# Master Informatics Eng.

2020/21

*A.J.Proença*

**Beyond Vector Extensions** *(online)*
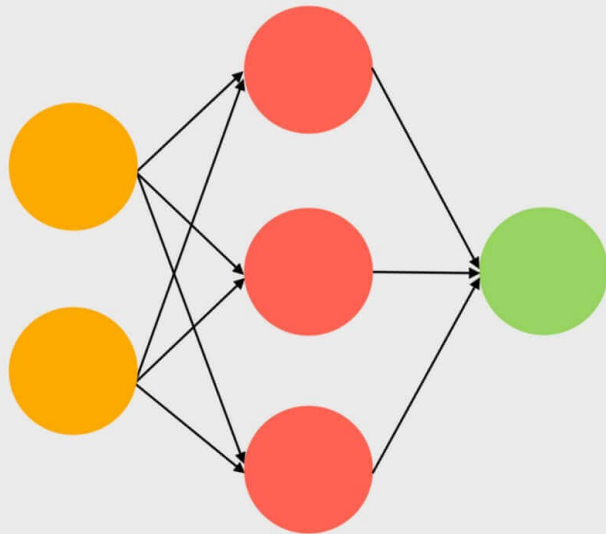
**(most slides are borrowed)**

- Evolution of vector/SIMD-extended architectures

  – **accelerators optimized for number crunching (GPU)**

  – **support for matrix <u>multiply + accumulate</u> operations**

    - most <u>scientific, engineering, AI & finance</u> applications use matrix computations, namely the dot product: multiply and accumulate the elements in a row of a matrix by the elements in a column from another matrix

    - typically these extensions are **Tensor Processing Unit** (TPU)

  – **support for half-precision FP & 8-bit integer**

    - <u>**machine learning using neural nets**</u> is becoming very popular; to compute the model parameters during training phase, intensive matrix products are used and with very low precision (is adequate!)
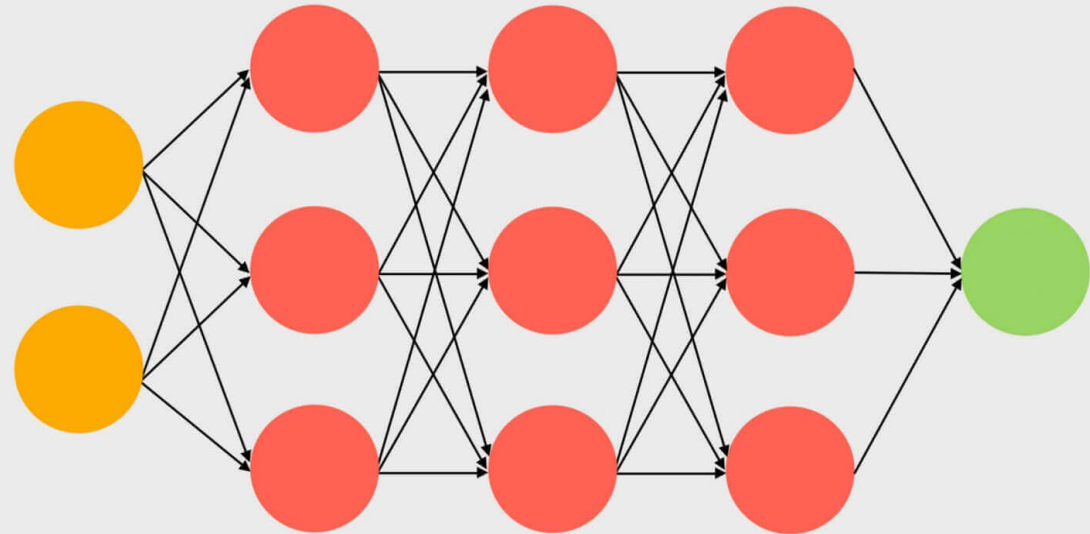
# Machine learning w/ neural nets & deep learning



**Artificial Neural Network (Single Layer ML)**

**Deep Neural Network (Multiple Layer ML)**
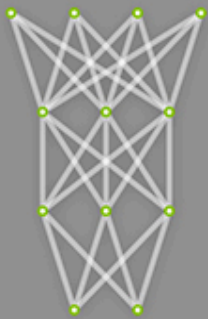
— Input Layer    — Hidden Layer    — Output Layer

https://www.techspot.com/article/2049-what-are-tensor-cores/

# Deep Learning workflow



Key algorithms to train & classify use matrix dot products,
but do not require high precision numbers!

*During model **training**, labeled data samples flow from input to output (I to O) through all layers of parametrized transformations — a forward pass. At the output end, the outpu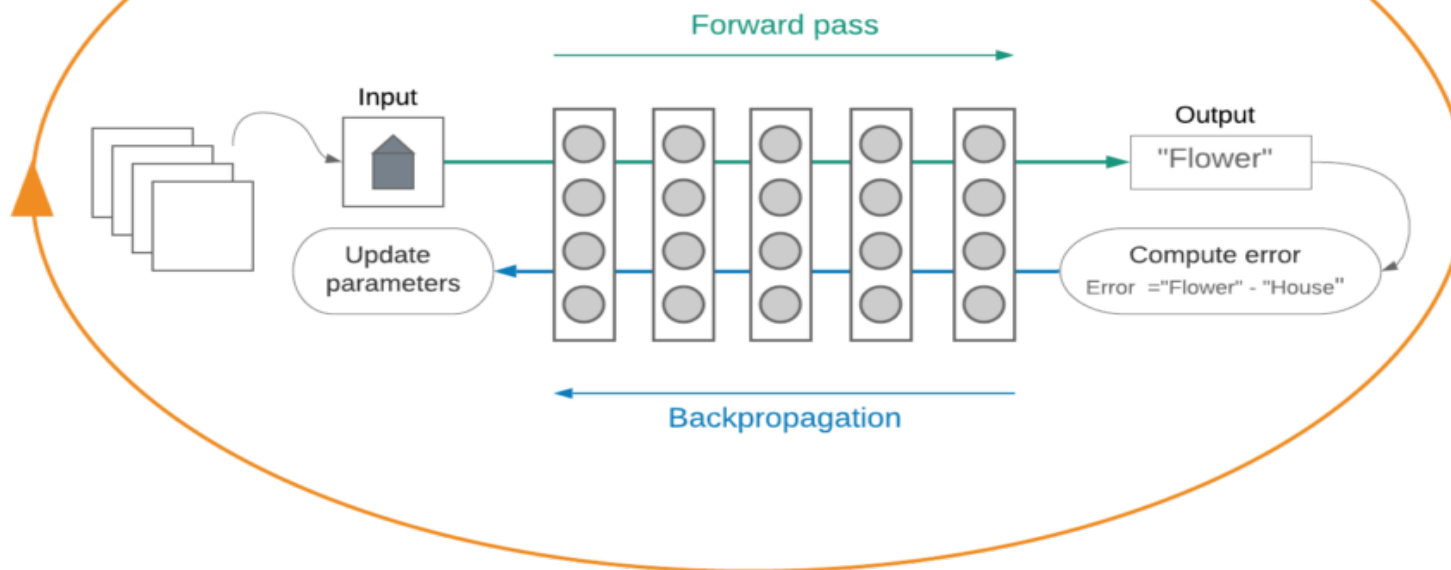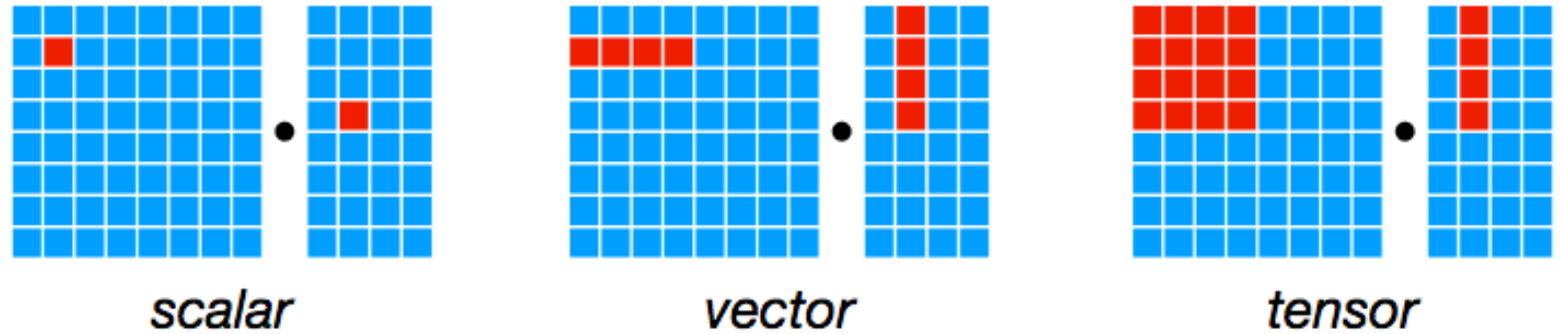t, or prediction, is compared to the correct answer for that particular input. Prediction error is computed; error being the difference between the predicted output and the correct one. Then, the error begins to work its way backwards in the O-to-I direction, via the backpropagation algorithm. As the error flows through each layer, it interacts with the I-to-O data that produced it (that data has been parked there, waiting for the error to come back) and, together, they determine how to change the layer's parameters to most effectively reduce the error. The parameters are then adjusted, and this process of forward-backpropagation steps continues for numerous passes over the set of training examples, until the error becomes insignificant or doesn't decrease anymore.*



https://www.cerebras.net/cerebras-wafer-scale-engine-why-we-need-big-chips-for-deep-learning/

# *Required hardware operations & data types to train & classify neural nets*

**Compute primitives**

scalar          vector          tensor

**Data type**

| Sign | Range | Precision |

FP32   8 BITS    23 BITS

int8

TF32 Range

TENSOR FLOAT 32 (TF32)   8 BITS   10 BITS

**IEEE 754**

TF32 Precision

FP16   5 BITS   10 BITS

**Google Brain** Team

BFLOAT16   8 BITS   7 BITS

# *Approaches to operations on tensors*

- **Tensor**: a mathematical object that describes the relationship between other mathematical objects that are all linked together; they are commonly shown as a <u>multidimensional array</u>

- Different approaches followed by chip manufacturers:
  - add <u>new extensions</u> to existing HPC vector devices
    - **NVidia**: tensor core units in HPC GPUs
    - **Intel**: AVX-512VNNI & AMX
  - develop SoC devices for embedded/specific application fields
    - neural net devices: **Google** TPU, **Intel** Habana, …
    - autonomous driving: **Tesla** FSD, **NVidia** Orin, …
    - smartphones: **Apple** A14 Bionic, **Huawei** Kirin 9000, Qualcomm Snapdragon, Samsung Exynos, …
    - gaming: …

