Mozart: Designing for Software Maturity and the Next Paradigm for Chip Architectures

Karu Sankaralingam | UW-Madison and SimpleMachines Tony Nowatzki | UCLA Greg Wright, Poly Palamuttam, Jitu Khare, Vinay Gangadhar, Preyas Shah | SimpleMachines

Where does AI hardware/software stand today?

- 1. The computational diversity needed to support AI is increasing
- 2. The software user experience expectations is increasing
- 3. GPU software maturity* is unrivalled in completeness and hence allows near complete dominance among AI industry deployment and researchers.
- 4. This support for model diversity is fuelling these trends and increasing GPU adoption!



* NVIDIA DL stack - cuDNN, TensorRT, etc.

Is there a problem?

- 1. GPU chip utilization is quite low (10% of peak)
- 2. The software stack is turgidly tied to hand-tuned/auto-tuned libraries, resulting in ML systems getting caught in the rut.
- 3. Tail wagging the dog: Because only some styles of kernels achieve high performance on GPUs, ML practitioners are forced to create ridiculously large models to get new capability, rather than use better algorithms.
- 4. Simpler models exist but don't run fast on existing HW or SW
 - a. Lottery Hypothesis shows simpler models have the same/superior learning ability.
 - b. Depth-wise separable convolutions have 10X fewer ops, but provide no speedups on today's HW.

New chip architectures must have software maturity to challenge GPU dominance and serve as an alternative

Designing for Software Maturity

Goal: automated realization of the compiler and software stack utilizing programsynthesis vs *human coding of high-perf libraries or non-automated/imperative compilers*.

Benefits:

- 1. Allows day-0 SW maturity
- 2. Avoiding the pit-fall of building software after silicon arrives
- 3. Overcome the problem of SW dev cycle that far exceeds the HW dev cycle and cannot match the rate of application change
- 4. Allows co-design of the architecture to match needs of software, because software stack exposes what applications are doing

Other IRs and Compiler:

- 1. TVM, GLOW and other IRs solve part of the problem
- 2. Don't get sufficiently down to object-code or sufficiently up the DLstack to serve as turn-key deployment

Composable Computing: An Architecture Paradigm That Allows SW Maturity

Overview of CPU Execution



Source code has rich semantic information

Machine code



Compute resources devoted to overhead instead of app perf.

Composable Computing Paradigm



Architecture Overview and Basic Tile



SoC Organization



Architecture Specification: Performance View



Mozart Chip, Board, System





Technology Node	TSMC 16FFC
Operating Frequency	1 GHz
Die Area	404 mm ²
Package	45mmx45mm FCBGA
Peak INT8 performance	48 TOPS
DRAM Channels	2x HBM2
DRAM Bandwidth	512 GB/s
Host Interface	PCle Gen3 x16

Deep Learning Performance

	Mozart @1GHz, 16nm, 75W PCIe Card	Bach (7nm) (projection)
Resnet	2.6X	6.6X
BERT	1.2X	3.1X
SSD-Resnet	1.0X	2.7X
RNN-T*	5.8X	15.1X
DLRM	0.3X	1.8X

* RNN-T includes many subtelities and optimizatin opportunities on batch-size in a sample

Relative Inf/Sec/Watt to NVIDIA A100

HW/SW Stack



High-level Attributes

Programmability Algorithm independence

Physical attributes

Physical scalability Modularity Efficiency

Software

C/C++ and Python SDK Resnet, BERT, RNN-T, DLRM Universal AI/ML accelerator

Program Synthesis and SW Maturity in Mozart

Compiler Flow



DFG MATMUL

quant6x6x64(N0[6], W0[6], N1[6], W1[6], reset, Z -> out0, out1) {

reset_counter = Acc64(0x04000000, reset, reset); reset_delayed = CmpEQ32x2(reset_counter, 0xfffffff04000001);

