

Parallelization of 0th Order Generalised Mode Acceleration Method

Implementation of OpenMP, MPI and CUDA Parallelizing Algorithms

Abhinandan Kumbhar- 193109007 R. Nitin Iyer - 19310R001 Jalaj Gupta - 19310R002 Saurabh Pal - 193104001

Abstract—For large systems, it is a common practice to implement model order reduction to minimize computational effort. In this project a displacement response of the linear elastic cantilever beam subjected to harmonic point load at the point of application is estimated. Various parallelization technique has been used to expedite the calculation process. Model order reduction has been done using Guyan condensation on the 3782×3782 system to get reduced system of 500×500 . A 0th order Generalized mode acceleration method proposed by J. Rixen [1] was applied to the reduced system. First 10 modes were used to approximate the response, and GMAM was done using OpenMP, MPI and GPGPU programming, and the results were compared with the serial code and MATLAB codes to test for consistency and time reduction.

Index Terms—GMAM, OpenMP, MPI, GPGPU, CUDA etc.

I. INTRODUCTION

DESCRIBE: As shown in adjacent figure. A cantilever beam is rigidly held in a rigid wall at one end and at free end at one edge a harmonic force of frequency comparable to first natural frequency of the beam is being applied longitudinally and transversally. This problem has been modelled as 2D problem. So in figure above one plane section is being depicted which has been considered for further analysis. The discretization of the plane in FEM resulted in 3782×3782 mass and stiffness matrices. As per modal analysis it has been found that modal participation factors for first ten modes are significant enough for this analysis. So instead of carrying out our computation on such a large mass and stiffness matrices, these matrices has been condensed to smaller 500×500 order using Guyan condensation. This model order reduction has been carried out accurately on MATLAB. Then these reduced matrices has been used for further computation, on C language, for estimating response of the system. Zeroth order generalized mode acceleration method has been used in the code for response calculation. This serial code written for zeroth order mode acceleration method is further parallelized using open MP, MPI and CUDA programming language to reduce execution time. In this report comparative timing study has been done to check code performance and results are also compared with MATLAB results to check for data consistency and coherency.

II. THEORY FOR GENERALIZED EIGENVALUE PROBLEM

Consider following generalized eigenvalue

$$\mathbf{K}\Phi = \lambda\mathbf{M}\Phi$$

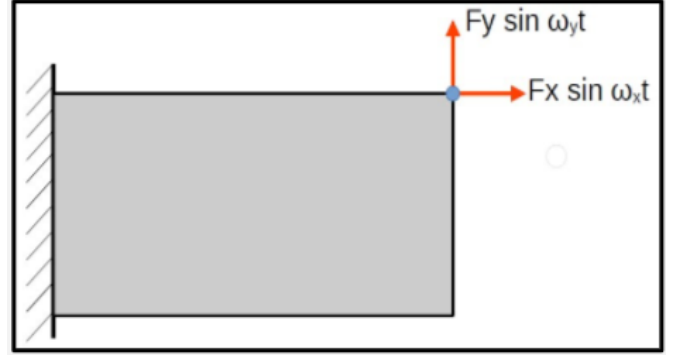


Fig. 1: Cantilever beam subjected to sinusoidal loading at the specified node

To solve above eigenvalue problem we will use the method proposed by [2], consider first following standard eigenvalue problem for \mathbf{M} :

$$\mathbf{M}\Phi_{\mathbf{M}} = \Phi_{\mathbf{M}}\lambda_{\mathbf{M}}$$

Where $\Phi_{\mathbf{M}}$, and $\lambda_{\mathbf{M}}$ are the eigenvector and eigenvalue matrices of \mathbf{M} , respectively. Then, we have:

$$\begin{aligned} \mathbf{M}\Phi_{\mathbf{M}} = \Phi_{\mathbf{M}}\lambda_{\mathbf{M}} &\implies \Phi_{\mathbf{M}}^{-1}\mathbf{M}\Phi_{\mathbf{M}} = \Phi_{\mathbf{M}}^{-1}\Phi_{\mathbf{M}}\lambda_{\mathbf{M}} \\ &\implies \lambda_{\mathbf{M}} = \Phi_{\mathbf{M}}^T\mathbf{M}\Phi_{\mathbf{M}} \end{aligned}$$