# Approaches to operations on tensors

- **Tensor**: a mathematical object that describes the relationship between other mathematical objects that are all linked together; they are commonly shown as a <u>multidimensional array</u>

- Different approaches followed by chip manufacturers:
  - add <u>new extensions</u> to existing HPC vector devices
    - **NVidia**: tensor core units in HPC GPUs
    - **Intel**: AVX-512VNNI & AMX
  - develop SoC devices for embedded/specific application fields
    - neural net devices: **Google** TPU, **Intel** Habana, …
    - autonomous driving: **Tesla** FSD, **NVidia** Orin, …
    - smartphones: **Apple** A14 Bionic, **Huawei** Kirin 9000, Qualcomm Snapdragon, Samsung Exynos, …
    - gaming: …

$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32     FP16             FP16             FP16 or FP32

Figure 8.     Tensor Core 4x4 Matrix Multiply and Accumulate

FP16 storage/input     Full precision product     Sum with FP32 accumulator     Convert to FP32 result

F16
F16
×
more products
+
F32
F32

**For each SM:**
**8x 64 FMA ops/cycle**
**1 KFLOP/cycle!**

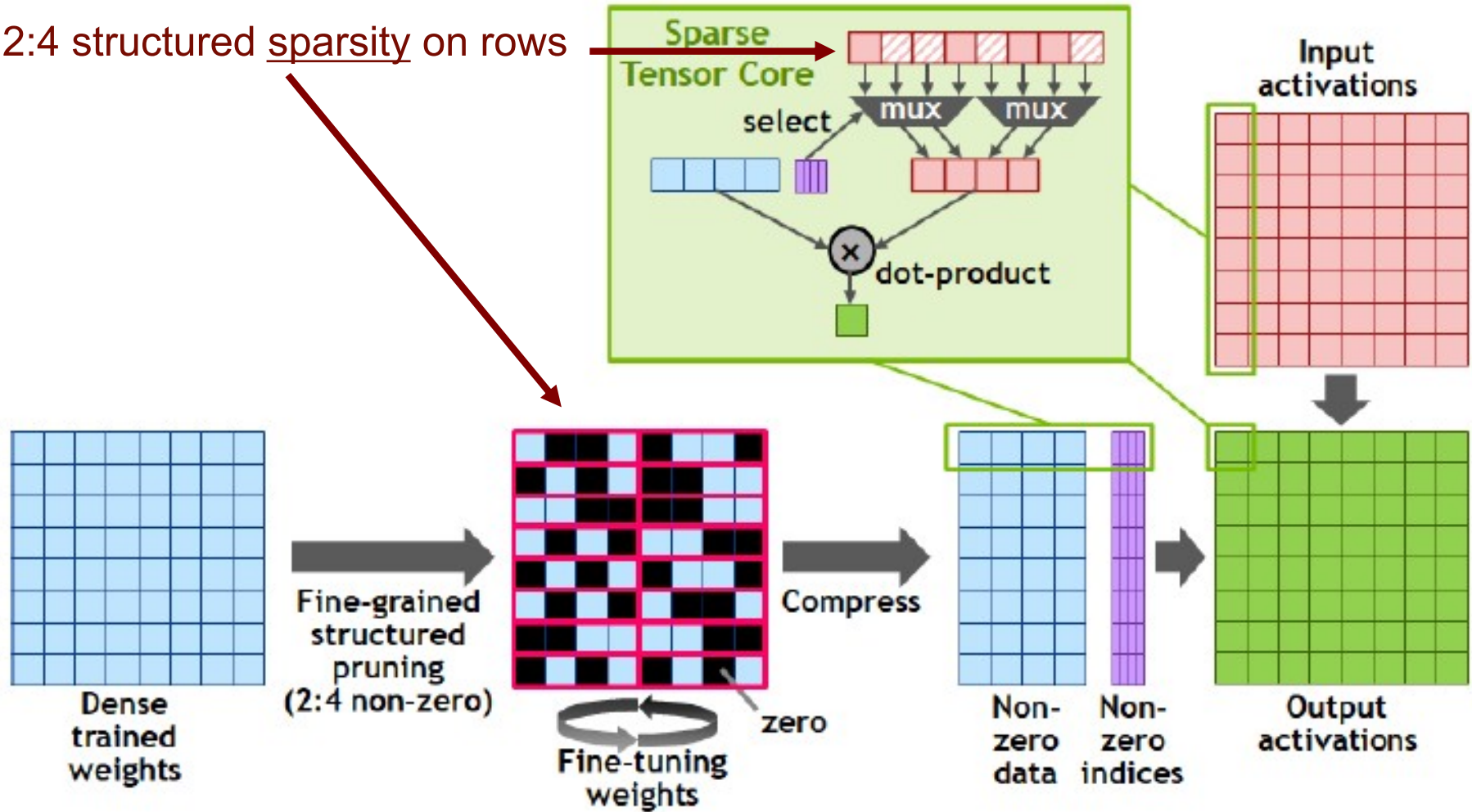Figure 9.     Mixed Precision Multiply and Accumulate in Tensor Core

http://www.nvidia.com/content/gated-pdfs/Volta-Architecture-Whitepaper-v1.1.pdf
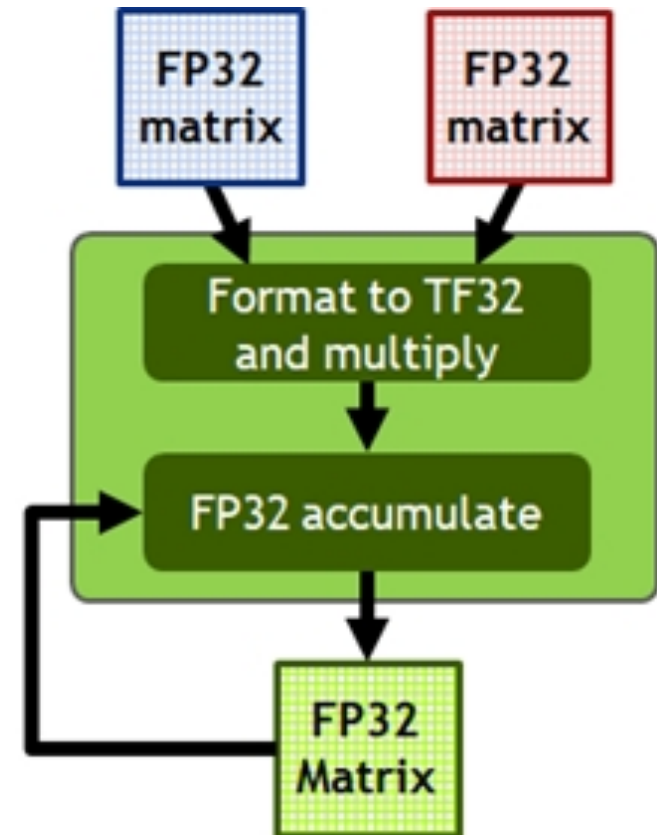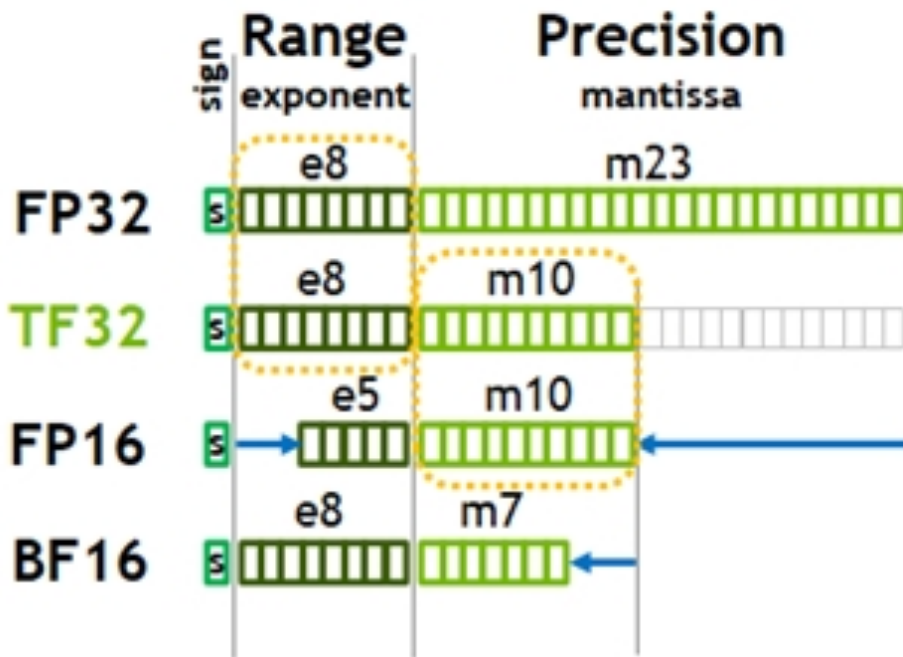
- GEMM *(Generic Matrix Multiplication)* computes **D = A * B + C**
  - **A** can be up to 8x8 matrix (mixed-precision ith **B**)
  - **B, C, D** can be up to 8x4 matrix

- Each SM in A100 has 4 Tensor Cores
  - 4x 256 FMA ops/cycle (FP16), **2 KFLOP/cycle**

- A100 with Fine-Grained Structured <u>Sparsity</u>
  - 2:4 structured sparsity on rows (2 non-zero values in every 4-entry vector)



Initialize : Dense Matrix → Train : Dense Matrix → Prune : Sparse Matrix → Retrain : Sparse Matrix → Sparsity Optimized Tensor Core

https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/nvidia-ampere-architecture-whitepaper.pdf

