UF Herbert Wertheim College of Engineering UNIVERSITY of FLORIDA

Coarse-Grained Floorplanning for streaming CNN applications on Multi-Die FPGAs

Danielle Tchuinkou Kwadjo

ECE Department, University of Florida

Outline



- Rale Balancing
- Multi-level Graph Partitioning
- Coarse-grained Floor planning

Conclusion

INTRODUCTION AND MOTIVATION

- Introduction
- DL Evolution
- Challenges
- \circ FINN

1

2

3

4

Results

- Granularity Exploration
- Resource Utilization
- \circ Productivity





- DNN workloads have become significant.
- As the complexity of machine learning algorithms increases, more data becomes available.





Figure 2: Evolution of depth, error-rate, and number of parameters over the years



FPGA Overview and Benefits

- Both GPU and FPGA are growing fast in artificial intelligence acceleration area.
- GPU is now dominating the market as it has less engineering.
 - Perform poorly when it comes to inference as requests often arrive one at a time.
- However, compared to GPU, FPGA has several outstanding features:
 - **Flexibility:** FPGA allows engineers to reconfigure underlying hardware architecture.
 - Efficient with lower precision deep learning algorithms, such as binary neural network
 [1] and ternary neural network[2].
 - **Low latency:** FPGAs are capable of data parallelism, but also pipeline parallelism
 - High Power Efficiency:
 - Xilinx Virtex Ultrascale+, FPGA board has general purpose compute efficiency of 277 GOP/s/W,
 - NVidia Tesla P4, GPU produced by Nvidia, the efficiency is of 208 GOP/s/W



Figure 4: Xilinx Alveo U50



Department of Electrical & Computer Engineering





Related Work

Challenges

Benchma

CNN Accelerators

- Most FPGA inference accelerators are based on overlay architectures [1], [2], i.e., GPP MM circuits onto which compute is scheduled to execute the CNN layers in sequence.
 - This approach is flexible, as it enables potentially any CNN topology to be executed by a single accelerator,
 - However, it is not efficient \rightarrow frequent transfers of weights and activations from/to the memory







Related Work

Challenges

Benchma

CNN FPGA Accelerators

- An alternative FPGA accelerator architecture is streaming accelerator
 - CNN inference has achieved the lowest latency, highest throughput, and lowest power dissipation.
 - It is limited by the resources on the FPGA.
 - Better fit for cloud environment

















Introduction CNN & FPGA Related Work Challenges Benchmark **Proposed Framework** 1) **Rate Balancing**: we propose a MIP model to find the optimal parameters for a CNN accelerator.





Herbert Wertheim College of Engineering

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

UF

1)

2)

Challenges 1×3×224×224 C₁ ω_1 **Proposed Framework** ω1 C_1 W (24×3×3×3) B (24) C₂ ω_2 **Rate Balancing**: we propose a model C_2 ω2 Relu to find the optimal parameters for the ω_3 C₃ *C*₃ configuration of a CNN accelerator ω4 Conv C_4 ω_4 **Computational Graph** G = C_4 03 W (58×24×1×1) B (58) ω5 (V, E, ω, φ) C_5 C₅ ω_5 set of Vertices V, vertices weights ω, *ω*6 edge set E, edges weight φ. ω_6 C₆ C_7 C_6 W (24×1×3×3) W (58×1×3×3) C7 ω_7 ω8 B (24) B (58) ω7 ω_8 C₈ C9 ωg W (58×24×1×1) W (58×58×1×1) Cg C_8 B (58) B (58) ω10 ω9 ω_{10} *C*₁₁ C₁₀ C₁₁ C_{10}

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Challenges

Proposed Framework

- **Rate Balancing**: we propose a model 1) to find the optimal parameters for the configuration of a CNN accelerator.
- **Computational Graph** G =2) (V, E, ω, φ)
 - set of Vertices V, vertices weights ω,
 - edge set E, edges weight φ.
- **Multi-level Graph Partitioning** 3)





Herbert Wertheim College of Engineering

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Introduction CNN & FPGA Related Work Challenges Benchmark

Proposed Framework

- 1) Rate Balancing: we propose a model to find the optimal parameters for the configuration of a CNN accelerator.
- 2) Computational Graph $G = (V, E, \omega, \varphi)$
 - set of Vertices V, vertices weights ω,
 - edge set E, edges weight φ.
- 3) Multi-level Graph Partitioning
- 4) Coarse-grained Floor planning





Introduction CNN & FPGA Related Work Challenges Benchmark

FINN

- 1) FINN enables the design of heterogeneous custom streaming architecture for a given topology
- 2) Separate compute engines are dedicated to each layer, communicating via on-chip data streams.
- 3) It has two main units:
 - 1) The **Sliding Window Unit** (SWU): Supplies the MVU with the image by applying interleaving and implementing the im2col algorithm.
 - 2) The **matrix-vector unit** (MVU): input and output buffers and an array of PEs., each with several SIMD lanes.





