedaría:

[1 **2 3** -**2** 5 6 **1** 4 0 0 4 3 0 6 7 2 9 9 6 9]

Como se puede ver, dos hilos consecutivos han escrito en ubicaciones en memoria no consecutivas ya que se saltaron los elementos 5 y 6, es esta falta de continuidad a la hora de escribir en memoria la que perjudica considerablemente el rendimiento, en este ejemplo, la distancia entre estos dos elementos es solamente dos, pero trabajando en matrices de ordenes más grandes, esta distancia puede ser mucho mayor, y no solo entre dos hilos consecutivos, sino entre más hilos de varios grupos de trabajo, que en conjunto conllevarían una considerable baja en el rendimiento.

Este comportamiento no ocurre al almacenar en formato SKS ya que los ceros entre los elementos 2 y 3 de la primera columna son almacenados en este formato, con estos ceros es posible mantener la memoria fusionada a la hora de escribir en memoria.

Otro aspecto a tomar en cuenta es el ensamblado de la matriz de rigidez, para hacer más efectivo el uso del formato SKS, la matriz de rigidez debería ser ensamblada tal que los elementos distintos de cero estén más concentrados cerca de la diagonal, esto está relacionado con la numeración de los nudos; en estructuras complejas, hacer esto podría no ser tan fácil, para estos casos el reordenamiento de filas y columnas por el algoritmo de Cuthill–McKee podría ser beneficioso.

CONCLUSIONES

- Entre las técnicas consideradas en esta tesis, las que arrojan mejores resultados (mejor balance en memoria y tiempo de ejecución), son los métodos con factorización (factorización de Cholesky y la factorización LDL^T) aplicadas sobre la matriz en formato SKS, no existiendo diferencias significativas entre estas y, en cambio, con diferencias significativas en comparación con las demás.
- La matriz sin comprimir alcanza rápidamente el límite de memoria del dispositivo (GPU). Para el dispositivo usado en esta tesis (GT 630M), el máximo orden alcanzado fue de 142422 (Ver Tabla 15 pág. 142), la memoria utilizada usando el formato SKS fue de 1483.67 MB (casi 1.5 GB); la memoria que se utilizaría en formato denso llegaría a 162273.35 MB (más de 160 GB!), capacidad que no ha sido alcanzada hasta la fecha por ninguna tarjeta gráfica, un problema de estas dimensiones ha sido resuelto por una tarjeta gráfica de apenas 2 GB de memoria, gracias a la compresión de la matriz.
- El formato de almacenamiento CSC, resulta más efectivo que el SKS para una matriz determinada, esto porque en el formato CSC se ignora la totalidad de ceros en la matriz, cosa que en el formato SKS no pasa, ya que dependiendo de la posición de los elementos distintos de cero, podría ser necesario almacenar algunos ceros, sin embargo, para los métodos con factorización, almacenar la matriz en formato CSC conlleva un problema, y es que la memoria que se necesita para almacenar la factorización es mayor a la memoria utilizada en la matriz original en formato CSC. Según los resultados, aplicar una factorización aumenta la memoria utilizada de la matriz en formato CSC a un valor incluso mayor que la que utiliza el formato SKS, la ventaja de este último es que no requiere memoria adicional para almacenar la

factorización. Al usar el método de los gradientes conjugados no se requiere ninguna factorización, por lo que la matriz original en formato CSC se mantiene, y es por tanto el método que mejor optimiza la memoria.

- Aunque el método de los gradientes conjugados, según la conclusión anterior, aplicada sobre la matriz en formato CSC, es la que mejor optimiza la memoria, el problema está en el tiempo de ejecución (Ver Tabla 11 pág. 135 y Tabla 13 pág. 140), es por eso que los métodos con factorización, aplicadas sobre la matriz en formato SKS, por tener el mejor balance de memoria y tiempo de ejecución, son las que sobresalen de las demás.
- El índice de dispersión para el formato SKS para el soporte central de puente analizado, disminuye cuando el orden aumenta; el índice de dispersión para el orden cercano a 142422 es de 0.009, lo que significa un 0.09% de la memoria que se utilizaría si almacenáramos en formato denso (Ver Gráfico 7 pág. 143).
- Las pruebas fueron realizadas en una **GeForce GT 630M**, cuyas características principales son: número de multiprocesadores 2, número de núcleos CUDA 96, con capacidad de cálculo (*compute capability*) 2.1, cantidad total de memoria global 2048 MB, cantidad máxima de hilos (o work ítems) en un grupo de trabajo 1024, cantidad total de memoria local por grupo de trabajo 48 KB. En general, al usarse otros dispositivos, el rendimiento mejorará según aumenten estos valores.
- Adicionalmente para notar el impacto que tiene usar otro dispositivo con mejores características se ha usado la tarjeta gráfica GTX 1070, cuyas características en contraste con la GT 630M puede observarse en la Tabla 17 (pág. 145). La diferencia encontrada al usar otro dispositivo, usando las mismas funciones (es decir los mismos kernels), resolviendo el mismo sistema de ecuaciones, es bastante significativa, para el orden de 142422 la tarjeta gráfica GT 630M arrojó un tiempo de 225.44 segundos para la factorización

de Cholesky sobre la matriz SKS, mientras que la GTX 1070, para el mismo método y formato arrojó un tiempo de 23.11 segundos.

- Visto los resultados en tiempo de ejecución con los arrojados por el solver incorporado del OpenFEM, concluimos que la codificación aún está sujeta a mejoras, tanto para los métodos con factorización sobre la matriz SKS, y en mayor grado para las demás técnicas.
- Los tiempos de ejecución para los métodos con factorización sobre la matriz CSC, no incluyen el tiempo necesario para realizar la factorización simbólica, es decir, si incluyéramos este tiempo, los resultados para éste formato de almacenamiento serían algo mayores.
- Para el método de los gradientes conjugados, lo que al principio tomamos como una ventaja, en la que éste método solo requería de operaciones sencillas en cada iteración, y que podían hacerse haciendo uso de las funciones incorporadas en ArrayFire, sin necesidad de escribir un kernel; resultaron teniendo inconvenientes en el tiempo de ejecución, podemos conjeturar que esto es debido a las demasiadas llamadas, construcciones y ejecuciones de kernels que se hacen para cada función incorporada utilizada, en nuestro caso, estos podrían haberse hecho dentro de una misma llamada a un kernel único (dentro de lo que fuera posible).

RECOMENDACIONES

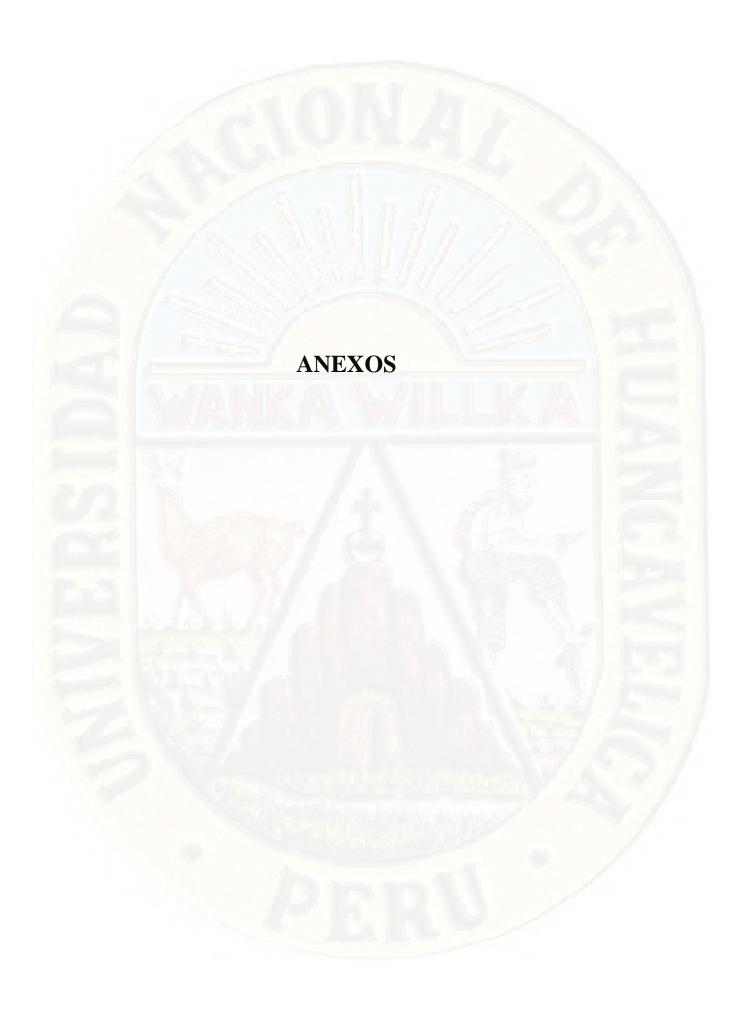
- Limitándonos a las técnicas aplicadas en ésta tesis, y a los códigos correspondientes a cada una de estas, se recomienda el uso de los métodos con factorización (factorización de Cholesky y factorización LDL^T), sobre la matriz en formato SKS, ya que son los que arrojan resultados más balanceados en memoria y tiempo de ejecución.
- Para hacer más efectivo el uso del formato SKS e incluso mejorar el uso del formato CSC, la matriz de rigidez debería ser ensamblada tal que los elementos distintos de cero estén más concentrados cerca de la diagonal, esto está relacionado con la numeración de los nudos; en estructuras complejas hacer esto podría no ser tan fácil, para estos casos se recomienda el uso del algoritmo de Cuthill–McKee para reenumerar filas y columnas, sin perder la simetría, y conseguir una matriz con elementos distintos de cero concentrados más cerca de la diagonal, obviamente el vector de fuerzas también debe pivotarse consistentemente con el pivotado de la matriz de rigidez; luego de resolver el sistema puede hacerse un pivote inverso para conseguir la solución del sistema original.
- Trabajar las lecturas o escrituras sobre memoria fusionada es uno de los temas más importantes a la hora de escribir kernels, la gran diferencia en el rendimiento entre la aplicación de los métodos sobre matrices en formato SKS y CSC recae justamente en este aspecto, básicamente se trata de que hilos consecutivos lean o escriban en ubicaciones en memoria consecutiva, como se vio en la sección de discusión de resultados, cuando este no es el caso, el rendimiento decae considerablemente.

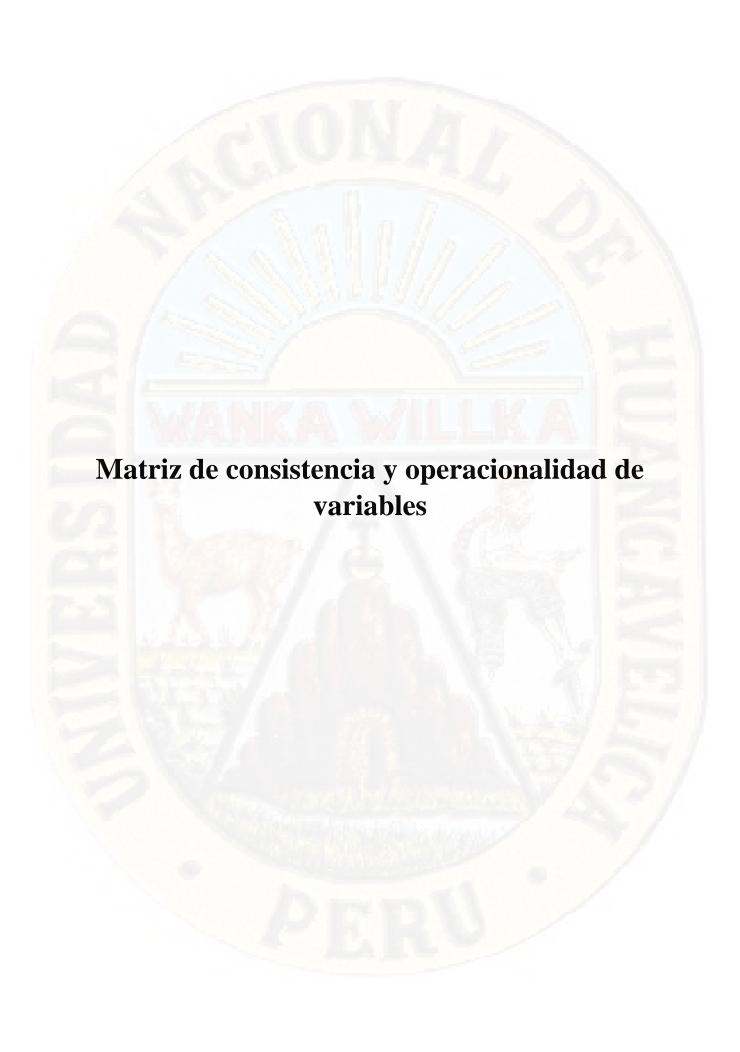
Usar un dispositivo con más multiprocesadores mejora el rendimiento, en general, mientras más multiprocesadores tenga un dispositivo, más instrucciones en paralelo podrán realizarse y por lo tanto el trabajo terminará más rápido, algo que se evidenció al comparar los resultados arrojadas por las tarjetas gráficas GT 630M y GTX 1070 en esta tesis.

Bibliografía

- ArrayFire company. (s.f.). *ArrayFire documentation*. Obtenido de http://arrayfire.org/docs/index.htm
- B.P.Demidovich, I. (1977). *Cálculo numérico fundamental*. (M. M. Calvo, Trad.) Madrid: PARANINFO.
- Celigüeta Lizarza, J. T. (2008). *Método de los Elementos Finitos para análisis estructural* (Tercera ed.). Navarra: Tecnun.
- Félix, M. V. (2015). Factorización Cholesky simbólica.
- Hernández Sampieri, R., Fernández Collado, C., & Baptista Lucio, P. (2010). *Metodología de la investigación* (Quinta ed.). México: The McGraw-Hill.
- Juan Carlos Camacho Puello, M. d. (2012). Análisis estructural con el método de elementos finitos asistido por computadora. Cartagena de Indias: Universidad Tecnológica de Bolivar.
- Juliacomputing.com. (s.f.). *GPU Programming in Julia*. Obtenido de https://juliacomputing.com/blog/2016/06/09/julia-gpu.html
- JuliaLang. (s.f.). *Julia 1.1 Documentation*. Obtenido de https://docs.julialang.org/en/v1.1/
- Niño Rojas, V. M. (2011). *Metodología de la investigación*. Bogotá: Ediciones de la U.
- NVIDIA Corporation. (n.d.). CUDA C Programming Guide. Retrieved from https://docs.nvidia.com/cuda/archive/8.0/cuda-c-programming-guide/index.html
- NVIDIA. (s.f.). OpenCL Programming for the CUDA Architecture.
- O'Connor, J. L. (1997). Técnicas de cálculo para sistemas de ecuaciones, programación lineal y programación entera. Madrid.

- Sánchez Meza, R. A. (2000). Procedimientos de gradiente conjugada en el análisis estructural. (U. N. PERÚ, Ed.) Obtenido de http://cybertesis.uni.edu.pe/handle/uni/465
- Santamaría, M. J. (1999). Factorización de Cholesky modificada de matrices dispersas sobre multiprocesadores. Universidad de Santiago de Compostela, Facultad de física, Departamento de Electrónica y Computación.
- Shucai Xiao, W.-c. F. (2009). Inter-Block GPU Communication via Fast Barrier Synchronization.
- Steven C. Rennich, D. S. (2016). Accelerating sparse Cholesky factorization on GPUs.
- stuff, W. (s.f.). *Skyline Storage Format for sparse matrices*. Obtenido de https://234satwant.wordpress.com/2014/08/29/skyline-method-for-sparse-matrices/
- Valdivia, F. A. (2014). Procedimiento para el análisis dinámico de estructuras usando el método de los elementos finitos. Lima: Pontifica Universidad Católica del Perú.
- Vargas-Félix, M. (2015). Gradiente conjugado.
- Watkins, D. S. (2002). Fundamentals of matrix computations (Second Edition ed.). (I. John Wiley & Sons, Ed.) New York: Wiley-Interscience. Obtenido de https://books.google.com.pe
- Wikipedia. (26 de Octubre de 2013). *GPGPU*. Obtenido de https://es.wikipedia.org/wiki/GPGPU
- Wikipedia. (s.f.). Sparse matrix. Obtenido de https://en.wikipedia.org/wiki/Sparse_matrix



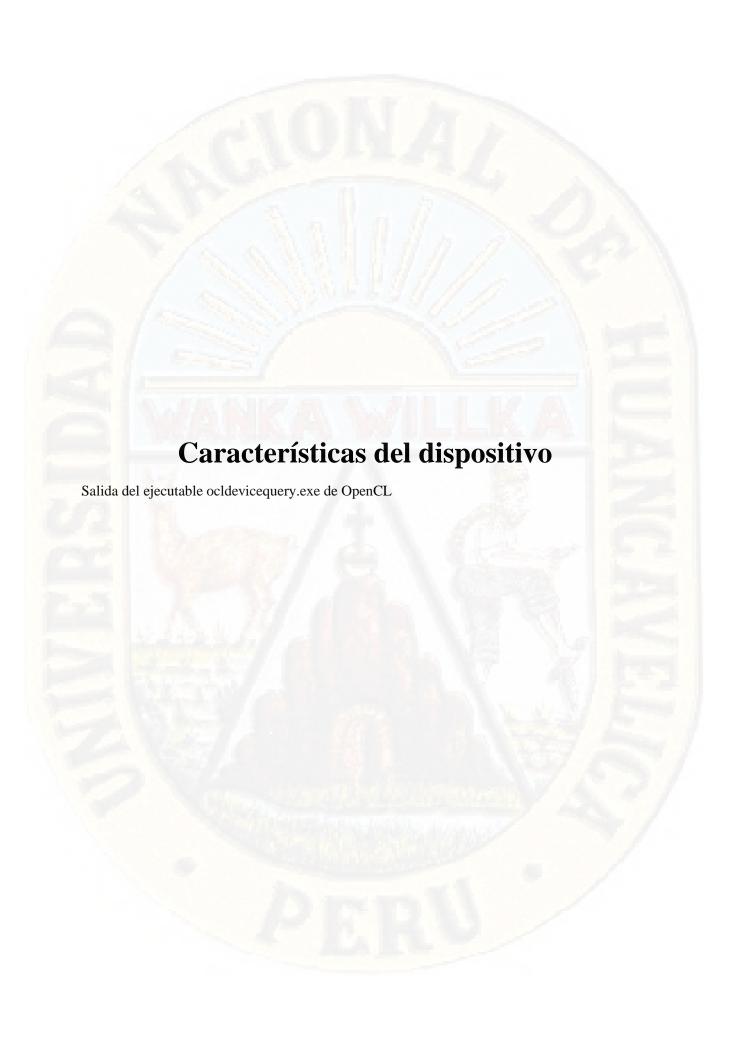


MATRIZ DE CONSISTENCIA

"TÉCNICAS NUMÉRICAS SOBRE PROCESADOR GRÁFICO, PARA LA SOLUCIÓN DE SISTEMAS DE ECUACIONES LINEALES EN ANÁLISIS ESTRUCTURAL"

OPERACIONALIDAD DE VARIABLES

Variables	Dimensiones	Metodología
V. Independiente • Técnicas numéricas sobre procesador gráfico.	 Método (formato de almacenamiento) Gauss Jordan (Matriz densa) Factorización de Cholesky (matriz densa) Factorización de Cholesky (matriz CSC) Factorización de Cholesky (matriz SKS) Factorización LDLT (Matriz densa). Factorización LDLT (Matriz CSC). Factorización LDLT (matriz SKS) Gradientes conjugados (Matriz densa). Gradientes conjugados (matriz CSC). Gradientes conjugados (matriz CSC). Gradientes conjugados (matriz SKS). 	Revisión bibliográfica
V. Dependiente • Solución de sistemas de ecuaciones lineales en análisis estructural	 Tiempo de ejecución (s) Memoria utilizada (MB) 	Obtención del sistema de ecuaciones: OpenFEM, herramienta de elementos finitos Memoria: medición manual de acuerdo al número de elementos de la matriz y la precisión de los elementos. Tiempo: Macro @time de Julia Macro @elapsed de Julia



```
G:\oclDeviceQuery.exe Starting...
OpenCL SW Info:
 CL_PLATFORM_NAME:
                     NVIDIA CUDA
 CL_PLATFORM_VERSION:
                       OpenCL 1.2 CUDA 9.1.84
 OpenCL SDK Revision:
                           7027912
OpenCL Device Info:
 1 devices found supporting OpenCL:
 Device GeForce GT 630M
 CL_DEVICE_NAME:
                                GeForce GT 630M
  CL_DEVICE_VENDOR:
                               NVIDIA Corporation
  CL_DRIVER_VERSION:
                                391.35
                               OpenCL 1.1 CUDA
 CL_DEVICE_VERSION:
  CL_DEVICE_OPENCL_C_VERSION:
                                     OpenCL C 1.1
  CL_DEVICE_TYPE:
                                 CL_DEVICE_TYPE_GPU
  CL_DEVICE_MAX_COMPUTE_UNITS:
                                     2
  CL_DEVICE_MAX_WORK_ITEM_DIMENSIONS: 3
  CL_DEVICE_MAX_WORK_ITEM_SIZES: 1024 / 1024 / 64
  CL_DEVICE_MAX_WORK_GROUP_SIZE: 1024
  CL_DEVICE_MAX_CLOCK_FREQUENCY: 950 MHz
  CL_DEVICE_ADDRESS_BITS:
                                      512 MByte
  CL_DEVICE_MAX_MEM_ALLOC_SIZE:
  CL_DEVICE_GLOBAL_MEM_SIZE:
                                      2048 MByte
  CL_DEVICE_ERROR_CORRECTION_SUPPORT: no
  CL_DEVICE_LOCAL_MEM_TYPE: local
  CL DEVICE LOCAL MEM SIZE:
                                48 KByte
  CL_DEVICE_MAX_CONSTANT_BUFFER_SIZE: 64 KByte
  CL_DEVICE_QUEUE_PROPERTIES:
     CL_QUEUE_OUT_OF_ORDER_EXEC_MODE_ENABLE
  CL_DEVICE_QUEUE_PROPERTIES:
     CL_QUEUE_PROFILING_ENABLE
  CL_DEVICE_IMAGE_SUPPORT:
  CL_DEVICE_MAX_READ_IMAGE_ARGS: 128
  CL_DEVICE_MAX_WRITE_IMAGE_ARGS:
  CL DEVICE SINGLE FP CONFIG:
                                      denorms INF-quietNaNs
round-to-nearest round-to-zero round-to-inf fma
                                      2D_MAX_WIDTH
  CL_DEVICE_IMAGE <dim>
                                                        16384
                           2D MAX HEIGHT
                                            16384
                           3D_MAX_WIDTH
                                             2048
                           3D_MAX_HEIGHT
                                            2048
                           3D_MAX_DEPTH
                                            2048
  CL_DEVICE_EXTENSIONS:
                                      cl_khr_global_int32
_base_atomics
                           cl_khr_global_int32
extended atomics
                           cl_khr_local_int32_base_atomics
                           cl_khr_local_int32_extended_atomics
```

cl_khr_fp64

```
cl_khr_byte_addressable_store
cl_khr_icd
cl_khr_gl_sharing
cl_nv_compiler_options
cl_nv_device_attribute_query
cl_nv_pragma_unroll
cl_nv_d3d10_sharing
cl_khr_d3d10_sharing
cl_nv_d3d11_sharing
cl_nv_copy_opts
```

```
CL_DEVICE_COMPUTE_CAPABILITY_NV:
                                       2.1
 NUMBER OF MULTIPROCESSORS:
                                       2
  NUMBER OF CUDA CORES:
                                       96
  CL_DEVICE_REGISTERS_PER_BLOCK_NV:
                                       32768
  CL_DEVICE_WARP_SIZE_NV:
                                 32
  CL_DEVICE_GPU_OVERLAP_NV:
                                 CL_TRUE
  CL_DEVICE_KERNEL_EXEC_TIMEOUT_NV:
                                       CL_TRUE
  CL DEVICE INTEGRATED MEMORY NV:
                                       CL FALSE
  CL_DEVICE_PREFERRED_VECTOR_WIDTH_<t>CHAR 1, SHORT 1, INT 1,
LONG 1, FLOAT 1, DOUBLE 1
```

```
2D Image Formats Supported (71)
      Channel Order
                       Channel Type
1
      CL_R
                       CL_FLOAT
2
      CL_R
                       CL_HALF_FLOAT
3
      CL_R
                       CL_UNORM_INT8
4
      CL R
                       CL UNORM INT16
5
                       CL_SNORM_INT16
      CL_R
                       CL_SIGNED_INT8
6
      CL_R
7
                       CL_SIGNED_INT16
      CL_R
      CL_R
8
                       CL_SIGNED_INT32
9
      CL_R
                       CL_UNSIGNED_INT8
10
      CL_R
                       CL_UNSIGNED_INT16
11
                       CL_UNSIGNED_INT32
      CL_R
12
      CL_A
                       CL_FLOAT
                       CL_HALF_FLOAT
13
      CL A
14
                       CL_UNORM_INT8
      CL_A
                       CL_UNORM_INT16
15
      CL_A
                       CL_SNORM_INT16
16
      CL_A
17
                       CL_SIGNED_INT8
      CL_A
18
                       CL_SIGNED_INT16
      CL_A
19
      CL_A
                       CL_SIGNED_INT32
20
      CL_A
                       CL_UNSIGNED_INT8
21
                       CL_UNSIGNED_INT16
      CL_A
22
      CL_A
                       CL_UNSIGNED_INT32
23
                       CL_FLOAT
      CL_RG
24
      CL_RG
                       CL_HALF_FLOAT
25
      CL RG
                       CL UNORM INT8
26
      CL_RG
                       CL_UNORM_INT16
27
      CL_RG
                       CL_SNORM_INT16
28
                       CL_SIGNED_INT8
      CL_RG
```

```
29
      CL_RG
                       CL_SIGNED_INT16
                       CL_SIGNED_INT32
30
      CL_RG
31
      CL_RG
                       CL_UNSIGNED_INT8
                       CL_UNSIGNED_INT16
32
      CL_RG
33
      CL_RG
                       CL_UNSIGNED_INT32
34
      CL_RA
                       CL_FLOAT
35
      CL_RA
                       CL_HALF_FLOAT
36
      CL_RA
                       CL UNORM INT8
      CL_RA
37
                       CL_UNORM_INT16
                       CL_SNORM_INT16
38
      CL_RA
39
      CL_RA
                       CL_SIGNED_INT8
40
      CL_RA
                       CL_SIGNED_INT16
41
      CL_RA
                       CL_SIGNED_INT32
42
                       CL_UNSIGNED_INT8
      CL_RA
43
      CL_RA
                       CL_UNSIGNED_INT16
44
      CL RA
                       CL_UNSIGNED_INT32
45
      CL_RGBA
                       CL_FLOAT
46
      CL_RGBA
                       CL_HALF_FLOAT
47
      CL_RGBA
                       CL_UNORM_INT8
48
      CL RGBA
                       CL UNORM INT16
49
      CL_RGBA
                       CL_SNORM_INT16
50
      CL_RGBA
                       CL_SIGNED_INT8
51
      CL_RGBA
                       CL_SIGNED_INT16
52
      CL_RGBA
                       CL_SIGNED_INT32
      CL_RGBA
53
                       CL_UNSIGNED_INT8
54
      CL_RGBA
                       CL_UNSIGNED_INT16
      CL_RGBA
                       CL_UNSIGNED_INT32
55
56
      CL_BGRA
                       CL_UNORM_INT8
57
      CL_BGRA
                       CL_SIGNED_INT8
58
      CL_BGRA
                       CL_UNSIGNED_INT8
59
      CL_ARGB
                       CL_UNORM_INT8
60
      CL_ARGB
                       CL_SIGNED_INT8
61
      CL ARGB
                       CL UNSIGNED INT8
62
      CL_INTENSITY
                       CL FLOAT
                       CL_HALF_FLOAT
      CL_INTENSITY
63
      CL_INTENSITY
                       CL_UNORM_INT8
64
      CL_INTENSITY
65
                       CL_UNORM_INT16
66
      CL_INTENSITY
                       CL_SNORM_INT16
67
      CL_LUMINANCE
                       CL_FLOAT
      CL_LUMINANCE
                       CL_HALF_FLOAT
68
69
      CL_LUMINANCE
                       CL_UNORM_INT8
                       CL UNORM INT16
70
      CL LUMINANCE
71
      CL_LUMINANCE
                       CL_SNORM_INT16
```

