

DOJO

The Microarchitecture of Tesla's Eye-Scale Computer

Emil Talpeanu, Jonathan Williams, Debjit Das Sarma

架构

AI 芯片 - NPU详解



Talk Overview

1. AI 计算体系

- 深度学习计算模式
- 计算体系与矩阵运算

2. AI 芯片基础

- 通用处理器 CPU
- 从数据看 CPU 计算
- 通用图形处理器 GPU
- AI专用处理器 NPU/TPU
- 计算体系架构的黄金10年

1. 华为昇腾 NPU

- 达芬奇架构
- 昇腾AI处理器

2. 谷歌 TPU

- TPU 核心脉动阵列
- TPU 系列架构

3. 特斯拉 DOJO

- DOJO 架构

4. 国内外其他AI芯片

- AI芯片的思考

Talk Overview

I. 基本内容

- DOJO 整体架构
- DOJO Core架构
- DOJO Core 前端处理
- DOJO Core 执行引擎
- SRAM 与内存
- 内核与物理实现
- 问题与思考



DOJO

整体架构介绍

DOJO 超级计算机系统

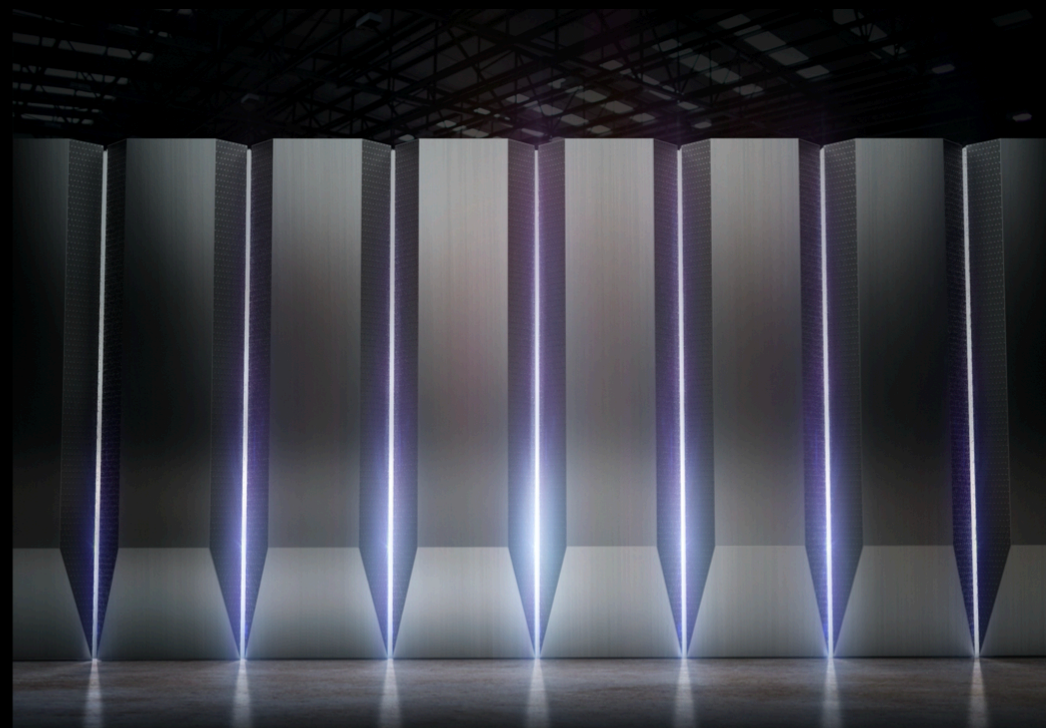
Tesla's in-house supercomputer for Machine Learning

Highly scalable and fully flexible distributed system

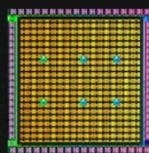
- Optimized for Neural Network training workloads
- General-purpose system capable of adapting to new algorithms and applications

Built from grounds up with large systems in mind

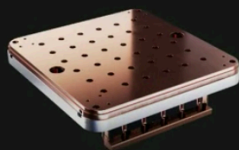
- Not evolved from existing small systems



DOJO 超级计算机系统



D2



TILE
V2



DIP
V2



DNIC
V2



PyTorch

DOJO EXTENSION

DOJO COMPILER ENGINE

LLVM

DOJO DRIVERS

DOJO 超级计算机系统

- 每个 DOJO 节点都有一个内核，具有CPU专用内存和 I/O接口。每个内核拥有一个 1.25MB 的 SRAM 作为主存。其中 SRAM 能以 400GB/s 速度加载，并以 270GB/s 存储。