2:4 structured <u>sparsity</u> on rows

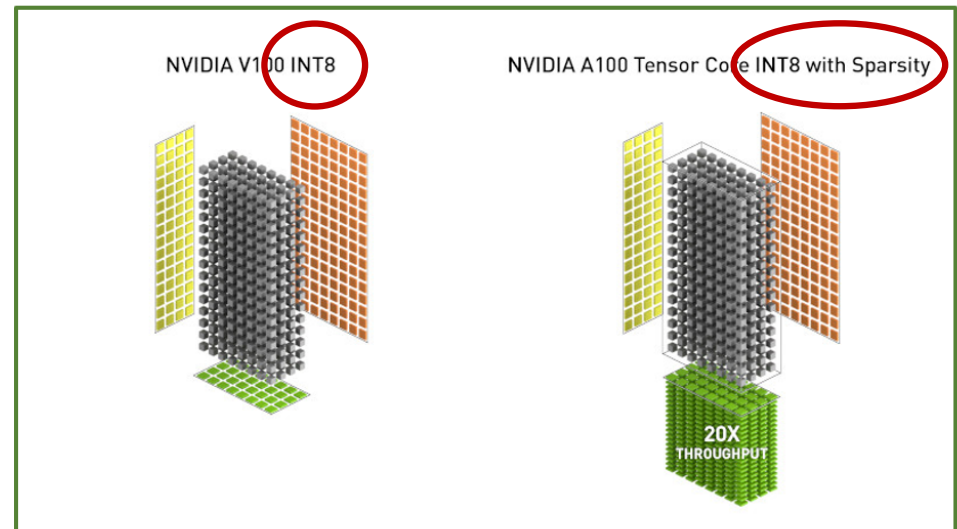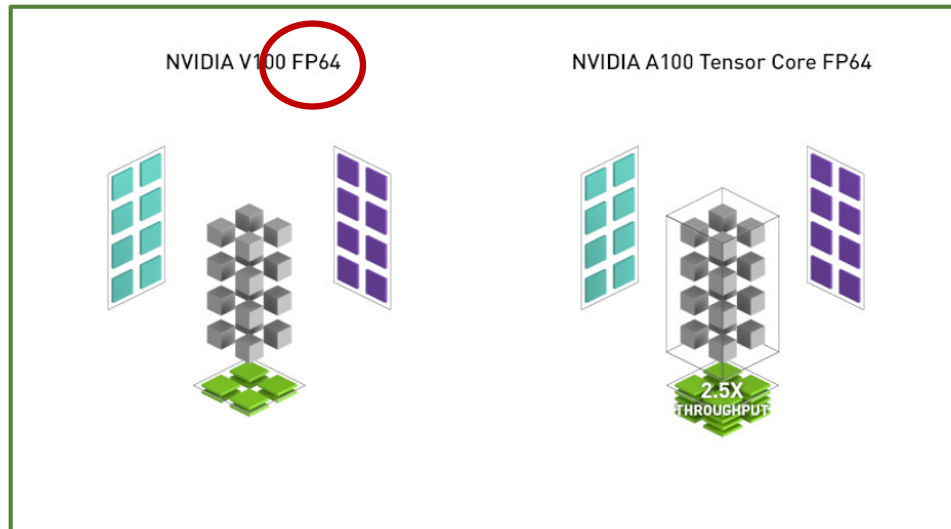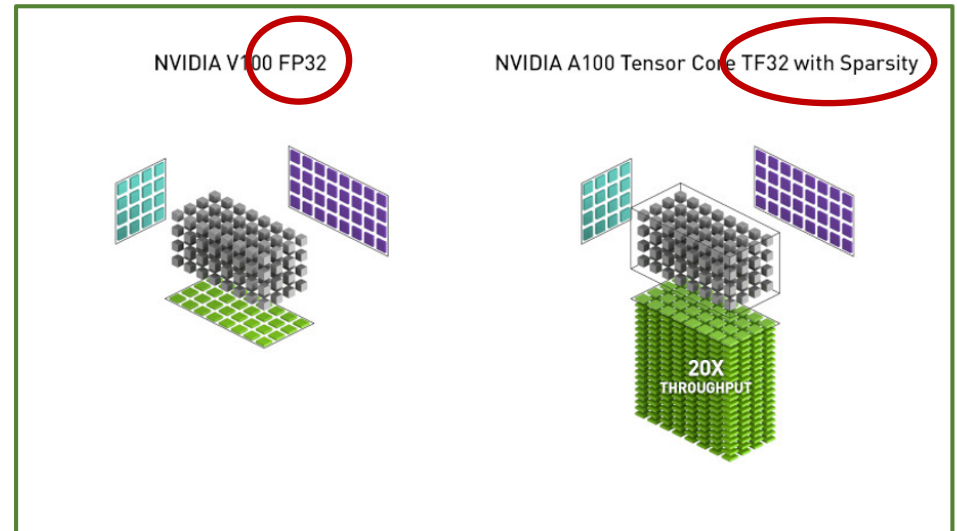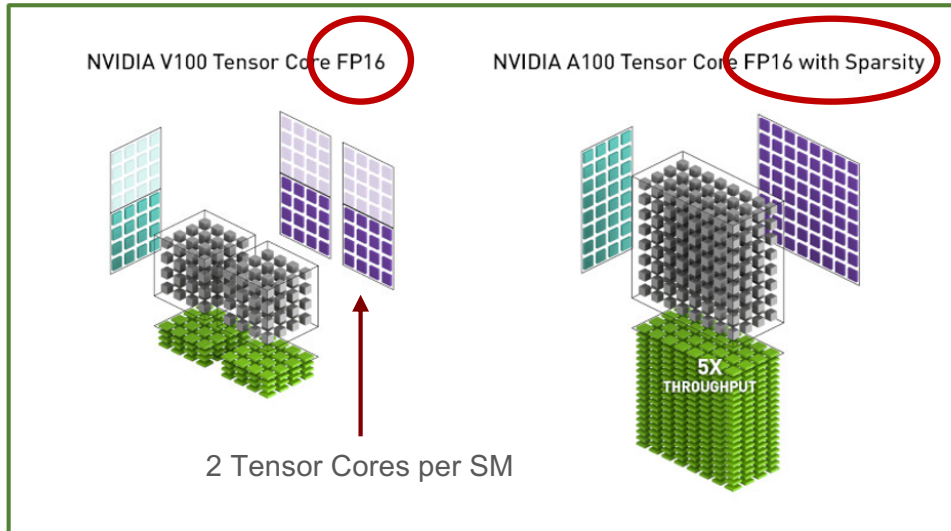https://www.nextplatform.com/2020/05/28/diving-deep-into-the-nvidia-ampere-gpu-architecture/

**TF32: same <u>range</u> as FP32 and same <u>precision</u> as FP16**
The FP multiplier scales with the square of the mantissa width ($8^2/11^2 \approx 0.5$)

https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/nvidia-ampere-architecture-whitepaper.pdf

NVIDIA V100 Tensor Core FP16     NVIDIA A100 Tensor Core FP16 with Sparsity

2 Tensor Cores per SM

5X THROUGHPUT

NVIDIA V100 FP32     NVIDIA A100 Tensor Core TF32 with Sparsity

20X THROUGHPUT

NVIDIA V100 FP64     NVIDIA A100 Tensor Core FP64

2.5X THROUGHPUT

NVIDIA V100 INT8     NVIDIA A100 Tensor Core INT8 with Sparsity

20X THROUGHPUT

https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/nvidia-ampere-architecture-whitepaper.pdf

| Tesla GPU | "Fermi" GF100 | "Fermi" GF104 | "Kepler" GK104 | "Kepler" GK110 | "Maxwell" GM200 | "Pascal" GP100 | "Volta" GV100 | "Turing" TU104 | "Ampere" GA100 |
|---|---|---|---|---|---|---|---|---|---|
| Compute Capability | 2.0 | 2.1 | 3.0 | 3.5 | 5.3 | 6.0 | 7.0 | 7.0 | 8.0 |
| Streaming Multiprocessors (SMs) | 16 | 16 | 8 | 15 | 24 | 56 | 84 | 72 | 128 |
| FP32 CUDA Cores / SM | 32 | 32 | 192 | 192 | 128 | 64 | 64 | 64 | 64 |
| FP32 CUDA Cores | 512 | 512 | 1,536 | 2,880 | 3,072 | 3,584 | 5,376 | 4,608 | 8,192 |
| FP64 Units | – | – | 512 | 960 | 96 | 1,792 | 2,688 | – | 4,096 |
| Tensor Core Units | | | | | | | 672 | 576 | 512 |
| Threads / Warp  **SIMT/SIMD instr** | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| Max Warps / SM  **SMT** | 48 | 48 | 64 | 64 | 64 | 64 | 64 | 64 | 64 |
| Max Threads / SM | 1,536 | 1,536 | 2,048 | 2,048 | 2,048 | 2,048 | 2,048 | 2,048 | 2,048 |
| Max Thread Blocks / SM | 8 | 8 | 16 | 16 | 32 | 32 | 32 | 32 | 32 |
| 32-bit Registers / SM | 32,768 | 32,768 | 65,536 | 65,536 | 65,536 | 65,536 | 65,536 | 65,536 | 65,536 |
| Max Registers / Thread | 63 | 63 | 63 | 255 | 255 | 255 | 255 | 255 | 255 |
| Max Threads / Thread Block | 1,024 | 1,024 | 1,024 | 1,024 | 1,024 | 1,024 | 1,024 | 1,024 | 1,024 |
| Shared Memory Size Configs | 16 KB | 16 KB | 16 KB | 16 KB | 96 KB | 64 KB | Config | Config | Config |
| | 48 KB | 48 KB | 32 KB | 32 KB | | | Up To | Up To | Up To |
| | | | 48 KB | 48 KB | | | 96 KB | 96 KB | 164 KB |

https://www.nextplatform.com/2020/05/28/diving-deep-into-the-nvidia-ampere-gpu-architecture/

