



HAL
open science

Leveraging Subgraph Extraction for Performance Portable Programming Frameworks on DL Accelerators

Xiao Zhang, Huiying Lan, Tian Zhi

► **To cite this version:**

Xiao Zhang, Huiying Lan, Tian Zhi. Leveraging Subgraph Extraction for Performance Portable Programming Frameworks on DL Accelerators. 15th IFIP International Conference on Network and Parallel Computing (NPC), Nov 2018, Muroran, Japan. pp.179-184, 10.1007/978-3-030-05677-3_21 . hal-02279540

HAL Id: hal-02279540

<https://inria.hal.science/hal-02279540v1>

Submitted on 5 Sep 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Leveraging Subgraph Extraction for Performance Portable Programming Frameworks on DL Accelerators

Xiao Zhang^{1,2,3}, Huiying Lan¹, and Tian Zhi¹

¹ Intelligent Processor Research Center, Institute of Computing Technology (ICT), CAS, China

² University of Chinese Academy of Sciences (UCAS), China

³ Cambricon Tech. Ltd.

Abstract. Deep learning framework plays an important role in connecting hardware platform and algorithm. In recent years, some domain-specific deep learning accelerators with better performance and energy efficiency were proposed by researchers. However, current frameworks lack enough considerations about how to better support the possible new features brought by accelerators. In this paper, we propose to build a performance portable programming framework with subgraph extraction. The intuition is that increasing ratio of optimizations are taken from the top-level framework to the low-level software stack of accelerator. In response to this development trend, framework needs to pay more attention to the splitting strategy of computation graph for the heterogeneous computation.

1 Introduction

In recent years, we have witnessed many significant breakthroughs of deep learning algorithm in a multitude of domains. This superior accuracy, however, comes at the cost of high computational complexity. Researchers try to design more efficient architectures based on the features of deep learning algorithm and get some promising results [3, 5, 10, 4, 7–9]. These results show that domain-specific accelerators outstand in both speed and energy efficiency compared to traditional solutions.

On the other hand, in order to explore and deploy deep learning algorithm conveniently, both academia and industry have developed several deep learning frameworks, such as MXNet [2], TensorFlow [1] and Caffe [6]. Those frameworks automatically optimize the computation flow, generate high-performance kernels and schedule kernels in parallel if possible.

However, there is a gap between emerging DL accelerators and existing programming frameworks. In order to run deep learning algorithm with the highest performance, some accelerators and its software stacks have tried to break the wall and search optimal solution in a large space. Unfortunately, current deep learning frameworks only provide limited adaptations for this new feature.

2 Motivation

2.1 DLA and Graph Fusion

We designed and implemented a deep learning accelerator and its software stack, and we call the accelerator DLA in following sections. The design of DLA is concluded from multiple deep learning accelerators, including NVidia DLA, DaDianNao [4] and TPU. There are multiple cores in DLA. Each core in DLA can complete a computation task independently, which makes it actually a parallel model with shared global memory.

Compared to traditional limited method that fusing some specific sequence composed of element-wise operators issued by framework, software stack of DLA offers a more radical solution. It optimizes and fuses the total graph (see the Figure 1). This strategy has several benefits. First, the experts developed lower stack can give better solution because they know more about hardware architecture. Also, fusing a large graph into a single node greatly saves the kernel launch cost, which is important for inference task.

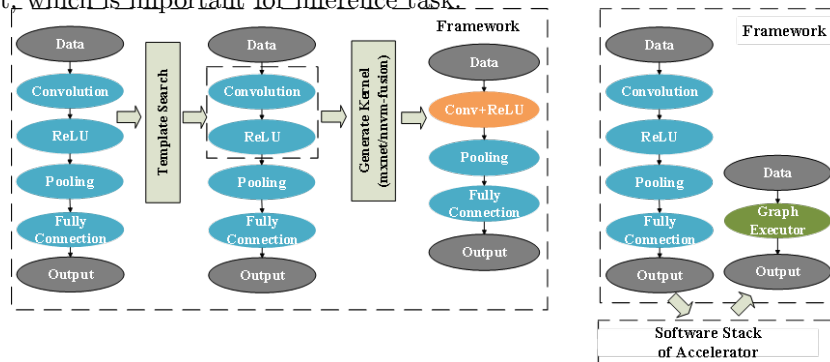


Fig. 1: In left part, the framework searches limited templates and generates new kernels to replace them. In right part, optimization stack of accelerator receives the whole graph, optimizes and generates a new executor back to framework.

