

# A Unique Hastening Communication-Bound Sparse Repeatabe Solver

<sup>1</sup>P.AKSHAY TEJA, <sup>2</sup>B.BALAJI, <sup>3</sup>B.BALA KRISHNA, <sup>4</sup>Dr.B.S.R.MURTHY

<sup>1</sup> M.Tech Student, DEPT OF ELECTRONICS & COMMUNICATION ENGINEERING. GANDHI ACADEMY OF TECHNICAL EDUCATION, Ramapuram (Katamommu Gudem), Chilkur(M), Kodad, Telangana 508206

<sup>2</sup> M.TECH, Assistant Professor, DEPT OF ELECTRONICS & COMMUNICATION ENGINEERING. GANDHI ACADEMY OF TECHNICAL EDUCATION, Ramapuram (Katamommu Gudem), Chilkur(M), Kodad, Telangana 508206

<sup>3</sup> M.TECH, Assistant Professor, HOD, DEPT OF ELECTRONICS & COMMUNICATION ENGINEERING. GANDHI ACADEMY OF TECHNICAL EDUCATION, Ramapuram (Katamommu Gudem), Chilkur(M), Kodad, Telangana 508206

<sup>4</sup> Phd, Professor, & Principal. GANDHI ACADEMY OF TECHNICAL EDUCATION, Ramapuram (Katamommu Gudem), Chilkur(M), Kodad, Telangana 508206

## ABSTRACT:

This paper presents an organized procedure to pick this algorithmic parameter  $k$ , which supplies communication-computation compromise on hardware accelerators like FPGA and GPU. While  $k$  iterations from the iterative solver could be unrolled to provided decrease in communication cost, the extent of the unrolling is dependent around the underlying architecture, its memory model, and also the development in redundant computation. Buying and selling communication with redundant computation can boost the plastic efficiency of FPGAs and GPUs in speeding up communication-bound sparse iterative solvers. With an NVidia C2050 GPU, we demonstrate a speedup over standard iterative solvers for a variety of benchmarks which this speedup is restricted through the development in redundant computation. In comparison, for FPGAs, we produce an architecture-aware formula that limits off-nick communication but enables communication between your processing cores. This reduces redundant computation and enables large  $k$  and therefore greater speedups. Our method for FPGA supplies a speedup over same-generation GPU products where  $k$  is selected carefully for architectures for a variety of benchmarks. We offer predictive models to know this compromise and show how careful choice of  $k$  can result in performance improvement that otherwise demands significant rise in memory bandwidth.

**Keywords:** *Iterative numerical methods, sparse matrix-vector multiply, matrix powers kernel, field programmable gate arrays (FPGAs), graphics processing units (GPUs).*

## I. INTRODUCTION

Among the communication-intensive scientific computations is definitely an iterative solver employed for fixing large-scale sparse straight line system of equations and eigenvalue problems. The price of a higher performance scientific computation operating on large datasets includes two factors i) computation price of carrying out floating-point procedures ii) communication cost (both latency and bandwidth) of moving data inside the memory hierarchy in consecutive situation or between processors in parallel situation [1]. The answer of those problems is calculated from the Kriol subspace span, in which a new vector is produced in every iteration. Iterative solvers are difficult to accelerate because they spend more often than not in communication-bound computations, like sparse matrix-vector multiply (SpMV) and vector-vector procedures (us dot items and vector additions). For big scale problems in which the matrix A does not fit on-nick, regardless of how much parallelism could be used to accelerate SpMV, the performance from the iterative solver is bounded through the available off-nick memory bandwidth, e.g.

with 2 flops per 4 bytes (single-precision) in SpMV, the utmost theoretical performance is 71 GFLOPs with an NvidiaC2050GPU and 17GFLOPs on the Virtex6 FPGA. Furthermore, the information dependency between these procedures hinders overlapping communication with computation. This leads to fewer than 7 % and 4 % efficiency of GPU and FPGA correspondingly. The communication issue is attached to the memory wall problem. Because of technology scaling, computation performance is growing in a dramatic rate whereas communication performance is enhancing but in a reduced rate. It's a well-known idea to formulate algorithmic improvements that hide memory latency and optimize memory bandwidth. For iterative solvers, Demmel et al. trade communication with redundant computation by changing k SpMVs using the matrix forces kernel. The important thing idea would be to partition the matrix into blocks and performs k SpMVs on blocks without fetching the block again within the consecutive situation and carrying out redundant computation to prevent communication along with other