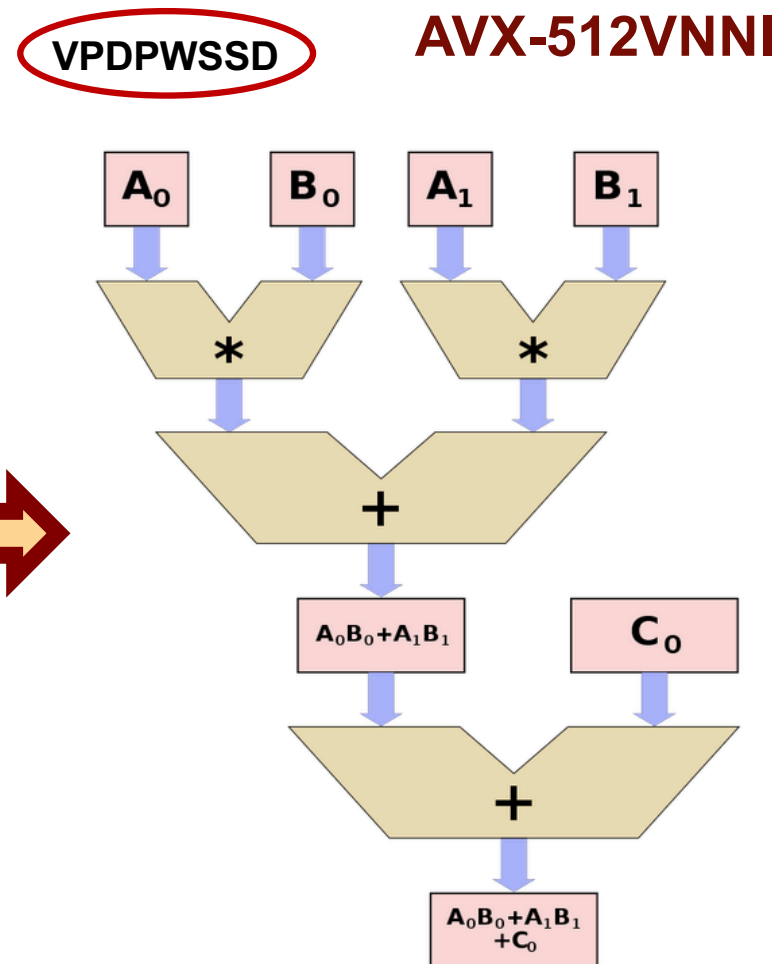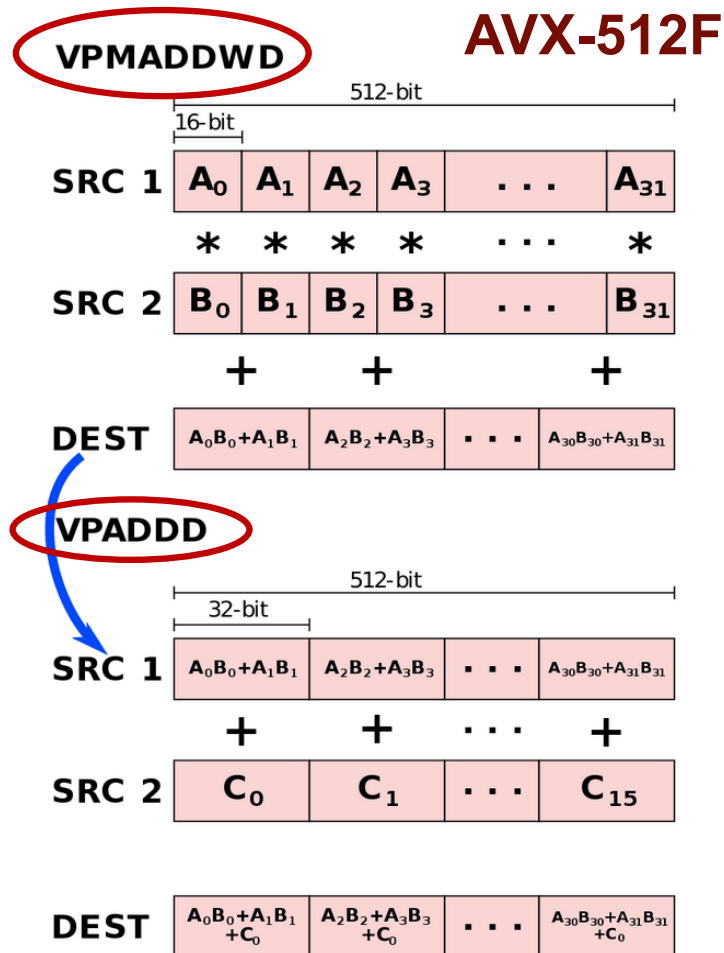# *Approaches to operations on tensors*

- **Tensor**: a mathematical object that describes the relationship between other mathematical objects that are all linked together; they are commonly shown as a <u>multidimensional array</u>

- Different approaches followed by chip manufacturers:
  - add <u>new extensions</u> to existing HPC vector devices
    - **NVidia**: tensor core units in HPC GPUs
    - **Intel**: AVX-512VNNI & AMX
  - develop SoC devices for embedded/specific application fields
    - neural net devices: **Google** TPU, **Intel** Habana, …
    - autonomous driving: **Tesla** FSD, **NVidia** Orin, …
    - smartphones: **Apple** A14 Bionic, **Huawei** Kirin 9000, Qualcomm Snapdragon, Samsung Exynos, …
    - gaming: …

**VNNI:** 2 new instr + 2 new extensions to merge previous set of 2 & 3 instr; ex.:



AVX-512F

AVX-512VNNI

VPMADDWD

VPADDD

VPDPWSSD

https://en.wikichip.org/wiki/x86/avx512_vnni

# The Intel Advanced Matrix Extension (AMX)
## *(expected in 2021)*

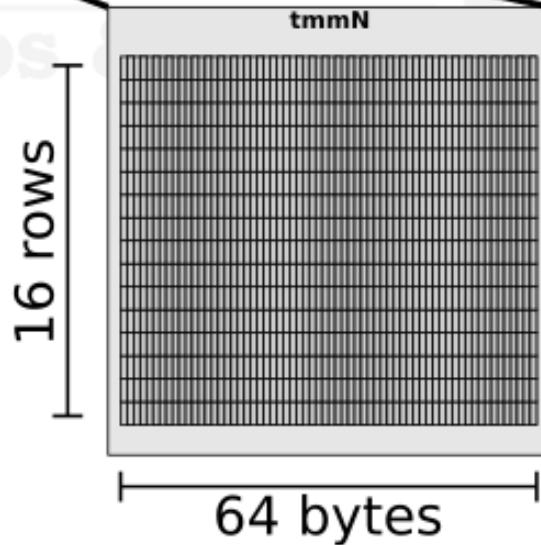**IA Host** ⟷ **Tiles and Accelerator Commands**

**Accelerator 1 (TMUL)** ← FMA

• • •

**Accelerator N**

**TILECFG**

**Tile RegFile**

| tmm0 |
| tmm1 |
| tmm2 |
| tmm3 |
| tmm4 |
| tmm5 |
| tmm6 |
| tmm7 |

**8 matrix registers**

**Each matrix reg is 1024 bytes long ( = 16 x 64B )**

Advanced Matrix Extension (AMX) is an x86 extension that introduces a matrix register file and new instructions for operating on matrices.

**tmmN**

16 rows

64 bytes

# AMX instructions in Accelerator 1

## Tile A

K

M

## Tile B

N

K

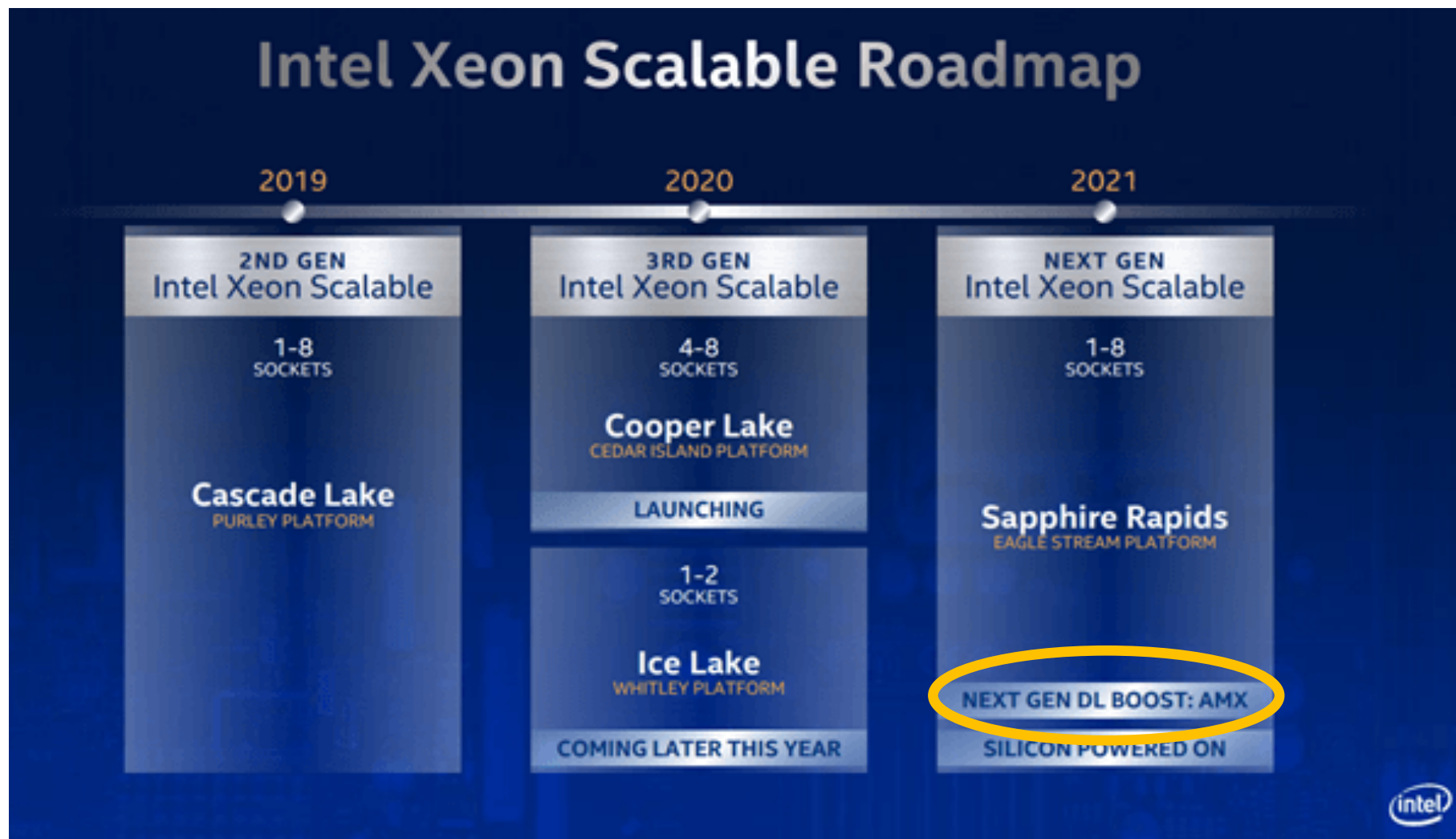| AMX Extensions | | |
|---|---|---|
| **Feature Set** | **Description** | **Instructions** |
| AMX-TILE | The base matrix tile architecture support. | 7 instructions |
| AMX-INT8 | Dot-product of Int8 tiles. | 4 instructions |
| AMX-BF16 | Dot-product of BF16 tiles. | 1 instruction |

Data types

## TMUL
### C += A * B

## Tile C

N

M

| | Sign | Range | Precision |
|---|---|---|---|
| FP32 | | 8 BITS | 23 BITS |
| | | TF32 Range | |
| TENSOR FLOAT 32 (TF32) | | 8 BITS | 10 BITS |
| | | | TF32 Precision |
| FP16 | | 5 BITS | 10 BITS |
| BFLOAT16 | | 8 BITS | 7 BITS |

