

welcome back!!

CS 744: TPU

Shivaram Venkataraman

Spring 2024

ADMINISTRIVIA

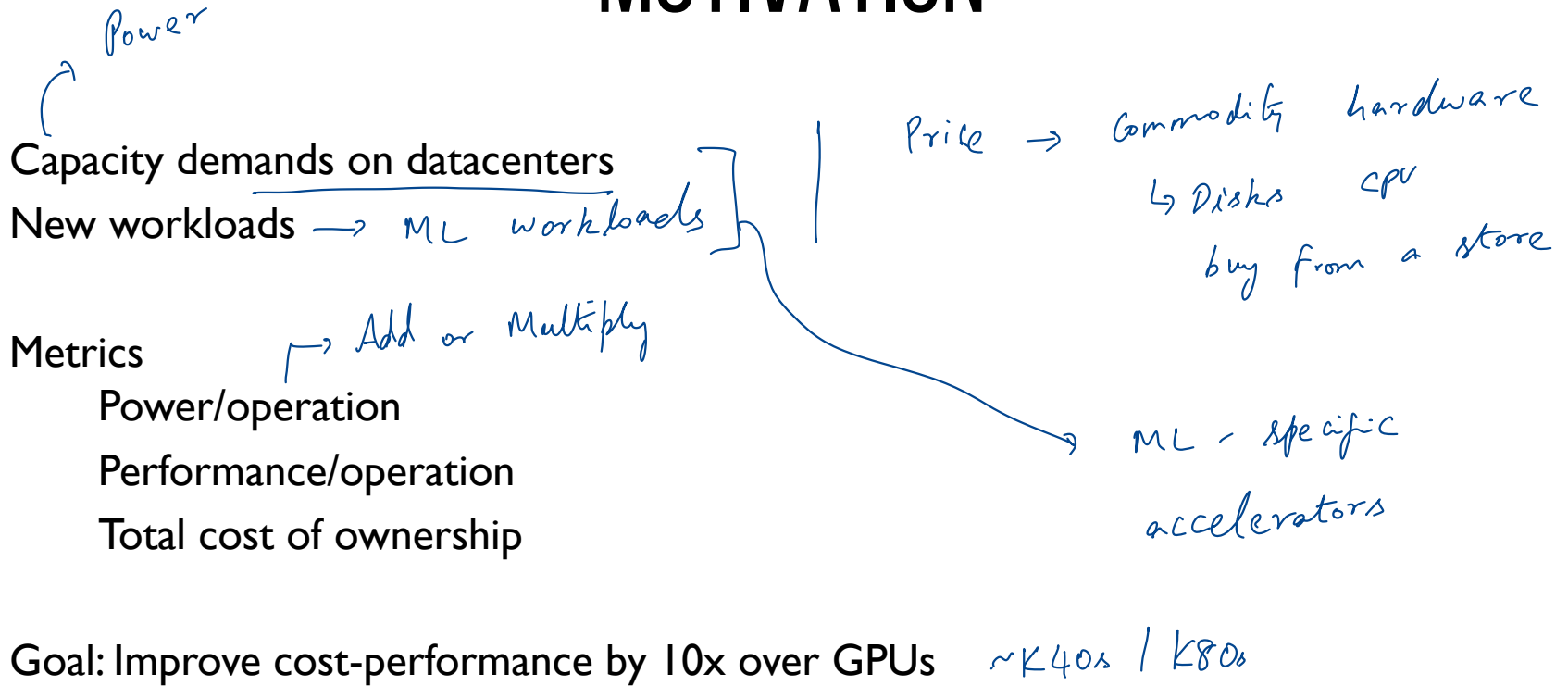
Midterm 2, April 25th

- Papers from SCOPE to HeMem *~ similar count*
- Similar format as first midterm
- Details on Piazza → *sample mid terms*