High throughput, general purpose CPU

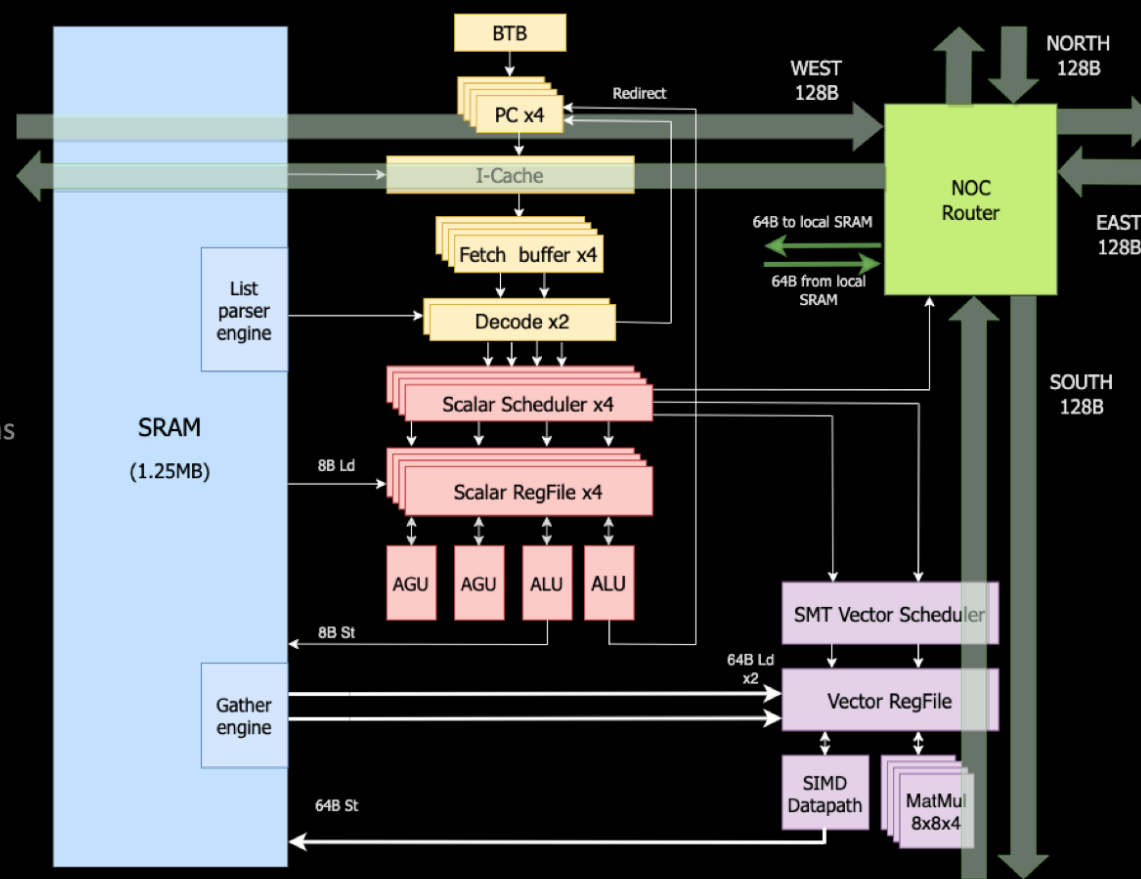
DOJO nodes are full-fledged computers

- Dedicated CPU, local memory, communication interface

Superscalar, multi-threaded organization

- Optimized for high-throughput math applications rather than control heavy code

Custom ISA optimized for ML kernels



DOJO 超级计算机系统

- DI 采用台积电 7nm 制程工艺，645mm² 面积上拥有 500 亿颗晶体管，BF16、CFP8 算力可达 362TFLOPS，FP32 算力可达 22.6TFLOPS，TDP（热设计功耗）为 400W。

