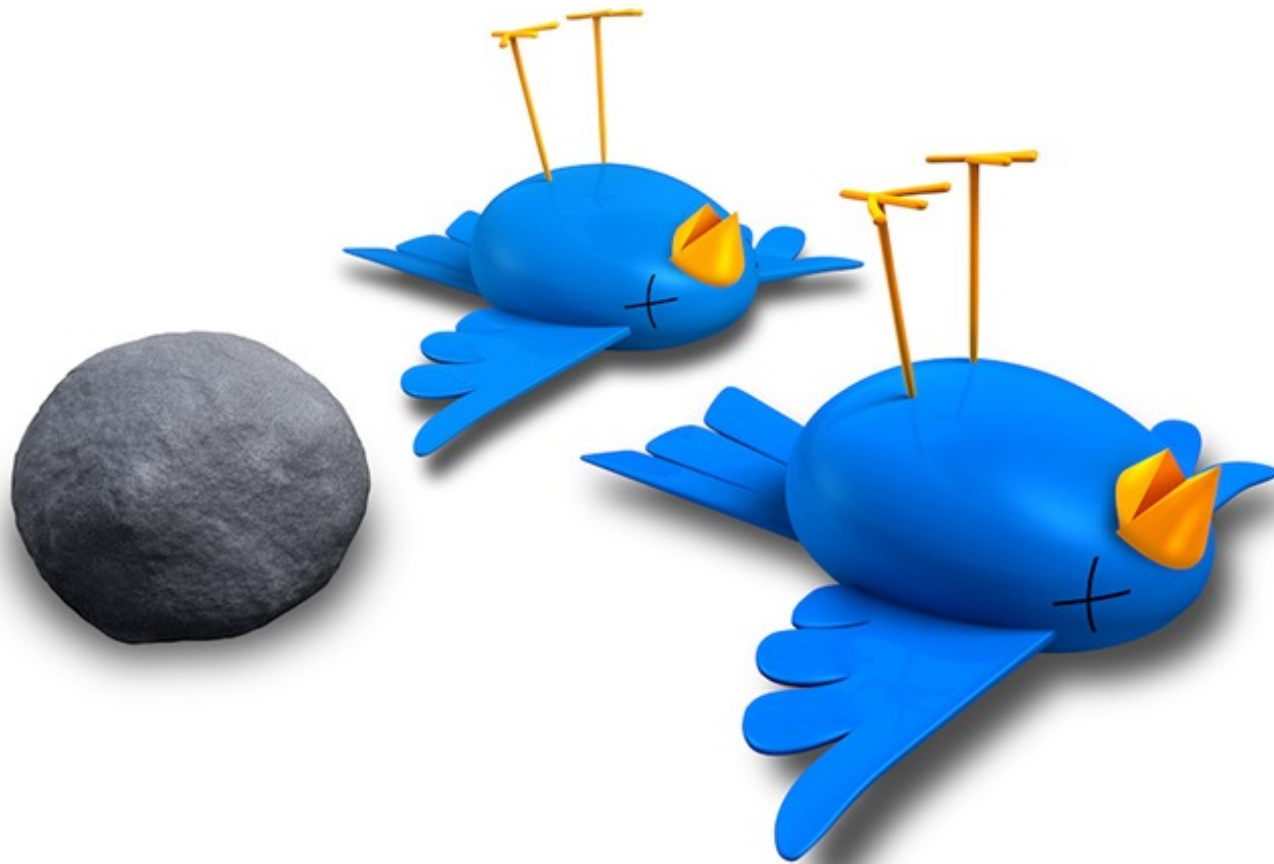


# EPFL CS723

## Ten Lessons From Three Generations Shaped Google's TPUv4i

# Elevator speech



# What is the problem? How important is it?

- **Growing demand** for AI services with ML workloads that become more **compute intensive**.
- ML workloads are **expensive** (infrastructure + energy), especially at **inference** to serve large-scale AI applications
- Need for **specialized hardware** for ML that is:
  - **Cost effective** and **low-latency**
  - **Flexible** and **scalable**
- Important problem since specialized hardware is a key to:
  - Reduce **infrastructure** and **operational costs**
  - Reduce **environmental** footprint
  - Making AI more **accessible**

# (1) What are the insights?

- **Insights** from 5+ years building DSAs for DDN at Google, summarized into 10 **lessons**.
- Lessons split into **3 categories**:
  1. Lessons applying to **any DSA**
  2. Lessons focusing on **DNN DSAs**
  3. Lessons about **DNN applications**

## (2) Insights: 3 lessons applying to any DSA

#	Insight	Solution
1	<b>Logic</b> improves much faster than wires and SRAM	⇒ DSAs architecture should mitigate the impact of <b>interconnect latency</b> and <b>wire scaling challenges</b>
2	DSA <b>architectures</b> and <b>compilers co-evolve</b> , and developing optimized compilers is expensive and time-consuming	⇒ New DSA should <b>leverage existing compiler optimizations</b> .
3	Different from benchmarking tools, big companies care about the <b>overall lifetime efficiency</b> of DSAs	⇒ Optimize for <b>perf/TCO</b> is a better target when designing DSA than perf/CapEx.

$$\text{TCO} = \text{CapEx} + 3 \times \text{OpEx}$$

### (3) Insights: 3 lessons focusing on DNN DSAs

#	Insight	Solution
4	<b>Backwards ML compatibility</b> enables rapid deployment of trained DNNs	⇒ New DSA version should look like previous version from <b>compiler's perspective</b> and support the same precision.
5	Inference DSA requires <b>global deployment</b> to reduce latency, but <b>liquid cooling</b> might not be possible everywhere	⇒ Inference DSAs should support <b>air cooling</b> .
6	Some inference apps need <b>floating point</b> arithmetic	⇒ DSA should support <b>multiple precision</b> (bf16 and int8), i.e. <b>quantization</b> must be <b>optional</b> .

## (4) Insights: 4 lessons about DNN applications

#	Insight	Solution
7	Production inference normally needs <b>multi-tenancy</b> (for lower costs, reduced latency and good software engineering practices)	⇒ DSAs need <b>local memory</b> (i.e. DRAM , since all weights can't fit in SRAM) for fast <b>switching</b> time between models.
8	DNNs <b>grow</b> ~ <b>1.5x</b> annually in memory and compute	⇒ DSAs architects should provide <b>headroom</b> to support DNN growth.
9	DNN workloads <b>evolve</b> with DNN breakthroughs	⇒ <b>Programmability</b> and <b>flexibility</b> are important factors for inference DSAs to track DNN progress.
10	The inference <b>SLO</b> is <b>P99 latency</b> , not batch size	⇒ DSAs should take advantage of <b>larger batch sizes</b> .

# What is the solution? Is it feasible?

These 10 lessons helped shape the new TPUv4i's design:

1. Based on **TPUv3**
2. One core chip for **inference** (two core TPUv4 for training)

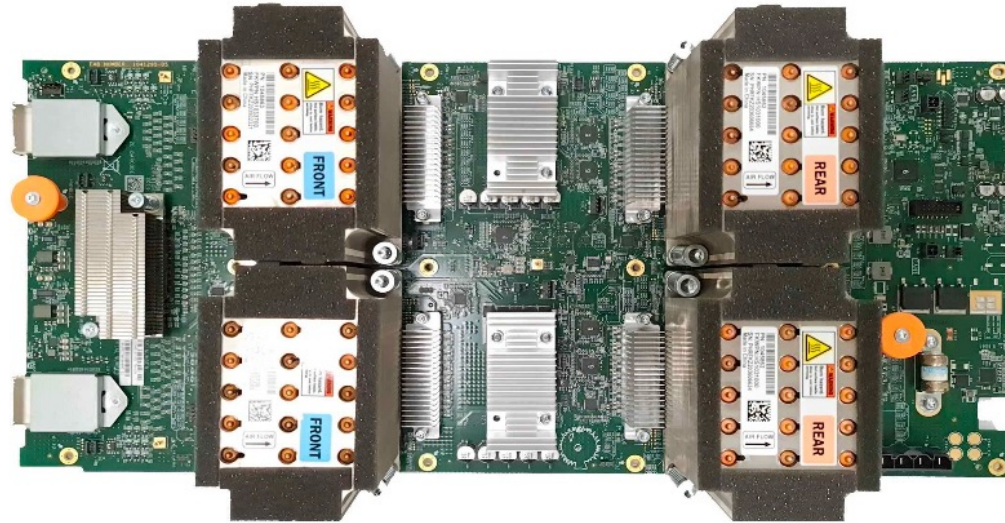


Figure 7. TPUv4i board with 4 chips that are connected by ICI.