2.2 Heterogeneous Computation

Heterogeneous computation is unavoidable for DLA and other accelerators. Some operators in new algorithms are hard to parallelize or to abstract to the tensor operators offered by accelerators, and the frequency of embedding accelerator in mobile device might be reduced to save energy. As a result, assigning some parts on CPU might bring better total performance. Thus, before we use lower software stack to optimize graph, we need to extract a subgraph composed by operators assigned on DLA. In other words, framework should have a clever split strategy and method to extract appropriate subgraph from the original deep networks.

3 Subgraph Extraction

When we try to extract a subgraph based on whether each operator is well-supported by accelerator, the direct intuition is to make it a maximum connected

convex subgraph. Connectivity guarantees data relation between operators which is necessary for most optimizing methods. Maximum grants the largest searching space and reduces kernel launch overheads. Convexity is used as a constraint to avoid circle which leads to dead lock when scheduling. A subgraph \mathbf{S} of a directed acyclic graph \mathbf{G} is convex if and only if there is no directed path between two vertices of \mathbf{S} which contains an arch not in \mathbf{S} (see the Figure 2).

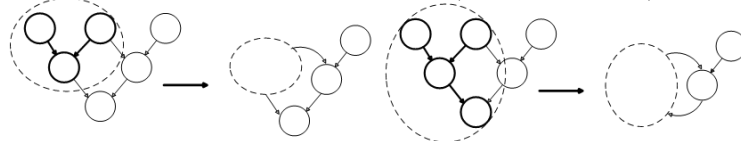


Fig. 2: Example of convex and non-convex subgraph.

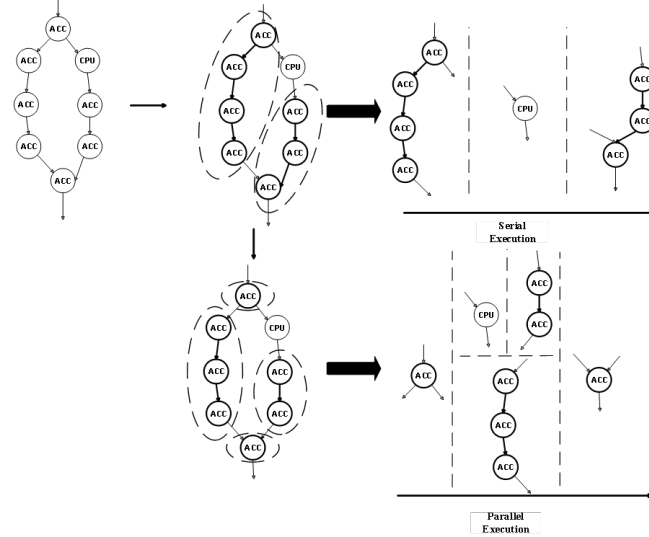


Fig. 3: Post-prune strategy. The ACC node represents operator assigned on accelerator, and the CPU node represents operator assigned on CPU.

Merging a large subgraph into a single node helps the corresponding computation to run faster, however, it may hinder scheduler to get maximum parallelism in some case. As Figure 3 shows, the fused graph must wait for all its input to be ready even though some inputs are not necessary at the early stage of its computation. Similarly, although not all the outputs of a subgraph are generated at the final stage, all descendants must keep waiting until computation of total subgraph finishes. So, we append a post-prune process to split each subgraph into smaller parts, each of which has only one input and output operator.

4 Evaluation

The experiment platform is DLA, a multi-core deep learning accelerator as we mentioned before. We first evaluate the performance before and after the graph fusion to demonstrate the validation of graph fusion. As shown in Figure 4, performance of all six entire-network benchmarks are improved, which achieves a

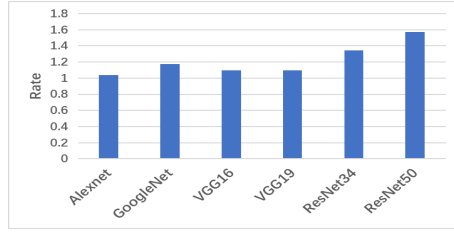


Fig. 4: Relative speedup of graph.

