CS 744: BIG DATA SYSTEMS

Shivaram Venkataraman Fall 2018

ADMINISTRIVIA

- Final class on Dec 6th, No class on Dec 11th
- Poster session Dec 13th, details on Piazza
- Course Project final report
- Course feedback form

ACCELERATORS

MOTIVATION

Capacity demands on datacenters New workloads

Metrics

Total cost of ownership (Depends on price ?) Power / operation Performance / operation

Goal: Improve cost-performance by 10x over GPUs

WORKLOAD: ML INFERENCE

Name	LOC	Layers					Nonlinear	Weights	TPU Ops /	TPU Batch	% of Deployed TPUs
		FC	Conv	Vector	Pool	Total	function	weights	Weight Byte	Size	in July 2016
MLP0	100	5				5	ReLU	20M	200	200	610/
MLP1	1000	4				4	ReLU	5M	168	168	01%
LSTM0	1000	24		34		58	sigmoid, tanh	JZIVI	64	64	2007
LSTM1	1500	37		19		56	sigmoid, tanh	34M	96	96	29%
CNN0	1000		16			16	ReLU	8M	2888	8	50%
CNN1	1000	4	72		13	89	ReLU	100M	1750	32	5%

DNN: RankBrain, LSTM: subset of GNM Translate CNNs: Inception, DeepMind AlphaGo

What are some notable points from this table ?

WORKLOAD: ML INFERENCE

Quantization \rightarrow Lower precision, energy use

8-bit integer multiplies (unlike training), 6X less energy and 6X less area

Need for predictable latency and not throughput e.g., 7ms at 99th percentile

OUTLINE

Motivation

TPU Design Comparison with CPU, GPU Selected Lessons ?

TPU DESIGN









INSTRUCTIONS

CISC format (why ?)

- I. Read_Host_Memory
- 2. Read_Weights
- 3. MatrixMultiply/Convolve
- 4. Activate
- 5. Write_Host_Memory

SYSTOLIC EXECUTION

Problem: Reading a large SRAM uses much more power than arithmetic!

Wave-like propagation of data

Not Von Neumann (fixed compute)

No memory access / cache / bus etc.



PERFORMANCE: ROOFLINE MODEL



COMPARISON WITH CPU, GPU

	Die									
Model	mm^2	nm	MHz	TDP	Measured		TOPS/s		CP/a	On-Chip
					Idle	Busy	8b	FP	GD/S	Memory
Haswell E5-2699 v3	662	22	2300	145W	41W	145W	2.6	1.3	51	51 MiB
NVIDIA K80 (2 dies/card)	561	28	560	150W	25W	98W		2.8	160	8 MiB
TPU	<331*	28	700	75W	28W	40W	92		34	28 MiB

COMPARISON WITH GPU, CPU

Туре	Batch	99th% Response	Inf/s (IPS)	% Max IPS
CPU	16	7.2 ms	5,482	42%
CPU	64	21.3 ms	13,194	100%
GPU	16	6.7 ms	13,461	37%
GPU	64	8.3 ms	36,465	100%
TPU	200	7.0 ms	225,000	80%
TPU	250	10.0 ms	280,000	100%

SELECTED LESSONS

- Latency more important than throughput for inference
- LSTMs and MLPs are more common than CNNs
- Performance counters are helpful
- Remember architecture history

SUMMARY

New workloads \rightarrow new hardware requirements

Domain specific design (understand workloads!)

- No features to improve the average case
- No caches, branch prediction, out-of-order execution etc.
- Simple design with MACs, Unified Buffer gives efficiency

Drawbacks

- No sparse support, training support (TPU v2, v3)
- Vendor specific ?