



Alexis Balayre

Artificial Intelligence Assignment

School of Aerospace, Transport and Manufacturing
Computational Software of Techniques Engineering

MSc
Academic Year: 2023 - 2024

Supervisor: Dr Jun Li
18th March 2024

Abstract

High-Performance Computing (HPC) systems are pivotal in solving complex computational problems. Leveraging such systems, this report investigates the optimization of sparse matrix-vector multiplication (SpMV) - a crucial operation in scientific computations. Two prevalent storage formats, Compressed Sparse Row (CSR) and ELLPACK, are parallelized using OpenMP and CUDA to enhance SpMV's efficiency on the CRES-CENT2 HPC cluster at Cranfield University.

The study reveals that CSR format, when parallelized with OpenMP, achieves superior performance across a majority of matrices due to efficient memory management and task distribution. CUDA parallelization exhibits significant speed-ups, especially for matrices with regular structures, indicating an intricate relationship between matrix properties and the efficiency of CSR and ELLPACK formats.

Experimental results suggest that the parallelization strategy should be selected based on matrix characteristics. For matrices with high density or regularity, CUDA outperforms OpenMP, whereas for less regular matrices, OpenMP provides sufficient speed-up. Sequential execution remains competitive for small matrices or those with unfavourable characteristics for parallelization.

This report underscores the necessity of aligning parallel programming strategies with matrix properties to fully exploit HPC capabilities, providing a foundation for future research towards more nuanced and effective computational techniques in scientific computing.

Table of Contents

Abstract	ii
Table of Contents	iii
List of Figures	vi
List of Tables	vii
1 Introduction	1
1.1 Dataset	1
1.2 Significance of the Competition	1
1.3 Evaluation Metric	1
2 literature Review	3
2.1 Basic principles	3
2.2 Object Detection Technologies	3
2.2.1 Convolutional Neural Networks (CNN)	3
2.2.2 Region-based approaches	3
2.2.2.1 R-CNN	3
2.2.2.2 Fast R-CNN	3
2.2.2.3 Faster R-CNN	4
2.2.3 Single Processing Methods	4
2.2.3.1 YOLO (You Only Look Once)	4
2.2.3.2 SSD (Single Shot MultiBox Detector)	4
2.3 Ethical Considerations	4
3 Methodology	5
3.1 Problem Statement	5
3.2 Data Structures	5
3.2.1 Sparse Matrix	5
3.2.1.1 Compressed Sparse Row (CSR) Format	5
3.2.1.2 ELLPACK	6
3.2.2 Fat Vector	6
3.3 Sparse Matrix - Fat Vector Multiplication	7
3.3.1 Sequential Algorithm using CSR Format	7
3.3.1.1 Algorithm Flow	7
3.3.1.2 Temporal Complexity Analysis	8

3.3.2	Sequential Algorithm using ELLPACK Format	9
3.3.2.1	Algorithm Flow	9
3.3.2.2	Temporal Complexity Analysis	10
3.4	Parallel Algorithms	11
3.4.1	OpenMP (Open Multi-Processing)	11
3.4.1.1	Workflow and Optimisation Strategy	11
3.4.1.2	Expected Performance Improvements	12
3.4.1.3	Challenges and Considerations	12
3.4.2	CUDA	13
3.4.2.1	Workflow and Kernel Design	13
3.4.2.2	Expected Performance Improvements	13
3.4.2.3	Challenges and Considerations	14
3.5	Implementation and Testing	15
3.5.1	Workflow	15
3.5.2	Comprehensive Testing Strategy	15
3.5.2.1	Test Parameters Configuration	15
3.5.2.1.1	OpenMP Parameters:	15
3.5.2.1.2	CUDA Parameters:	16
3.5.2.2	Correctness Verification	16
3.5.2.3	Performance Evaluation	16
4	Results and Discussion	18
4.1	Results	18
4.1.1	Sequential Algorithms	18
4.1.2	OpenMP	20
4.1.2.1	Chunk Size Analysis	20
4.1.2.2	Thread Count Analysis	20
4.1.2.3	Performance Comparison	20
4.1.3	CUDA	23
4.1.3.1	X Block Size Analysis	23
4.1.3.2	Y Block Size Analysis	23
4.1.3.3	Performance Comparison	23
4.1.4	Overall Performance Comparison	25
5	Conclusion	27
	References	28
A	Documentation	29
A.A	Project tree	29
A.B	Getting Started	30
A.C	Methods Overview	30
A.C.1	Utils.h	30
A.C.1.1	convertCRStoELLPACK	30
A.C.1.2	areMatricesEqual	30
A.C.1.3	readMatrixMarketFile	30
A.C.1.4	generateLargeFatVector	31

A.C.2	matrixMultivectorProductCRS.h	31
A.C.2.1	matrixMultivectorProductCRS	31
A.C.3	matrixMultivectorProductCRSOpenMP.h	31
A.C.3.1	matrixMultivectorProductCRSOpenMP	31
A.C.4	matrixMultivectorProductCRSCUDA.h	32
A.C.4.1	matrixMultivectorProductCRSCUDA	32
A.C.5	matrixMultivectorProductELLPACK.h	32
A.C.5.1	matrixMultivectorProductELLPACK	32
A.C.6	matrixMultivectorProductELLPACKOpenMP.h	32
A.C.6.1	matrixMultivectorProductELLPACKOpenMP	32
A.C.7	matrixMultivectorProductELLPACKCUDA.h	33
A.C.7.1	matrixMultivectorProductELLPACKCUDA	33

B Source Codes 34

List of Figures

List of Tables

3.1	Summary of sparse matrices (?)	17
4.1	Performance Comparison of CRS and ELLPACK Sequential Algorithms .	18
4.2	CRS vs ELLPACK using OpenMP	20
4.3	CRS vs ELLPACK using CUDA	24
4.4	Overall Performance Comparison	25

Chapter 1

Introduction

Accurately diagnosing thoracic abnormalities from radiographs (X-rays) represents a considerable challenge, even for experienced radiologists. The complexity and critical nature of diagnosing these anomalies requires increasingly sophisticated decision support tools. In this context, the competition organised by Vingroup's Big Data Institute, supported by Vingroup JSC and launched in August 2018, aims to promote fundamental research in data science and artificial intelligence, with a particular focus on medical image processing.