*AJProença, Advanced Architectures, MiEI, UMinho, 2020/21*

Intel Xeon Scalable Roadmap

| 2019 | 2020 | 2021 |
|---|---|---|
| **2ND GEN** Intel Xeon Scalable | **3RD GEN** Intel Xeon Scalable | **NEXT GEN** Intel Xeon Scalable |
| 1-8 SOCKETS | 4-8 SOCKETS | 1-8 SOCKETS |
| Cascade Lake PURLEY PLATFORM | Cooper Lake CEDAR ISLAND PLATFORM — LAUNCHING | Sapphire Rapids EAGLE STREAM PLATFORM |
| | 1-2 SOCKETS — Ice Lake WHITLEY PLATFORM — COMING LATER THIS YEAR | NEXT GEN DL BOOST: AMX — SILICON POWERED ON |

# *Approaches to operations on tensors*

- **Tensor**: a mathematical object that describes the relationship between other mathematical objects that are all linked together; they are commonly shown as a <u>multidimensional array</u>

- Different approaches followed by chip manufacturers:
  - add <u>new extensions</u> to existing HPC vector devices
    - **NVidia**: tensor core units in HPC GPUs
    - **Intel**: AVX-512VNNI & AMX
  - develop SoC devices for embedded/specific application fields
    - neural net devices: **Google** TPU, **Intel** Habana,
    - autonomous driving: **Tesla** FSD, **NVidia** Orin, …
    - smartphones: **Apple** A14 Bionic, **Huawei** Kirin 9000, Qualcomm Snapdragon, Samsung Exynos, …
    - gaming: …

Google
**Tensor Processing Unit**

- The Matrix Unit: 65,536 (256x256) 8-bit multiply-accumulate units ← **INT8**
- 700 MHz clock rate **F M A**
- Peak: 92T operations/second
  - 65,536 * 2 * 700M → **92 TOPS**
- >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) **SRAM**
- 3.5X as much on-chip memory vs GPU
- Two 2133MHz DDR3 DRAM channels
- 8 GiB of off-chip weight DRAM memory

## TPU: High-level Chip Architecture

## Chip floor plan

TPU: a Neural Network Accelerator Chip



**Unified Buffer for Local Activations** (96Kx256x8b = 24 MiB) 29% of chip

**Matrix Multiply Unit** (256x256x8b=64K MAC) 24%

DRAM port ddr3 3%

Host Interf. 2%

**Accumulators** (4Kx256x32b = 4 MiB) 6%

Control 2%

Activation Pipeline 6%

PCIe Interface 3%

Misc. I/O 1%

DRAM port ddr3 3%



TPUs are intensively used by Google, namely in Google Photos, RankBrain, StreetView & Google Translate

https://cloud.google.com/blog/big-data/2017/05/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu

# Google TPUv2 (May'17)

## TPUv2 Chip

- 16 GB of HBM
- 600 GB/s mem BW
- Scalar unit: 32b float
- MXU: 32b float accumulation but reduced precision for multipliers
- 45 TFLOPS

bfloat

HBM 8 GB

core

scalar unit

MXU 128x128

core

scalar unit

MXU 128x128

HBM 8 GB

Google

**Tensor Processing Unit**

https://www.anandtech.com/show/16005/hot-chips-2020-live-blog-google-tpuv2-and-tpuv3-230pm-pt

TPU Core: Vector Unit (Lane)

https://www.anandtech.com/show/16005/hot-chips-2020-live-blog-google-tpuv2-and-tpuv3-230pm-pt

# TPU Core: Matrix Multiply Unit

- 128 x 128 systolic array
  - Streaming LHS and results
  - Stationary RHS (w/ optional transpose)
- Numerics
  - bfloat16 multiply
    - {s, e, m} = {1, 8, 7}
    - The original!
  - float32 accumulation

https://www.anandtech.com/show/16005/hot-chips-2020-live-blog-google-tpuv2-and-tpuv3-230pm-pt

# Google TPUv3 *(May'18)*

TPUv4 released Jun 2020
but no data available yet…

TPU v2 - 4 chips, 2 cores per chip

TPU v3 - 4 chips, 2 cores per chip

- Intel acquired **Nervana** Engine *(Aug 2016)*
- Intel launched Nervana NNP (Neural Net Processor) *(Oct 2017)*
- Key features: <u>matrix multiplication</u> & <u>convolution</u> *(for neural nets)*
- Intel discontinued Nervana NNP *(Jan 2020…)*

- Intel acquired the Israel chipmaker **Habana** Labs *(Dec 2019)*
  - Habana <u>training</u> chip **Gaudi**, with support to FP32, INT32, **BF16**, INT16, INT8, UINT32, UINT16, UINT8
  - Habana <u>inference</u> chip **Goya**, with support to FP32, INT32, INT16, INT8, UINT32, UINT16, UINT8

https://en.wikichip.org/wiki/habana/microarchitectures/gaudi

TPC: Tensor Processing Core
GEMM: General Matrix Multiply

# *Inference Processor: GOYA*

https://en.wikichip.org/wiki/habana/microarchitectures/goya

TPC: Tensor Processing Core
GEMM: General Matrix Multiply

# *Approaches to operations on tensors*

- **Tensor**: a mathematical object that describes the relationship between other mathematical objects that are all linked together; they are commonly shown as a multidimensional array

- Different approaches followed by chip manufacturers:
  - add new extensions to existing HPC vector devices
    - **NVidia**: tensor core units in HPC GPUs
    - **Intel**: AVX-512VNNI & AMX
  - develop SoC devices for embedded/specific application fields
    - neural net devices: **Google** TPU, **Intel** Habana, …
    - autonomous driving: **Tesla** FSD, **NVidia** Orin, …
    - smartphones: **Apple** A14 Bionic, **Huawei** Kirin 9000, Qualcomm Snapdragon, Samsung Exynos, …
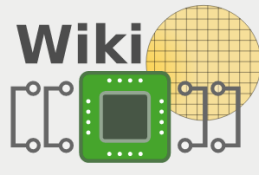    - gaming

# Tesla Full Self-Driving chip (FSD)



**ISP (24-bit)**

**Camera I/F**

**Safety System**

**Security System**

**GPU 1 GHz (600 GFLOPS)**

**Video Encode (H.265)**

**Quad-Core Cortex-A72 2.2 GHz**

**LPDDR4-4266 (64-bit)**

**NoC**

**NPU 2 GHz (36.86 TOPS)**

**NPU 2 GHz (36.86 TOPS)**

**Quad-Core Cortex-A72 2.2 GHz**

**Quad-Core Cortex-A72 2.2 GHz**

**LPDDR4-4266 (64-bit)**

https://en.wikichip.org/wiki/tesla_(car_company)/fsd_chip

# The Neural Processing Unit in FSD

# NVidia roadmap for Drive systems



DRIVE PX 2

DRIVE Xavier

DRIVE Pegasus

Orin

DRIVE PX Parker

Xavier SoC

**16 CSI**
109 Gbps
1gE & 10gE

**DLA**
5 TFLOPS FP16
10 TOPS INT8

**Video Processor**
1.2 GPIX/s Encode
1.8 GPIX/s Decode

**PVA**
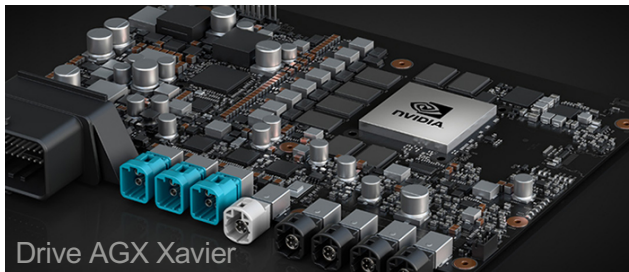1.6 TOPS
Stereo Disparity
Optical Flow
Image Processing

**ISP**
1.5 GPIX/s
Native Full-range HDR
Tile-based Processing

**Volta GPU**
FP32 / FP16 / INT8 Multi Precision
512 CUDA Cores
1.3 CUDA TFLOPS
20 Tensor Core TOPS

**Carmel ARM64 CPU**
8 Cores
10-wide Superscalar
2700 SpecInt2000
Functional Safety Features
Dual Execution Mode
Parity & ECC

**256-Bit LPDDR4**
137 GB/s

Jun'18

# NVidia Xavier SoC



Drive AGX Xavier

256-bit LPDDR4X

NVLink

PCIe 4

MM/DLA

PVA

Volta

Southbridge

Xavier SoC

Octa-core Carmel

dGPU

NVLink

Xavier SoC 2

PCIe4

## Deep Learning Accelerator

Control Bus

Configuration and Control

Convolutional Buffer (SRAM)
- Input activations
- Filter weights

Convolution Core

Post-processing

Memory Interface

SDRAM

Internal RAM

**D**eep **L**earning **A**ccelerator
Large array of multiply-accumulate units optimized for CNNs

### Volta GPU

| SM | SM | SM | SM |
|---|---|---|---|
| 128 KiB L1 | 128 KiB L1 | 128 KiB L1 | 128 KiB L1 |

512 KiB L2

| 128 KiB L1 | 128 KiB L1 | 128 KiB L1 | 128 KiB L1 |
|---|---|---|---|
| SM | SM | SM | SM |

### CPU Complex

| CPU | CPU | CPU | CPU |
|---|---|---|---|
| 2 MiB L2 | | 2 MiB L2 | |

4 MiB L3

| 2 MiB L2 | | 2 MiB L2 | |
|---|---|---|---|
| CPU | CPU | CPU | CPU |

https://fuse.wikichip.org/news/1618/hot-chips-30-nvidia-xavier-soc/

*AJProença, Advanced Architectures, MiEI, UMinho, 2020/21*

# NVidia Drive SoC's: Parker, Xavier, Orion

| NVIDIA ARM SoC Specification Comparison | | | |
|---|---|---|---|
| | **Orin 2021?** | **Xavier 2018** | **Parker 2016** |
| **CPU Cores** | 12x Arm "Hercules" Cortex-A78AE * | 8x NVIDIA Custom ARM "Carmel" | 2x NVIDIA Denver + 4x Arm Cortex-A57 |
| **GPU Cores** | Ampere iGPU (?? cores) | Xavier Volta iGPU (512 CUDA Cores) | Parker Pascal iGPU (256 CUDA Cores) |
| **INT8 DL TOPS** | 200 TOPS | 30 TOPS | N/A |
| **FP32 TFLOPS** | ? | 1.3 TFLOPs | 0.7 TFLOPs |
| **Manufacturing Process** | 7nm? | TSMC 12nm FFN | TSMC 16nm FinFET |
| **TDP** | ~5-45W | 30W | 15W |

* AE, *Automotive Enhanced*: improved functional security and safety

**NVidia Partners**

CARS

# *Approaches to operations on tensors*

- **Tensor**: a mathematical object that describes the relationship between other mathematical objects that are all linked together; they are commonly shown as a multidimensional array

