

CS 744: TPU

Shivaram Venkataraman Fall 2020

ADMINISTRIVIA

Midterm 2, Dec 3rd

- Papers from SCOPE to TPU
- Similar format etc.

Presentations Dec 8, 10

- Sign up sheet
- Presentation template
- Trial run?

MOTIVATION



```
Metrics

Power/operation

Performance/operation \longrightarrow hatency (or Tput) \rightarrow workloads are latency sensitive

Total cost of ownership

\longrightarrow Buy (Build) + Operate
```

Goal: Improve cost-performance by 10x over GPUs

© CNNS are only 5%, MCRS are 61%. () Number of weights is not borrelated with batch size and												
	CNNS Dþs /	hare 1. Byte	very	high				A M	op= -P & LSTM Ops/byte	tare and h	same batch size	
Name	LOC	Layers					Nonlinear	Weights	TPU Ops /	TPU Batch	% of Deployed TPUs	
		FC	Conv	Vector	Pool	Total	function	weignis	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	52M	64 🧲	→ 64	2007	
LSTM1	1500	37		19		56	sigmoid, tanh	34M	<u>96</u>	<i>→</i> 96	29%	
CNN0	1000		16			16	ReLU	8M	2888	8	50%	
CNN1	1000	4	72 ·		13	89	ReLU	100M	1750	32		

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

WORKLOAD: MI INFERNCE

Quantization \rightarrow Lower precision, energy use 32 bit float \rightarrow 8 bit integer

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

caches -> improve average branch prediction case scenario

- Focus on Inference Only!



INSTRUCTIONS

CISC format (why ?)

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

-> specialized instruction set -> CISC Ly Instructions encode operations that take many cycles to run

SYSTOLIC EXECUTION

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

HASWELL ROOFLINE

(memory officing										
Mem Compite Carbes, L1, L2, L3	L3 : 32 L2 : 16 = L1 : 8	MВ ~ 60	AKI 4 MB	SUN max T	WIIF power	HCPL	J, GP GPU H band	U as a lindth	higher	mem
			Die							
Model	mm ²	nm	MHz	ם תד ס תד	Меа	Measured		TOPS/s		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 v1	<331*	28	700	75W	28W	4 <u>0</u> W	92		34	28 MiB
Empared CPU and GPU				- GPUs bring down Power used when idle - TPUs not so much						

SELECTED LESSONS

- Latency more important than throughput for inference
- LSTMs and MLPs are more common than CNNs
- Performance counters are helpful -> to also improve compilers of DNN models
- 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 ?

DISCUSSION

https://forms.gle/tss99VSCMeMjZx7P6

O harger batches have higher tput but also higher tail latercy Inferences per sec

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	<u>6.7 ms</u>	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%

much higher 2 TPU Hputs are (3) TPU can be at higher util while meeting Ims latency target ④ GPU has higher IPS at some botch size compared to CPU
→ relates to arg. latency

How would TPUs impact serving frameworks like Clipper? Discuss what specific effects it could have on distributed serving systems architecture

() TPUS have 8GB to store many models La but this might break containers in Clipper (2) Mragglers are less frequent (3) Batching (Auto batching) an be very helpful (4)

NEXT STEPS

No class Thursday! Happy Thanksgiving!

Next week schedule: Tue: Fairness in ML, Summary Thu: Midterm 2

ENERGY PROPORTIONALITY

Haswell

Target Workload