Spatial Accelerator at the Edge

Tulika Mitra

Department of Computer Science National University of Singapore

NANDA Workshop, September 2024

Edge Intelligence



Real-time data processing and response locally on-device at the network's edge





Edge Intelligence Challenge





Short battery life

Power-hungry applications drain battery life very quickly.



Al power gap







Architecture-Compiler Codesign for < 1W Power

Microprocessor: von Neumann Model of Computation

- High programmability: Supports all application
- Inherently sequential: Burden is on hardware to extract parallelism
 - High energy overhead and limited parallelism
- 1980-2005: Unprecedented performance growth of microprocessors due to Moore's Law, Dennard Scaling, and micro-architectural innovations
- Hardware-driven parallelism extraction is too costly for edge devices

<pre>#include <stdlib.h> int sub(int x, int y){ return 2*x+y; } int main(int argc, char ** argv){ int a; a = atoi(argv[1]); return sub(argc,a); }</stdlib.h></pre>	.8ext:00000000 _sub: .8ext:0000001 .8ext:0000003 .8ext:0000000 .8ext:0000000 .8ext:0000000 .8ext:0000000 .8ext:00000010 .8ext:00000011 .8ext:00000014 .8ext:00000014 .8ext:0000001A .8ext:0000001A .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024 .8ext:00000024	push ebp mov ebp, esp mov eax, [ebp+8] mov ecx, [ebp+0Ch] lea eax, [ecx+eax*2] pop ebp retn push ebp mov ebp, esp push ecx mov eax, [ebp+0Ch] mov ecx, [ebp+0Ch] mov ecx, [eax+4] push ecx call dword ptr ds:imp_atoi add esp, 4 mov [ebp-4], eax mov edx, [ebp-4] push eax call _sub add esp, 8 mov esp, ebp pop ebp 88





AMD Ryzen 7000 ~65W, 3.6 GHz 12 cores, 24 threads



~0.8 mW @100MHz

Spatial Accelerators: Scaling Parallelism

Software-driven parallelism extraction





What can we do with 1W power?

At 100MHz, 0.9V, 45nm

- Compute only: ~43K MAC units
 4.3 TOPS/W, 8-bit integer
- Storage only: ~100MB, 100 access per cycle
- How much storage for data?
- How much storage for configurations?

Accelerator design choices are tradeoff between on-chip compute, data storage, and configuration storage

eger		FP		Memory	
ld		FAdd		Cache	(64bit)
8 bit	0.03pJ	16 bit	0.4pJ	8KB	10pJ
32 bit	0.1pJ	32 bit	0.9pJ	32KB	20pJ
ult		FMult		1MB	100pJ
8 bit	0.2pJ	16 bit	1.1pJ	DRAM	1.3-2.6nJ
32 bit	3.1pJ	32 bit	3.7pJ		

 25pJ
 6pJ
 Control
 70 pJ

 ↑

 ↑

 ↑

 ↑

 I-Cache Access
 Register File Access
 Add

Instruction Energy Breakdown

Mark Horowitz. ISSCC 2014



Custom Edge AI Accelerators



Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6	Layer 7	Layer 8
							-
							×.



35 TOPS at 4 TOPS/W

Processo

1

Mythic:

Domain-specific Accelerator: Close to zero configuration cost Very little storage cost



Hailo 26 TOPS at 3 TOPS/W **Programmability:** Introduces configuration storage and access cost



Accumulator Activation 167 GiB/s Normalize / Pool

Edge TPU: 4 TOPS at 2 TOPS/W



I FIFO Fetcher

CGRA Spatial Accelerator: Efficiency with Programmability

- Parallelism: Array of simple processing elements (PE)
- Operation of a PE and the data transfer between PEs can be configured per cycle
- Exposed to software: Compiler sets the configuration according to the dataflow of the application
- Simple hardware and hence low power: During execution, PE array and switches simply follow the pre-defined configurations
- As number of configurations is limited, CGRA can only support recurring computation, i.e., loops



