



Domain-Specific Architectures for Deep Neural Networks

David Patterson, Google AI and UC Berkeley
April 2019

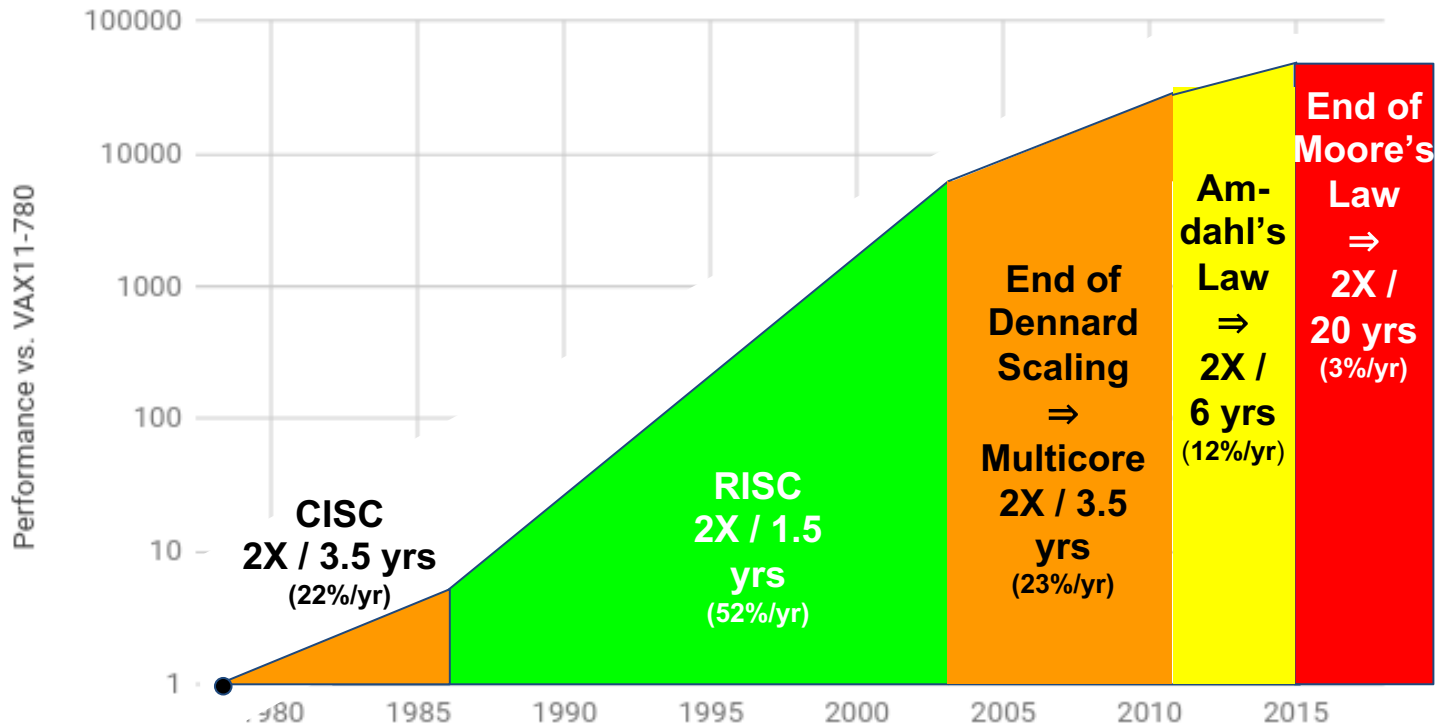
Based on Two Talks

In-Data Center Performance Analysis of a Tensor Processing Unit
ISCA 2017, Jouppi, Young, Patil, Patterson, et al.
(TPUv1, for inference)

Codesign in Google TPUs, by Cliff Young, HotChips 2017
(TPUv2 and later, for training and inference)

End of Growth of Performance?

40 years of Processor Performance



Based on SPECintCPU. Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

What's Left?

Since

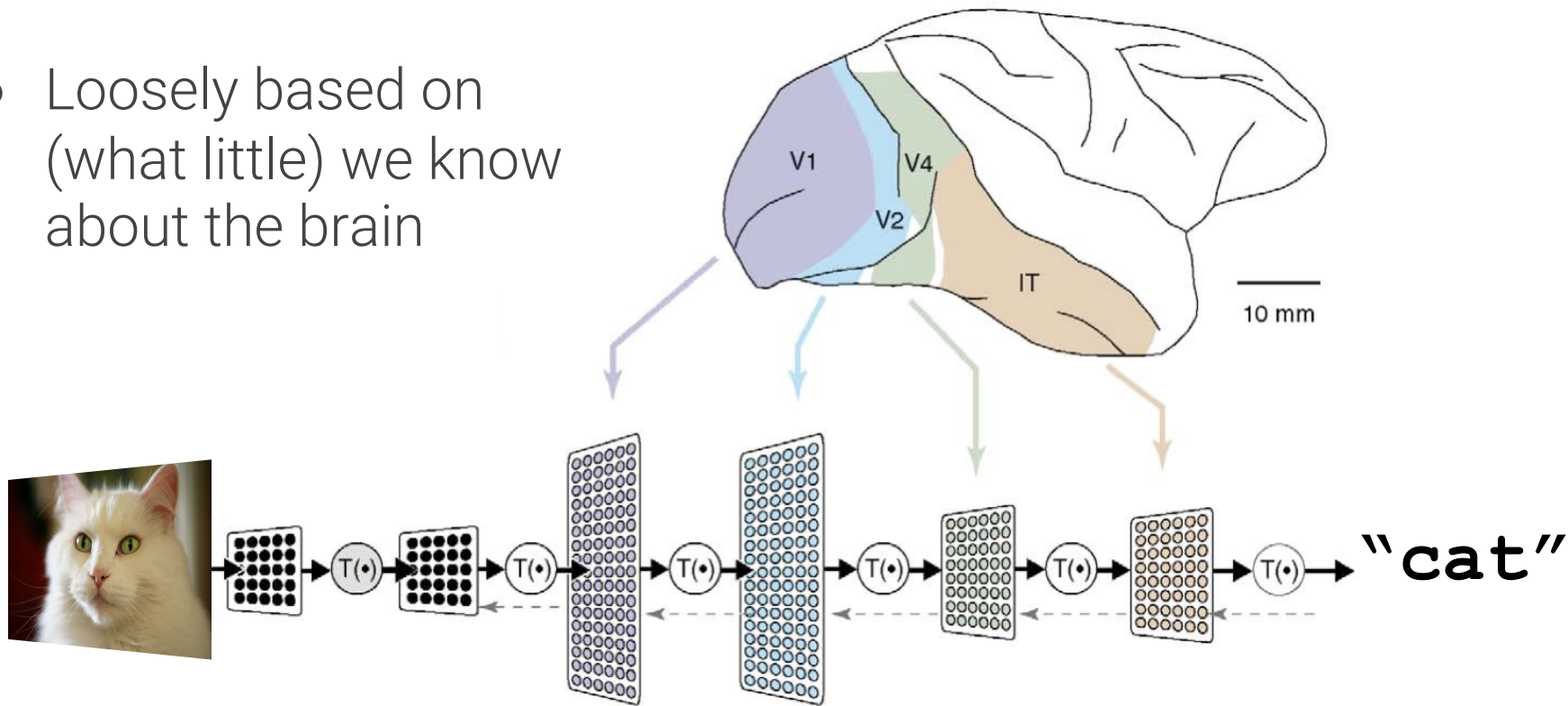
- Transistors not getting much better
- Power budget not getting much higher
- Already switched from 1 inefficient processor/chip to N efficient processors/chip

Only path left is *Domain Specific Architectures*

- Just do a few tasks, but extremely well

What is Deep Learning?

- Loosely based on (what little) we know about the brain



Key NN Concepts for Architects

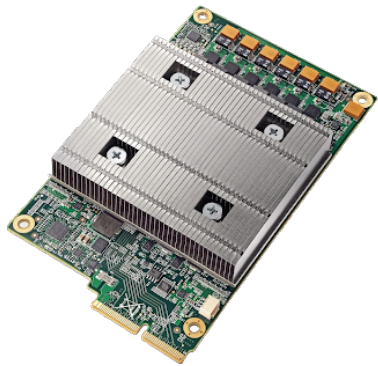
- *Training* or learning (development)
vs. *Inference* or prediction (production)
- *Batch size*
 - Problem: DNNs have millions of weights that take a long time to load from memory (DRAM)
 - Solution: Large batch \Rightarrow Amortize weight-fetch time by inferring (or training) many input examples at a time
- Floating-Point vs. Integer ("*Quantization*")
 - Training in Floating Point on GPUs popularized DNNs
 - Inferring in Integers faster, lower energy, smaller

- 2013: Prepare for success-disaster of new DNN apps
 - Scenario with users speaking to phones 3 minutes per day:
If only CPUs, need 2X-3X times whole fleet
 - Unlike some hardware targets, DNNs applicable to a wide range of problems, so can reuse for solutions in speech, vision, language, translation, search ranking, ...
- Custom hardware to reduce the TCO of DNN inference phase by 10X vs. CPUs
 - Must run existing apps developed for CPUs and GPUs
- A very short development cycle
 - Started project 2014, running in datacenter 15 months later:
Architecture invention, compiler invention, hardware design, build, test, deploy
- Google CEO Sundar Pichai reveals Tensor Processing Unit at Google I/O on May 18, 2016 as “10X performance/Watt”

cloudplatform.googleblog.com/2016/05/Google-supercharges-machine-learning-tasks-with-custom-chip.html

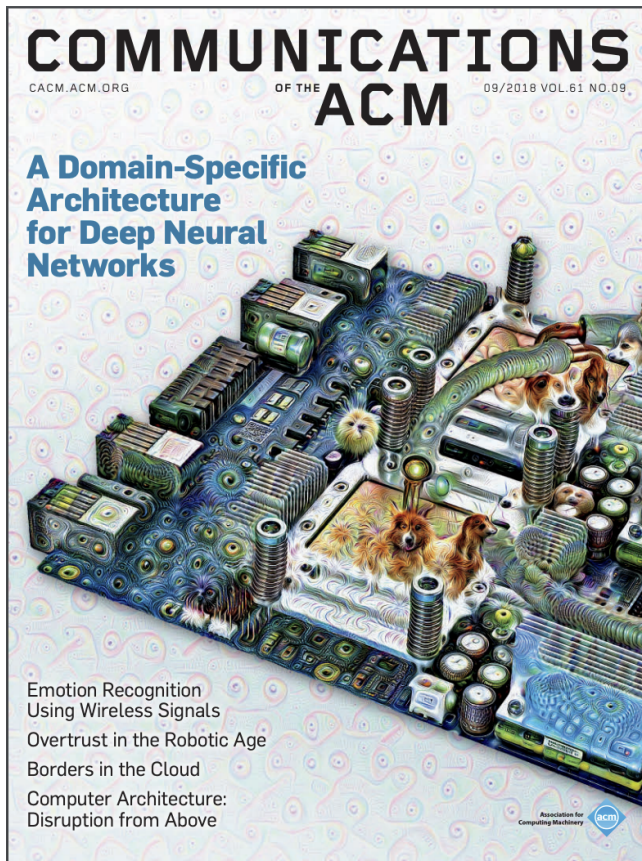
Tensor Processing Unit v1 (deployed 2015)

Google-designed chip for neural net inference



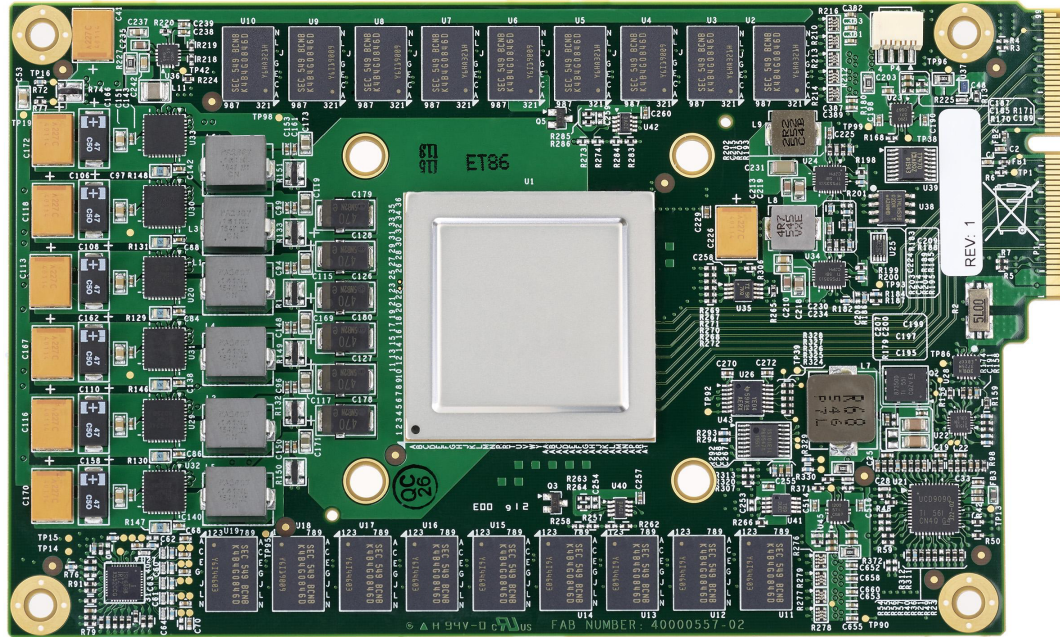
In production use for ≈ 4 years: used by billions on search queries, for neural machine translation, for AlphaGo match, ...