Poster session: May 2nd

- More details soon

MOTIVATION



WORKLOAD

→ ML inference workloads

handful of layers of compute kernels

how much compute for a byte of data (model weights)

model size

90% not CNN models

Name	LOC	Layers					Nonlinear function	Weights	TPU Ops / Weight Byte	TPU Batch Size	% of Deployed TPUs in July 2016
		FC	Conv	Vector	Pool	Total					
MLP0	100	5				5	ReLU	20M	200	200	61%
MLP1	1000	4				4	ReLU	5M	168	168	
LSTM0	1000	24		34		58	sigmoid, tanh	52M	64	64	29%
LSTM1	1500	37		19		56	sigmoid, tanh	34M	96	96	
CNN0	1000		16			16	ReLU	8M	2888	8	5%
CNN1	1000	4	72		13	89	ReLU	100M	1750	32	

larger → LLMs ~ GBs

DNN: RankBrain, LSTM: subset of GNM Translate

CNNs: Inception, DeepMind AlphaGo

WORKLOAD: ML INFERENCE

→ representation of weights to integers
Quantization → Lower precision, energy use

→ Goal: Perf/watt

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

InFaaS → SLO satisfaction rate

→ hardware plays a key role

TPU DESIGN CONTROL

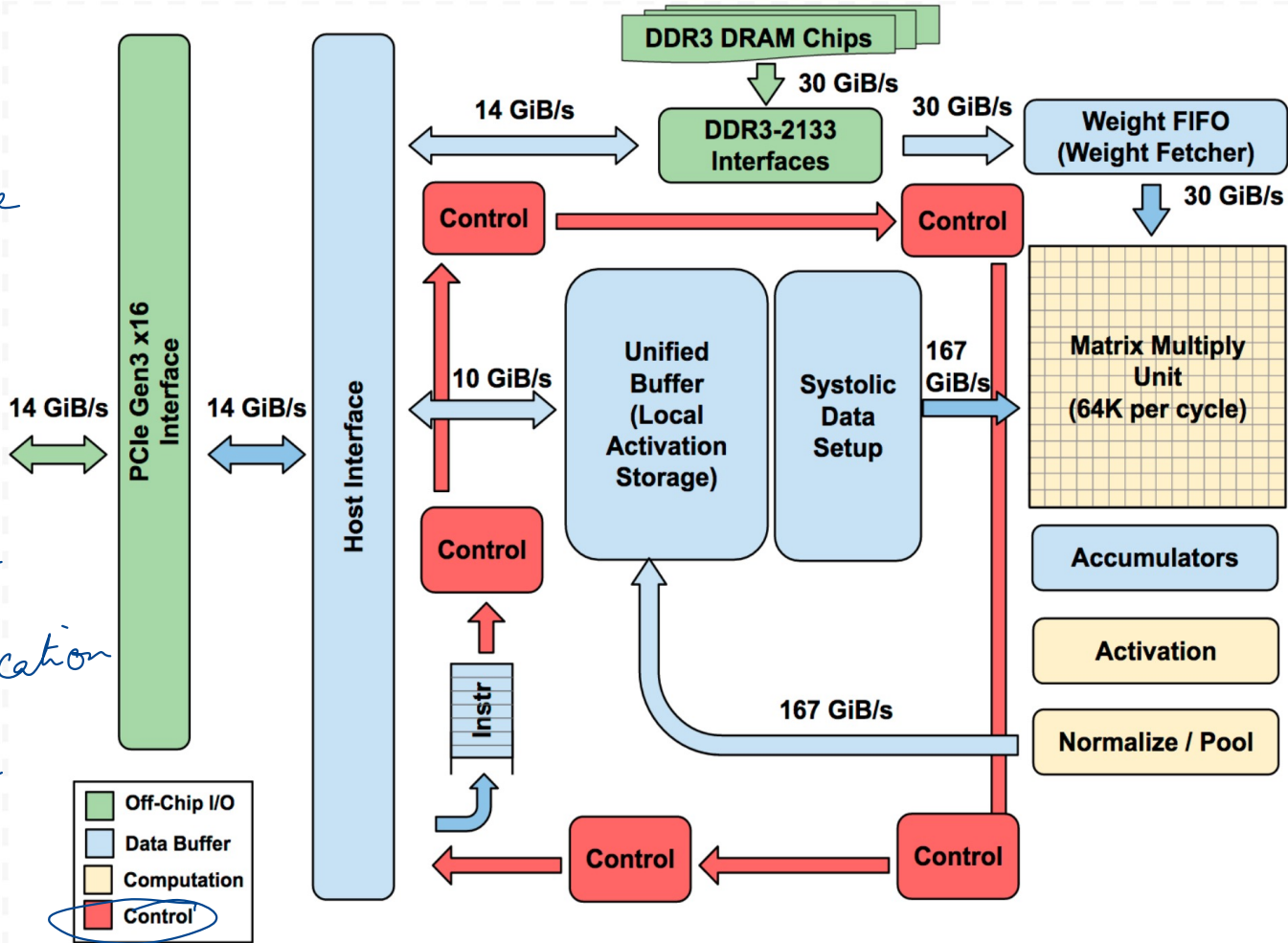
Mounted on a PCI-e

→ Compatibility

Low bandwidth when using PCI-e

↳ minimize host to device communication

↳ Instructions are coarse grained



COMPUTE

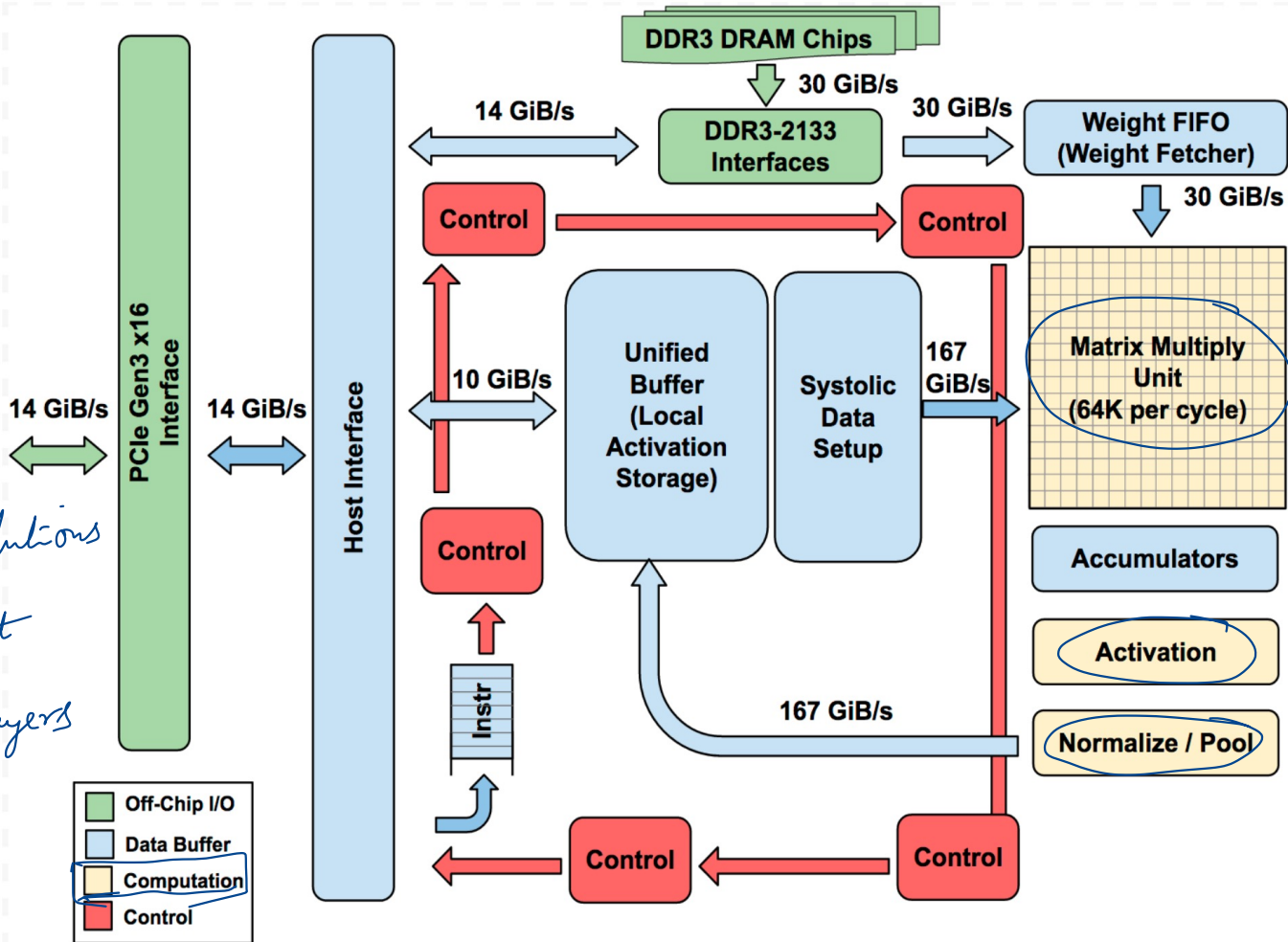
- ML specific
→ Activations / Pooling etc.