TSMC 7nm, 645mm²

Physically and logically arranged as a 2D array

- 354 DOJO processing nodes on die

Extremely modular design

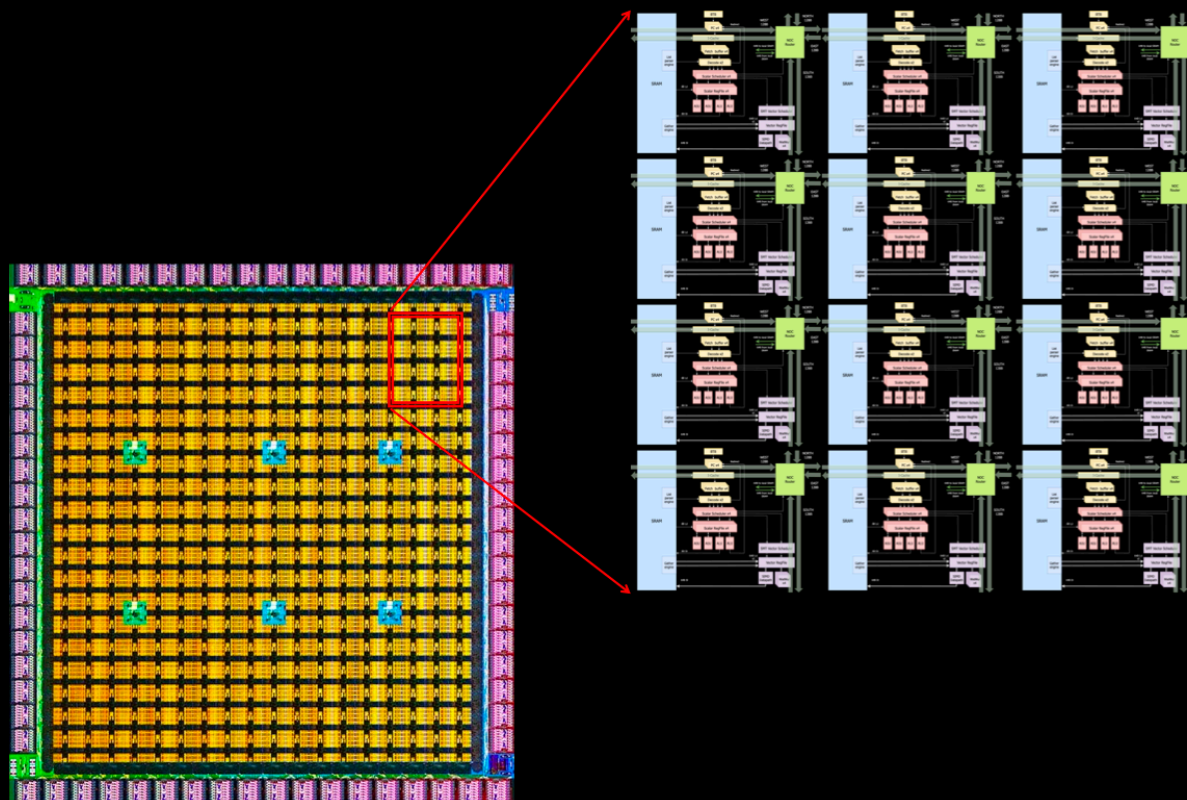
362 TFlops BF16/CFP8, 22 TFlops FP32 @2GHz

440 MB SRAM

Custom low power serdes channels on all edges

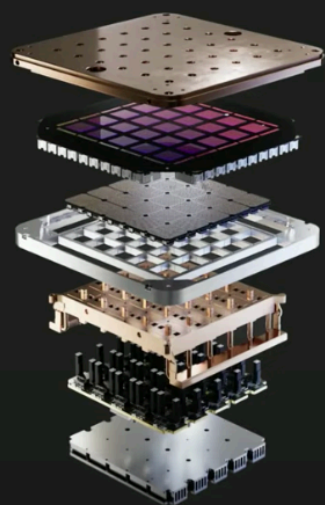
- 576 bidirectional channels
- 2 TB bandwidth on each edge

Seamless connection to neighboring dies



DOJO 超级计算机系统

基于DI芯片，特斯拉推出晶圆上系统级方案，通过应用台积电SoW封装技术，将所有25颗DI裸片都集成到一个训练Tile上，每个Dojo训练Tile消耗15kW。特斯拉Dojo训练Tile中有计算、I/O、功率和液冷模块。