[A Domain-Specific Architecture for Deep Neural Networks](#), Jouppi, Young, Patil, & Patterson, *Communications of the ACM*, September 2018



TPUv1 Card & Package

- TPUv1 Card to replace a disk
- Up to 4 cards / server



Inference Datacenter Workload (95%)

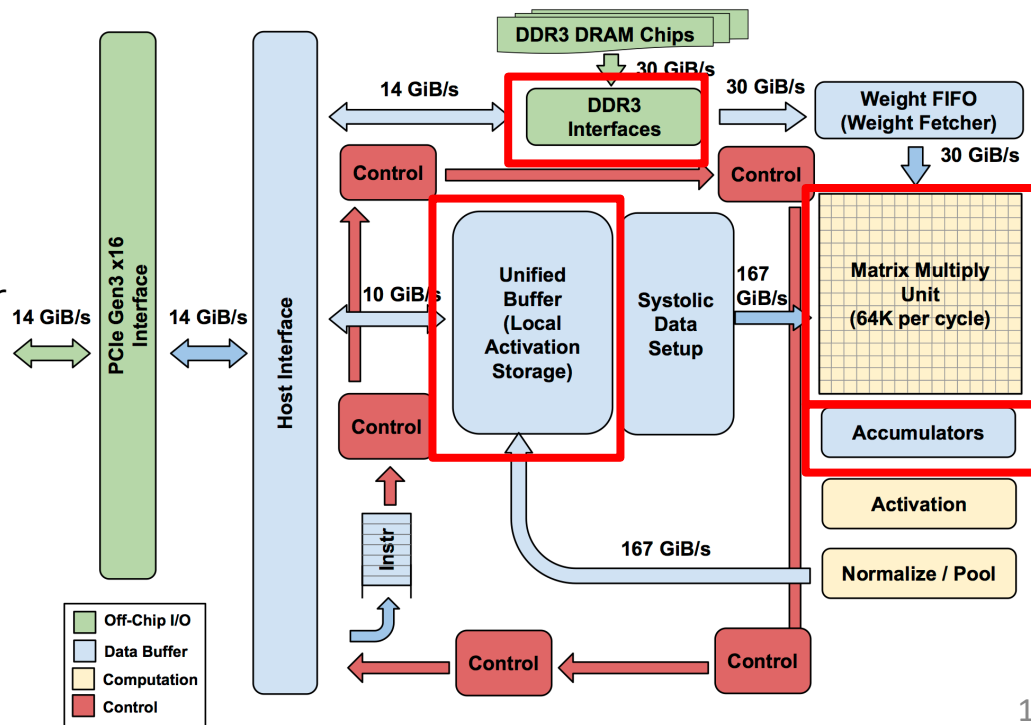
<i>Name</i>	<i>LOC</i>	<i>Layers</i>					<i>Nonlinear function</i>	<i>Weights</i>	<i>TPUv1 Ops / Weight Byte</i>	<i>TPUv1 Batch Size</i>	<i>% Deployed</i>
		<i>FC</i>	<i>Conv</i>	<i>Vector</i>	<i>Pool</i>	<i>Total</i>					
MLP0	0.1k	5				5	ReLU	20M	200	200	61%
MLP1	1k	4				4	ReLU	5M	168	168	
LSTM0	1k	24		34		58	sigmoid, tanh	52M	64	64	29%
LSTM1	1.5k	37		19		56	sigmoid, tanh	34M	96	96	
CNN0	1k		16			16	ReLU	8M	2888	8	5%
CNN1	1k	4	72		13	89	ReLU	100M	1750	32	

TPUv1 Architecture and Implementation

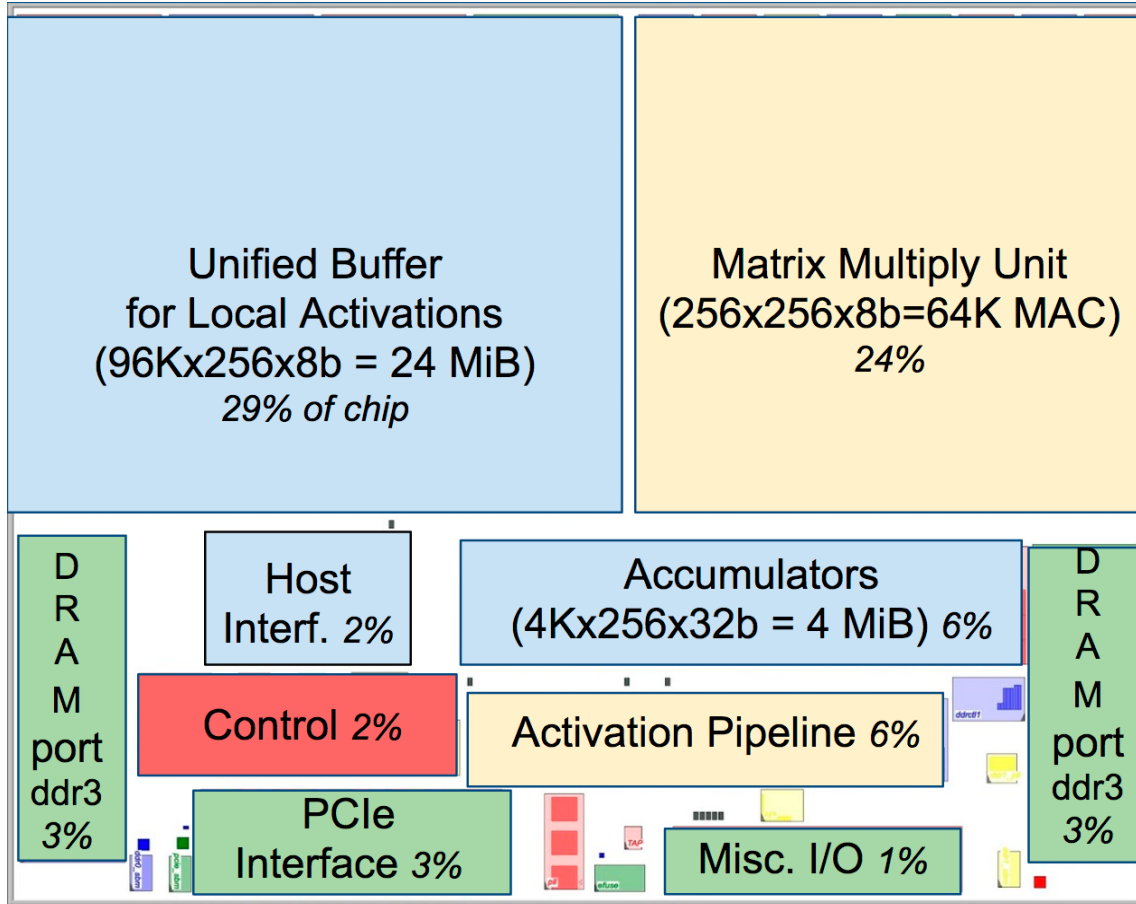
- Add as accelerators to existing servers
 - So connect over I/O bus (“PCIe”)
 - TPUv1 \approx matrix accelerator on I/O bus
- Host server sends it instructions like a Floating Point Unit
 - Unlike GPU that fetches and executes own instructions

- The Matrix Unit: 65,536 (256x256)
8-bit multiply-accumulate units
- 700 MHz clock rate
- Peak: 92T operations/second
 - $65,536 * 2 * 700M$
- >25X as many MACs vs GPU
- >100X as many MACs vs CPU
- 4 MiB of on-chip Accumulator memory
- 24 MiB of on-chip Unified Buffer (activation memory)
- 3.5X as much on-chip memory vs GPU
- Two 2133MHz DDR3 DRAM channels
- 8 GiB of off-chip weight DRAM memory

TPUv1: High-level Chip Architecture



TPUv1: a Neural Network Accelerator Chip



TPUv1 Architecture, programmer's view

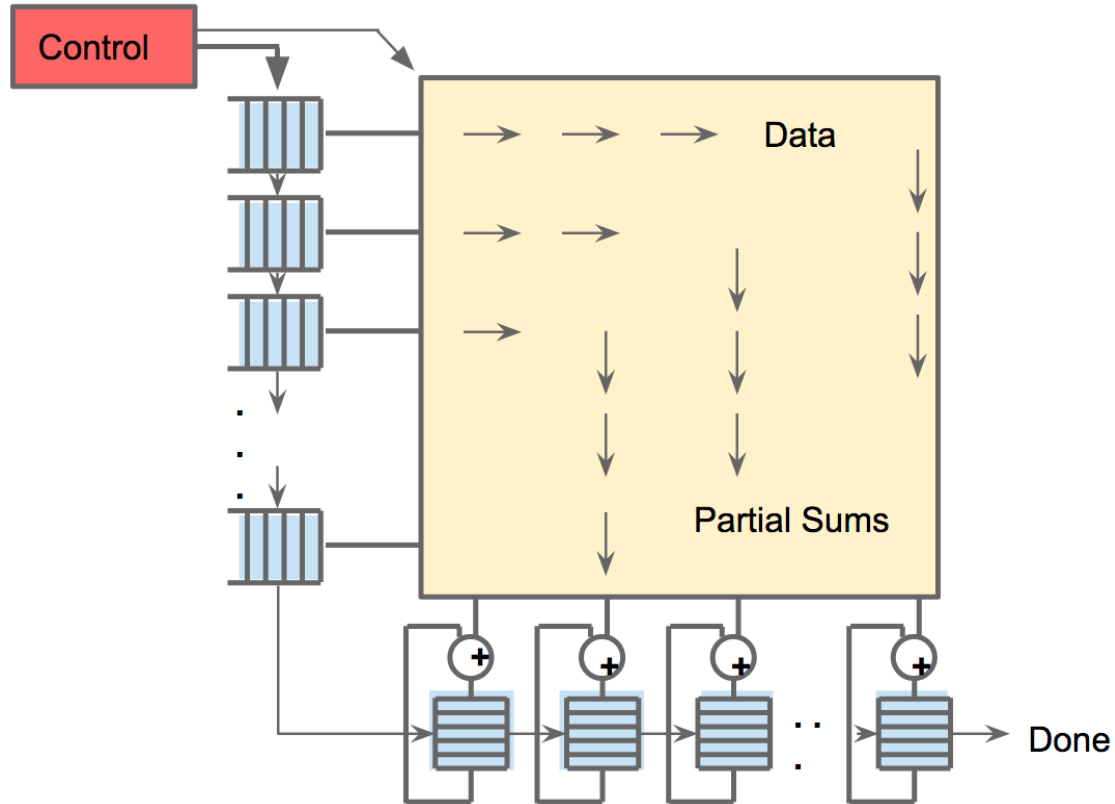
- 5 main (CISC) instructions
 - Read_Host_Memory
 - Write_Host_Memory
 - Read_Weights
 - MatrixMultiply/Convolve
 - Activate (ReLU, Sigmoid, Maxpool, LRN, ...)
- Average Clock cycles per instruction: >10
- 4-stage overlapped execution, 1 instruction type / stage
 - Execute other instructions while matrix multiplier busy
- Complexity in SW: No branches, in-order issue, SW controlled buffers, SW controlled pipeline synchronization

Systolic Execution in Matrix Array

- Problem: energy/ time for repeated SRAM accesses of matrix multiply
- Solution: “Systolic Execution” to compute data on the fly in buffers by pipelining control and data
 - Relies on data from different directions arriving at cells in an array at regular intervals and being combined

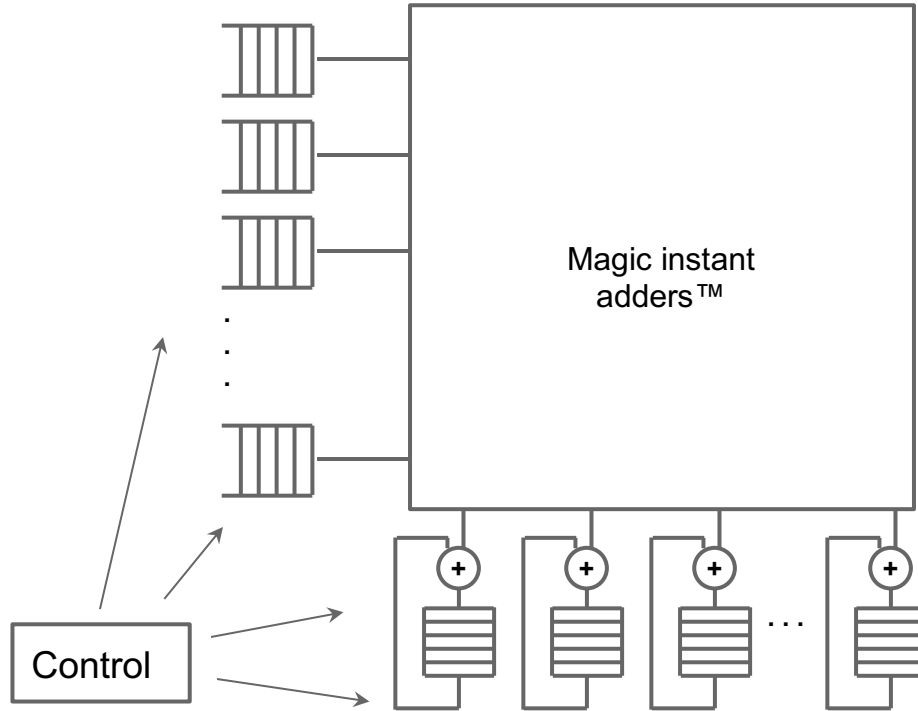
Systolic Execution:

Control and Data are pipelined



Can now ignore pipelining in matrix

Pretend each 256B input read at once, & they instantly update 1 location of each of 256 accumulator RAMs.



Relative Performance: 3 Contemporary Chips (2015)

<i>Processor</i>	<i>mm²</i>	<i>Clock MHz</i>	<i>TDP Watts</i>	<i>Idle Watts</i>	<i>Memory GB/sec</i>	<i>Peak TOPS/chip</i>	
						<i>8b int.</i>	<i>32b FP</i>
CPU: Haswell (18 core)	662	2300	145	41	51	2.6	1.3
GPU: Nvidia K80 (2 / card)	561	560	150	25	160	--	2.8
TPUv1	<331*	700	75	28	34	91.8	--

*TPUv1 is less than half die size of the Intel Haswell processor

K80 and TPUv1 in 28 nm process; Haswell fabbed in Intel 22 nm process

These chips and platforms chosen for comparison because widely deployed in Google data centers

GPUs and TPUs added to CPU server

Relative Performance: 3 Platforms

<i>Processor</i>	<i>Chips/ Server</i>	<i>DRAM</i>	<i>TDP Watts</i>	<i>Idle Watts</i>	<i>Observed Busy Watts in datacenter</i>
CPU: Haswell (18 cores)	2	256 GB	504	159	455
NVIDIA K80 (13 cores) (2 die per card; 4 cards per server)	8	256 GB (host) + 12GB x 8	1838	357	991
TPUv1 (1 core) (1 die per card; 4 cards per server)	4	256GB (host) + 8GB x 4	861	290	384

These chips and platforms chosen for comparison because widely deployed in Google datacenters

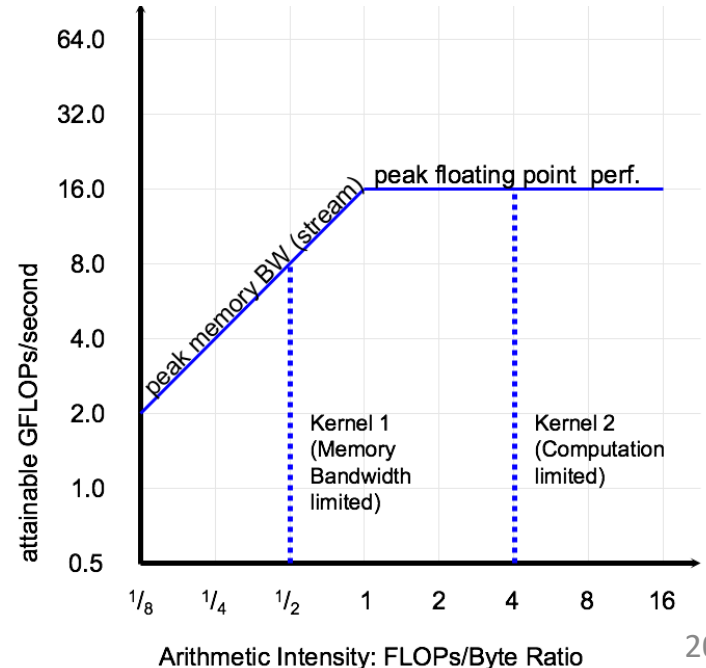
2 Limits to performance:

1. Peak Computation
2. Peak Memory Bandwidth
(For apps with large data that don't fit in cache)

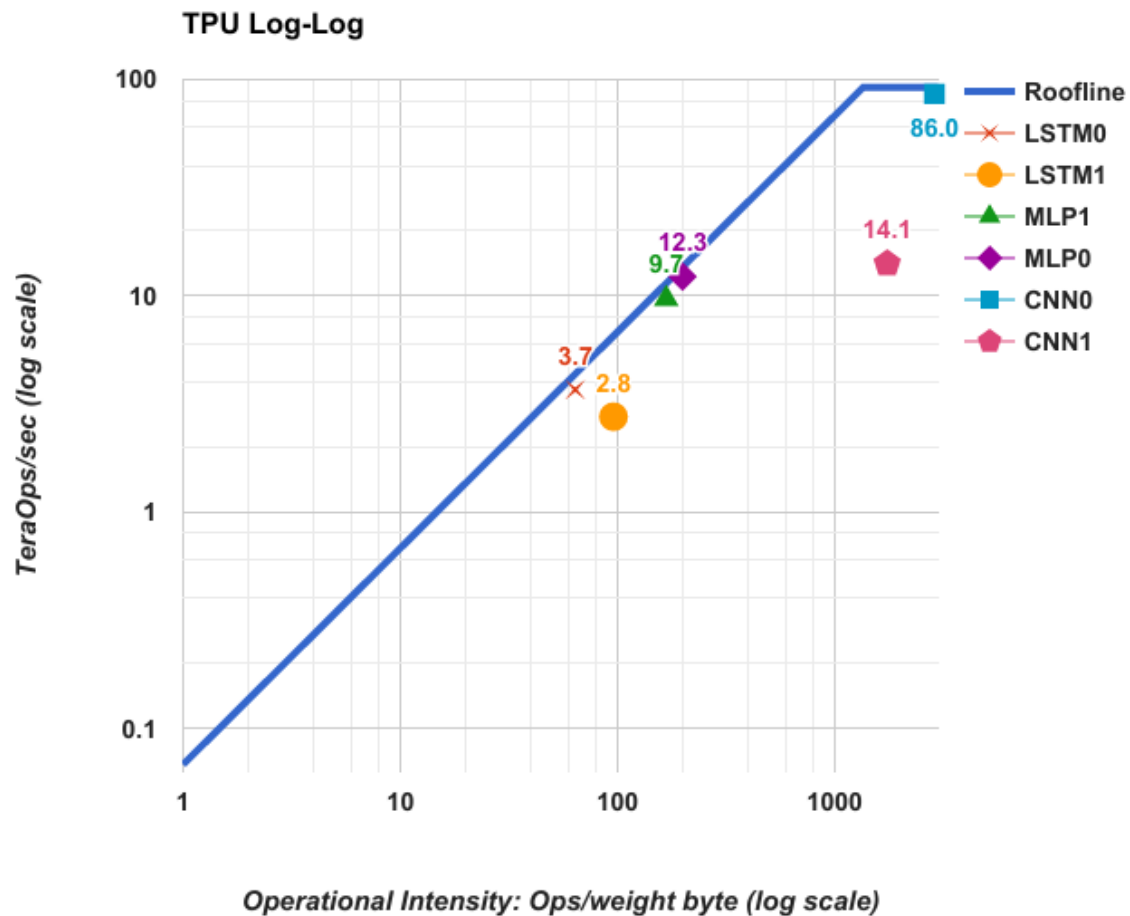
Arithmetic Intensity (FLOP/byte or reuse) determines which limit
Weight-reuse = Arithmetic Intensity for DNN roofline

Roofline Visual Performance Model

$$\text{GFLOP/s} = \text{Min}(\text{Peak GFLOP/s}, \text{Peak GB/s} \times \text{AI})$$

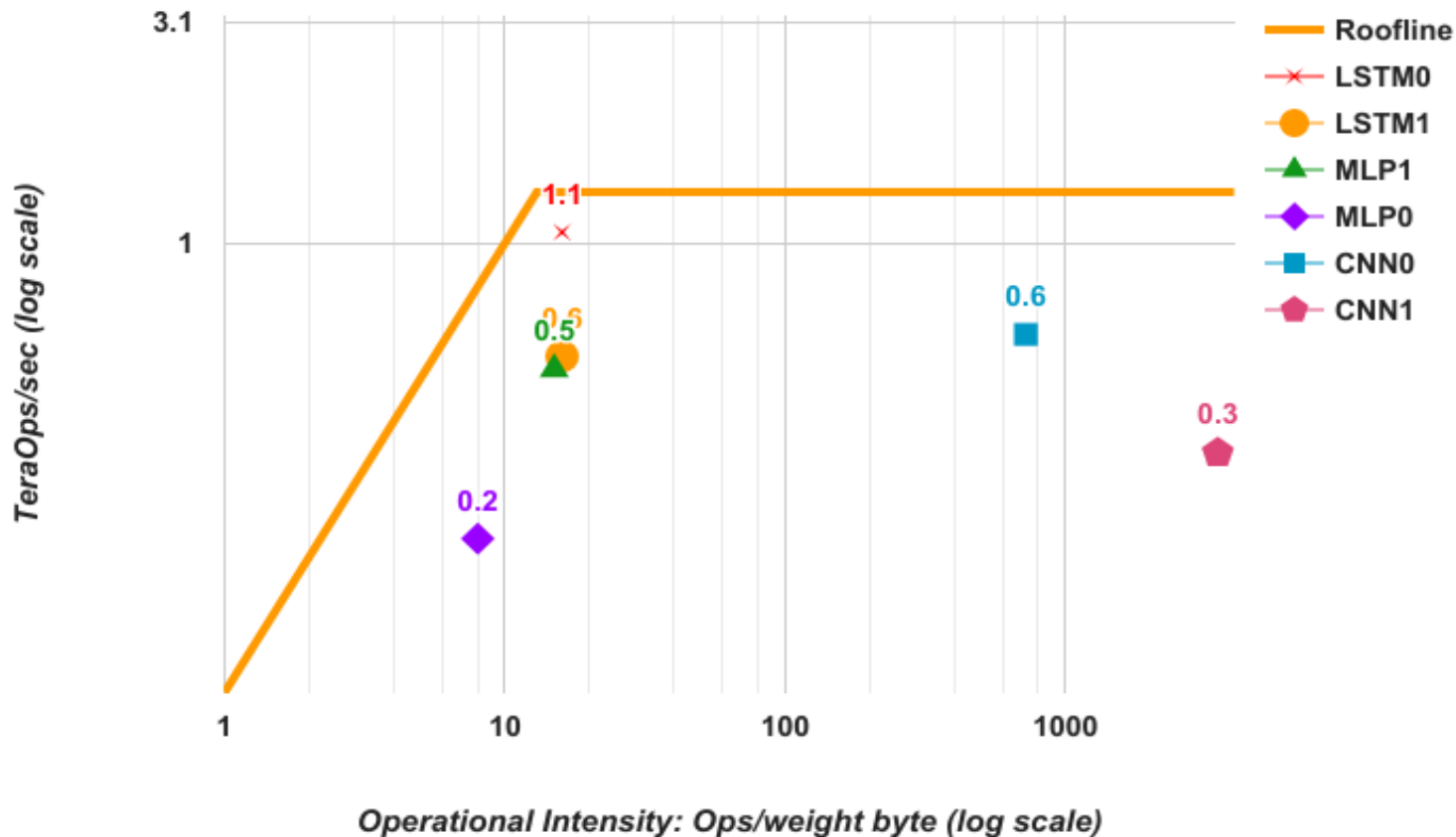


TPUv1 Die Roofline

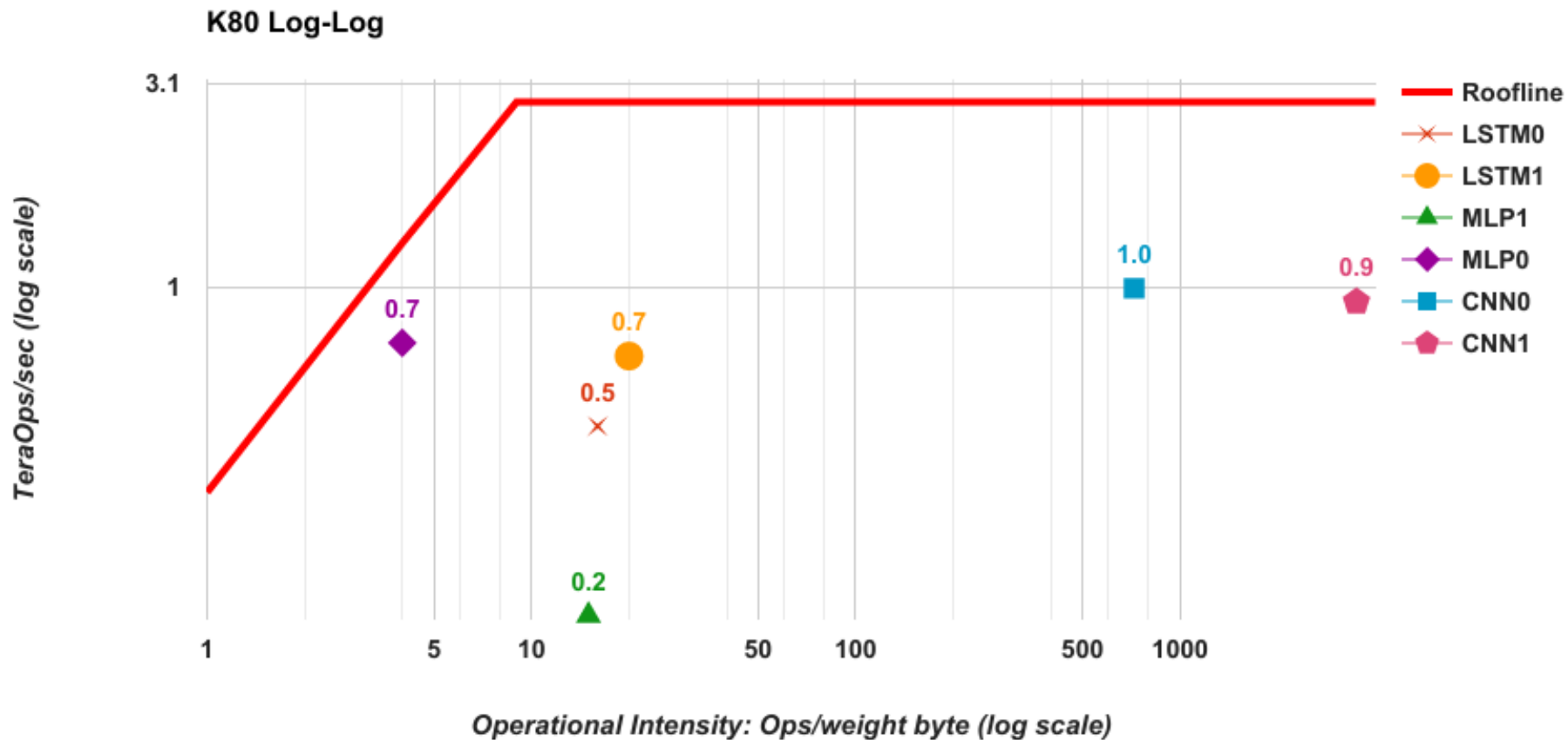


Haswell (CPU) Die Roofline

Haswell Log-Log



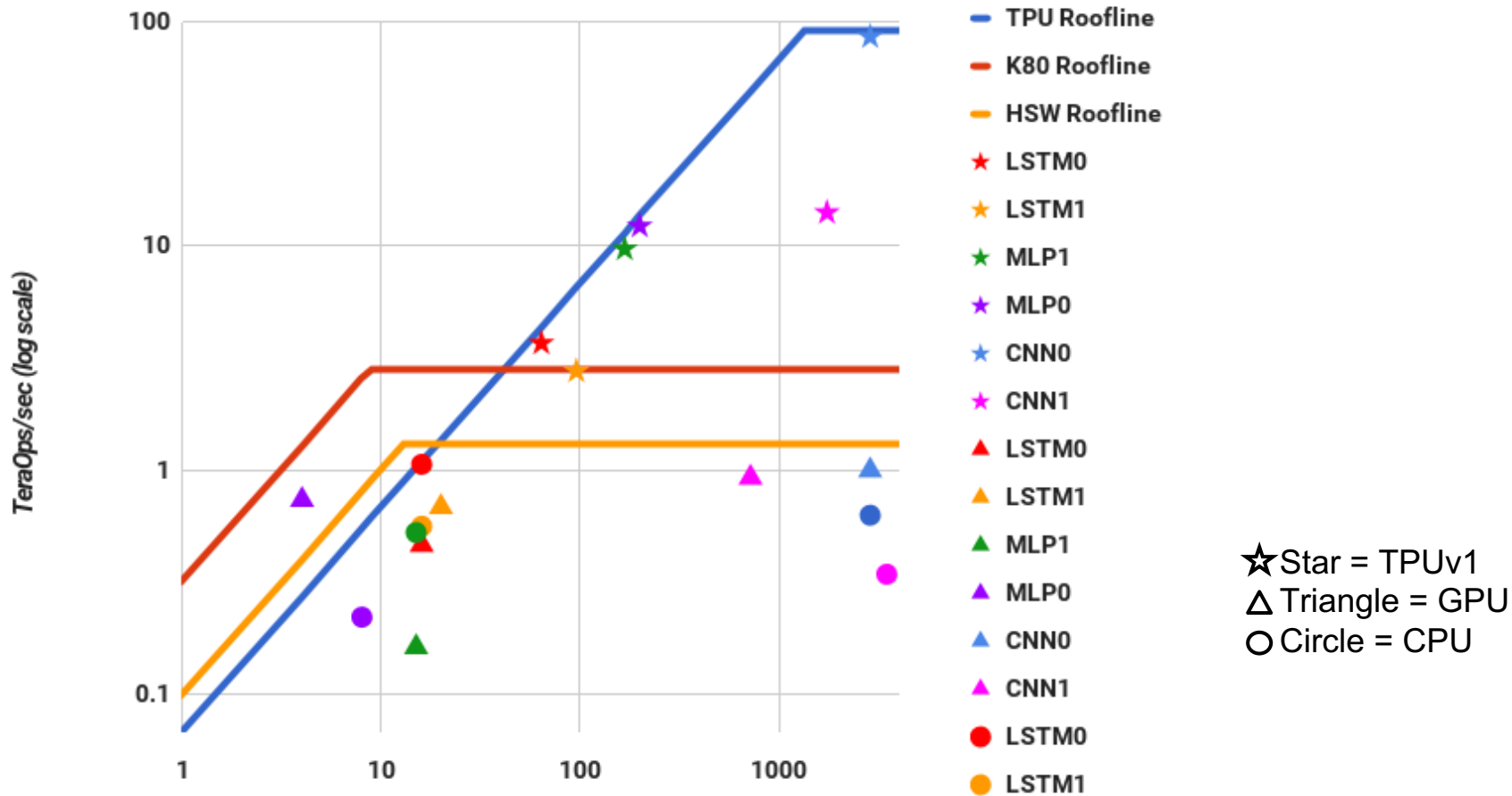
K80 (GPU) Die Roofline



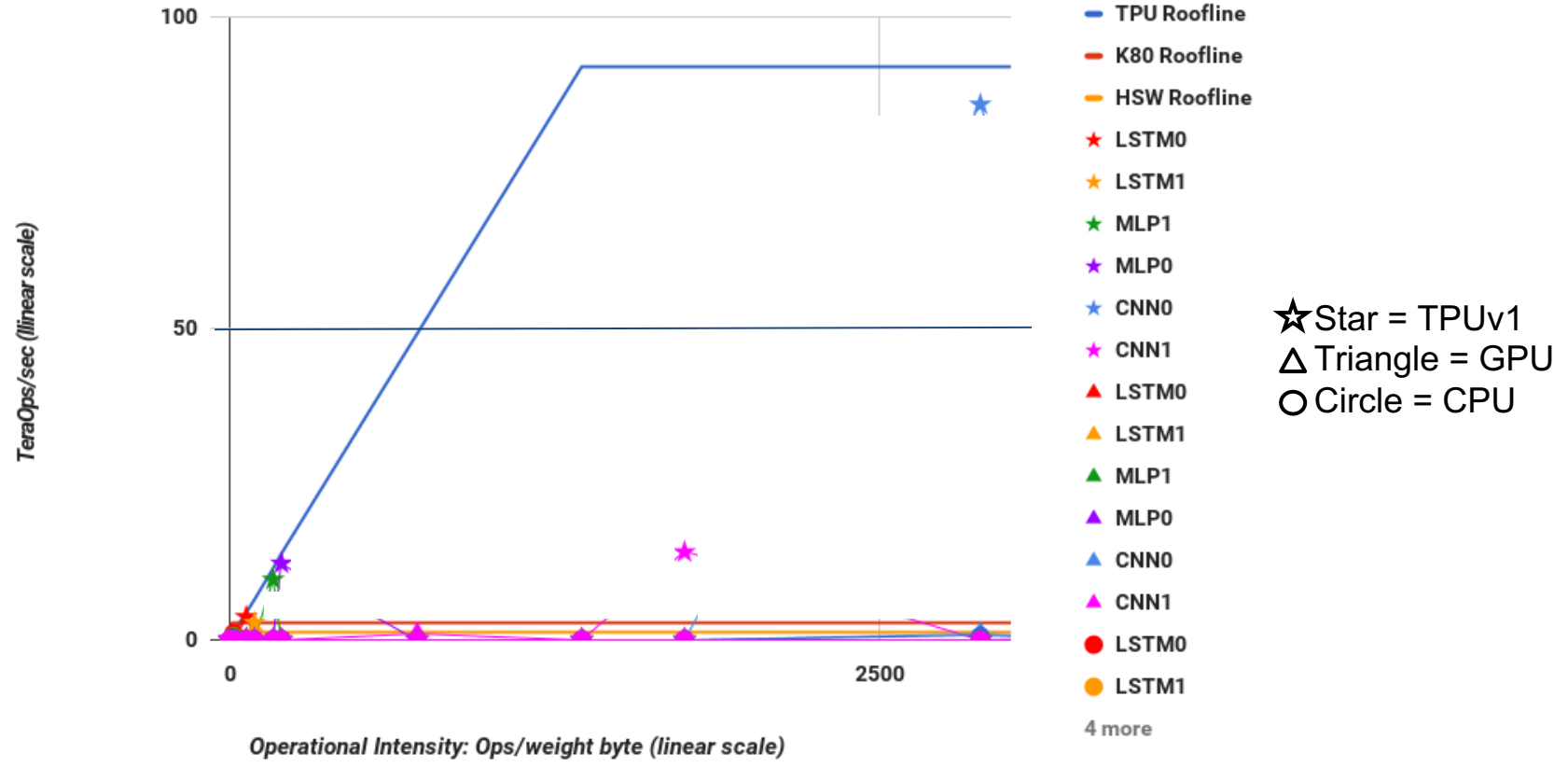
Why so far below Rooflines? (MLP0)