3D Image Formats Supported (71)

8

CL_R

Channel Type Channel Order 1 CL_R CL_FLOAT 2 CL_R CL_HALF_FLOAT 3 CL_UNORM_INT8 CL_R 4 CL_UNORM_INT16 CL_R CL SNORM INT16 5 CL R CL SIGNED INT8 6 CL_R 7 CL_R CL_SIGNED_INT16

```
9
      CL_R
                        CL_UNSIGNED_INT8
                        CL_UNSIGNED_INT16
10
      CL_R
      CL_R
                        CL_UNSIGNED_INT32
11
      CL_A
12
                        CL_FLOAT
13
      CL_A
                        CL_HALF_FLOAT
14
      CL_A
                        CL_UNORM_INT8
15
      CL_A
                        CL_UNORM_INT16
16
      CL A
                        CL SNORM INT16
      CL_A
17
                        CL_SIGNED_INT8
18
      CL_A
                        CL_SIGNED_INT16
19
      CL_A
                        CL_SIGNED_INT32
20
      CL_A
                        CL_UNSIGNED_INT8
21
      CL_A
                        CL_UNSIGNED_INT16
22
      CL_A
                        CL_UNSIGNED_INT32
23
      CL_RG
                        CL_FLOAT
24
      CL RG
                        CL_HALF_FLOAT
25
      CL_RG
                        CL_UNORM_INT8
26
                        CL_UNORM_INT16
      CL_RG
                        CL_SNORM_INT16
27
      CL_RG
28
      CL RG
                        CL SIGNED INT8
29
      CL_RG
                        CL_SIGNED_INT16
30
      CL_RG
                        CL_SIGNED_INT32
31
      CL_RG
                        CL_UNSIGNED_INT8
32
      CL_RG
                        CL_UNSIGNED_INT16
33
      CL_RG
                        CL_UNSIGNED_INT32
                        CL_FLOAT
34
      CL_RA
      CL_RA
                        CL_HALF_FLOAT
35
36
      CL_RA
                        CL_UNORM_INT8
37
      CL_RA
                        CL_UNORM_INT16
38
      CL_RA
                        CL_SNORM_INT16
39
      CL_RA
                        CL_SIGNED_INT8
40
      CL_RA
                        CL_SIGNED_INT16
41
      CL RA
                        CL SIGNED INT32
42
      CL_RA
                        CL_UNSIGNED_INT8
43
      CL_RA
                        CL_UNSIGNED_INT16
44
      CL_RA
                        CL_UNSIGNED_INT32
45
      CL_RGBA
                        CL_FLOAT
46
      CL_RGBA
                        CL_HALF_FLOAT
47
      CL_RGBA
                        CL_UNORM_INT8
48
      CL_RGBA
                        CL_UNORM_INT16
49
      CL_RGBA
                        CL_SNORM_INT16
50
      CL RGBA
                        CL SIGNED INT8
51
      CL_RGBA
                        CL_SIGNED_INT16
      CL_RGBA
52
                        CL_SIGNED_INT32
53
      CL_RGBA
                        CL_UNSIGNED_INT8
54
      CL_RGBA
                        CL_UNSIGNED_INT16
                        CL_UNSIGNED_INT32
55
      CL_RGBA
56
      CL_BGRA
                        CL_UNORM_INT8
57
      CL_BGRA
                        CL_SIGNED_INT8
58
      CL_BGRA
                        CL_UNSIGNED_INT8
59
      CL_ARGB
                        CL_UNORM_INT8
      CL_ARGB
60
                        CL_SIGNED_INT8
61
      CL_ARGB
                        CL_UNSIGNED_INT8
                        CL FLOAT
62
      CL INTENSITY
63
      CL INTENSITY
                        CL_HALF_FLOAT
64
      CL_INTENSITY
                        CL_UNORM_INT8
65
                        CL_UNORM_INT16
      CL_INTENSITY
```

```
66
     CL_INTENSITY
                      CL_SNORM_INT16
                      CL_FLOAT
67
     CL_LUMINANCE
68
     CL_LUMINANCE
                      CL_HALF_FLOAT
69
     CL_LUMINANCE
                      CL_UNORM_INT8
70
      CL_LUMINANCE
                      CL_UNORM_INT16
71
      CL_LUMINANCE
                      CL_SNORM_INT16
```

oclDeviceQuery, Platform Name = NVIDIA CUDA, Platform Version = OpenCL 1.2 CUDA 9.1.84, SDK Revision = 7027912, NumDevs = 1, Device = GeForce GT 630M

System Info:

Local Time/Date = 0:37:29, 12/11/2019

CPU Arch: 9
CPU Level: 6

of CPU processors: 8
Windows Build: 9200

Windows Ver: 6.2 (Windows Vista / Windows 7)

ocldevicequery Starting... OpenCL SW Info: CL_PLATFORM_NAME: NVIDIA CUDA CL_PLATFORM_VERSION: OpenCL 1.2 CUDA 9.2.239 OpenCL SDK Revision: 7027912 OpenCL Device Info: 1 devices found supporting OpenCL: Device GeForce GTX 1070 CL_DEVICE_NAME: GeForce GTX 1070 CL_DEVICE_VENDOR: NVIDIA Corporation CL_DRIVER_VERSION: 399.31 OpenCL 1.2 CUDA CL_DEVICE_VERSION: CL_DEVICE_OPENCL_C_VERSION: OpenCL C 1.2 CL_DEVICE_TYPE: CL_DEVICE_TYPE_GPU CL_DEVICE_MAX_COMPUTE_UNITS: CL_DEVICE_MAX_WORK_ITEM_DIMENSIONS: 3 CL_DEVICE_MAX_WORK_ITEM_SIZES: 1024 / 1024 / 64 CL_DEVICE_MAX_WORK_GROUP_SIZE: 1024 CL_DEVICE_MAX_CLOCK_FREQUENCY: 1645 MHz CL_DEVICE_ADDRESS_BITS: 2048 MByte CL_DEVICE_MAX_MEM_ALLOC_SIZE: CL_DEVICE_GLOBAL_MEM_SIZE: 8192 MByte CL_DEVICE_ERROR_CORRECTION_SUPPORT: no CL_DEVICE_LOCAL_MEM_TYPE: local CL DEVICE LOCAL MEM SIZE: 48 KByte CL_DEVICE_MAX_CONSTANT_BUFFER_SIZE: 64 KByte CL_DEVICE_QUEUE_PROPERTIES: CL_QUEUE_OUT_OF_ORDER_EXEC_MODE_ENABLE CL_DEVICE_QUEUE_PROPERTIES: CL_QUEUE_PROFILING_ENABLE CL_DEVICE_IMAGE_SUPPORT: CL_DEVICE_MAX_READ_IMAGE_ARGS: 256 CL_DEVICE_MAX_WRITE_IMAGE_ARGS: 16 CL_DEVICE_SINGLE_FP_CONFIG: denorms INF-quietNaNs round-to-nearest round-to-zero round-to-inf fma 2D_MAX_WIDTH CL_DEVICE_IMAGE <dim> 16384 2D MAX HEIGHT 32768 3D_MAX_WIDTH 16384 3D_MAX_HEIGHT 16384 3D_MAX_DEPTH 16384 CL_DEVICE_EXTENSIONS: cl_khr_global_int32 _base_atomics cl_khr_global_int32 extended atomics cl_khr_local_int32_base_atomics cl_khr_local_int32_extended_atomics

cl_khr_fp64

```
cl_khr_byte_addressable_store
cl_khr_icd
cl_khr_gl_sharing
cl_nv_compiler_options
cl_nv_device_attribute_query
cl_nv_pragma_unroll
cl_nv_d3d10_sharing
cl_khr_d3d10_sharing
cl_nv_d3d11_sharing
cl_nv_copy_opts
```

```
6.1
  CL_DEVICE_COMPUTE_CAPABILITY_NV:
                                      16
 NUMBER OF MULTIPROCESSORS:
 NUMBER OF CUDA CORES:
                                      4294967280
  CL_DEVICE_REGISTERS_PER_BLOCK_NV:
                                      65536
  CL_DEVICE_WARP_SIZE_NV:
  CL_DEVICE_GPU_OVERLAP_NV:
                                CL_TRUE
  CL_DEVICE_KERNEL_EXEC_TIMEOUT_NV:
                                      CL_TRUE
  CL_DEVICE_INTEGRATED_MEMORY_NV:
                                     CL FALSE
  CL_DEVICE_PREFERRED_VECTOR_WIDTH_<t>CHAR 1, SHORT 1, INT 1,
LONG 1, FLOAT 1, DOUBLE 1
```

2D I	mage Formats Sur	oported (75)
#	Channel Order	Channel Type
1 2 3 4 5 6 7 8	CL_R CL_R CL_R CL_R CL_R CL_R CL_R CL_R	CL_FLOAT CL_HALF_FLOAT CL_UNORM_INT8 CL_UNORM_INT16 CL_SNORM_INT16 CL_SIGNED_INT8 CL_SIGNED_INT16 CL_SIGNED_INT32 CL_UNSIGNED_INT8
10 11	CL_R CL_R	CL_UNSIGNED_INT16 CL_UNSIGNED_INT32
12 13 14 15 16 17 18 19 20 21 22 23	CL_A CL_A CL_A CL_A CL_A CL_A CL_A CL_A	CL_FLOAT CL_HALF_FLOAT CL_UNORM_INT8 CL_UNORM_INT16 CL_SNORM_INT16 CL_SIGNED_INT8 CL_SIGNED_INT16 CL_SIGNED_INT32 CL_UNSIGNED_INT8 CL_UNSIGNED_INT16 CL_UNSIGNED_INT16
24 25 26 27 28	CL_RG CL_RG CL_RG CL_RG CL_RG CL_RG	CL_FLOAT CL_HALF_FLOAT CL_UNORM_INT8 CL_UNORM_INT16 CL_SNORM_INT16 CL_SIGNED_INT8

```
29
      CL_RG
                       CL_SIGNED_INT16
                       CL_SIGNED_INT32
30
      CL_RG
31
      CL_RG
                       CL_UNSIGNED_INT8
                       CL_UNSIGNED_INT16
32
      CL_RG
33
      CL_RG
                       CL_UNSIGNED_INT32
34
      CL_RA
                       CL_FLOAT
35
      CL_RA
                       CL_HALF_FLOAT
36
      CL_RA
                       CL UNORM INT8
      CL_RA
37
                       CL_UNORM_INT16
                       CL_SNORM_INT16
38
      CL_RA
39
      CL_RA
                       CL_SIGNED_INT8
40
      CL_RA
                       CL_SIGNED_INT16
41
      CL_RA
                       CL_SIGNED_INT32
42
                       CL_UNSIGNED_INT8
      CL_RA
43
      CL_RA
                       CL_UNSIGNED_INT16
44
      CL RA
                       CL_UNSIGNED_INT32
45
      CL_RGBA
                       CL_FLOAT
46
      CL_RGBA
                       CL_HALF_FLOAT
47
      CL_RGBA
                       CL_UNORM_INT8
48
      CL RGBA
                       CL UNORM INT16
49
      CL_RGBA
                       CL_SNORM_INT16
50
      CL_RGBA
                       CL_SIGNED_INT8
51
      CL_RGBA
                       CL_SIGNED_INT16
52
      CL_RGBA
                       CL_SIGNED_INT32
      CL_RGBA
53
                       CL_UNSIGNED_INT8
54
      CL_RGBA
                       CL_UNSIGNED_INT16
                       CL_UNSIGNED_INT32
55
      CL_RGBA
56
      CL_BGRA
                       CL_UNORM_INT8
57
      CL_BGRA
                       CL_SIGNED_INT8
58
      CL_BGRA
                       CL_UNSIGNED_INT8
59
      CL_ARGB
                       CL_UNORM_INT8
60
      CL_ARGB
                       CL_SIGNED_INT8
61
      CL ARGB
                       CL UNSIGNED INT8
62
      CL_INTENSITY
                       CL FLOAT
      CL_INTENSITY
                       CL_HALF_FLOAT
63
                       CL_UNORM_INT8
      CL_INTENSITY
64
      CL_INTENSITY
65
                       CL_UNORM_INT16
66
      CL_INTENSITY
                       CL_SNORM_INT16
67
      CL_LUMINANCE
                       CL_FLOAT
      CL_LUMINANCE
                       CL_HALF_FLOAT
68
69
      CL_LUMINANCE
                       CL_UNORM_INT8
                       CL UNORM INT16
70
      CL LUMINANCE
71
      CL_LUMINANCE
                       CL_SNORM_INT16
72
      CL_BGRA
                       CL_SNORM_INT8
73
                       CL_SNORM_INT16
      CL_BGRA
74
      CL ARGB
                       CL SNORM INT8
75
      CL_ARGB
                       CL_SNORM_INT16
3D Image Formats Supported (75)
```

Channel Order Channel Type

1 CL_R CL_FLOAT
2 CL_R CL_HALF_FLOAT
3 CL_R CL_UNORM_INT8
4 CL_R CL_UNORM_INT16

```
5
      CL_R
                        CL_SNORM_INT16
6
      CL_R
                        CL_SIGNED_INT8
7
      CL_R
                        CL_SIGNED_INT16
8
      CL_R
                        CL_SIGNED_INT32
9
      CL_R
                        CL_UNSIGNED_INT8
10
      CL_R
                        CL_UNSIGNED_INT16
11
      CL_R
                        CL_UNSIGNED_INT32
12
      CL A
                        CL FLOAT
13
      CL_A
                        CL_HALF_FLOAT
14
      CL_A
                        CL_UNORM_INT8
15
      CL_A
                        CL_UNORM_INT16
16
      CL_A
                        CL_SNORM_INT16
17
      CL_A
                        CL_SIGNED_INT8
18
                        CL_SIGNED_INT16
      CL_A
19
      CL_A
                        CL_SIGNED_INT32
20
      CL A
                        CL_UNSIGNED_INT8
21
      CL_A
                        CL_UNSIGNED_INT16
22
                        CL_UNSIGNED_INT32
      CL_A
23
      CL_RG
                        CL_FLOAT
24
      CL RG
                        CL_HALF_FLOAT
25
      CL_RG
                        CL_UNORM_INT8
26
      CL_RG
                        CL_UNORM_INT16
27
      CL_RG
                        CL_SNORM_INT16
28
      CL_RG
                        CL_SIGNED_INT8
29
      CL_RG
                        CL_SIGNED_INT16
30
      CL_RG
                        CL_SIGNED_INT32
      CL_RG
                        CL_UNSIGNED_INT8
31
32
      CL_RG
                        CL_UNSIGNED_INT16
33
      CL_RG
                        CL_UNSIGNED_INT32
34
      CL_RA
                        CL_FLOAT
35
                        CL_HALF_FLOAT
      CL_RA
36
      CL_RA
                        CL_UNORM_INT8
37
      CL RA
                        CL UNORM INT16
38
      CL_RA
                        CL_SNORM_INT16
39
      CL_RA
                        CL_SIGNED_INT8
                        CL_SIGNED_INT16
40
      CL_RA
                        CL_SIGNED_INT32
41
      CL_RA
42
      CL_RA
                        CL_UNSIGNED_INT8
      CL_RA
                        CL_UNSIGNED_INT16
43
44
      CL_RA
                        CL_UNSIGNED_INT32
45
      CL_RGBA
                        CL_FLOAT
                        CL HALF FLOAT
46
      CL RGBA
47
      CL_RGBA
                        CL_UNORM_INT8
      CL_RGBA
                        CL_UNORM_INT16
48
                        CL_SNORM_INT16
49
      CL_RGBA
      CL_RGBA
50
                        CL SIGNED INT8
51
      CL_RGBA
                        CL_SIGNED_INT16
52
      CL_RGBA
                        CL_SIGNED_INT32
53
      CL_RGBA
                        CL_UNSIGNED_INT8
54
      CL_RGBA
                        CL_UNSIGNED_INT16
55
      CL_RGBA
                        CL_UNSIGNED_INT32
56
      CL_BGRA
                        CL_UNORM_INT8
                        CL_SIGNED_INT8
57
      CL_BGRA
58
      CL BGRA
                        CL UNSIGNED INT8
59
      CL ARGB
                        CL UNORM INT8
                        CL_SIGNED_INT8
60
      CL_ARGB
61
      CL_ARGB
                        CL_UNSIGNED_INT8
```

```
62
      CL_INTENSITY
                       CL_FLOAT
      CL_INTENSITY
                       CL_HALF_FLOAT
63
64
      CL_INTENSITY
                       CL_UNORM_INT8
65
                       CL_UNORM_INT16
      CL_INTENSITY
66
      CL_INTENSITY
                       CL_SNORM_INT16
67
      CL_LUMINANCE
                       CL_FLOAT
68
                       CL_HALF_FLOAT
      CL_LUMINANCE
69
      CL LUMINANCE
                       CL UNORM INT8
                       CL_UNORM_INT16
70
      CL_LUMINANCE
71
      CL_LUMINANCE
                       CL_SNORM_INT16
72
      CL_BGRA
                       CL_SNORM_INT8
73
      CL_BGRA
                       CL_SNORM_INT16
74
      CL_ARGB
                       CL_SNORM_INT8
75
      CL_ARGB
                       CL_SNORM_INT16
```

oclDeviceQuery, Platform Name = NVIDIA CUDA, Platform Version = OpenCL 1.2 CUDA 9.2.239, SDK Revision = 7027912, NumDevs = 1, Device = GeForce GTX 1070

System Info:

Local Time/Date = 15:55:36, 11/23/2019

CPU Arch: 9
CPU Level: 6

of CPU processors: 8
Windows Build: 9200

Windows Ver: 6.2 (Windows Vista / Windows 7)

Código fuente

Se muestra el código fuente del archivo donde están escritas las funciones usando la biblioteca de ArrayFire, que llaman a los kernels OpenCL, según el siguiente detalle:

Formato de la matriz	Método	Kernel	función	línea
Denso	Gauss Jordan	Gauss_Jordan_c	SELgj_c	807
	Factorización de Cholesky	Cholesky_c	SELchol_c, fac_chol_c	1681, 1828
	Factorización LDL [™]	ldlt_c	SELldlt_c, fac_ldlt_c	1930, 2071
	Gradientes conjugados		SEL_gc	1289
CSC	Factorización de Cholesky	Cholesky_sparse_c	<pre>fac_sparse_chol_c, SELchol_sparse_c</pre>	2165, 2549
	Factorización LDL ^T	ldlt_sparse_c	<pre>fac_sparse_ldlt_c, SELldlt_sparse_c</pre>	2920, 3090
	Gradientes conjugados		SELgc_sparse	1406
SKS Factor de Cho	Factorización de Cholesky	chol_sparse_sk	<pre>fac_sparse_chol_sks , SELchol_sparse_sks</pre>	2692, 2798
	Factorización LDL ^T	ldlt_sparse_sks	<pre>fac_sparse_ldlt_sks , SELldlt_sparse_sks</pre>	3235, 3335
	Gradientes conjugados		SELgc_sparse_sks	1527