The main objective of the competition is to develop automated systems capable of locating and classifying 14 types of thoracic anomalies from chest X-ray images. This initiative underlines the need for greater precision in medical diagnosis, where current methods struggle in particular to specify the location of findings on X-ray images, potentially leading to incorrect diagnoses.

1.1 Dataset

The dataset provided to participants includes 18,000 annotated chest scans, of which 15,000 images are for training and 3,000 for evaluation. These annotations were carefully collected via VinBigData's VinLab platform from de-identified studies provided by two Vietnamese hospitals.

1.2 Significance of the Competition

This competition promises significant advances in medical diagnosis by potentially providing radiologists with a reliable secondary opinion. By automating the detection and location of chest X-ray findings, the solution aims to lighten the workload of healthcare professionals and improve diagnostic accuracy for patients, particularly benefiting those in resource-limited settings.

1.3 Evaluation Metric

The competition uses the PASCAL VOC 2010 Mean Average Precision (mAP) metric with an Intersection on Union (IoU) threshold of ≥ 0.4 to evaluate submissions. This metric

highlights the importance of accuracy in the detection and classification process.

Chapter 2

literature Review

Object detection is a fundamental task in computer vision that involves identifying and locating objects of different categories in an image or video. Unlike image classification, which assigns a label to the entire image, object detection aims to provide a label and bounding box for each object of interest in the image.

2.1 Basic principles

Object detection generally involves two main tasks: object classification (knowing what objects are) and object localisation (knowing where objects are). To be successful, an object detection system must be able to recognise objects under a variety of conditions, such as different sizes, viewing angles, and occlusion levels.

2.2 Object Detection Technologies

2.2.1 Convolutional Neural Networks (CNN)

CNNs are at the heart of many advances in object detection. They are particularly effective at extracting hierarchical features from images thanks to their layered structure, which includes convolutional layers, pooling layers and fully connected layers.

2.2.2 Region-based approaches

2.2.2.1 R-CNN

R-CNN (Regions with CNN features) uses proposed regions to identify potentially interesting parts of the image, then applies a CNN to each of these proposed regions to classify the objects.

2.2.2.2 Fast R-CNN

Fast R-CNN improves on R-CNN by using a more efficient architecture that shares convolution calculations across the entire image, reducing processing time.

2.2.2.3 Faster R-CNN

Faster R-CNN introduces a Region Proposal Network (RPN) that generates region proposals directly from image features, further improving detection speed and accuracy.

2.2.3 Single Processing Methods

2.2.3.1 YOLO (You Only Look Once)

YOLO divides the image into a grid and predicts bounding boxes and class probabilities for each grid cell in a single pass, offering high processing speed.

2.2.3.2 SSD (Single Shot MultiBox Detector)

SSD combines the advantages of region-based approaches and single-shot methods, using bounding boxes at different scales and aspect ratios to predict the presence of objects in the image.

2.3 Ethical Considerations

Chapter 3

Methodology

3.1 Problem Statement

Consider a sparse matrix A of dimensions $m \times n$ and a fat vector X of dimensions $n \times k$. The objective is to perform the multiplication $A \times X$, yielding a result that is of dimensions $m \times k$.

The matrix A is defined as:

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix} \quad (3.1)$$

where most elements of A are zeros.

The vector X is defined as:

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \quad (3.2)$$

3.2 Data Structures

3.2.1 Sparse Matrix

Two data structures were used to represent sparse matrices: Compressed Sparse Row (CSR) and ELLPACK formats. Both formats are designed to store and manipulate sparse matrices efficiently, reducing the storage space and computational load associated with zeros in the matrix.

3.2.1.1 Compressed Sparse Row (CSR) Format

The CSR format represents sparse matrices efficiently by storing only the non-zero elements and their positions. This format uses three arrays:

- **values**: An array holding all the non-zero elements of the matrix, stored sequentially as they appear in the matrix from top to bottom and left to right.
- **colIndices**: An array storing the column indices of each non-zero element in the *values* array, indicating the exact column position of each element.
- **rowPtr**: An array where each entry marks the starting point in the *values* and *colIndices* arrays for the non-zero elements of a specific row, facilitating direct access to each row's data.

The SparseMatrixCRS structure is defined in appendix ??.

3.2.1.2 ELLPACK

The ELLPACK format optimises the storage and computation for matrices with a relatively uniform distribution of non-zero elements per row. It employs two 2D arrays for this purpose:

- **values**: A 2D array where each row corresponds to a row in the original sparse matrix, containing the non-zero values. Rows are padded with zeros to equalize the length across all rows, determined by `maxNonZerosPerRow`.
- **colIndices**: A 2D array parallel to *values*, storing the column indices for each non-zero value in the matrix. It uses padding (typically with an invalid index such as -1) to match the structure of *values*.
- **maxNonZerosPerRow**: Specifies the fixed number of elements in each row of *values* and *colIndices*, determined by the row with the maximum number of non-zero elements.
- **numRows** and **numCols**: Indicate the dimensions of the matrix, specifically the total number of rows and columns.

ELLPACK format is particularly advantageous for parallel computations on GPUs due to its consistent row length, which enables efficient memory access patterns and simplifies parallelization strategies. The SparseMatrixELLPACK structure is defined in appendix ??.

3.2.2 Fat Vector

The structure FatVector is designed to represent fat vectors or low-dimensional dense matrices:

- **values** : Contains the values of the fat vector or matrix.
- **numRows** and **numCols** : Indicate the dimensions of the vector or matrix, making it possible to represent both column vectors and matrices with several columns.

The FatVector structure is defined in appendix ??.

3.3 Sparse Matrix - Fat Vector Multiplication

3.3.1 Sequential Algorithm using CSR Format

Let A be a sparse matrix of size $m \times n$ with NZ non-zero elements, stored in CSR format, and X be a fat vector of size $n \times k$. The sequential algorithm for multiplying A by X using the CSR format is implemented in Appendix ??.