<i>Type</i>	<i>Batch</i>	<i><u>99th% Response</u></i>	<i>Inf/s (IPS)</i>	<i>% Max IPS</i>
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%
TPUv1	200	7.0 ms	225,000	80%
TPUv1	250	10.0 ms	280,000	100%

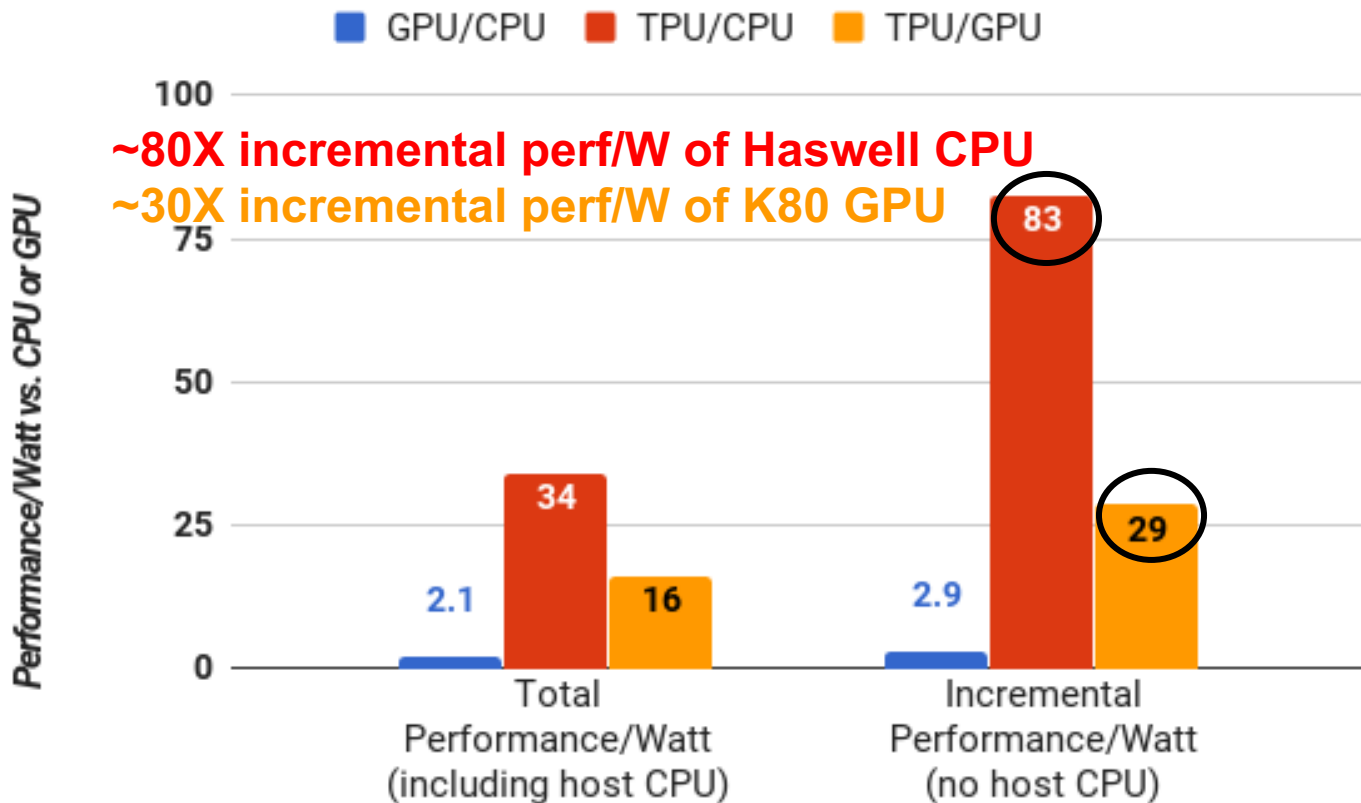
Log Rooflines for CPU, GPU, TPUv1



TPUv1



Perf/Watt TPUv1 vs CPU & GPU

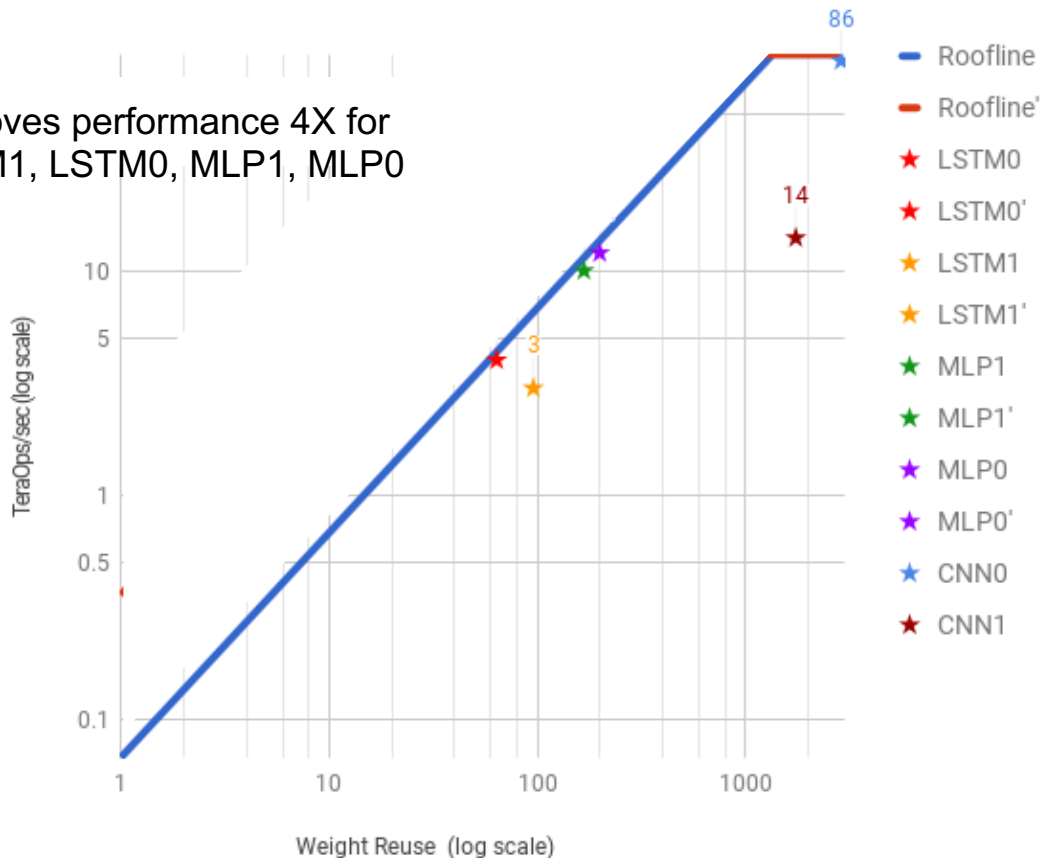


- Current DRAM
 - 2 DDR3 2133 \Rightarrow 34 GB/s
- Replace with GDDR5 like in K80 \Rightarrow 180 GB/s
 - Move Ridge Point from 1400 to 256

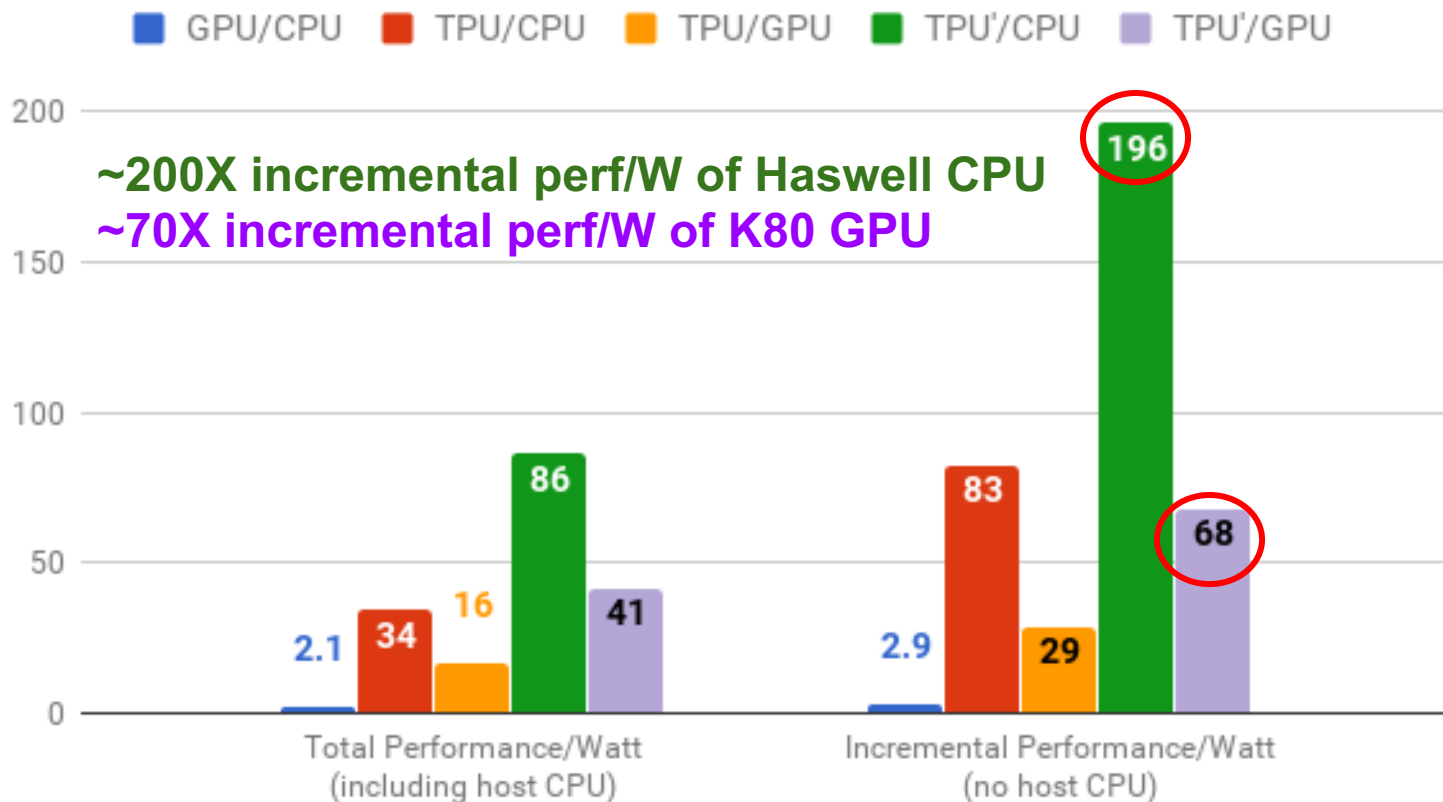
Improving TPUv1:
Move “Ridge Point” to
the left

Revised TPUv1 Raises Roofline

Improves performance 4X for LSTM1, LSTM0, MLP1, MLP0



Perf/Watt Original & Revised TPUv1



Related Work

Two survey articles document that custom NN ASICs go back at least 25 years [Int96][Aas02]. For example, CNAPS chips contained a 64 SIMD array of 16-bit by 8-bit multipliers, and several CNAPS chips could be connected together to a sequencer [Ham90]. The Synapse-1 system was based on a custom systolic multiply-accumulate chip called the MA-16, which performed sixteen 16-bit multiplies at a time [Ram91]. The system concatenated several MA-16 chips together and had custom hardware to do activation functions.

Twenty-five SPERT-II workstations, accelerated by the T0 custom ASIC, were deployed starting in 1995 to do both NN training and inference for speech recognition [Aas98]. The 40-MHz T0 added vector instructions to the MIPS instruction set architecture. The eight-lane vector unit could produce up to sixteen 32-bit arithmetic results per clock cycle based on 8-bit and 16-bit inputs, making it 25 times faster at inference and 20 times faster at training than a SPARC-20 workstation. They found that 16 bits were insufficient for training, so they used two 16-bit words instead, which doubled training time. To overcome that drawback, they introduced “bunches” (batches) of 32 to 1000 data points to reduce time spent updating weights, which made it faster than training with one word but no batches.

To more recent DaianNao Family of NN architectures minimizes memory accesses both on the chip and to external DRAM by having efficient architectural support for the memory access patterns that appear in NN applications [Keu16]. [Chai16a]. All use 16-bit integer operations and all designs dove down to layout, but no chips were fabricated. The original DaianNao uses an array of 16-bit integer multiply-accumulate units with 44 kB of on-chip memory and is estimated to be 3 mm² (65 nm), and run at 1 GHz, and to consume 0.5W [Chai16a]. Most of this energy went to DRAM accesses for weights, so one successor DaianNao (“big computer”) includes eDRAM of 36 MB of weights on chip [Chai16b]. The goal was to have enough memory in a multiplex system to avoid external DRAM accesses. The follow-on PuDanNao (“general computer”) is aimed at more traditional machine learning algorithms beyond DNNs, such as support vector machines [Liu15]. Another offshoot is ShiDanNao (“vision computer”) aimed at CNNs, which avoids DRAM accesses by connecting the accelerator directly to the sensor [Dai15].

The Convolution Engine is also focused on CNNs for image processing [Qad13]. This design deploys 64 10-bit multiply-accumulator units and customizes a Tensilica processor estimated to run at 800 MHz in 45 nm. It is projected to be 8X to 15X more energy efficient than an SIMD processor, and within 2X to 3X of custom hardware designed just for a specific kernel.

The Fatfish benchmark paper recently reports results contradictory to ours, with the GPU training inference much faster than the CPU [Ade16]. However, their CPU and GPU are not server-class, the CPU has only four cores, the applications do not use the CPU’s AVX instructions, and there is no response-time cutoff (see Table 4) [Bro16].