Nota: Los códigos (kernels y funciones) de la presente tesis están almacenadas en el repositorio: https://github.com/4lrdyD/oclLinealSol, los números de línea de la tabla podrían no mantenerse debido a posibles actualizaciones de código.

```
//revisión 0.9.6 01-12-2019, 18:10 VS 2017
 4
 5 #include "Header.h"
 6 #include "Common.h"
 7 char * matrixMul_source=AFire::kernel_src("Common.h",
       "matrixMul.cl");
 8
   char * GJordan_source = AFire::kernel_src("Common.h",
 9
10
       "Gauss_Jordan.cl");
11 char * Cholesky_source = AFire::kernel_src("Common.h",
       "Cholesky.cl");
12
13 char * ldlt source = AFire::kernel src("Common.h", "ldlt.cl");
14 char * gc_source = AFire::kernel_src("Common.h",
        "Gradiente_conjugado.cl");
15
   char * util_src = AFire::kernel_src("Common.h", "util.cl");
16
17
   using namespace af;
18
19 //----
20 //utilidades
21 //----
22 void AFire::copy(const af::array &A, af::array &B,size_t length) {
23
24
       // 2. Obtain the device, context, and queue used by ArrayFire
25
       static cl_context af_context = afcl::getContext();
       static cl_device_id af_device_id = afcl::getDeviceId();
26
27
       static cl command queue af queue = afcl::getQueue();
28
       // 3. Obtain cl_mem references to af::array objects
29
       cl_mem * d_A = A.device<cl_mem>();
       cl_mem * d_B = B.device<cl_mem>();
30
31
32
       // 4. Load, build, and use your kernels.
33
             For the sake of readability, we have omitted error checking.
34
       int status = CL SUCCESS;
       // A simple copy kernel, uses C++11 syntax for multi-line strings.
35
       const char * kernel_name = "copy_kernel";
36
       const char * source = R"(
37
           void kernel
38
39
           copy_kernel(__global float * gA, __global float * gB)
40
               int id = get global id(0);
41
42
               gB[id] = gA[id];
43
44
       // Create the program, build the executable, and extract the entry point
45
46
       // for the kernel.
47
       cl_program program = clCreateProgramWithSource(af_context, 1, &source,
         NULL, &status);
48
       status = clBuildProgram(program, 1, &af_device_id, NULL, NULL);
49
       cl_kernel kernel = clCreateKernel(program, kernel_name, &status);
       // Set arguments and launch your kernels
50
       clSetKernelArg(kernel, 0, sizeof(cl_mem), d_A);
51
       clSetKernelArg(kernel, 1, sizeof(cl_mem), d_B);
52
53
       clEnqueueNDRangeKernel(af_queue, kernel, 1, NULL, &length, NULL, 0, NULL, →
          NULL);
       // 5. Return control of af::array memory to ArrayFire
54
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
7
```

```
55
         A.unlock();
 56
         B.unlock();
 57
         // ... resume ArrayFire operations
         // Because the device pointers, d_x and d_y, were returned to ArrayFire's
 58
 59
         // control by the unlock function, there is no need to free them using
        // clReleaseMemObject()
 60
 61
    }
 62
    void AFire::sumaaf(af_array* out , af_array dA, af_array dB) {
 63
 64
         //to store the result
 65
         af array dC;
 66
 67
         af_copy_array(&dC, dA);
 68
         // 2. Obtain the device, context, and queue used by ArrayFire
 69
         static cl_context af_context = afcl::getContext();
 70
 71
         static cl_device_id af_device_id = afcl::getDeviceId();
 72
         static cl_command_queue af_queue = afcl::getQueue();
 73
 74
         dim t order[AF MAX DIMS];
 75
        af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
 76
         size_t order = _order[0];
 77
 78
        int status = CL SUCCESS;
 79
         // 3. Obtain cl_mem references to af_array objects
 80
 81
         cl mem *d A = (cl mem*)clCreateBuffer(af context,
 82
             CL_MEM_READ_ONLY, sizeof(float) * order,
 83
             NULL, &status);
 84
         af_get_device_ptr((void**)d_A, dA);
 85
 86
         cl_mem *d_B = (cl_mem*)clCreateBuffer(af_context,
             CL_MEM_READ_ONLY, sizeof(float) * order,
 87
 88
             NULL, &status);
         af_get_device_ptr((void**)d_B, dB);
 89
 90
         cl_mem *d_C = (cl_mem*)clCreateBuffer(af_context,
 91
             CL_MEM_WRITE_ONLY, sizeof(float) * order,
 92
 93
             NULL, &status);
         af_get_device_ptr((void**)d_C, dC);
 94
 95
        // 4. Load, build, and use your kernels.
 96
 97
               For the sake of readability, we have omitted error checking.
 98
         // A simple sum kernel, uses C++11 syntax for multi-line strings.
        const char * kernel_name = "sum_kernel";
 99
         const char * source = R"(
100
             void <u>kernel</u>
101
102
             sum_kernel(__global float * gC, __global float * pA, __global float * >
                gB)
103
                 int id = get_global_id(0);
104
105
                 gC[id] = gA[id]+gB[id];
106
         )";
107
        // Create the program, build the executable, and extract the entry point
108
        // for the kernel.
109
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
110
         cl program program = clCreateProgramWithSource(af context, 1, &source,
           NULL, &status);
         status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
111
112
         cl_kernel sumkernel = clCreateKernel(program, kernel_name, &status);
113
         // Set arguments and launch your kernels
114
         clSetKernelArg(sumkernel, 0, sizeof(cl_mem), d_C);
115
         clSetKernelArg(sumkernel, 1, sizeof(cl_mem), d_A);
         clSetKernelArg(sumkernel, 2, sizeof(cl_mem), d_B);
116
         clEnqueueNDRangeKernel(af_queue, sumkernel, 1, NULL, &order, NULL, 0,
117
           NULL, NULL);
118
119
         // 5. Return control of af::array memory to ArrayFire
120
         af unlock array(dA);
121
         af unlock array(dB);
122
         af_unlock_array(dC);
123
124
         //copy results to output argument
125
         af_copy_array(out, dC);
126
127
        // ... resume ArrayFire operations
128
        // Because the device pointers, d_x and d_y, were returned to ArrayFire's
        // control by the unlock function, there is no need to free them using
129
130
        // clReleaseMemObject()
131
    }
132
133 void AFire::printh(float *A,size_t length)
134 {
135
136
        for (int i = 0; i < length; i++)</pre>
137
138
             std::cout << A[i] << std::endl;</pre>
139
140
         std::cout << std::endl;</pre>
141 }
142
    void AFire::mulaf(af_array* dC, const af_array dA,
143
        const af_array dB) {
144
145
         //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
146
         static cl_context af_context = afcl::getContext();
         static cl_device_id af_device_id = afcl::getDeviceId();
147
148
         static cl command queue af queue = afcl::getQueue();
149
150
151
         //detalles importantes sobre el kernel
152
         //-----
         /*suponiendo que las matrices son compatibles para la multiplicación:
153
         El detalle aquí es el uso adecuado del Kernel en matrixMul.cl
154
155
         el kernel está escrito para multiplicar matrices cuyos datos están
           ordenados por filas,
156
         sin embargo las matrices de tipo af::array son ordenadas por columnas, si ₹
            establecemos los argumentos
         en el kernel sin tomar en cuenta este detalle, el resultado será erroneo.
157
         el Kernel acepta ocho argumentos: suponiendo C=A*B
158
159
         1.vector que guarda los datos de C
160
         2.vector que guarda los datos de A
161
         3.vector que guarda los datos de B
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
162
         4.cero con el tamaño necesario de memoria para un grupo en A
163
         5.cero con el tamaño necesario de memoria para un grupo en B
164
         6.ancho de A
165
         7.ancho de B
         8.alto de A
166
167
         sean 2 matrices del tipo af::array
168
169
170
      A=
171
      0
             8
                12 16
                        20
                             24
                                 28
         5
            9
                        21
                             25
                                 29
172
      1
                13
                    17
173
      2 6
            10 14
                    18
                        22
                            26
                                 30
174
         7
            11
                15
                    19
                        23
                            27
175
176 B=
          9
             17
                 25
177
      1
                     33
                         41
178
      2
         10
             18
                 26
                     34
                         42
179
      3
         11 19
                 27
                     35
                         43
180
      4
         12 20 28
                     36 44
181
      5
         13 21
                 29
                     37 45
182
         14 22
                 30
                     38 46
      6
         15 23
183
      7
                     39
                         47
                 31
184
      8
         16 24
                 32
                     40 48
185
         si queremos multiplicar estas 2 matrices usando este kernel, los datos
186
           serían:
         dA=\{0,1,2,...,31\}
187
188
         dB=\{1,2,3,\ldots,48\}
189
         así son extraídos con: A.device<cl_mem>();
         como podemos comprobar los están ordenados por columnas al construir las
190
           matrices correspondientes
191
         así como estan definidas las matrices las dimensiones serían:
192
193
194
         ancho de A=8
195
         ancho de B=6
         alto de A=4
196
197
198
         sin embargo si usamos el kernel en matrixMul.cl, con los argumentos así
           como están, en realidad se
         estarían multiplicando estas matrices:
199
200
201 A=
202
       0
               2
                            5
                                6
                                    7
           1
                   3
203
       8
           9
              10
                  11
                       12
                           13
                               14
                                   15
              18
                  19
                               22
                                   23
204
      16
          17
                       20
                           21
      24
              26
          25
                           29
205
                  27
                      28
                               30
                                   31
206
207
     B=
208
           2
               3
                   4
                       5
                           6
       1
209
       7
               9
           8
                  10 11
                           12
              15
210
      13 14
                  16
                      17
                           18
      19
          20
              21
                  22
                       23
                           24
211
212
      25
          26
              27
                  28
                      29
                           30
213
              33
                  34
                      35
                           36
      31
          32
214
      37
         38
              39
                  40
                      41
                           42
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
215
          44 45
                       47
                           48
                  46
216
217
      que son diferentes a las originales, ya que como se dijo al comienzo el
        kernel asume que los valores estan
      ordenados por filas, si se realiza la multiplicación el resultado será
218
        erroneo.
219
      Sin embargo es posible usar el mismo kernel para obtener un resultado
220
        válido, y es invirtiendo el orden de
221
      los argumentos de entrada, invirtiendo también para cada unos sus
        dimensiones, es decir colocando B de 6x8 como
      primer argumento y A de 8x4 como segundo argumento, como se podrá comprobar →
222
        el resultado final es el esperado.
223
224
      B=
            2
                3
                                7
225
       1
                    4
                        5
                            6
                                     8
       9
          10
               11
                   12
                           14
                               15
 226
                       13
                                   16
                   20
          18
               19
                       21
                           22
                               23
                                   24
227
      17
228
      25
          26
               27
                   28 29
                           30
                               31
                                   32
229
      33
               35
                   36
                           38
                               39
           34
                       37
 230
      41
          42
              43
                   44
                       45
                           46
                               47
                                   48
231
232
     A=
 233
       0
           1
                2
                    3
       4
                    7
234
           5
                6
       8
           9
               10
 235
                   11
236
      12
          13
               14
                   15
 237
      16
          17
               18
                   19
238
      20
          21
               22
                   23
239
      24
          25
               26
                   27
      28
          29
               30
                   31
240
 241
      B*A=
242
243
      672
            708
                   744
                         780
 244
      1568
            1668
                   1768
                         1868
245
      2464
            2628
                         2956
                   2792
 246
      3360
            3588
                   3816
                         4044
247
      4256
            4548
                   4840
                         5132
248
      5152
            5508
                   5864
                         6220
249
250
      A*B= (A y B originales)
251
       672 1568 2464
                        3360
                              4256
                                      5152
                        3588 4548
252
      708 1668 2628
                                    5508
253
      744 1768
                  2792
                              4840
                                     5864
                        3816
254
           1868
                  2956
                        4044
                              5132
                                     6220
255
      no es necesario transponer B*A (cambiados) para obtener A*B (originales) ya →
256
        que el resultado saldrá
257
      en forma de vector, estos se ordenarán automáticamente por columna una vez
        vuelvan a ser parte del campo de C
258
      o matriz resultado.
259
260
      */
261
         dim_t Bdims[AF_MAX_DIMS];
262
         af_get_dims(&Bdims[0], &Bdims[1], &Bdims[2], &Bdims[3], dB);
263
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
264
          size t wB = Bdims[1];
265
          size_t hB = Bdims[0];
266
267
         dim_t Adims[AF_MAX_DIMS];
         af_get_dims(&Adims[0], &Adims[1], &Adims[2], &Adims[3], dA);
268
269
          size_t hA = Adims[0];
270
         af_dtype typef;
271
272
         af_get_type(&typef, dA);
273
274
         int msize = 0;
          if (typef == f64)
275
276
             msize = sizeof(double);
277
         else if (typef == f32)
278
             msize = sizeof(float);
279
         else;
280
         //matriz resultado
281
282
         af_array Rs;
283
         unsigned ndims = 2;
 284
         dim_t dim[] = { Adims[0],Bdims[1] };
         af_randu(&Rs, ndims, dim, typef);
285
286
 287
          size_t program_length = strlen(matrixMul_source);
288
          int status = CL_SUCCESS;
 289
290
         //obteniendo las referencias cl_mem de los objetos af::array
 291
292
         cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_ONLY, msize * Adims[0] * Adims[1],
293
294
              NULL, &status);
 295
         af_get_device_ptr((void**)d_A, dA);
296
297
          cl mem *d B = (cl mem*)clCreateBuffer(af context,
              CL MEM READ ONLY, msize * Bdims[0] * Bdims[1],
298
299
              NULL, &status);
         af_get_device_ptr((void**)d_B, dB);
 300
301
302
          cl_mem *d_C = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_WRITE_ONLY, msize * Adims[0] * Bdims[1],
303
304
              NULL, &status);
305
```

```
af_get_device_ptr((void**)d_C, Rs);
306
307
         //creando el programa, construyendo el ejecutable y extrayendo el punto
           de entrada
308
         // para el Kernel
309
         cl_program program = clCreateProgramWithSource(
310
             af_context, 1, (const char **)&matrixMul_source,
311
             &program_length, &status);
312
         status = clBuildProgram(program, 1, &af_device_id,
313
             NULL, NULL, NULL);
314
315
316
         char* kernelName;
317
         if (typef == f64)
             kernelName = "matrixMul";
318
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
319
         else if (typef == f32)
             kernelName = "matrixMul_sp";
320
321
         else;
322
         cl_kernel mulkernel = clCreateKernel(program,
323
324
             kernelName, &status);
325
         // estableciendo los argumentos, ver detalles importantes sobre el kernel >
326
             (más arriba)
327
         int i = 0;
         clSetKernelArg(mulkernel, i++, sizeof(cl_mem), d_C);
328
         clSetKernelArg(mulkernel, i++, sizeof(cl_mem), d_B);
329
330
         clSetKernelArg(mulkernel, i++, sizeof(cl_mem), d_A);
331
         clSetKernelArg(mulkernel, i++,
             sizeof(double) * BLOCK_SIZE *BLOCK_SIZE, 0);
332
         clSetKernelArg(mulkernel, i++,
333
334
              sizeof(double) * BLOCK_SIZE *BLOCK_SIZE, 0);
335
         clSetKernelArg(mulkernel, i++, sizeof(cl_int), &hB);
336
         clSetKernelArg(mulkernel, i++, sizeof(cl_int), &hA);
337
         clSetKernelArg(mulkernel, i++, sizeof(cl_int), &wB);
338
         size_t localWorkSize[] = { BLOCK_SIZE, BLOCK_SIZE };
339
         size_t globalWorkSize[] = { shrRoundUp(BLOCK_SIZE, hA),
340
             shrRoundUp(BLOCK_SIZE, wB) };
 341
342
         //ejecutando el Kernel
 343
344
         clEnqueueNDRangeKernel(af queue, mulkernel, 2, 0,
             globalWorkSize, localWorkSize, 0, NULL, NULL);
345
346
         //devolviendo el control de memoria af::array a ArrayFire
347
         af unlock array(dA);
348
 349
         af_unlock_array(dB);
350
         af_unlock_array(Rs);
351
         // ... reanudando las operaciones ArrayFire
352
     //copiando el resultado en el argumento de salida
353
         af_copy_array(dC, Rs);
354
355
         af_release_array(Rs);
356
357
358 char * AFire::find kernel test(const char* file)
359 {
         char * cl_files = AFire::get_env_path(CL_VAR_ENT);
360
361
         std::string cd_files_path = cl_files;
362
         std::string file name = file;
         std::string path_to_file = cd_files_path + "/" + file_name;
363
364
         size_t program_length;
365
         char *header = oclLoadProgSource(path_to_file.c_str(),
           &program_length);
366
         return header;
367
     }
368
369 char * AFire::helpk(const char * string)
370 {
         char * str = (char *)malloc(10 * sizeof(char));
371
372
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
8
```

```
373
         strcpy_s(str, sizeof(char) * 10, "hola ");
         strcat_s(str, sizeof(char) * 10, "will");
374
375
         return str;
376
         //return (char*)string;
377
378 }
379
380 //Obtiene la ruta explicita de una variable de
     //entorno, si no se encuentra la variable
382
    //devolverá "./"
383 char * AFire::get_env_path(const char * env)
384 {
385
         std::string env s = env;
386
         env_s += "=";
         int cont = 0;
387
388
         bool bool_s = 0;
389
         size_t pos_s;
390
         std::string path_s;
391
392
         extern char ** environ;
393
         while (environ[cont] != NULL && !bool s)
394
             path s = std::string(environ[cont]);
395
396
             pos s = path s.find(env s);
             if (pos_s != std::string::npos)
397
398
399
                 path_s.replace(pos_s, env_s.length(), "");
400
                 bool_s = 1;
401
             }
402
             cont += 1;
403
         }
404
         if (bool_s)
405
406
             char* file path = (char*)malloc(path s.length() + 1);
407
             strcpy_s(file_path, path_s.length() + 1, path_s.c_str());
408
409
             return file_path;
410
         }
411
         else
412
             return "./";
413
414
         }
415
     }
416
     /*Obtiene el código fuente de un kernel (OpenCL) y lo concatena
417
     con el correspondiente encabezado, el programa estará listo para
418
     construirse con el código devuelto
419
420 */
421 char * AFire::kernel_src(const char* hfile, const char * clfile)
422
423
         size_t program_length;
         //obteniendo el valor de la variable de entorno
424
425
         char * cl files = AFire::get env path(CL VAR ENT);
426
         std::string cl_files_path = cl_files;
427
428
         //encabezado comun a los Kernels
```

```
...sual Studio 2013\Projects\AF_func\AF_sample.cpp
429
         std::string file name = hfile;
         std::string header path = cl files path + "/" + file name;
430
431
         char *header = oclLoadProgSource(header_path.c_str(), "",
           &program_length);
432
433
         //kernel
434
         file_name = clfile;
         std::string source_path = cl_files_path + "/" + file_name;
435
436
         char * source = oclLoadProgSource(source_path.c_str(), header,
           &program_length);
437
         return source;
438 }
439
     void AFire::sparse_mat_vec_mul(af_array* dC, af_array elmA,
440
441
         af_array colA, af_array rowA, af_array dB) {
         //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
442
443
         //cl_context af_context;
444
         static cl_context af_context = afcl::getContext();
445
         static cl_device_id af_device_id = afcl::getDeviceId();
446
         static cl_command_queue af_queue = afcl::getQueue();
447
448
         //creando copia de dB
449
         af array c;
450
         af_copy_array(&c, dB);
451
         //2. Obteniendo parámetros necesarios
452
453
         //longitud de los vectores
454
455
         dim_t _order[AF_MAX_DIMS];
         af_get_dims(&_order[0], &_order[1], &_order[2],
456
             &_order[3], elmA);
457
458
         cl_int size_elmA = _order[0];
459
460
         af get dims(& order[0], & order[1], & order[2],
461
             &_order[3], colA);
         cl_int size_colA = _order[0];
462
463
464
         af_get_dims(&_order[0], &_order[1], &_order[2],
465
             &_order[3], rowA);
466
         cl_int size_rowA = _order[0];
467
         size t localWorkSize = BLOCK SIZE * BLOCK SIZE;
468
469
         size t globalWorkSize = localWorkSize * BLOCK SIZE;
470
471
         int status = CL SUCCESS;
472
473
         af_dtype typef;
474
         af_get_type(&typef, elmA);
475
476
         int msize = 0;
477
         if (typef == f64)
             msize = sizeof(double);
478
479
         else if (typef == f32)
480
             msize = sizeof(float);
481
         else;
482
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
483
         //3.obteniendo las referencias cl mem de los objetos af::array
484
         cl mem *d elmA = (cl mem*)clCreateBuffer(af context,
485
             CL_MEM_READ_ONLY, msize*size_elmA,
486
             NULL, &status);
         af_get_device_ptr((void**)d_elmA, elmA);
487
488
489
        cl_mem *d_colA = (cl_mem*)clCreateBuffer(af_context,
             CL_MEM_READ_ONLY, sizeof(int)*size_colA,
490
491
             NULL, &status);
492
         af_get_device_ptr((void**)d_colA, colA);
493
494
         cl_mem *d_rowA = (cl_mem*)clCreateBuffer(af_context,
495
             CL MEM READ ONLY, sizeof(int)*size rowA,
496
             NULL, &status);
497
         af_get_device_ptr((void**)d_rowA, rowA);
498
499
         cl_mem *d_b = (cl_mem*)clCreateBuffer(af_context,
500
             CL_MEM_READ_ONLY, msize*size_colA,
501
             NULL, &status);
502
         af_get_device_ptr((void**)d_b, dB);
503
504
        cl mem *d c = (cl mem*)clCreateBuffer(af context,
             CL MEM READ WRITE, msize*size colA,
505
506
             NULL, &status);
        af_get_device_ptr((void**)d_c, c);
507
508
509
         size_t program_length = strlen(gc_source);
510
511
         //4.creando el programa, construyendo el ejecutable y extrayendo el punto →
            de entrada
512
         // para el Kernel
513
         cl program program = clCreateProgramWithSource(af context,
514
             1, (const char **)&gc_source, &program_length,
515
             &status);
         status = clBuildProgram(program, 1, &af_device_id,
516
             NULL, NULL, NULL);
517
518
519
         char* kernelName;
520
         if (typef == f64)
             kernelName = "sparse_mat_vec_mul";
521
522
         else if (typef == f32)
             kernelName = "sparse mat vec mul sp";
523
524
525
         cl kernel kernel = clCreateKernel(program, kernelName,
526
             &status);
527
528
         cl_int step = 0;
529
         // 5.estableciendo los argumentos
530
         int i = 0;
531
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_elmA);
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_colA);
532
         clSetKernelArg(kernel, i++, sizeof(cl mem), d rowA);
533
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_b);
534
535
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_c);
536
         clSetKernelArg(kernel, i++, msize*localWorkSize, 0);
537
         clSetKernelArg(kernel, i++, sizeof(cl int), &size colA);
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
11
```

```
538
         clSetKernelArg(kernel, i++, sizeof(cl_int), &step);
539
540
        //6. ejecutando el kernel
541
542
         //transpose(L)*b
543
         clEnqueueNDRangeKernel(af_queue, kernel, 1, 0,
544
             &globalWorkSize, &localWorkSize, 0, NULL,
545
             NULL);
546
547
         //L*b
548
         step++;
         clSetKernelArg(kernel, 7, sizeof(cl_int), &step);
549
550
         clEnqueueNDRangeKernel(af queue, kernel, 1, 0,
551
             &globalWorkSize, &localWorkSize, 0, NULL,
             NULL);
552
553
554
         //D*b
555
         step++;
556
         clSetKernelArg(kernel, 7, sizeof(cl_int), &step);
557
         clEnqueueNDRangeKernel(af_queue, kernel, 1, 0,
558
             &globalWorkSize, &localWorkSize, 0, NULL,
559
             NULL);
560
561
         //7. devolviendo el control de memoria af::array a ArrayFire
         af_unlock_array(elmA);
562
         af_unlock_array(colA);
563
564
         af unlock array(rowA);
565
        af_unlock_array(dB);
566
        af_unlock_array(c);
567
         //copiando al argumento de salida
568
569
        af_copy_array(dC, c);
570
571
        af release array(c);
572
    }
573
    void AFire::sparse_sks_mat_vec_mul(af_array* dC,
574
575
         af_array elmA, af_array idxA, af_array dB) {
576
         //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
         //cl_context af_context;
577
         static cl context af context = afcl::getContext();
578
579
         static cl device id af device id = afcl::getDeviceId();
         static cl_command_queue af_queue = afcl::getQueue();
580
581
         //creando copia de dB
582
583
         af array c;
584
        af_copy_array(&c, dB);
585
586
        //2. Obteniendo parámetros necesarios
587
588
         //longitud de los vectores
         dim t order[AF MAX DIMS];
589
590
         af_get_dims(&_order[0], &_order[1], &_order[2],
591
             &_order[3], elmA);
592
         cl_int size_elmA = _order[0];
593
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
594
          af get dims(& order[0], & order[1], & order[2],
595
             &_order[3], idxA);
596
         cl_int size_idxA = _order[0];
597
598
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
599
         size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
600
         int status = CL_SUCCESS;
601
602
603
         af_dtype typef;
         af_get_type(&typef, elmA);
604
605
606
         int msize = 0;
607
          if (typef == f64)
             msize = sizeof(double);
608
609
          else if (typef == f32)
610
             msize = sizeof(float);
611
         else;
612
613
         //3.obteniendo las referencias cl mem de los objetos af::array
614
         cl_mem *d_elmA = (cl_mem*)clCreateBuffer(af_context,
             CL_MEM_READ_ONLY, msize*size_elmA,
615
616
             NULL, &status);
617
         af_get_device_ptr((void**)d_elmA, elmA);
618
         cl_mem *d_idxA = (cl_mem*)clCreateBuffer(af_context,
619
620
             CL_MEM_READ_ONLY, sizeof(int)*size_idxA,
621
             NULL, &status);
622
         af_get_device_ptr((void**)d_idxA, idxA);
623
         cl mem *d b = (cl mem*)clCreateBuffer(af context,
624
625
             CL_MEM_READ_ONLY, msize*size_idxA,
626
             NULL, &status);
627
         af get device ptr((void**)d b, dB);
628
         cl_mem *d_c = (cl_mem*)clCreateBuffer(af_context,
629
             CL_MEM_READ_WRITE, msize*size_idxA,
630
631
             NULL, &status);
632
         af_get_device_ptr((void**)d_c, c);
633
634
          size t program length = strlen(gc source);
635
          //4.creando el programa, construyendo el ejecutable y extrayendo el punto 🤊
636
            de entrada
637
          // para el Kernel
          cl_program program = clCreateProgramWithSource(af_context,
638
639
             1, (const char **)&gc_source, &program_length,
640
             &status);
641
          status = clBuildProgram(program, 1, &af_device_id,
642
             NULL, NULL, NULL);
643
         char* kernelName;
644
645
          if (typef == f64)
```

kernelName = "sparse_mat_vec_mul_sks";

kernelName = "sparse mat vec mul sks sp";

else if (typef == f32)