# What is the solution? Is it feasible?

*Unequal hardware improvements*

*Leverage compiler optimizations*

*Design for perf/TCO*

*Backwards ML*

*Air cooling*

*Floating point support*

*Multi-tenancy*

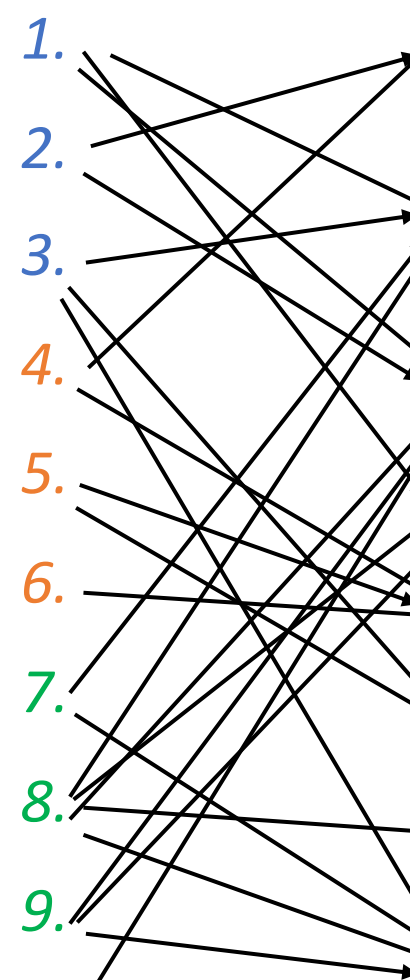
*Memory/compute 1.5x grow*

*DNN workloads evolve*

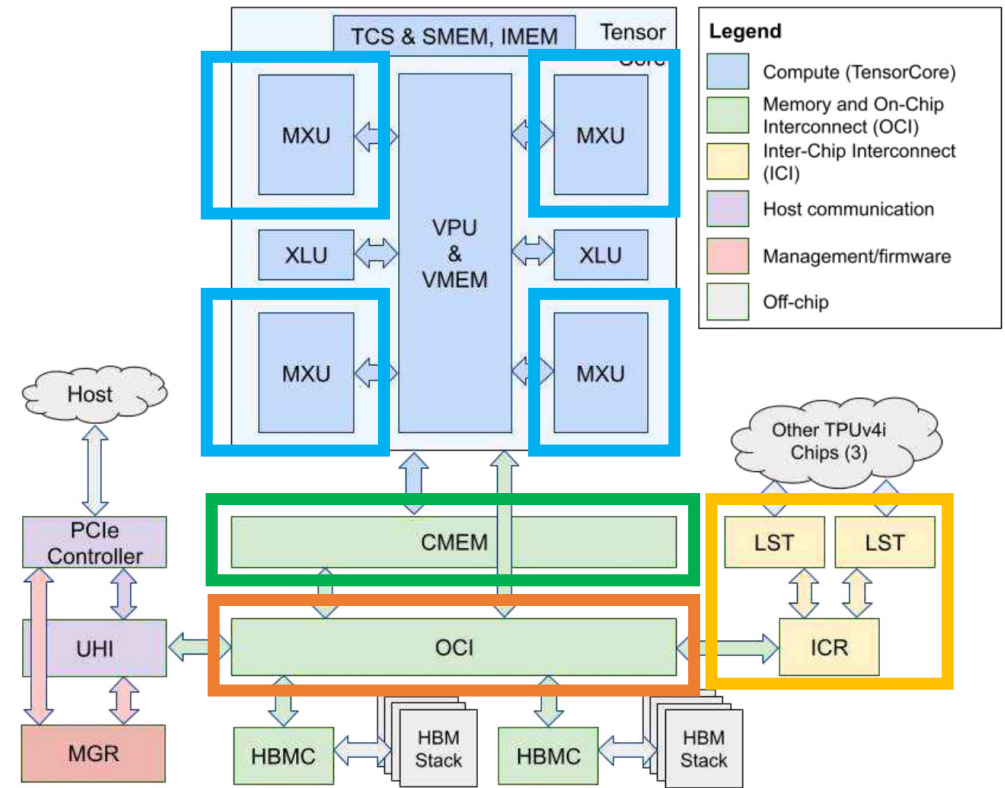
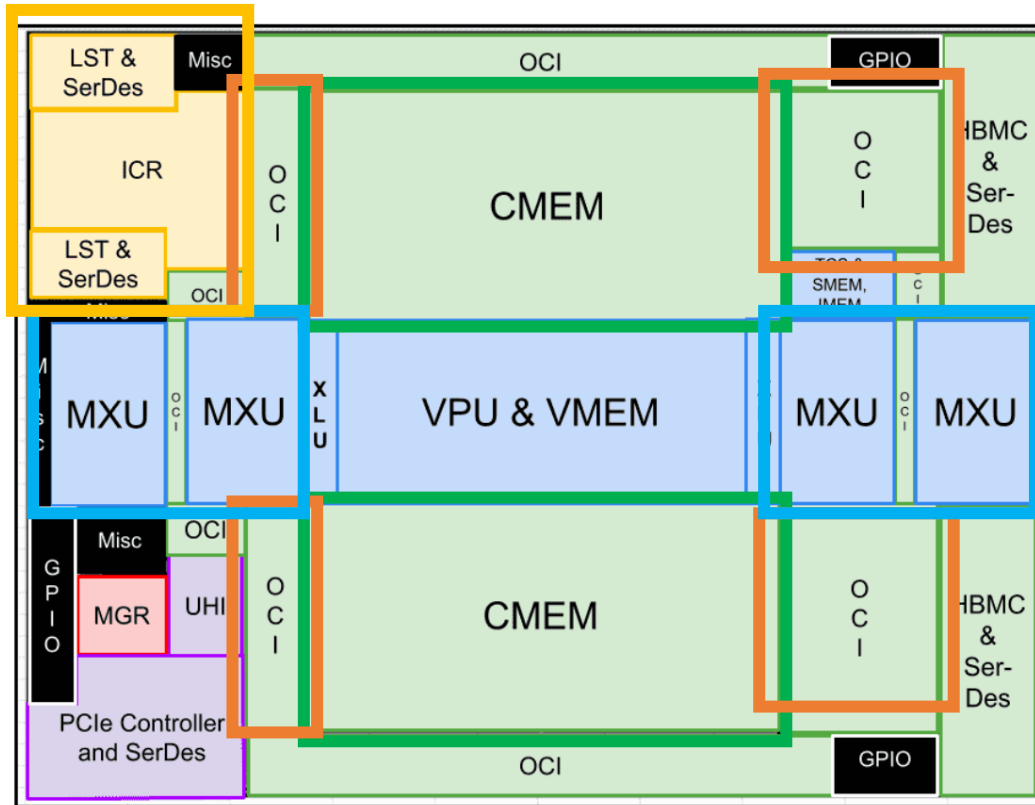
*Inference SLO is P99 latency*

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.
- 8.
- 9.
- 10.