NW000 = QuantOp(N0[0], W0[0], Z);NW001 = QuantOp(N0[1], W0[1], Z);NW002 = QuantOp(N0[2], W0[2], Z);NW003 = QuantOp(N0[3], W0[3], Z); NW004 = QuantOp(N0[4], W0[4], Z);NW005 = QuantOp(N0[5], W0[5], Z);AS0 = SAdd32x2(NW000, NW001);AS1 = SAdd32x2(NW002, NW003);AS2 = SAdd32x2(NW004, NW005);AT0 = SAdd32x2(AS0, AS1);AU0 = SRedAcc32x2(AT0, AS2, reset_delayed); NW100 = QuantOp(N0[0], W1[0], Z);NW101 = QuantOp(N0[1], W1[1], Z); NW102 = QuantOp(N0[2], W1[2], Z); NW103 = QuantOp(N0[3], W1[3], Z);NW104 = QuantOp(N0[4], W1[4], Z);NW105 = QuantOp(N0[5], W1[5], Z);AS3 = SAdd32x2(NW100, NW101);AS4 = SAdd32x2(NW102, NW103);AS5 = SAdd32x2(NW104, NW105);AT1 = SAdd32x2(AS3, AS4);AU1 = SRedAcc32x2(AT1, AS5, reset delayed);

out0 = Concat32(AU0, AU1);



Software Capability

Complete deep-learning compiler including model characterization, graph optimization (fusion, splitting, rewriting), quantization (including quantization on non-linear operators like Gelu, Softmax, Erf), tensor-mapping to memory, backend optimized object code generation, and lightweight high-speed custom software runtime.

Application	Original trained FP32 Optimized & quantized model Mozart model		trained FP32 Optimized & quantized model Mozart model			Resne	et — BE	RT —	SSD-Resr	net
	Total Ops	#uniq Ops	Total Ops	#uniq Ops	80% 60%	66%				
Resnet50	461	14	84	12	40%	51%	40% 33%			
SSD-Resnet34	3320	46	111	10	20% 0%		6%		3%	
BERT	757	24	322	17		1	2	3	4	
SSD-Mobilenet	4104	52	91	7	Appl	Application		Compiler runtime to transform FP32 model		
3D-Unet	587	19	100	15	Res	Resnet50		8 seconds		
					SSD-R	SSD-Resnet34		19 seconds		
DLRM	47	10	19	11	В	BERT		29 seco	onds	
						RM		3 seco	nds	

Day-0 Software Maturity is Necessary for all Future Chips

- 1. Architecture **definition** is the key challenge for future chips and new architectures **must allow day-0 SW maturity**
 - a. Can you bootstrap software for that architecture (before/while) defining it?
 - b. Can you define an architecture that is correct and meets some figures-of-merit requirements?
- 2. Chip implementation of a defined architecture is deterministic
 - a. But includes about 8 months of exposure: 4 months for "final" physical design, 4 months manufacture, 1.5 months bringup
- 3. Program Synthesis and Auto-generation of SW stacks is necessary for all future chips
- 4. Mozart and Composable Computing Paradigm is a compelling path forward

The Key Research Ideas

1) Spatial compiler that was retargetable to "any" spatial hardware. PLDI 2013 paper. Distinguished Paper award, CACM Research Highlights Nomination. <u>"A General Constraint-centric Scheduling Framework for Spatial Architectures"</u>

2) Five behaviors that capture hardware and software interactions. HPCA 2016 paper. IEEE Micro Top Picks. <u>"Pushing the Limits of Accelerator Efficiency While Retaining General-Purpose Programmability"</u>

3) Hybrid Von-Neuman Dataflow and Stream-Dataflow: Practical blend of above two ideas with a concrete arch, microarchitecture and small enough design. ISCA-2015, ISCA-2017. Top Picks, CACM Research Highlights.

"Heterogeneous Von Neumann/Dataflow Microprocessors"

"Stream-Dataflow Acceleration"

Thank You