- Different approaches followed by chip manufacturers:
  - add new extensions to existing HPC vector devices
    - **NVidia**: tensor core units in HPC GPUs
    - **Intel**: AVX-512VNNI & AMX
  - develop SoC devices for embedded/specific application fields
    - neural net devices: **Google** TPU, **Intel** Habana, …
    - autonomous driving: **Tesla** FSD, **NVidia** Orin, …
    - smartphones: **Apple** A14 Bionic, **Huawei** Kirin 9000, Qualcomm Snapdragon, Samsung Exynos, …
    - gaming: …

# Evolution of Apple A series

**Bionic** => with a neural engine

**A14**
**6-core CPU**
2 FireStorm cores (high-performance)
4 IceStorm cores (energy-efficient)

**A14**
**Neural Engine**
16 cores

**A14**
ML accelerators    AMX blocks

**A14**    4-core GPU
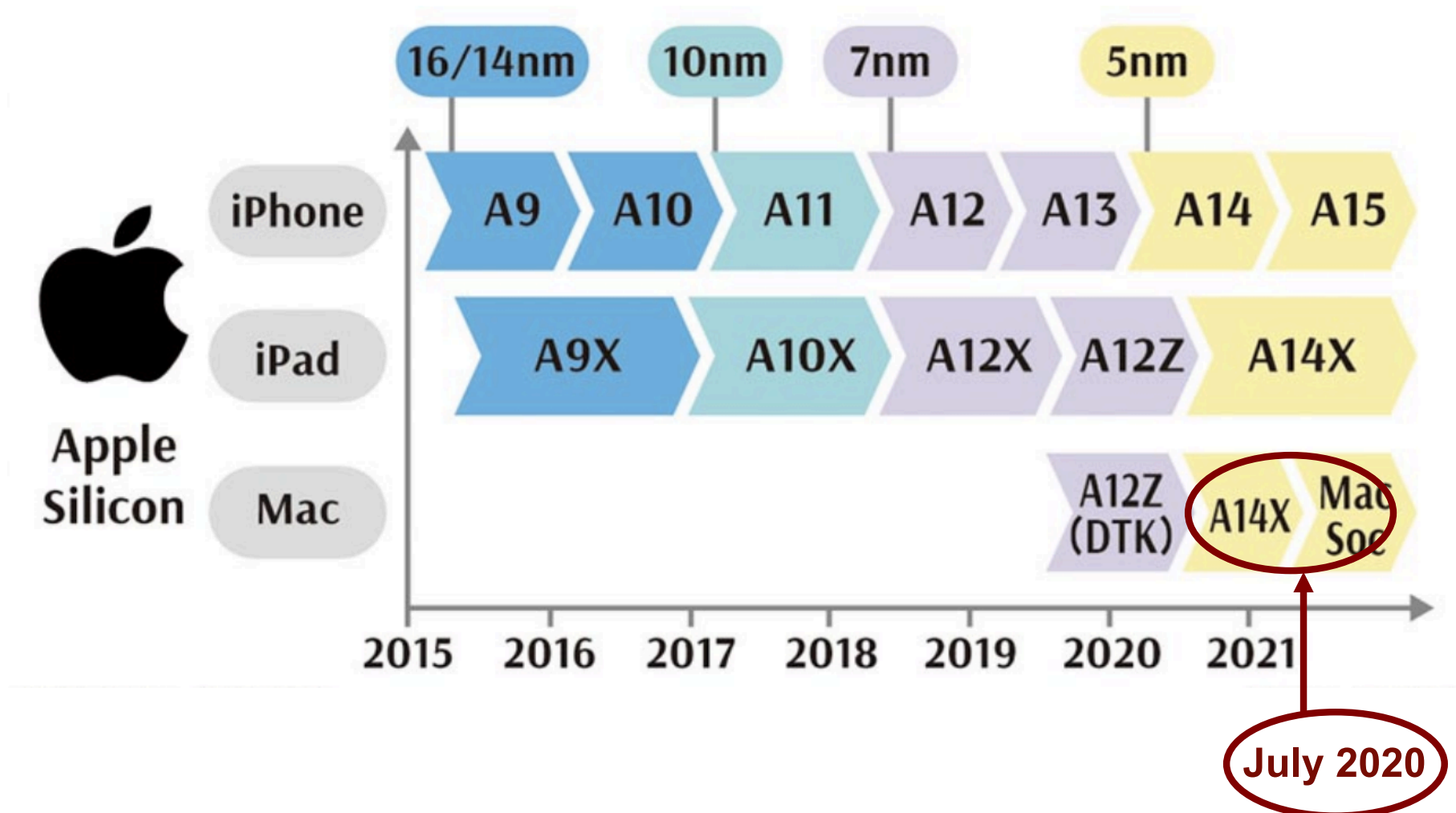4 cores (Apple-designed)

Apple-designed
64-bit six-core CPU
implementing ARMv8.4 ISA
8 MiB L2 cache

16-core Neural Engine +
2nd gen ML accelerators (AMX) +
high-performance 4-core GPU  =>
powerful image recognition, natural language learning, motion analysis, …

*Apple M1:*
*an extended version of A14*

Nov'20

Top panel (A14):
- Machine learning controller
- New 6-core CPU
- Next-generation ML accelerators
- 16-core NEURAL ENGINE
- 5 nanometer process
- A14
- 11 trillion Operations per second
- 11.8 billion Transistors
- Advanced image signal processor
- New 4-core GPU
- Secure Enclave

Bottom panel (M1):
- 5 nanometer process
- Machine learning accelerators
- 16-core Neural Engine — 11 trillion operations per second
- Thunderbolt / USB 4 controller
- Media encode and decode engines
- M1
- Up to 8-core GPU
- 8-core CPU
- 4 to 8
- 6 to 8
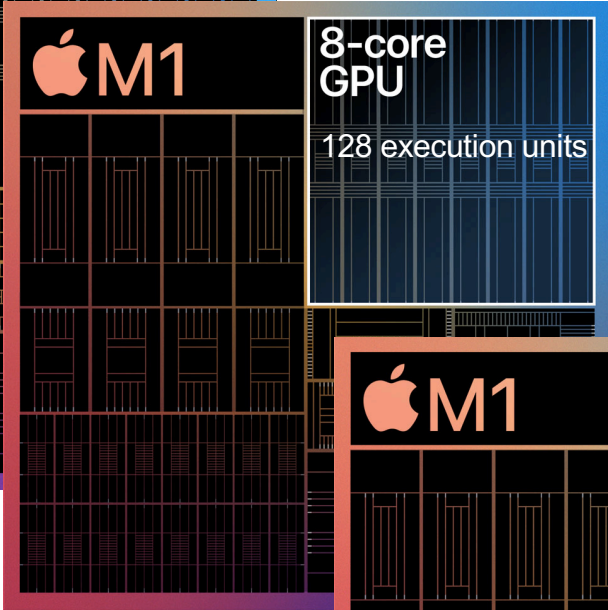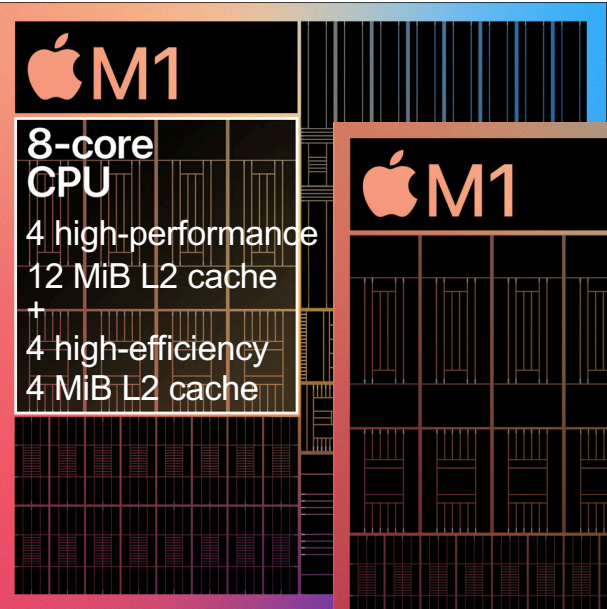- 16 billion transistors
- Advanced image signal processor
- Secure Enclave
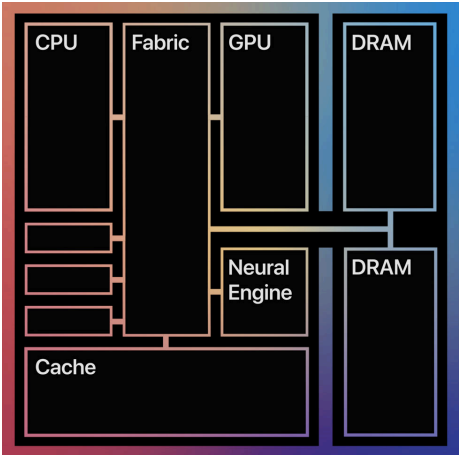- Unified memory architecture
- Industry-leading performance per watt

Apple M1 SoC

M1

8-core CPU
4 high-performance
12 MiB L2 cache
+
4 high-efficiency
4 MiB L2 cache

8-core GPU
128 execution units

16-core Neural Engine

5-nanometer process

16 billion transistors

8-Core GPU

8x 16b LPDDR4X Channels

SLC Cache

4 Firestorm Perf Cores +12MB L2

16-Core Neural Engine

4 Icestorm Efficiency Cores +4MB L2

ANANDTECH

CPU | Fabric | GPU | DRAM

Neural Engine | DRAM

Cache

*AJProença, Advanced Architectures, MiEI, UMinho, 2020/21*

Will the first Apple silicon laptop be faster than every Intel MacBook Pro?