1. *Compiler compatibility, not binary compatibility*
2. *Increased on-chip SRAM storage with CMEM*
3. *4D tensor DMA*
4. *Custom on-chip interconnect (OCI)*
5. *Arithmetic improvements*
6. *Clock Rate, TDP*
7. *Inter-Chip Interconnect scaling*
8. *Workload analysis features*



# What is the solution? Is it feasible?



2. Increased on-chip SRAM storage with common memory (CMEM)

5. Arithmetic improvements with four MXUs

4. Custom on-chip interconnect (OCI)

7. Inter-chip Interconnect (ICI) scaling

3. 4D tensor Direct Memory Access (DMA)

# What is the solution? Is it feasible?

## **1. Compiler compatibility, not binary compatibility:**

XLA compiler compatible with multiple TPUs

⇒ only LLO should be optimized per TPU

## **8. Clock Rate, TDP (Thermal Design Power):**

Clock rate of 1.05 GHz + chip TDP of 175W

⇒ Support air cooling and reduce TCO

## **6. Workload analysis features:**

Tracing and performance counter hardware features

⇒ Analyze system-level bottlenecks and support continuous system-level performance improvements

# What is the takeaway message?

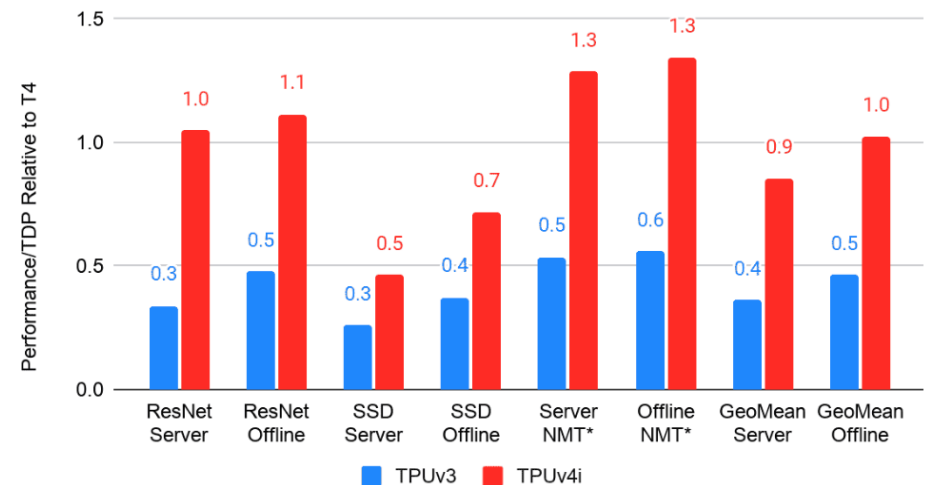
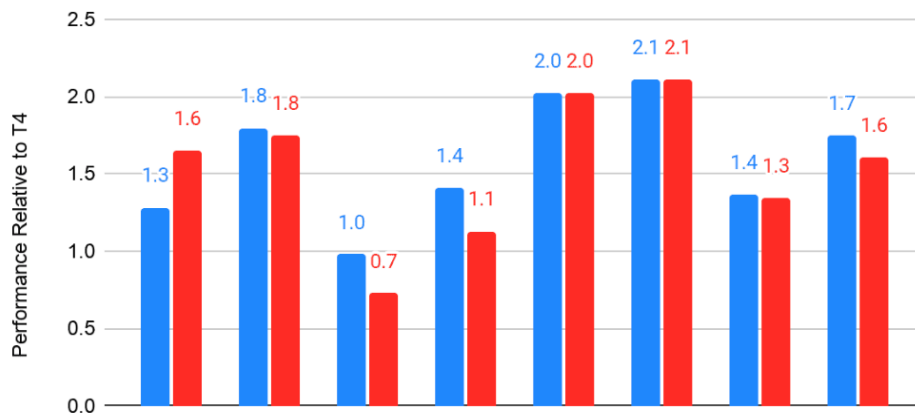
Key takeaway messages in designing DSA for DNN:

1. **Design** should be driven by the specific needs of **real-world workloads**.
2. **Iterative improvement** is the key to keeping up with the fast-evolving demand.
3. DSA should be **programmable** and **flexible** to cover previous/current/future needs.
4. **Moore's Law** is diminishing and **Dennard scaling** is dead  $\Rightarrow$  **hardware/software/DNN co-design** is the best chance for DNN DSAs to continue making progress.

# What is the takeaway message?

Takeaway messages regarding the TPUv4i chip itself:

1. TPUv4i has **similar performance** to TPUv3 (1.9x TPUv2), but is much **more energy efficient** (2.3x perf/TDP).
2. "Killing 2 birds with 1 stone": two separate chips for inference/training is key to achieve **better perf/TCO**.
3. **Thermal Design Power** offers a good proxy for DSA **Total cost of ownership**



# Will this paper win the test of time award?

NO

1. **Architectural choices** made for the TPUv4i **obsolete** in 15+ years?
2. New **computing paradigms** are emerging.  
⇒ DNN DSAs may be very different in 15+ years.
3. Most insights are related to ML **challenges from 2015-2020**  
⇒ ML community's challenges will be different in 15+.

Name one reason why this paper should have not appeared in MLSYS, NeurIPS, ICML, OSDI, ASPLOS, etc.?

The paper was presented to ISCA21 (103 citations on Semantic Scholar)

1. Valuable for the hardware designer community.
2. Scope limited to Google infrastructure?
3. Limited comparison to existing (similar) DSA for DNN?
4. Architectural description broad, with few details
5. Concerns regarding writing:
  1. Paper is very dense, with a lot of jargon/abbreviations, hard to read.
  2. Some motivations are obscured by referencing internal feedback from other design teams, without providing further insights.

# Preliminary

<i>Feature</i>	<i>TPUv1</i>	<i>TPUv2</i>	<i>TPUv3</i>	<i>TPUv4i</i>	<i>NVIDIA T4</i>
Peak TFLOPS / Chip	92 (8b int)	<u>46 (bf16)</u>	<u>123 (bf16)</u>	<u>138 (bf16/8b int)</u>	65 (ieee fp16)/130 (8b int)
First deployed (GA date)	Q2 2015	<u>Q3 2017</u>	<u>Q4 2018</u>	<u>Q1 2020</u>	Q4 2018
DNN Target	Inference only	<u>Training &amp; Inf.</u>	<u>Training &amp; Inf.</u>	<u>Inference only</u>	Inference only
Network links x Gbits/s / Chip	--	4 x 496	<u>4 x 656</u>	<u>2 x 400</u>	--
Max chips / supercomputer	--	256	<u>1024</u>	--	--
Chip Clock Rate (MHz)	700	700	<u>940</u>	<u>1050</u>	585 / (Turbo 1590)
Idle Power (Watts) Chip	28	<u>53</u>	<u>84</u>	<u>55</u>	36
TDP (Watts) Chip / System	75 / 220	<u>280 / 460</u>	<u>450 / 660</u>	<u>175 / 275</u>	70 / 175
Die Size (mm <sup>2</sup> )	< 330	<u>&lt; 625</u>	<u>&lt; 700</u>	<u>&lt; 400</u>	545
Transistors (B)	3	<u>9</u>	<u>10</u>	<u>16</u>	14
Chip Technology	28 nm	<u>16 nm</u>	16 nm	<u>7 nm</u>	12 nm
Memory size (on-/off-chip)	28MB / 8GB	<u>32MB / 16GB</u>	32MB / 32GB	<u>144MB / 8GB</u>	18MB / 16GB
Memory GB/s / Chip	34	<u>700</u>	<u>900</u>	<u>614</u>	320 (if ECC is disabled)
MXU Size / Core	1 256x256	<u>1 128x128</u>	<u>2 128x128</u>	<u>4 128x128</u>	8 8x8
Cores / Chip	1	<u>2</u>	2	<u>1</u>	40
Chips / CPUHost	4	4	4	<u>8</u>	8



Operation		Picojoules per Operation		
		45 nm	7 nm	45 / 7
+	Int 8	0.03	0.007	4.3
	Int 32	0.1	0.03	3.3
	BFloat 16	--	0.11	--
	IEEE FP 16	0.4	0.16	2.5
	IEEE FP 32	0.9	0.38	2.4
×	Int 8	0.2	0.07	2.9
	Int 32	3.1	1.48	2.1
	BFloat 16	--	0.21	--
	IEEE FP 16	1.1	0.34	3.2
	IEEE FP 32	3.7	1.31	2.8
SRAM	8 KB SRAM	10	7.5	1.3
	32 KB SRAM	20	8.5	2.4
	1 MB SRAM <sup>1</sup>	100	14	7.1
GeoMean <sup>1</sup>		--	--	2.6
DRAM		Circa 45 nm	Circa 7 nm	
	DDR3/4	1300 <sup>2</sup>	1300 <sup>2</sup>	1.0
	HBM2	--	250-450 <sup>2</sup>	--
	GDDR6	--	350-480 <sup>2</sup>	--

Table 2. Energy per Operation: 45 nm [16] vs 7 nm. Memory is pJ per 64-bit access.