Eigen-vector matrix has to be orthogonal because \mathbf{M} is a symmetric matrix. If above equation is pre and post multiplied by $\lambda_{\mathbf{M}}^{-0.5}$. Then,

$$\lambda_{\mathbf{M}}^{-0.5}\Phi_{\mathbf{M}}^T\mathbf{M}\Phi_{\mathbf{M}}\lambda_{\mathbf{M}}^{-0.5} = \lambda_{\mathbf{M}}^{-0.5}\lambda_{\mathbf{M}}\lambda_{\mathbf{M}}^{-0.5} = \mathbf{I}$$

$$\hat{\Phi}_{\mathbf{M}} = \Phi_{\mathbf{M}}\lambda_{\mathbf{M}}^{-0.5}$$

Let $\hat{\mathbf{K}}$ be defined as

$$\hat{\mathbf{K}} = \hat{\Phi}_{\mathbf{M}}^T\mathbf{K}\hat{\Phi}_{\mathbf{M}}$$

Since \mathbf{K} is symmetric this matrix would also be symmetric. The eigenvalue problem for $\hat{\mathbf{K}}$ is

$$\hat{\mathbf{K}}\Phi_{\mathbf{K}} = \lambda_{\mathbf{K}}\Phi_{\mathbf{K}}$$

where, $\Phi_{\mathbf{K}}$ and $\lambda_{\mathbf{K}}$ are the eigenvector and eigenvalue matrices of $\hat{\mathbf{K}}$. Pre-multiply above matrix by $\Phi_{\mathbf{K}}^T$

$$\Phi_{\mathbf{K}}^T\hat{\mathbf{K}}\Phi_{\mathbf{K}} = \lambda_{\mathbf{K}}$$

Now plugging back expression for $\hat{\mathbf{K}}$ in above expression

$$\begin{aligned}\Phi_K^T \hat{\Phi}_M^T K \hat{\Phi}_M^T \Phi_K &= \lambda_K \\ \Phi_K^T (\Phi_M \lambda_M^{-0.5})^T K \Phi_M \lambda_M^{-0.5} \Phi_K &= \lambda_K \\ \implies \Phi^T K \Phi &= \lambda_K\end{aligned}$$

where,

$$\Phi = \hat{\Phi}_M \Phi_K = \Phi_M \lambda_M^{-0.5} \Phi_K$$

This expression also implies that Φ also diagonalizes \mathbf{B} and gives identity matrix. So,

$$\Phi^T M \Phi = \mathbf{I} \implies \Phi^T M \Phi \lambda_K = \lambda_K$$

$$\Phi^T M \Phi \lambda_K = \Phi^T K \Phi$$

$$M \Phi \lambda_K = K \Phi$$

This is expression for generalized eigenvalue problems. So, the following algorithm has been devised for calculation of the eigenvalues and eigenvectors-

- 1) $\Phi_M, \lambda_M \leftarrow M \Phi_M = \lambda_M \Phi_M$
- 2) $\hat{\Phi}_M \leftarrow \hat{\Phi}_M = \Phi_M \lambda_M^{-0.5} \approx \Phi_M (\lambda_M^{-0.5} + \epsilon \mathbf{I})^{-1}$
- 3) $\hat{\mathbf{K}} \leftarrow \hat{\mathbf{K}} = \hat{\Phi}_M^T K \hat{\Phi}_M$
- 4) $\Phi_K \lambda_K \leftarrow \hat{\mathbf{K}} \Phi_K = \lambda_K \Phi_K$
- 5) $\lambda \leftarrow \lambda = \lambda_K$
- 6) $\Phi \leftarrow \Phi = \hat{\Phi}_M \Phi_K$
- 7) Φ and λ

While solving generalized eigenvalue problem (\mathbf{K}, \mathbf{M}) , we have to solve two simple eigenvalue problems of \mathbf{M} and $\hat{\mathbf{K}}$. We have used QR decomposition for solving this problem. We will now study QR decomposition in details, followed by GMAM method

