

### PipeFuser: Building Flexible Pipeline Architecture for DNN Accelerators via Layer Fusion

Xilang Zhou<sup>1</sup>, Shuyang Li<sup>1</sup>, Haodong Lu<sup>2</sup>, Kun Wang<sup>1,\*</sup>

<sup>1</sup>State Key Lab of ASIC & System, Fudan University, Shanghai, China <sup>2</sup>Nanjing University of Posts and Telecommunications, Nanjing, China. kun.wang@ieee.org

### Introduction

### • Preliminaries

- Background
- Design Challenges
- PipeFuser Design
  - Microarchitecture and Computation Flow
  - Co-design Engine
- Experiment
- Conclusion

### Introduction



Three architectural paradigms of DNN accelerators:

(a) non-pipeline: a unified accelerator processing all layers

(b) full-pipeline: Dedicated accelerators tailored for each individual DNN layer, fully pipelined

(c) fused-pipeline : Pipelined accelerators of fused layer groups

### Introduction

**Non-pipeline**: need to accommodate diverse DNN layers, frequent off-chip memory access

**Full-pipeline**: scarce resource and large on-chip buffering demand with increasing model depth; pipeline imbalance

**Fused-pipeline**: improved pipeline balance; alleviate resource constraints; less on-chip storage demand

### Introduction

### **Contribution:**

- We leverage the layer fusion technique and fused-pipeline architecture for DNN accelerator deployment.
- We develop an end-to-end automation framework, **PipeFuser**, for fused-pipeline deployment.
- A GA-based co-design engine for fast fusion strategy determination.
- Competitive performance with significant performance speedup and higher flexibility compared to non-pipeline and full-pipeline architectures.

• Introduction

### Preliminaries

- Background
- Design Challenges
- PipeFuser Design
  - Microarchitecture and Computation Flow
  - Co-design Engine
- Experiment
- Conclusion

# **Background on Layer Fusion Technique**

First proposed in [1] Combined with **tiling** technique Reduces intermediate data transfer between on-chip and off-chip within the fusion group



[1]M. Alwani, H. Chen, M. Ferdman, and P. Milder, "Fused-layer cnn accelerators," in MICRO, 2016.

# **Performance analysis via Roofline Model**

- Compute-To-Communication (CTC) ratio: # of MAC operations per byte
- Full-pipeline and fused-pipeline: computation-bounded
- Fused-pipeline arechitecture is more resource-efficient and performance-optimized



## **Design Challenges**

- Challenge 1: How to formulate an efficient fusion strategy.
- Challenge 2: How to explore the expansive mapping space.

↓ Co-optimization problem

- Introduction
- Preliminaries
  - Background
  - Design Challenges

### PipeFuser Design

- Microarchitecture and Computation Flow
- Co-design Engine
- Experiment
- Conclusion

### **Framework Overview**



### Three main components of

PipeFuser:

Model/Hardware Parser

**Co-design Engine** 

Hardware generator

# **Microarchitecture and Computation Flow**





**Computation Flow** 

### Microarchitecture of Pipefuser Modules:

- Computation
- Memory
- Control Unit

# **Co-design Engine**

The co-design Engine is based on a Genetic Algorithm GA is a search heuristic inspired by the process of natural selection

Four phases are considered in the GA:

- Initialization
- Crossover
- Mutation
- Evaluation and Selection

# **Co-design Engine**

The fusion and mapping parameters (template parameters) need to be encoded into genome form first



Notation:

 $T_r$  and  $T_c$ : dimensions of the tile size (global)  $T_m$  and  $T_n$ : parallelism dimensions G: array of layers that are fused together

Total number of chromosomes = The number of fusion groups = the number of accelerators

# **Co-design Engine**

### Algorithm Flow

Algorithm 1: Explore Fusion and Mapping Spaces

```
Data: the probability of the crossover operation P_c.
         the probability of the mutation operation P_m.
         the number of the solutions within the
         population T, the number of generation G, the
         total layer of the net N_{layer}
  Result: optimal number of accelerator n_{acc}, optimal
           solution S
1 while n_{acc} \leq N_{layer} do
      initialize population X of T individuals;
2
      while q \leq = G do
3
          evaluate fitness of each individual:
4
          select parents for next generation;
5
          if random(0,1) < P_c then
6
              create offspring through crossover;
7
              exclude ineligible offspring:
 8
9
          end
          if random(0,1) < P_m then
10
              create offspring through mutation;
11
              exclude ineligible offspring;
12
          end
13
      end
14
15 end
```

**Initialization:** initialize population X of T individuals **Crossover:** 

- directly swapping T<sub>r</sub> and T<sub>c</sub> between two sets of solutions.
- T<sub>m</sub> and T<sub>n</sub> are swapped within the same fusion group