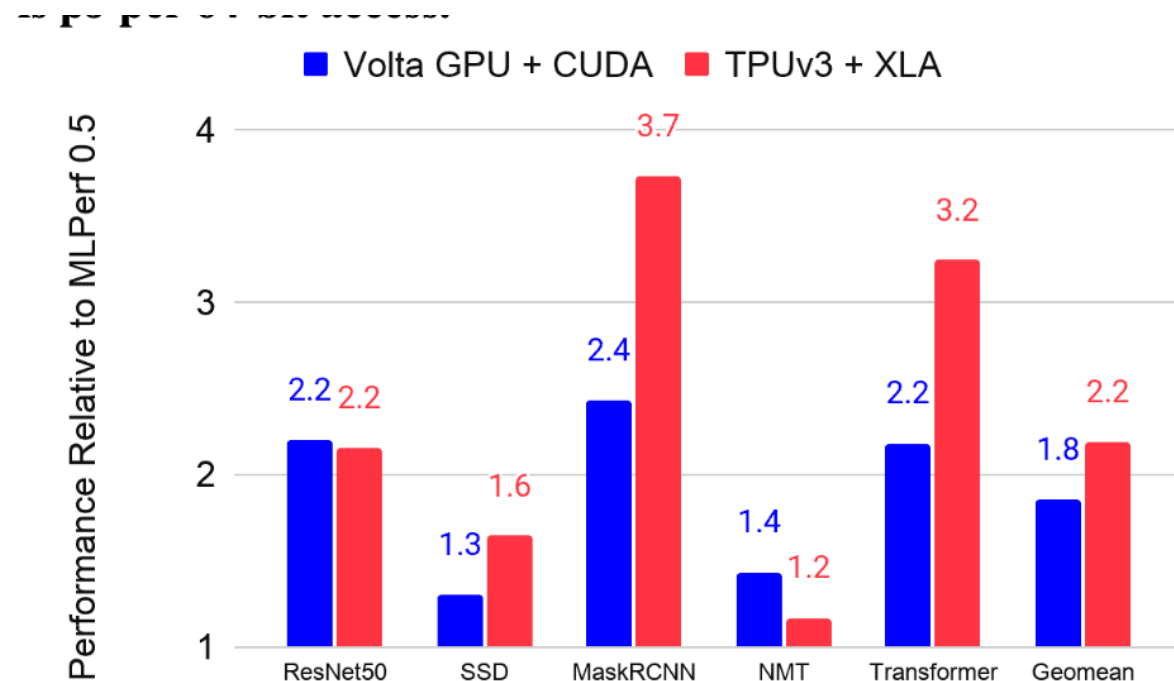
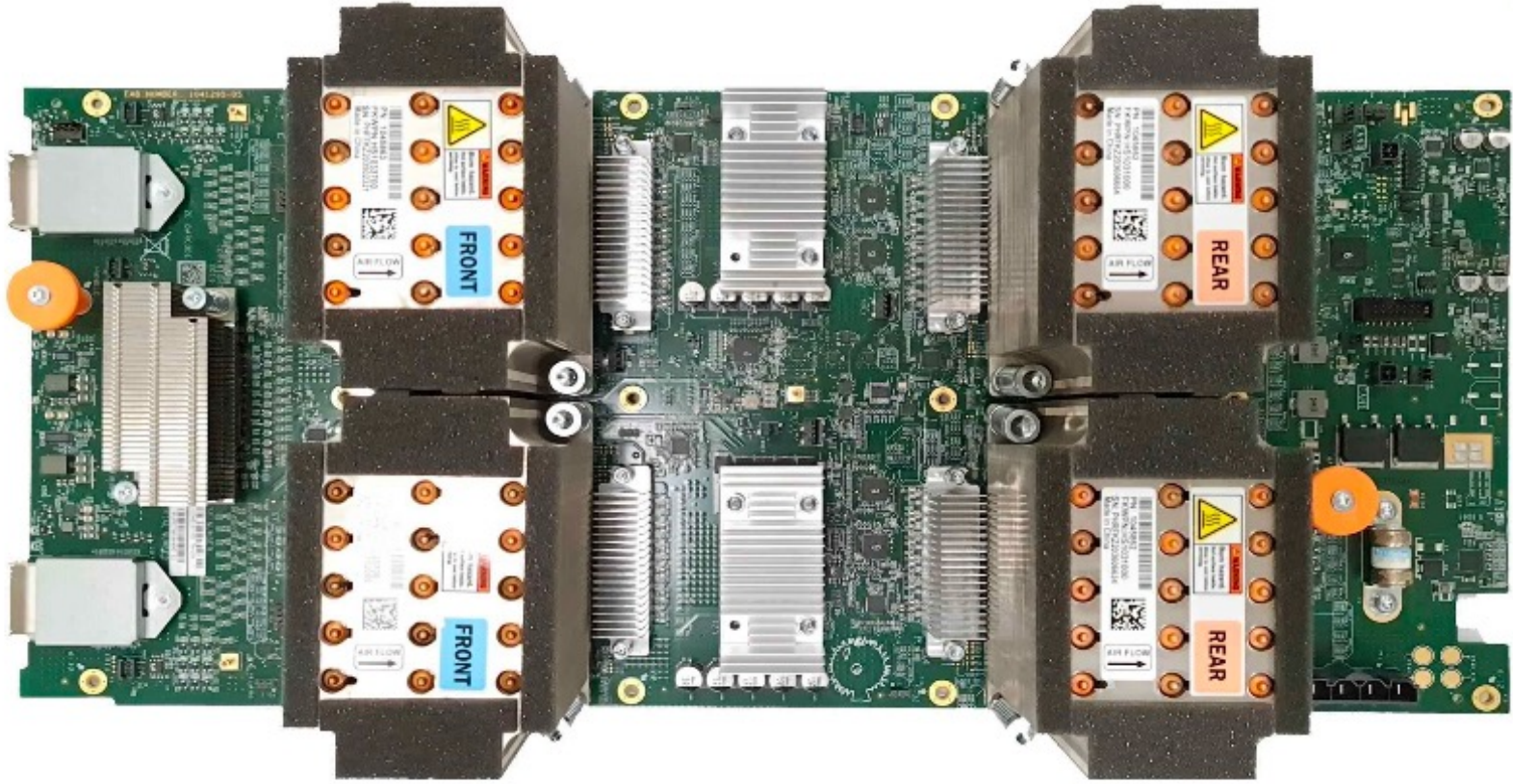
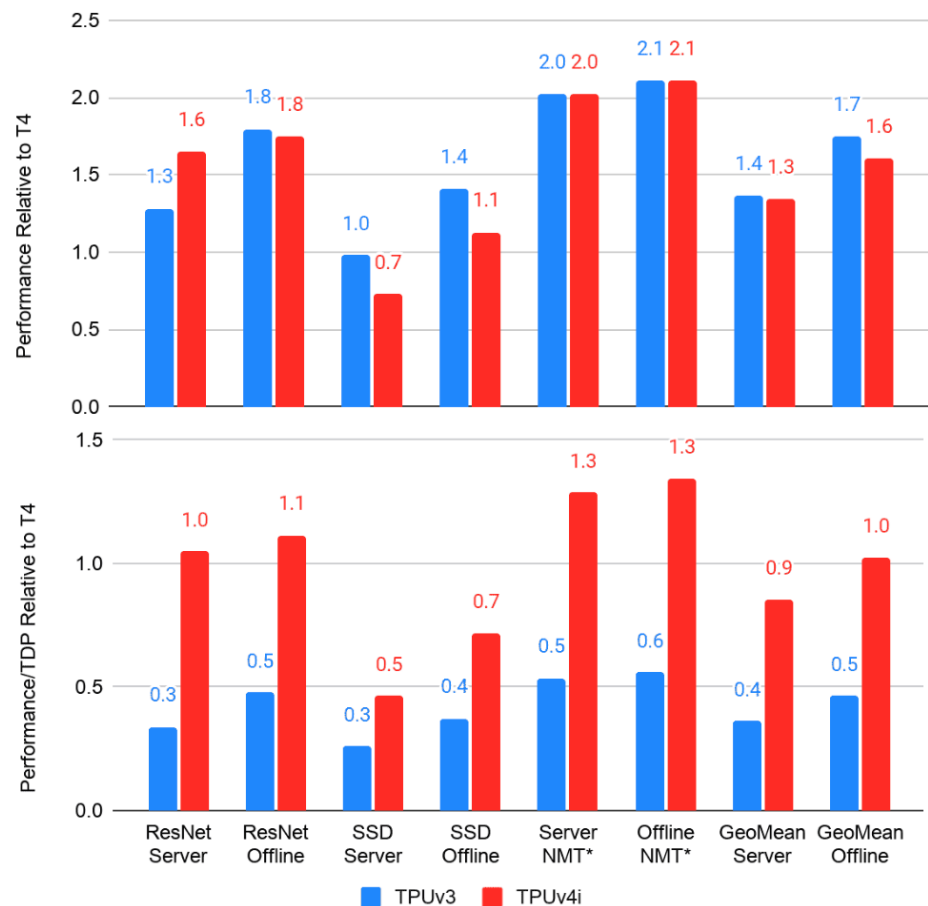