III. THEORY OF QR DECOMPOSITION WITH GRAM-SCHMIDT

We now perform QR decomposition of a matrix \mathbf{A} as \mathbf{QR} , where \mathbf{Q} is orthogonal matrix and \mathbf{R} is upper triangular matrix. The computation of \mathbf{Q} is iterative process in which each iteration deals with calculation of k^{th} column of \mathbf{Q} by using below formula

$$\mathbf{Q}_k = \mathbf{A}_k - \sum_{n=0}^{k-1} (\mathbf{A}_k \cdot \mathbf{Q}_n) \mathbf{Q}_n$$

The \mathbf{Q}_k is then used to rotate \mathbf{A}_k and proceed to next iteration.

IV. THEORY OF GMAM FOR 0^{th} ORDER

After obtaining the eigenvalues and eigenvectors from the algorithm explained in the above section, the response can be obtained as-

$$q_n = \sum_{r=1}^k \frac{x_r x_r^T F}{\omega_r^2 - \omega^2} \left[\sin \omega t - \frac{\omega}{\omega_r} \sin \omega_r t \right]$$

where, x_r is the r th eigenvector, k is the number of modes considered & being 10 in our case. ω is forcing frequency, ω_r is square root of r^{th} eigenvalue. The F is vector of force amplitudes.

V. POTENTIAL REGIONS FOR PARALLELIZATION

After studying the algorithm of solving generalized eigenvalue problem and GMAM, we found out that major portion of code deals with matrix multiplication. So, we decided to perform parallelization of matrix multiplication.

While solving generalized eigenvalue problem (\mathbf{K}, \mathbf{M}) , we have to solve two simple eigenvalue problems of \mathbf{M} and $\hat{\mathbf{K}}$. We have used QR decomposition for solving this problem, Which is iterative process. Suppose we perform L iterations of QR decomposition for solving single simple eigenvalue problem. The QR decomposition in itself needs N i.e size of matrix, iterations. Then it means we need to do $2 * L * N$ iterations in total for solving two simple eigenvalue problems. We can parallelize this QR decomposition algorithm and save a lot of time.

Depending upon above study, we have decided to parallelize QR decomposition and matrix multiplication using OpenMP, OpenMPI and Cuda.

VI. IMPLEMENTATION OF OPENMP, MPI AND GPGPU ALGORITHMS

OpenMP :

The OpenMp parallelization is done by observing the 'for' loops and making the appropriate variables as private to each thread.

MPI :

Parallelization of QR decomposition

$$\mathbf{Q}_k = \mathbf{A}_k - \sum_{n=0}^{k-1} (\mathbf{A}_k \cdot \mathbf{Q}_n) \mathbf{Q}_n$$

Formation of \mathbf{Q} is parallelized by distribution of

$$\sum_{n=0}^{k-1} (\mathbf{A}_k \cdot \mathbf{Q}_n) \mathbf{Q}_n$$

among the processors. Each process has to perform computation for k/cores columns. Root process will normalize and assemble the columns of \mathbf{Q} . The \mathbf{Q} will then be used to rotate the mass of transformed matrix, and then next iteration will be performed.

Using matrix multiplication, we'll transform the \mathbf{K} matrix, and perform QR decomposition on transformed \mathbf{K} matrix.

Parallelization of Matrix Multiplication ($C = AB$) :- Share \mathbf{B} matrix with all the processors. Send sets of rows to different processors, and perform multiplication of those rows with the \mathbf{B} matrix. Receive the computed rows at root processor, and assemble all the rows at appropriate positions.

$$\text{Setsofrows} = \frac{\text{No.ofrowsofA}}{\text{No.ofProcessors}}$$

CUDA parallelization

As explained in the previous section, Two major areas are identified where we can apply parallelization. Kernel functions

are written for below two codes and used them wherever they are required.

- 1) Matrix multiplication
- 2) QR decomposition

CUDA kernel for matrix multiplication ($AB=C$)

Suppose we have matrix C of size $N \times N$, We have launched a kernel having grid of $N \times N$ size, each grid block having N threads.