646

647

648

```
649
         else;
650
         cl kernel kernel = clCreateKernel(program, kernelName,
651
             &status);
652
         cl_int step = 0;
653
654
         // 5.estableciendo los argumentos
655
         int i = 0;
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_elmA);
656
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_idxA);
657
658
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_b);
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_c);
659
         clSetKernelArg(kernel, i++, msize*localWorkSize, 0);
660
661
         clSetKernelArg(kernel, i++, sizeof(cl_int), &size_idxA);
         clSetKernelArg(kernel, i++, sizeof(cl_int), &step);
662
663
        //6. ejecutando el kernel
664
665
         //U*b
666
667
         clEnqueueNDRangeKernel(af_queue, kernel, 1, 0,
             &globalWorkSize, &localWorkSize, 0, NULL,
668
669
             NULL);
670
         //L*b
671
672
         step++;
         clSetKernelArg(kernel, 6, sizeof(cl_int), &step);
673
         clEnqueueNDRangeKernel(af_queue, kernel, 1, 0,
674
675
             &globalWorkSize, &localWorkSize, 0, NULL,
676
             NULL);
677
678
         //7. devolviendo el control de memoria af::array a ArrayFire
679
         af unlock array(elmA);
680
         af_unlock_array(idxA);
681
682
         af unlock array(dB);
683
         af_unlock_array(c);
684
         //copiando al argumento de salida
685
686
         af_copy_array(dC, c);
687
         af_release_array(c);
688
689
690
691
    //Algebra Lineal
692 //----
693
694 void AFire::SELgj_f(af_array* dC, af_array dA, af_array dB) {
695
         //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
696
         //cl_context af_context;
697
         static cl_context af_context = afcl::getContext();
698
         static cl_device_id af_device_id = afcl::getDeviceId();
         static cl_command_queue af_queue = afcl::getQueue();
699
700
         //acoplando A con B
701
702
         af array Ac;
703
         af_join(&Ac, 1, dA, dB);
704
```

```
705
         //ya que el kernel asume la matriz como ordenada por filas
706
         //debemos transponer la matriz af::array A para obtener la
707
         //solución correcta, ya que ArrayFire ordena lo elementos
708
         //por columna
709
         af_array At;
710
         af_transpose(&At, Ac, false);
711
         dim_t _order[AF_MAX_DIMS];
712
         af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
713
714
         size_t order = _order[0];
715
         size t localWorkSize = BLOCK SIZE * BLOCK SIZE;
716
717
         size t globalWorkSize = localWorkSize * BLOCK SIZE;
718
719
         int status = CL_SUCCESS;
720
721
         af_dtype typef;
         af_get_type(&typef, dA);
722
723
724
        int msize = 0;
725
         if (typef == f64)
             msize = sizeof(double);
726
         else if (typef == f32)
727
728
            msize = sizeof(float);
729
        else;
730
731
         //obteniendo las referencias cl mem de los objetos af::array
         cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
732
733
             CL_MEM_READ_WRITE, msize * order*(order + 1),
734
             NULL, &status);
735
        af_get_device_ptr((void**)d_A, At);
736
        size t program length = strlen(GJordan source);
737
738
739
740
         //creando el programa, construyendo el ejecutable y extrayendo el punto
           de entrada
741
         // para el Kernel
742
         cl_program program = clCreateProgramWithSource(af_context, 1, (const char →
            **)&GJordan_source, &program_length, &status);
         status = clBuildProgram(program, 1, &af device id, NULL, NULL, NULL);
743
744
745
         char* kernelName;
746
         if (typef == f64)
             kernelName = "Gauss Jordan f";
747
748
         else if (typef == f32)
749
             kernelName = "Gauss_Jordan_f_sp";
750
         else;
751
         cl_kernel kernel = clCreateKernel(program, kernelName,
752
             &status);
753
         // estableciendo los argumentos
754
755
         int i = 0;
756
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_A);
757
         clSetKernelArg(kernel, i++, sizeof(cl_int), &order);
758
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
759
          for (int j = 0; j < order; j++)</pre>
760
              clSetKernelArg(kernel, 2, sizeof(cl_int), &j);
761
762
              //ejecutando el Kernel
763
              clEnqueueNDRangeKernel(af_queue, kernel, 1, 0, &globalWorkSize,
                &localWorkSize,
764
                  0, NULL, NULL);
765
766
767
          //devolviendo el control de memoria af::array a ArrayFire
         af_unlock_array(At);
768
769
770
         //hasta aqui At contiene en su última fila
771
772
         //y en su diagonal principal, los valores
         //finales que deben dividirse para obtener
773
774
          //la solución final
775
776
         //extrayendo la última fila
777
         af release array(Ac);
         af index t* indexers = 0;
 778
         af create indexers(&indexers);
779
         af set seq param indexer(indexers, order, order, 1,
780
 781
              0, false);
         af_set_seq_param_indexer(indexers, 0, order - 1, 1,
782
783
              1, false);
784
         af_index_gen(&Ac, At, 2, indexers);
785
786
         //transponiendo
         af_array Atr;
787
         af_transpose(&Atr, Ac, false);
788
 789
790
         //extrayendo la diagonal
791
         af release array(Ac);
792
         af_diag_extract(&Ac, At, 0);
793
```

```
//dividiendo
794
795
         af release array(At);
796
         af_div(&At, Atr, Ac, false);
797
798
         // copiando el resultado en dC
799
         af_copy_array(dC, At);
800
801
         af_release_array(Atr);
802
         af release array(Ac);
803
         af_release_array(At);
804
         af_release_indexers(indexers);
805
806
807
    void AFire::SELgj_c(af_array* dC, af_array dA, af_array dB) {
         //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
808
         //cl context af context;
809
810
         static cl_context af_context = afcl::getContext();
811
         static cl_device_id af_device_id = afcl::getDeviceId();
         static cl_command_queue af_queue = afcl::getQueue();
812
813
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
814
         //acoplando A con B
815
         af array Ac;
816
         af_join(&Ac, 1, dA, dB);
817
         dim_t _order[AF_MAX_DIMS];
818
819
         af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
820
         size_t order = _order[0];
821
         size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
822
823
         size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
824
825
         int status = CL_SUCCESS;
826
827
         af_dtype typef;
828
         af_get_type(&typef, dA);
829
830
         int msize = 0;
         if (typef == f64)
831
832
             msize = sizeof(double);
833
         else if (typef == f32)
834
             msize = sizeof(float);
835
         else;
836
837
         //obteniendo las referencias cl mem de los objetos af::array
         cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
838
             CL_MEM_READ_WRITE, msize * order*(order + 1),
839
840
             NULL, &status);
841
         af_get_device_ptr((void**)d_A, Ac);
842
         size_t program_length = strlen(GJordan_source);
843
844
845
         //creando el programa, construyendo el ejecutable y extrayendo el punto
           de entrada
846
         // para el Kernel
847
         cl program program =
             clCreateProgramWithSource(af_context, 1,
848
             (const char **)&GJordan_source, &program_length,
849
850
                 &status);
851
         status = clBuildProgram(program, 1, &af_device_id,
             NULL, NULL, NULL);
852
853
         char* kernelName;
854
855
         if (typef == f64)
             kernelName = "Gauss_Jordan_c";
856
857
         else if (typef == f32)
             kernelName = "Gauss_Jordan_c_sp";
858
859
860
         cl_kernel kernel = clCreateKernel(program, kernelName,
861
             &status);
862
         // estableciendo los argumentos
863
864
         int i = 0;
         clSetKernelArg(kernel, i++, sizeof(cl_mem), d_A);
865
866
         clSetKernelArg(kernel, i++, sizeof(cl_int), &order);
867
868
         for (int j = 0; j < order; j++)
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
17
```

```
869
             clSetKernelArg(kernel, 2, sizeof(cl_int), &j);
870
871
             //ejecutando el Kernel
             clEnqueueNDRangeKernel(af_queue, kernel, 1, 0,
872
                 &globalWorkSize, &localWorkSize, 0, NULL,
873
874
                 NULL);
875
         }
876
         //devolviendo el control de memoria af::array a ArrayFire
877
878
         af_unlock_array(Ac);
879
         //hasta aqui Ac contiene en su última columna
880
881
         //y en su diagonal principal, los valores
         //finales que deben dividirse para obtener
882
883
         //la solución final
884
885
         //extrayendo la última columna
         af_array Au;
886
887
         af_index_t* indexers = 0;
888
         af_create_indexers(&indexers);
889
         af_set_seq_param_indexer(indexers, 0, order - 1, 1,
890
             0, false);
         af set seq param indexer(indexers, order, order, 1,
891
892
             1, false);
         af_index_gen(&Au, Ac, 2, indexers);
893
894
895
         //extrayendo la diagonal
896
         af array Ad;
897
         af_diag_extract(&Ad, Ac, 0);
898
         //dividiendo
899
900
         af_release_array(Ac);
901
         af_div(&Ac, Au, Ad, false);
902
         // copiando el resultado en dC
903
         af_copy_array(dC, Ac);
904
905
906
         af release array(Au);
907
         af_release_array(Ad);
908
         af_release_array(Ac);
909
         af release indexers(indexers);
910 }
911
912 void AFire::SELgj_fshr(af_array* dC, af_array dA,
913
         af array dB) {
914
         //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
915
         static cl_context af_context = afcl::getContext();
916
         static cl_device_id af_device_id = afcl::getDeviceId();
917
         static cl_command_queue af_queue = afcl::getQueue();
918
919
         /*creando copia de B*/
         af_array Bc;
920
921
         af_copy_array(&Bc, dB);
922
923
         /*ya que el kernel asume la matriz como ordenada por filas
924
         debemos transponer la matriz af::array A para obtener la
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
18
```

```
925
         solución correcta, ya que ArrayFire ordena lo elementos
926
         por columna*/
927
         af_array At;
928
         af_transpose(&At, dA, false);
929
         //para almacenar el resultado
930
931
         af_array Rs;
932
         af_copy_array(&Rs, dB);
933
934
         dim_t Adims[AF_MAX_DIMS];
935
         af_get_dims(&Adims[0], &Adims[1], &Adims[2], &Adims[3],
936
             dA);
937
         size t order = Adims[0];
938
939
         af_dtype typef;
         af_get_type(&typef, dA);
940
941
942
         int msize = 0;
943
         if (typef == f64)
944
             msize = sizeof(double);
945
         else if (typef == f32)
946
             msize = sizeof(float);
947
         else;
948
949
         size_t program_length = strlen(GJordan_source);
         int status = CL_SUCCESS;
950
951
         //obteniendo las referencias cl_mem de los objetos af::array
952
953
         cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
             CL_MEM_READ_WRITE, msize * order * order,
954
955
             NULL, &status);
956
         af_get_device_ptr((void**)d_A, At);
957
958
         cl mem *d B = (cl mem*)clCreateBuffer(af context,
959
             CL MEM READ WRITE, msize * order, NULL, &status);
         af_get_device_ptr((void**)d_B, Bc);
960
961
962
         cl mem *d C = (cl mem*)clCreateBuffer(af context,
963
             CL_MEM_WRITE_ONLY, msize * order, NULL, &status);
964
         af_get_device_ptr((void**)d_C, Rs);
965
         //creando el programa, construyendo el ejecutable y extrayendo el punto
966
           de entrada
967
         // para el Kernel
         cl program program = clCreateProgramWithSource(af context, 1, (const char →
968
            **)&GJordan_source, &program_length, &status);
969
         status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
970
971
         char* kernelName;
972
         if (typef == f64)
             kernelName = "Gauss_Jordan_fshr";
973
974
         else if (typef == f32)
975
             kernelName = "Gauss_Jordan_fshr_sp";
976
977
         cl_kernel kernel = clCreateKernel(program, kernelName,
978
             &status);
```

```
979
 980
          // estableciendo los argumentos
 981
          int i = 0;
 982
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_C);
 983
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_A);
 984
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_B);
 985
          clSetKernelArg(kernel, i++, sizeof(double)*BLOCK_SIZE*BLOCK_SIZE, 0);
          clSetKernelArg(kernel, i++, sizeof(cl_int), &order);
 986
 987
 988
          size_t localWorkSize[] = { BLOCK_SIZE,BLOCK_SIZE };
 989
          size_t globalWorkSize[] =
 990
          { shrRoundUp(localWorkSize[0],order),
 991
              shrRoundUp(localWorkSize[1], order) };
 992
          //ejecutando el Kernel
 993
          clEnqueueNDRangeKernel(af_queue, kernel, 2, 0,
 994
 995
              globalWorkSize, localWorkSize,0, NULL, NULL);
 996
 997
          //devolviendo el control de memoria af::array a ArrayFire
 998
          af unlock array(At);
 999
          af unlock array(Bc);
1000
          af_unlock_array(Rs);
1001
1002
          //copiando el resultado
1003
          af_copy_array(dC, Rs);
1004
1005
          af release array(At);
1006
          af_release_array(Bc);
1007
          af_release_array(Rs);
1008
1009
1010
      void AFire::SELgj_f2d(af_array* dC, af_array dA,
          af array dB) {
1011
1012
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
1013
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
1014
          static cl_command_queue af_queue = afcl::getQueue();
1015
1016
1017
          /*creando copia de B*/
1018
          af_array Bc;
1019
          af copy array(&Bc, dB);
1020
1021
          //para almacenar el resultado
1022
          af array Rs;
1023
          af_copy_array(&Rs, dB);
1024
          /*ya que el kernel asume la matriz como ordenada por filas
1025
1026
          debemos transponer la matriz af::array A para obtener la
1027
          solución correcta, ya que ArrayFire ordena lo elementos
1028
          por columna*/
1029
          af_array At;
1030
          af transpose(&At, dA, false);
1031
1032
          dim t Adims[AF MAX DIMS];
          af_get_dims(&Adims[0], &Adims[1], &Adims[2], &Adims[3],
1033
1034
              dA);
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
20
```

```
1035
          size t order = Adims[0];
1036
1037
          af dtype typef;
1038
          af_get_type(&typef, dA);
1039
1040
          int msize = 0;
          if (typef == f64)
1041
              msize = sizeof(double);
1042
1043
          else if (typef == f32)
1044
              msize = sizeof(float);
1045
          else;
1046
1047
          size t program length = strlen(GJordan source);
1048
          int status = CL_SUCCESS;
1049
          //obteniendo las referencias cl_mem de los objetos af::array
1050
1051
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
1052
              CL_MEM_READ_WRITE, msize * order * order,
1053
              NULL, &status);
1054
          af_get_device_ptr((void**)d_A, At);
1055
          cl mem *d B = (cl mem*)clCreateBuffer(af context,
1056
              CL_MEM_READ_WRITE, msize * order, NULL, &status);
1057
          af_get_device_ptr((void**)d_B, Bc);
1058
1059
          cl_mem *d_C = (cl_mem*)clCreateBuffer(af_context,
1060
1061
              CL_MEM_WRITE_ONLY, msize * order, NULL, &status);
1062
          af_get_device_ptr((void**)d_C, Rs);
1063
          //creando el programa, construyendo el ejecutable y extrayendo el punto
1064
            de entrada
1065
          // para el Kernel
1066
          cl program program = clCreateProgramWithSource(af context, 1, (const char →
             **)&GJordan source, &program length, &status);
          status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
1067
1068
1069
          char* kernelName;
1070
          if (typef == f64)
1071
              kernelName = "Gauss_Jordan_f2d";
          else if (typef == f32)
1072
1073
              kernelName = "Gauss Jordan f2d sp";
1074
          cl_kernel kernel = clCreateKernel(program, kernelName,
1075
1076
              &status);
1077
          // estableciendo los argumentos
1078
1079
          int i = 0;
1080
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_C);
1081
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_A);
1082
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_B);
          clSetKernelArg(kernel, i++, sizeof(cl_int), &order);
1083
1084
1085
          size t localWorkSize[] = { BLOCK SIZE, BLOCK SIZE };
1086
          size_t globalWorkSize[] = { BLOCK_SIZE,BLOCK_SIZE };
1087
1088
          //ejecutando el Kernel
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
1089
          clEnqueueNDRangeKernel(af_queue, kernel, 2, 0,
1090
              globalWorkSize, localWorkSize, 0, NULL, NULL);
1091
1092
          //devolviendo el control de memoria af::array a ArrayFire
          af_unlock_array(At);
1093
1094
          af unlock array(Bc);
1095
          af_unlock_array(Rs);
1096
1097
          //copiando el resultado
1098
          af_copy_array(dC, Rs);
1099
1100
          af_release_array(At);
1101
          af release array(Bc);
1102
          af_release_array(Rs);
1103 }
1104
1105
      void AFire::SELgj_c2d(af_array* dC, af_array dA,
          af_array dB) {
1106
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
1107
1108
          static cl_context af_context = afcl::getContext();
1109
          static cl_device_id af_device_id = afcl::getDeviceId();
          static cl_command_queue af_queue = afcl::getQueue();
1110
1111
1112
          /*creando copia de B*/
1113
          af_array Bc;
          af_copy_array(&Bc, dB);
1114
1115
          //para almacenar el resultado
1116
1117
          af_array Rs;
          af_copy_array(&Rs, dB);
1118
1119
1120
          /*creando copia de A*/
1121
          af array At;
1122
          af_copy_array(&At, dA);
1123
1124
          dim t Adims[AF MAX DIMS];
          af_get_dims(&Adims[0], &Adims[1], &Adims[2], &Adims[3],
1125
1126
              dA);
1127
          size_t order = Adims[0];
1128
1129
          af dtype typef;
          af_get_type(&typef, dA);
1130
1131
1132
          int msize = 0;
1133
          if (typef == f64)
              msize = sizeof(double);
1134
1135
          else if (typef == f32)
1136
              msize = sizeof(float);
1137
          else;
1138
          size_t program_length = strlen(GJordan_source);
1139
          int status = CL SUCCESS;
1140
1141
1142
          //obteniendo las referencias cl_mem de los objetos af::array
1143
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
1144
              CL MEM READ WRITE, msize * order * order,
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
22
```

```
1145
              NULL, &status);
1146
          af_get_device_ptr((void**)d_A, At);
1147
1148
          cl_mem *d_B = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, msize * order, NULL, &status);
1149
1150
          af_get_device_ptr((void**)d_B, Bc);
1151
          cl_mem *d_C = (cl_mem*)clCreateBuffer(af_context,
1152
              CL_MEM_WRITE_ONLY, msize * order, NULL, &status);
1153
          af get_device_ptr((void**)d_C, Rs);
1154
1155
1156
          //creando el programa, construyendo el ejecutable y extrayendo el punto
            de entrada
1157
          // para el Kernel
1158
          cl_program program = clCreateProgramWithSource(af_context, 1, (const char →
             **)&GJordan_source, &program_length, &status);
1159
          status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
1160
1161
          char* kernelName;
1162
          if (typef == f64)
              kernelName = "Gauss Jordan c2d";
1163
          else if (typef == f32)
1164
              kernelName = "Gauss_Jordan_c2d_sp";
1165
1166
          else:
          cl_kernel kernel = clCreateKernel(program, kernelName,
1167
1168
              &status);
1169
          // estableciendo los argumentos
1170
1171
          int i = 0;
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_C);
1172
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_A);
1173
1174
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_B);
1175
          clSetKernelArg(kernel, i++, sizeof(cl_int), &order);
1176
          size t localWorkSize[] = { BLOCK SIZE, BLOCK SIZE };
1177
          size_t globalWorkSize[] = { BLOCK_SIZE,BLOCK_SIZE };
1178
1179
1180
          //ejecutando el Kernel
1181
          clEnqueueNDRangeKernel(af_queue, kernel, 2, 0,
              globalWorkSize, localWorkSize, 0, NULL, NULL);
1182
1183
          //devolviendo el control de memoria af::array a ArrayFire
1184
          af unlock array(At);
1185
          af_unlock_array(Bc);
1186
1187
          af_unlock_array(Rs);
1188
1189
          //copiando el resultado
1190
          af_copy_array(dC, Rs);
1191
1192
          af_release_array(At);
1193
          af_release_array(Bc);
1194
          af release array(Rs);
1195
      }
1196
      void AFire::prueba_shr(af_array* dC, af_array dA,
1197
1198
          af array dB) {
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
23
```

```
1199
          //Obteniendo el dispositivo, contexto y la cola usada
1200
          //por ArrayFire
1201
          static cl_context af_context = afcl::getContext();
1202
          static cl_device_id af_device_id = afcl::getDeviceId();
1203
          static cl_command_queue af_queue = afcl::getQueue();
1204
1205
          dim_t Adims[AF_MAX_DIMS];
          af_get_dims(&Adims[0], &Adims[1], &Adims[2], &Adims[3], dA);
1206
1207
          size_t order = Adims[0];
1208
1209
          af_dtype typef;
1210
          af_get_type(&typef, dA);
1211
1212
          int msize = 0:
1213
          if (typef == f64)
              msize = sizeof(double);
1214
1215
          else if (typef == f32)
1216
              msize = sizeof(float);
1217
          else;
1218
1219
          size t program length = strlen(GJordan source);
          int status = CL SUCCESS;
1220
1221
1222
          //creando copia de dA y dB, para poder modificarlos
1223
          af_array Bc;
          af_copy_array(&Bc, dB);
1224
1225
          af array Ac;
1226
          af_copy_array(&Ac, dA);
1227
          //para almacenar el resultado
1228
1229
          af array Rs;
1230
          af_copy_array(&Rs, dB);
1231
1232
          //obteniendo las referencias cl mem de los objetos af::array
          cl mem *d A = (cl mem*)clCreateBuffer(af context,
1233
1234
              CL_MEM_READ_WRITE, msize * order * order,
1235
              NULL, &status);
1236
          af_get_device_ptr((void**)d_A, Ac);
1237
1238
          cl_mem *d_B = (cl_mem*)clCreateBuffer(af_context,
1239
              CL MEM READ WRITE, msize * order, NULL, &status);
          af get device ptr((void**)d B, Bc);
1240
1241
1242
          cl mem *d C = (cl mem*)clCreateBuffer(af context,
1243
              CL_MEM_WRITE_ONLY, msize * order, NULL, &status);
1244
          af_get_device_ptr((void**)d_C, Rs);
1245
1246
1247
          //creando el programa, construyendo el ejecutable y extrayendo el punto
            de entrada
1248
          // para el Kernel
          cl program program = clCreateProgramWithSource(af context, 1, (const char →
1249
             **)&GJordan_source, &program_length, &status);
1250
          status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
1251
1252
          char* kernelName;
```

```
...sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
24
```

```
1253
          if (typef == f64)
              kernelName = "prueba shr";
1254
1255
          else if (typef == f32)
              kernelName = "prueba_shr_sp";
1256
1257
          else;
          cl_kernel kernel = clCreateKernel(program, kernelName,
1258
1259
              &status);
1260
          // estableciendo los argumentos
1261
1262
          int i = 0;
1263
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_C);
1264
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_A);
1265
          clSetKernelArg(kernel, i++, sizeof(cl_mem), d_B);
          clSetKernelArg(kernel, i++, sizeof(double)*BLOCK_SIZE*BLOCK_SIZE, 0);
1266
1267
          clSetKernelArg(kernel, i++, sizeof(cl_int), &order);
1268
1269
          size_t localWorkSize[] = { BLOCK_SIZE, BLOCK_SIZE };
1270
          size_t globalWorkSize[] =
1271
              { shrRoundUp(localWorkSize[0], order), shrRoundUp(localWorkSize[1],
                order) };
1272
          //ejecutando el Kernel
          clEnqueueNDRangeKernel(af queue, kernel, 2, 0,
1273
              globalWorkSize, localWorkSize, 0, NULL, NULL);
1274
1275
          //devolviendo el control de memoria af::array a ArrayFire
1276
          af unlock array(Ac);
1277
1278
          af unlock array(Bc);
1279
          af_unlock_array(Rs);
1280
          //copiando el resultado
1281
1282
          af copy array(dC, Rs);
1283
          af release array(Ac);
1284
1285
          af release array(Bc);
1286
          af_release_array(Rs);
1287
1288
1289
      void AFire::SEL_gc(af_array* C, af_array A, af_array
1290
          double lerr) {
1291
          dim_t _order[AF_MAX_DIMS];
1292
          af_get_dims(&_order[0], &_order[1], &_order[2],
1293
1294
              &_order[3], A);
1295
          size_t order = _order[0];
1296
1297
          af_dtype typef;
1298
          af_get_type(&typef, A);
1299
1300
          double normr;
1301
          //af_array de ayuda
          af_array zero;
1302
          dim_t d_order[] = { 1 };
1303
1304
          af_constant(&zero, 0, 1, d_order, typef);
1305
          af_array Ax0;
1306
          af_array rtxp;
1307
          af_array ptxz;
```

```
1308
           af_array axp;
1309
          af_array B;
1310
          af_array axz;
1311
          af_array copyr;
1312
          af_array rtxz;
1313
          af_array rsp;
1314
          af_array Bxp;
1315
1316
          //x0=b
1317
          af_array x0;
1318
          af_copy_array(&x0,
                               );
1319
          //r=b-A*x0
1320
1321
          af_array r;
          af_matmul(&AxO, A, xO, AF_MAT_NONE, AF_MAT_NONE);
1322
1323
          af_sub(&r, , Ax0, false);
1324
1325
          //p = r
1326
          af_array p;
1327
          af_copy_array(&p, r);
1328
1329
          //z=A*p
          af array z;
1330
1331
          af_matmul (&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
1332
          //a = (r' *p)/(p' *z)
1333
1334
          af_array a;
1335
          af_matmul (&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
1336
          af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
1337
          af_div(&a, rtxp, ptxz, false);
1338
1339
          //x = x0 + a.*p
1340
          af_array x;
1341
          af_mul(&axp, a, p, true);
1342
          af_add(&x, x0, axp, false);
1343
          int contin;
1344
           int I = 0;
1345
1346
           for (int i = 0; i < order; i ++)
1347
               //r -= a. *z
1348
               af_copy_array(&copyr, r);
1349
               af_mul(&axz, a, z, true);
1350
               af_sub(&r, copyr, axz, false);
1351
1352
               af_norm(&normr, r, AF_NORM_EUCLID, 1, 1);
1353
               if (normr <= lerr)</pre>
1354
1355
                   break;
1356
1357
               //B = -(r' *z)/(p' *z)
1358
               af_matmul(&rtxz, r, z, AF_MAT_TRANS, AF_MAT_NONE);
               af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
1359
1360
               af_div(&rsp, rtxz, ptxz, false);
1361
               af_sub(&B, zero, rsp, false);
1362
               //p = r + B.*p
1363
```

```
1364
               af_mul(&Bxp, B, p, true);
 1365
               af_add(&p, r, Bxp, false);
 1366
               //z = A*p
 1367
1368
               af_matmul(&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
1369
               //a = (r' *p)/(p' *z)
1370
               af_matmul(&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
1371
               af_matmul(&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
 1372
 1373
               af_div(&a, rtxp, ptxz, false);
 1374
               //x += a.*p
1375
1376
               af_mul(&axp, a, p, true);
 1377
               af_copy_array(&x0, x);
 1378
               af_add(&x, x0, axp, false);
 1379
 1380
               I + +;
 1381
           }
 1382
 1383
           //copi ando el resultado en el argumento de salida
 1384
           af_copy_array(C, x);
 1385
           //liberando objetos af_array usados
 1386
 1387
           af_rel ease_array(zero);
 1388
           af_rel ease_array(Ax0);
           af_rel ease_array(rtxp);
 1389
 1390
           af_rel ease_array(ptxz);
           af_rel ease_array(axp);
 1391
 1392
           af_rel ease_array(B);
 1393
           af_rel ease_array(axz);
 1394
           af_rel ease_array(copyr);
 1395
           af_release_array(rtxz);
1396
           af_release_array(rsp);
 1397
           af_release_array(Bxp);
 1398
           af_release_array(x0);
1399
           af_release_array(r);
 1400
           af_rel ease_array(p);
 1401
           af_release_array(z);
 1402
           af_release_array(a);
 1403
           af_release_array(x);
1404 }
 1405
      void AFire::SELgc_sparse(af_array* C, af_array elmA,
1406
           af_array col A, af_array rowA, af_array ,
 1407
 1408
           double lerr) {
 1409
           dim_t _order[AF_MAX_DIMS];
 1410
           af_get_dims(&_order[0], &_order[1], &_order[2],
1411
1412
               &_order[3], col A);
1413
           size_t order = _order[0];
1414
           af_dtype typef;
 1415
 1416
           af_get_type(&typef, elmA);
 1417
 1418
           double normr;
           //af_array de ayuda
 1419
```

```
1420
           af_array zero;
           dim_t d_order[] = { 1 };
 1421
 1422
           af_constant(&zero, 0, 1, d_order, typef);
 1423
           af_array Ax0;
1424
           af_array rtxp;
1425
           af_array ptxz;
 1426
           af_array axp;
 1427
           af_array B;
 1428
           af_array axz;
 1429
           af_array copyr;
 1430
           af_array rtxz;
 1431
           af_array rsp;
1432
           af_array Bxp;
 1433
 1434
           //x0=b
           af_array x0;
 1435
 1436
           af_copy_array(&x0,
 1437
 1438
           //r=b-A*x0
 1439
           af_array r;
 1440
           //af_matmul (&AxO, A, xO, AF_MAT_NONE, AF_MAT_NONE);
 1441
           AFire::sparse_mat_vec_mul(&Ax0, elmA, colA, rowA, x0);
           af_sub(&r, , Ax0, false);
 1442
 1443
           //p = r
 1444
 1445
           af_array p;
 1446
           af_copy_array(&p, r);
 1447
 1448
           //z=A*p
 1449
           af_array z;
           //af_matmul (&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
 1450
 1451
           AFire::sparse_mat_vec_mul(&z, elmA, colA, rowA, p);
 1452
 1453
           //a = (r' *p)/(p' *z)
 1454
           af array a;
           af_matmul (&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
 1455
           af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
 1456
 1457
           af_div(&a, rtxp, ptxz, false);
 1458
 1459
           //x = x0 + a.*p
1460
           af_array x;
 1461
           af_mul(&axp, a, p, true);
1462
           af_add(&x, x0, axp, false);
1463
 1464
           int contin;
 1465
           int I = 0;
           for (int i = 0; i < order; i++)
 1466
 1467
1468
               //r -= a. *z
 1469
               af_copy_array(&copyr, r);
 1470
               af_mul(&axz, a, z, true);
 1471
               af_sub(&r, copyr, axz, false);
 1472
               af_norm(&normr, r, AF_NORM_EUCLID, 1, 1);
 1473
               if (normr <= lerr)</pre>
 1474
                   break;
 1475
```

```
1476
1477
              //B = -(r' *z)/(p' *z)
              af_matmul (&rtxz, r, z, AF_MAT_TRANS, AF_MAT_NONE);
1478
1479
              af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
1480
              af_div(&rsp, rtxz, ptxz, false);
1481
              af_sub(&B, zero, rsp, false);
1482
              //p = r + B.*p
1483
              af_mul(&Bxp, B, p, true);
1484
1485
              af_add(&p, r, Bxp, false);
1486
              //z = A*p
1487
              //af_matmul(&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
1488
1489
              AFire::sparse_mat_vec_mul(&z, elmA, colA, rowA, p);
1490
              //a = (r' *p)/(p' *z)
1491
              af_matmul (&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
1492
1493
              af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
              af_div(&a, rtxp, ptxz, false);
1494
1495
1496
              //x += a.*p
1497
              af_mul(&axp, a, p, true);
1498
              af_copy_array(&x0, x);
1499
              af_add(&x, x0, axp, false);
1500
1501
              1++;
1502
          }
1503
1504
          //copi ando el resultado en el argumento de salida
1505
          af_copy_array(C, x);
