#### ECE 6775 High-Level Digital Design Automation Fall 2024

# Deep Learning Acceleration on FPGAs



**Cornell University** 



#### Announcements

- Final project
  - In-depth exploration of a topic related to high-level design automation
    - (1) Designing new accelerators with HLS; <u>OR</u>
      (2) Developing new automation algorithms/tools
  - 3 or 4 students per team, up to 12 teams
    - 12 teams = 10 \* 4 students + 2 \* 3 students
  - Weekly meetings start next week
    - A Google sheet is created for meeting scheduling, sign up by Saturday noon
  - More project guidelines coming up in the next lecture
  - Project abstract due Friday 11/8

#### Agenda

- CNN acceleration on FPGAs
- Potential of FPGA-based LLM inference

### **Recap: Convolutional Neural Network (CNN)**



- Front: convolutional layers learn visual features
- Back: fully-connected layers perform classification using the visual features

#### **Convolutional Layer**



- An output pixel is connected to its neighboring region on each input feature map (fmap)
- All pixels on an output feature map use the same filter weights



- Four main sources of parallelism
  - 1. Across input feature maps (i.e., input channels)



- Four main sources of parallelism
  - 1. Across input feature maps (i.e., input channels)
  - 2. Across output feature maps (i.e., output channels)



- Four main sources of parallelism
  - 1. Across input feature maps (i.e., input channels)
  - 2. Across output feature maps (i.e., output channels)
  - 3. Across different output pixels (i.e., filter positions)



- Four main sources of parallelism
  - 1. Across input feature maps (i.e., input channels)
  - 2. Across output feature maps (i.e., output channels)
  - 3. Across different output pixels (i.e., filter positions)
  - 4. Across filter pixels

#### **Parallelism in the Code**



## Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks

Cheng Zhang<sup>1</sup>, Peng Li<sup>2</sup>, Guangyu Sun<sup>1,3</sup>, Yijin Guan<sup>1</sup>, Bingjun Xiao<sup>2</sup>, Jason Cong<sup>2,3,1</sup>

<sup>1</sup> Peking University

<sup>2</sup> UCLA

<sup>3</sup> PKU/UCLA Joint Research Institute in Science and Engineering

FPGA 2015

#### **Main Contributions**

- 1. Analysis of the different sources of parallelism in the convolution kernel of a CNN
- 2. Quantitative performance modeling of the hardware design space using the Roofline method
- 3. Design and implementation of a CNN accelerator for FPGA using Vivado HLS, evaluated on AlexNet

#### **Challenges to FPGA Acceleration**

- We can't just unroll all the loops due to limited FPGA resources
- Must choose the right code transformations to exploit the parallelism in a resource efficient way



## **Loop Tiling**



Offloading just the inner loops requires only a small portion of the data to be stored on FPGA chip



Loop Tiling

## **Code with Loop Tiling**



for (tii=ti; tii<min(ti+Ti, I); tii++) {
 for (ki=0; ki<K; ki++) {
 for (kj=0; kj<K; kj++) {
 output\_fm[too][trr][tcc] +=
 weights[too][tii][ki][kj]\*input\_fm[tii][S\*trr+ki][S\*tcc+kj];
 }}
}</pre>

// software: read output feature map

10

}}}

5	for (trr=row; trr <min(row+tr, r);="" th="" trr++)="" {<=""></min(row+tr,>
6	for (tcc=col; tcc <min(col+tc, c);="" tcc++)="" th="" {<=""></min(col+tc,>
7	<b>for</b> (too=to; too <min(to+to, o);="" th="" too++)="" {<=""></min(to+to,>
8	for (tii=ti; tii <min(ti+ti, i);="" td="" tii++)="" {<=""></min(ti+ti,>
9	<b>for</b> (ki=0; ki <k; ki++)="" td="" {<=""></k;>
10	<b>for</b> (kj=0; kj <k; kj++)="" td="" {<=""></k;>
	output_fm[too][trr][tcc] +=
	weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
	}}}}

9	for $(ki=0; ki {$
10	for $(kj=0; kj {$
5	for (trr=row; trr <min(row+tr, r);="" th="" trr++)="" {<=""></min(row+tr,>
6	<b>for</b> (tcc=col; tcc <min(col+tc, c);="" tcc++)="" th="" {<=""></min(col+tc,>
7	<b>for</b> (too=to; too <min(to+to, o);="" td="" too++)="" {<=""></min(to+to,>
8	for (tii=ti; tii <min(ti+ti, i);="" td="" tii++)="" {<=""></min(ti+ti,>
	output_fm[too][trr][tcc] +=
	weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
	}}}}