↳ efficient

→ Matrix Multiply

- MLPs & Convolutions

- one unit that handles these layers



DATA

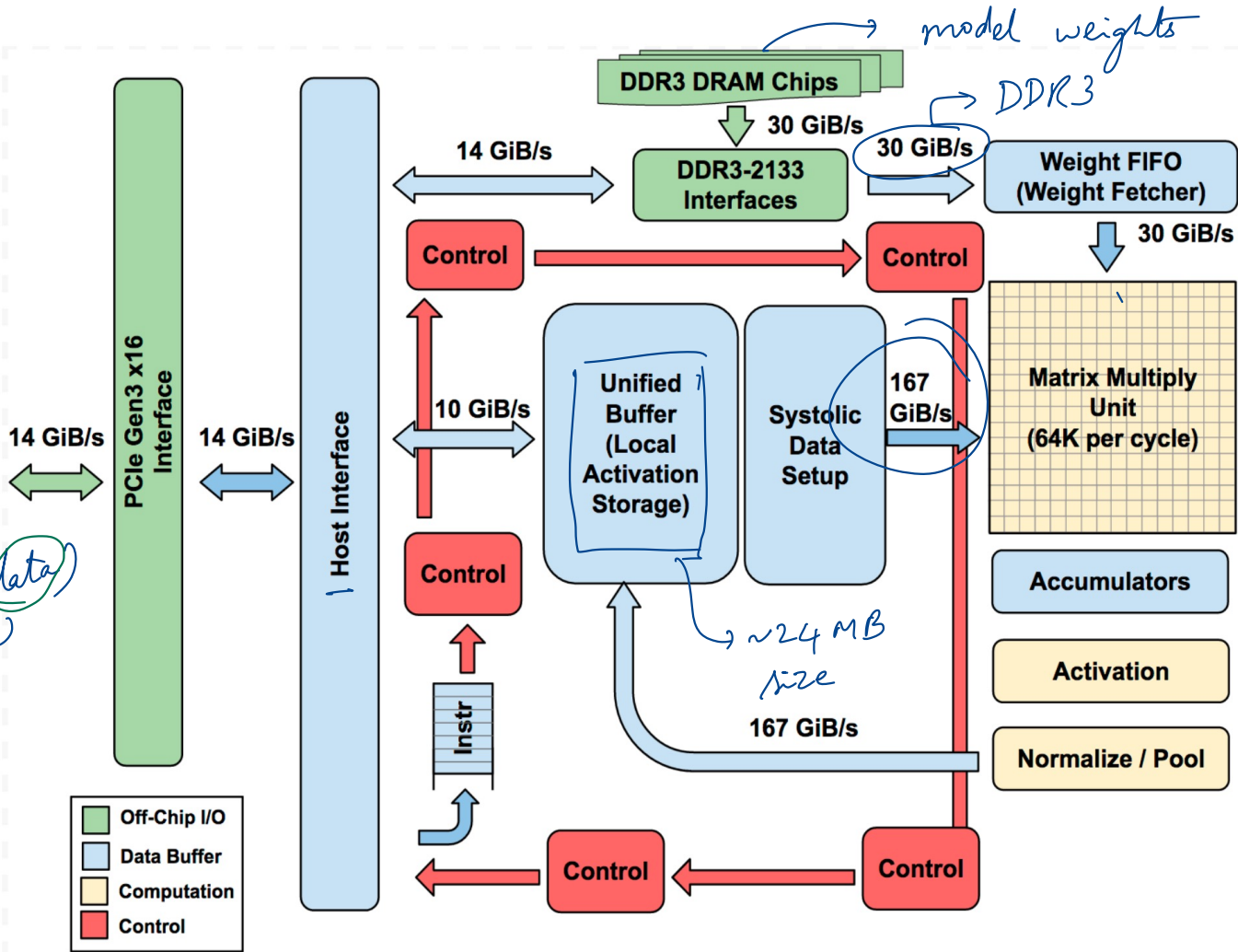
- Two kinds of data

① Model weights
→ static, load
them ahead of
time

② model forward (data)

input examples

Activations & layer
outputs also go into
