

大模型系列 - 集合通信

片内互联



ZOMI

大模型业务全流程

数据 & 模型算法

模型训练 & 微调

模型验证 & 推理部署

2. 数据处理

开源
数据
数据
预处理
向量
数据库

3. 模型算法

LLM模型
架构
多模态
一切皆Tokens

4. 模型训练

混合
精度
梯度
检查
梯度
累积
...

5. 分布式并行

训练集群稳定性

6. 模型微调

全参微调
低参微调
指令微调

7. 模型验证

下游
任务
测评
标准
bench
mark

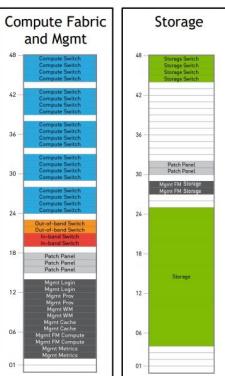
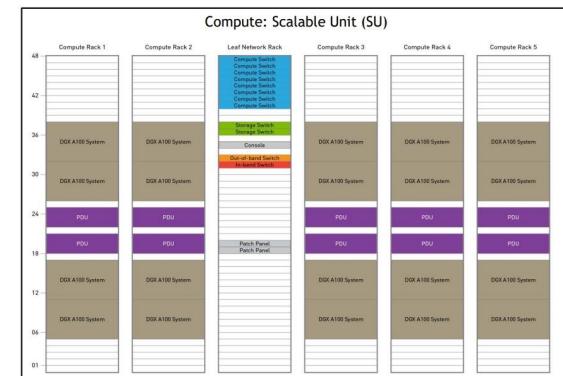
8. 推理与智能体