9	for (ki=0; ki <k; ki++)="" th="" {<=""></k;>
10	for (kj=0; kj <k; kj++)="" td="" {<=""></k;>
5	for (trr=row; trr <min(row+tr, r);="" td="" trr++)="" {<=""></min(row+tr,>
6	<pre>for (tcc=col; tcc<min(col+tc, c);="" pre="" tcc++)="" {<=""></min(col+tc,></pre>
	#pragma HLS pipeline
7	<b>for</b> (too=to; too <min(to+to, o);="" td="" too++)="" {<=""></min(to+to,>
8	for (tii=ti; tii <min(ti+ti, i);="" td="" tii++)="" {<=""></min(ti+ti,>
	output_fm[too][trr][tcc] +=
	weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
	}}}}

9	for (ki=0; ki <k; ki++)="" th="" {<=""></k;>
10	for (kj=0; kj <k; kj++)="" td="" {<=""></k;>
5	for (trr=row; trr <min(row+tr, r);="" td="" trr++)="" {<=""></min(row+tr,>
6	<pre>for (tcc=col; tcc<min(col+tc, c);="" pre="" tcc++)="" {<=""></min(col+tc,></pre>
	#pragma HLS pipeline
7	<b>for</b> (too=to; too <min(to+to, o);="" td="" too++)="" {<=""></min(to+to,>
	#pragma HLS unroll
8	for (tii=ti; tii <min(ti+ti, i);="" td="" tii++)="" {<=""></min(ti+ti,>
	output_fm[too][trr][tcc] +=
	weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];
	}}}

9	for (ki=0; ki <k; ki++)="" th="" {<=""><th></th></k;>		
10	for (kj=0; kj <k; kj++)="" td="" {<=""><td></td></k;>		
5	<pre>for (trr=row; trr<min(row+tr, pre="" r);="" trr++)="" {<=""></min(row+tr,></pre>		
6	<pre>for (tcc=col; tcc<min(col+tc, c);="" pre="" tcc++)="" {<=""></min(col+tc,></pre>		
	#pragma HLS pipeline		
7	<b>for</b> (too=to; too <min(to+to, o);="" td="" too++)="" {<=""></min(to+to,>		
	#pragma HLS unroll	/ <b>T</b> :\	
8	for (tii=ti; tii <min(ti+ti, i);="" td="" tii++)="" {<=""><td>(1)</td></min(ti+ti,>	(1)	
	#pragma HLS unroll and output (To) channels		
	output_fm[too][trr][tcc] +=		
	weights[too][tii][ki][kj]*input_fm[tii][S*trr+ki][S*tcc+kj];		
	<pre>}}}}</pre>		

*L* is the pipeline depth (# of pipeline stages, II=1)

Number of cycles to execute the above loop nest  $\approx K \times K \times Tr \times Tc + L \approx Tr \times Tc \times K^2$ 



#### **Generated Hardware**

Performance and size of the accelerator determined by tile factors Ti and To

Number of data transfers determined by Tr and Tc

#### **Overall Accelerator Architecture**



#### **Design Space Complexity**

- Challenge: Number of available optimizations present a huge space of possible designs
  - What is the optimal loop order?
  - What tile size to use for each loop?
- Implementing and testing each design by hand will be slow and error-prone
  - Some designs will exceed the on-chip compute/memory capacity
- Solution: Performance modeling + automated design space exploration

#### **Performance Modeling**

- Three design metrics are estimated:
  - Total number of operations (FLOP)
    - Depends on the CNN model parameters
  - Total external memory access (Byte)
    - Depends on the CNN weight and activation size

#### – Total execution time (Sec)

- Depends on the hardware architecture (e.g., tile factors To and Ti)
- Ignore resource constraints for now

#### **Performance Modeling**

- Total operations FLOPs  $\approx 2 \times 0 \times I \times R \times C \times K^2$
- Execution time = Number of Cycles × Clock Period

- Number of cycles 
$$\approx \left[\frac{O}{To}\right] \times \left[\frac{I}{Ti}\right] \times \left[\frac{R}{Tr}\right] \times \left[\frac{C}{Tc}\right] \times (Tr \times Tc \times K^2)$$
  
 $\approx \left[\frac{O}{To}\right] \times \left[\frac{I}{Ti}\right] \times R \times C \times K^2$ 

- External memory accesses =  $ai \times Bi + aw \times Bw + ao \times Bo$ 
  - Size of input fmap buffer:  $Bi = Ti \times (Tr + K 1)(Tc + K 1)$  with stride=1
  - Size of output fmap buffer:  $Bo = To \times Tr \times Tc$
  - Size of weight buffer:  $Bw = Ti \times To \times K^2$
  - External access times:  $ao = \left[\frac{O}{To}\right] \times \left[\frac{R}{Tr}\right] \times \left[\frac{C}{Tc}\right]$ ,  $ai = aw = \left[\frac{I}{Ti}\right] \times ao$ 
    - Input features and weights are reused across multiple output tiles

### **Design Space Exploration with Roofline**

Which design point do you prefer?