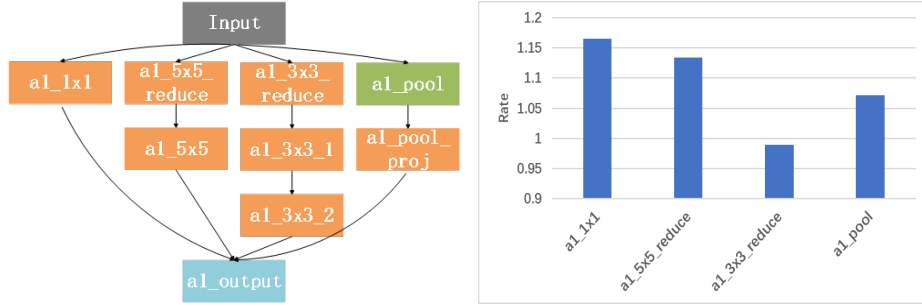


Fig. 5: Left figure shows the structure of the inception-v3 block. Right figure shows speedup after the post-prune strategy. Horizontal axis label represents part of the block assigned to CPU

speedup of $1.18\times$ on average compared with the baseline, which we do not implement the graph fusion. Specifically, the improvement of ResNet34 and ResNet50 is clearly higher than other four networks.

Then we evaluate the speedup of the post prune process. We use the intuitive maximum connected convex subgraph extraction strategy as the baseline. In order to accurately evaluate the prune strategy, we choose a basic block of operators with multiple branches from inception-v3 networks for its enough braches. To trigger subgraph extraction, we seperately assign operators on different branch to CPU and evaluate the speedup. As the result shown in the figure 5, except for assigning operator on the critical path to CPU, performance of the other three heterogeneous computation get a speedup of $1.1\times$ on average, which is an obvious improvement.

5 Conclusion

In this paper, we propose a performance portable programming framework. The key motivation is that framework needs a subgraph extraction strategy to better balance schedule parallelism and fusion efficiency. We implement such a framework by migrating MXNet. This strategy is designed to cooperate framework with lower software stack in heterogeneous computation task, because none of them can complete the whole task independently. This strategy can be used in a wider field if accelerators choose to take over framework to optimize the computation graph by themselves.

Acknowledgment

This work is partially supported by the National Key Research and Development Program of China (under Grant 2017YFA0700902, 2017YFB1003101), the NSF of China (under Grants 6147239, 61432016, 61473275, 61522211, 61532016, 61521092, 61502446, 61672491, 61602441, 61602446, 61732002, 61702478), the 973 Program of China (under Grant 2015CB358800), National Science and Technology Major Project (2018ZX01031102) and Strategic Priority Research Program of Chinese Academy of Sciences (XDBS01050200).

References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M.: Tensorflow: a system for large-scale machine learning (2016)
2. Chen, T., Li, M., Li, Y., Lin, M., Wang, N., Wang, M., Xiao, T., Xu, B., Zhang, C., Zhang, Z.: Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *Statistics* (2015)
3. Chen, T., Du, Z., Sun, N., Wang, J., Wu, C., Chen, Y., Temam, O.: Diannao: a small-footprint high-throughput accelerator for ubiquitous machine-learning. In: *Proceedings of the 19th international conference on Architectural support for programming languages and operating systems*. pp. 269–284. ACM (2014)
4. Chen, Y., Luo, T., Liu, S., Zhang, S., He, L., Wang, J., Li, L., Chen, T., Xu, Z., Sun, N.: Dadiannao: a machine-learning supercomputer. In: *Ieee/acm International Symposium on Microarchitecture*. pp. 609–622 (2014)
5. Du, Z., Fasthuber, R., Chen, T., Ienne, P., Li, L., Luo, T., Feng, X., Chen, Y., Temam, O.: Shidiannao. *Acm Sigarch Computer Architecture News* 43(3), 92–104 (2015)
6. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guaradarama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093* (2014)
7. Liu, D., Chen, T., Liu, S., Zhou, J., Zhou, S., Teman, O., Feng, X., Zhou, X., Chen, Y.: Pudiannao: A polyvalent machine learning accelerator. In: *Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems*. pp. 369–381 (2015)
8. Liu, S., Du, Z., Tao, J., Han, D., Luo, T., Xie, Y., Chen, Y., Chen, T.: Cambricon: An instruction set architecture for neural networks. In: *Proceedings of the 43rd International Symposium on Computer Architecture*. pp. 393–405. IEEE Press (2016)
9. Reagen, B., Whatmough, P., Adolf, R., Rama, S., Lee, H., Lee, S.K., Hernández-Lobato, J.M., Wei, G.Y., Brooks, D.: Minerva: Enabling low-power, highly-accurate deep neural network accelerators. In: *ACM SIGARCH Computer Architecture News*. vol. 44, pp. 267–278. IEEE Press (2016)
10. Zhang, S., Du, Z., Zhang, L., Lan, H., Liu, S., Li, L., Guo, Q., Chen, T., Chen, Y.: Cambricon-x: An accelerator for sparse neural networks. In: *Ieee/acm International Symposium on Microarchitecture*. pp. 1–12 (2016)