3.3.1.1 Algorithm Flow

Algorithm 1 Sparse Matrix-Fat Vector Multiplication (CRS)

Require: A is an $m \times n$ sparse matrix in CRS format

Require: X is an $n \times k$ fat vector

Ensure: Y is an $m \times k$ fat vector, result of $A \times X$

```

for  $i = 0$  to  $A.numRows - 1$  do
  for  $j = A.rowPtr[i]$  to  $A.rowPtr[i + 1] - 1$  do
     $colIndex \leftarrow A.colIndices[j]$ 
     $value \leftarrow A.values[j]$ 
    for  $k = 0$  to  $X.numCols - 1$  do
       $yIndex \leftarrow i \times X.numCols + k$ 
       $xIndex \leftarrow colIndex \times X.numCols + k$ 
       $Y.values[yIndex] \leftarrow Y.values[yIndex] + value \times X.values[xIndex]$ 
    end for
  end for
end for

```

The algorithmic flow can be more explicitly detailed by:

1. **Initialisation:** Prepare the result vector with the same number of rows as the sparse matrix and initialise all elements to zero.
2. **Iteration Over Rows:** For each row in the sparse matrix, use the rowPtr array to find the starting and ending indices of non-zero elements in that row.
3. **Iteration Over Non-Zero Elements:** For each non-zero element identified in the previous step, retrieve the column index and value from colIndices and values arrays, respectively.
4. **Multiplication and Accumulation:** Use the column index to identify the corresponding element(s) in the fat vector. Multiply each non-zero element by the corresponding element in the fat vector and accumulate the product in the appropriate position of the result vector. This step is repeated for each column in the fat vector if it has more than one column.
5. **Result Compilation:** After iterating through all rows and their non-zero elements, the result vector contains the product of the sparse matrix and the fat vector.

3.3.1.2 Temporal Complexity Analysis

Given a sparse matrix A of size $m \times n$ with NZ non-zero elements and a fat vector X of size $n \times k$, the serial algorithm for multiplying $A \times X$ iterates through each non-zero element of the matrix A to compute the product.

The algorithm performs two operations (a multiplication and an addition) for each non-zero element with respect to each column of X , resulting in a total of $2 \times NZ \times k$ operations.

Hence, the time complexity of the sparse matrix-fat vector multiplication algorithm can be expressed as:

$$T(n) = O(NZ \times k) \quad (3.3)$$

The CRS format is more efficient in terms of computation when the distribution of non-zero elements is irregular, as it only iterates over these elements. However, it may not be as efficient for parallel processing due to the irregular memory access patterns.

3.3.2 Sequential Algorithm using ELLPACK Format

Let A be a sparse matrix of size $m \times n$ with NZ non-zero elements, stored in CSR format, and X be a fat vector of size $n \times k$. The sequential algorithm for multiplying A by X using the ELLPACK format is implemented in Appendix ??.

3.3.2.1 Algorithm Flow

Algorithm 2 Sparse Matrix-Fat Vector Multiplication (ELLPACK)

Require: A is an $m \times n$ sparse matrix in ELLPACK format

Require: X is an $n \times k$ fat vector

Ensure: Y is an $m \times k$ fat vector, result of $A \times X$

```

for  $i = 0$  to  $A.numRows - 1$  do
  for  $j = 0$  to  $A.maxNonZerosPerRow - 1$  do
     $colIndex \leftarrow A.colIndices[i \times A.maxNonZerosPerRow + j]$ 
     $value \leftarrow A.values[i \times A.maxNonZerosPerRow + j]$ 
    if  $colIndex \neq -1$  then
      for  $k = 0$  to  $X.numCols - 1$  do
         $yIndex \leftarrow i \times X.numCols + k$ 
         $xIndex \leftarrow colIndex \times X.numCols + k$ 
         $Y.values[yIndex] \leftarrow Y.values[yIndex] + value \times X.values[xIndex]$ 
      end for
    end if
  end for
end for

```

The algorithmic flow can be more explicitly detailed by:

1. **Initialisation:** Similarly, prepare the result vector with an appropriate size and initialise all values to zero.
2. **Iteration Over Rows:** Iterate over each row of the sparse matrix, given that the ELLPACK format stores a fixed number of elements (equal to the maximum number of non-zero elements in any row) for every row.
3. **Iteration Over Elements:** For each element in a row (up to the maximum number of non-zero elements per row), check if the column index is valid (not a padding indicator, such as -1).
4. **Multiplication and Accumulation:** For each valid non-zero element, perform multiplication with the corresponding element(s) in the fat vector, similar to the CRS algorithm. This involves using the column index to locate the correct element in the fat vector and accumulating the product in the result vector.
5. **Handling Padding:** Ignore any padding indicators (e.g., column index -1) during multiplication to ensure that they do not affect the result.
6. **Result Compilation:** After processing all rows, the result vector is fully populated with the product of the sparse matrix in ELLPACK format and the fat vector.

3.3.2.2 Temporal Complexity Analysis

Given a sparse matrix A of size $m \times n$ with NZ_{\max} as the maximum number of non-zero elements in any row and a fat vector X of size $n \times k$. The sequential algorithm for multiplying $A \times X$ iterates over each row and each column index/value pair within the row, resulting in a total of $2 \times m \times NZ_{\max} \times k$ operations.

Hence, the time complexity of the sparse matrix-fat vector multiplication algorithm can be expressed as:

$$T(n) = O(m \times NZ_{\max} \times k) \quad (3.4)$$

The ELLPACK format can lead to faster execution times in architectures that favour regular memory access patterns, despite potentially higher computational complexity due to padding, particularly when the sparse matrix is relatively dense or the maximum number of non-zero elements per row is close to the average number over all rows. However, its efficiency decreases with increasing filling (i.e. when the difference between the maximum and average number of non-zero elements per row is large).

3.4 Parallel Algorithms

3.4.1 OpenMP (Open Multi-Processing)

OpenMP provides a powerful framework for parallelising computational tasks in shared-memory architectures. When applied to sparse matrix-vector multiplication, OpenMP enables significant performance enhancements for both CSR and ELLPACK formats by distributing the computation across multiple CPU threads. The parallel algorithms for both formats are implemented in the appendix ??.