unified buffer



INSTRUCTIONS

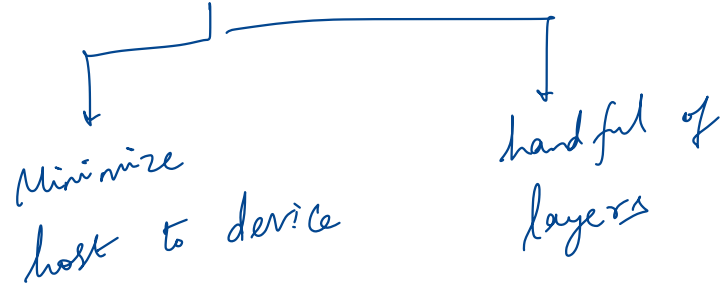
CISC format (why ?)

→ specialized with
Complex instructions

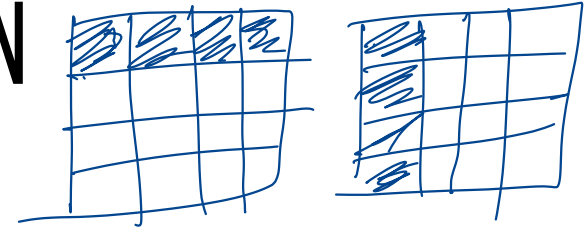
1. Read_Host_Memory_input
2. Read_Weights
3. MatrixMultiply/Convolve
4. Activate
5. Write_Host_Memory