TRAINING TILE
Compute + I/O + Power + Cooling

CHIP

PACKAGE

BOARD

SYSTEM

DOJO 超级计算机系统

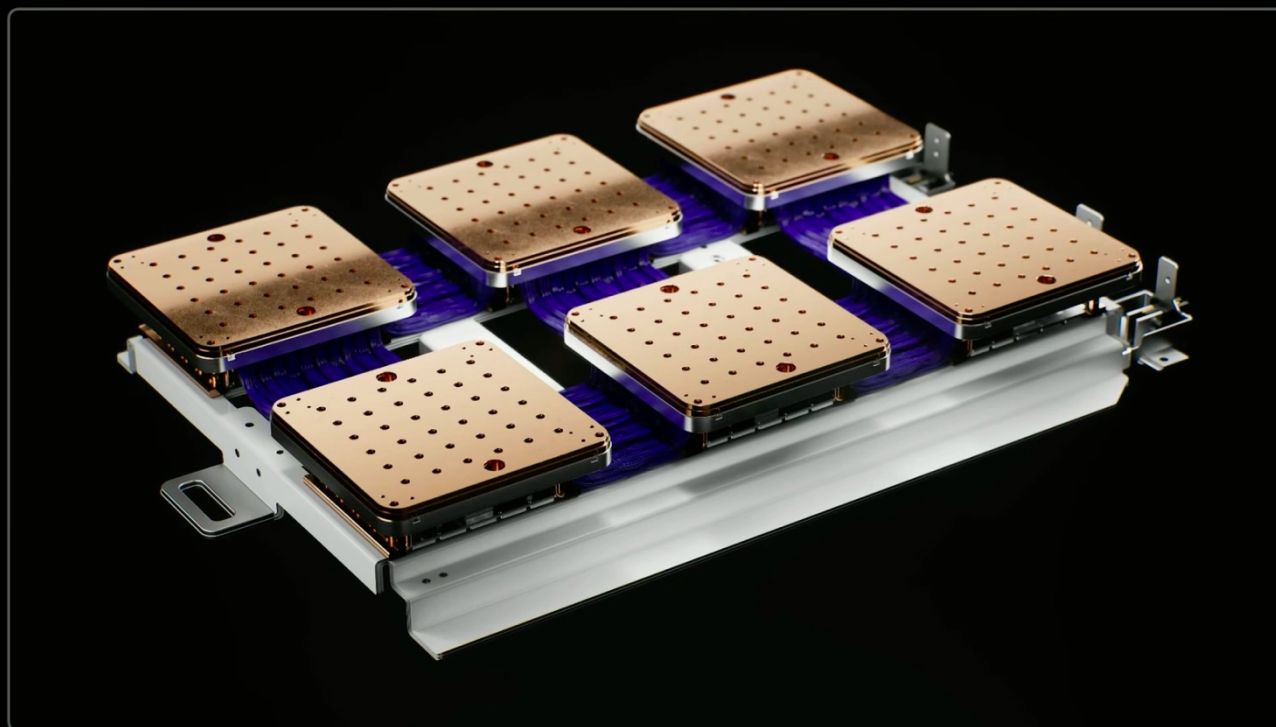
- DOJO System Tray 有高速连接、密集集成等特性，75mm高度能支持135kg。其BF16/CFP8峰值算力可达到54TFLOPS，功耗100+kW。

HIGH-SPEED CONNECTIVITY

POWER + MECHANICAL + THERMAL
2000A at 52VDC

DENSE INTEGRATION
75mm height to support 135kg

54 PFLOPS BF16/CFP8
13.4 TB/S BISECTION BW
100+ KW POWER



DOJO 超级计算机系统

- DOJO 接口处理器是一个具有高带宽内存PCIe卡，使用特斯拉传输协议TTP，主机 CPU 和训练Tile 之间的桥梁。每个 DIP 有 32GB 的 HBM。

HIGH-BANDWIDTH MEMORY FOR TRAINING

800 GB/s Total Memory Bandwidth

FULL BANDWIDTH MEMORY + INGEST TO TILE

Tesla Transport Protocol (TTP) -
Full custom protocol

TTP OVER ETHERNET (TTPOE)