Dataflow Processor

A Preliminary Architecture for a Basic Data-Flow Processor"

- Expose instruction-level parallelism in dataflow graph at compile time
- Dataflow graph is directly mapped onto hardware
- Operation is triggered when input data is ready
- Dataflow computing dominates most accelerator designs today

Jack B. Dennis and David P. Misunas Project MAC Massachusetts Institute of Technology





Dataflow synthesis on CGRA

Mapping the dataflow graph (DFG) of the loop body

- **Placement:** assign DFG operations to PEs
- **Routing:** map data dependencies through links between PEs rather than through cache or scratchpad

Image: space of the loop bodyState of the loop bodyState of the loop body

N2 N3 N4 DFG

F1 \leftrightarrow F2 f fF3 \leftrightarrow F4

Placement and routing



Spatio-temporal mapping

2x2 CGRA



CGRA Spatial Accelerator: Efficiency + Programmability

Instantiate different task-specific virtual or logical accelerators on the same hardware fabric through software-directed reconfiguration

Data

PE



Google TPU: Specialized

Hardware architecture is designed specially to support matrix multiplication dataflow

CGRA: Generalized

Hardware architecture is general and reconfigured to support different dataflows



CGRA: Mapping Diverse Dataflows



NUS National University of Singapore

CGRA for Sustainability

Embodied footprint: Raw material extraction, manufacturing, assembly, transportation, end-of-life processing

Operational footprint: Device use during its entire lifetime



Embodied footprint per unit chip area is increasing with technology nodes

Embodied footprint dominated most consumer computing devices



Concurrency

Carbon Footprint Improvement



HyCUBE: a 0.9V 26.4 MOPS/mW, 290 pJ/cycle, Power Efficient Accelerator for IoT Applications ASSCC'19

HyCUBE : A CGRA with Reconfigurable Single-cycle Multi-hop Interconnect_DAC 2017

- Reconfigurable single-cycle •
- multi-hop interconnect
- Creates large dynamic neighborhood reachable within single-cycle through asynchronous bypassing of intermediate PEs



HyCUBE CGRA @ NUS



F3

N5

N2

N7

N3

Ν4

Ν6



CGRA Architectural Innovations





streaming

application

CASCADE **Decoupled Access-Execute CGRA** 3x speedup, 2.2x performance/watt over iso-area CGRA

FLEX: Spatio-temporal vector dataflow 45% less energy. 1.9× power efficiency at same performance over CGRA

ICED: **DVFS-aware CGRA** 1.32x energy-efficiency over CGRA at same performance

Pipeline

16

LDO

6 x 7

DCO ADPLL

6-0-0

CLK ISLAND

vnc Bypass FIFO

CASCADE: High Throughput Data Streaming via Decoupled Access/Execute CGRA. CASES 2019 FLEX : Introducing FLEXible Execution on CGRA with Spatio-Temporal Vector Dataflow. ICCAD 2023 ICED: An Integrated CGRA Framework Enabling DVFS-Aware Acceleration. MICRO 2024

PACE CGRA+RISC-V SoC: 2023





SDRAN

controller

Onchip SRAM

...........

Memory

Fabric

PACE SoC Specifications (Simulation)

Tech. Node

UMC 40nm ULP

Micron

SDRAM

Chip (32MB)

PACE: A Scalable and Energy Efficient CGRA in a RISC-V SoC for Edge Computing Applications. Hot Chips 2024