3.4.1.1 Workflow and Optimisation Strategy

The application of OpenMP in sparse matrix-vector multiplication encompasses a series of steps designed to efficiently parallelise the computation and optimise resource utilisation:

1. **Memory Allocation and Data Initialization:** Initially, the sparse matrix and fat vector are allocated in memory. OpenMP does not require explicit memory allocation on a separate device but operates directly on data in the process's address space.
2. **Kernel Execution:**
 - *For CSR Format:* An OpenMP parallel for loop iterates over the rows of the matrix, where each thread processes a subset of rows, computing dot products between the non-zero elements and the corresponding entries of the fat vector (See Appendix ??).
 - *For ELLPACK Format:* Similarly, an OpenMP parallel for loop is employed, with threads iterating over rows. The fixed-length rows of the ELLPACK format potentially offer more regular memory access patterns, which can be advantageous for performance (See Appendix ??).

Both implementations utilise `#pragma omp parallel` for directives, with dynamic scheduling to manage workload distribution among threads.

3. **Performance Measurement:** The execution time for the parallel region is captured using `omp_get_wtime()`, allowing for the calculation of performance metrics such as execution time and GFLOPS (Giga Floating-Point Operations Per Second).
4. **Data Retrieval and Cleanup:** Upon completion of the parallel computation, the result vector is available for further processing. Unlike CUDA, there is no need for explicit data transfer between host and device memory spaces.
5. **Resource Management:** OpenMP abstracts much of the resource management, simplifying the parallelization process. However, developers should still consider the optimal number of threads and scheduling strategies to maximise performance.

3.4.1.2 Expected Performance Improvements

Leveraging OpenMP for sparse matrix-vector multiplication offers potential for substantial performance gains, attributed to:

- **Parallel Processing:** Utilising multiple CPU cores to perform computations in parallel significantly reduces overall execution time.
- **Efficient Workload Distribution:** The ability to dynamically schedule work among threads can lead to more balanced computation, especially for matrices with uneven distributions of non-zero elements.
- **Reduced Overhead:** Compared to CUDA, OpenMP operates within the existing memory space of the application, eliminating the overhead associated with data transfer between host and device.

3.4.1.3 Challenges and Considerations

Despite the advantages, several considerations must be addressed to optimise performance:

- **Thread Overhead:** The benefits of adding more threads diminish beyond a certain point, where the overhead of thread management can outweigh performance gains.
- **Memory Access Patterns:** Especially relevant for the ELLPACK format, ensuring efficient memory access patterns can enhance performance.
- **Optimal Use of Resources:** Balancing the computational load across available CPU cores and selecting the appropriate chunk size for dynamic scheduling are key factors in optimizing performance.

3.4.2 CUDA

CUDA (Compute Unified Device Architecture) is a parallel programming model developed by NVIDIA. It enables graphics processing units (GPUs) to be used for general-purpose computing outside the traditional graphics context. With CUDA, developers can exploit the massively parallel computing power of NVIDIA GPUs for computationally intensive computing more efficiently than with traditional CPU approaches. The parallel algorithms for both CSR and ELLPACK formats are implemented in the appendix ??.

3.4.2.1 Workflow and Kernel Design

The CUDA-based implementation follows a structured workflow, incorporating specialised kernels to exploit parallel processing capabilities of GPUs:

1. **Memory Allocation and Data Transfer:** Initially, necessary memory for matrix data (values, column indices, and row pointers or ELLPACK structures), the fat vector, and the result vector is allocated on the GPU. Data is then transferred from the host to these allocated spaces.
2. **Kernel Execution:**
 - *CSR Kernel:* Processes the sparse matrix in CSR format. Each thread calculates a single element of the result vector by iterating over non-zero elements of a particular row, multiplying each by the corresponding vector element, and accumulating the result (See Appendix ??).
 - *ELLPACK Kernel:* Similarly, processes the matrix in ELLPACK format. Threads are assigned to matrix rows, with each thread iterating over the fixed-size row data, performing multiplication and accumulation operations (See Appendix ??).

Both kernels use atomic operations to safely add results to the output vector in cases of potential write conflicts.

3. **Performance Measurement:** CUDA events are used to record the start and stop times of kernel execution, allowing for precise measurement of computation time and subsequent performance analysis.
4. **Data Retrieval:** After kernel execution, the resulting vector is transferred back to the host for further processing or analysis.
5. **Resource Cleanup:** Finally, allocated GPU memory is freed, and CUDA events are destroyed to clean up resources.

3.4.2.2 Expected Performance Improvements

CUDA parallelization offers significant speedups for sparse matrix-vector multiplication by leveraging the massive parallelism of GPU cores. Performance gains are realised through:

- **Fine-grained Parallelism:** Each non-zero element or row of the matrix can be processed in parallel, significantly reducing computation time.
- **Memory Bandwidth Utilisation:** Efficient use of memory bandwidth by coalescing memory accesses and minimising global memory transactions.
- **Load Balancing:** Dynamic assignment of work to threads helps mitigate performance degradation due to imbalanced non-zero element distribution across rows.

3.4.2.3 Challenges and Considerations

While CUDA accelerates sparse matrix operations, developers must consider factors such as the sparsity pattern of the matrix, the optimal configuration of CUDA threads and blocks, and the overhead of memory transfers between host and device. Proper tuning and optimization strategies, including choosing appropriate block sizes and utilizing shared memory, are crucial for maximizing performance on GPUs.

3.5 Implementation and Testing

In order to implement all the algorithms and test them, 2 main programs were created:

- `runOpenMP.cpp` to test the OpenMP implementations (See Appendix ??).
- `runCuda.cpp` to test the CUDA implementations (See Appendix ??).

3.5.1 Workflow

Both programs follow the same workflow:

1. Reading matrices from files in Matrix Market format (??).
2. Conversion of matrices from CRS format to ELLPACK format.
3. Generation of fat vectors with random values for multiplication.
4. Perform matrix-vector multiplications using both sequential and parallel algorithms.
5. Validate the results of the parallel algorithms against the sequential ones.
6. Measuring and displaying performance.

3.5.2 Comprehensive Testing Strategy

To validate and benchmark the efficiency of parallel implementations using OpenMP and CUDA for sparse matrix-vector multiplication, a systematic testing strategy is employed. This strategy encompasses a range of test parameters, correctness verification, and performance evaluation methods.