1506
1507
          //liberando objetos af_array usados
          af_release_array(zero);
1508
1509
          af_release_array(Ax0);
1510
          af_release_array(rtxp);
1511
          af_rel ease_array(ptxz);
1512
          af_release_array(axp);
1513
          af_rel ease_array(B);
1514
          af_rel ease_array(axz);
1515
          af_rel ease_array(copyr);
1516
          af_release_array(rtxz);
1517
          af_release_array(rsp);
1518
          af_release_array(Bxp);
1519
          af_release_array(x0);
1520
          af_release_array(r);
1521
          af_release_array(p);
1522
          af_rel ease_array(z);
1523
          af_rel ease_array(a);
1524
          af_release_array(x);
1525
     }
1526
     void AFire::SELgc_sparse_sks(af_array* C, af_array elmA,
1527
          af_array idxA, af_array , double lerr) {
1528
1529
1530
          dim_t _order[AF_MAX_DIMS];
          af_get_dims(&_order[0], &_order[1], &_order[2],
1531
```

```
&_order[3], idxA);
1532
1533
          size_t order = _order[0];
1534
1535
          af_dtype typef;
1536
          af_get_type(&typef, elmA);
1537
1538
          double normr;
1539
          //af_array de ayuda
1540
          af_array zero;
1541
          dim_t d_order[] = { 1 };
1542
          af_constant(&zero, 0, 1, d_order, typef);
1543
          af_array Ax0;
1544
          af_array rtxp;
1545
          af_array ptxz;
1546
          af_array axp;
1547
          af_array B;
1548
          af_array axz;
1549
          af_array copyr;
1550
          af_array rtxz;
1551
          af_array rsp;
1552
          af_array Bxp;
1553
          //x0=b
1554
1555
          af array x0;
1556
          af_copy_array(&x0, );
1557
          //r=b-A*x0
1558
1559
          af_array r;
          //af_matmul (&AxO, A, xO, AF_MAT_NONE, AF_MAT_NONE);
1560
1561
          AFire::sparse_sks_mat_vec_mul(&Ax0, elmA, idxA, x0);
          af_sub(&r, , Ax0, false);
1562
1563
1564
          //p = r
1565
          af_array p;
1566
          af_copy_array(&p, r);
1567
1568
          //z=A*p
1569
          af_array z;
1570
          //af_matmul (&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
1571
          AFire::sparse_sks_mat_vec_mul(&z, elmA, idxA, p);
1572
          //a = (r' *p)/(p' *z)
1573
1574
          af_array a;
1575
          af_matmul (&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
          af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
1576
1577
          af_div(&a, rtxp, ptxz, false);
1578
          //x = x0 + a.*p
1579
1580
          af_array x;
1581
          af_mul(&axp, a, p, true);
1582
          af_add(&x, x0, axp, false);
1583
          int contin;
1584
1585
          int I = 0;
1586
          for (int i = 0; i < order; i + +)
1587
```

```
1588
               //r -= a. *z
1589
               af_copy_array(&copyr, r);
1590
               af_mul(&axz, a, z, true);
1591
               af_sub(&r, copyr, axz, false);
1592
               af_norm(&normr, r, AF_NORM_EUCLID, 1, 1);
1593
1594
               if (normr <= lerr)</pre>
1595
                   break;
1596
               //B = -(r' *z)/(p' *z)
1597
               af_matmul(&rtxz, r, z, AF_MAT_TRANS, AF_MAT_NONE);
1598
               af_matmul (&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
1599
               af_div(&rsp, rtxz, ptxz, false);
1600
               af_sub(&B, zero, rsp, false);
1601
1602
               //p = r + B.*p
1603
               af_mul(&Bxp, B, p, true);
1604
1605
               af_add(&p, r, Bxp, false);
1606
1607
               //z = A*p
1608
               //af_matmul (&z, A, p, AF_MAT_NONE, AF_MAT_NONE);
1609
               AFire::sparse_sks_mat_vec_mul(&z, elmA, idxA, p);
1610
               //a = (r' *p)/(p' *z)
1611
1612
               af_matmul(&rtxp, r, p, AF_MAT_TRANS, AF_MAT_NONE);
               af_matmul(&ptxz, p, z, AF_MAT_TRANS, AF_MAT_NONE);
1613
1614
               af_div(&a, rtxp, ptxz, false);
1615
1616
               //x += a.*p
1617
               af_mul(&axp, a, p, true);
1618
               af_copy_array(&x0, x);
1619
               af_add(&x, x0, axp, false);
1620
1621
               1++;
1622
          }
1623
1624
          //copi ando el resultado en el argumento de salida
1625
          af_copy_array(C, x);
1626
1627
          //liberando objetos af_array usados
1628
          af_rel ease_array(zero);
1629
          af release array(Ax0);
1630
          af_release_array(rtxp);
1631
          af_release_array(ptxz);
1632
          af_release_array(axp);
1633
          af_release_array(B);
1634
          af_rel ease_array(axz);
1635
          af_rel ease_array(copyr);
1636
          af_release_array(rtxz);
1637
          af_release_array(rsp);
1638
          af_rel ease_array(Bxp);
1639
          af_release_array(x0);
1640
          af_release_array(r);
1641
          af_release_array(p);
1642
          af_release_array(z);
1643
          af_release_array(a);
```

```
1644
          af_release_array(x);
1645 }
1646
1647 array AFire:: SEL_gc(array A, array , double lerr) {
1648
1649
          bool gfor_stat = gforGet();
1650
          gforSet(true);
          si ze_t order = A. di ms(0);
1651
1652
          array x0 = ;
1653
1654
          array r = b - matmul(A, x0);
1655
          array p = r;
1656
          array z = matmul(A, p);
          array a = matmul (r, p, AF_MAT_TRANS) /
1657
1658
             matmul (p, z, AF_MAT_TRANS);
1659
          array x = x0 + a * p;
1660
          array B;
1661
1662
          int I = 0;
1663
          for (int i = 0; i < order; i ++)
1664
1665
              r -= a * z;
              if (norm(r) <= lerr)</pre>
1666
1667
                  break;
1668
              B = -matmul(r, z, AF\_MAT\_TRANS) /
                  matmul (p, z, AF_MAT_TRANS);
1669
1670
              p = r + B * p;
              z = matmul(A, p);
1671
1672
              a = matmul(r, p, AF_MAT_TRANS) /
1673
                  matmul (p, z, AF_MAT_TRANS);
              x += a * p;
1674
1675
              1++;
1676
1677
          gforSet(gfor_stat);
1678
          return x;
1679
      }
1680
      void AFire::SELchol_c(af_array* dC, af_array dA, af_array dB) {
1681
1682
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
1683
          //cl_context af_context;
          static cl_context af_context = afcl::getContext();
1684
1685
          static cl device id af device id = afcl::getDeviceId();
1686
          static cl_command_queue af_queue = afcl::getQueue();
1687
          //copi a de A y B
1688
          af_array Ac;
1689
1690
          af_array Bc;
1691
          af_copy_array(&Ac, dA);
1692
          af_copy_array(&Bc, dB);
1693
1694
          dim_t _order[AF_MAX_DIMS];
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
1695
1696
          size_t order = _order[0];
1697
1698
          size_t local WorkSize = BLOCK_SIZE * BLOCK_SIZE;
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
1699
```

```
1700
1701
          int status = CL_SUCCESS;
1702
1703
          af_dtype typef;
1704
          af_get_type(&typef, dA);
1705
          int msize = 0;
1706
1707
          if (typef == f64)
1708
              msi ze = si zeof(double);
1709
          else if (typef == f32)
1710
              msi ze = si zeof(float);
1711
          el se;
1712
1713
          //obteniendo las referencias cl_mem de los objetos af::array
1714
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
1715
              CL_MEM_READ_WRITE, msize*order*order,
1716
              NULL, &status);
          af_get_device_ptr((void**)d_A, Ac);
1717
1718
1719
          cl_mem *d_B = (cl_mem*)clCreateBuffer(af_context,
1720
              CL_MEM_READ_WRITE, msi ze*order,
1721
              NULL, &status);
          af_get_device_ptr((voi d**)d_B, Bc);
1722
1723
1724
          size_t program_l ength = strl en(Chol esky_source);
1725
1726
          //creando el programa, construyendo el ejecutable y extrayendo el punto
1727
            de entrada
1728
          // para el Kernel
          cl_program program = cl CreateProgramWi thSource(af_context, 1, (const char >
1729
              **)&Cholesky_source, &program_length, &status);
          status = cl BuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
1730
1731
1732
          char* kernel Name;
          if (typef == f64)
1733
              kernel Name = "Chol esky_c";
1734
          else if (typef == f32)
1735
              kernel Name = "Chol esky_c_sp";
1736
1737
          el se;
          cl_kernel kernel = clCreateKernel(program, kernelName,
1738
1739
              &status);
1740
1741
          // estableciendo los argumentos
1742
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_A);
1743
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_B);
1744
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
1745
1746
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
1747
1748
          for (int j = 0; j < order; j ++)
1749
              //Se modifican los elementos de la columna i a partir
1750
1751
              //de la fila j + 1.
1752
              //A(n, j) = A(n, j) / sqrt(A(j, j)) n > j
              int sstep = 0;
1753
```

```
cl SetKernel Arg(kernel, 4, sizeof(cl_int), &i);
1754
               cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
1755
               cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1756
                   &global WorkSize, &local WorkSize, O, NULL,
1757
1758
                   NULL);
1759
               //para una columna s, s > j, se modifican las filas n >= s
1760
               //A(n, s) = A(n, s) - A(s, j)*A(n, j) s > j, n >= s
1761
1762
               sstep++;
               cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
1763
               cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1764
                   &global WorkSize, &local WorkSize, O, NULL,
1765
                   NULL):
1766
1767
1768
          //cada elemento diagonal es reemplazado por su raiz cuadrada
1769
           //asi termina la factorización, la matriz Ac, tendra en su
1770
           //parte triangular inferior, los elementos de la factorización
1771
1772
          //de chol esky
1773
          int fstep = 2;
1774
          cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
1775
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
               &global WorkSize, &local WorkSize, O, NULL,
1776
1777
               NULL);
1778
          //luego de la factorización se tiene el sistema equivalente
1779
          //L*transpose(L)x=b que debe resol verse separadamente para
1780
1781
          //Ly=b
1782
          //y luego
1783
          //transpose(L)x=y
1784
1785
           //resol vi endo Ly=b
1786
          fstep++;
           cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
1787
1788
           for (int j = 0; j < order; j ++)
1789
1790
               //para un paso j, modificará el vector b tal que :
               \label{eq:bounds} $$ //b[n] = b[n] - b[j] * L[n, j] / L[j, j] n > j $$ cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j); $$
1791
1792
1793
               cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1794
                   &global WorkSize, &local WorkSize, O, NULL,
1795
                   NULL);
1796
           }
1797
1798
           fstep++;
           cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
1799
           cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1800
               &global WorkSize, &local WorkSize, O, NULL,
1801
1802
               NULL);
1803
1804
           //resolviendo transpose(L)x=y
1805
          fstep++;
          cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
1806
1807
           for (int j = 0; j < order; j++)
1808
           {
               //para un paso j, modifica b[order-j], este valor
1809
```

```
//es un componente final de la solución del sistema
1810
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j);
1811
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1812
1813
                  &global WorkSize, &local WorkSize, O, NULL,
1814
                  NULL);
1815
1816
          //devolviendo el control de memoria af::array a ArrayFire
1817
1818
          af_unlock_array(Ac);
1819
          af_unlock_array(Bc);
1820
          // copi ando el resultado en dC
1821
1822
          af_copy_array(dC, Bc);
1823
1824
          af_release_array(Ac);
1825
          af_rel ease_array(Bc);
1826 }
1827
     void AFire::fac_chol_c(af_array* L, af_array A) {
1828
1829
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
1830
          //cl_context af_context;
1831
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
1832
1833
          static cl_command_queue af_queue = afcl::getQueue();
1834
          //2. Obteniendo parámetros necesarios
1835
          //copi a de A
1836
1837
          af_array Ac;
1838
          af_copy_array(&Ac, A);
1839
          dim_t _order[AF_MAX_DIMS];
1840
1841
          af_get_dims(&_order[0], &_order[1], &_order[2],
1842
              &_order[3], A);
1843
          size_t order = _order[0];
1844
          size t local WorkSize = BLOCK SIZE * BLOCK SIZE;
1845
1846
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
1847
1848
          int status = CL_SUCCESS;
1849
1850
          af_dtype typef;
1851
          af_get_type(&typef, A);
1852
1853
          int msize = 0;
1854
          if (typef == f64)
1855
              msi ze = si zeof(doubl e);
1856
          else if (typef == f32)
              msi ze = si zeof(float);
1857
1858
          el se;
1859
1860
          //3. obteniendo las referencias cl_mem de los objetos af::array
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
1861
              CL_MEM_READ_WRITE, msize*order*order,
1862
1863
              NULL, &status);
1864
          af_get_device_ptr((voi d**)d_A, Ac);
1865
```

```
1866
          size_t program_length = strlen(Cholesky_source);
1867
1868
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto →
             de entrada
1869
          // para el Kernel
          cl_program program = clCreateProgramWithSource(af_context, 1, (const char →
1870
              **) & Cholesky_source, & program_length, & status);
          status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
1871
1872
1873
          char* kernel Name;
1874
          if (typef == f64)
              kernel Name = "Chol esky_c";
1875
          else if (typef == f32)
1876
              kernel Name = "Chol esky_c_sp";
1877
1878
          el se:
1879
          cl_kernel kernel = clCreateKernel(program, kernelName,
1880
              &status);
1881
1882
          // 5. estableciendo los argumentos
1883
          int i = 0;
1884
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_A);
1885
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), 0);
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
1886
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
1887
1888
          //6. Ej ecutando el kernel
1889
          for (int j = 0; j < order; j++)
1890
1891
1892
              //Se modifican los elementos de la columna j a partir
1893
              //de la fila j + 1.
              //A(n, j) = A(n, j) / sqrt(A(j, j)) n > j
1894
1895
              int sstep = 0;
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j);
1896
1897
              cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
1898
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1899
                   &global WorkSize, &local WorkSize, O, NULL,
1900
                  NULL);
1901
1902
              //para una columna s, s > j, se modifican las filas n >= s
1903
              //A(n, s) = A(n, s) - A(s, j)*A(n, j) s > j, n >= s
1904
              sstep++;
              cl SetKernel Arg(kernel, 5, sizeof(cl int), &sstep);
1905
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1906
                   &global WorkSize, &local WorkSize, O, NULL,
1907
1908
                   NULL);
1909
          }
1910
          //cada el emento di agonal es reemplazado por su rai z cuadrada
1911
1912
          //asi termina la factorización, la matriz Ac, tendra en su
1913
          //parte triangular inferior, los elementos de la factorización
1914
          //de chol esky
          int fstep = 2;
1915
          cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
1916
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
1917
1918
              &global WorkSize, &local WorkSize, O, NULL,
1919
              NULL);
```

```
1920
          //7. devolviendo el control de memoria af::array a ArrayFire
1921
1922
          af_unlock_array(Ac);
1923
1924
          // copi ando el resultado en dC
1925
          af_copy_array(L, Ac);
1926
          af_release_array(Ac);
1927
1928
1929
1930 void AFire:: SELIdIt_c(af_array* dC, af_array dA, af_array dB) {
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
1931
1932
          //cl_context af_context;
1933
          static cl_context af_context = afcl::getContext();
1934
          static cl_device_id af_device_id = afcl::getDeviceId();
1935
          static cl_command_queue af_queue = afcl::getQueue();
1936
          //copi a de A y B
1937
1938
          af_array Ac;
1939
          af_array Bc;
1940
          af_copy_array(&Ac, dA);
1941
          af_copy_array(&Bc, dB);
1942
1943
          dim t order[AF MAX DIMS];
1944
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
1945
          size_t order = _order[0];
1946
1947
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
1948
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
1949
          int status = CL_SUCCESS;
1950
1951
1952
          af_dtype typef;
1953
          af_get_type(&typef, dA);
1954
          int msize = 0;
1955
1956
          if (typef == f64)
              msi ze = si zeof(doubl e);
1957
1958
          else if (typef == f32)
1959
              msi ze = si zeof(float);
1960
          el se;
1961
          //obteniendo las referencias cl_mem de los objetos af::array
1962
1963
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
1964
              CL_MEM_READ_WRITE, msize*order*order,
              NULL, &status);
1965
1966
          af_get_devi ce_ptr((voi d**)d_A, Ac);
1967
1968
          cl_mem *d_B = (cl_mem*)clCreateBuffer(af_context,
1969
              CL_MEM_READ_WRITE, msize*order,
1970
              NULL, &status);
          af_get_device_ptr((void**)d_B, Bc);
1971
1972
1973
          size_t program_length = strlen(ldlt_source);
1974
1975
```

```
1976
          //creando el programa, construyendo el ejecutable y extrayendo el punto
            de entrada
1977
          // para el Kernel
          cl_program program = clCreateProgramWi thSource(af_context, 1, (const char >)
1978
             **)&IdIt_source, &program_length, &status);
          status = cl BuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
1979
1980
          char* kernel Name;
1981
          if (typef == f64)
1982
              kernelName = "Idlt_c";
1983
1984
          else if (typef == f32)
              kernel Name = "IdIt_c_sp";
1985
1986
          el se:
          cl_kernel kernel = clCreateKernel (program, kernel Name,
1987
1988
              &status);
1989
1990
          // estableciendo los argumentos
1991
          int i = 0:
1992
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_A);
1993
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_B);
1994
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
1995
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
1996
          for (int j = 0; j < order; j ++)
1997
1998
              //Se modifican los elementos de la columna j a partir
1999
2000
              //de la fila j + 1.
              //A(n,j)=A(n,j)/A(j,j) n>j
2001
2002
              int sstep = 0;
2003
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j);
              cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
2004
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2005
                  &global WorkSize, &local WorkSize, O, NULL,
2006
                  NULL);
2007
2008
2009
              //para una columna s, s > j, se modifican las filas
2010
              //n >= s
              //A(n, s) = A(n, s) - A(s, j)*A(n, j)*A(j, j)
2011
2012
              //s > j, n >= s
2013
              sstep++;
              cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
2014
2015
              cl EnqueueNDRangeKernel (af queue, kernel, 1, 0,
                  &global WorkSize, &local WorkSize, O, NULL,
2016
2017
                  NULL);
2018
2019
          //luego de la factorización se tiene el sistema equivalente
2020
2021
          //L*D*transpose(L)x=b que debe resolverse separadamente para
2022
          //Ly=b, luego para Dz=y y finalmente para transpose(L)x=z
2023
2024
          //resol vi endo Ly=b
2025
          int fstep = 2;
          cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
2026
2027
          for (int j = 0; j < order; j ++)
2028
          {
              //para un paso j, modificará el vector b tal que :
2029
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
38
```

```
//b[n] = b[n] - b[i] * L[n, i] n > i
2030
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &i);
2031
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2032
2033
                  &global WorkSize, &local WorkSize, O, NULL,
2034
                  NULL);
2035
2036
          //resol vi endo Dz=y
2037
2038
          fstep++;
          cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
2039
2040
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2041
              &global WorkSize, &local WorkSize, O, NULL,
              NULL):
2042
2043
2044
          //resol vi endo transpose(L)x=z
2045
          fstep++;
          cl SetKernel Arg(kernel, 5, sizeof(cl_int), &fstep);
2046
2047
          for (int j = 0; j < order; j ++)
2048
2049
              //para un paso j modificará el elemento b[i]
2050
              //tal que i = n - j:
2051
              //b[i] = b[i] - sumatoria{ j = i + 1, n }(z[j] * L[j, i])
              //este valor será la solución final x[i] del
2052
2053
              //sistema Ax = b
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j);
2054
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2055
                  &global WorkSize, &local WorkSize, O, NULL,
2056
                  NULL);
2057
2058
          }
2059
          //devolviendo el control de memoria af::array a ArrayFire
2060
2061
          af unlock array(Ac);
          af_unlock_array(Bc);
2062
2063
2064
          // copi ando el resul tado en dC
2065
          af_copy_array(dC, Bc);
2066
2067
          af_release_array(Ac);
2068
          af_release_array(Bc);
2069
     }
2070
     void AFire::fac Idlt c(af array* L, af array A) {
2071
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2072
2073
          //cl_context af_context;
2074
          static cl_context af_context = afcl::getContext();
2075
          static cl_device_id af_device_id = afcl::getDeviceId();
2076
          static cl_command_queue af_queue = afcl::getQueue();
2077
2078
          //2. Obteniendo parámetros necesarios
2079
          //copi a de A
2080
          af_array Ac;
2081
          af_copy_array(&Ac, A);
2082
          dim t order[AF MAX DIMS];
2083
2084
          af_get_dims(&_order[0], &_order[1], &_order[2],
              &_order[3], A);
2085
```

```
2086
          size_t order = _order[0];
2087
2088
          size_t local WorkSize = BLOCK_SIZE * BLOCK_SIZE;
2089
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
2090
2091
          int status = CL_SUCCESS;
2092
2093
          af_dtype typef;
2094
          af_get_type(&typef, A);
2095
2096
          int msize = 0;
2097
          if (typef == f64)
              msi ze = si zeof(doubl e);
2098
2099
          else if (typef == f32)
2100
             msi ze = si zeof(float);
2101
          el se;
2102
          //3. obteni endo las referencias cl_mem de los objetos af::array
2103
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
2104
2105
              CL_MEM_READ_WRITE, msize*order*order,
2106
              NULL, &status);
          af_get_device_ptr((void**)d_A, Ac);
2107
2108
2109
          size_t program_l ength = strl en(ldl t_source);
2110
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto →
2111
             de entrada
2112
          // para el Kernel
          cl_program program = clCreateProgramWithSource(af_context, 1, (const char >
2113
             **)&IdIt_source, &program_Iength, &status);
          status = cl BuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
2114
2115
          char* kernel Name;
2116
2117
          if (typef == f64)
2118
              kernel Name = "IdIt c";
          else if (typef == f32)
2119
2120
              kernel Name = "Idlt_c_sp";
2121
2122
          cl_kernel kernel = clCreateKernel (program, kernel Name,
2123
              &status);
2124
2125
          // 5. estableciendo los argumentos
2126
          int i = 0;
2127
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_A);
2128
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
          cl SetKernel Arg(kernel, i++, msize*local WorkSize, 0);
2129
2130
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
2131
2132
          //6. Ej ecutando el kernel
2133
          for (int j = 0; j < order; j ++)
2134
2135
              //Se modifican los elementos de la columna j a partir
2136
              //de la fila i + 1.
2137
              //A(n,j)=A(n,j)/A(j,j) n>j
              int sstep = 0;
2138
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j);
2139
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
40
```

```
cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
2140
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2141
2142
                  &global WorkSize, &local WorkSize, O, NULL,
2143
                  NULL);
2144
              //para una columna s, s > j, se modifican las filas
2145
2146
              //A(n, s) = A(n, s) - A(s, j)*A(n, j)*A(j, j)
2147
2148
              //s > j, n >= s
2149
              sstep++;
2150
              cl SetKernel Arg(kernel, 5, sizeof(cl_int), &sstep);
2151
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
                  &global WorkSize, &local WorkSize, O, NULL,
2152
                  NULL);
2153
2154
2155
          //7. devolviendo el control de memoria af::array a ArrayFire
2156
2157
          af_unlock_array(Ac);
2158
2159
          // copi ando el resultado en dC
2160
          af_copy_array(L, Ac);
2161
          af_release_array(Ac);
2162
2163
      }
2164
2165 void AFire::fac_sparse_chol_c(af_array elmA, af_array colA,
          af_array rowA, af_array elmL, af_array colL,
2166
          af_array rowL) {
2167
2168
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2169
          //cl_context af_context;
          static cl_context af_context = afcl::getContext();
2170
          static cl_device_id af_device_id = afcl::getDeviceId();
2171
          static cl_command_queue af_queue = afcl::getQueue();
2172
2173
2174
          //2. Obteni endo parámetros necesarios
2175
2176
          //longitud de los vectores
          dim_t _order[AF_MAX_DIMS];
2177
2178
          af_get_dims(&_order[0], &_order[1], &_order[2],
2179
              &_order[3], elmA);
2180
          size_t size_elmA = _order[0];
2181
          af_get_dims(&_order[0], &_order[1], &_order[2],
2182
              &_order[3], col A);
2183
2184
          size_t size_colA = _order[0];
2185
2186
          af_get_dims(&_order[0], &_order[1], &_order[2],
              &_order[3], rowA);
2187
2188
          size_t size_rowA = _order[0];
2189
2190
          af_get_dims(&_order[0], &_order[1], &_order[2],
2191
              &_order[3], elmL);
2192
          size_t size_elmL = _order[0];
2193
2194
          af_get_dims(&_order[0], &_order[1], &_order[2],
              &_order[3], coll);
2195
```

```
size_t size_colL = _order[0];
2196
2197
2198
          af_get_dims(&_order[0], &_order[1], &_order[2],
2199
              &_order[3], rowL);
2200
          size_t size_rowL = _order[0];
2201
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
2202
2203
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
2204
2205
          int status = CL_SUCCESS;
2206
2207
          af_dtype typef;
2208
          af_get_type(&typef, elmA);
2209
2210
          int msize = 0;
          if (typef == f64)
2211
2212
              msize = sizeof(double);
          else if (typef == f32)
2213
2214
              msi ze = si zeof(fl oat);
2215
          el se;
2216
2217
          //3. obteni endo las referencias cl_mem de los objetos af::array
          cl_mem *d_el mA = (cl_mem*)clCreateBuffer(af_context,
2218
              CL_MEM_READ_ONLY, msize*size_elmA,
2219
2220
              NULL, &status);
          af_get_device_ptr((void**)d_elmA, elmA);
2221
2222
          cl_mem *d_colA = (cl_mem*)clCreateBuffer(af_context,
2223
2224
              CL_MEM_READ_ONLY, sizeof(int)*size_colA,
2225
              NULL, &status);
2226
          af_get_devi ce_ptr((voi d**)d_col A, col A);
2227
          cl_mem *d_rowA = (cl_mem*)clCreateBuffer(af_context,
2228
2229
              CL_MEM_READ_ONLY, sizeof(int)*size_rowA,
2230
              NULL, &status);
          af_get_device_ptr((void**)d_rowA, rowA);
2231
2232
2233
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, msize*size_elmL,
2234
2235
              NULL, &status);
          af_get_device_ptr((void**)d_elmL, elmL);
2236
2237
          cl_mem *d_colL = (cl_mem*)clCreateBuffer(af_context,
2238
2239
              CL_MEM_READ_WRITE, sizeof(int)*size_colL,
2240
              NULL, &status);
2241
          af_get_device_ptr((void**)d_colL, colL);
2242
2243
          cl_mem *d_rowL = (cl_mem*)clCreateBuffer(af_context,
2244
              CL_MEM_READ_WRITE, sizeof(int)*size_rowL,
2245
              NULL, &status);
2246
          af_get_device_ptr((void**)d_rowL, rowL);
2247
          size_t program_l ength = strlen(Chol esky_source);
2248
2249
2250
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto →
             de entrada
```

```
2251
          // para el Kernel
          cl_program program = clCreateProgramWithSource(af_context, 1, (const char >
2252
             **) & Chol esky_source, & program_l ength, & status);
2253
          status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
2254
2255
          char* kernel Name;
          if (typef == f64)
2256
2257
              kernel Name = "Chol esky_sparse_c";
2258
          else if (typef == f32)
              kernel Name = "Chol esky_sparse_c_sp";
2259
2260
          el se;
2261
          cl_kernel kernel = clCreateKernel (program, kernel Name,
2262
              &status);
2263
2264
          int step = 0;
2265
          int sstep = -1;
2266
          // 5. estableciendo los argumentos
2267
2268
          int i = 0;
2269
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mA);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_colA);
2270
2271
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowA);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
2272
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_colL);
2273
2274
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_el mA);
2275
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_colA);
2276
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), 0);
2277
2278
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
2279
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &step);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &sstep);
2280
2281
2282
          //6. ej ecutando el kernel
2283
2284
          //primera ejecución con sstep=-1 copiará los elementos
2285
          //de elmA en el lugar correspondiente de elmL, tomar
2286
          //en cuenta que los elementos iniciales del elmL son
2287
          //todos cero.
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2288
2289
              &global WorkSize, &local WorkSize, O, NULL,
2290
              NULL);
2291
2292
          //6. Ej ecutando el kernel
          for (int j = 0; j < size_colA-1; j++)
2293
2294
2295
              //Se modifican los elementos de la columna j a partir
2296
              //de la fila j + 1.
2297
              //A(n, j) = A(n, j) / sqrt(A(j, j)) n > j
2298
              int sstep = 0;
2299
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
2300
              cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2301
2302
                   &global WorkSize, &local WorkSize, O, NULL,
2303
                   NULL);
2304
              //para una columna s, s > j, se modifican las filas n >= s
2305
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
43
```

```
//A(n, s) = A(n, s) - A(s, j)*A(n, j) s > j, n >= s
2306
2307
              sstep++;
              cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
2308
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2309
2310
                  &global WorkSize, &local WorkSize, O, NULL,
2311
                  NULL);
2312
          }
2313
2314
          //cada elemento diagonal es reemplazado por su raiz cuadrada
2315
          //asi termina la factorización, la matriz Ac, tendra en su
2316
          //parte triangular inferior, los elementos de la factorización
2317
          //de chol esky
          int fstep = 2;
2318
2319
          cl SetKernel Arg(kernel, 11, si zeof(cl_int), &fstep);
2320
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
              &global WorkSize, &local WorkSize, O, NULL,
