













Example Domain: Deep Neural Networks Inpired by neuron of Cutput Innut Neuron the brain Computes non-linear "activiation" function of the weighted sum of input values Neurons arranged in layers Name **DNN** layers Weights **Operations/Weight** MLPO 5 20M 200 MLP1 4 5M 168 LSTM0 52M 58 64 LSTM1 56 34M 96 CNN0 16 8M 2888 CNN1 89 100M 1750 Figure 7.5 Six DNN applications that represent 95% of DNN workloads for inference at Google in 2016, which we use in Section 7.9. The columns are the DNN name, the number of layers in the DNN, the number of weights, and operations per weight (operational intensity). Figure 7.41 on page 595 goes into more detail on these DNNs.

xample I	e Domain: Deep Neural Networks						
■ Most – To	practioners will choose an existing design ology and Data type						
■ Train – Ca – Su	ing (learning culate weights pervised learn	(learning): ate weights using backpropagation algorithm vised learning: stocastic gradient descent					
Type of dat	a Problem area	Size of benchmark's training set	DNN architecture	Hardware	Training time		
text [1]	Word prediction (word2vec)	100 billion words (Wikipedia)	2-layer skip gram	1 NVIDIA Titan X GPU	6.2 hours		
audio [2]	Speech recognition	2000 hours (Fisher Corpus)	11-layer RNN	1 NVIDIA K1200 GPU	3.5 days		
images [3	Image classification	1 million images (ImageNet)	22-layer CNN	1 NVIDIA K20 GPU	3 weeks		
video [4]	activity recognition	1 million videos (Sports-1M)	8-layer CNN	10 NVIDIA GPUs	1 month		
Figure 7.6 T	raining set sizes and trai	ning time for several D	NNs (landola, 20	16).			

Inferrence: use neural network for classification









2) Convolutional Neural Network (CNN)

- Batches:
 - Reuse weights once fetched from memory across multiple inputs
 - Increases operational intensity
- Quantization
 - Use 8- or 16-bit fixed point
- Summary:
 - Need the following kernels:
 - Matrix-vector multiply
 - Matrix-matrix multiply
 - Stencil
 - ReLU
 - Sigmoid
 - Hyperbolic tangent











Guidelines for DSAs

Guideline	TPU	Catapult	Crest	Pixel Visual Core
Design target	Data center ASIC	Data center FPGA	Data center ASIC	PMD ASIC/SOC IP
1. Dedicated memories	24 MiB Unified Buffer, 4 MiB Accumulators	Varies	N.A.	Per core: 128 KiB line buffer, 64 KiB P.E. memory
2. Larger arithmetic unit	65,536 Multiply- accumulators	Varies	N.A.	Per core: 256 Multiply- accumulators (512 ALUs)
3. Easy parallelism	Single-threaded, SIMD, in-order	SIMD, MISD	N.A.	MPMD, SIMD, VLIW
 Smaller data size 	8-Bit, 16-bit integer	8-Bit, 16-bit integer 32-bit Fl. Pt.	21-bit Fl. Pt.	8-bit, 16-bit, 32-bit integer
 Domain- specific lang. 	TensorFlow	Verilog	TensorFlow	Halide/TensorFlow

Examples of DSAs

Tensor Processing Unit
 Microsoft Catapult
 Intel Crest
 Pixel Visual Core









TPU Implementation TPU chip fabricated using the 28-nm process, 700 MHz clock. - Less than half size of an Intel Haswell CPU, which is 662 mm². Local Unified Buffer for Matrix multiply unit (256x256x8b = 64K MAC) activations (96Kx256x8b = 24 MiB) 24% 29% of chip Host Accumulators DRAM RAM Interf. 2% (4Kx256x32b = 4 MiB) 6% Control 2% Activation pipeline 6% port ddr3 3% Figure 7.16 TPU printed circuit board. It can be inserted into the slot for an SATA disk in a server, but the card uses the PCIe bus. PCle Misc. I/O 1% Interface 3% Figure 7.15 Floor plan of TPU die. The shading follows Figure 7.14. The light data buffers are 37%, the light computation units are 30%, the medium I/O is 10%, and the dark control is just 2% of the die. Control is much larger (and much more difficult to design) in a CPU or GPU. The unused white space is a consequence of the emphasis on time to tape-out for the TPU. 10



TPU and the 5 Guidelines

- Use dedicated memories
 - 24 MiB dedicated buffer, 4 MiB accumulator buffers
- Invest resources in arithmetic units and dedicated memories
 - 60% of the memory and 250X the arithmetic units of a server-class CPU
- Use the easiest form of parallelism that matches the domain
 - Exploits 2D SIMD parallelism
- Reduce the data size and type needed for the domain
 - Primarily uses 8-bit integers
- Use a domain-specific programming language
 - Uses TensorFlow







































47

Tensor Processing Unit (TPU) – V1, ..., V5, ...

	TPUv1	TPUv2	TPUv3	TPUv4[14][16]	TPUv5 ^[17]	Edge v
Date introduced	2016	2017	2018	2021	2023	2018
Process node	28 nm	16 nm	16 nm	7 nm	Unstated	
Die size (mm ²)	331	< 625	< 700	< 400	Unstated	
On-chip memory (MiB)	28	32	32	32	48	
Clock speed (MHz)	700	700	940	1050	Unstated	
Memory	8 GIB DDR3	16 GIB HBM	32 GIB HBM	32 GIB HBM	16 GB HBM	
Memory bandwidth	34 GB/s	600 GB/s	900 GB/s	1200 GB/s	819 GB/s	
TDP (W)	75	280	220	170	Not Listed	2
TOPS (Tera Operations Per Second)	92	45	123	275	393	4
TOPS/W	0.31	0.16	0.56	1.62	Not Listed	2

Chip feature	Cloud TPU v3	Cloud TPU v4		
Peak compute per chip	123 teraflops (bf16)	275 teraflops (bf16 or int8)		
HBM2 capacity and bandwidth	32 GiB, 900 GB/s	32 GiB, 1200 GB/s		
Measured min/mean/max power	123/220/262 W	90/170/192 W		
TPU pod size	1024 chips	4096 chips		
Interconnect topology	2D torus	3D torus		
Peak compute per pod	126 petaflops (bf16)	1.1 exaflops (bf16 or int8)		
All-reduce bandwidth per pod	340 TB/s	1.1 PB/s		
Bisection bandwidth per pod	6.4 TB/s	24 TB/s		

Cloud TPU v4 pods performance table (source: Google)

https://en.wikipedia.org/wiki/Tensor_Processing_Unit https://cloud.google.com/tpu

https://cloud.google.com/blog/topics/systems/tpu-v4-enables-performance-energy-and-co2e-efficiency-gains







Conclusion: A New Golden Age turing lecture End of Dennard Scaling and Moore's Law > architecture innovation to improve performance/cost/energy • Security \Rightarrow architecture innovation too Domain Specific Languages ⇒Domain Specific A New Golden Architectures Age for Computer • Free, open architectures and open-source implementations Architecture everyone can innovate and contribute • Cloud FPGAs \Rightarrow all can design and deploy custom "HW" • Agile HW development \Rightarrow all can afford to make (small) chips • Like 1980s, great time for architects in academia & in industry!