`matmult <<< grid(N,N), N >>> (A,B,C)`

Each block corresponds to one element of C matrix and calculates one entry of C matrix $C[i][j]$. Within each block, N threads will help compute dot product i^{th} row of A and j^{th} column of B.

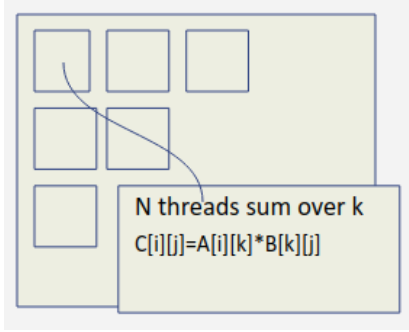


Fig. 2: Cuda Kernel for matrix multiplication

CUDA kernel for QR decomposition ($A=QR$)

QR decomposition is performed in n iterations. k^{th} iteration corresponds to calculation of k^{th} column of Q matrix. So, below explained kernels are launched for n times for computing each column of Q. Formula for calculating k^{th} column of Q is as follows:

$$\text{Temp}_k = A_k - \sum_{n=0}^{k-1} (A_k \cdot Q_n) Q_n$$

$$Q_k = \text{Temp}_k / \text{Sqrt}(\text{Temp}_k \cdot \text{Temp}_k)$$

As we can see, k dot products have to be performed for k^{th} iteration. Dot product of A_k and Q_n are then multiplied with Q_n i.e. $(A_k \cdot Q_n) Q_n$ and sum them all. Then, subtracted the summation from A_k to get vector Temp_k . Unit vector of Temp_k is stored in k^{th} column of Q. All this is done in three steps.

- 1) Step 1: Calculation of k dot products i.e. $A_k Q_n$.

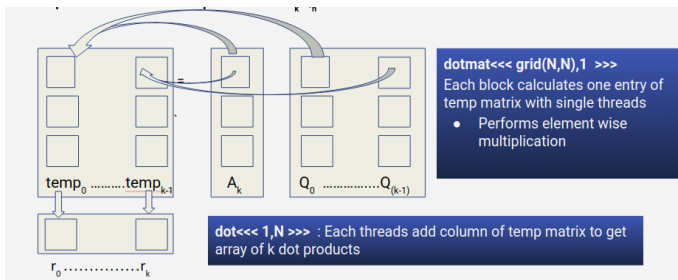


Fig. 3: Cuda Kernels for step 1

- 2) Step 2: Calculate matrix $(A_k \cdot Q_n) Q_n$.
- 3) Step 3: Calculate Temp_k and Q_k .

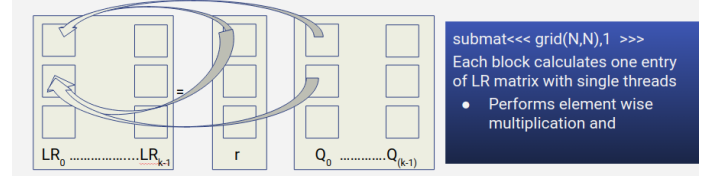


Fig. 4: Cuda Kernel for step 2

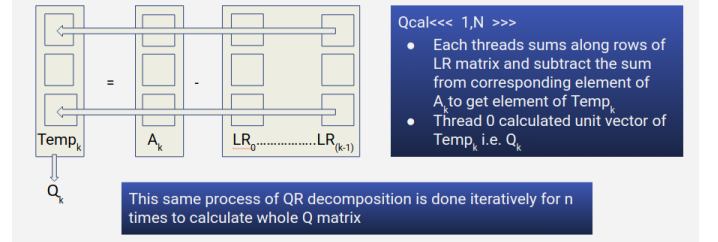


Fig. 5: Cuda Kernel for step 3

VII. RESULTS AND DISCUSSION

System specifications : The code were run on the CDAC Param Sanganak. For OpenMp, number of threads used were equal to number of cores allocated. Similarly, for OpenMPI, number of cores assigned were same as number of processes.

Time study : The serial code took 48 mins to run and below is the time study for different techniques.