Catapult is the most widely deployed example of using reconfigurability to support DNNs, which many have proposed [Far09][Chai10][Fu11][Pee13][Cav15][Zha15]. They chose FPGAs over GPUs to reduce power as well as the risk that latency-sensitive applications wouldn’t map well to GPUs. FPGAs can also be re-spurposed, such as for search, compression, and network interface cards [Put15]. The TPU project actually began with FPGAs, but we abandoned them when we saw that the FPGAs of that time were not competitive in performance compared to the GPUs of that time, and the TPU could be much lower power than GPUs while being as fast or faster, giving it potentially significant benefits over both of FPGAs and GPUs.

Although first published in 2014 [Put14], Catapult is a TPU contemporary since it deployed 28-nm Stratix V FPGAs into datacenters concurrently with the TPU in 2015. Catapult has a 200 MHz clock, 3,926 18-bit MACs, 5 MB of on-chip memory, 11 GBs memory bandwidth, and uses 25 Watts. The TPU has a 900 MHz clock, 65,536 8-bit MACs, 28 MB, 34 GBs, and typically uses 40 Watts. A revised version of Catapult uses newer FPGAs and was deployed at large scale in 2016 [Liu16].

Catapult V1 runs CNNs—using a systolic matrix multiplier—2.3X as fast as a 2.1 GHz, 16-core, dual-socket server [Ovi15a]. Using the next generation of FPGAs (14-nm Artix 10) of Catapult V2, performance might go up to 7X, and perhaps even 17X with more careful floorplanning [Ovi15b]. Although it’s apples versus oranges, a current TPU die runs its CNNs 40X to 70X versus a somewhat faster server (Tables 2 and 6). Perhaps the biggest difference is that to get the best performance the user must write long programs in the low-level hardware-design-language Verilog [Met16][Put16] versus writing short programs using the high-level TensorFlow framework. That is, reprogrammability comes from software for the TPU rather than from firmware for the FPGA.

Recent research, which appeared after the TPU was deployed, accelerates DNNs by optimizing the cases when weights and data are very small or zero. Our tight schedule precluded such optimizations in the TPU, but we saw the same opportunity in our studies. The Efficient Inference Engine is based on a first pass that reduces the number of weights by about a factor of 10 [Han15] as a separate step by filtering out very small values and then uses Huffman encoding to shrink the data even further to improve inference performance [Han16]. Cvnutils [Aib16] avoids multiplications when an activation input is zero—which it is 44% of the time, presumably in part due to ReLU nonlinear function that transforms negative values to zero—to improve performance by an average 1.4 times.

Eyeris is a novel, low-power dataflow architecture that takes advantage of zeros by run-length encoding data to reduce the memory footprint and saves power by avoiding computations when an input is zero [Chai16a]. Using Eyeris terminology, a TPU convolutional layer maps C and M to the rows and columns of the matrix unit, taking HWN cycles to perform one pass. With high C/M, it takes RS passes to process the layer, for low C/M, a number of techniques reduce passes and improve utilization. (More can be found in the online references [Ros11][Ros15][Ros15c][Ros15d][Tho15][Liu15].)

Eyeris is a design system that crosses algorithm, architecture, and circuit disciplines to reduce power by 8X to 1.6x by pruning activation data with small values and in part by quantizing the data [Real16]. [Gup15] looks at 16-bit fixed-point arithmetic for training instead of for inference. Others leverage the lower precision of DNN calculations by utilizing analog circuits during the computation to improve energy and performance [LiK16] [Sha16]. By tailoring an instruction set to DNNs, Cambridge reduces code size [Liu16]. Recent work looked at processor-in-memory architectures for NNs [Chai16][Kim16].

Related Work

Comparing the TPU to some of these architectures:

- [Chai16a] DMAAs data from DRAM to input and weight buffers. They are read by the 3-stage pipelined NFU that performs multiplies, adds, and non-linear-functions; the results go to the output buffer, and then to DRAM. The NFU has no storage and isn’t systolic.
- [Gup15] appears to stream both matrix inputs while storing partial sums in the systolic array; the TPU stores the weights in a stochastic random file streaming the other input and the pre-activation partial sums. The TPU doesn’t support stochastic rounding.
- [Zha15] is built out of computation units equivalent to a 4x2 version of the TPU matrix unit. In an ASIC, the wiring cost of the crossbars that connect input and output buffers to these compute engines would be significant. We are surprised that we didn’t see architectural support for additional reductions to combine results from compute engines in [Zha15].

All three of [Gup15][Chai16a][Zha15] store activations in DRAM during computation; the TPU’s Unified Buffer is sized so that no DRAM spilling or reloading happens during normal operation.

References

[Aba16] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., 2016. TensorFlow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*.

[Aib16] Alberio, J., Judd, P., Hehringer, T., Aamodt, T., Jørgen, N.E. and Moshovos, A. 2016. Cvnutils: Inefficient-but-Cheap Deep Neural Network Computing. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Ade16] Ade, E., Naraia, S., Rogers, B., Vici, G.V. and Broussard, D., 2016, September. Fatfish: reference workloads for modern deep learning methods. *IEEE International Symposium on Workload Characterization (ISWC)*.

[Aas02] Aasawale, K. 2002. Programmable Neurocomputers. In *The Handbook of Brain Theory and Neural Networks: Second Edition*, M. A. Arbib (Ed.), MIT Press. ISBN 0-262-01197-2. November 2002. <http://people.csis.yorku.ca/abib/abib-handbook-neurocomp.pdf>

[Aas98] Aasawale, K. 1998. Aasawale, K., Beck, Johnson, J., Wazarynek, J., Kingbirdy, B. and Morgan, N., November 1998. Training Neural Networks with Spars-H. Chapter 11 in *Parallel Architectures for Artificial Neural Networks: Paradigms and Implementations*, N. Sandhu and P. Sarachandran (Eds.), Elsevier, Amsterdam. ISBN 0-444-82644-2.

[Bar16] Barua, L. and Haldol, U., 2007. The case for energy-proportional computing. *IEEE Computer*, vol. 40, (Bar16). Jan 6, September 29, 2016. Private. No. 2, February Issue for Amazon EC2 up-to 16 GPUs. <http://amazon.com/aws/whitepapers/ec2-q2-amazon-ec2-up-to-16-gpu/>

[Bro16] Brooks, D. November 4, 2016. Private communication.

[Cav16] Caulfield, A.M., Chang, E.S., Patnam, A., Haselman, H.A.J.F.M., Humphrey, S.H.M., Daniqi, P.K.J.Y.K., Onvichaov, I.T.M.K., Lanka, M.P., E.S. and Bharg, D.C.D., 2016. A 16-bit Scale Acceleration Architecture. *MFGO conference*.

[Cai15] Cavigelli, L., Gschwend, D., Mayer, C., Wildi, S., Muehle, B. and Benini, L., 2015, May. On-chip convolutional neural networks. *Proceedings of the 25th edition on Great Lakes Symposium on VLSI*.

[Chai16a] Chai, D., Liu, S., Sun, Z., Wang, J., Wu, C., Chen, Y. and Temam, O., 2016. A dynamically configurable coprocessor for convolutional neural networks. *Proceedings of the 37th International Symposium on Computer Architecture*.

[Chai16b] Chen, Y., Dao, Z., Sun, Z., Wang, J., Wu, C., Chen, Y. and Temam, O., 2014, December. DaianNao: A machine-learning supercomputer. *Proceedings of the 7th Annual International Symposium on Microarchitectures*.

[Chen16] Chen, Y.H., Chen, J., Liu, S., Wang, J., Wu, C., Chen, Y., and Temam, O., 2016. Efficient Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Chen17] Chen, Y., Chen, Y., Xiao, Z., Sun, Z., Zhang, T., Zhou, J., Wang, Y., Wang, Y., and Xie, Y., 2016. DaianNao Family: Efficient Hardware Accelerators for Machine Learning. *ICML Highlights, Communications of the ACM*, 9(9).

[Chih16] Chi, P., Li, S., Gu, P., Xu, C., Zhang, T., Zhou, J., Wang, Y., Wang, Y., and Xie, Y., 2016. PRIME: A Novel Processing-in-Memory Architecture for Neural Network Computation in ReRAM-Based Main Memory. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Choi15] Clark, J. October 26, 2015. Google Using Its Lucrative Web Search Over to AI Machines. *Bloomberg Technology*. www.bloomberg.com.

[Dal16] Dalby, W. February 9, 2016. High Performance Hardware for Machine Learning. *Cadence ESN Summit*.

[Dea13] Dean, J., Ghemawat, S., and Sun, J., L., 2013. The tall and wide. *Communications of the ACM*, 56(2).

[Dea14] Dean, J. July 7, 2014. Large-Scale Distributed Learning with TensorFlow for Building Intelligent Systems. ACM Webinar.

[Dai15] Dai, Z., Fadiher, R., Chen, T., Isene, P., Li, L., Luo, T., Feng, X., Chen, Y. and Temam, O., 2015. June. ShiDanNao: Shifting system design to the user. *Proceedings of the 42nd International Symposium on Computer Architecture*.

[Far09] Farabet, C., Puaud, C., Han, J.Y. and Lecun, Y., 2009. August. Gena: An FPGA-based processor for convolutional networks, 2009 International Conference on Field Programmable Logic and Applications.

[Far11] Farabet, C., Martini, B., Corda, B., Akelrod, P., Culicovski, I. and LeCun, Y., 2011. June. *Neuflow: A runtime reconfigurable dataflow processor for vision*. In *CIPR 2011 Workshop*.

[Gup15] Gupta, S., Agrawal, A., Goldkornik, K. and Narayanan, P., 2015, July. Deep Learning with Limited Numerical Precision. *ICML*.

[Ham90] Hammar, D., 1990. June. A VLSI architecture for high-performance, low-cost, on-chip learning. *1990 IJCNN International Joint Conference on Neural Networks*.

[Han15] Han, S., Pool, J., Tran, J., and Dally, W. 2015. Learning both weights and connections for efficient neural networks. In *Advances in Neural Information Processing Systems*.

[Han16] Han, S., Liu, X., Pool, J., Pfister, A., Han, A., Hosoi, M. and Dally, W., 2016. IEEE. Efficient inference engine on compressed deep neural networks. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[He16] He, K., Zhang, S., Ren, S. and Sun, J., 2016. Identity mappings in deep residual networks. Also in *arXiv preprint arXiv:1603.06027*.

[Hos15] Hosoya, J.I. and Patterson, D.A., 2015. Computer architecture: a quantitative approach, 4th edition, Elsevier.

[Hol99] Holdre, U. and Harrison, L., 2009. *The datascaper: a computer. Morgan and Claypool*.

[Int96] Imiti, P., Corne, T. and Kahn, G., 1996. Special-purpose digital hardware for neural networks: An architectural survey. *Journal of VLSI signal processing systems for technology*, 1(1).

[Int16] Intel, 2016. Intel® Xeon® Processor E3-8600 v5. http://ark.intel.com/products/576601/Intel-Xeon-Processor-E3-8600-v5-15M-Cache-2_16GHz

[Jou16] Joppa, N. May 18, 2016. Google supercharges machine learning tasks with TPU custom chip. <http://cloudplatform.googleblog.com/2016/05/18/tpu.html>

[Keu16] Kratoch, K., 2016. On-chip design of a neural network architecture: a quantitative approach. *Communications of the ACM*, 9(11).

[Kim16] Kim, D., Kang, H.H., Cha, S., Yamahatachi, S. and Makhadmeh, S., 2016. Neurocube: A Programmable Digital Neuromorphic Architecture with High-Density 3D Memory. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Kim17] Kishikawa, A., Saito, E., and Hirono, G., 2017. Efficient classification with deep convolutional neural networks. *Advances in neural information processing systems*.

[Kim20] King, H.T. and Levinson, C.E., 1980. Algorithms for VLSI speech arrays. *Introduction to VLSI systems*.

[Lang16] Lang, J.D., 2009. Identifying shades of gray. *The SPIE conference*, 7213.

[Lar16] Laruffi, M. March 10, 2016. Google Looks To Open Up Streaming Access To Make GPU Programming Easier, Phoronix, http://www.phoronix.com/scan.php?page=news_item&catid=google-Stream-Exec-Parallel.

[LiK16] Li, Kan-Wei, R. Hou, Y., Guo, J., Polansky, M. and Zhong, L., 2016. RedEye: An Analog CVNet Image Sensor Architecture for Continuous Mobile Vision. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Liu15] Liu, D., Dao, Z., Liu, S., Zhao, J., Chen, S., Temam, O., Feng, X., Zhou, X. and Chen, Y., 2015. March. PuDanNao: A polyvalent machine learning accelerator. *Proceedings of the 42nd International Symposium on Computer Architecture*.

[Liu16] Liu, S., Dao, Z., Tao, J.H., Han, D., Luo, T., Xie, Y., Chen, Y. and Chen, T., 2016. Cambridge: An instruction set architecture for neural networks. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Met16] Metz, C., September 26, 2016. Microsoft Bet's Future On A Reprogrammable Computer Chip. *Wired Magazine*, <http://www.wired.com/2016/09/microsoft-beta-future-chip-program-4/>

[Nvi15] Nvidia, January 2015. Tesla K80 GPU Accelerator. *Round Specification* <http://nvidia.com/page/tesla-k80-gpu-features.html#specifications-071316-003-005.pdf>

[Nvi16] Nvidia, 2016. Tesla GPU Accelerators For Servers. <http://www.nvidia.com/object/tesla-servers.html>.

[Ovi15a] Oviichaov, K., Rawate, O., Kim, J.Y., Fowers, J., Strauss, K. and Chung, E.S., February 2, 2015. Accelerating deep convolutional neural networks using specialized hardware. Microsoft Research Whitepaper. <http://www.microsoft.com/en-us/research/publication/accelerating-deep-convolutional-neural-networks-using-specialized-hardware/>

[Ovi15b] Oviichaov, K., Rawate, O., Kim, J.Y., Fowers, J., Strauss, K. and Chung, E.S., 2015, August. Toward accelerating deep learning at scale using specialized hardware in the datacenter. *2015 IEEE Hot Chips 27 Symposium (HOTS)*.

[Put14] Puttenber, D.A., 2004. Latency lags handbook. *Communications of the ACM*, 47(10).

[Pee13] Peeren, M., Setis, A.A., Mesnam, B. and Corporal, H., 2013, October. Memory-centric design for convolutional neural networks. In *2013 IEEE/ACM International Conference on Computer Design (ICCD)*.

[Put14] Patnam, A., Caulfield, A.M., Chang, E.S., Chio, D., Constantinescu, K., Demme, J., Farnazizadeh, H., Fowers, J., Gopal, G.P., Gray, J., Haselman, M., Hauck, S., Heil, S., Hornsby, A., Kim, J.Y., Lanka, S., Laria, J., Peterson, E., Pope, S., Smith, A., Thong, J., Xiao, P.Y., Burger, D., 2014. June. A reconfigurable fabric for accelerating large-scale datacenter services. *41st International Symposium on Computer Architecture*.

[Put15] Patnam, A., Caulfield, A.M., Chang, E.S., Chio, D., Constantinescu, K., Demme, J., Farnazizadeh, H., Fowers, J., Gopal, G.P., Gray, J., Haselman, M., Hauck, S., Heil, S., Hornsby, A., Kim, J.Y., Lanka, S., Laria, J., Peterson, E., Pope, S., Smith, A., Thong, J., Xiao, P.Y., Burger, D., 2015. A Reconfigurable Fabric for Accelerating Large-Scale Datacenter Services. *IEEE Micro*, 35(3).