2321
2322
              NULL);
2323
2324
          //7. devolviendo el control de memoria af::array a ArrayFire
2325
          af_unlock_array(elmA);
          af_unlock_array(col A);
2326
2327
          af_unlock_array(rowA);
2328
          af_unlock_array(elmL);
2329
          af unlock array(colL);
2330
          af_unl ock_array(rowL);
2331
     }
2332
2333 void AFire:: SELchol_sparse_c(af_array* dC, af_array elmA,
2334
          af_array col A, af_array rowA, af_array el mL,
2335
          af_array coll, af_array rowL, af_array dB) {
2336
2337
          AFire:: fac_sparse_chol_c(elmA, colA, rowA, elmL,
2338
              colL, rowL);
2339
          AFire:: SELchol_sparse_c(dC, elmL, colL, rowL, dB);
2340
          /*//1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2341
          //cl_context af_context;
2342
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
2343
2344
          static cl_command_queue af_queue = afcl::getQueue();
2345
          //creando copi a de dB
2346
2347
          af array b;
2348
          af_copy_array(&b, dB);
2349
2350
          //2. Obteni endo parámetros necesarios
2351
2352
          //longitud de los vectores
2353
          dim_t _order[AF_MAX_DIMS];
2354
          af_get_dims(&_order[0], &_order[1], &_order[2],
2355
              &_order[3], elmA);
2356
          cl_int size_elmA = _order[0];
2357
2358
          af_get_dims(&_order[0], &_order[1], &_order[2],
2359
              &_order[3], colA);
2360
          cl_int size_colA = _order[0];
2361
```

```
2362
          af_get_dims(&_order[0], &_order[1], &_order[2],
2363
              &_order[3], rowA);
2364
          cl_int size_rowA = _order[0];
2365
2366
          af_get_dims(&_order[0], &_order[1], &_order[2],
2367
              &_order[3], elmL);
2368
          cl_int size_elmL = _order[0];
2369
2370
          af_get_dims(&_order[0], &_order[1], &_order[2],
2371
              &_order[3], coll);
2372
          cl_int size_colL = _order[0];
2373
          af_get_dims(&_order[0], &_order[1], &_order[2],
2374
              &_order[3], rowL);
2375
2376
          cl_int size_rowL = _order[0];
2377
2378
          size_t local WorkSize = BLOCK_SIZE * BLOCK_SIZE;
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
2379
2380
2381
          int status = CL_SUCCESS;
2382
2383
          af_dtype typef;
          af_get_type(&typef, elmA);
2384
2385
2386
          int msize = 0;
          if (typef == f64)
2387
              msi ze = si zeof(doubl e);
2388
          else if (typef == f32)
2389
2390
              msi ze = si zeof(float);
2391
          el se;
2392
2393
          //3. obteni endo las referencias cl mem de los objetos af::array
2394
          cl_mem *d_el mA = (cl_mem*)clCreateBuffer(af_context,
2395
              CL_MEM_READ_ONLY, msi ze*si ze_el mA,
2396
              NULL, &status);
2397
          af_get_device_ptr((void**)d_elmA, elmA);
2398
          cl_mem *d_colA = (cl_mem*)clCreateBuffer(af_context,
2399
2400
              CL_MEM_READ_ONLY, sizeof(int)*size_colA,
2401
              NULL, &status);
          af_get_device_ptr((void**)d_colA, colA);
2402
2403
2404
          cl_mem *d_rowA = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_ONLY, sizeof(int)*size_rowA,
2405
2406
              NULL, &status);
2407
          af_get_device_ptr((void**)d_rowA, rowA);
2408
2409
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
2410
              CL_MEM_READ_WRITE, msize*size_elmL,
2411
              NULL, &status);
2412
          af_get_device_ptr((void**)d_elmL, elmL);
2413
2414
          cl mem *d colL = (cl mem*)clCreateBuffer(af context,
              CL_MEM_READ_WRITE, sizeof(int)*size_colL,
2415
2416
              NULL, &status);
          af_get_devi ce_ptr((voi d**)d_col L, col L);
2417
```

```
2418
2419
          cl_mem *d_rowL = (cl_mem*)clCreateBuffer(af_context,
2420
              CL_MEM_READ_WRITE, sizeof(int)*size_rowL,
2421
              NULL, &status);
2422
          af_get_devi ce_ptr((voi d**)d_rowL, rowL);
2423
          cl_mem *d_b = (cl_mem*)clCreateBuffer(af_context,
2424
              CL_MEM_READ_WRITE, msize*size_colA,
2425
              NULL, &status);
2426
2427
          af_get_devi ce_ptr((voi d**)d_b, b);
2428
2429
          size_t program_length = strlen(Cholesky_source);
2430
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto →
2431
             de entrada
2432
          // para el Kernel
2433
          cl_program program = clCreateProgramWithSource(af_context,
2434
              1, (const char **) & Chol esky_source, & program_l ength,
2435
              &status);
2436
          status = cl BuildProgram(program, 1, &af_device_id,
              NULL, NULL, NULL);
2437
2438
2439
          char* kernel Name;
2440
          if (typef == f64)
2441
              kernel Name = "Chol esky_sparse_c";
2442
          else if (typef == f32)
2443
              kernel Name = "Chol esky_sparse_c_sp";
2444
2445
          cl_kernel kernel = clCreateKernel (program, kernel Name,
2446
              &status);
2447
2448
          cl int step = 0;
2449
          cl_int sstep = -1;
2450
2451
          // 5. estableciendo los argumentos
2452
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mA);
2453
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_col A);
2454
2455
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowA);
2456
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_elmL);
2457
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_col L);
2458
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_el mA);
2459
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_col A);
2460
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_b);
2461
          cl SetKernel Arg(kernel, i++, msi ze*local WorkSi ze, 0);
2462
2463
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &step);
2464
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &sstep);
2465
2466
          //6. ej ecutando el kernel
2467
2468
          //primera ejecución con sstep=-1 copiará los elementos
2469
          //de el mA en el lugar correspondiente de el mL, tomar
2470
          //en cuenta que los elementos iniciales del elmL son
2471
          //todos cero.
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2472
```

```
&global WorkSize, &local WorkSize, O, NULL,
2473
2474
              NULL);
2475
2476
          for (cl_int_j = 0; j < size_colA - 1; j++)
2477
              //Se modifican los elementos de la columna j a partir
2478
2479
              //de la fila i + 1.
              //A(n, j) = A(n, j) / sqrt(A(j, j)) n > j
2480
2481
              sstep = 0;
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
2482
2483
              cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2484
                  &global WorkSize, &local WorkSize, O, NULL,
2485
                  NULL);
2486
2487
              //para una columna s, s > j, se modifican las filas n >= s
2488
2489
              //A(n, s) = A(n, s) - A(s, j)*A(n, j) s > j, n >= s
              sstep++;
2490
2491
              cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
2492
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2493
                  &global WorkSize, &local WorkSize, O, NULL,
2494
                  NULL);
2495
2496
2497
          //cada el emento di agonal es reemplazado por su rai z cuadrada
          //asi termina la factorización, la matriz Ac, tendra en su
2498
          //parte triangular inferior, los elementos de la factorización
2499
2500
          //de chol esky
          sstep++;
2501
2502
          cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2503
2504
              &global WorkSize, &local WorkSize, 0, NULL,
2505
              NULL);
2506
2507
          //termina la factorización
          //resolviendo el sistema Ly=b (y=transpose(L)*x)
2508
2509
          sstep++;
2510
          cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
2511
          for (cl_int_j = 0; j < size_colA - 1; j++) {
2512
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
2513
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2514
                  &global WorkSize, &local WorkSize, O, NULL,
2515
                  NULL);
2516
2517
          cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
2518
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2519
2520
              &global WorkSize, &local WorkSize, O, NULL,
2521
              NULL);
2522
2523
          //resolviendo el sistema transpose(L)*x=y
2524
          cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
2525
2526
2527
          for (cl_int_j = 0; j < size_colA; j++) {
2528
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
47
```

```
cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2529
2530
                  &global WorkSize, &local WorkSize, O, NULL,
2531
                  NULL):
2532
2533
2534
          //7. devolviendo el control de memoria af::array a ArrayFire
          af_unlock_array(elmA);
2535
          af_unlock_array(col A);
2536
2537
          af_unlock_array(rowA);
2538
          af_unlock_array(elmL);
2539
          af_unlock_array(col L);
2540
          af_unlock_array(rowL);
2541
          af_unlock_array(b);
2542
2543
          //copi ando al argumento de salida
2544
          af_copy_array(dC, b);
2545
2546
          af_rel ease_array(b); */
2547
     }
2548
2549
     void AFire::SELchol_sparse_c(af_array* dC, af_array elmL,
2550
          af_array coll, af_array rowL, af_array dB) {
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2551
2552
          //cl context af context;
2553
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
2554
2555
          static cl_command_queue af_queue = afcl::getQueue();
2556
2557
          //creando copia de dB
2558
          af_array b;
2559
          af_copy_array(&b, dB);
2560
2561
          //2. Obteniendo parámetros necesarios
2562
2563
          //longitud de los vectores
2564
          dim_t _order[AF_MAX_DIMS];
          af_get_dims(&_order[0], &_order[1], &_order[2],
2565
2566
              &_order[3], elmL);
2567
          size_t size_elmL = _order[0];
2568
2569
          af_get_dims(&_order[0], &_order[1], &_order[2],
2570
              & order[3], coll);
          size_t size_colL = _order[0];
2571
2572
2573
          af_get_dims(&_order[0], &_order[1], &_order[2],
2574
              &_order[3], rowL);
2575
          size_t size_rowL = _order[0];
2576
2577
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
2578
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
2579
2580
          int status = CL SUCCESS;
2581
2582
          af dtype typef;
2583
          af_get_type(&typef, elmL);
2584
```

```
2585
          int msize = 0;
          if (typef == f64)
2586
2587
              msi ze = si zeof(doubl e);
2588
          else if (typef == f32)
2589
              msi ze = si zeof(float);
2590
          el se:
2591
2592
          //3. obteniendo las referencias cl_mem de los objetos af::array
2593
          cl_mem *d_el mL = (cl_mem*)clCreateBuffer(af_context,
2594
              CL_MEM_READ_WRITE, msize*size_elmL,
2595
              NULL, &status);
2596
          af_get_device_ptr((void**)d_elmL, elmL);
2597
2598
          cl_mem *d_colL = (cl_mem*)clCreateBuffer(af_context,
2599
              CL_MEM_READ_WRITE, sizeof(int)*size_colL,
2600
              NULL, &status);
2601
          af_get_device_ptr((void**)d_colL, colL);
2602
2603
          cl_mem *d_rowL = (cl_mem*)clCreateBuffer(af_context,
2604
              CL_MEM_READ_WRITE, sizeof(int)*size_rowL,
2605
              NULL, &status);
2606
          af_get_device_ptr((void**)d_rowL, rowL);
2607
          cl mem *d b = (cl mem*)clCreateBuffer(af context,
2608
2609
              CL_MEM_READ_WRITE, msize*size_colL,
              NULL, &status);
2610
2611
          af_get_device_ptr((void**)d_b, b);
2612
2613
          size_t program_l ength = strlen(Cholesky_source);
2614
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto P
2615
             de entrada
          // para el Kernel
2616
2617
          cl_program program = cl CreateProgramWi thSource(af_context,
2618
              1, (const char **) & Chol esky_source, & program_l ength,
2619
              &status);
          status = cl BuildProgram(program, 1, &af_device_id,
2620
              NULL, NULL, NULL);
2621
2622
2623
          char* kernel Name;
2624
          if (typef == f64)
              kernel Name = "Chol esky sparse c";
2625
          else if (typef == f32)
2626
2627
              kernel Name = "Chol esky_sparse_c_sp";
2628
2629
          cl_kernel kernel = clCreateKernel(program, kernelName,
2630
              &status);
2631
2632
          int step = 0;
2633
          int sstep = 0;
2634
          int kk = 0;
2635
          // 5. estableciendo los argumentos
2636
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
2637
2638
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
2639
```

```
... sual Studio 2013\Projects\AF_func\AF_sample.cpp
```

```
49
```

```
2640
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_colL);
2641
2642
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowL);
2643
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &kk);
2644
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_colL);
2645
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_b);
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
2646
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &step);
2647
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &sstep);
2648
2649
2650
          //6. ejecutando el kernel
2651
          //resolviendo el sistema Ly=b (y=transpose(L)*x)
2652
2653
          sstep=3:
2654
          cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
2655
          for (int j = 0; j < size_colL - 1; j++) {
2656
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
2657
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2658
2659
                   &global WorkSize, &local WorkSize, O, NULL,
2660
                   NULL);
2661
          }
2662
          sstep++;
2663
          cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
2664
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
              &global WorkSize, &local WorkSize, O, NULL,
2665
2666
              NULL);
2667
2668
          //resolviendo el sistema transpose(L)*x=y
2669
          cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
2670
2671
          for (int j = 0; j < size_colL; j++) {
2672
2673
2674
              cl SetKernel Arg(kernel, 10, si zeof(cl_int), &i);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2675
                   &global WorkSize, &local WorkSize, O, NULL,
2676
2677
                  NULL);
2678
2679
          //7. devolviendo el control de memoria af::array a ArrayFire
2680
2681
          af unlock array(elmL);
2682
          af unlock array(col L);
          af_unlock_array(rowL);
2683
2684
          af_unlock_array(b);
2685
2686
          //copi ando al argumento de salida
2687
          af_copy_array(dC, b);
2688
2689
          af_rel ease_array(b);
2690 }
2691
2692 void AFire::fac_sparse_chol_sks(af_array elmA,
2693
          af_array i dxA) {
2694
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2695
          //cl_context af_context;
```

```
... sual Studio 2013\Projects\AF_func\AF_sample.cpp
```

```
50
```

```
2696
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
2697
2698
          static cl_command_queue af_queue = afcl::getQueue();
2699
2700
          //2. Obteni endo parámetros necesarios
2701
          //longitud de los vectores
          dim_t _order[AF_MAX_DIMS];
2702
2703
          af_get_dims(&_order[0], &_order[1], &_order[2],
2704
              &_order[3], elmA);
2705
          size_t size_elmL = _order[0];
2706
2707
          af_get_dims(&_order[0], &_order[1], &_order[2],
2708
              &_order[3], i dxA);
          size_t size_idxL = _order[0];
2709
2710
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
2711
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
2712
2713
2714
          int status = CL_SUCCESS;
2715
          af_dtype typef;
2716
2717
          af_get_type(&typef, elmA);
2718
          int msize = 0;
2719
          if (typef == f64)
2720
              msi ze = si zeof(double);
2721
          else if (typef == f32)
2722
2723
              msi ze = si zeof(float);
2724
          el se;
2725
          //3. obteni endo las referencias cl_mem de los objetos af::array
2726
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
2727
              CL_MEM_READ_WRITE, msize*size_elmL,
2728
2729
              NULL, &status);
2730
          af_get_device_ptr((void**)d_elmL, elmA);
2731
          cl_mem *d_i dxL = (cl_mem*)clCreateBuffer(af_context,
2732
2733
              CL_MEM_READ_ONLY, sizeof(int)*size_idxL,
2734
              NULL, &status);
2735
          af_get_device_ptr((void**)d_idxL, idxA);
2736
2737
          size_t program_l ength = strlen(Cholesky_source);
2738
2739
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto →
             de entrada
2740
          // para el Kernel
2741
          cl_program program = cl CreateProgramWi thSource(af_context,
2742
              1, (const char **) & Chol esky_source, & program_l ength,
2743
2744
          status = cl BuildProgram(program, 1, &af_device_id,
2745
              NULL, NULL, NULL);
2746
          char* kernel Name;
2747
2748
          if (typef == f64)
              kernel Name = "chol_sparse_sks";
2749
          else if (typef == f32)
2750
```

```
2751
              kernel Name = "chol_sparse_sks_sp";
2752
2753
          cl_kernel kernel = clCreateKernel (program, kernel Name,
2754
              &status);
2755
          // 5. estableciendo los argumentos
2756
2757
          int i = 0:
2758
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_i dxL);
2759
          cl SetKernel Arg(kernel, i++, si zeof(cl_int), &si ze_el mL);
2760
2761
          cl SetKernel Arg(kernel, i++, si zeof(cl_int), &si ze_i dxL);
2762
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
          cl SetKernel Arg(kernel, i++, msi ze*local WorkSi ze, 0);
2763
2764
2765
          cl_int sstep = 0;
          //6. Ej ecutando el kernel
2766
2767
          for (cl_int_j = 0; j < size_idxL - 1; j++)
2768
2769
              //Se modifican los elementos de la columna j a partir
2770
              //de la fila j + 1.
              //A(n, j) = A(n, j) / sqrt(A(j, j)) n > j
2771
2772
              sstep = 0;
              cl SetKernel Arg(kernel, 6, sizeof(cl_int), &j);
2773
              cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
2774
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2775
                   &global WorkSize, &local WorkSize, O, NULL,
2776
                  NULL);
2777
2778
2779
              //para una columna s, s > j, se modifican las filas n >= s
2780
              //A(n, s) = A(n, s) - A(s, j) * A(n, j) s > j, n >= s
2781
              sstep++;
2782
              cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2783
2784
                   &global WorkSize, &local WorkSize, O, NULL,
2785
                  NULL);
2786
          }
2787
          sstep++;
2788
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2789
2790
              &global WorkSize, &local WorkSize, O, NULL,
2791
              NULL);
2792
2793
          //7. devolviendo el control de memoria af::array a ArrayFire
2794
          af_unlock_array(elmA);
2795
          af_unlock_array(idxA);
2796 }
2797
2798 void AFire:: SELchol_sparse_sks(af_array* dC, af_array elmL,
2799
          af_array idxL, af_array dB) {
2800
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2801
          //cl_context af_context;
          static cl_context af_context = afcl::getContext();
2802
2803
          static cl device id af device id = afcl::getDeviceId();
2804
          static cl_command_queue af_queue = afcl::getQueue();
2805
          //creando copia de dB
2806
```

```
2807
          af_array b;
2808
          af_copy_array(&b, dB);
2809
2810
          //2. Obteni endo parámetros necesarios
2811
2812
          //longitud de los vectores
          dim_t _order[AF_MAX_DIMS];
2813
          af_get_dims(&_order[0], &_order[1], &_order[2],
2814
2815
              &_order[3], elmL);
2816
          size_t size_elmL = _order[0];
2817
2818
          af_get_dims(&_order[0], &_order[1], &_order[2],
2819
              &_order[3], i dxL);
          size_t size_idxL = _order[0];
2820
2821
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
2822
2823
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
2824
2825
          int status = CL_SUCCESS;
2826
          af_dtype typef;
2827
2828
          af_get_type(&typef, elmL);
2829
          int msize = 0;
2830
2831
          if (typef == f64)
              msi ze = si zeof(double);
2832
2833
          else if (typef == f32)
              msi ze = si zeof(float);
2834
2835
          el se;
2836
          //3. obteni endo las referencias cl_mem de los objetos af::array
2837
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
2838
2839
              CL_MEM_READ_ONLY, msize*size_elmL,
2840
              NULL, &status);
2841
          af_get_device_ptr((void**)d_elmL, elmL);
2842
          cl_mem *d_i dxL = (cl_mem*)clCreateBuffer(af_context,
2843
              CL_MEM_READ_ONLY, sizeof(int)*size_idxL,
2844
2845
              NULL, &status);
2846
          af_get_device_ptr((void**)d_idxL, idxL);
2847
2848
          cl mem d b = (cl mem)clCreateBuffer(af context,
              CL_MEM_READ_WRITE, msize*size_idxL,
2849
2850
              NULL, &status);
2851
          af_get_devi ce_ptr((voi d**)d_b, b);
2852
2853
          size_t program_length = strlen(Cholesky_source);
2854
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto >
2855
             de entrada
2856
          // para el Kernel
          cl_program program = clCreateProgramWi thSource(af_context,
2857
2858
              1, (const char **) & Chol esky_source, & program_l ength,
2859
              &status);
2860
          status = cl BuildProgram(program, 1, &af_device_id,
              NULL, NULL, NULL);
2861
```

```
2862
2863
          char* kernel Name;
2864
          if (typef == f64)
              kernel Name = "chol_sparse_sks";
2865
2866
          else if (typef == f32)
2867
              kernel Name = "chol_sparse_sks_sp";
2868
          el se:
2869
          cl_kernel kernel = clCreateKernel(program, kernelName,
2870
              &status);
2871
          cl_int_kk = 0;
2872
          // 5. estableciendo los argumentos
2873
2874
          int i = 0:
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
2875
2876
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_idxL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &kk);
2877
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_idxL);
2878
2879
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_b);
2880
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
2881
2882
          //6. ej ecutando el kernel
2883
          //resolviendo el sistema Ly=b (y=transpose(L)*x)
2884
          cl int sstep = 3;
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
2885
2886
          for (cl_int_j = 0; j < size_idxL - 1; j++) {
2887
              cl SetKernel Arg(kernel, 6, sizeof(cl_int), &j);
2888
2889
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2890
                   &global WorkSize, &local WorkSize, O, NULL,
2891
                   NULL);
2892
          }
2893
          sstep++;
2894
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
2895
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2896
              &global WorkSize, &local WorkSize, O, NULL,
2897
              NULL);
2898
2899
          //resolviendo el sistema transpose(L)*x=y
2900
          sstep++;
2901
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
          for (cl_int j = 0; j < size_idxL; j++) {
2902
2903
              cl SetKernel Arg(kernel, 6, sizeof(cl int), &i);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
2904
2905
                   &global WorkSize, &local WorkSize, O, NULL,
2906
                   NULL);
2907
          }
2908
2909
          //7. devolviendo el control de memoria af::array a ArrayFire
2910
          af_unlock_array(elmL);
2911
          af_unlock_array(idxL);
2912
          af_unl ock_array(b);
2913
2914
          //copi ando al argumento de salida
2915
          af_copy_array(dC, b);
2916
2917
          af_rel ease_array(b);
```

```
2918 }
2919
2920 void AFire:: fac_sparse_ldlt_c(af_array elmA, af_array colA,
2921
          af_array rowA, af_array elmL, af_array colL,
2922
          af_array rowL) {
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
2923
2924
          //cl_context af_context;
2925
          static cl_context af_context = afcl::getContext();
2926
          static cl_device_id af_device_id = afcl::getDeviceId();
2927
          static cl_command_queue af_queue = afcl::getQueue();
2928
2929
          //2. Obteni endo parámetros necesarios
2930
2931
          //Iongi tud de los vectores
2932
          dim_t _order[AF_MAX_DIMS];
          af_get_dims(&_order[0], &_order[1], &_order[2],
2933
2934
              &_order[3], elmA);
2935
          size_t size_elmA = _order[0];
2936
2937
          af_get_dims(&_order[0], &_order[1], &_order[2],
2938
              &_order[3], col A);
2939
          size_t size_colA = _order[0];
2940
2941
          af_get_dims(&_order[0], &_order[1], &_order[2],
2942
              &_order[3], rowA);
2943
          size_t size_rowA = _order[0];
2944
2945
          af_get_dims(&_order[0], &_order[1], &_order[2],
2946
              &_order[3], elmL);
2947
          size_t size_elmL = _order[0];
2948
2949
          af_get_dims(&_order[0], &_order[1], &_order[2],
2950
              &_order[3], coll);
2951
          size_t size_colL = _order[0];
2952
          af_get_dims(&_order[0], &_order[1], &_order[2],
2953
2954
              &_order[3], rowL);
2955
          size_t size_rowL = _order[0];
2956
2957
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
2958
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
2959
          int status = CL_SUCCESS;
2960
2961
2962
          af_dtype typef;
2963
          af_get_type(&typef, elmA);
2964
2965
          int msize = 0;
2966
          if (typef == f64)
2967
              msi ze = si zeof(doubl e);
2968
          else if (typef == f32)
2969
              msi ze = si zeof(float);
2970
          el se:
2971
2972
          //3. obteni endo las referencias cl_mem de los objetos af::array
2973
          cl_mem *d_elmA = (cl_mem*)clCreateBuffer(af_context,
```

```
... sual Studio 2013\Projects\AF_func\AF_sample.cpp
```

```
55
```

```
2974
              CL_MEM_READ_ONLY, msize*size_elmA,
2975
              NULL, &status);
2976
          af_get_device_ptr((void**)d_elmA, elmA);
2977
2978
          cl_mem *d_colA = (cl_mem*)clCreateBuffer(af_context,
2979
              CL_MEM_READ_ONLY, sizeof(int)*size_colA,
2980
              NULL, &status);
          af_get_device_ptr((void**)d_colA, colA);
2981
2982
2983
          cl_mem *d_rowA = (cl_mem*)clCreateBuffer(af_context,
2984
              CL_MEM_READ_ONLY, sizeof(int)*size_rowA,
2985
              NULL, &status);
          af_get_devi ce_ptr((voi d**)d_rowA, rowA);
2986
2987
2988
          cl_mem *d_el mL = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, msize*size_elmL,
2989
2990
              NULL, &status);
2991
          af_get_device_ptr((void**)d_elmL, elmL);
2992
2993
          cl_mem *d_colL = (cl_mem*)clCreateBuffer(af_context,
2994
              CL_MEM_READ_WRITE, sizeof(int)*size_colL,
2995
              NULL, &status);
2996
          af_get_device_ptr((void**)d_colL, colL);
2997
2998
          cl_mem *d_rowL = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, sizeof(int)*size_rowL,
2999
3000
              NULL, &status);
          af_get_device_ptr((void**)d_rowL, rowL);
3001
3002
3003
          size_t program_l ength = strlen(ldlt_source);
3004
3005
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto P
             de entrada
3006
          // para el Kernel
3007
          cl_program program = clCreateProgramWi thSource(af_context,
              1, (const char **)&IdIt_source, &program_I ength,
3008
3009
              &status);
          status = cl BuildProgram(program, 1, &af_device_id,
3010
3011
              NULL, NULL, NULL);
3012
          char* kernel Name;
3013
3014
          if (typef == f64)
              kernel Name = "Idlt_sparse_c";
3015
          else if (typef == f32)
3016
              kernel Name = "Idlt_sparse_c_sp";
3017
3018
3019
          cl_kernel kernel = clCreateKernel (program, kernel Name,
3020
              &status);
3021
3022
          int step = 0;
3023
          int sstep = -1;
3024
3025
          // 5. estableciendo los argumentos
3026
          int i = 0:
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mA);
3027
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_col A);
3028
```

```
... sual Studio 2013\Projects\AF_func\AF_sample.cpp
```

```
56
```

```
3029
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowA);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
3030
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_colL);
3031
3032
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowL);
3033
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_el mA);
3034
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_colA);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
3035
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
3036
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &step);
3037
3038
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &sstep);
3039
3040
          //6. ej ecutando el kernel
3041
          //primera ejecución con sstep=-1 copiará los elementos
3042
3043
          //de el mA en el lugar correspondiente de el mL, tomar
3044
          //en cuenta que los elementos iniciales del elmL son
3045
          //todos cero.
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3046
              &global WorkSize, &local WorkSize, O, NULL,
3047
3048
              NULL);
3049
3050
          //6. Ej ecutando el kernel
          for (int j = 0; j < size_col A - 1; j++)
3051
3052
3053
              //Se modifican los elementos de la columna j a partir
              //de la fila j + 1.
3054
3055
              //A(n,j)=A(n,j)/A(j,j) n>j
              int sstep = 0;
3056
3057
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
3058
              cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3059
                   &global WorkSize, &local WorkSize, O, NULL,
3060
3061
                   NULL);
3062
3063
              //para una columna s, s > j, se modifican las filas n >= s
3064
              //A(n, s) = A(n, s) - A(s, j) * A(n, j) * A(j, j) s > j, n >= s
3065
              sstep++;
              cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
3066
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3067
3068
                   &global WorkSize, &local WorkSize, O, NULL,
3069
                   NULL);
3070
3071
3072
          //7. devolviendo el control de memoria af::array a ArrayFire
3073
          af_unlock_array(elmA);
3074
          af_unlock_array(col A);
3075
          af_unlock_array(rowA);
3076
          af_unl ock_array(el mL);
3077
          af_unlock_array(coll);
3078
          af_unlock_array(rowL);
3079
      }
3080
     void AFire::SELIdIt_sparse_c(af_array* dC, af_array elmA,
3081
3082
          af_array colA, af_array rowA, af_array elmL,
3083
          af_array coll, af_array rowL, af_array dB) {
3084
```

```
3085
          AFire::fac_sparse_ldlt_c(elmA, colA, rowA, elmL,
3086
              colL, rowL);
3087
          AFire:: SELIdIt_sparse_c(dC, elmL, colL, rowL, dB);
3088 }
3089
3090 void AFire:: SELIdIt_sparse_c(af_array* dC, af_array elmL,
          af_array coll, af_array rowL, af_array dB) {
3091
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3092
3093
          //cl_context af_context;
3094
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
3095
3096
          static cl_command_queue af_queue = afcl::getQueue();
3097
3098
          //creando copi a de dB
3099
          af_array b;
3100
          af_copy_array(&b, dB);
3101
          //2. Obteni endo parámetros necesarios
3102
3103
3104
          //longitud de los vectores
3105
          dim_t _order[AF_MAX_DIMS];
3106
          af_get_dims(&_order[0], &_order[1], &_order[2],
              &_order[3], elmL);
3107
3108
          size_t size_elmL = _order[0];
3109
          af_get_dims(&_order[0], &_order[1], &_order[2],
3110
3111
              &_order[3], coll);
          size_t size_colL = _order[0];
3112
3113
3114
          af_get_dims(&_order[0], &_order[1], &_order[2],
3115
              &_order[3], rowL);
3116
          size_t size_rowL = _order[0];
3117
3118
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
3119
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
3120
          int status = CL_SUCCESS;
3121
3122
3123
          af_dtype typef;
3124
          af_get_type(&typef, elmL);
3125
3126
          int msize = 0;
3127
          if (typef == f64)
              msi ze = si zeof(doubl e);
3128
3129
          else if (typef == f32)
3130
              msi ze = si zeof(float);
3131
          el se;
3132
3133