Enables extending communication over standard Ethernet
Native hardware support

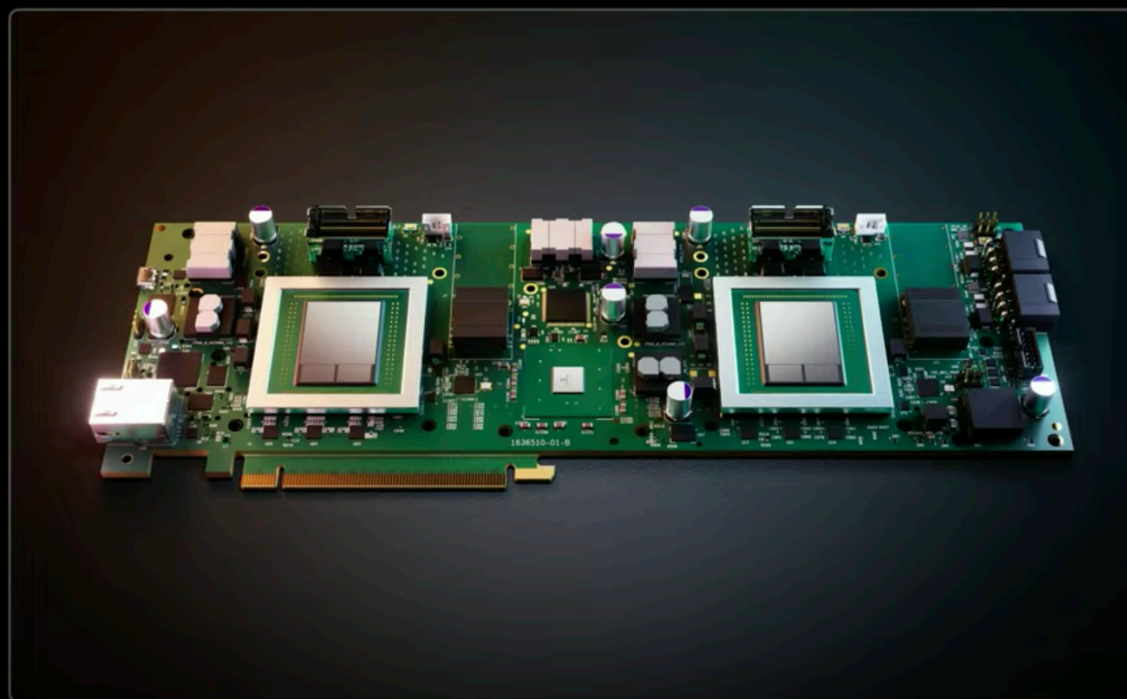
STANDARD PCIE HOST INTERFACE

32 GB HIGH-BANDWIDTH DRAM

900 TB/S TTP BANDWIDTH

50 GB/S ETHERNET BANDWIDTH

32 GB/S GEN4 PCIE BANDWIDTH



DOJO 超级计算机系统

- TTPOE 可将标准以太网转换至 Z 平面拓扑，拥有高 Z 平面拓扑连。最多可以将 5 个 DIP 以 900GB/s 的速度连接到一个训练瓦片上，以达到 4.5TB/s 的总量，每个训练 Tile 共有 160GB 的 HBM。