[Put16] Patnam, A., Caulfield, A.M., Chang, E.S., Chio, D., Constantinescu, K., Demme, J., Farnazizadeh, H., Fowers, J., Gopal, G.P., Gray, J., Haselman, M., Hauck, S., Heil, S., Hornsby, A., Kim, J.Y., Lanka, S., Laria, J., Peterson, E., Pope, S., Smith, A., Thong, J., Xiao, P.Y., Burger, D., 2016. A Reconfigurable Fabric for Accelerating Large-Scale Datacenter Services. *Communications of the ACM*.

[Qad13] Qader, W., Hamed, R., Shachan, O., Venkatesan, P., Kozyrakis, C. and Horowitz, M.A., 2013, June. Convolution engine: balancing efficiency and flexibility in specialized computing. *Proceedings of the 40th International Symposium on Computer Architecture*.

[Rau11] Ramesher, U., Reichter, J., Raab, W., Andlas, J., Bruehl, N., Hachmann, U. and Westering, M., 1991. Design of a 1st Generation Neurocomputer. In *VLSI Design of Neural Networks*. Springer U.S.

[Real16] Reagan, B., Whittington, P., Adfoll, R., Rama, S., Lee, H. S., Hernandez-Lobato, J.M., Wei, G.Y. and Brooks, D., 2016. Minerva: Enabling low-power, highly-accurate deep neural network acceleration. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Ros15a] Ross, J., Joppa, N., Pheasant, A., Young, G., Norman, A., Thorogood, G., Liu, D., 2015. Neural Network Processor. Patent Application No. 62/164,931.

[Ros15] Ross, J., Pheasant, A., 2015. Computing Convolutions Using a Neural Network Processor. Patent Application No. 62/164,902.

[Ros16] Ross, J., 2015. Prefetching Weights for a Neural Network Processor. Patent Application No. 62/164,901.

[Ros17] Ross, J., Thomson, G., 2015. Rotating Data for Neural Network Computations. Patent Application No. 62/164,908.

[Ros18] Roushkovsky, O., Deng, J., Su, H., Krause, S., Saftschek, S., M. S., Huang, Z., Kapurthy, A., Khosla, A., Bernstein, M. and Berg, A.C., 2015. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(1).

[Sch16] Schumacher, E. and Brubaker, J., 2009, June. The user and business impact of server delays, additional bytes, and HTTP chunking in web search. In *Web Performance and Operations Conference*.

[Sha16] Shafer, A., Nag, A., Marulamban, N., Balasubramanian, R., Strachan, J.P., Hu, M., Williams, R.S. and Srikanar, V., 2016. ISAAC: A Convolutional Neural Network Accelerator with In-Situ Analog Arithmetic in Crossbars. *Proceedings of the 43rd International Symposium on Computer Architecture*.

[Sil16] Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneerselvam, V., Lanctot, M. and Dieleman, S., 2016. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587).

[Smith16] Smith, J.E., 1982, April. Decoupled accelerators for convolutional architectures. *Proceedings of the 11th International Symposium on Computer Architecture*.

[Ste15] Steinberg, D., 2015. Full-Chip Simulations, Keys to Success. *Proceedings of the Synopsys User Group (SUG) Silicon Valley 2015*.

[Ste15a] Steingy, C., Liu, W., Liu, Y., Serment, P., Reed, S., Angulo, L., Fernh, D., Vanhucque, V. and Rabinovich, A., 2015. Going deep with deep learning: an end-to-end system for image classification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

[Tho15] Thorson, G., Clark, C., Liu, D., 2015. Vector Computation Unit in a Neural Network Processor. Patent Application No. 62/165,022.

[Wol10] Williams, S., Waterman, A. and Patterson, D., 2009. RedEye: an insightful visual processor module for multicore architectures. *Communications of the ACM*.

TPUv1 succeeded because of

Conclusions (1/2)

- Large matrix multiply unit
- Substantial software-controlled on-chip memory
- Run whole inference models to reduce host CPU
- Single-threaded, deterministic execution model
good match to 99th-percentile response time
- Enough flexibility to match NNs of 2017 vs. 2013
- Omission of GP features \Rightarrow small, low power die
- Use of 8-bit integers in the quantized apps
- Apps in TensorFlow, so easy to port at speed

Conclusions (2/2)

- Inference prefers latency over throughput
- K80 GPU relatively poor at inference (vs. training)
- Small redesign improves TPUv1 at low cost
- 15-month design & live on I/O bus yet TPUv1 15X-30X faster Haswell CPU, K80 GPU (inference), <math>< \frac{1}{2}</math> die size, $\frac{1}{2}$ Watts
 - 65,536 (8-bit) TPUv1 MACs cheaper, lower energy, & faster 576 (32-bit) CPU MACs, 2496 GPU (32-bit) MACs
- 10X difference in computer products are rare

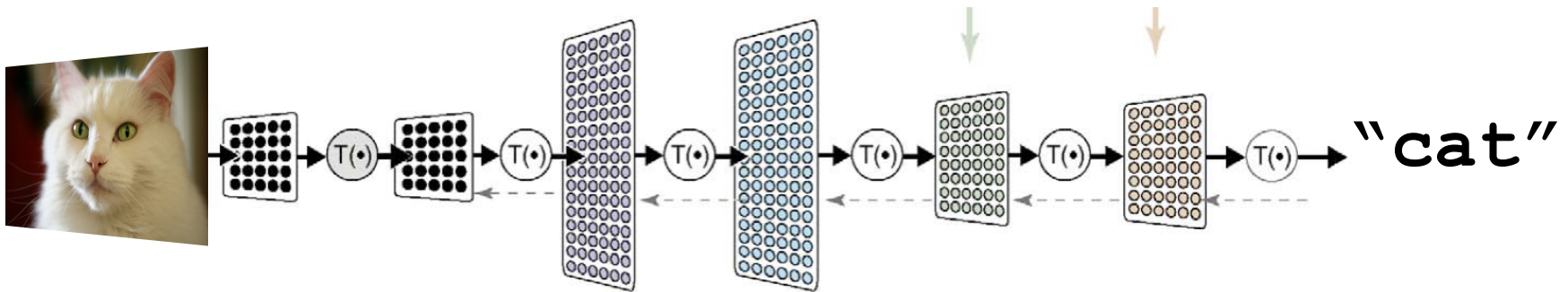
Two Talks

In-Data Center Performance Analysis of a Tensor Processing Unit
ISCA 2017, Jouppi, Young, Patil, Patterson, et al.
(TPUv1, for inference)

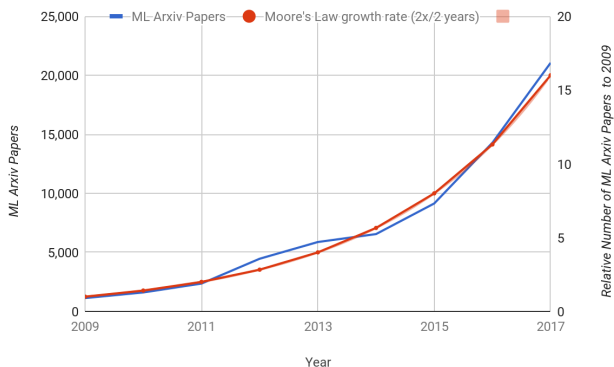
Codesign in Google TPUs
(TPUv2 and later, for training and inference)

Observation: Training >> Inference

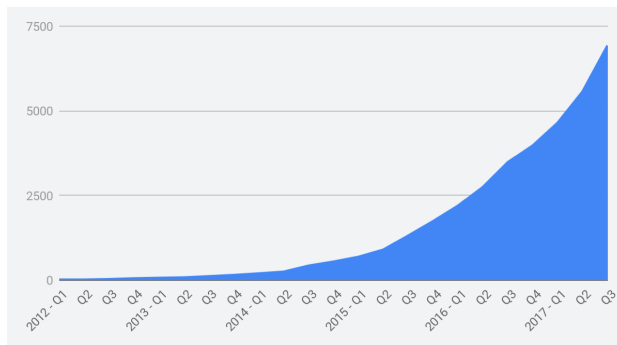
- 3x** the computation: forward propagation, backward propagation, and weight update.
- Much longer data storage lifetimes: **memory** capacity and bandwidth.
- Huge **training datasets** for training, versus scale-out to serve inference.
- Changes to algorithms and model structure require more **flexibility**.
- Many more potential **Amdahl's Law** bottlenecks.



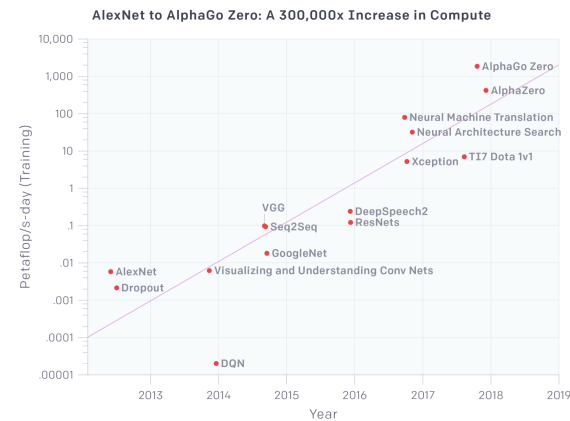
Exponential Growth in Deep Learning



ArXiv papers about ML
~18 months



Google project directories
~18 months



FLOPs to train a model
~3.5 months
(~10X per year)

Why all this growth? Because Deep Learning works.

Classic Codesign at the HW/SW Interface



Definition: design spanning two fields for a common goal.

- Classic version is between architecture and compiler.

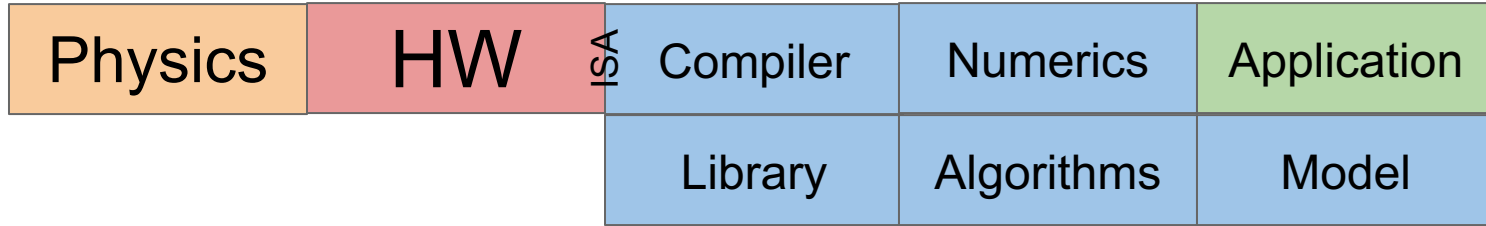
Instruction Set Architecture (ISA) as interface/contract between levels.

Example of pushing things back and forth: instruction scheduling.

- VLIW (static scheduling)
- OoO (dynamic scheduling)
- Answer today=both.

Ultimately ISA is a **single** thin layer between the hardware and software domains.

Codesign for Domain-Specific Architectures



(conceptual, not rigorous diagram)

Now, there are **many** different layers, with **many** different interfaces.

TPUs are still digital (for now).

- Some startups are pushing into physics (NVRAM, Flash, optical).
- Need to do codesign from physics to application: hard!

Fallacy: TPUs are ASICs, so they are not Programmable

ASIC: **A**pplication-**S**pecific **I**ntegrated **C**ircuit

Means only “build whatever you want into the chip.”

ASICs include general-purpose cores, SoCs, and fixed-function designs.

TPUs are Domain-Specific Architectures (DSAs) for Machine Learning.

We designed them to meet our current and future needs.

They include the flexibility to handle future models.

Choosing the **right** amount of flexibility is central to our codesign process.

For the technically nitpicky:

TPUv1 is a coprocessor, controlled by the host.

TPUv2 and successors are Turing-complete.

TPUs power both Google **research** and Google **production** applications.

Training: DNN Supercomputer or Cluster of CPUs with DNN Accelerators?

- Single-chip system—built as coprocessor to a CPU like TPUv1—would work fine for inference and standard cluster networks
 - AlphaGo used cluster of 64 TPUv1 chips
- Went instead with large supercomputer because
 - Training takes weeks to months on single chip for our production training runs
 - Deep neural network wisdom was bigger datasets + bigger machines led to breakthroughs
- Build a NN **supercomputer** (TPU v2/v3) vs build a NN **coprocessor chip** (TPU v1)

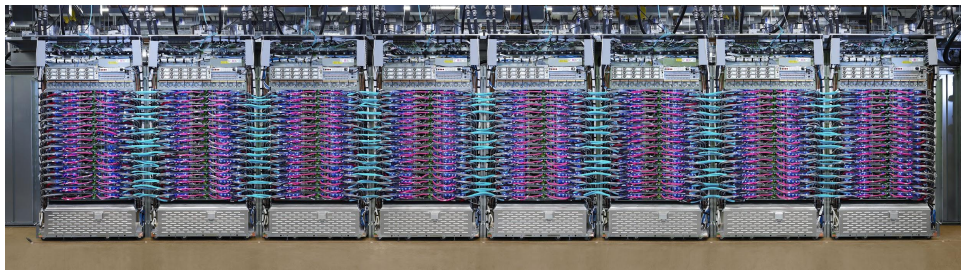
Codesign: Datacenter-scale Supercomputer

Three ways to spend on computer architectural resources:

- Compute
- Memory
- Interconnect

Cloud TPUs were designed from the beginning to be networked.

- Data parallelism through increased batch size scales seamlessly today.
- Model parallelism is underway.



ML Quality \approx Correctness: Fast but Incorrect Uninteresting

- 1% quality loss to ML practitioners can be like getting wrong answer
 - Aiming for intelligent app for a billion people, so lower quality can mean worse experience for millions of people / loss of income
- For datacenter production apps, training has to be in floating point
 - Researchers exploring fixed point for training but at cost in quality
 - Production remains floating point (but FP32 sufficient, no need for FP64)

Codesign: Reduced-precision numerics

Neural-network researchers showed that inference did not need float32 precision.

Vanhoucke, V., Senior, A. and Mao, M.Z., 2011, December. Improving the speed of neural networks on CPUs.
In *Proc. Deep Learning and Unsupervised Feature Learning NIPS Workshop* (Vol. 1, p. 4).

TPUv1 deployed with int8 hardware and support for int16 through software.
int16 intended as “insurance”; used mostly in LSTM-based models.

Training has traditionally been done in floating point.

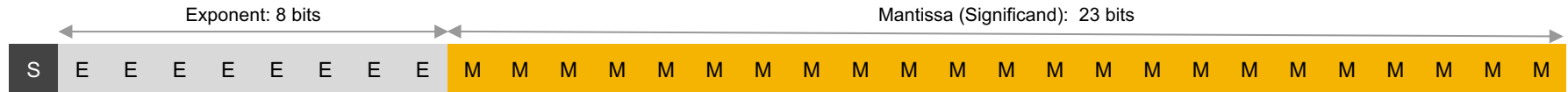
Q: Does training require full float32 resolution?

A: Not everywhere. But where matters.

A brief guide to Floating Point Formats

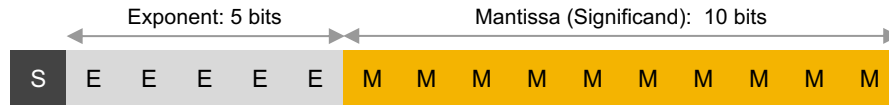
fp32: Single-precision IEEE Floating Point Format

Range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



fp16: Half-precision IEEE Floating Point Format

Range: $\sim 5.96e^{-8}$ to 65504



bfloat16: Brain Floating Point Format

Range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



Bfloat16 as Good Codesign

Hardware: small mantissa reduces multiplier power, area

$$\text{float32: } 23^2=529$$

$$\text{float16: } 10^2=100$$

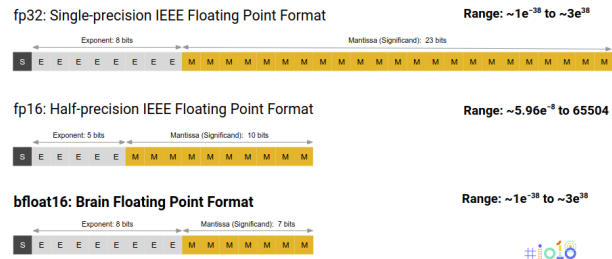
$$\text{bfloat16: } 7^2=49$$

Software: same dynamic range on number line, same Inf/NaN behavior as float

Numerics: Unlike IEEE fp16, bfloat16 trains without loss scaling [Micikevicius 2017]

System: bfloat16 as an implementation technique inside the matrix multiplier.