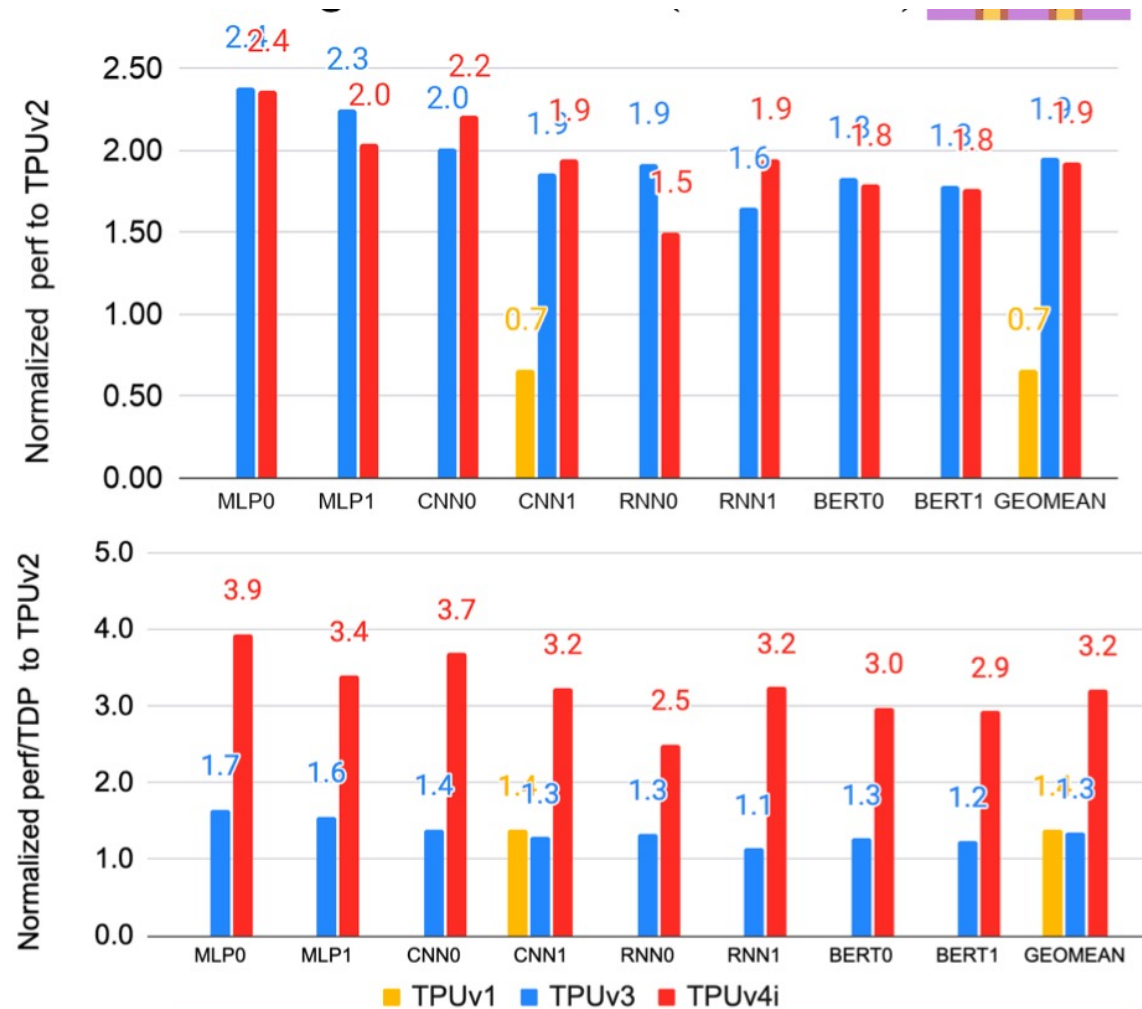
Figure 2. DSA gains per chip for MLPerf Training 0.5 to 0.7 over 20 months for the same compilers. The unverified TPUv3 MLPerf 0.5 scores for Mask R-CNN and Transformer are from [23]; all other results are from [28].



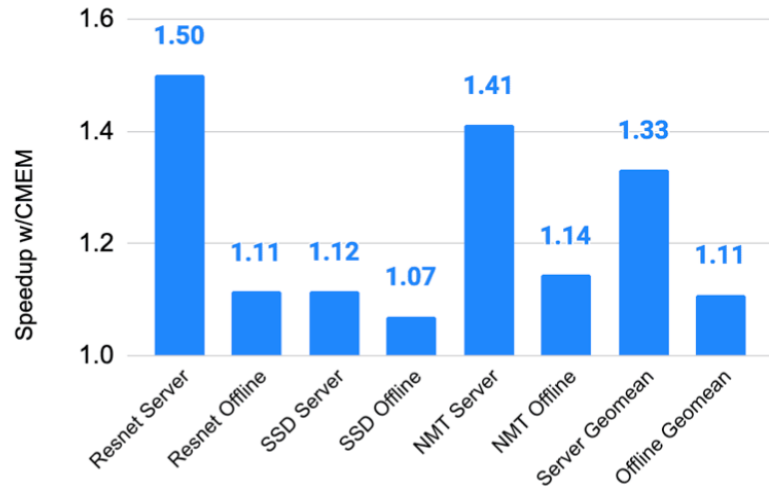
**Figure 7. TPUv4i board with 4 chips that are connected by ICI.**



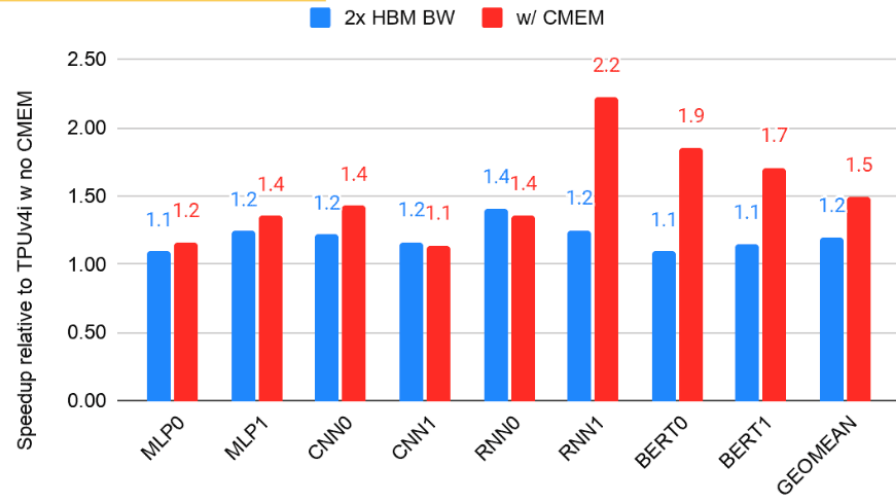
**Figure 9. Performance (top) and performance/system TDP** ③ relative to T4 for TPUv3/v4i in our datacenter (\$7.B). Note that the T4 uses int8 for ResNet and SSD and fp16 for NMT. The TPUs use bf16 for all three to maintain backwards ML compatibility ④ with TPUv3/v4. MLPerf Inference 0.7 omits NMT, so we use MLPerf Inference 0.5 code for it and 0.7 code for ResNet and SSD in Figures 9, 10, and 13. (These results are unofficial, as they have not been verified by MLPerf.)