processors within the parallel situation. In this manner the communication price is reduced at the fee for rise in redundant computation [2]. They reveal that this kind of approach can minimize latency inside a grid and both latency and bandwidth on the multi-core CPUs to provide speedup correspondingly over  $k$  SpMV's for banded matrices. As the communication staying away from approach is promising, there's two primary challenges connected with this particular communication-staying away from approach on parallel architectures i) how you can keep your redundant computation to a minimum to reduce the computation cost and ii) how to decide on the optimal worth of the algorithmic parameter  $k$ , which minimizes overall runtime by supplying a compromise between computation and communication cost. Within this paper, we show the way we can boost the plastic efficiency of FPGAs and GPUs in speeding up communication-bound sparse iterative solver. Like a motivation, we show a compromise between computation and communication cost for FPGA and GPU to reduce overall runtime. The primary contributions of the paper are: i) Communication optimization inside the memory hierarchy of merely one stream

multiprocessor (SM) in addition to between different SMs while mapping the matrix forces kernel to some GPU. ii) An architecture-aware matrix forces kernel that suits the effectiveness of the FPGAs to prevent redundant computation along with a resource-restricted methodology to choose  $k$  for the FPGA. iii) A unified predictive type of the matrix forces kernel for GPU and FPGA, which will help us understanding communication computation tradeoffs in choosing the algorithmic parameter  $k$ . While using steepest ascent approach, we show which facet of future GPU and FPGA architectures have to be enhanced to attain greater performance. iv) For a variety of problem dimensions, a quantitative comparison from the matrix forces kernel on FPGA.

## II. RELATED WORK

Within this work, we target large structured sparse banded matrices. The banded matrices are kept in band storage format where matrix of band size  $b$  is stored like a matrix. We decide this band structure for 2 primary reasons. First, computations on such matrices happen to be utilized as an architectural evaluation benchmark because of high parallelism and occasional

computational intensity, offering possibilities to take advantage of on-nick parallelism and challenges with connected memory systems. Next, they naturally arise in several scientific computations like stencils in partial differential equation (PDE) solvers, semi-definite optimization problems and model predictive control.

### III. EXISTED SYSTEM

We first present the present GF100 GPU architecture after which discuss different optimization techniques that cause high throughput [3]. Then we produce an analytical model to calculate and comprehend the performance from the matrix forces kernel on any GPU device (model parameters are acquired using micro-benchmarks). We make use of the same model to pick  $k$  for current GPU architectures as well as make architectural projections to acquire a preferred performance with future products. The GPU comprises 14 streaming multiprocessors (SMs) each operating at 1.15 GHz. Each SM has 32 floating-point cores able to carrying out 1 single-precision flop/cycle reaching an optimum throughput. We feature out sensitivity analysis of GPU performance with regards to the algorithmic parameters

with constant architectural parameters. We highlight that by carefully picking the algorithmic parameters we are able to achieve greater performance over  $k$  SpMV's that otherwise requires significant architectural modifications when it comes to global memory bandwidth and latency.

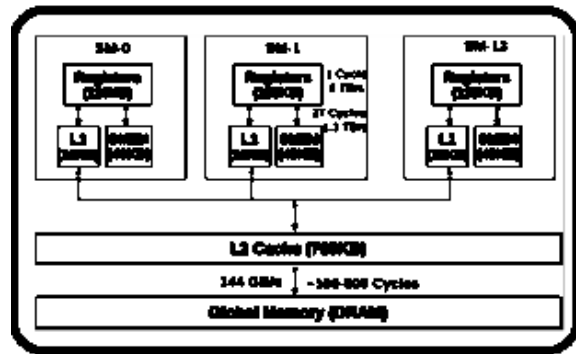


Fig.1. Nvidia GPU Structure

### IV. PROPOSED MODEL

The possibility to make use of FPGAs in high-performance computing comes from the truth that computer architecture could be specialized to accelerate a specific task solvers. We currently introduce the architecture-aware hybrid matrix forces kernel that exploits these architectural features to obtain high throughput. In connection with this, we present an origin-restricted methodology for choosing an ideal  $k$  for any target FPGA device. The suggested formula loads the blocks from the matrix in the slow memory into large on-nick memory utilizing a consecutive formula