Geekbench 5 data from Primate Labs

A77: old ARM generation…



## Kirin 9000

| CPU | NPU | GPU |
|---|---|---|
| **1** fast<br>Cortex-A77@ 3.13 GHz | **2**<br>Big-Core | **24 core**<br>Mali-G78 |
| **3** medium<br>Cortex-A77@ 2.54 GHz | **+** | |
| **4** efficient<br>Cortex-A55@ 2.05 GHz | **1**<br>Tiny-Core | |

5nm EUV

| | |
|---|---|
| Balong 5000 Modem | ISP 6.0 |
| LPDDR 5/4X | UFS |
| HiFi Audio | 4K HDR Video |
| Mobile Secure Processor | |

And after this…

Yet some more odd approaches:
- SoC with reconfigurable components: **Xilinx** ACAP
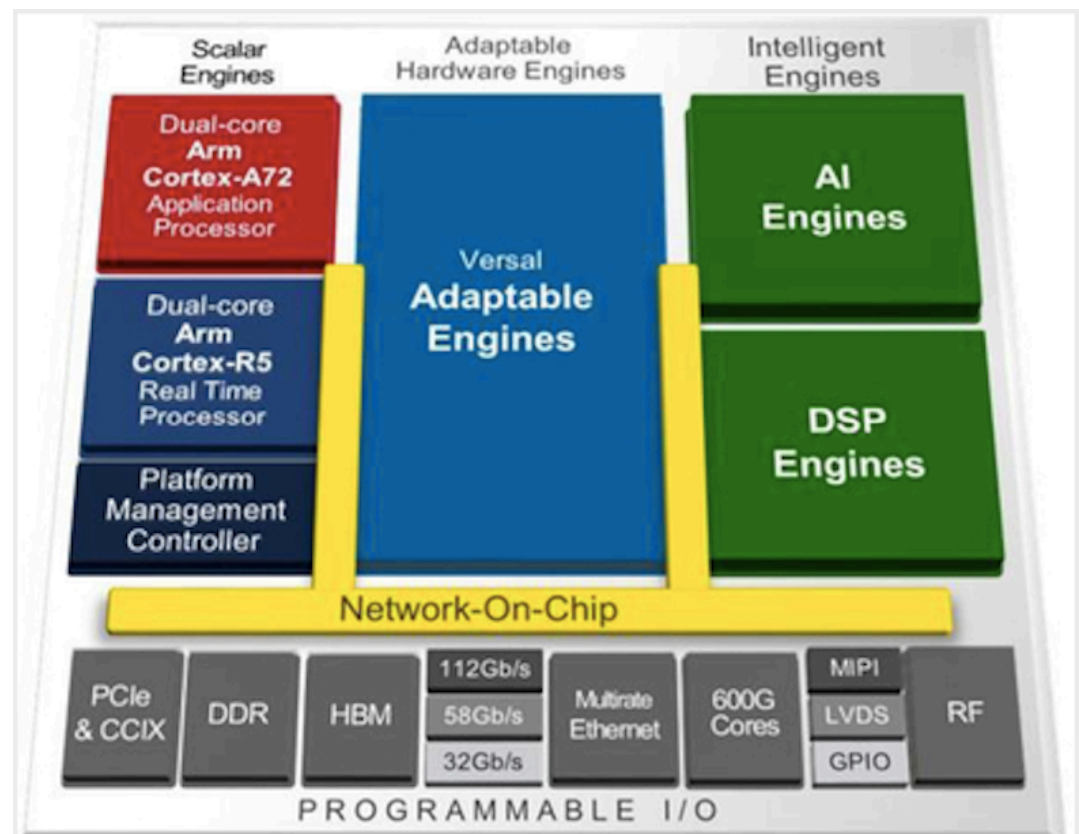- a system based on a **very large** "chip": **Cerebras**

## ACAP die, an adaptable accelerator-fabric ecosystem with:

- multicore ARM SoC's
- an FPGA fabric with distributed memory
- hw-programmable DSP engines
- AI Engines with vector units
- other specialized accelerators
- a flexible NoC interconnection

AI Engine: Xilinx Reinvents Multi-Core Compute

**Traditional Multi-core** (cache-based architecture)

**AI Engine Array** (intelligent engine)

Fixed, shared Interconnect
- Blocking limits compute
- Timing not deterministic

Data Replicated
- Robs bandwidth
- Reduces capacity

Dedicated Interconnect
- Non-blocking
- Deterministic

Local, Distributed Memory
- No cache misses
- Higher bandwidth
- Less capacity required

Page 16

https://forums.xilinx.com/t5/Adaptable-Advantage-Blog/Xilinx-Unveiled-the-Secret-Sauce-of-the-Ultimate-AI-Inference/ba-p/898197

AI Engine: Tile-Based Architecture

**Non-Blocking Interconnect**
Up to 200+ GB/s bandwidth per tile

**Local Memory**
Multi-bank implementation
Shared across neighbor cores

**ISA-based Vector Processor**
Software Programmable
(e.g., C/C++)

**Cascade Interface**
Partial results to next core

**Data Mover**
Non-neighbor data communication
Integrated synchronization primitives

https://forums.xilinx.com/t5/Adaptable-Advantage-Blog/Xilinx-Unveiled-the-Secret-Sauce-of-the-Ultimate-AI-Inference/ba-p/898197

# AI Inference Mapping on Versal ACAP

A = Activations
W = Weights

$$\begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \times \begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bma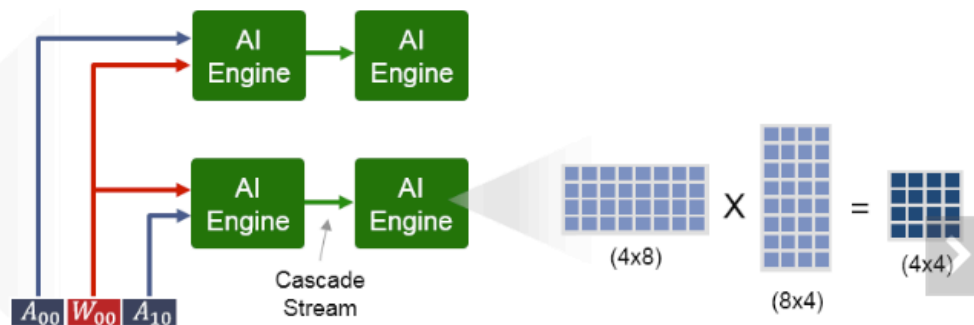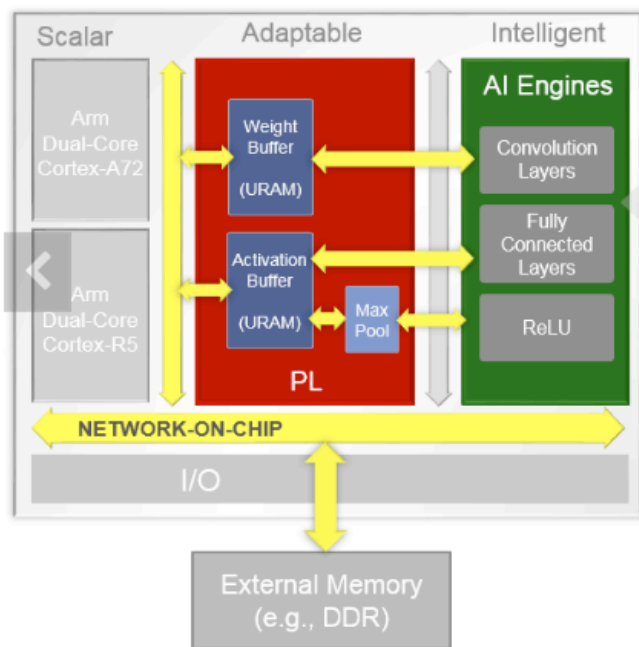trix} = \begin{bmatrix} A_{00} \times W_{00} + A_{01} \times W_{10} & \dots \\ A_{10} \times W_{00} + A_{11} \times W_{10} & \dots \end{bmatrix}$$

> Custom memory hierarchy
> > Buffer on-chip vs off-chip; Reduce latency and power
> Broadcast on AI interconnect (Weights and Activations)
> > Read once: reduce memory bandwidth
> AI-optimized vector instructions (128 INT8 mults/cycle)

Page 25

https://forums.xilinx.com/t5/Adaptable-Advantage-Blog/Xilinx-Unveiled-the-Secret-Sauce-of-the-Ultimate-AI-Inference/ba-p/898197

**Cerebras** Wafer Scale Engine (WSE):
the largest chip ever built



**46,225 mm$^2$ chip**
56x larger than the biggest GPU ever made
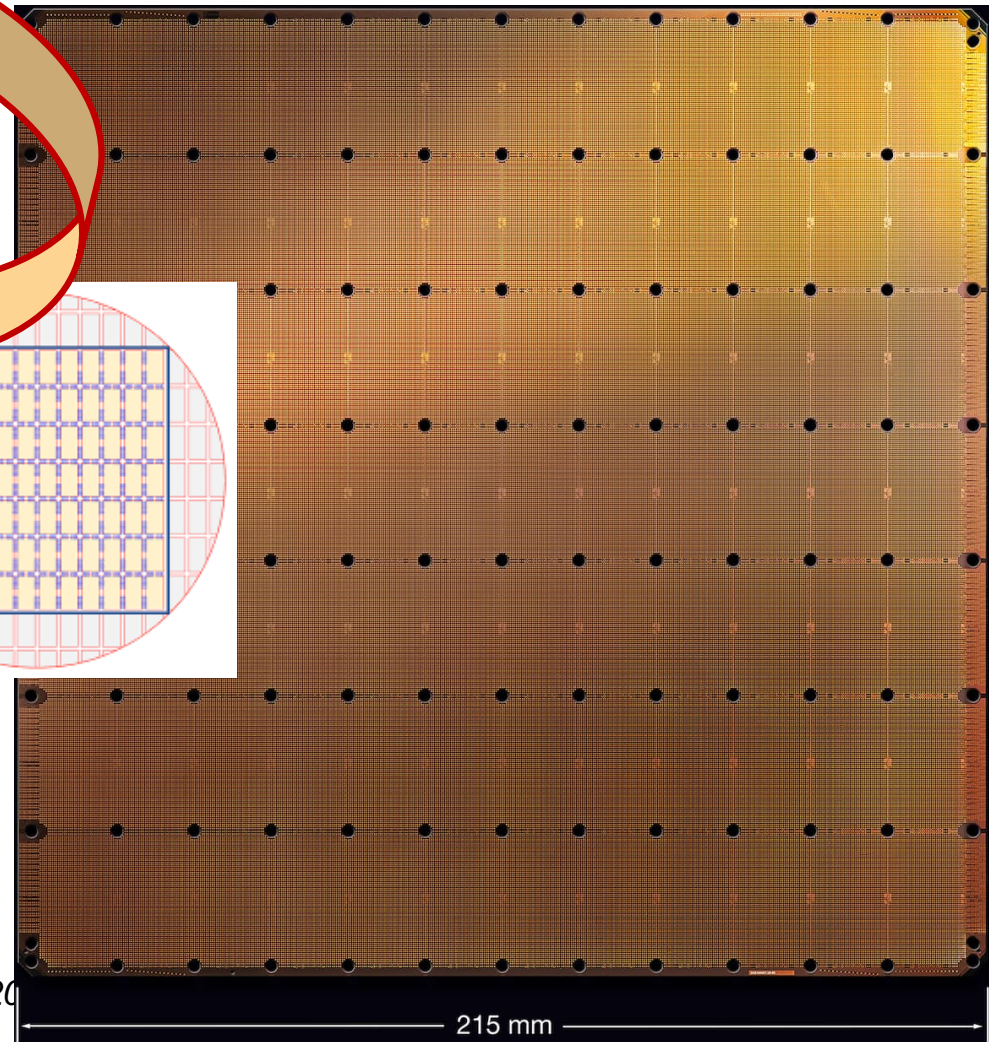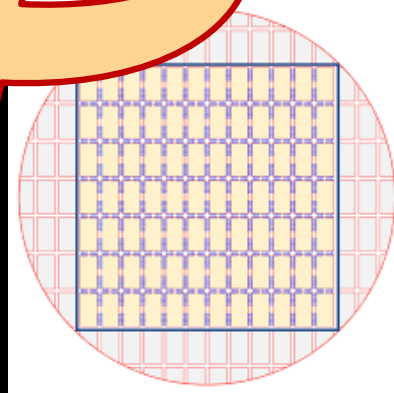
**400,000 core**
78x more cores