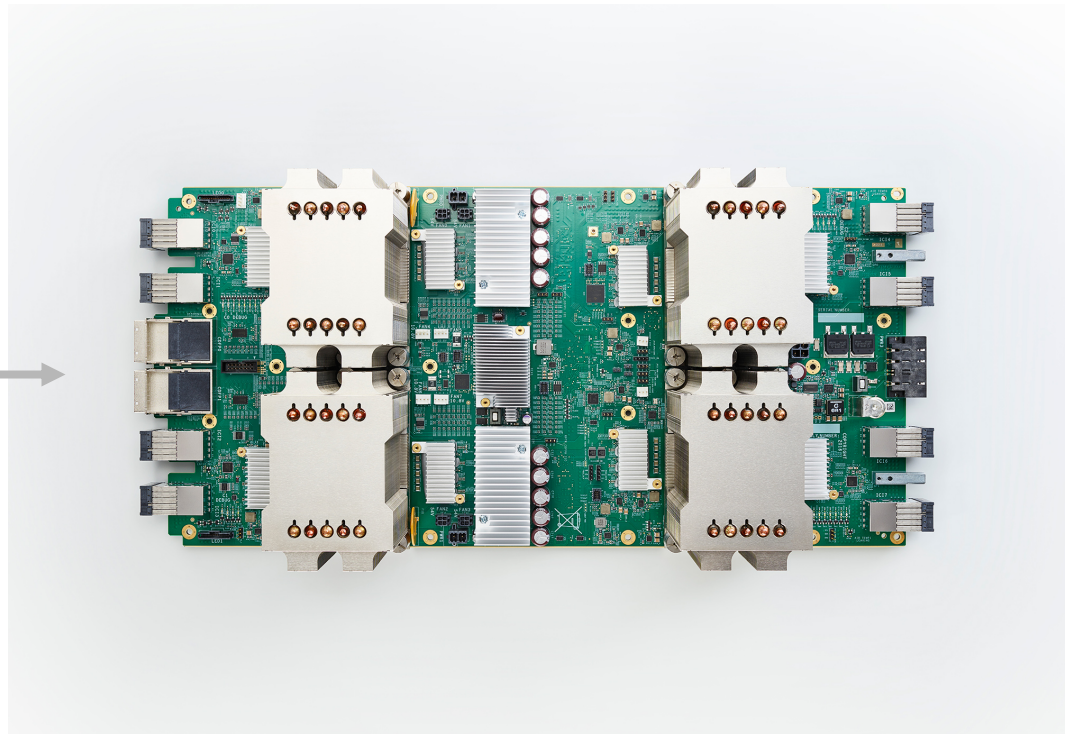
Can also expose it to save memory capacity and bandwidth, with more work



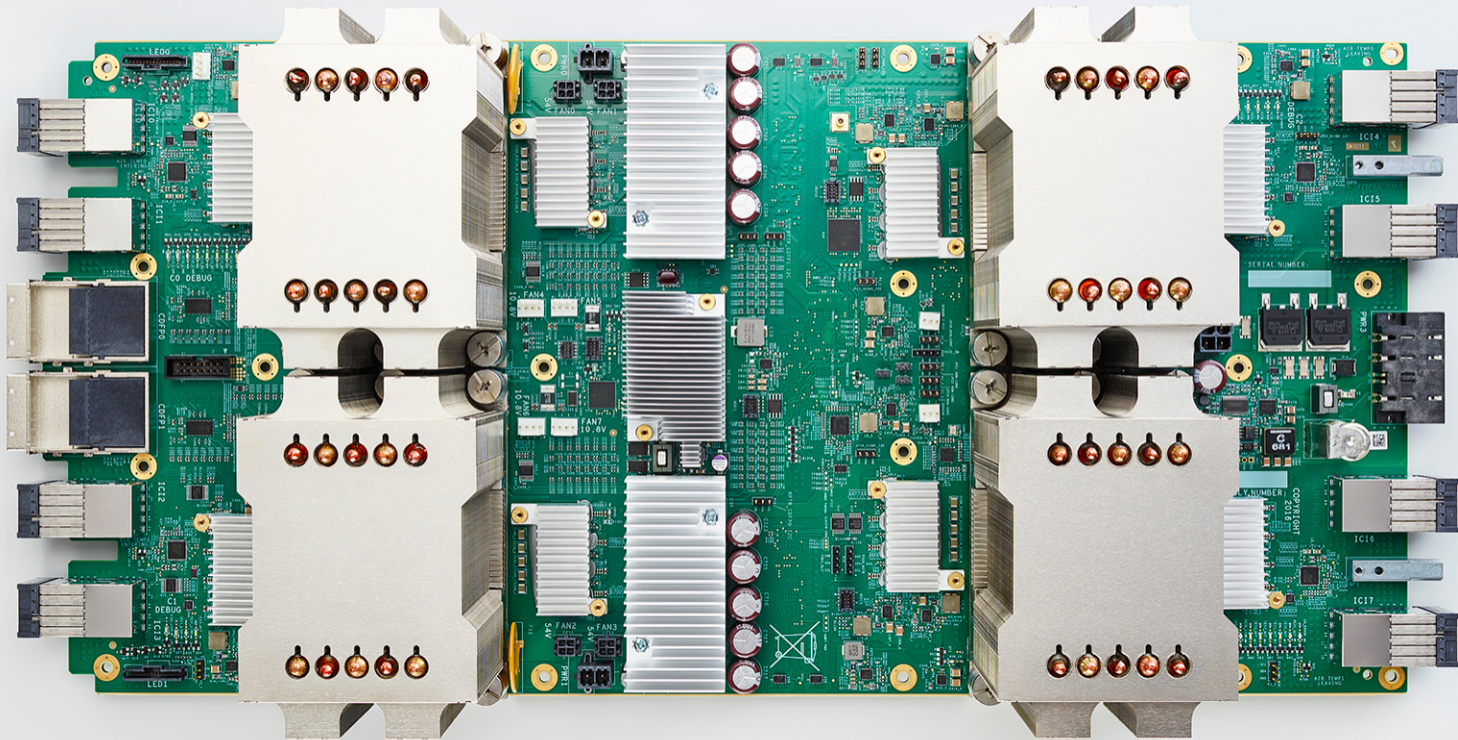


Host Server

PCI v3 x 32

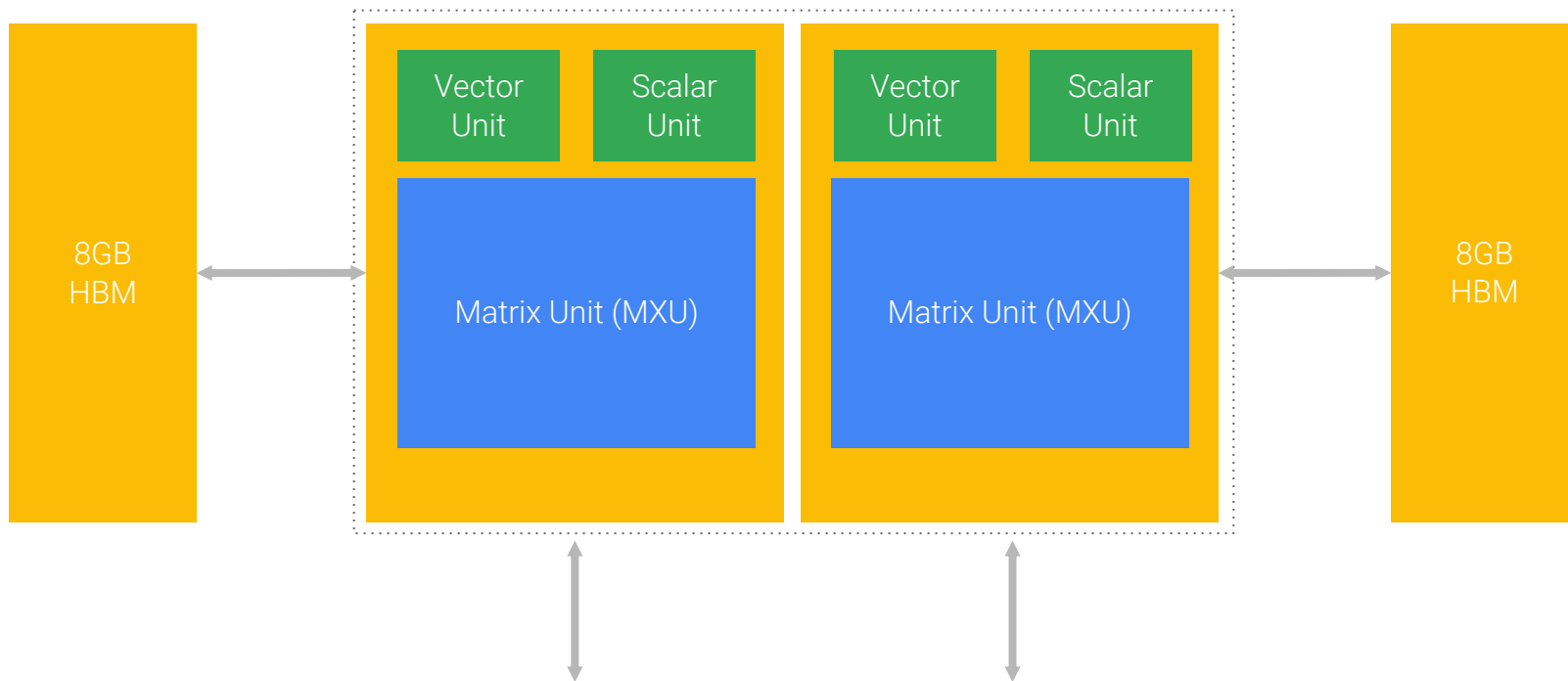


Cloud TPU (v2), Generally Available

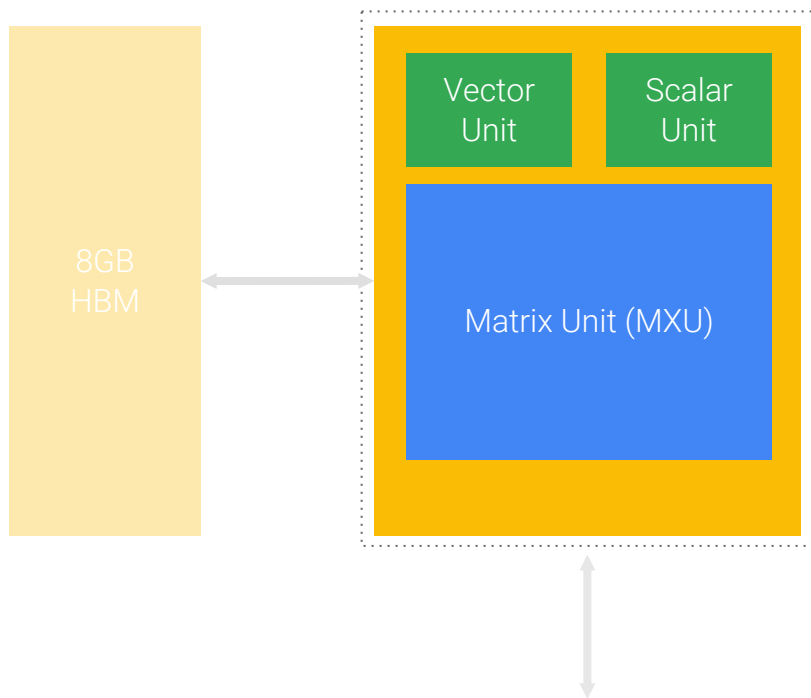


- 180 teraflops of computation, 64 GB of HBM memory, 2400 GB/s mem. BW
- Designed to be connected together into larger configurations

Cloud TPU chip layout



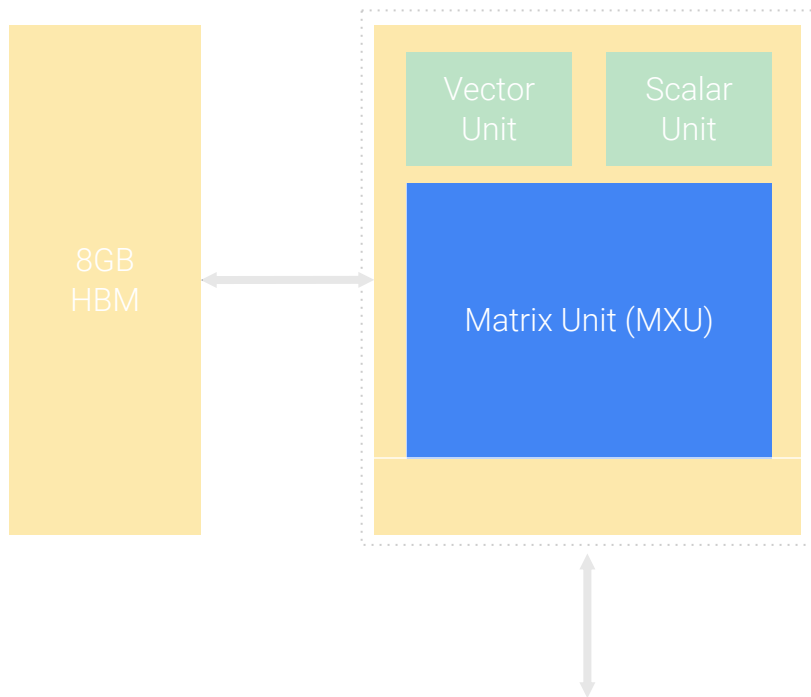
Cloud TPU chip layout



22.5 TFLOPS per core

- 2 cores per chip
- 4 chips per 180 TFLOP Cloud TPU
- Scalar Unit
- Vector Unit
- Matrix Unit
- Mostly float32