3.5.2.1 Test Parameters Configuration

3.5.2.1.1 OpenMP Parameters:

- **Matrix Sparsity:** Varied densities of hollow matrices are tested to simulate real-world scenarios and understand performance across different sparsity levels.
- **Vector Sizes:** Fat vector sizes are varied (1, 2, 3, 6) to assess the impact on performance, reflecting different use cases.
- **Thread Count:** Tests are conducted with 1 to 16 threads to identify the optimal concurrency level that maximizes performance.
- **Chunk Sizes:** To fine-tune dynamic workload distribution, chunk sizes are varied (2, 4, 8, 16, 32, 64, 128, 256), allowing for in-depth analysis of the parallel loop scheduling efficiency.

3.5.2.1.2 CUDA Parameters:

- **Matrix Sparsity:** Similar to OpenMP, different densities of hollow matrices are evaluated to gauge CUDA's effectiveness across varying sparsity patterns.
- **Vector Sizes:** A range of fat vector sizes are tested to examine CUDA's adaptability to different data scales.
- **X Block Size:** CUDA block sizes along the X dimension are varied (8, 16, 32, 64) to optimize thread block configuration for the GPU architecture.
- **Y Block Size:** Similarly, Y block sizes (1, 2, 4, 8, 16) are adjusted to explore the impact on parallel execution efficiency.

The sparse matrices used for testing are summarised in Table 3.1.

3.5.2.2 Correctness Verification

To ensure the accuracy of both OpenMP and CUDA parallel implementations, a two-step verification process is adopted:

1. **Result Comparison:** The outcomes of the sequential (baseline) and parallel implementations are compared using the `areMatricesEqual` function. This function evaluates if the two matrices are identical within a predefined tolerance level, ensuring computational integrity.
2. **Tolerance Threshold:** A small tolerance is allowed for floating-point operations to account for numerical precision variances inherent in parallel computations.

3.5.2.3 Performance Evaluation

Performance testing is structured to provide a comprehensive understanding of the efficiency gains from parallelization:

1. **Execution Time Measurement:** The time taken by each implementation to complete the matrix-vector multiplication is precisely measured, using built-in timing functionalities like `omp_get_wtime()` for OpenMP and CUDA events for GPU measurements.
2. **GFLOPS Calculation:** Performance is quantitatively compared in terms of Giga Floating-Point Operations Per Second (GFLOPS), offering a normalized metric to evaluate computational speed. The formula used for GFLOPS calculation is:

$$\text{GFLOPS} = \frac{2 \times \text{NZ} \times k}{\text{Execution Time} \times 10^9} \quad (3.5)$$

where NZ is the number of non-zero elements in the sparse matrix, k is the number of columns in the fat vector and the execution time is in seconds.

3. **Repetitions for Accuracy:** Each test configuration is executed 20 times. This repetition ensures statistical significance, allowing for the calculation of average execution times and minimizing the impact of outliers.

Through this comprehensive testing approach, the parallel implementations' correctness and performance are meticulously evaluated, ensuring reliability and efficiency of the sparse matrix-vector multiplication algorithms.

Matrix Name	m	NZ	AVG NZR	MAX NZR	Symmetric
cavity10	2597	76367	29.4	62	False
PR02R	161070	8185136	50.8	92	False
nlpkkt80	1062400	28704672	27.0	28	True
Cube_Coup_dt0	2164760	127206144	58.8	68	True
roadNet-PA	1090920	3083796	2.8	9	True
ML_Laplace	377002	27689972	73.4	74	False
bcsstk17	10974	428650	39.1	150	True
mhda416	416	8562	20.6	33	False
af_1_k101	503625	17550675	34.8	35	True
thermal1	82654	574458	7.0	11	True
thermomech_TK	102158	711558	7.0	10	True
cage4	9	49	5.4	6	False
cant	62451	4007383	64.2	78	True
dc1	116835	766396	6.6	114190	False
raefsky2	3242	294276	90.8	108	False
rdist2	3198	56934	17.8	61	False
mcfe	765	24382	31.9	81	False
olm1000	1000	3996	4.0	6	False
lung2	109460	492564	4.5	8	False
webbase-1M	1000005	3105536	3.1	4700	False
mhd4800a	4800	102252	21.3	33	False
west2021	2021	7353	3.6	12	False
thermal2	1228045	8580313	7.0	11	True
adder_dcop_32	1813	11246	6.2	100	False
mac_econ_fwd500	206500	1273389	6.2	44	False
FEM_3D_thermal1	17880	430740	24.1	27	False
amazon0302	262111	1234877	4.7	5	False
cop20k_A	121192	2624331	21.7	81	True
olafu	16146	1015156	62.9	89	True
af23560	23560	484256	20.6	21	False

Table 3.1: Summary of sparse matrices (?)

Chapter 4

Results and Discussion

4.1 Results

4.1.1 Sequential Algorithms

The following table shows the performance comparison of the sequential algorithms using the CRS and ELLPACK formats.

Table 4.1: Performance Comparison of CRS and ELLPACK Sequential Algorithms

Matrix	CRS	ELLPACK	Best Structure
Cube_Coup_dt0	6.313540	2.659870	CRS
FEM_3D_thermal1	5.939650	2.543850	CRS
ML_Laplace	6.629070	2.775640	CRS
PR02R	6.577900	2.755250	CRS
adder_dcop_32	4.481260	2.360730	CRS
af23560	5.804440	2.509250	CRS
af_1_k101	6.139610	2.608660	CRS
amazon0302	2.344020	0.968039	CRS
bcsstk17	6.256790	2.654560	CRS
cage4	0.871628	0.849791	CRS
cant	6.419710	2.687030	CRS
cavity10	6.032170	2.625250	CRS
cop20k_A	4.771770	2.155300	CRS
dc1	4.686860	2.349130	CRS
lung2	4.634780	2.319800	CRS
mac_econ_fwd500	3.684580	1.851230	CRS
mcfe	6.069120	2.634280	CRS
mhd4800a	5.474340	2.549000	CRS
mhda416	5.524470	2.516310	CRS
nlpkkt80	5.932140	2.567120	CRS
olafu	6.481010	2.724440	CRS
olm1000	4.775320	2.192310	CRS