DISAGGREGATED HIGH SPEED MEMORY

Standard PCIe form factor
20x cards per tray & 32GB HBM per card

HIGH-BANDWIDTH INGEST

PCIe & Ethernet connectivity

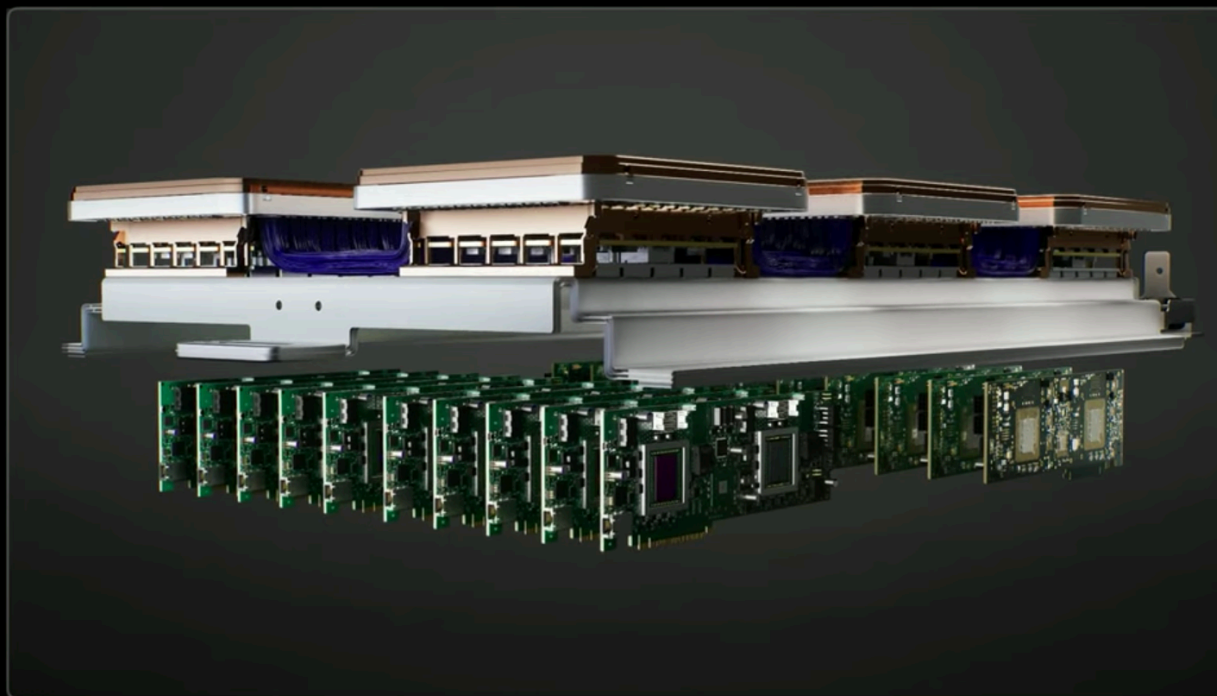
HIGH RADIX Z-PLANE CONNECTIVITY

Shortcuts across compute plane

640 GB HIGH-BANDWIDTH DRAM

1 TB/S ETHERNET BANDWIDTH

18 TB/S AGGREGATE BANDWIDTH TO TILES



DOJO 超级计算机系统

- DOJO 主机接口

INGEST PROCESSING

PCIe connectivity to Interface Processors
Hardware Video Decoder Support

USER APPLICATIONS

x86 Linux Environment
User-Scheduled jobs

512 TOTAL X86 CORES

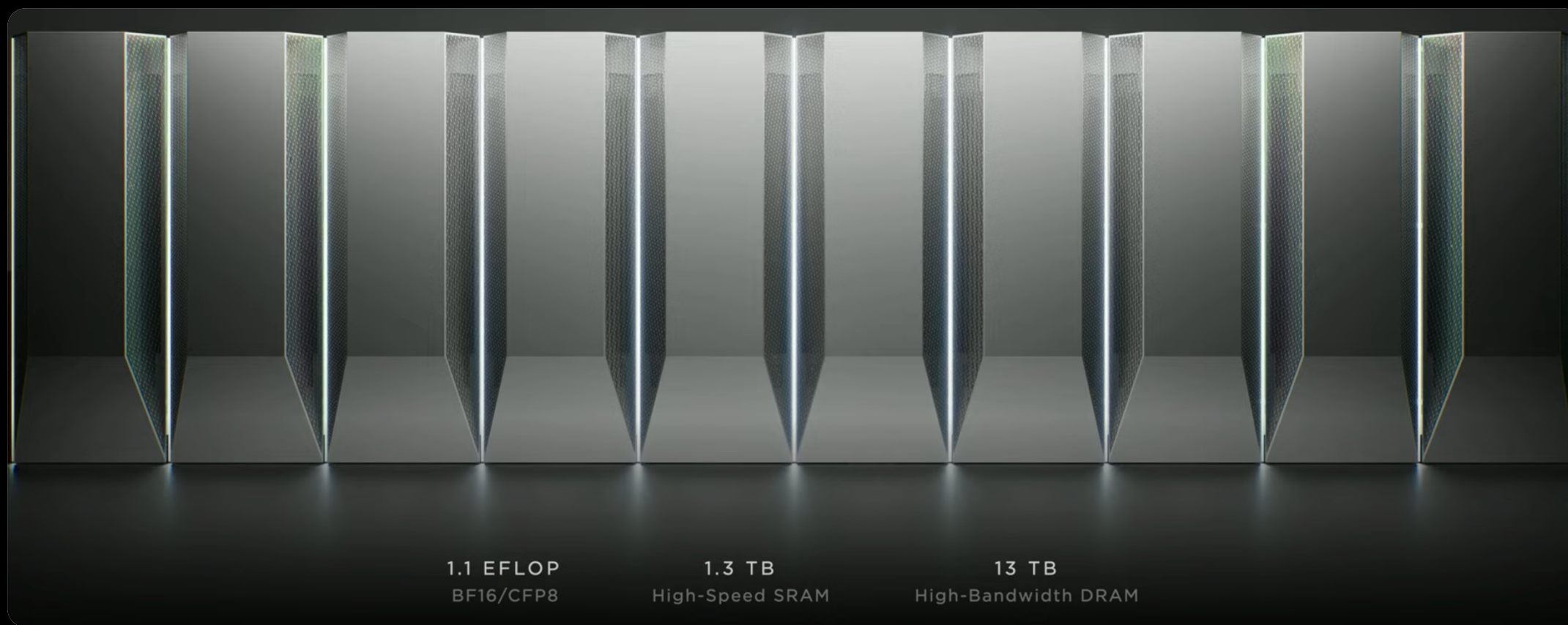
8 TB TOTAL MEMORY

640 GB/S PCIe BANDWIDTH



DOJO ExaPOD 集群将突破E级算力

- 每个DOJO ExaPOD 集成了120个训练模块，内置3000个DI芯片，拥有超过100万个训练节点，BF16/CFP8 峰值算力达到 1.1EFLOPS（百亿亿次浮点运算），拥有 1.3TB 高速 SRAM 和 13TB 高带宽 DRAM。



