

POSTER: GPU based Near Data Processing for Image Processing with Pattern Aware Data Allocation and Prefetching

Keshav Mishra

Indian Institute of Technology Bhilai
Computer Science and Engineering
12341140

Abstract-This paper proposes a GPU-based Near-data- processing (NDP) architecture as well as a well-matched programming model considering both the characteristics of image applications and NDP constraints. First, data allocation to the processing unit is handled to keep the data locality considering the memory access pattern. Second, this predictable allocation enables to design a compact but efficient NDP architecture. By applying a prefetcher that leverages the pattern aware data allocation, the number of active warps and on-chip SRAM size of NDP are significantly reduced. This allows to satisfy the NDP constraints and increases the opportunity to integrate more processing units on a memory logic die. The evaluation results for various image processing benchmarks show that the proposed NDP GPU improves the performance compared to the baseline GPU.

I INTRODUCTION

Recent applications with high resolution and frame rates have very large energy and processing time overhead because of the data movement between the processor and the off-chip memory [6]. Near-data-processing (NDP) is a promising alternative for reducing the overhead caused by data movement. Programmable GPU-based NDP is attractive because it is capable of processing a wider range of algorithms than hardware accelerators. Recent studies have mainly been conducted to accelerate specific applications such as deep neural networks (DNN) [3] or 3D rendering [7]. Therefore, previous NDP studies may not be extensively applied to various image processing algorithms. NDP for image processing has

not been sufficiently studied despite the strong demand for efficient image processing in many fields.

In this paper, a GPU-based NDP architecture and a programming model are proposed. The two approaches are complementary and are optimized together. In a general GPU, a large register size is needed for latency hiding. However, a GPU dedicated to image processing can be made more compact than general GPUs while maintaining a high performance level. The proposed programming model allocates data to the processing unit by considering the memory access pattern and data locality of the image processing algorithms. Taking advantage of this predictable allocation, an access pattern-aware prefetcher is added, and the hardware resources are greatly reduced.

II PROPOSED NDP GPU

The proposed GPU-based NDP allows not only to reduce the energy consumption by decreasing the memory access overhead but also to exploit the wide memory bandwidth. Also, careful design efforts are made to consider the 3D stacked memory environment, including the limited area and power constraints of the logic die.

A. Access Pattern Aware Data Allocation

To take advantage of the data access pattern in image processing applications, cooperative thread array (CTA)-to-streaming multiprocessor (SM) allocation is performed based on the data area information processed by CTAs. The allocation schemes are classified into row-major, column-major, and tile-based. Software-based data allocation uses the feature that CTA is not deallocated before the job assigned to that CTA is finished. Only a single CTA is allocated to each SM, and each CTA processes their data through loop iteration. Fig. 1 shows an example of software-based realization of row-major data allocation. The left side of Fig. 1 represents that the entire image is divided into four areas, which are allocated to four SMs. The right side of Fig. 1 magnifies the image area allocated to SM1, where each of the boxes represents the area processed by one CTA, with the width and height of $blockDim.x$ and $blockDim.y$, respectively. The area surrounded by a gray dotted line is called a row block.

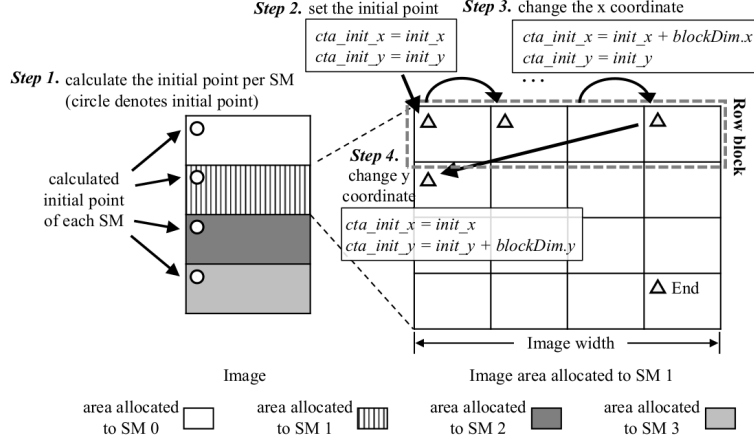


Figure 1: Software-based row major pattern aware data allocation example.

First, the initial coordinate of the area allocated to each SM is calculated (*Step 1*). Then, the top-left coordinate for the first CTA execution denoted by a triangle is set to be the initial coordinate (*Step 2*). Following the CTA execution, the x-coordinate is moved in the row-major direction (*Step 3*). *Step 3* is repeated until the execution is finished in the first row block, and then the process is moved to the next row block (*Step 4*). *Steps 3* and *Step 4* are repeated until the execution is finished for the entire data area allocated to SM0.

B. NDP GPU Architecture

SM Structure: The number of warps and the size of on-chip SRAM that are appropriate for the SM structure of NDP GPU are proposed in this section.