Matrix	CRS	ELLPACK	Best Structure
raefsky2	6.556030	2.767780	CRS
rdist2	5.762990	2.465300	CRS
roadNet-PA	2.721200	1.240050	CRS
thermal1	4.599780	2.228430	CRS
thermal2	3.206760	1.437750	CRS
thermomech_TK	3.558870	1.592880	CRS
webbase-1M	3.750440	1.888800	CRS
west2021	4.184400	2.181630	CRS

A few key observations can be made:

- **CRS Dominance:** For the majority of matrices tested, the CRS format systematically outperforms the ELLPACK format. This is evident for matrices such as Cube_Coup_dt0, FEM_3D_thermal1, and ML_Laplace, where CRS not only handles large matrices efficiently but also matrices with high average and maximum numbers of non-zero elements per row.
- **Efficiency in Dense Matrices:** The CRS format appears to be particularly efficient at handling high-density matrices (higher AVG NZR and MAX NZR), which could be attributed to its storage efficiency and the way it streamlines the multiplication process for rows with varying lengths of non-zero elements.

4.1.2 OpenMP

4.1.2.1 Chunk Size Analysis

As shown in Figures ?? and ??, the performance of the OpenMP parallel algorithms is influenced by the chunk size. For lower density matrices, the optimal chunk size is 64 for the CRS format and 32 for the ELLPACK format. For higher density matrices, the optimal chunk size is 128 or 256 for both formats.

4.1.2.2 Thread Count Analysis

As shown in Figures ?? and ??, the performance of the OpenMP parallel algorithms is influenced by the number of threads. As more threads are used, the performance improves up to a certain point, after which the performance starts to stagnate due to the overhead of managing the threads.

4.1.2.3 Performance Comparison

Table 4.2: CRS vs ELLPACK using OpenMP

Matrix	CRS	ELLPACK	Best Structure	Speedup
Cube_Coup_dt0	29.0243	27.0415	CRS	4.597152
FEM_3D_thermal1	22.5277	19.9711	CRS	3.792766
ML_Laplace	33.2766	31.9121	CRS	5.019799
PR02R	32.5086	30.2968	CRS	4.942094
adder_dcop_32	5.12204	4.83417	CRS	1.142991
af23560	21.7606	18.9078	CRS	3.748958
af_1_k101	26.067	23.4946	CRS	4.245709
amazon0302	7.29883	5.23835	CRS	3.113809
bcsstk17	25.5293	23.3597	CRS	4.080255
cage4	0.385164	0.334653	CRS	0.441890
cant	29.7794	27.6956	CRS	4.638745
cavity10	19.0148	16.7258	CRS	3.152232
cop20k_A	19.5918	17.8729	CRS	4.105772
dc1	10.4524	9.39085	CRS	2.230150
lung2	8.832	6.85328	CRS	1.905592
mac_econ_fwd500	10.5515	8.84879	CRS	2.863691
mcfe	12.6097	11.8456	CRS	2.077682
mhd4800a	18.3501	16.1582	CRS	3.352021
mhda416	7.06099	6.80262	CRS	1.278130
nlpkkt80	24.4853	21.694	CRS	4.127566
olafu	29.966	27.9111	CRS	4.623662
olm1000	2.95895	2.1558	CRS	0.619634
raefsky2	27.8548	25.7504	CRS	4.248730
rdist2	15.5283	12.2527	CRS	2.694487

Matrix	CRS	ELLPACK	Best Structure	Speedup
roadNet-PA	4.63584	3.01227	CRS	1.703601
thermal1	11.7254	8.83925	CRS	2.549122
thermal2	8.84723	6.75174	CRS	2.758931
thermomech_TK	10.9908	8.10205	CRS	3.088284
webbase-1M	6.11469	5.10783	CRS	1.630393
west2021	3.89256	2.92762	CRS	0.930255

The performance evaluation of the CRS and ELLPACK formats, when parallelized with OpenMP, shows a preference for CRS in a majority of cases. This trend, indicated by a frequent superiority of the CRS format, highlights its suitability for OpenMP's parallelization capabilities, probably due to more optimal memory management and task distribution between threads.

The observed speed-up factor varies significantly between matrices, with particularly remarkable speed-ups for matrices such as *Cube_Coup_dt0* and *FEM_3D_thermal1*, highlighting the potential for improving performance via OpenMP's parallel optimisation.

Special cases such as *cage4* and *olm1000* show lower speed-ups, revealing the limits of parallelization for certain matrix structures. The influence of matrix density and size is also notable, with high densities and large fat vectors favouring the CRS format more, illustrating its effectiveness in processing large workloads under OpenMP.

4.1.3 CUDA

4.1.3.1 X Block Size Analysis

As shown in Figures ??, ??, the performance of the CUDA parallel algorithms is influenced by the X block size. For the CRS format, the optimal X block size is 16 most matrices, while for the ELLPACK format, the optimal X block size is 16 or even 8 for some matrices with higher density.

4.1.3.2 Y Block Size Analysis

As shown in Figures ?? and ??, the performance of the CUDA parallel algorithms is influenced by the Y block size. For the CRS format, the optimal Y block size is 16 for most matrices, while for the ELLPACK format, the optimal Y block size is 16 or 4 for some matrices with higher density.