          //3. obteni endo las referencias cl_mem de los objetos af::array
3134
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
3135
              CL_MEM_READ_WRITE, msize*size_elmL,
3136
              NULL, &status);
          af_get_device_ptr((void**)d_elmL, elmL);
3137
3138
3139
          cl_mem *d_colL = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, sizeof(int)*size_colL,
3140
```

```
3141
              NULL, &status);
          af_get_device_ptr((void**)d_colL, colL);
3142
3143
3144
          cl_mem *d_rowL = (cl_mem*)clCreateBuffer(af_context,
3145
              CL_MEM_READ_WRITE, sizeof(int)*size_rowL,
3146
              NULL, &status);
          af_get_device_ptr((void**)d_rowL, rowL);
3147
3148
          cl_mem *d_b = (cl_mem*)clCreateBuffer(af_context,
3149
3150
              CL_MEM_READ_WRITE, msize*size_colL,
3151
              NULL, &status);
          af_get_device_ptr((voi d**)d_b, b);
3152
3153
          size_t program_length = strlen(ldlt_source);
3154
3155
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto >
3156
             de entrada
3157
          // para el Kernel
3158
          cl_program program = cl CreateProgramWi thSource(af_context,
3159
              1, (const char **) & I dl t_source, & program_I ength,
3160
              &status);
          status = cl BuildProgram(program, 1, &af_device_id,
3161
              NULL, NULL, NULL);
3162
3163
3164
          char* kernel Name;
          if (typef == f64)
3165
              kernel Name = "Idlt_sparse_c";
3166
          else if (typef == f32)
3167
              kernel Name = "Idlt_sparse_c_sp";
3168
3169
3170
          cl_kernel kernel = clCreateKernel(program, kernelName,
3171
              &status);
3172
3173
          int step = 0;
          int sstep = 0;
3174
3175
          int kk = 0;
          // 5. estableciendo los argumentos
3176
3177
          int i = 0;
3178
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
3179
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
3180
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
3181
          cl SetKernel Arg(kernel, i++, sizeof(cl mem), d el mL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_colL);
3182
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowL);
3183
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &kk);
3184
3185
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_colL);
3186
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_b);
          cl SetKernel Arg(kernel, i++, msize*local WorkSize, 0);
3187
3188
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &step);
3189
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &sstep);
3190
3191
          //6. ejecutando el kernel
3192
3193
          //resolviendo el sistema Ly=b (y=D*transpose(L)*x)
3194
          sstep = 2;
          cl SetKernel Arg(kernel, 11, sizeof(cl_int), &sstep);
3195
```

```
for (int j = 0; j < size_colL - 1; j++) {
3196
3197
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
3198
3199
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3200
                  &global WorkSize, &local WorkSize, O, NULL,
3201
                  NULL);
3202
          }
3203
3204
          //resolviendo el sistema Dz=y, z=transpose(L)*x
3205
          sstep++;
3206
          cl SetKernel Arg(kernel, 11, si zeof(cl_int), &sstep);
3207
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
              &global WorkSize, &local WorkSize, O, NULL,
3208
3209
              NULL);
3210
3211
          //resolviendo el sistema transpose(L)*x=z
3212
          cl SetKernel Arg(kernel, 11, si zeof(cl_i nt), &sstep);
3213
3214
3215
          for (int j = 0; j < size_colL; j ++) {
3216
3217
              cl SetKernel Arg(kernel, 10, sizeof(cl_int), &j);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3218
                  &global WorkSize, &local WorkSize, O, NULL,
3219
3220
                  NULL);
3221
          }
3222
3223
          //7. devolviendo el control de memoria af::array a ArrayFire
3224
          af_unlock_array(elmL);
3225
          af_unlock_array(col L);
3226
          af_unlock_array(rowL);
          af_unlock_array(b);
3227
3228
3229
          //copi ando al argumento de salida
3230
          af_copy_array(dC, b);
3231
3232
          af_release_array(b);
3233
      }
3234
3235
     void AFire:: fac_sparse_IdIt_sks(af_array elmA,
3236
          af_array idxA) {
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3237
          //cl_context af_context;
3238
3239
          static cl_context af_context = afcl::getContext();
3240
          static cl_device_id af_device_id = afcl::getDeviceId();
3241
          static cl_command_queue af_queue = afcl::getQueue();
3242
3243
          //2. Obteniendo parámetros necesarios
3244
          //longitud de los vectores
3245
          dim_t _order[AF_MAX_DIMS];
3246
          af_get_dims(&_order[0], &_order[1], &_order[2],
3247
              &_order[3], elmA);
3248
          size_t size_elmL = _order[0];
3249
3250
          af_get_dims(&_order[0], &_order[1], &_order[2],
              &_order[3], idxA);
3251
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
3252
          size_t size_idxL = _order[0];
3253
3254
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
3255
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
3256
          int status = CL_SUCCESS;
3257
3258
3259
          af_dtype typef;
          af_get_type(&typef, el mA);
3260
3261
3262
          int msize = 0;
3263
          if (typef == f64)
              msi ze = si zeof(doubl e);
3264
          else if (typef == f32)
3265
3266
              msi ze = si zeof(float);
3267
          el se;
3268
          //3. obteniendo las referencias cl_mem de los objetos af::array
3269
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
3270
3271
              CL_MEM_READ_WRITE, msize*size_elmL,
3272
              NULL, &status);
3273
          af_get_device_ptr((void**)d_elmL, elmA);
3274
3275
          cl mem *d idxL = (cl mem*)clCreateBuffer(af context,
3276
              CL_MEM_READ_ONLY, sizeof(int)*size_idxL,
              NULL, &status);
3277
          af_get_device_ptr((void**)d_idxL, idxA);
3278
3279
3280
          size_t program_length = strlen(ldlt_source);
3281
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto P
3282
             de entrada
3283
          // para el Kernel
3284
          cl_program program = cl CreateProgramWi thSource(af_context,
3285
              1, (const char **) &I dI t_source, &program_I ength,
3286
              &status);
          status = cl BuildProgram(program, 1, &af_device_id,
3287
              NULL, NULL, NULL);
3288
3289
3290
          char* kernel Name;
3291
          if (typef == f64)
              kernel Name = "Idlt sparse sks";
3292
          else if (typef == f32)
3293
3294
              kernel Name = "Idlt_sparse_sks_sp";
3295
3296
          cl_kernel kernel = clCreateKernel(program, kernelName,
3297
              &status);
3298
3299
          // 5. estableciendo los argumentos
3300
          int i = 0;
3301
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_i dxL);
3302
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_elmL);
3303
3304
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_idxL);
3305
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), 0);
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
3306
```

```
3307
3308
          //6. Ej ecutando el kernel
3309
          for (int j = 0; j < size_i dxL - 1; j++)
3310
3311
               //Se modifican los elementos de la columna j a partir
3312
               //de la fila j + 1.
               //A(n,j)=A(n,j)/A(j,j) n>j
3313
3314
               int sstep = 0;
               cl SetKernel Arg(kernel, 6, si zeof(cl_int), &j);
cl SetKernel Arg(kernel, 7, si zeof(cl_int), &sstep);
3315
3316
3317
               cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
                   &global WorkSize, &local WorkSize, O, NULL,
3318
3319
                   NULL):
3320
3321
               //para una columna s, s > j, se modifican las filas n >= s
3322
               //A(n, s) = A(n, s) - A(s, j) * A(n, j) * A(j, j) s > j, n > = s
3323
               cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
3324
               cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3325
3326
                   &global WorkSize, &local WorkSize, O, NULL,
3327
                   NULL);
3328
          }
3329
3330
          //7. devolviendo el control de memoria af::array a ArrayFire
3331
          af_unl ock_array(el mA);
          af_unl ock_array(i dxA);
3332
3333 }
3334
3335
      void AFire:: SELIdIt_sparse_sks(af_array* dC, af_array elmL,
3336
          af_array idxL, af_array dB) {
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3337
3338
          //cl context af context;
          static cl_context af_context = afcl::getContext();
3339
3340
          static cl_device_id af_device_id = afcl::getDeviceId();
3341
          static cl_command_queue af_queue = afcl::getQueue();
3342
          //creando copia de dB
3343
3344
          af_array b;
3345
          af_copy_array(&b, dB);
3346
          //2. Obteniendo parámetros necesarios
3347
3348
3349
          //longitud de los vectores
          dim_t _order[AF_MAX_DIMS];
3350
3351
          af_get_dims(&_order[0], &_order[1], &_order[2],
3352
               &_order[3], elmL);
3353
          size_t size_elmL = _order[0];
3354
3355
          af_get_dims(&_order[0], &_order[1], &_order[2],
3356
               &_order[3], idxL);
3357
          size_t size_idxL = _order[0];
3358
          size t local WorkSize = BLOCK SIZE * BLOCK SIZE;
3359
3360
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
3361
          int status = CL_SUCCESS;
3362
```

```
3363
3364
          af_dtype typef;
3365
          af_get_type(&typef, elmL);
3366
3367
          int msize = 0;
3368
          if (typef == f64)
              msi ze = si zeof(double);
3369
          else if (typef == f32)
3370
              msi ze = si zeof(float);
3371
3372
          el se;
3373
3374
          //3. obteni endo las referencias cl_mem de los objetos af::array
          cl_mem *d_elmL = (cl_mem*)clCreateBuffer(af_context,
3375
              CL_MEM_READ_ONLY, msi ze*si ze_el mL,
3376
3377
              NULL, &status);
          af_get_device_ptr((void**)d_elmL, elmL);
3378
3379
          cl_mem *d_i dxL = (cl_mem*)clCreateBuffer(af_context,
3380
3381
              CL_MEM_READ_ONLY, sizeof(int)*size_idxL,
3382
              NULL, &status);
          af_get_device_ptr((void**)d_idxL, idxL);
3383
3384
          cl_mem d_b = (cl_mem)clCreateBuffer(af_context,
3385
3386
              CL_MEM_READ_WRITE, msize*size_idxL,
3387
              NULL, &status);
          af_get_devi ce_ptr((voi d**)d_b, b);
3388
3389
3390
          size_t program_l ength = strlen(ldlt_source);
3391
3392
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto
             de entrada
3393
          // para el Kernel
3394
          cl_program program = cl CreateProgramWi thSource(af_context,
3395
              1, (const char **) &I dI t_source, &program_I ength,
3396
              &status);
3397
          status = cl BuildProgram(program, 1, &af_device_id,
              NULL, NULL, NULL);
3398
3399
3400
          char* kernel Name;
3401
          if (typef == f64)
              kernel Name = "I dl t_sparse_sks";
3402
3403
          else if (typef == f32)
              kernel Name = "I dl t_sparse_sks_sp";
3404
3405
          el se:
3406
          cl_kernel kernel = clCreateKernel(program, kernelName,
3407
              &status);
3408
3409
          cl_int_kk = 0;
3410
          // 5. estableciendo los argumentos
3411
          int i = 0;
3412
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_el mL);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_i dxL);
3413
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &kk);
3414
3415
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_idxL);
3416
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_b);
          cl SetKernel Arg(kernel, i++, msi ze*l ocal WorkSi ze, 0);
3417
```

```
3418
3419
          //6. ej ecutando el kernel
3420
          //resolviendo el sistema Ly=b (y=D*transpose(L)*x)
3421
          cl_int sstep = 2;
3422
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
3423
          for (cl_int_j = 0; j < size_idxL - 1; j++) {
3424
              cl SetKernel Arg(kernel, 6, sizeof(cl_int), &j);
3425
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3426
                  &global WorkSize, &local WorkSize, O, NULL,
3427
3428
                  NULL);
3429
          }
3430
3431
          //resolviendo el sistema Dz=y, z=transpose(L)*x
3432
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
3433
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3434
              &global WorkSize, &local WorkSize, O, NULL,
3435
3436
              NULL);
3437
3438
          //resolviendo el sistema transpose(L)*x=z
3439
          sstep++;
          cl SetKernel Arg(kernel, 7, sizeof(cl_int), &sstep);
3440
          for (cl_int_j = 0; j < size_idxL; j++) {
3441
3442
              cl SetKernel Arg(kernel, 6, sizeof(cl_int), &j);
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3443
3444
                  &global WorkSize, &local WorkSize, O, NULL,
3445
                  NULL);
3446
          }
3447
          //7. devolviendo el control de memoria af::array a ArrayFire
3448
3449
          af unlock array(elmL);
          af_unlock_array(idxL);
3450
3451
          af_unlock_array(b);
3452
3453
          //copi ando al argumento de salida
3454
          af_copy_array(dC, b);
3455
3456
          af_release_array(b);
3457
     }
3458 //----
3459 //Pruebas y otros
3460 //----
3461 void AFire::fenceTest_sp(af::array &A, af::array &B) {
3462
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3463
          static cl_context af_context = afcl::getContext();
3464
          static cl_device_id af_device_id = afcl::getDeviceId();
3465
          static cl_command_queue af_queue = afcl::getQueue();
3466
3467
          //obteniendo las referencias cl_mem de los objetos af::array
3468
          cl_mem * d_A = A. device<cl_mem>();
          cl_mem * d_B = B. device<cl_mem>();
3469
3470
3471
          //detalles importantes sobre el kernel
3472
          /**/
3473
```

```
3474
3475
          size_t order = (int)A.dims(0);
3476
3477
          size_t program_l ength = strlen(GJordan_source);
3478
          int status = CL_SUCCESS;
3479
          //creando el programa, construyendo el ejecutable y extrayendo el punto
3480
            de entrada
          // para el Kernel
3481
3482
          cl_program program = clCreateProgramWithSource(af_context, 1, (const char →
             **)&GJordan_source, &program_length, &status);
3483
          status = clBuildProgram(program, 1, &af_device_id, NULL, NULL, NULL);
          cl_kernel kernel = clCreateKernel(program, "fenceTest_sp", &status);
3484
3485
3486
          // estableciendo los argumentos
3487
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_A);
3488
3489
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_B);
3490
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
3491
3492
          size_t local WorkSize[] = { BLOCK_SIZE, BLOCK_SIZE };
3493
          size_t global WorkSize[] =
          { shrRoundUp(I ocal WorkSi ze[0], order), shrRoundUp(I ocal WorkSi ze[1],
3494
            order) };
3495
          //ejecutando el Kernel
          cl EnqueueNDRangeKernel (af_queue, kernel, 2, 0,
3496
3497
              global WorkSize, local WorkSize, 0, NULL, NULL);
3498
3499
          //devolviendo el control de memoria af::array a ArrayFire
3500
          A. unlock();
3501
          B. unlock();
3502
     }
3503
3504 void AFire::test_1(af::array &A, af::array &B)
3505 {
3506
          //B = A;
3507
          size_t len = A. dims(0);
          seq idx(0, len - 1, 1);
3508
3509
          af::copy(B, A, idx);
3510 }
3511
3512 void AFire::test_2(af_array dA)
3513 {
3514
          //1. Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3515
          //cl_context af_context;
3516
          static cl_context af_context = afcl::getContext();
3517
          static cl_device_id af_device_id = afcl::getDeviceId();
          static cl_command_queue af_queue = afcl::getQueue();
3518
3519
3520
          //2. Obteniendo parámetros necesarios
3521
3522
          //longitud de los vectores
3523
          dim_t _order[AF_MAX_DIMS];
          af_get_dims(&_order[0], &_order[1], &_order[2],
3524
3525
              &_order[3], dA);
3526
          size_t size_elmA = _order[0];
```

```
3527
3528
          af_array zeros;
3529
          dim_t size[] = { size_elmA };
3530
          af_constant(&zeros, 0, 1, size, s32);
3531
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
3532
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
3533
3534
          int status = CL_SUCCESS;
3535
3536
3537
          af_dtype typef;
3538
          af_get_type(&typef, dA);
3539
3540
          int msize = 0:
3541
          if (typef == f64)
              msi ze = si zeof(double);
3542
3543
          else if (typef == f32)
3544
              msi ze = si zeof(float);
3545
          el se;
3546
3547
          //3. obteni endo las referencias cl_mem de los objetos af::array
3548
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, msize*size_elmA,
3549
3550
              NULL, &status);
3551
          af_get_devi ce_ptr((voi d**)d_A, dA);
3552
          cl_mem *d_zeros = (cl_mem*)clCreateBuffer(af_context,
3553
3554
              CL_MEM_READ_WRITE, sizeof(int)*size_elmA,
              NULL, &status);
3555
3556
          af_get_devi ce_ptr((voi d**)d_zeros, zeros);
3557
3558
          size_t program_l ength = strl en(gc_source);
3559
3560
          //4. creando el programa, construyendo el ejecutable y extrayendo el punto 🤊
             de entrada
3561
          // para el Kernel
          cl_program program = clCreateProgramWi thSource(af_context,
3562
3563
              1, (const char **) &gc_source, &program_l ength,
3564
              &status);
3565
          status = cl BuildProgram(program, 1, &af_device_id,
3566
              NULL, NULL, NULL);
3567
          char* kernel Name;
3568
3569
          if (typef == f64)
              kernel Name = "prueba_1";
3570
3571
          else if (typef == f32)
3572
              kernel Name = "prueba_1_sp";
3573
          el se;
3574
          cl_kernel kernel = clCreateKernel(program, kernelName,
3575
              &status);
3576
          // 5. estableciendo los argumentos
3577
3578
          int i = 0;
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_A);
3579
3580
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_zeros);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &size_el mA);
3581
```

```
3582
3583
          //6. ej ecutando el kernel
3584
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3585
              &global WorkSize, &local WorkSize, O, NULL,
3586
              NULL);
3587
          af_print_array(zeros);
3588
          //7. devolviendo el control de memoria af::array a ArrayFire
3589
3590
          af_unlock_array(dA);
3591
          af_unlock_array(zeros);
3592
3593
          af_rel ease_array(zeros);
3594 }
3595
3596 void AFire::global_sync_test(af_array* dC, af_array dA,
3597
          af_array dB) {
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3598
3599
          //cl_context af_context;
3600
          static cl_context af_context = afcl::getContext();
3601
          static cl_device_id af_device_id = afcl::getDeviceId();
          static cl_command_queue af_queue = afcl::getQueue();
3602
3603
          //obteni endo el número de multiprocesadores, necesario
3604
3605
          //para determinar el número de blogues para los kernels,
3606
          //y hacer posible el uso de las funciones de sincronización
          //ver artículo:
3607
          //Inter-Block GPU Communication via Fast Barrier Synchronization
3608
          //Shucai Xi ao and Wu-chun Feng
3609
3610
          //Department of Computer Science Virginia Tech
3611
          //2009/9/19
3612
          //Pág. 5
3613
          cl_uint compute_units;
          cl GetDevi cel nfo(af_devi ce_i d,
3614
3615
              CL_DEVICE_MAX_COMPUTE_UNITS, sizeof(cl_uint),
3616
              &compute units, NULL);
3617
          compute_units = BLOCK_SIZE;
3618
          //acopl ando A con B
3619
          af_array Ac;
3620
          af_join(&Ac, 1, dA, dB);
3621
3622
          //vectores de sincronización
3623
          const dim_t nBlks[] = { compute_units };
3624
          af_array syncln;
3625
          af_array syncOut;
          af_constant(&syncIn, 0, 1, nBlks, s32);
3626
3627
          af_constant(&syncOut, 0, 1, nBlks, s32);
3628
3629
          dim_t _order[AF_MAX_DIMS];
3630
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], dA);
3631
          size_t order = _order[0];
3632
          size t local WorkSize = BLOCK SIZE * BLOCK SIZE;
3633
3634
          size_t globalWorkSize = localWorkSize * compute_units;
3635
3636
          int status = CL_SUCCESS;
3637
```

```
3638
          af_dtype typef;
3639
          af_get_type(&typef, dA);
3640
3641
          int msize = 0;
3642
          if (typef == f64)
3643
              msi ze = si zeof(doubl e);
          else if (typef == f32)
3644
3645
              msi ze = si zeof(float);
3646
          el se;
3647
3648
          //obteniendo las referencias cl_mem de los objetos af::array
3649
          cl_mem *d_A = (cl_mem*)clCreateBuffer(af_context,
3650
              CL_MEM_READ_WRITE, msize * order*(order + 1),
3651
              NULL, &status);
3652
          af_get_device_ptr((voi d**)d_A, Ac);
3653
3654
          cl_mem *d_syncin = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, sizeof(int)*compute_units,
3655
3656
              NULL, &status);
3657
          af_get_device_ptr((void**)d_syncIn, syncIn);
3658
3659
          cl_mem *d_syncOut = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_WRITE, sizeof(int)*compute_units,
3660
3661
              NULL, &status);
3662
          af_get_device_ptr((void**)d_syncOut, syncOut);
3663
          size_t program_l ength = strlen(GJordan_source);
3664
3665
3666
          //creando el programa, construyendo el ejecutable y extrayendo el punto
            de entrada
3667
          // para el Kernel
3668
          cl_program program =
3669
              cl CreateProgramWi thSource(af_context, 1,
              (const char **)&GJordan_source, &program_length,
3670
3671
                   &status);
3672
          status = clBuildProgram(program, 1, &af_device_id,
              NULL, NULL, NULL);
3673
3674
3675
          char* kernel Name;
3676
          if (typef == f64)
              kernel Name = "gl obal _sync_test";
3677
3678
          else if (typef == f32)
              kernel Name = "gl obal_sync_test_sp";
3679
3680
          el se:
          cl_kernel kernel = clCreateKernel(program, kernelName,
3681
3682
              &status);
3683
3684
          // estableciendo los argumentos
3685
          int i = 0;
3686
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_A);
3687
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_syncln);
3688
3689
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_syncOut);
3690
3691
          int device;
          af_get_devi ce(&devi ce);
3692
```

```
3693
          for (int | = 0; | < order; | ++)
3694
3695
              cl SetKernel Arg(kernel, 4, sizeof(cl_int), &j);
3696
3697
              //ej ecutando el Kernel
              cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3698
                  &global WorkSize, &local WorkSize, O, NULL,
3699
3700
                  NULL);
              //af_sync(device);
3701
3702
3703
3704
          //devolviendo el control de memoria af::array a ArrayFire
          af_unlock_array(Ac);
3705
3706
3707
          //hasta aqui Ac contiene en su última columna
          //y en su diagonal principal, los valores
3708
3709
          //finales que deben dividirse para obtener
          //la solución final
3710
3711
3712
          //extrayendo la última columna
3713
          af_array Au;
3714
          af_i ndex_t* i ndexers = 0;
3715
          af_create_i ndexers(&i ndexers);
3716
          af_set_seq_param_indexer(indexers, 0, order - 1, 1,
3717
              0, false);
          af_set_seq_param_indexer(indexers, order, order, 1,
3718
3719
              1, false);
          af_index_gen(&Au, Ac, 2, indexers);
3720
3721
3722
          //extrayendo la diagonal
3723
          af_array Ad;
3724
          af_diag_extract(&Ad, Ac, 0);
3725
3726
          //di vi di endo
3727
          af_release_array(Ac);
3728
          af_div(&Ac, Au, Ad, false);
3729
3730
          // copi ando el resultado en dC
3731
          af_copy_array(dC, Ac);
3732
          af_release_array(Au);
3733
3734
          af release array(Ad);
3735
          af_release_array(Ac);
          af_release_i ndexers(i ndexers);
3736
3737
      }
3738
3739 void AFire::sks_util_1(af_array* index, af_array row,
3740
          af_array col) {
3741
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3742
          //cl_context af_context;
3743
          static cl_context af_context = afcl::getContext();
          static cl_device_id af_device_id = afcl::getDeviceId();
3744
3745
          static cl_command_queue af_queue = afcl::getQueue();
3746
3747
          double orderd;
          double helper;
3748
```

```
3749
          af_max_all(&orderd, &helper, col);
3750
          int order = (int)orderd;
3751
          order++:
3752
3753
          af_dtype typef;
3754
          af_get_type(&typef, row);
3755
3756
          af_array idx;
          dim_t d_order[] = { order };
3757
3758
          af_constant(&idx, 0, 1, d_order, typef);
3759
3760
          dim_t _order[AF_MAX_DIMS];
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], row);
3761
3762
          size_t lennz = _order[0];
3763
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
3764
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
3765
3766
3767
          int status = CL_SUCCESS;
3768
3769
          //obteniendo las referencias cl_mem de los objetos af::array
3770
          cl_mem *d_col = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_ONLY, sizeof(int)*Iennz,
3771
3772
              NULL, &status);
3773
          af_get_device_ptr((void**)d_col, col);
3774
          cl_mem *d_row = (cl_mem*)clCreateBuffer(af_context,
3775
              CL_MEM_READ_ONLY, sizeof(int)*Iennz,
3776
3777
              NULL, &status);
3778
          af_get_device_ptr((void**)d_row, row);
3779
3780
          cl_mem *d_i dx = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_WRITE_ONLY, sizeof(int)*order,
3781
3782
              NULL, &status);
3783
          af_get_device_ptr((void**)d_idx, idx);
3784
3785
          size_t program_length = strlen(util_src);
3786
3787
3788
          //creando el programa, construyendo el ejecutable y
3789
          //extrayendo el punto de entrada
3790
          // para el Kernel
3791
          cl_program program =
              cl CreateProgramWi thSource(af_context, 1,
3792
              (const char **) & util_src, & program_length,
3793
3794
                  &status):
          status = cl BuildProgram(program, 1, &af_device_id,
3795
3796
              NULL, NULL, NULL);
3797
3798
          char* kernel Name;
3799
          kernel Name = "sks_util_1";
3800
3801
          cl_kernel kernel = clCreateKernel(program, kernelName,
3802
              &status):
3803
          // estableciendo los argumentos
3804
```

```
... sual Studio 2013\Projects\AF_func\AF_sample.cpp
```

```
70
```

```
3805
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_idx);
3806
3807
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_row);
3808
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_col);
3809
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &lennz);
3810
3811
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
3812
3813
                  &global WorkSize, &local WorkSize, O, NULL,
3814
                  NULL);
3815
3816
          //devolviendo el control de memoria af::array a ArrayFire
          af_unlock_array(col);
3817
3818
          af_unlock_array(row);
3819
3820
          // copi ando el resultado en dC
3821
          af_copy_array(index, idx);
3822
3823
          af_release_array(idx);
3824 }
3825
3826 void AFire::sks_util_2(af_array sks, af_array ptr,
          af_array nz, af_array row, af_array col) {
3827
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3828
3829
              //cl_context af_context;
          static cl_context af_context = afcl::getContext();
3830
          static cl_device_id af_device_id = afcl::getDeviceId();
3831
3832
          static cl_command_queue af_queue = afcl::getQueue();
3833
3834
          double orderd;
3835
          double helper;
3836
          af_max_all(&orderd, &helper, col);
          int order = (int)orderd;
3837
3838
          order++;
3839
          af_dtype typef;
3840
3841
          af_get_type(&typef, nz);
3842
3843
          dim_t _order[AF_MAX_DIMS];
3844
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], row);
3845
          size_t lennz = _order[0];
3846
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], sks);
3847
3848
          size_t lensks = _order[0];
3849
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
3850
          size_t globalWorkSize = localWorkSize * BLOCK_SIZE;
3851
3852
3853
          int status = CL_SUCCESS;
3854
3855
          int msize = 0;
          if (typef == f32)
3856
              msi ze = si zeof(float);
3857
          else if (typef == f64)
3858
3859
              msi ze = si zeof(doubl e);
3860
          el se:
```

```
3861
          //obteniendo las referencias cl_mem de los objetos af::array
3862
3863
          cl_mem *d_sks = (cl_mem*)clCreateBuffer(af_context,
3864
              CL_MEM_WRITE_ONLY, msi ze*lensks,
3865
              NULL, &status);
          af_get_device_ptr((void**)d_sks, sks);
3866
3867
          cl_mem *d_ptr = (cl_mem*)clCreateBuffer(af_context,
3868
3869
              CL_MEM_READ_ONLY, sizeof(int)*order,
3870
              NULL, &status);
3871
          af_get_device_ptr((void**)d_ptr, ptr);
3872
          cl_mem *d_nz = (cl_mem*)clCreateBuffer(af_context,
3873
3874
              CL_MEM_READ_ONLY, msize*lennz,
3875
              NULL, &status);
          af_get_device_ptr((void**)d_nz, nz);
3876
3877
          cl_mem *d_row = (cl_mem*)clCreateBuffer(af_context,
3878
3879
              CL_MEM_READ_ONLY, sizeof(int)*Iennz,
3880
              NULL, &status);
3881
          af_get_device_ptr((void**)d_row, row);
3882
          cl_mem *d_col = (cl_mem*)clCreateBuffer(af_context,
3883
3884
              CL_MEM_READ_ONLY, sizeof(int)*Iennz,
3885
              NULL, &status);
3886
          af_get_device_ptr((void**)d_col, col);
3887
3888
3889
          size_t program_length = strlen(util_src);
3890
3891
3892
          //creando el programa, construyendo el ejecutable y
3893
          //extrayendo el punto de entrada
3894
          // para el Kernel
3895
          cl_program program =
3896
              cl CreateProgramWi thSource(af_context, 1,
              (const char **) & util_src, & program_length,
3897
3898
                  &status);
3899
          status = clBuildProgram(program, 1, &af_device_id,
3900
              NULL, NULL, NULL);
3901
3902
          char* kernel Name;
3903
          if (typef == f32)
              kernel Name = "sks_util_2_sp";
3904
3905
          else if (typef == f64)
              kernelName = "sks_util_2";
3906
3907
3908
          cl_kernel kernel = clCreateKernel (program, kernel Name,
3909
              &status);
3910
3911
          // estableciendo los argumentos
3912
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_sks);
3913
3914
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_ptr);
3915
          cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_nz);
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_row);
3916
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
                                                                                      72
3917
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_col);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &l ennz);
3918
3919
3920
3921
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
              &global WorkSize, &local WorkSize, O, NULL,
3922
3923
              NULL);
3924
3925
          //devolviendo el control de memoria af::array a ArrayFire
3926
          af_unlock_array(sks);
3927
          af_unlock_array(ptr);
3928
          af_unlock_array(nz);
3929
          af_unlock_array(row);
3930
          af_unlock_array(col);
3931 }
3932
3933 void AFire::csc_util_1(af_array* index, af_array row,
3934
          af_array col) {
3935
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
3936
          //cl_context af_context;
3937
          static cl_context af_context = afcl::getContext();
3938
          static cl_device_id af_device_id = afcl::getDeviceId();
          static cl_command_queue af_queue = afcl::getQueue();
3939
3940
3941
          double orderd;
3942
          double helper;
          af_max_all(&orderd, &helper, col);
3943
3944
          int order = (int)orderd;
3945
          order++;
3946
```