- (1) **Pattern-aware prefetcher:** A stride-based prefetcher is mounted inside the SM of NDP GPU. Because the data access pattern of the algorithm is preserved through the proposed data allocation, accurate prefetching is provided with a simple stride-based prefetcher.
- (2) **The number of active warps:** The access pattern-aware data allocation and prefetching make the necessary data readily available. Accordingly, as the stall decreases, the latency may be hidden with a smaller number of warps compared to the baseline GPU.
- (3) **On-chip SRAM size:** The size of on-chip SRAM, including the register, L1 data cache, and shared memory, is reduced by using a smaller number of warps.

Off-chip L2 Cache: The proposed NDP GPU uses no L2 cache. Data allocation and the prefetcher in the proposed NDP GPU help to increase the

hit rate of the L1 cache to nearly 100%, reducing the need for an L2 cache, which typically mitigates the miss penalty of the L1 cache in the memory hierarchy.

The number of SMs: Area constraints of the logic die should be considered along with performance and energy consumption to determine the optimal number of SMs for the NDP GPU.

III EVALUATION

A cycle-accurate simulator, GPGPU-Sim [1], is modified to verify the impact of the proposed NDP architecture. Four configurations are evaluated as listed below.

- Host: This is a conventional GPU. It is similar to the NVIDIA Maxwell GPU [2] and has 64 SMs.
- B-NDP-L2: This denotes a baseline NDP with L2 cache. It has the same SM architecture as Host, but it is located on the logic die.
- B-NDP-woL2: This denotes a baseline NDP where the L2 cache is not used. Other configurations are the same as with B-NDP-L2.
- IP-NDP: This is the proposed NDP specialized for image processing. It has a prefetcher and a reduced-size SRAM.

Six commonly used image processing algorithms are tested [4, 5]. For the evaluation of IP-NDP, the code is modified according to the proposed method.

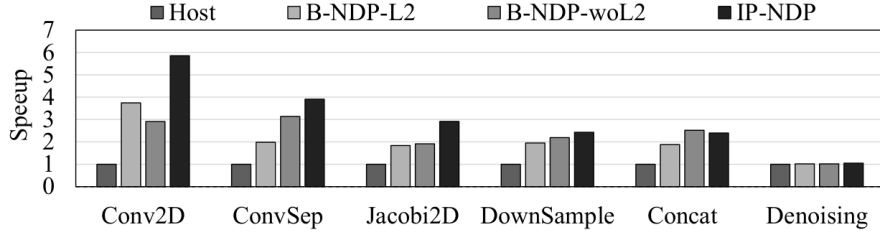


Figure 2: Speedup for B-NDP-L2, B-NDP-woL2, IP-NDP with varying the number of SMs in the NDP stack. All results are normalized to Host.

Fig. 2 shows the increase in speed for the four configurations normalized to Host. Two observations from this result are explained below. First, image processing benefits from the increased bandwidth. B-NDP-L2 shows a $2.06\times$ performance improvement on average over Host. Second, the proposed IP-NDP outperforms other configurations. IP-NDP is $3.09\times$ (up to $5.85\times$), $1.44\times$ (up to $1.97\times$), and $1.31\times$ (up to $2.01\times$) faster than Host, B-NDP-L2, and B-NDP-woL2, respectively.

Thus, it can be concluded that the proposed pattern-aware data allocation and NDP GPU works very well. Although not shown in this paper, the proposed scheme consumes 52% less energy than Host due to lower static power and faster execution, while requiring only 57.3% of the area compared to Host.

ACKNOWLEDGMENT

This paper was the result of the research project supported by SK hynix Inc. and by the R&D program of MOTIE/KEIT [No. 10077609, Developing Processor-Memory-Storage Integrated Architecture for Low Power, High Performance Big Data Servers]

REFERENCES

- [1] A. Bakhoda, G. L. Yuan, W. W. Fung, H. Wong, and T. M. Aamodt. Analyzing cuda workloads using a detailed gpu simulator. In *2009 IEEE International Symposium on Performance Analysis of Systems and Software*, pages 163–174. IEEE, 2009.
- [2] NVIDIA Corporation. Gtx 980: Featuring maxwell, the most advanced gpu ever made. White paper, NVIDIA Corporation, 2014.
- [3] M. Gao, J. Pu, X. Yang, M. Horowitz, and C. Kozyrakis. Tetris: Scalable and efficient neural network acceleration with 3d memory. In *ACM SIGARCH Computer Architecture News*, volume 45, pages 751–764. ACM, 2017.
- [4] S. Grauer-Gray, L. Xu, R. Searles, S. Ayalasomayajula, and J. Cavazos. Auto-tuning a high-level language targeted to gpu codes. In *2012 Innovative Parallel Computing (InPar)*, pages 1–10. IEEE, 2012.
- [5] NVIDIA. Nvidia cuda sdk 4.2, 2011.
- [6] I. K. Park, N. Singhal, M. H. Lee, S. Cho, and C. Kim. Design and performance evaluation of image processing algorithms on gpus. *IEEE Transactions on Parallel and Distributed Systems*, 22(1):91–104, 2011.
- [7] C. Xie, S. L. Song, J. Wang, W. Zhang, and X. Fu. Processing-in-memory enabled graphics processors for 3d rendering. In *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 637–648. IEEE, 2017.