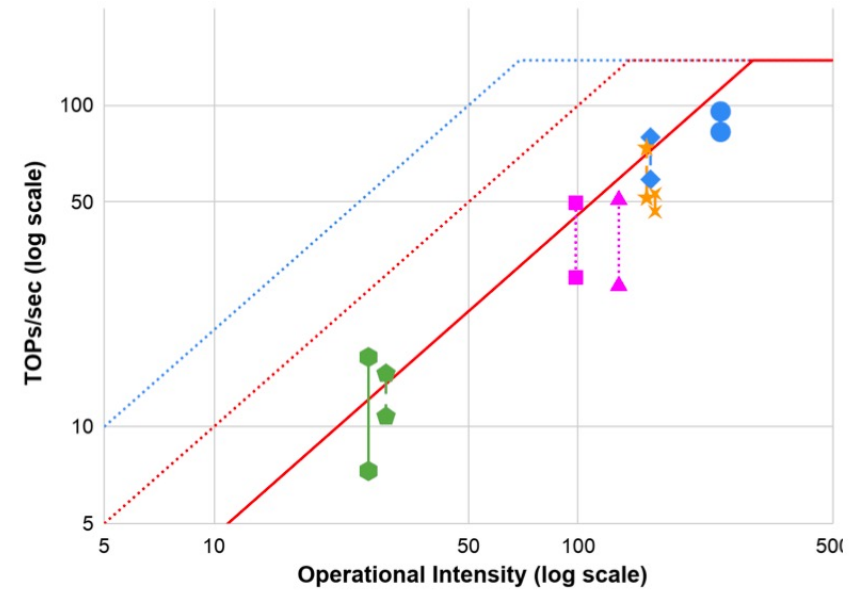
**Figure 8. Performance (top) and performance/system Watt** ③ for production apps ⑨ relative to TPUv2 for the other TPUs.



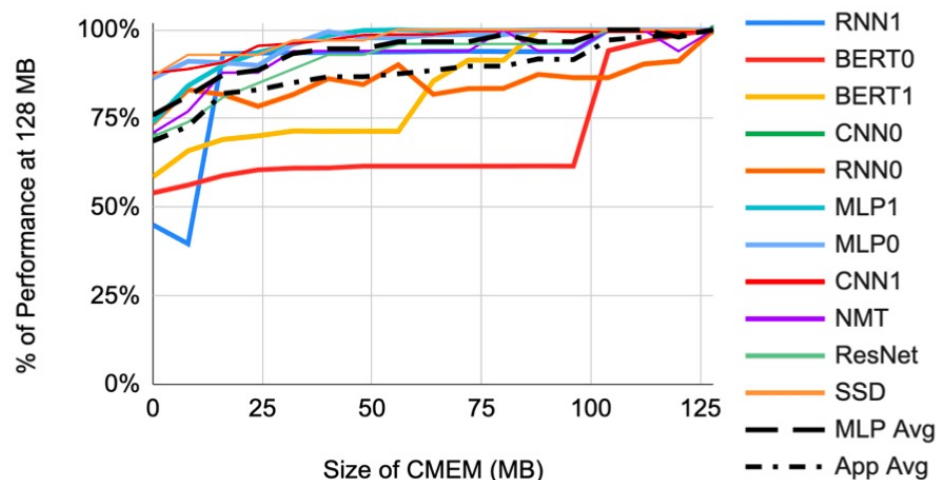
**Figure 10. Unverified MLPerf Inference impact of turning on CMEM vs no CMEM.**



**Figure 11. Performance of 2X HBM bandwidth with no CMEM relative to CMEM with standard HBM bandwidth.**



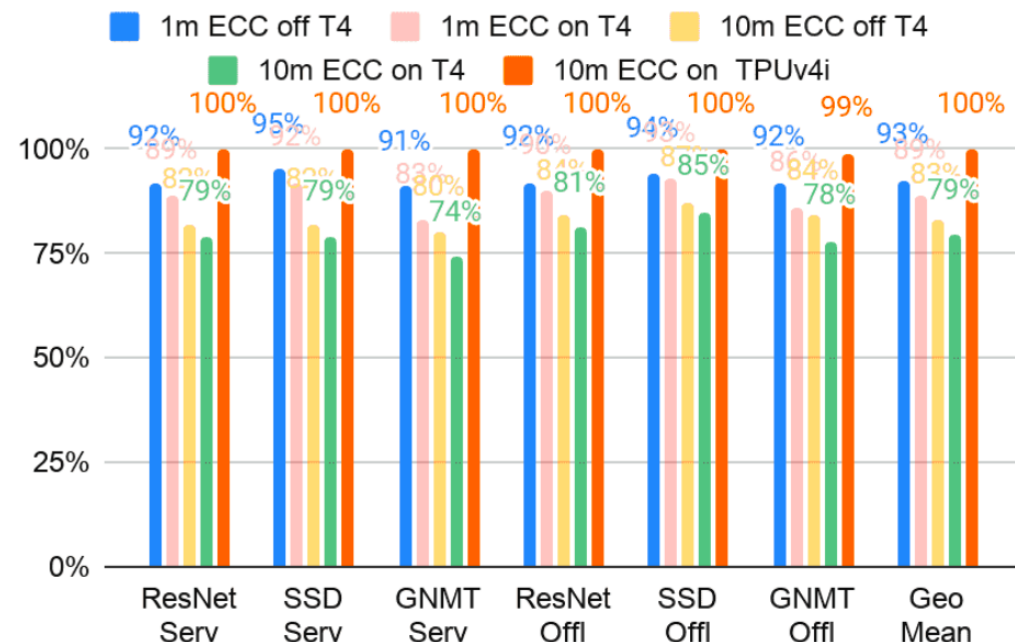
**Figure 12. Roofline model showing apps without CMEM (low point) vs with CMEM (high point). Operational intensity (OI) here is operations divided by memory accesses to HBM or to CMEM. If OI were relative to HBM only, CMEM would increase OI and move the points to the right as well as up.**



**Figure 13. Percent of 128 MB speed as CMEM varies 0–128 MB for the apps and MLPerf Inference 0.5-0.7 server code.**

	<i>perf/TDP</i> <i>vs T4</i>	<i>TDP</i> <i>(W)</i>	Cores	<i>GHz</i>	<i>MiB</i> <i>SRAM</i>	<i>Type</i> <i>DRAM</i>
T4	1.00	70	40	0.6	20	GDDR6
Goya	0.46	200	8	2.1	≈32	DDR4
Nervana	0.69	100	24	1.1	60	LPDDR4
Zebra	0.25	225	--	--	54	DDR4
HanGuang	2.17	280	4	0.7	192	none

**Table 6. Five DSAs that ran the MLPerf Inference 0.5 (Goya, Nervana, HanGuang) or 0.7 (T4, Zebra) server scenarios. Goya ran only SSD, and the last three ran only ResNet. Performance is relative to T4 for MLPerf Inference 0.5. TDP in this table is per chip rather than per system, since the latter was unavailable. All have 16–32GB of DRAM except HanGuang, which has none. All results are from [29].**



**Figure 14. T4 and TPUv4 running unverified MLPerf Inference benchmarks 0.5/0.7 at Google with memory ECC off and on. NVIDIA’s MLPerf score is 100% for T4 and 100% for TPUv4i is unverified MLPerf in our datacenter. For the 1 minute case, the T4 was idle initially. It then ran MLPerf with ECC off and on for 1 minute at the fastest clock rate. (T4 offers inline ECC, which uses memory bandwidth.) For the 10 minute case, the machines are not idle beforehand. TPUv4i’s speed is unchanged whether ECC is on or off or how long it runs.**

TPUv4i FLOPS/s would drop to 22% and Roofline to 97%.