after which performs computations inside the block in parallel without having done redundant computations. The significant from the formula, having a toy example, where we've 2 outer blocks which are loaded sequentially in the slow memory into FPGA on-chip memory [4]. Each outer block is further partitioned into two sub blocks that may be processed in parallel by a range of processing elements (PEs) your SIMD fashion All of the vertices in the sub-block are calculated inside a pipelined fashion utilizing a reduction circuit inside the PE. After each level within the graph, PEs has to communicate dependencies for their nearest neighbors. However, unlike GPU, where this communication is just possible using shared global memory and it is therefore prevented using redundant computation, the PEs inside the FPGA utilize low-latency FIFOs and therefore steer clear of the redundant computation. This enables bigger values of  $k$  and therefore greater possible speedups. Speedup over  $k$  SpMV's utilized in standard iterative solvers. This speedup factor is nearly two times of the items we accomplished with parallel matrix forces kernel on GPU. Even though the internet performance of FPGA (85 GFLOPs) is under those of GPU (123 GFLOPs) with this

band size, however, we've better plastic efficiency with FPGA (18.8 percent) in comparison to GPU. The model is exact because of the highly predictive nature of FPGAs like a computing platform [5]. As FPGA has relatively bigger on-chip memory in comparison to GPU, we plan to keep  $k$  vectors on-chip to get used by subsequent modules in communication staying away from iterative solver. You will find three procedures in the matrix forces kernel on FPGA: loading the block, computing the sub-blocks in parallel as well as an optional stage for storing the  $k$  vectors to the off-chip memory if they don't fit on-chip.

## V. CONCLUSION

Although unrolling  $k$  iterations while using matrix forces kernel provides significant performance improvement in comparison to plain  $k$  SpMV's on the GPU, the performance is restricted because of quadratic development in redundant computations. This enables us large worth of  $k$  and therefore superior plastic efficiency in comparison to GPU. Buying and selling communication with computation boosts the plastic efficiency of hardware accelerators like FPGAs and GPU for speeding up communication-bound sparse iterative

solver. Our suggested hybrid matrix forces kernel for FPGA exploits the architectural options that come with this significantly different platform to reduce redundant computations. Our architectural insight shows a good formula-architecture interaction can offer similar performance, which otherwise requires significant enhancements in memory bandwidth. For a variety of at random produced banded matrices, we shown speedup over GPU for big and small band dimensions correspondingly.

## REFERENCES

- [1] E. Klerk, “Exploiting Special Structure in Semi definite Programming: A Survey of Theory and Applications,” *Eur. J. Oper. Res.*, vol. 201, no. 1, pp. 1-10, Feb. 2010.
- [2] E. Anderson, Z. Bai, C. Bischof, L.S. Blackford, J. Demmel, J.J. Dongarra, J. Du Croz, S. Hammarling, A. Greenbaum, A. McKenney, and D. Sorensen, *LAPACK Users’ Guide (Third Edition)*. Philadelphia, PA, USA: SIAM, 1999, .
- [3] M. Strout, L. Carter, and J. Ferrante, “Sparse Tiling for Stationary Iterative

Methods,” *Int’l J. High Perform. Comput. Appl.*, vol. 18, no. 1, pp. 95-113, Feb. 2004.

[4] D. Culler, R. Karp, D. Patterson, A. Sahay, E.K. Schauer, E. Santos, R. Subramonian, and T.V. Eicken, “LogP: Towards a Realistic Model of ParallelComputation,” *ACMSigplan notices*, vol. 28, no. 7, pp. 1-12, July 1993.

[5] Xilinx DS816 Floating-Point Operator v6.0, 2012. [Online]. Available: [http://www.xilinx.com/support/documentation/ip\\_documentation/floating\\_point/v6\\_0/ds816\\_floating\\_point.pdf](http://www.xilinx.com/support/documentation/ip_documentation/floating_point/v6_0/ds816_floating_point.pdf)