↳ output

many ops to do to run
this instruction



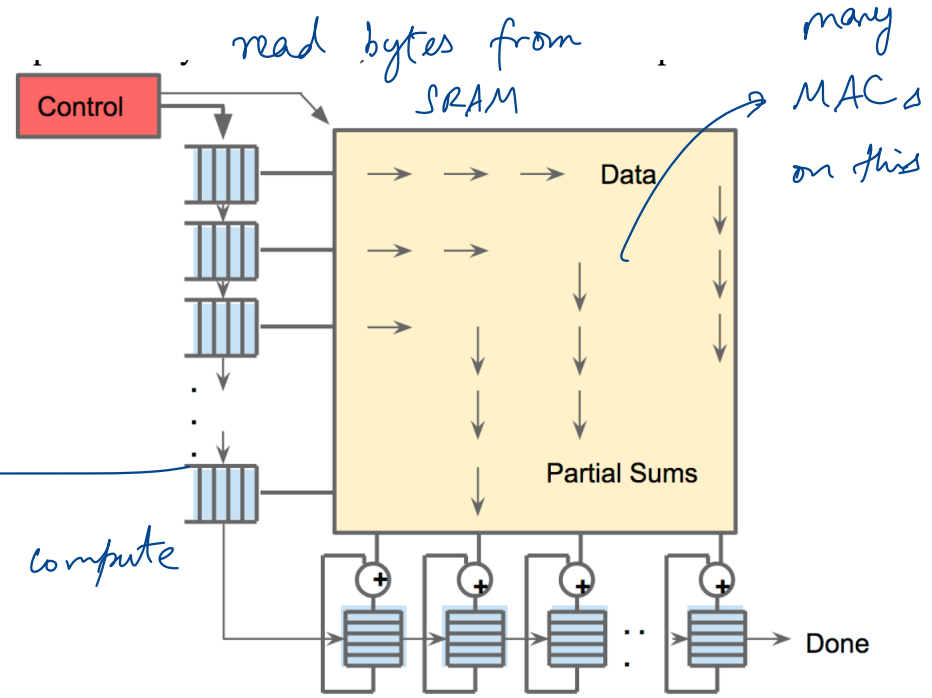
SYSTOLIC EXECUTION



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

$c[i] = a[i] * b[j]$

read SRAM
multiply
write SRAM



Systolic arrays
- Data streamed through compute units

x-axis: Computational intensity

y-axis: tput or FLOPs that you get

Blue line

↳ peak performance for given operation intensity

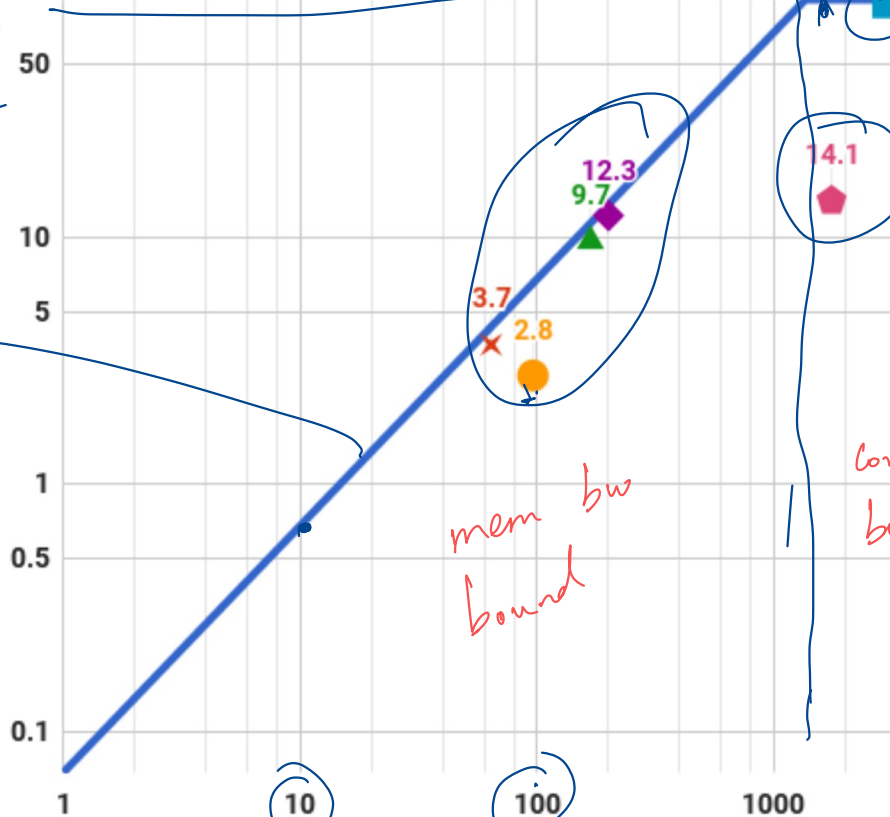
weight fetcher bottleneck
↳ most mem BW models