Convolution Kernel: 3x3, IFM: 56 x56, C=64, OFM 56 x56



14

UF Herbert Wertheim College of Engineering UNIVERSITY of FLORIDA

Details of the proposed framework

0

1)

1)

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Rate Balancing $batency_i \simeq (1 + \epsilon) \times Latency_{i+1} \quad \forall i = 1, ..., N$ Inference model with N vertices and a platform with $Latency_i = F_i^n \times F_i^s$ with M SLRs, with $F_i^n = \frac{OFM_{H_i} \times OFM_{W_i}}{P_i}$, and $F_i^s = \frac{K_i^2 \times IFM_{Ch_i}}{S_i}$ For a layer $i \rightarrow$ Maximize (S_i, P_i) such that $P_{i} = \sum_{p=1}^{64} \sigma_{i,p} \times x_{i,p} \quad and \quad \sum_{p=1}^{64} \sigma_{i,p} = 1$ $S_{i} = \sum_{p=1}^{64} \gamma_{i,p} \times y_{i,p} \quad and \quad \sum_{p=1}^{64} \gamma_{i,p} = 1$ $= \sum_{p=1}^{2 \times \epsilon * 100} \delta_{i,p} \times z_{i,p} \quad and \quad \sum_{p=1}^{2 \times \epsilon * 100} \delta_{i,p} = 1$ Variables Constraints: $\square P_i \rightarrow \sigma_{i,n}$ $\bullet S_i \rightarrow \gamma_{i,p}$ $\bullet \varepsilon \rightarrow z_{i,p}$ **Resources Constraints:** $\begin{cases} \sum_{i=1}^{N} F_{lut_i}(P_i, S_i) \leq LUT_{VRs}, & \forall i = 1, ..., M\\ \sum_{i=1}^{N} F_{dsp_i}(P_i, S_i) \leq DSP_{VRs}, & \forall i = 1, ..., M\\ \sum_{i=1}^{N} F_{bram_i}(P_i, S_i) \leq BRAM_{VRs}, & \forall i = 1, ..., M \end{cases}$ • $F_{t_i}(P_i, S_i)$, a linear function that estimates the number of resources of type *t* demanded by the *ith* layer.

Herbert Wertheim College of Engineering

Partitioning

POWERING THE NEW ENGINEER TO TRANSFORM THE FUTURE

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Graph

- We implement a recursive balanced bipartitioning
- The weight of the heaviest partition $\leq \varepsilon \times$ $\omega(V)$
- **Stop**(k ≥ #F P Gas
- Considering bi-partitioning produce 2^n partitions per iterations:
 - inbalanced partitions
 - or too many partitions.



Herbert Wertheim College of Engineering

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

UF



Herbert Wertheim College of Engineering

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

POWERING THE NEW ENGINEER TO TRANSFORM THE FUTURE

Rate BalancingGraph
PartitioningMotivation
ExampleDesign
Generatio
nPlacementResults

Motivation Example



	Ċ.	9	LR.	2
	3			
	E.			
12.00	8			Т
136/0	8			
	R.			1
Bashin	Å.			
	λ,			
1 3	Ð	G	EQ.	P.
	R.		N.	۶.
16 MIN	2 N		5	25
Stat Jrd	-2		×	8
6	Ъ.		8	8
	ê		ê.	8
18. Li	2		Å.	SX.
58 A 10	12		5	E.
	×		×	×
ŝ	6	1	धुर	<u>اچ</u>
e	ê		e	ê
X	Ş.		X	5X
E	R		Ę	8
×	8		8	8
Ê –	8		8	5
ê	ê		ē	ē
Š.	8		3	S.

 Graph partitioning + Floor planning improves the frequency by 28%:

325 MHz to 416 MHz.



DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Rate Balancing Graph Motivation Example Design Generatio Results Results



They consider the number of cells in a design, their connections, and the target FPGA device's physical architecture to generate a circuit according to specified constraints







lead to overall QoR improvement in a design



DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Design Generatio

OCC

- It allows to implement and analyze (resource analysis, timing analysis, power analysis, etc) a module independently of the rest of the design.
- It enables reusing and preserving the characteristics of placed and routed modules within a top-level design.

2 **Floor planning**

- Utilizing pblock constraints allows carefully selecting the FPGA resources that will be used by each design component.
- Port planning with PartPins



Clock routing

• Accurately run the timing analysis on the OOC modules



• Once a module attains a desirable performance (Fmax, area, power, etc)

Checkpoint file generation



Herbert Wertheim College of Engineering

D FPGAs with *k* Dies

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

UF

Placement D_{H_i} h_i Wi $\bar{D}_{\{W_i\}}$ Ц S Ξ Modules must respect linear localization constraints: • Two adjacent M_i and M_j are non overlapping.

The localization must respect placement and congestion costs.