OpenMP (min)	OpenMPI (min)	Cuda (min)
8 thds - 9.1	8 PEs - 9.4	8 cores - 11.02
12 thds - 10.4	12 PEs - 9.3	12 cores - 10.9
16 thds - 12.5	16 PEs - 10.7	16 cores - 9.7
20 thds - 15.03	20 PEs - 11.1	20 cores - 10.1

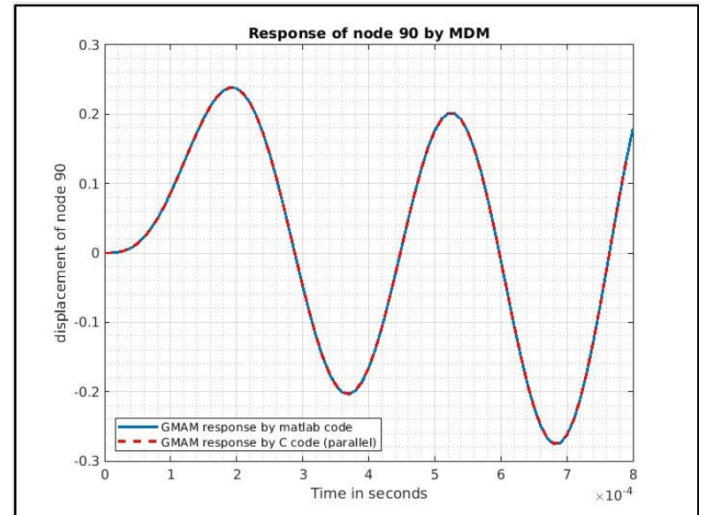


Fig. 6: Response obtained from code plotted in Matlab

The response of node 3781 (90 in reduced matrix) is plotted and its error value is also plotted. we can see that the error is of the order 10^{-4} (0.1%)

The same problem was solved using Matlab code to confirm that the parallelized code are having data consistency and coherency.

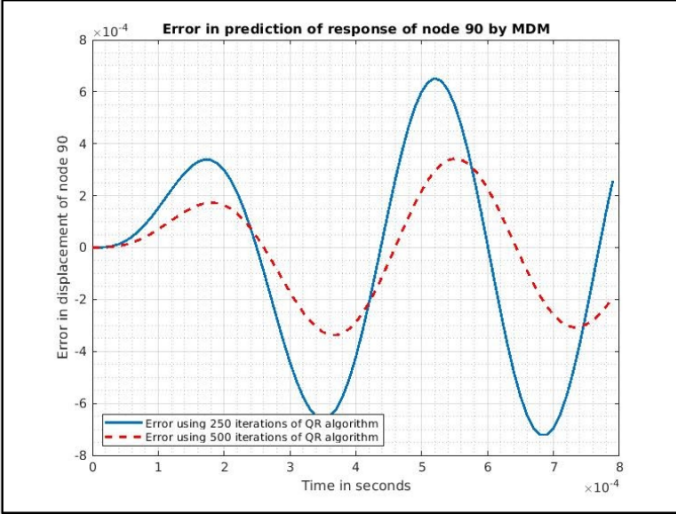


Fig. 7: Error variation with changing number of iterations used for predicting response

In case of MPI programming for parallelization, the increase in time as number of processes goes on increasing is counter intuitive, but it may have happened because decrease in time of computation by each process is not significant as compared to time increase due to increase in data sharing. Since after distributing the workload equally amongst all processes the extra remainder work has to be done by the root process, the time taken by root process is bottle-neck.

We can observe that for CUDA program, time remains almost same even if we increase the number of cores. It may have happened because we are launching the same kernels in every case. We are launching same threads and same numbers of blocks for all various numbers of cores.

As explained by Amdahl's law, We do not get speed-up as expected. Since the every iteration of QR decomposition is dependent on previous iteration, we have limitation on speed up, still we could see speed-up upto 4-5 times.

REFERENCES

- [1] D. Rixen, "Generalized mode acceleration methods and modal truncation augmentation," in *19th AIAA Applied Aerodynamics Conference*, p. 1300, 2001.
- [2] B. Ghogh, F. Karray, and M. Crowley, "Eigenvalue and generalized eigenvalue problems: Tutorial," *arXiv preprint arXiv:1903.11240*, 2019.