量化
压缩
推理
加速
9. Agent
智能体

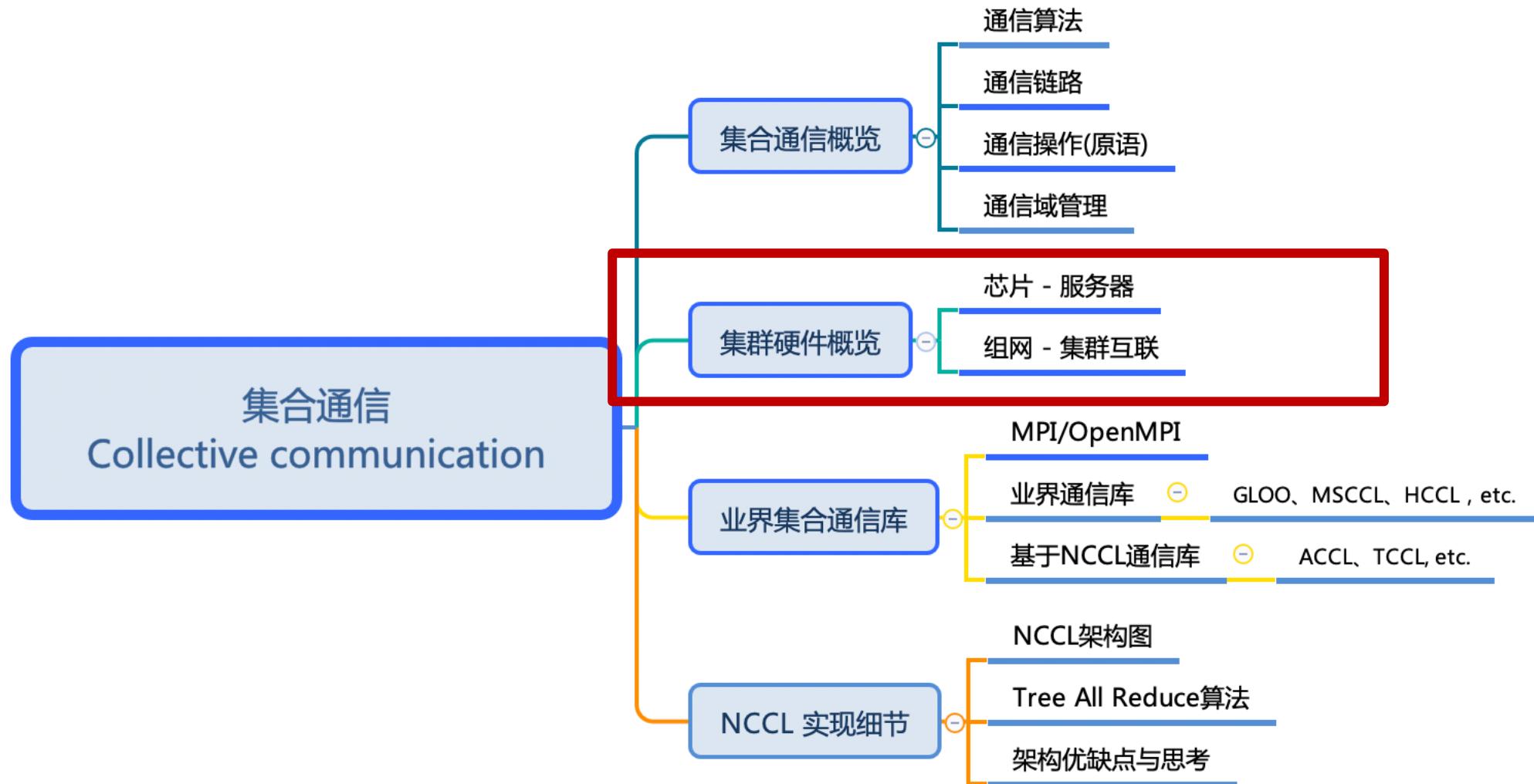
1. AI 集群建设

计算、存储、
网络
AI 集群机房建
设
AI 集群上线与
运维

集群算力准备



思维导图 XMind

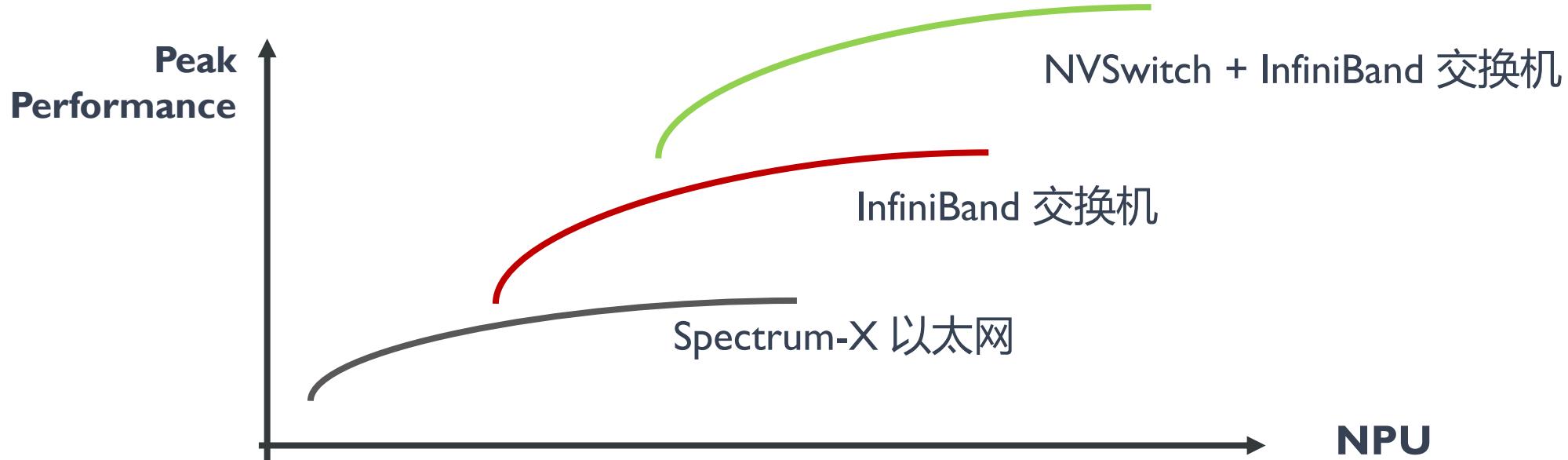


Background

01. 集群背景介绍