### **Mutation:**

- Mutation-Parallel-Dim
- Mutation-Tile Size
- Mutation-Group

#### Evaluation and Selection: tournament selection

- Introduction
- Preliminaries
  - Background
  - Design Challenges
- PipeFuser Design
  - Microarchitecture and Computation Flow
  - Co-design Engine

### • Experiment

• Conclusion

# **Experiment Setup**

### Platform

Xilinx Alveo U200 FPGA hosted on Intel Xeon Server CPU (Gold 5218R@2.10GHz), synthesized with Vitis 2022.1

### Hardware Baseline

non-pipeline: HybridDNN<sup>[1],</sup> FlexCNN<sup>[2]</sup> full-pipeline: DNNbuilder<sup>[3]</sup>, TGPA<sup>[4]</sup> hybrid architecture: DNNExplorer<sup>[5]</sup> segment-grained pipeline: Deepburning-SEG<sup>[6]</sup>

### DNN Workload

VGGNet, ResNet

### 8-bit fixed-point data representation

H. Ye, X. Zhang, Z. Huang, G. Chen, and D. Chen, "Hybriddnn: A framework for high-performance hybrid dnn accelerator design and implementation," in DAC, 2020.
 S. Basalama, A. Sohrabizadeh, J. Wang, L. Guo, and J. Cong, "Flexcnn: An end-to-end framework for composing cnn accelerators on fpga," in TRETS, 2023.
 X. Zhang, J. Wang, C. Zhu, Y. Lin, J. Xiong, W.-m. Hwu, and D. Chen, "Dnnbuilder: An automated tool for building high-performance dnn hardware accelerators for fpgas," in ICCAD, 2018.
 X. Wei, Y. Liang, X. Li, C. H. Yu, P. Zhang, and J. Cong, "Tgpa: Tile-grained pipeline architecture for low latency cnn inference," in ICCAD, 2018.

[5] X. Zhang and D. Chen, "Dnnexplorer: a framework for modeling and exploring a novel paradigm of fpga-based dnn accelerator," in ICCAD, 2020.
 [6] X. Cai, Y. Wang, X. Ma, Y. Han, and L. Zhang, "Deepburning-seq: Generating dnn accelerators of segment-grained pipeline architecture," in MICRO, 2022.

### **Experiment: Comparison with Previous Designs**

Design HybridDNN [15] FlexCNN [3] **DNNBuilder** [21] **TGPA** [13] **DNNExplorer** [20] **DeepBurning-SEG** [4] Ours Device VU9P U250 KU115 VU9P KU115 KU115 U200 Model VGG16 VGG16 VGG16 VGG19 VGG16 VGG16 VGG16 16bit (8bit) Precision 12bit 16bit (8bit) 16bit 16bit 8bit 8bit 167 235 235 250 Freq.(MHz) 241 210 200 5136 4736 4318 4096 4444 5128 5026 **DSPs** (75.9%)(37.98%)(78%)(60%)(80.5%)(92.9%)(73.5%)2692 1578 1690 1648 1486 1310 BRAM NA (45.93%)(81.8%)(81%)(78%)(76.3%)(76.9%)Throughput 3376 1543 (2329) 2011 (4022) 1702 4778 7684 1510 (GOPS) **DSP Density** 0.66 0.33 0.93 0.37 0.38 0.93 1.52 (GOPS / DSPs)

### TABLE I Performance Comparison with Previous Works

### **Experiment: Full-pipeline vs Fused-pipeline**



### **Experiment: Scalability Evaluation**



- With Limited Resources
- With Increasing Layer Depth

#### TABLE II Scalability Evaluation

Model	VGG11	VGG13	VGG16	VGG19	ResNet50	ResNet152
Group.number	3	3	4	5	4	6
Freq.(MHz)	250	250	250	235	220	220
DSPs	3982	4390	5026	5621	3560	4496
BRAM	982	1073	1310	1582	1244	1543
Throughput (GOPS)	3560	5231	7684	8431	2106	5415

- Introduction
- Preliminaries
  - Background
  - Design Challenges
- PipeFuser Design
  - Microarchitecture and Computation Flow
  - Co-design Engine
- Experiment
- Conclusion

# Conclusion

PipeFuser

- Optimized Fused-Pipeline Architecture
- Automated End-to-End Framework

GA-based Co-design Engine

- Evaluation
- Higher Throughput than Non-pipeline and Full-pipeline Architectures
- Higher DSP density
- Better Deployment Flexibility

### Acknowlegdement

This work was financially supported in part by National Key Research and Development Program of China under Grant 2021YFA1003602, and in part by Shanghai Pujiang Program under Grant 22PJD003.





# Thank you!

kun.wang@ieee.org