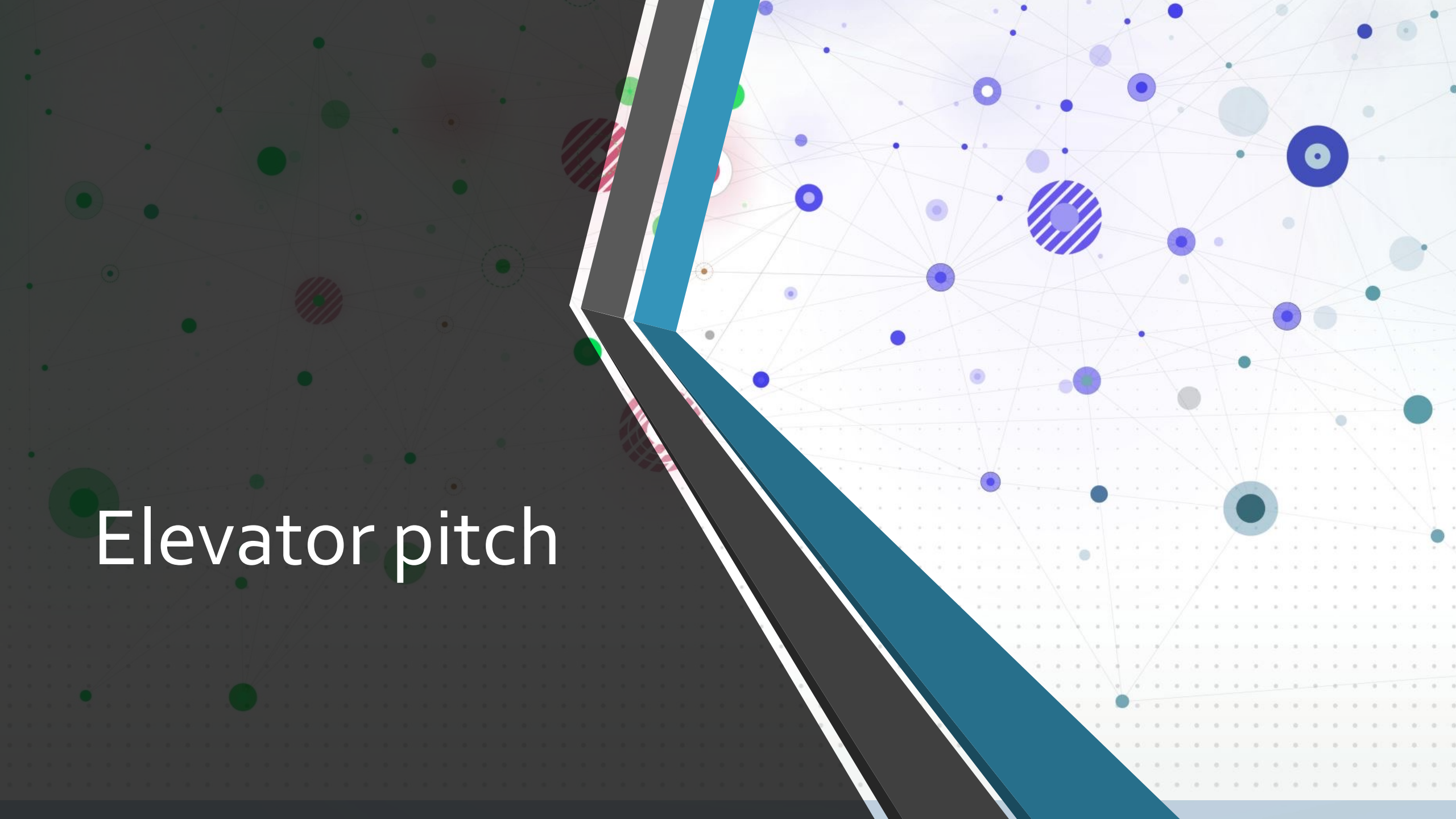
TPU	TPUv1	TPUv2	TPUv3	TPUv4i
MXUs/Chip	1 256 x256	2 128x128	4 128x128	4 128x128
MXUs % Die Area	24%	8%	11%	11%
FLOPS/s Utilization	20%	51%	38%	33%
HBM Roofline Util.	20%	66%	63%	99%

**Table 7. Average utilization of peak performance and of roofline for our eight production applications.**

# RaPiD:

AI Accelerator for Ultra-low Precision Training and Inference

# Elevator pitch





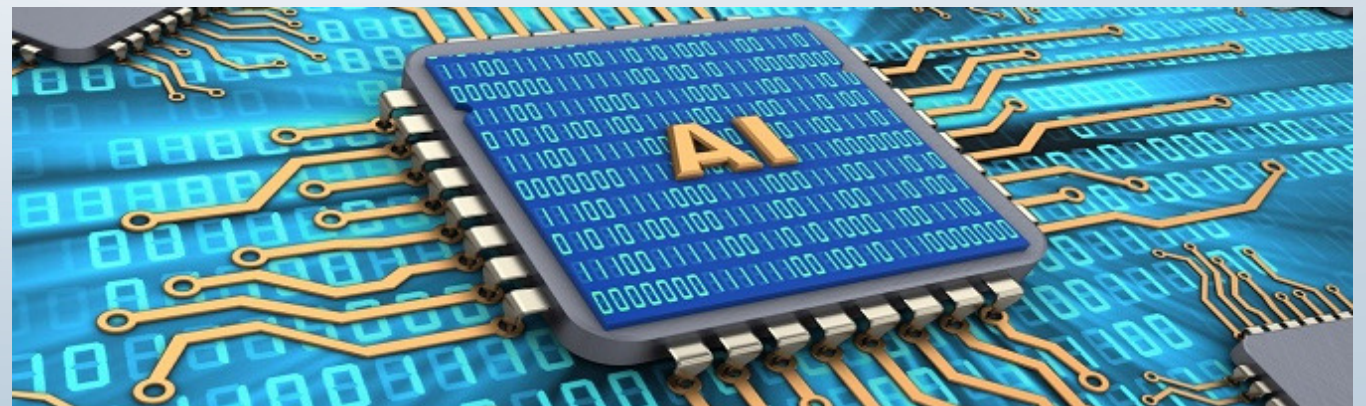
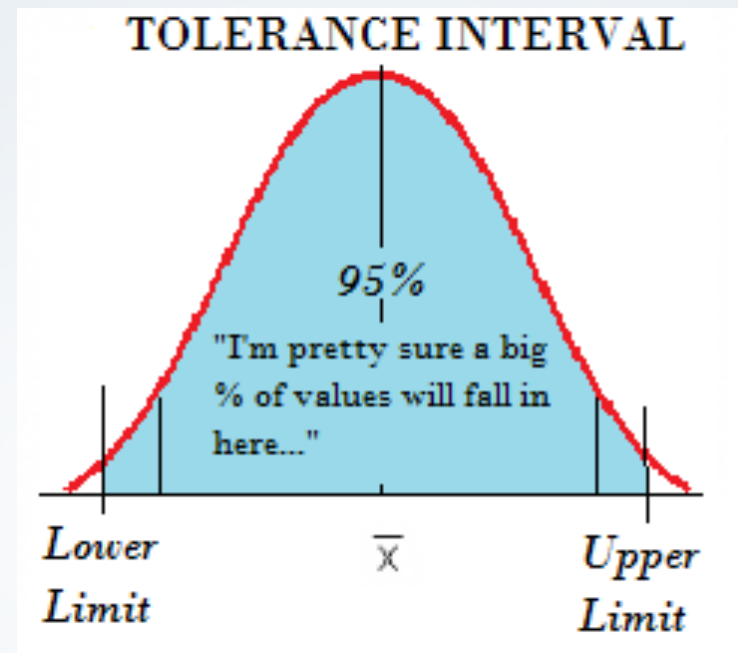


AI =

or



Move forward?



Then,  
hopefully...

FASTER

CHEAPER

BETTER

# Write-up questions

x

# 1. What is the problem?

- AI/DNN have different paradigm of computation workloads
- the high accuracy of these tasks have a high computational cost.
- => they severely stress the capabilities of traditional computing platforms.

## 2. What are the insights?

- AI workloads are static dataflow graphs
- the computations only need a small number of primitives
- specialized hardware accelerators can improve system performance
- AI-workloads are error-resilient, so precision scaling can be beneficial.

### 3. What is the solution?

- varying level of precision is designed and supported;
- both performance and energy efficiency should be improved;
- several operations, like data-shuffle and polling, by special function units;
- zero-gating logic to boost the performance of sparse AI models;
- memory neighbor interface to incorporate core-to-core and core-to-memory communication and synchronization.