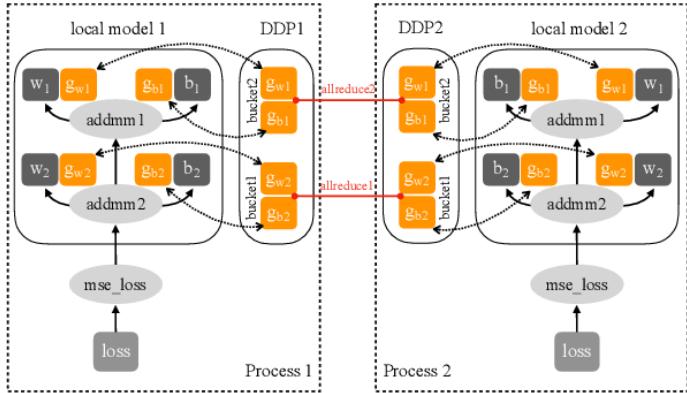
不同 AI 时期对通信的诉求

- **AI 集群:**
 - 单一化业务，整个 AI 系统只为大模型（LLM、LMM等）或者搜广推服务，几乎没有其他业务的复用性；
 - 用于超大规模的模型（百/千亿参数量）的训练、推理，Lo 基础大模型算法研究的探索；
 - 训练大模型走极致性能优化路线 vs 虚拟化云服务和 AI 通用算力服务化走性价比路线；



