

### CS 744: TPU

Shivaram Venkataraman Fall 2021

#### ADMINISTRIVIA

Midterm 2. Dec 7<sup>rd</sup>

- Papers from SCOPE to TPU ->
- Similar format etc. -> questions more specified
  Details on Piazza writing your assumptions etc.

Poster session: Dec 14th

More details soon

# MOTIVATION



Goal: Improve cost-performance by 10x over GPUs

For every lyte of data (model) WORKLOAD - ML inference											
		0		(	í Á			7 ×	maller	8	mixture of models ith low?
				Iguang			Nonlingan	$\bigwedge$	TDUOng	TDUPatch	of CNNS
Name	LOC	FC	Conv	Vector	Pool	Total	function	Weights	Weight Byte	Size	in July 2016
MLP0	100	5				5	ReLU	20M	200	200	
MLP1	1000	4				4	ReLU	5M	168	168	61%
LSTM0	1000	24		34		58	sigmoid, tanh	52M	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	570
	Jarge								L		

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

#### WORKLOAD: MI INFERNCE

Quantization -> Lower precision, energy use -> model weights could be Converted to integers

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

Particulary important for serving - not just los ang latency - but we want predictability!







# INSTRUCTIONS

CISC format (why ?)

- I. Read\_Host\_Memory ~ input
- 2. Read\_Weights ~ mold
- 3. MatrixMultiply/Convolve
- 4. Activate
- 5. Write\_Host\_Memory

- Complex Instruction Set La PCie bus not that high La Each instruction can take a long time! -> maps very closely to the workload

#### SYSTOLIC EXECUTION

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







### COMPARISON WITH CPU, GPU

	Die									
Model	2	nm	MHz	TDP	Measured		TOPS/s		$CD/\pi$	On-Chip
	mm				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
idle power is nuch lower power performance ra							. <del>0</del> -			

### SELECTED LESSONS



- Latency more important than throughput for inference
- LSTMs and MLPs are more common than CNNs
- Performance counters are helpful  $\rightarrow$  for benchmarking 1 debugging
- 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/LFeaeME4pFdHZdMV6

For all he	ardwarc							
incr batch	tize		GPU CPU sacrifice utilization					
> incr		7	The reach 7ms SLD					
util &	Туре	Batch	99th% Response	Inf/s (IPS)	% Max IPS			
tput	CPU	16	7.2 ms	5,482	42% 100%			
but at	CPU	64	21.3 ms	13,194				
cost of	GPU	16	<b>6.7</b> ms	13,461	37% 100%			
latercy	GPU	64	8.3 ms	36,465				
	TPU	(200)	7.0 ms	225,000	80%			
	TPU 250		10.0 ms	280,000	100%			
	C	an ha with la batch	alle rger mize	J t Able to meet 1 ms while 72	ook Hput			

How would TPUs impact serving frameworks like Nexus? What specific effects it could have on distributed serving systems architecture

#### **NEXT STEPS**

Next week schedule

Tue: Midterm 2

Thu: Last class! (Fairness in ML, Summary)