4.1.3.3 Performance Comparison

Table 4.3: CRS vs ELLPACK using CUDA

Matrix	CRS	ELLPACK	Best Structure	Speedup
amazon0302	12.3120	10.4822	CRS	5.252515
cage4	0.033586	0.031649	CRS	0.038533
lung2	10.2714	9.81548	CRS	2.216157
mac_econ_fwd500	15.6609	15.4014	CRS	4.250389
mhd4800a	21.9362	21.4274	CRS	4.007095
olm1000	1.89504	1.48348	CRS	0.396840
raefsky2	70.6624	69.7615	CRS	10.778230
roadNet-PA	8.07317	5.84257	CRS	2.966768
thermal1	15.5753	13.318	CRS	3.386097
thermal2	19.2541	17.0523	CRS	6.004222
thermomech_TK	15.3241	14.2953	CRS	4.305889
west2021	2.66929	2.26954	CRS	0.637915
Cube_Coup_dt0	111.963	113.348	ELLPACK	17.953161
FEM_3D_thermal1	32.0555	33.4671	ELLPACK	5.634524
ML_Laplace	176.774	188.886	ELLPACK	28.493590
PR02R	135.828	146.651	ELLPACK	22.294501
adder_dcop_32	1.20824	4.12316	ELLPACK	0.920089
af23560	29.2923	29.4083	ELLPACK	5.066518
af_1_k101	79.1365	80.4176	ELLPACK	13.098161
bcsstk17	42.2982	45.3715	ELLPACK	7.251562
cant	99.3766	110.127	ELLPACK	17.154513
cavity10	21.571	23.0165	ELLPACK	3.815625
cop20k_A	46.062	47.6264	ELLPACK	9.980867
dc1	0.733577	14.7532	ELLPACK	3.147779
mcfe	10.5186	11.2408	ELLPACK	1.852130
mhda416	4.46372	4.6875	ELLPACK	0.848498
nlpkkt80	67.4344	67.5926	ELLPACK	11.394303
olafu	66.3819	69.052	ELLPACK	10.654512
rdist2	15.0779	15.7503	ELLPACK	2.733008
webbase-1M	8.04237	8.68042	ELLPACK	2.314507

The results of parallelizing the CRS and ELLPACK algorithms with CUDA show a strong preference for the CRS format on certain matrices (such as *amazon0302*, *lung2*, and *thermal1*), attributable to better sparsity management and memory access patterns optimized for GPUs. In contrast, ELLPACK shows superior performance for matrices with regular structures, such as *Cube_Coup_dt0* and *ML_Laplace*, due to more uniform memory access.

The variability of the speed-up factor between matrices highlights the importance of the specific structure of the matrix in the relative efficiency of CRS and ELLPACK under CUDA. In particular, symmetric and dense arrays tend to favour ELLPACK, highlighting the significant role of density and symmetry in the performance of the formats.

4.1.4 Overall Performance Comparison

Table 4.4: Overall Performance Comparison

Matrix	Best Performance	Best Structure	Best Method	Speedup
amazon0302	12.3120	CRS	CUDA	5.252515
lung2	10.2714	CRS	CUDA	2.216157
mac_econ_fwd500	15.6609	CRS	CUDA	4.250389
mhd4800a	21.9362	CRS	CUDA	4.007095
raefsky2	70.6624	CRS	CUDA	10.778230
roadNet-PA	8.07317	CRS	CUDA	2.966768
thermal1	15.5753	CRS	CUDA	3.386097
thermal2	19.2541	CRS	CUDA	6.004222
thermomech_TK	15.3241	CRS	CUDA	4.305889
Cube_Coup_dt0	113.348	ELLPACK	CUDA	17.953161
FEM_3D_thermal1	33.4671	ELLPACK	CUDA	5.634524
ML_Laplace	188.886	ELLPACK	CUDA	28.493590
PR02R	146.651	ELLPACK	CUDA	22.294501
af23560	29.4083	ELLPACK	CUDA	5.066518
af_1_k101	80.4176	ELLPACK	CUDA	13.098161
bcsstk17	45.3715	ELLPACK	CUDA	7.251562
cant	110.127	ELLPACK	CUDA	17.154513
cavity10	23.0165	ELLPACK	CUDA	3.815625
cop20k_A	47.6264	ELLPACK	CUDA	9.980867
dc1	14.7532	ELLPACK	CUDA	3.147779
nlpkkt80	67.5926	ELLPACK	CUDA	11.394303
olafu	69.0520	ELLPACK	CUDA	10.654512
rdist2	15.7503	ELLPACK	CUDA	2.733008
webbase-1M	8.68042	ELLPACK	CUDA	2.314507
adder_dcop_32	5.12204	CRS	OpenMP	1.142991
mcfe	12.6097	CRS	OpenMP	2.077682
mhda416	7.06099	CRS	OpenMP	1.278130
cage4	0.871628	CRS	Serial	1.000000
olm1000	4.77532	CRS	Serial	1.000000
west2021	4.1844	CRS	Serial	1.000000

Exceptionally high performance is observed for certain matrices under CUDA, such as *Cube_Coup_dt0* and *ML_Laplace*, revealing the match between certain data structures and GPU optimisation. On the other hand, OpenMP shows a moderate advantage for specific matrices, such as *adder_dcop_32* and *mcfe*, offering a significant performance improvement without the need for specialised hardware, thanks to parallelization on shared memory architectures.

However, matrices of small size or with characteristics less favourable to parallelization, such as *cage4*, *olm1000*, and *west2021*, show no significant improvement with

OpenMP or CUDA, indicating that for some cases sequential execution remains the most suitable method.

These observations lead to the conclusion that the choice between CRS and ELLPACK formats, as well as the decision to use sequential execution, OpenMP or CUDA, should be informed by the specific properties of the matrices in question. Optimisation of the calculation parameters and careful evaluation of the data characteristics are essential to maximise the efficiency of operations on hollow matrices.

Chapter 5

Conclusion

In summary, this report explored the performance of the CSR and ELLPACK formats for fat vector multiplication of hollow matrices using the OpenMP and CUDA parallel programming paradigms. The results show a marked preference for the CSR format when parallelized with OpenMP, which is probably due to more efficient memory management and optimal task distribution between threads. On the other hand, CUDA performance is strongly influenced by the sparsity and structural regularity of matrices, with symmetric and dense matrices seemingly favouring the ELLPACK format.

The benefits of using CUDA on GPU architectures have been demonstrated, particularly for arrays that align well with memory access models optimised for these devices. However, it is clear that the benefits of parallelization with CUDA or OpenMP are closely related to the specific characteristics of the arrays in question, and that a sequential approach may be preferable for small arrays or those with patterns less conducive to parallelization.

These findings underline the importance of a thorough analysis of data structures and computational parameters to maximise the efficiency of operations on hollow matrices. Ultimately, the choice between CSR and ELLPACK formats, as well as the decision to use sequential execution, OpenMP or CUDA, must be informed by the specific properties of the matrices. Such a nuanced understanding will further optimise the performance of the scientific and engineering applications that depend on these intensive computations.