CGRA Comparison

PACE has highest power efficiency even at 40nm

	Academic					Commercial			
	Post synthesis simulation/Post P&R			Silicon		Post synthesis simulation/Post P&R	Silicon		
CGRA	UE_CGRA HPCA'21	SNAFU ISCA'21	TRANSPIRE DATE'20	RIPTIDE MICRO'22	AMBER <u>VLSI'22</u>	PACE	ULP-SRP <u>FPT'12</u>	Renesas DRP <u>VLSI'18</u>	Cardinal SN10 (Improved version of Plasticine for ML training) <u>ISSCC'22</u>
PE array	8x8	бхб		бхб	384 PEs	8x8	3x3	96 PEs	640 PMU 640 PCU
Node	28nm TSMC	Intel 22FFL	28nm FDSOI	Intel 22FFL	16nm FINFET	40nm UMC	40nm TSMC	28 nm	7nm TSMC
Voltage (V)	-	-	0.6	_	1.29	0.45	1.1	-	_
Frequency (MHz)	750	50	50	50	955	10	100	333 MHz	-
Area (mm ²)	0.25^{1}	0.27^{1}	0.27	0.25^{1}	20.1	3.94	_	_	-
Performance (MOPS)	625	71	-	62	367000	640	467 ²	960 gigaflops	2.6 petaflops
Power (mW)	14	0.54	-	0.24	682	1.1	11.3	-	-
Power efficiency (GOPS/W)	45	134	224	254	538	582	23	143 ³	_



¹ Only the PE array included

² Estimated from published 467 MIPS

³Estimated value from experiments

Compilation: Dataflow Synthesis on CGRA

- Dataflow synthesis on CGRA provides a fertile ground for research
- Graph Minor (2012): Theoretically optimal based on graph minor modeling
- ChordMap (2020): Streaming dataflow application synthesis
- HiMap (2021), Panorama (2022): Scalable mapping through hierarchical abstractions
- LISA (2022): Automated compiler synthesis using GNN



Number of sub-CGRA:4

5200 cycles

5300 cycles

2 5400 cycles

5300 cvcle

19

Graph Minor Approach for Application Mapping on CGRAs. FPT 2012 <u>TRETS 2014</u> Best paper award
ChordMap: Automated Mapping of Streaming Applications onto CGRA. <u>TCAD 2021</u>
HiMap: Fast and Scalable High-Quality Mapping on CGRA via Hierarchical Abstraction. <u>DATE 2021</u> <u>TCAD 2022</u>
Panorama: Divide-and-Conquer Approach for Mapping Complex Loop Kernels on CGRA. <u>DAC 2022</u> Publicity Paper
LISA: Graph Neural Network based Portable Mapping on Spatial Accelerators. <u>HPCA 2022</u> Distinguished Artifact

Scaling Dataflow Synthesis on CGRA

- Multiple abstractions to divide and conquer
- Clustering at both dataflow graph and architecture level
- Dramatic Improvement
 - 17x performance, 5x energy-efficiency
 - Reduce compilation time from days to15-min for 64x64 CGRA









Automated GNN-based Mapping



Instantiating TPU-like dataflow for GEMM



Google TPU Matrix Multiply Unit (MXU) Dataflow





HiMap GEMM Schedule & Dataflow on CGRA



NUS Morpher Open-Source CGRA Toolchain

14 3 10 3

return 01



23

https://github.com/ecolab-nus/morpher

NUS MLIR-Based End-to-End CGRA Compiler

CGRA SoC Emulation on FPGA

Domain-Specific CGRA

- Task-specific CGRA design
- RVAMP: Automated exploration of heterogeneous CGRA (compute, memory, network) design space for an application domain
- Architectural specification and compiler generation

REVAMP: A Systematic Framework for Heterogeneous CGRA Realization. ASPLOS 2022 Toolchain: <u>https://zenodo.org/record/5848404#.YgyrPTFByUk</u>

Emergence of Native Multimodal LLMs Unifying AI across modalities

set were on the desk near a red

- > Seamless interaction with physical world via integration of multimodal data
- Process and generate text, images, audio, and video within a single model
- Support complex tasks requiring multi-modal reasoning
- Perfect opportunity for Edge AI

Mind the Gap: From Smaller to True Edge AI

Size in billions of parameters (active)

MatMult is all you need

Attention is the cornerstone of LLMs enabling them to intelligently focus on the most relevant parts of any input