### 3. solution – cont'd

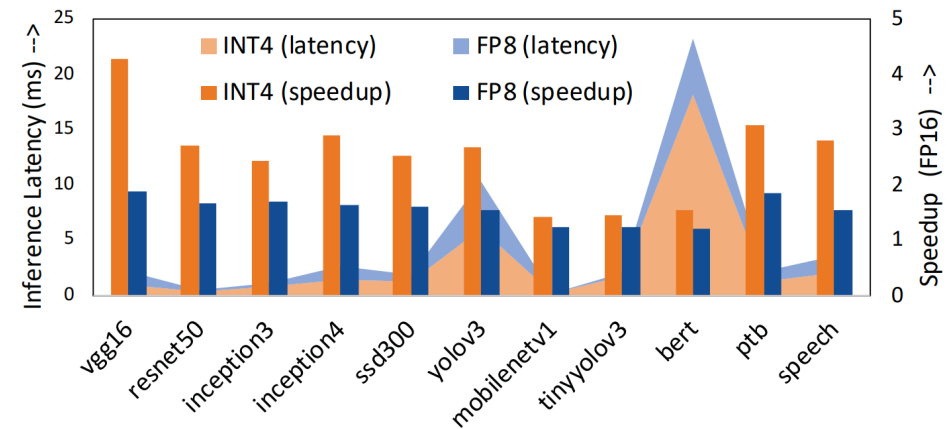


Figure 13: Classifications per second using 4-core RAPiD chip

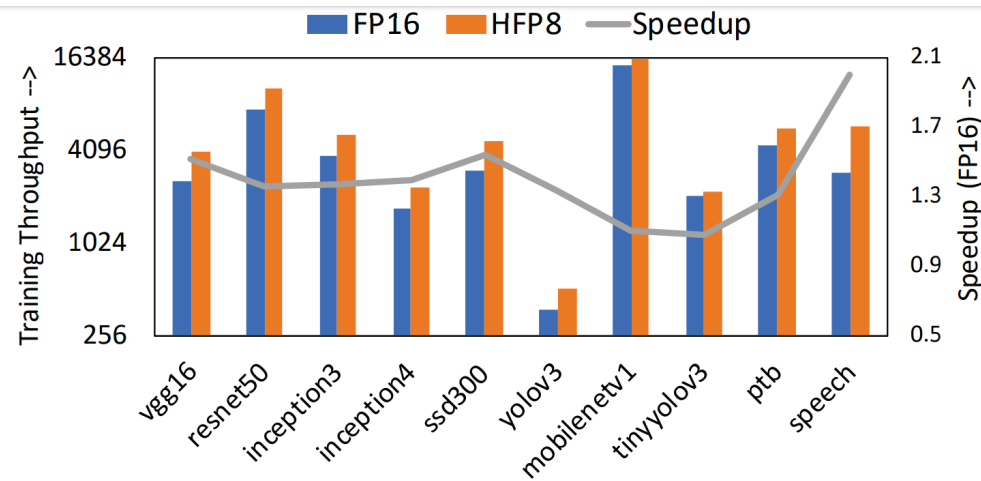
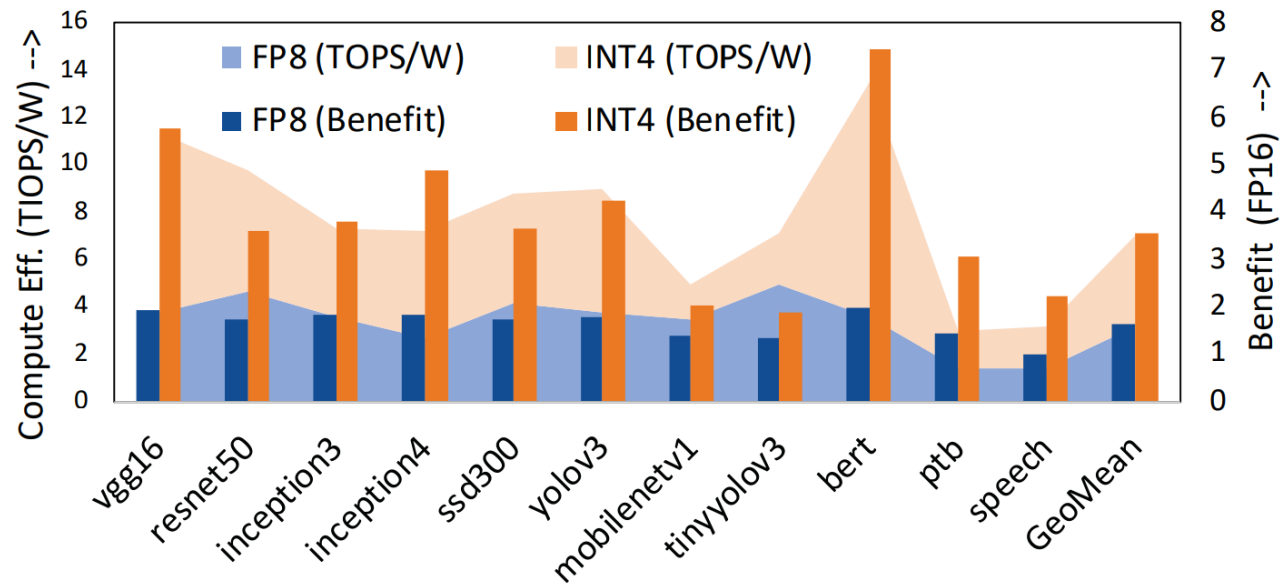


Figure 15: Throughput with 4-chip RAPiD training system





**Figure 14: Sustained TOPS/W on 4-core RAPID chip**

3. solution –  
cont'd

## 4. What is the takeaway message?

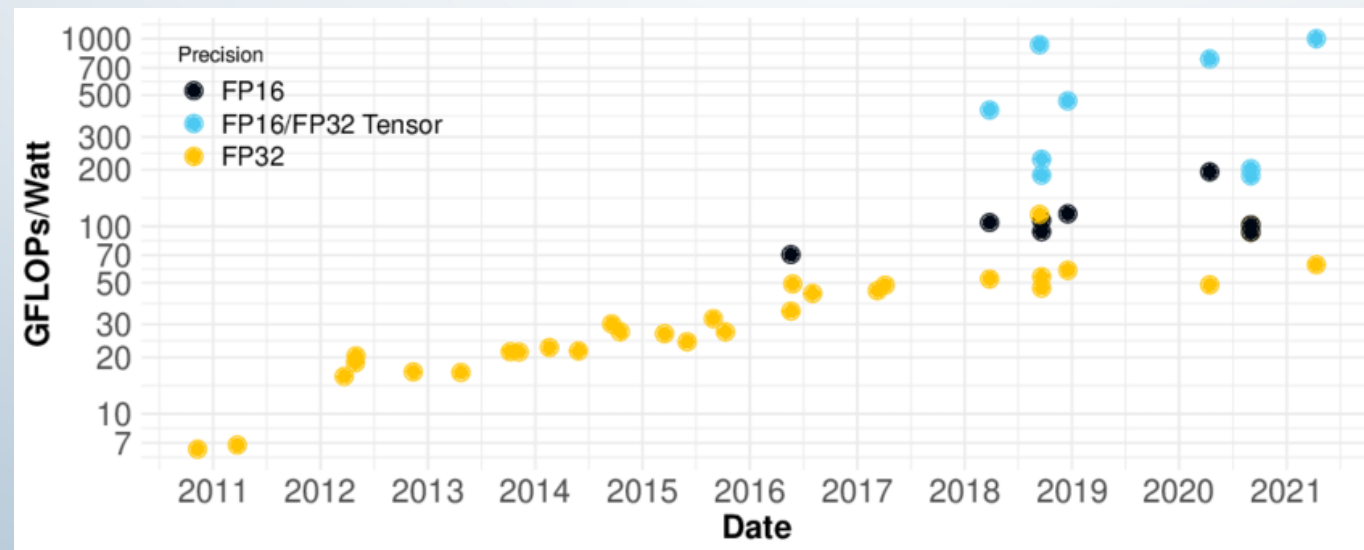
- Precision scaling is an effective foundation for designing hardware accelerators for AI.

## 5. Will this paper win the test of time award?

- Possible, from the perspective that the chip could be used for edge or battery-operated devices for AI applications with less energy consumption.
- However, in terms of generic AI accelerators, I don't think the design has achieved a great leap.

## 6. Reason why it isn't top conference quality?

- No critical reasons to reject this paper.
- However, how much better this accelerator performs with/compared-to different CPU/GPU/TPUs? Below btw.



Theoretically, for Nvidia GPU



# Thank you

Q&A