ROOFLINE MODEL

TPUs

86
Tops/sec

TeraOps/sec



Compute bound

- Roofline
- LSTM0
- LSTM1
- MLP1
- MLP0
- CNN0
- CNN1

mem bw bound

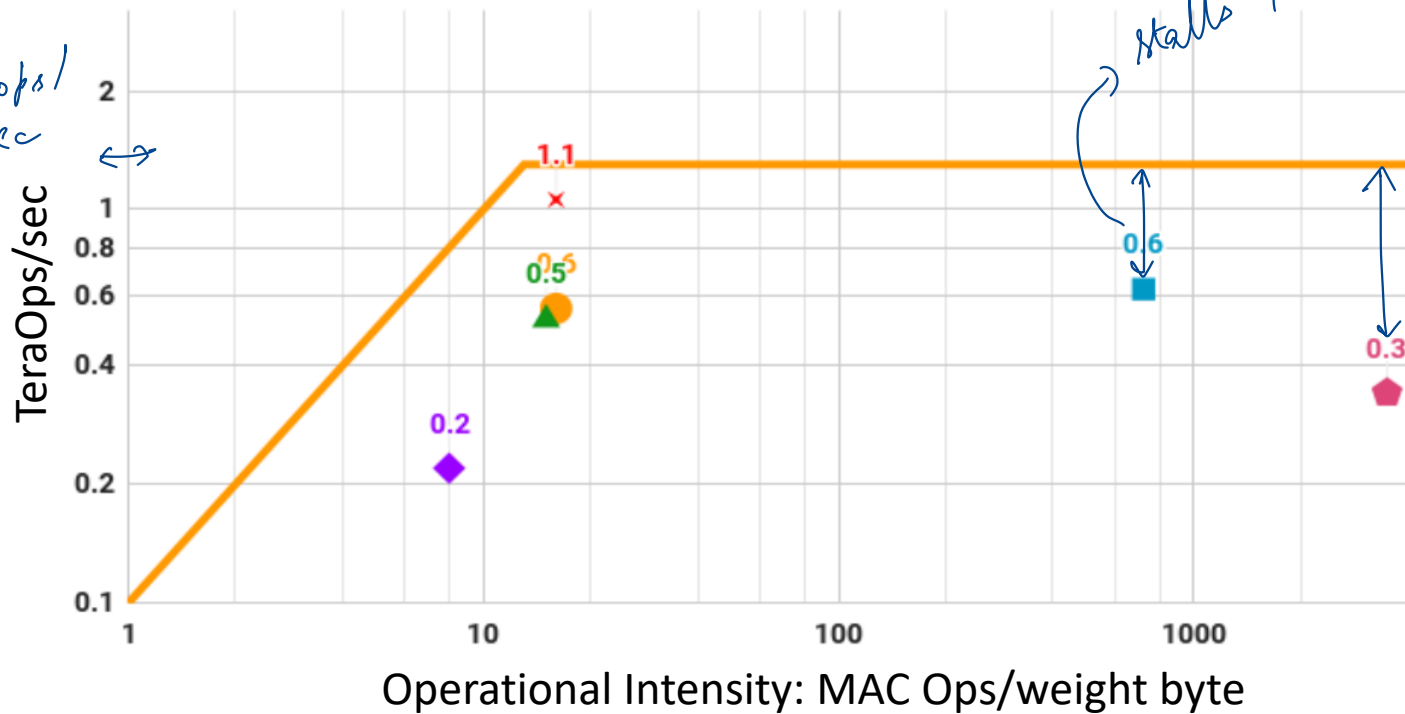
compute bound

Operational Intensity: MAC Ops/weight byte

CPU

HASWELL ROOFLINE

vb. 86 TeraOps/sec
on TPU



stalls / inefficiency

- Roofline
- LSTM0
- LSTM1
- MLP1
- MLP0
- CNN0
- CNN1

COMPARISON WITH CPU, GPU

<i>Model</i>	<i>Die</i>									
	<i>mm²</i>	<i>nm</i>	<i>MHz</i>	<i>TDP</i>	<i>Measured</i>		<i>TOPS/s</i>		<i>GB/s</i>	<i>On-Chip Memory</i>
					<i>Idle</i>	<i>Busy</i>	8b	FP		
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	<u>8 MiB</u>
TPU	<331*	28	700	75W	28W	40W	92	--	34	28 MiB

50%
better
than GPU

8-bit
integers

low
mem
bus

unified
buffer

SELECTED LESSONS

predictable

- Latency more important than throughput for inference
- LSTMs and MLPs are more common than CNNs → *Surprising in 2016!*
- Performance counters are helpful
- Remember architecture history → *Systolic arrays*

SUMMARY

New workloads → 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/P7mZsfK44PemjkXa7>

Type	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%

7ms SLO

able to support larger batch size

slightly higher than GPU latency at 100%.

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

As you add more hardware types

↳ need more profiling data

TPUs → quantization

↳ maintain more model variants

Sharing TPU might be tricky

coarse grained comp. units

- Additional need for software maps from PyTorch
↳ TPU

NEXT STEPS

Next week schedule

Tue: HeMem

Thu: Midterm 2