Herbert Wertheim College of Engineering

POWERING THE NEW ENGINEER TO TRANSFORM THE FUTURE

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

FROM MASTER SLIDE IF N/A

23



Herbert Wertheim College of Engineering

POWERING THE NEW ENGINEER TO TRANSFORM THE FUTURE

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Rate Bala				ation nple	Design Generatio n	Placement	Results
	Top Component	VOY VOX	22.114 24714 20 257143	Top Component	X2Y14 X2Y14 CmA	S. 02 Top Component	X2714 X2714 X3714 X3714
	Component ₁₋₁	XoY13	xm3 13	Component ₁₋₁	2 X0Y13 xm1 2 X2Y13 xm3 2 X4Y13	Component _{1.1}	2 X2Y13
(w_i, y_i)	Component ₁₋₂	11 X0/12 xmz	xmz xmz 11 X4Y12 1 X5V12	Component ₁₋₂	11 X0/11 2012 2012 2012 2012 2012 2012	Component ₁₋₂	1 X2Y1
	Component ₁₋₂		TASX 014	Component ₁₋₂	0710 X07 min min 2710 X27 min min xmii	Component ₁₋₂	7710 X077 710 X277 110 X277 1710 X477 1710 X477
	Component ₂₋₁	X 9/0X X 9/1X	x 19 x 19 x 19 x 19 x 10 x 10 x 10 x 10 x 10 x 10 x 10 x 10	Component _{2.1}	X0Y9 X X1Y9 x X2Y9 X X3Y9 X X3Y9 X	ि Component _{2.1}	019 XC 219 XC 219 XC 219 XC 419 XC
	Component _{2.2}	X0Y8 X1Y8 X1Y0	x3% X3% X4Y8 X5Y8	Component _{2.2}	X078 X178 X278 X376 X478	Component _{2.2}	X078 3 X178 2 X378 2 X378 2 X378 2 X578 2
	•	X0Y7 X1Y7 X1Y7	x3\7 X3\7 X4\7 X5\ <u>7</u>	•	x0Y7 x1Y7 x2Y7 x3Y7 x3Y7		хоү7 х1ү7 х3ү7 х4ү7 х5ү <u>7</u>
	•	X0Y6 X1Y6	x3Y6 X4Y6 X5Y <u>6</u>	•	 X0Y6 X1Y6 X2Y6 X3Y6 X4Y5 	5.sz	X0Y6 X1Y6 X3Y6 X3Y6 X3Y6
		13	2 2 2 5 5 5 5		175 375 375	273	75 75 75



Herbert Wertheim College of Engineering

POWERING THE NEW ENGINEER TO TRANSFORM THE FUTURE

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Placement h_i Wi (w_i, y_i) Component₁₋₂ 1. Create inter-connected nets uter1 pblock Component₁₋₂ 10 between components 2. Inter-component routing 9 25



- Experiment has been conducted on ResNet-50
- The hardware is generated using Vivado v2020.2 and RapidWright v2020.
- The components are implemented with vivado HLS.
- The MIP is solved with LocalSolver.



Herbert Wertheim College of Engineering

POWERING THE NEW ENGINEER TO TRANSFORM THE FUTURE

DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A



DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Rate Balancing Graph Motivation Example Generatio Placement Results

Granularity Exploration:





DEPARTMENT OR UNIT NAME. DELETE FROM MASTER SLIDE IF N/A

Rate Balancing Graph Partitionin		Design Generatio n		Re	sults
-------------------------------------	--	--------------------------	--	----	-------

	Layer-based ResNet			Block-based ResNet			Baseline ResNet		
	KFF	KLUTs	BRAM	KFF	KLUTs	BRAM	KFF	KLUTs	BRAM
Resources	741 ((↓ 26%)	421 (↓ 24%)	821.5 (≅)	801 (↓ 16%)	479 (↓ 10%)	761.6 (↓ 7%)	935	526	822
Latency (ms)	4.8 (↓ 31%)			4.2 (↓ 33 %)			6.3		
Frequency (MHz)	276 († 37%)			252 († 25.3%)			201		
Avg. Power (W)	208			225			235		
Energy Efficiently	29.79			27.98			22.15		

	Layer Granu	ResNet				
	Custom API	Inter-node Routing	Synthesis		P&R	
Time (hours)	1.14 4.82		4.16		8.9	
Ratio	~ 17.4%	~ 17.4% 82.6%		6%	64.4%	
Total (hours)	5.96 (2.1	13.02				
	Block					
	Custom API	Inter-node Ro				
Time (hours)	0.32	3.63				
Ratio	8.10%	91.8%				
Total (hours)	3.95	(3.21×↓)				



Conclusion

- We propose a framework to accelerate model inference on a multi-die Cloud FPGA Platform.
- It takes as input the model definition and performs an intensive search in the form of a MIP problem to determine each layer's highest degree of parallelism considering the platform constraints.
- The graph is then partitioned, and the resulting sub-graphs are allocated to the FPGAs' dies.
- Experiments and results show that our approach improves latency and maximum frequency, with little to no impact on the number of resources used.
- Future works we intend to extend to deeper neural networks.



UF Herbert Wertheim College of Engineering UNIVERSITY of FLORIDA

THANK YOU

0