af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], row);

//obteniendo las referencias cl_mem de los objetos af::array

3947

3948

3949 3950

3951

3952

39533954

3955

3956

3957

3958

3959 3960

3961 3962 3963

3964

3965

3966

3967

3968 3969

3970

3971

3972

af_dtype typef;

af_array idx;

af_get_type(&typef, row);

dim_t d_order[] = { order };

dim_t _order[AF_MAX_DIMS];

size_t lennz = _order[0];

int status = CL_SUCCESS;

NULL, &status);

NULL, &status);

af_constant(&idx, 0, 1, d_order, typef);

size t local WorkSize = BLOCK SIZE * BLOCK SIZE;

size_t global WorkSize = local WorkSize * BLOCK_SIZE;

cl_mem *d_col = (cl_mem*)clCreateBuffer(af_context,

cl mem *d row = (cl mem*)clCreateBuffer(af context,

CL_MEM_READ_ONLY, sizeof(int)*Iennz,

CL_MEM_READ_ONLY, sizeof(int)*Iennz,

af_get_device_ptr((void**)d_col, col);

af_get_devi ce_ptr((voi d**)d_row, row);

```
3973
          cl_mem *d_i dx = (cl_mem*)clCreateBuffer(af_context,
3974
3975
              CL_MEM_WRITE_ONLY, sizeof(int)*order,
3976
              NULL, &status);
3977
          af_get_device_ptr((void**)d_idx, idx);
3978
          size_t program_length = strlen(util_src);
3979
3980
3981
          //creando el programa, construyendo el ejecutable y
3982
3983
          //extrayendo el punto de entrada
3984
          // para el Kernel
3985
          cl_program program =
              cl CreateProgramWi thSource(af_context, 1,
3986
              (const char **)&util_src, &program_length,
3987
3988
                  &status);
          status = clBuildProgram(program, 1, &af_device_id,
3989
3990
              NULL, NULL, NULL);
3991
3992
          char* kernel Name;
3993
          kernel Name = "csc_util_1";
3994
3995
          cl_kernel kernel = clCreateKernel(program, kernelName,
3996
              &status);
3997
3998
          // estableciendo los argumentos
3999
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_idx);
4000
4001
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_row);
4002
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_col);
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &l ennz);
4003
4004
4005
4006
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
4007
              &global WorkSize, &local WorkSize, O, NULL,
4008
              NULL);
4009
          //devolviendo el control de memoria af::array a ArrayFire
4010
4011
          af_unlock_array(col);
4012
          af_unlock_array(row);
4013
4014
          // copi ando el resultado en dC
4015
          af_copy_array(index, idx);
4016
4017
          af_release_array(idx);
4018 }
4019
4020 void AFire::csc_util_2(af_array* del, af_array index,
4021
          af_array row) {
4022
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
4023
          //cl_context af_context;
4024
          static cl_context af_context = afcl::getContext();
4025
          static cl_device_id af_device_id = afcl::getDeviceId();
4026
          static cl_command_queue af_queue = afcl::getQueue();
4027
          dim_t _order[AF_MAX_DIMS];
4028
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
7
```

```
4029
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], index);
4030
          size_t order = _order[0];
4031
4032
          af_get_dims(&_order[0], &_order[1], &_order[2], &_order[3], row);
4033
          size_t lennz = _order[0];
4034
          af_dtype typef;
4035
          af_get_type(&typef, index);
4036
4037
4038
          af_array idx;
          dim_t d_order[] = { order };
4039
4040
          af_constant(&idx, 0, 1, d_order, typef);
4041
4042
          size t local WorkSize = BLOCK SIZE * BLOCK SIZE:
4043
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
4044
4045
          int status = CL_SUCCESS;
4046
4047
          //obteniendo las referencias cl_mem de los objetos af::array
4048
4049
          cl_mem *d_i ndex = (cl_mem*)clCreateBuffer(af_context,
4050
              CL_MEM_READ_ONLY, sizeof(int)*order,
4051
              NULL, &status);
4052
          af_get_device_ptr((void**)d_index, index);
4053
          cl_mem *d_row = (cl_mem*)clCreateBuffer(af_context,
4054
              CL_MEM_READ_ONLY, sizeof(int)*Iennz,
4055
4056
              NULL, &status);
4057
          af_get_device_ptr((void**)d_row, row);
4058
          cl_mem *d_i dx = (cl_mem*)clCreateBuffer(af_context,
4059
              CL_MEM_WRITE_ONLY, sizeof(int)*order,
4060
4061
              NULL, &status);
4062
          af_get_device_ptr((void**)d_idx, idx);
4063
4064
          size_t program_length = strlen(util_src);
4065
          //creando el programa, construyendo el ejecutable y
4066
4067
          //extrayendo el punto de entrada
4068
          // para el Kernel
4069
          cl_program program =
4070
              cl CreateProgramWi thSource(af_context, 1,
              (const char **) &util_src, &program_length,
4071
4072
                  &status):
          status = cl BuildProgram(program, 1, &af_device_id,
4073
4074
              NULL, NULL, NULL);
4075
4076
          char* kernel Name;
4077
          kernelName = "csc_util_2";
4078
4079
          cl_kernel kernel = clCreateKernel(program, kernelName,
4080
              &status);
4081
          // estableciendo los argumentos
4082
4083
          int i = 0;
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_idx);
4084
```

```
... sual Studio 2013\Projects\AF_func\AF_func\AF_sample.cpp
```

```
cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_index);
4085
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_row);
4086
4087
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
4088
4089
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
4090
              &global WorkSize, &local WorkSize, O, NULL,
4091
4092
              NULL);
4093
4094
          //devolviendo el control de memoria af::array a ArrayFire
4095
          af_unlock_array(index);
4096
          af_unlock_array(row);
4097
4098
          // copi ando el resultado en dC
4099
          af_copy_array(del, idx);
4100
          af_release_array(idx);
4101
4102 }
4103
4104 void AFire::csc_util_3(af_array del, af_array ptr) {
4105
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
4106
          //cl_context af_context;
          static cl_context af_context = afcl::getContext();
4107
4108
          static cl_device_id af_device_id = afcl::getDeviceId();
4109
          static cl_command_queue af_queue = afcl::getQueue();
4110
          dim_t _order[AF_MAX_DIMS];
4111
          af_get_dims(&_order[0], &_order[1], &_order[2],
4112
4113
              &_order[3], ptr);
4114
          size_t order = _order[0];
4115
          size_t localWorkSize = BLOCK_SIZE * BLOCK_SIZE;
4116
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
4117
4118
          int status = CL_SUCCESS;
4119
4120
4121
          //obteniendo las referencias cl_mem de los objetos af::array
4122
4123
          cl_mem *d_del = (cl_mem*)clCreateBuffer(af_context,
4124
              CL_MEM_READ_ONLY, sizeof(int)*order,
4125
              NULL, &status);
4126
          af_get_device_ptr((void**)d_del, del);
4127
4128
          cl_mem *d_ptr = (cl_mem*)clCreateBuffer(af_context,
4129
              CL_MEM_READ_WRITE, sizeof(int)*order,
4130
              NULL, &status);
4131
          af_get_devi ce_ptr((voi d**)d_ptr, ptr);
4132
          size_t program_length = strlen(util_src);
4133
4134
4135
          //creando el programa, construyendo el ejecutable y
          //extrayendo el punto de entrada
4136
4137
          // para el Kernel
4138
          cl_program program =
4139
              cl CreateProgramWi thSource(af_context, 1,
              (const char **) &util_src, &program_length,
4140
```

```
... sual Studio 2013\Projects\AF_func\AF_sample.cpp
```

```
76
```

```
4141
                  &status);
4142
          status = cl BuildProgram(program, 1, &af_device_id,
4143
              NULL, NULL, NULL);
4144
4145
          char* kernel Name;
          kernel Name = "csc_util_3";
4146
4147
4148
          cl_kernel kernel = clCreateKernel (program, kernel Name,
4149
              &status);
4150
4151
          // estableciendo los argumentos
4152
          int i = 0:
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_del);
4153
4154
          cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_ptr);
4155
          cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
          cl SetKernel Arg(kernel, i++,
4156
4157
              sizeof(int)*local WorkSize, 0);
4158
4159
4160
          cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
              &global WorkSize, &local WorkSize, O, NULL,
4161
4162
              NULL);
4163
4164
          //devolviendo el control de memoria af::array a ArrayFire
4165
          af_unlock_array(del);
          af_unl ock_array(ptr);
4166
4167 }
4168
4169
      void AFire::csc_util_4(af_array csc, af_array rowc,
4170
          af_array ptr, af_array nz, af_array row, af_array del) {
4171
          //Obteniendo el dispositivo, contexto y la cola usada por ArrayFire
4172
              //cl context af context;
          static cl_context af_context = afcl::getContext();
4173
4174
          static cl_device_id af_device_id = afcl::getDeviceId();
4175
          static cl_command_queue af_queue = afcl::getQueue();
4176
4177
          dim_t _order[AF_MAX_DIMS];
4178
          af_get_dims(&_order[0], &_order[1], &_order[2],
4179
              &_order[3], ptr);
4180
          size_t order = _order[0];
4181
4182
          af_get_dims(&_order[0], &_order[1], &_order[2],
4183
              &_order[3], row);
4184
          size_t lennz = _order[0];
4185
4186
          af_get_dims(&_order[0], &_order[1], &_order[2],
4187
              &_order[3], csc);
4188
          size_t lencsc = _order[0];
4189
4190
          af_dtype typef;
4191
          af_get_type(&typef, nz);
4192
          size t local WorkSize = BLOCK SIZE * BLOCK SIZE;
4193
4194
          size_t global WorkSize = local WorkSize * BLOCK_SIZE;
4195
          int status = CL_SUCCESS;
4196
```

```
4197
4198
          int msize = 0;
4199
          if (typef == f32)
4200
              msi ze = si zeof(float);
4201
          else if (typef == f64)
4202
              msi ze = si zeof(double);
4203
          el se;
4204
4205
          //obteniendo las referencias cl_mem de los objetos af::array
4206
          cl_mem *d_csc = (cl_mem*)clCreateBuffer(af_context,
4207
              CL_MEM_WRITE_ONLY, msi ze*lencsc,
4208
              NULL, &status);
          af_get_device_ptr((void**)d_csc, csc);
4209
4210
4211
          cl_mem *d_rowc = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_WRITE_ONLY, sizeof(int)*(lencsc - order),
4212
4213
              NULL, &status);
4214
          af_get_device_ptr((void**)d_rowc, rowc);
4215
4216
          cl_mem *d_ptr = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_ONLY, sizeof(int)*order,
4217
4218
              NULL, &status);
4219
          af_get_device_ptr((void**)d_ptr, ptr);
4220
4221
          cl_mem *d_nz = (cl_mem*)clCreateBuffer(af_context,
              CL_MEM_READ_ONLY, msize*lennz,
4222
4223
              NULL, &status);
          af_get_device_ptr((void**)d_nz, nz);
4224
4225
4226
          cl_mem *d_row = (cl_mem*)clCreateBuffer(af_context,
4227
              CL_MEM_READ_ONLY, sizeof(int)*Iennz,
4228
              NULL, &status);
4229
          af_get_device_ptr((void**)d_row, row);
4230
4231
          cl mem *d del = (cl mem*)clCreateBuffer(af context,
4232
              CL_MEM_READ_ONLY, sizeof(int)*order,
              NULL, &status);
4233
          af_get_device_ptr((void**)d_del, del);
4234
4235
4236
          size_t program_length = strlen(util_src);
4237
4238
          //creando el programa, construyendo el ejecutable y
4239
          //extrayendo el punto de entrada
4240
          // para el Kernel
4241
          cl_program program =
4242
              cl CreateProgramWi thSource(af_context, 1,
              (const char **)&util_src, &program_length,
4243
4244
                  &status);
4245
          status = cl BuildProgram(program, 1, &af_device_id,
4246
              NULL, NULL, NULL);
4247
4248
          char* kernel Name;
4249
          if (typef == f32)
              kernelName = "csc_util_4_sp";
4250
4251
          else if (typef == f64)
              kernel Name = "csc_util_4";
4252
```

```
4253
           el se;
4254
           cl_kernel kernel = clCreateKernel (program, kernel Name,
4255
               &status);
4256
           // estableciendo los argumentos
4257
4258
           int i = 0;
4259
           cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_csc);
           cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_rowc);
4260
           cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_ptr);
cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_nz);
4261
4262
4263
           cl SetKernel Arg(kernel, i++, si zeof(cl_mem), d_row);
4264
           cl SetKernel Arg(kernel, i++, sizeof(cl_mem), d_del);
4265
           cl SetKernel Arg(kernel, i++, sizeof(cl_int), &order);
4266
4267
4268
           cl EnqueueNDRangeKernel (af_queue, kernel, 1, 0,
4269
               &global WorkSize, &local WorkSize, O, NULL,
4270
               NULL);
4271
4272
           //devolviendo el control de memoria af::array a ArrayFire
4273
           af_unlock_array(csc);
4274
           af_unl ock_array(ptr);
4275
           af_unlock_array(nz);
           af_unlock_array(row);
4276
4277
           af_unlock_array(del);
4278
      }
4279
```