Cloud TPU chip layout

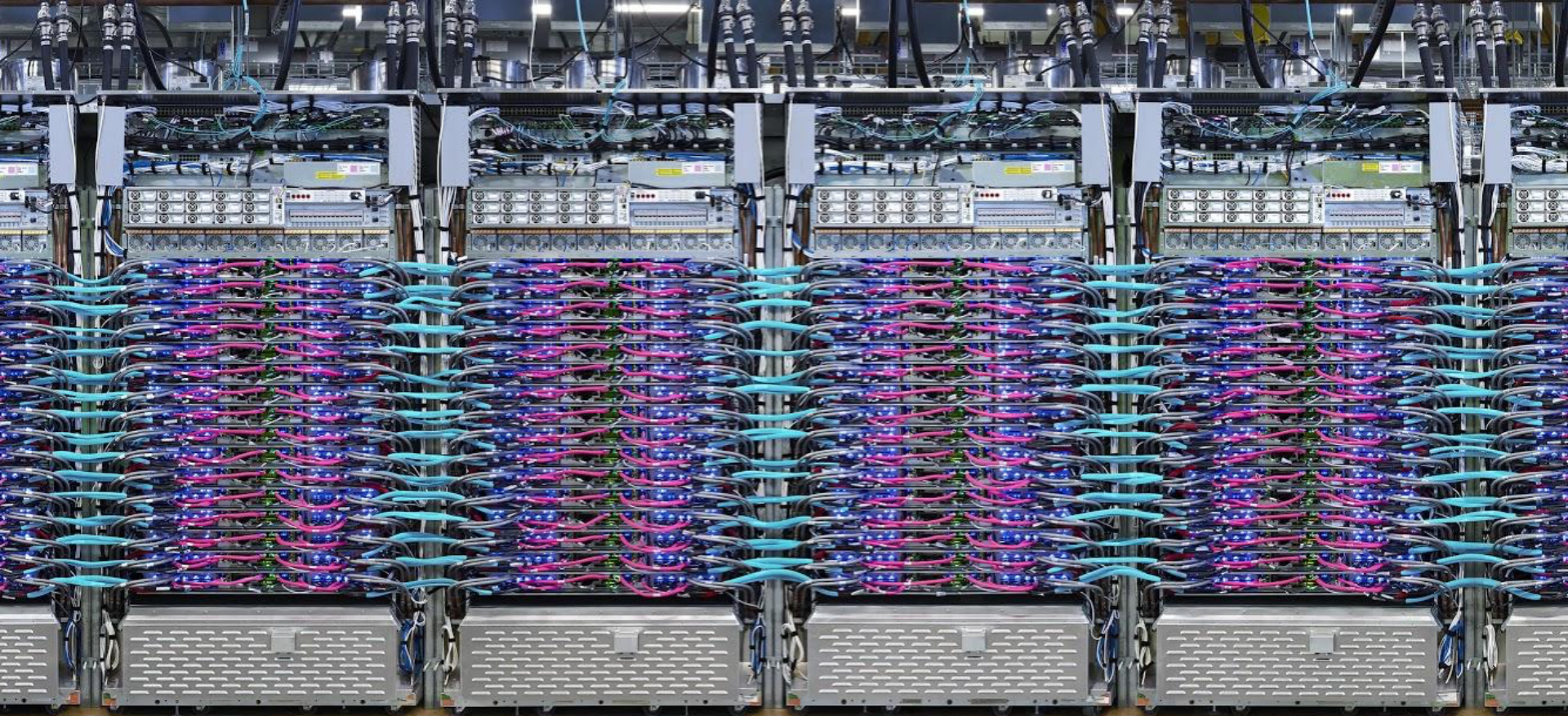


Matrix Unit (MXU)

- 128 x 128 systolic array
- bfloat16 multiplies
- float32 accumulate



Cloud TPU v2 Pod



TPU v3 Pod: Revealed at Google I/O May 2017

Relentless progress

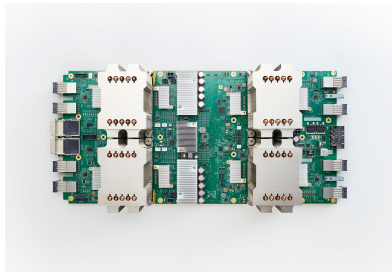
g.co/cloudtpu

TPU v1
(deployed 2015)



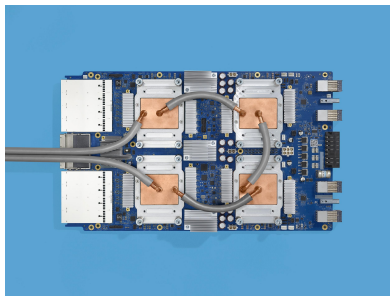
92 teraops
Inference only

Cloud TPUv2



180 teraflops
64 GB HBM
Training and inference
Generally available (GA)

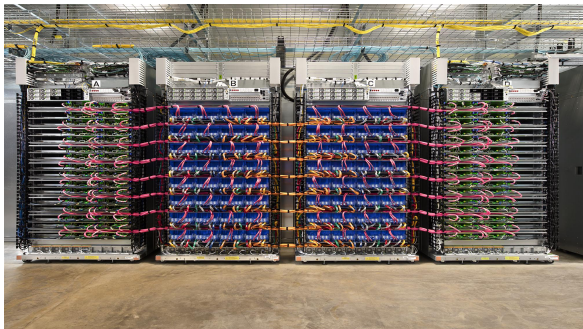
Cloud TPUv3



420 teraflops
128 GB HBM
Training and inference
Beta

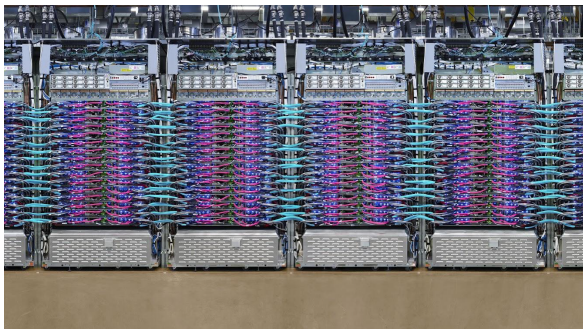
Relentless progress

g.co/cloudtpu



Cloud TPU Pod (v2, 2017)

11.5 petaflops
4 TB HBM
2-D toroidal mesh network
Training and inference
Alpha

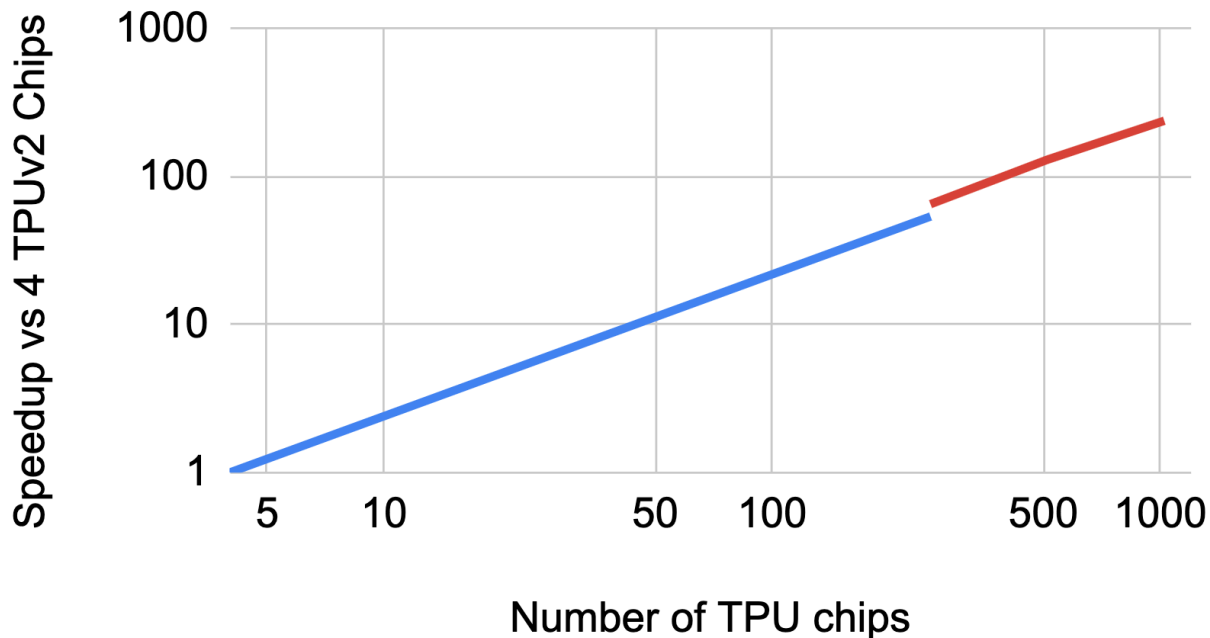


TPU v3 Pod (2018)

> 100 petaflops!
32 TB HBM
Liquid cooled
New chip architecture + larger-scale system

ResNet-50 Speedup

— TPUv2 — TPUv3



Ying, C., Kumar, S., Chen, D., Wang, T. and Cheng, Y., 2018. Image Classification at Supercomputer Scale. *arXiv preprint arXiv:1811.06992*.

ML performance, carefully measured

Focus on real-world data, time-to-accuracy, and cost

- Measure performance with **real input data**, not just synthetic
- Ensure that models **converge to expected accuracy** while achieving high performance
- Measure **total cost** via public clouds

Make ML benchmarks reproducible via open-source implementations

JUST LAUNCHED!



MLPerf

[mlperf.org]

A broad ML benchmark suite for measuring performance of ML frameworks, ML hardware accelerators, and ML cloud platforms

Researchers from Harvard, Stanford, University of California, Berkeley, University of Illinois, University of Minnesota, University of Texas, and University of Toronto

Supporting companies:

Alibaba, AMD, Arm, Baidu, Cadence, Cerebras, Cisco, Cray, Dividiti, Enflame Tech, Esperanto, Facebook, Google, Groq, Huawei, Intel, MediaTek, Mentor Graphics, Microsoft, Myrtle, Mythic, NetApp, NVIDIA, One Convergence, Qualcomm, Rpa2ai, Sambanova, Samsung S.LSI, Sigopt, Synopsys, Tensyr, Teradyne, Wave Computing, WekaIO

MLPerf (mlperf.org) in One Slide

Goal: Build “SPEC for Machine Learning”.

Consortium of companies and universities.

Philosophy:

- Agile development because ML is changing rapidly.
- Serve both the commercial and research communities.
- Enforce replicability to ensure reliable results.
- Use representative workloads, reflecting production use-cases.
- Keep benchmarking effort affordable (so all can play).

v0.5 Published December 2019

MLPerf 0.5 Training Results

<https://mlperf.org/results/> has the raw data and detailed submission information.

Of the 7 benchmarks:

- NVIDIA submitted 6/7 V100-based results
- Google submitted 3/7 TPUv2 and TPUv3 results
- Intel submitted 1/7 SkyLake (CPU) results

Some observations:

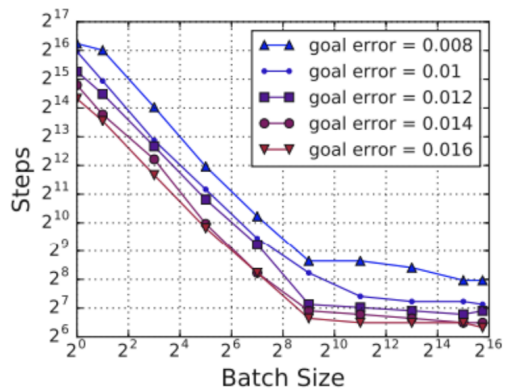
- What do the submissions say about generality and software maturity?
- For similar numbers of chips, V100 and TPUv3 look comparable.
- Both NVIDIA and Google showed huge scale (640- and 260-chip entries)

Next deadline May 2019 MLPerf 0.6

Some open codesign questions in Machine Learning

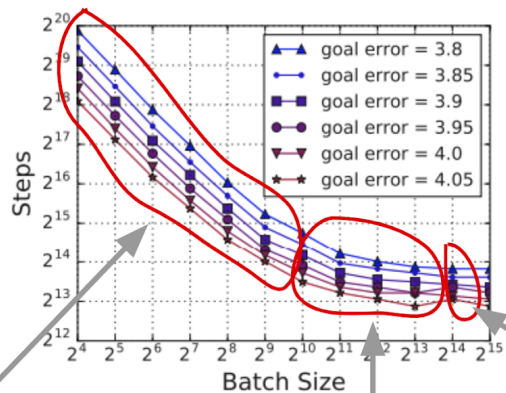
- What's the "best" architecture?
 - Will the market be the final arbiter?
 - At the end of Moore's Law, perhaps architectural efficiency matters more
- Software may matter more than hardware:
 - MultiFlow's Compiler as most important artifact.
 - Ease of use takes time: typically a decade for compilers to mature.
- What's the lower limit on numerics?
- How much more is sparsity going to matter?
 - Embeddings, attention, compute and memory savings. What else? Brains are sparse.
- When does batch=1 matter? Definitely for inference. For training?
 - Shallue, C.J., Lee, J., Antognini, J., Sohl-Dickstein, J., Frostig, R. and Dahl, G.E., 2018. Measuring the effects of data parallelism on neural network training. arXiv preprint arXiv:1811.03600.
- How can we use more weights, but touch fewer of them? Mixture of Experts.

The effects of data parallelism on neural network training



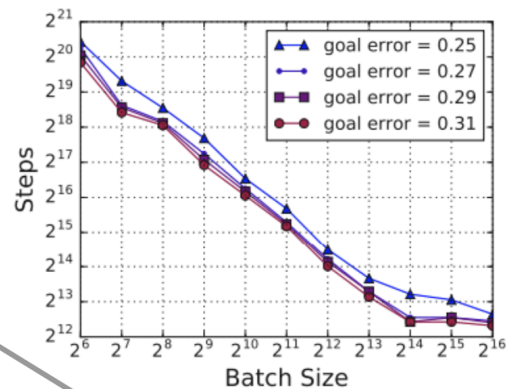
(a) Simple CNN on MNIST

Perfect scaling region



(b) Transformer on LM1B

Diminishing returns



(c) ResNet-50 on ImageNet

Maximal data parallelism

Measuring the Effects of Data Parallelism on Neural Network Training, Christopher J. Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig, George E. Dahl, arxiv.org/abs/1811.03600

ML Crisis as both Danger and Opportunity

Danger: the end of Moore's Law, Dennard Scaling, and standard CPU performance

- Limits of CMOS in sight
- Intel 10nm woes, Global Foundries 7nm exit

Opportunity: the revolution in ML

- Economic demand for ML accelerators
- Architectural and codesign experimentation and transformation
- Can we use ML to design better accelerators?

Irony: exponential demand for ML computation, just at the end of Moore's Law

- Efficiency is going to matter a lot
- Huge opportunities for HW/SW codesign in building TPUs and other DSAs

A New Golden Age

Hennessy and Patterson,
“A New Golden Age for Computer
Architecture,”

Communications of the ACM,
February 2019



Questions?

*4/5/17 Google published a blog on the TPU. A 17-page technical paper with same title will be on arXiv.org. (Paper will also appear at the *International Symposium on Computer Architecture* on June 26, 2017.)

<https://cloudplatform.googleblog.com/2017/04/quantifying-the-performance-of-the-TPU-our-first-machine-learning-chip.html>

My Story: Accidental Berkeley CS Professor*

- 1st college graduate in family; no CS/grad school plan
 - Wrestler, Math major in high school and college
- Accidental CS Undergrad
- Accidental PhD Student
 - New UCLA PhD (Jean-Loup Baer) took pity on me as undergrad
- Wife + 2 sons in Married Students Housing on 20 hour/week RAship
 - Lost RA-ship after ≈4 years because grant ended
 - Part time at nearby Hughes Aircraft Company ≈3 more years (7.5 years to PhD)
- Accidental Berkeley Professor
 - Wife forced me to call UC Berkeley CS Chair to check on application
- 1st project as Assistant Prof with an Associate Prof too ambitious & no resources
 - Took leave of absence at Boston computer company to rethink career; 3rd year Ass't Prof
- Tenure not easy (Conference papers vs journal papers, RISC too recent)
Full video: see "Closing Remarks", www2.eecs.berkeley.edu/patterson/ 2016

What Worked Well for Me*

- Maximize Personal Happiness vs. Personal Wealth
- Family First!
- Passion & Courage
 - Swing for the fences vs. Bunt for singles
- “Friends may come and go, but enemies accumulate”
- Winning as Team vs. Winning as Individual
 - “No losers on a winning team, no winners on a losing team”
- Seek Out Honest Feedback & Learn From It
 - Guaranteed Danger Sign: “I’m smartest person in the room”
- One (Big) Thing at a Time
 - “It’s not how many projects you start; It’s how many you finish”
- Natural Born Optimist

* Full video: see “Closing Remarks”, www2.eecs.berkeley.edu/patterson2016

9 Magic Words for a Long Relationship

“I Was Wrong.”

“You Were Right.”

“I Love You.”