LLM inference requires huge memory to store the model weights and key-value cache (e.g., 3 GB for the smallest Phi-3 Mini model: 3.8B parameters, 4 bit quantized, 4K context length)

states

Vaswani, Ashish. "Attention is all you need." NeurIPS 2017 Image credit: Damien Benveniste

A Diverse Landscape of Sparse Attentions

Attention computation scales quadratically with sequence length, requiring O(n²) operations for n tokens

- Advanced attention mechanisms leverage diverse **sparse matrix computations** with specialized patterns for enhanced efficiency
- Each type of attention mechanism requires **unique dataflow**

Memory transfer requires 200x more energy than compute

Quantized Attention: Trim Bits, Enable Edge AI

4-bit Mini-Floats provide comparable accuracy to
 32-bit floating point in transformer models

Complete on-chip storage (BRAM + URAM) of weights on U280 FPGA MatMult-Free Ternary LLM (1.58bits): *Work-in-progress* 370M parameters: ~5K tokens/sec @ 21W vs. ~100 tokens/sec @ 200W

Shedding the Bits: Pushing the Boundaries of Quantization with Minifloats on FPGAs. FPL 2024 * Scalable MatMul-free Language Modeling, Rui-Jie Zhu et al. arXiv 2024

FPGA Acceleration of
Sliding-Window AttentionASIC
Spars

ASIC accelerator for Sparse Matrix ML Workload

SWAT: FPGA Sliding-Window Attention Design achieves 15x energy-efficiency over GPU

ZeD: ASIC accelerator achieves 3.2x performance-per-area improvement over SOTA for sparse ML workloads

SWAT: Scalable and Efficient Window Attention-based Transformers Acceleration on FPGAs. DAC 2024 ZeD: A Generalized Accelerator for Variably Sparse Matrix Computations in ML. PACT 2024

Challenge: CGRA for Arbitrary Sparse Dataflows

- Minimalist reconfigurable array that adapts to diverse attentions without high area, energy, speed costs
- Smart algorithm that maps arbitrary attention dataflows on the fabric
 - Maintain high parallel processing while handling complex dataflow dependencies
 - Achieve software-level accuracy with lowprecision hardware technology

Acknowledgments

• Li-Shiuan Peh

- Anh Tuan Do
- Chen Liang
- Tan Cheng
- Manupa Karunaratne
- Wang Bo
- Aditi Kulkarni
- Dhananjaya Wijerathne
- Zhaoying Li
- Thilini Kaushalya Bandara
- V Vanchinathan
- Dan Wu
- Vishnu Paramasivam
- Chong Yi Sheng
- Pranav Dangi
- Zhenyu Bai
- Huize Li

NATIONAL Research Foundation

PRIME MINISTER'S OFFICE SINGAPORE

Ministry of Education SINGAPORE

Institute of Microelectronics

BLACK

AMD XILINX

THANK YOU

CGRA versus FPGA

- CGRA offers word-level reconfiguration while FPGA offers bit-level reconfiguration
 - Reduces reconfiguration overhead providing better performance and lower power
- Built-in arithmetic functions improves energy-efficiency/frequency, makes mapping simpler
 - Closer to compilation for VLIW architecture than HLS for ASIC/FPGA
 - Easier programmability for software developers compared to FPGA
- CGRA offers spatio-temporal mapping while FPGA offers only spatial mapping
 - Reduces area cost, which is important for edge devices

			Normalized w.r.t. FPGA				
	CGRA	Application	Power	Perf.	Power Efficiency		
ne		CNN	1.20	95.1	76.9		
Tir		GEMM	1.40	33.0	24.4		
		BFS	0.60	7.3	11.4		
	Amber (VLSI'22)	Resnet Conv2_x	0.11	1.0	9.2		
		Resnet Conv3_1	0.11	1.0	9.1		

Spatial Mapping

Spatio-temporal Mapping

36

CGRA vs FPGA