**18 GB on-chip SRAM**
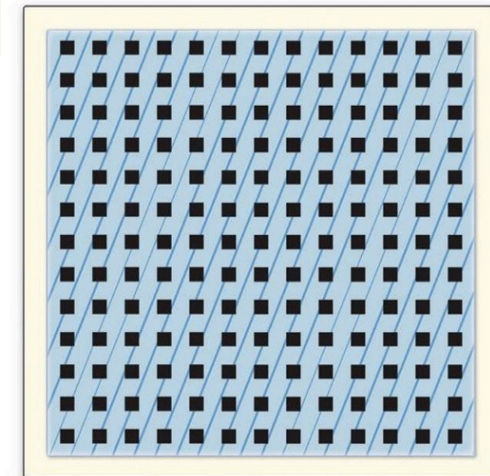3000x more on-chip memory

**100 Pb/s interconnect**
33,000x more bandwidth

215 mm

*AJProença, Advanced Architectures, MiEI, UMinho, 2020*

## Some additional data:

- 16 nm process

- 215 mm x 215 mm, ~15 kW consumption !

- 84 individual chips (12 wide by 7 tall)

- each chip:
  - 225 MiB SRAM
  - 54 x 94 = 5,076 Sparse Linear Algebra cores (SLA)
    (2 cores per row/column unused due to repair scheme leaving 4,888 usable cores)

- each core:
  - 47 kiB SRAM
  - Zeros not loaded from memory and zeros not multiplied
  - FP32 precision and scalar execution (can't filter zeros from memory with SIMD)
  - FMAC datapath with peak 8 operations per cycle)
  - Tensor control unit feeds the FMAC datapath with strided accesses
    (from memory or inbound data from links)
  - 4x 8 GB/s bidirectional links to its neighbours



Memory uniformly distributed across cores

■ Core   ▨ Memory

# Architecture Designed for Deep Learning

## Each component optimized for AI compute

### Compute
- Fully-programmable core, ML-optimized extensions
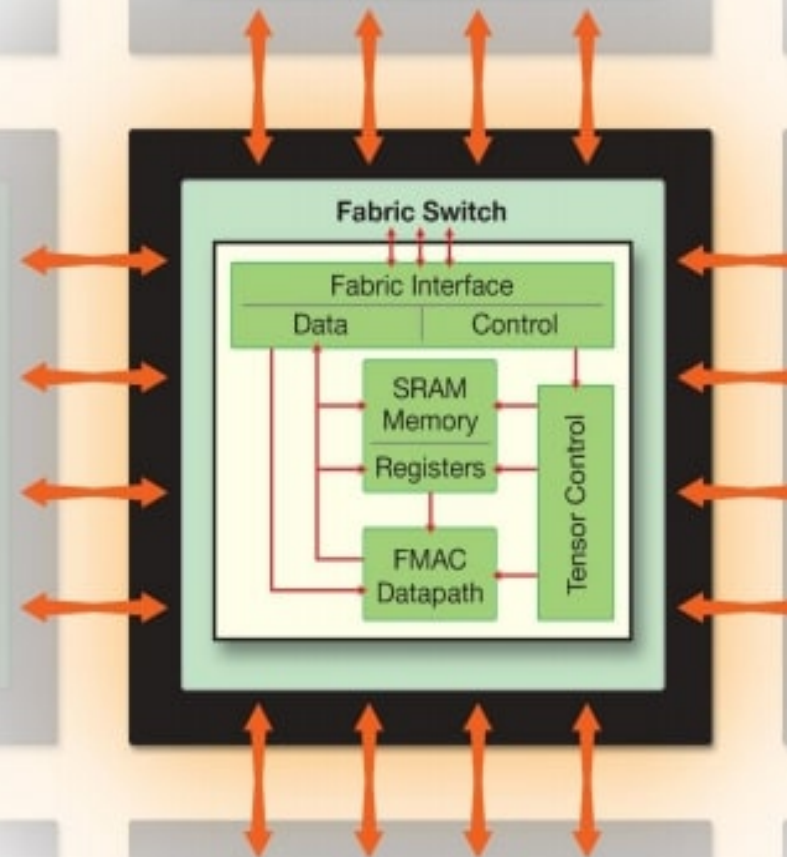- Dataflow architecture for sparse, dynamic workloads

### Memory
- Distributed, high performance, on-chip memory

### Communication
- High bandwidth, low latency fabric
- Cluster-scale networking on chip
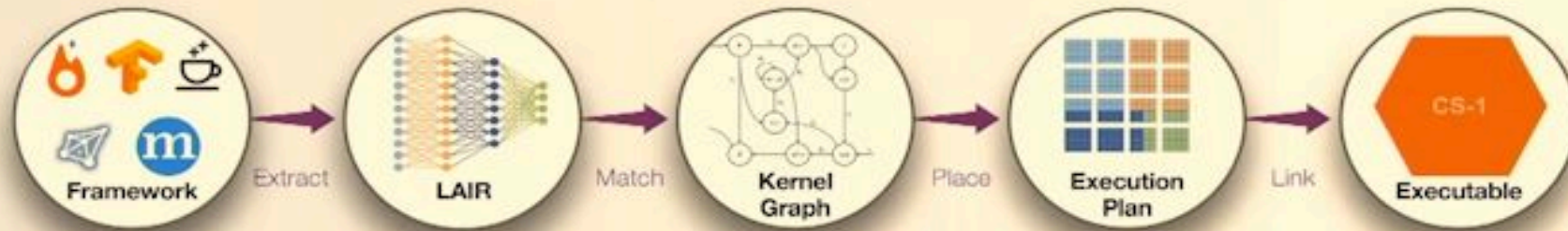- Fully-configurable to user-specified topology

**Together**, orders of magnitude performance and efficiency gain
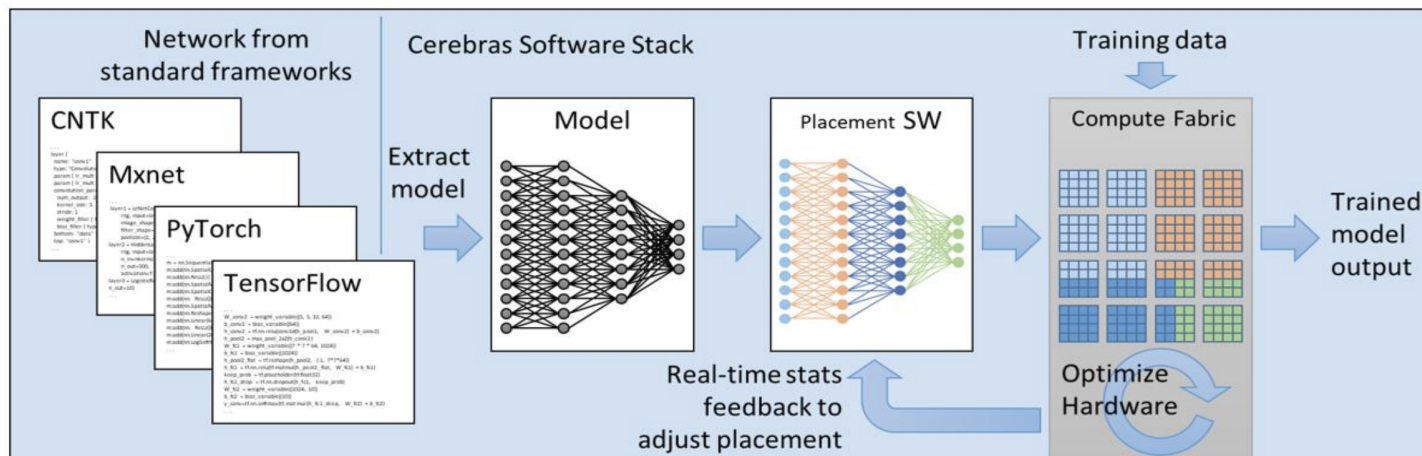
**Linear cluster-scale performance** on a single chip



https://www.anandtech.com/show/16006/hot-chips-2020-live-blog-cerebras-wse-programming-300pm-pt

Programming the Wafer-Scale Engine

Framework → Extract → LAIR → Match → Kernel Graph → Place → Execution Plan → Link → Executable (CS-1)

Users **program the WSE using standard ML frameworks**, e.g. TensorFlow, PyTorch

**Cerebras Graph Compiler automatically compiles the DNN graph**

• Extracts from Framework, converts to Cerebras IR, performs matching to Cerebras kernels
• Place & Route allocates compute and memory, configures on-chip network



Network from standard frameworks — CNTK, Mxnet, PyTorch, TensorFlow → Extract model → Cerebras Software Stack: Model → Placement SW → Compute Fabric → Trained model output

Training data

Real-time stats feedback to adjust placement ← Optimize Hardware

https://www.anandtech.com/show/16006/hot-chips-2020-live-blog-cerebras-wse-programming-300pm-pt

**SOFTWARE PLATFORM**

ML Frameworks — TensorFlow, PyTorch

Cerebras Software Platform
Cerebras Intermediate Representation
Cerebras Graph Compiler

Wafer-Scale Engine

TensorFlow → Extract → LAIR → Match → Kernel Graph → Place & Route → Execution Plan → Link → WSE Executable

Figure 8: A high-level overview of the compilation process for the WSE

**Cerebras WSE**

**2. Engine Block**

Power, cooling and packaging solution for the Wafer-Scale Engine

Wafer Scale Engine – Generation 2

850,000 **AI-optimized cores**

2.6 Trillion **Transistors**

TSMC 7nm **Process**

Cerebras executive Sean Lie

https://www.anandtech.com/show/16006/hot-chips-2020-live-blog-cerebras-wse-programming-300pm-pt