#### **Hardware Implementation**



### Understanding the Potential of FPGA-Based Spatial Acceleration for Large Language Model Inference

Hongzheng Chen<sup>1</sup>, Jiahao Zhang<sup>2,1</sup>, Yixiao Du<sup>1</sup>, Shaojie Xiang<sup>1</sup>, Zichao Yue<sup>1</sup>, Niansong Zhang<sup>1</sup>, Yaohui Cai<sup>1</sup>, Zhiru Zhang<sup>1</sup>

<sup>1</sup> Cornell University <sup>2</sup> Tsinghua University

ACM TRETS 2024 (FCCM Journal Track)

#### The Era of Large Language Models (LLMs)





Is FPGA suitable for efficient LLM inference?

#### **Transformer GPU Execution Breakdown**

- LLMs aren't just about matrix multiplications (MatMul)
- Non-linear & element-wise operators also play a significant role
  - Low compute-to-memory ratio (namely, low OI)
  - High kernel launch overheads

<b>Operator Class</b>	Representative Op	% FLOP	% Run Time
Tensor contraction	MM, MV	99.80	61.0
Stat. normalization	softmax, layernorm	0.17	25.5
Element-wise	bias, dropout	0.03	13.5

**Proportions for operator classes in the BERT model,** implemented in PyTorch, profiled with an NVIDIA V100 GPU (MM: matrix-matrix multiply; MV: matrix-vector multiply)



Data source: Andrei Ivanov et al., "Data Movement is All You Need: A Case Study on Optimizing Transformers", arXiv:2007.00072, 2021.

#### **Two-Stage LLM Generative Inference**

#### Stage 1: Prefill

Takes in user prompts and generates the first token (seq\_len > 1)

#### Stage 2: Decode

Processes the previously generated token to produce new tokens one at a time in an autoregressive manner (seq\_len = 1)



#### ΜV



compute	mk		
memory footprint	mk + k + m		

#### **Typical GPU Execution Cycle for a Single Operator**

- Analogy Compute is like a factory, requiring instructions (kernel launch overhead) and a steady data supply (memory bandwidth) to run efficiently
- If compute efficiency grows faster than data supply, it limits the system's ability to operate at peak performance



#### **Typical GPU Execution Cycle for Multiple Operators**

- Moving data to and from GPU compute units incurs a high memory bandwidth cost
- For memory-bound operations, more time is spent moving data than performing computation



#### **The Potential of Dataflow Execution**

- Model-specific dataflow accelerator
  - Pass intermediate results directly to the next operator
  - Reduced memory traffic leads to higher performance and higher energy efficiency



#### **Contributions of This Work**

- An analytical model to estimate performance and resource/bandwidth usage for model-specific dataflow accelerator on FPGAs
  - **Compute resource**: M is compute power in MACs/cycle and C is the # of layers per FPGA

$$\sum M_i C < M_{\text{tot}}, i \in \{q, k, v, a_1, a_2, p, f_1, f_2\}$$

• Memory capacity: S is buffer size

 $S_{\text{param}}C \! < \! \text{DRAM}_{\text{tot}}, \\ \sum S_iC \! < \! \text{SRAM}_{\text{tot}}, i \! \in \! \{ \text{tile}, \! \text{KV}, \! \text{FIFO} \}$ 

Memory port: s is tensor size and b is bitwidth  

$$R_{i} = \left[\frac{s_{i}b_{BRAM}}{M_{i}/r_{i} \times S_{BRAM}}\right] \times \frac{M_{i}/r_{i}}{k}$$

$$\sum_{i} CR_{i} + 2C(R_{a_{1}} + R_{a_{2}}) < \mathsf{SRAM}_{\mathsf{tot}}, i \in \{q, k, v, p, f_{1}, f_{2}\}$$

- Memory bandwidth: B is bandwidth  $\sum_i \! CB_i \! < \! B_{\mathsf{tot}}, i \! \in \! \{q,\!k,\!v,\!p,\!f_1,\!f_2\}$
- An optimized HLS kernel library and compose accelerator designs for LLM on FPGAs



#### **Estimation on Different Models and Devices**

#### Is FPGA suitable for efficient LLM inference?

- Latency estimation of GPT2 and LLaMA2 on different FPGAs
- GPU results are obtained through actual profiling

#### Compute-bound



Existing FPGAs are inferior in the compute-intensive prefill stage, but can outperform GPUs in the memory-intensive <u>decode</u> stage.

#### Acknowledgements

- This tutorial contains/adapts materials developed by
  - Ritchie Zhao (Cornell ECE PhD, now NVIDIA)
  - Authors of the following papers
    - Optimizing FPGA-based Accelerator Design for Deep Convolutional Neural Networks (FPGA'15, PKU-UCLA)
    - Understanding the Potential of FPGA-Based Spatial Acceleration for Large Language Model Inference (ACM TRETS, FCCM'24 Journal Track)

#### **Next Lecture**

Project Guidelines