Training Tile

Heat Out



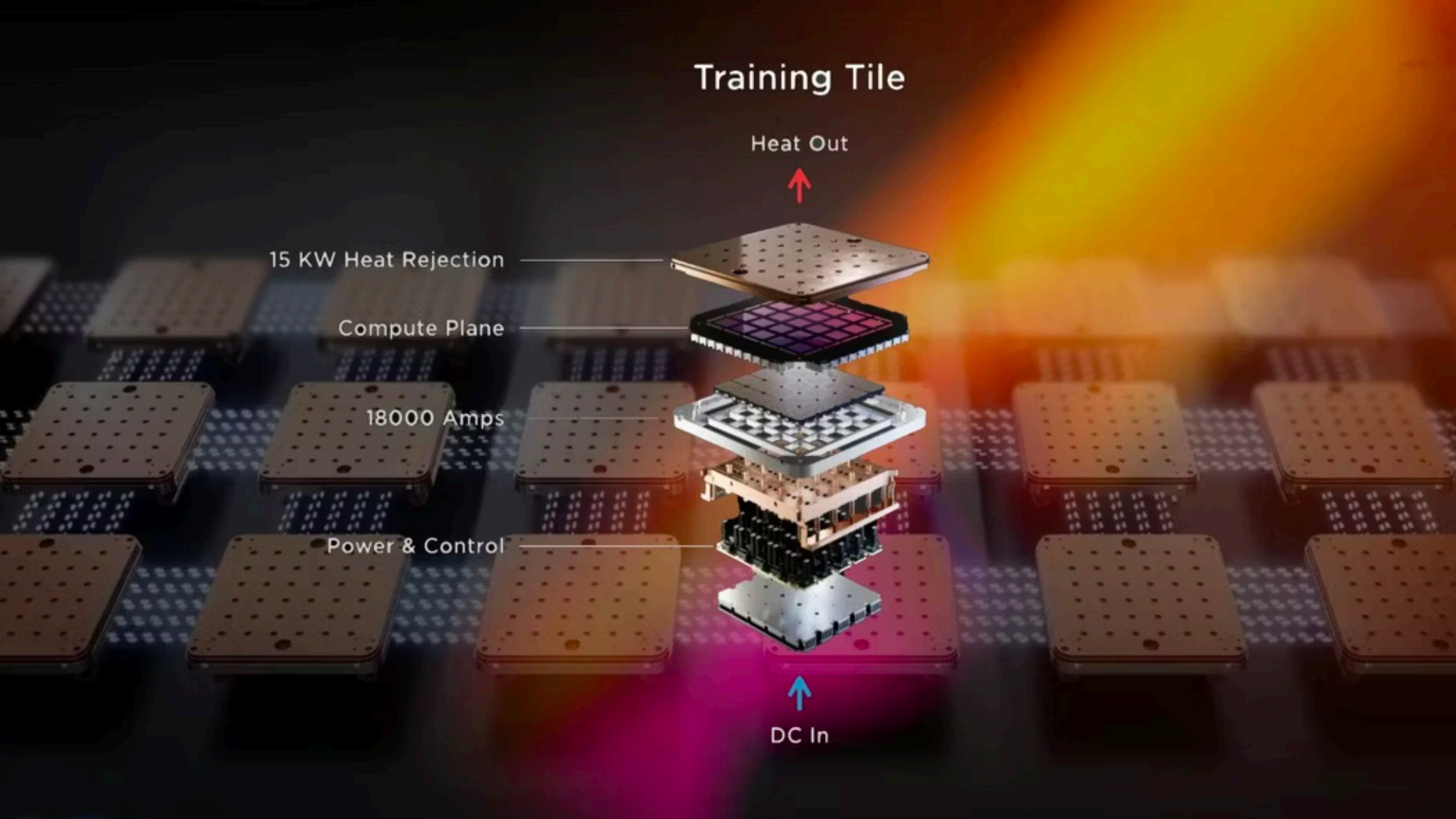
15 KW Heat Rejection

Compute Plane

18000 Amps

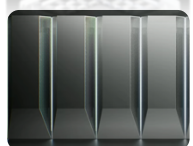
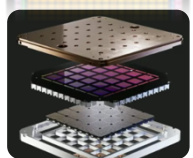
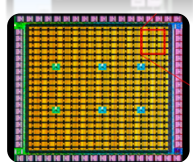
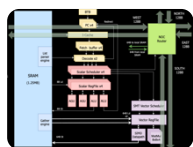
Power & Control

DC In



DOJO 超级计算机系统：D1 芯片、训练 Tile 和 ExaPOD 集群

- 每354个Dojo核心组成一块D1芯片，而每25颗芯片组成一个训练模组。最后120个训练模组组成一组ExaPOD计算集群，共计3000颗D1芯片。



分层	名称	片上SRAM	算力	备注
内核	DOJO Core	1.25MB	1.024 TFLOPS	单个计算核心，64bit，4个8x8x4矩阵计算核心，2GHz主频
芯片	DOJO D1	440MB	362 TFLOPS	单芯片，354核心数，654mm ²
格点	DOJO Tile	11GB	9050 TFLOPS	单个训练模组，每5x5个芯片组成一个训练模组
集群	ExaPOD	1320GB	1.1 EFLOPS	训练集群，每12个训练模组组成一个机柜，每10个机柜组成一个ExaPOD，一共3000个D1芯片