References

Appendix A

Documentation

Appendix A.A Project tree

```
Source Code/  
  CRS/  
    matrixMultivectorProductCRS .cpp  
    matrixMultivectorProductCRS .h  
    matrixMultivectorProductCRSCUDA .cu  
    matrixMultivectorProductCRSCUDA .h  
    matrixMultivectorProductCRSOpenMP .cpp  
    matrixMultivectorProductCRSOpenMP .h  
  ELLPACK/  
    matrixMultivectorProductELLPACK .cpp  
    matrixMultivectorProductELLPACK .h  
    matrixMultivectorProductELLPACKCUDA .cu  
    matrixMultivectorProductELLPACKCUDA .h  
    matrixMultivectorProductELLPACKOpenMP .cpp  
    matrixMultivectorProductELLPACKOpenMP .h  
  scripts /  
    cuda .sub  
    openMP .sub  
    parseCudaResults .sh  
    parseOpenMPResults .sh  
  cudaUtils .cu  
  makefile  
  MatrixDefinitions .h  
  runCuda .cpp  
  runOpenMP .cpp  
  utils .h  
  utils .cpp  
results /  
  images /  
  CUDA.csv  
  CUDA.ipynb  
  OpenMP.csv  
  OpenMP.ipynb
```

Appendix A.B Getting Started

To run the program, follow these steps:

1. Install the required compilers and libraries:
 - **OpenMP:** Install the GNU Compiler Collection (GCC) and OpenMP (?).
 - **CUDA:** Install the NVIDIA CUDA Toolkit (?).
2. Compile the files using the following command: `make all`.
3. Run the programs:
 - **OpenMP:** `./runOpenMP.o`
 - **CUDA:** `./runCuda.o`

Appendix A.C Methods Overview

A.C.1 Utils.h

A.C.1.1 convertCRStoELLPACK

Description: Read a sparse matrix from a Matrix Market file.

Parameters:

- `SparseMatrixCRS &crsMatrix`: The CRS matrix to convert.
- `SparseMatrixELLPACK &ellpackMatrix`: The ELLPACK matrix to convert to.

A.C.1.2 areMatricesEqual

Description: Compares two matrices for equality within a specified tolerance.

Parameters:

- `FatVector &mat1`: First matrix.
- `FatVector &mat2`: Second matrix.
- `double tolerance`: Tolerance for comparison.

Returns: `bool`: True if matrices are equal within the tolerance, false otherwise.

A.C.1.3 readMatrixMarketFile

Description: Reads a matrix from a Matrix Market file into a sparse matrix format.

Parameters:

- `std::string &filename`: Name of the Matrix Market file.
- `SparseMatrixCRS &matrix`: Sparse matrix to read into.

A.C.1.4 generateLargeFatVector

Description: Generates a random Fat Vector with specified dimensions.

Parameters:

- `FatVector &fatVector`: Fat vector to generate.
- `int n`: Number of rows.
- `int k`: Number of columns.

A.C.2 matrixMultivectorProductCRS.h**A.C.2.1 matrixMultivectorProductCRS**

Description: Perform the matrix-vector multiplication in the CRS format.

Parameters:

- `SparseMatrixCRS &sparseMatrix`: Sparse matrix in CRS format.
- `FatVector &fatVector`: Fat vector.
- `FatVector &result`: Result of the multiplication.
- `int testNumber`: Number of iterations for the performance measurement

A.C.3 matrixMultivectorProductCRSOpenMP.h**A.C.3.1 matrixMultivectorProductCRSOpenMP**

Description: Perform the matrix-vector multiplication in the CRS format using OpenMP

Parameters:

- `SparseMatrixCRS &sparseMatrix`: Sparse matrix in CRS format.
- `FatVector &fatVector`: Fat vector.
- `FatVector &result`: Result of the multiplication.
- `int testNumber`: Number of iterations for the performance measurement
- `int numThreads`: Number of threads to use.
- `int chunkSize`: Chunk size for the parallelization.

A.C.4 `matrixMultivectorProductCRSCUDA.h`

A.C.4.1 `matrixMultivectorProductCRSCUDA`

Description: Perform the matrix-vector multiplication in the CRS format using CUDA

Parameters:

- `SparseMatrixCRS &sparseMatrix`: Sparse matrix in CRS format.
- `FatVector &fatVector`: Fat vector.
- `FatVector &result`: Result of the multiplication.
- `int testNumber`: Number of iterations for the performance measurement
- `int xBlockSize`: X block size for the parallelization.
- `int yBlockSize`: Y block size for the parallelization.

A.C.5 `matrixMultivectorProductELLPACK.h`

A.C.5.1 `matrixMultivectorProductELLPACK`

Description: Perform the matrix-vector multiplication in the ELLPACK format.

Parameters:

- `SparseMatrixELLPACK &sparseMatrix`: Sparse matrix in ELLPACK format.
- `FatVector &fatVector`: Fat vector.
- `FatVector &result`: Result of the multiplication.
- `int testNumber`: Number of iterations for the performance measurement

A.C.6 `matrixMultivectorProductELLPACKOpenMP.h`

A.C.6.1 `matrixMultivectorProductELLPACKOpenMP`

Description: Perform the matrix-vector multiplication in the ELLPACK format using OpenMP

Parameters:

- `SparseMatrixELLPACK &sparseMatrix`: Sparse matrix in ELLPACK format.
- `FatVector &fatVector`: Fat vector.
- `FatVector &result`: Result of the multiplication.
- `int testNumber`: Number of iterations for the performance measurement
- `int numThreads`: Number of threads to use.
- `int chunkSize`: Chunk size for the parallelization.

A.C.7 `matrixMultivectorProductELLPACKCUDA.h`

A.C.7.1 `matrixMultivectorProductELLPACKCUDA`

Description: Perform the matrix-vector multiplication in the ELLPACK format using CUDA

Parameters:

- `SparseMatrixELLPACK &sparseMatrix`: Sparse matrix in ELLPACK format.
- `FatVector &fatVector`: Fat vector.
- `FatVector &result`: Result of the multiplication.
- `int testNumber`: Number of iterations for the performance measurement
- `int xBlockSize`: X block size for the parallelization.
- `int yBlockSize`: Y block size for the parallelization.

Appendix B

Source Codes