DOJO 架构设计哲学

- DOJO 采用存算一体架构（“存内计算”或者“近存计算”），单个可扩展计算平面、全局寻址快速存储器和统一的高带宽+低延迟。
- **面积精简**：将大量计算内核集成到芯片中，最大限度提高AI计算的吞吐量，因此需要在保障算力的情况下使单个内核的面积尽可能小，更好的折衷超算系统中算力堆叠和延迟的矛盾。
- **延迟精简**：为了实现其区域计算效率最大化，内核以 2 GHz 运行，只使用基本的分支预测器和小指令缓存，只保留必要部件架构，其余面积留给向量计算和矩阵计算单元。
- **功能精简**：通过削减对运行内部不是必须处理器功能，来进一步减少功耗和面积使用。DOJO 核心不进行数据端缓存，不支持虚拟内存，也不支持精确异常。

DOJO 架构设计哲学

- DOJO 采用存算一体架构（“存内计算”或者“近存计算”），单个可扩展计算平面、全局寻址快速存储器和统一的高带宽+低延迟。

名称	片上SRAM	算力	备注
DOJO DI	440MB	362 TFLOPS	单芯片，354核心数，654mm ²
A100	40MB	312 TFLOPS	单芯片，128核心数

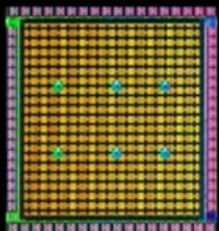
DOJO 超级计算机系统

D1 Chip

362 TFLOPs BF16/CFP8
22.6 TFLOPs FP32

10TBps/dir. On-Chip Bandwidth
4TBps/edge. Off-Chip Bandwidth

400W TDP



645mm²
7nm Technology

50 Billion Transistors

11+ Miles Of Wires

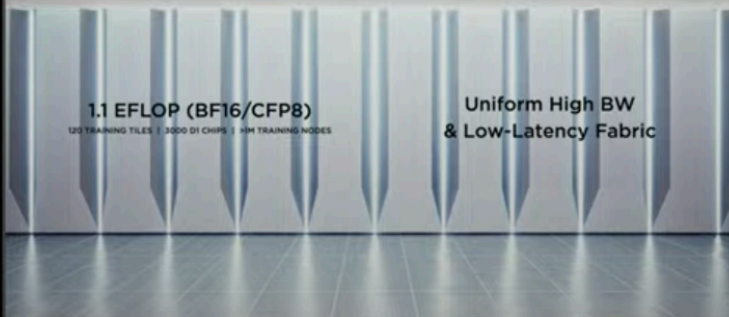
Training Tile



9 PFLOPs
36TB/s I/O BW
< 1 cu Ft

High-Performance
Extremely High-Bandwidth
Low Latencies
Lower Energy Communication

ExaPOD



1.1 EFLOP (BF16/CFP8)
120 TRAINING TILES | 3000 D1 CHIPS | 4M TRAINING NODES

Uniform High BW
& Low-Latency Fabric

Reference 引用&参考

1. <https://www.youtube.com/watch?v=uE2f7kiRhmw>
2. <https://www.youtube.com/watch?v=QurtwJdb5Ew>
3. <https://www.youtube.com/watch?v=DSw3lwsgNnc>
4. <https://chipsandcheese.com/2022/09/01/hot-chips-34-teslas-dojo-microarchitecture/>
5. https://en.wikipedia.org/wiki/Tesla_Dojo
6. <https://www.qbitai.com/2022/08/37209.html>
7. <https://zhidx.com/p/347884.html>
8. <https://www.cnblogs.com/wujianming-110117/p/17115152.html>
9. <https://mp.weixin.qq.com/s/xGytSXTW7-CL-OIV7y7QRw>
10. <https://mp.weixin.qq.com/s/uBL4x4MjzIGiCf2a0gnnUQ>
11. https://mp.weixin.qq.com/s/N_sMmndpyiq0_qbdwq3MuA
12. <https://www.51cto.com/article/717372.html>
13. <https://mp.weixin.qq.com/s/vmNsRVwmi3Azo-bPDanc2g>

BUILDING A BETTER CONNECTED WORLD

THANK YOU



Copyright©2014 Huawei Technologies Co., Ltd. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.