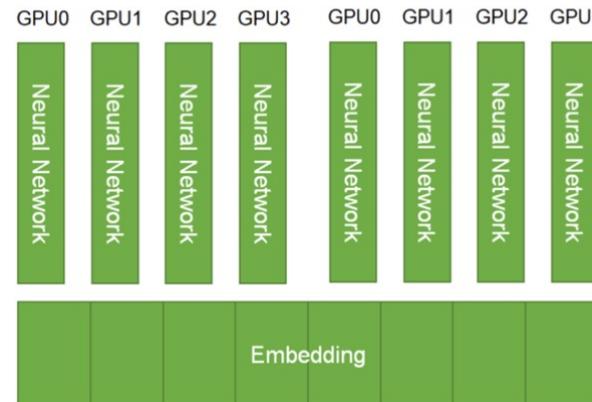
不同 AI 时期对通信的诉求

CV 模型--进入成熟期



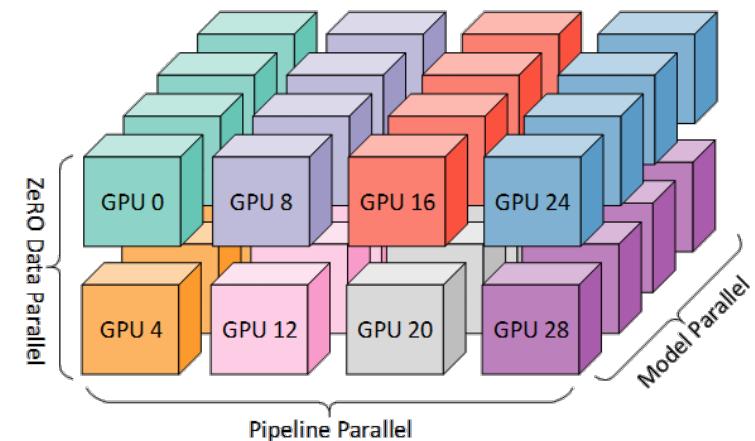
数据并行，模型参数少，All reduce 为主，参数面规模<4K

推荐模型--大规模应用



Embedding 并行 All2All，参数面规模<4K

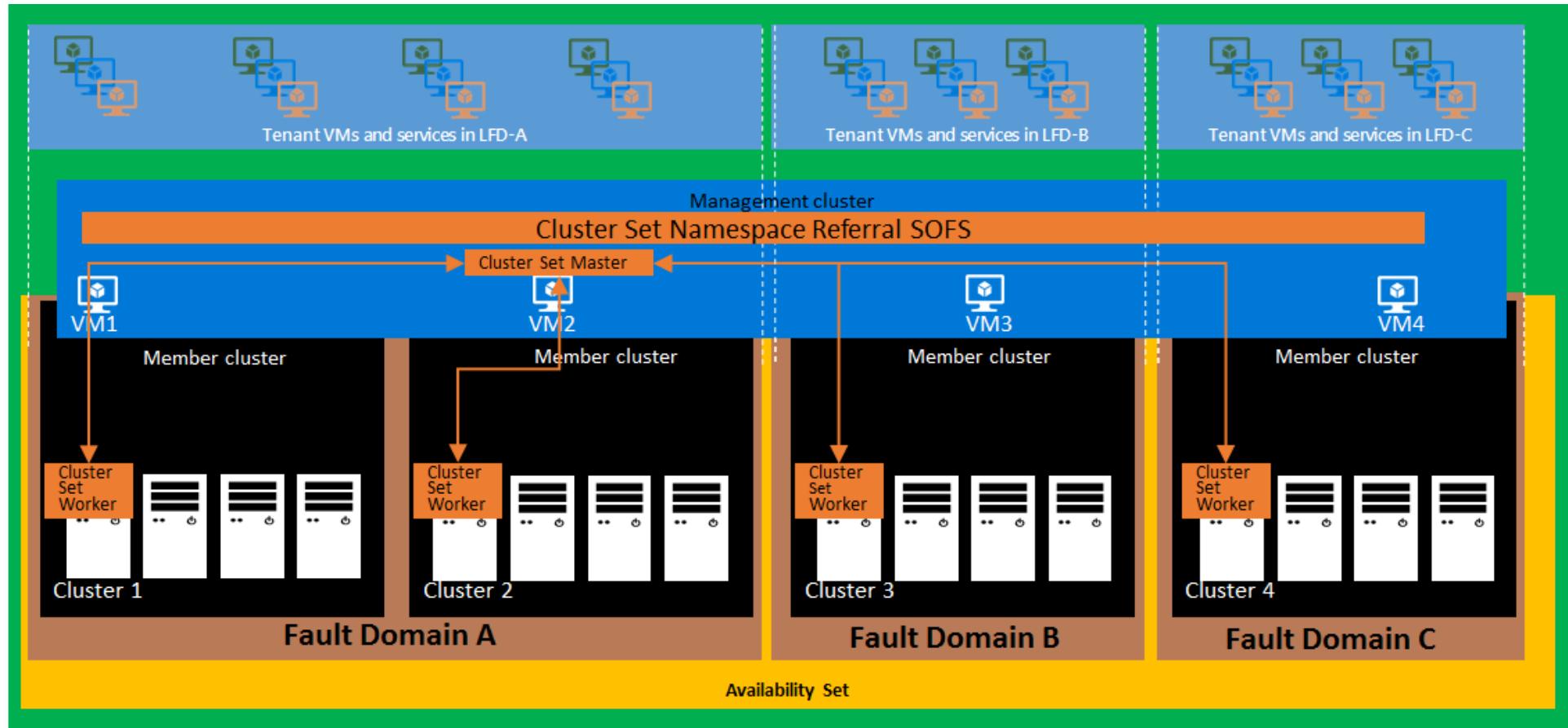
大模型—快速发展期



PTDES混合并行，集合通信都有使用，参数规模<32K

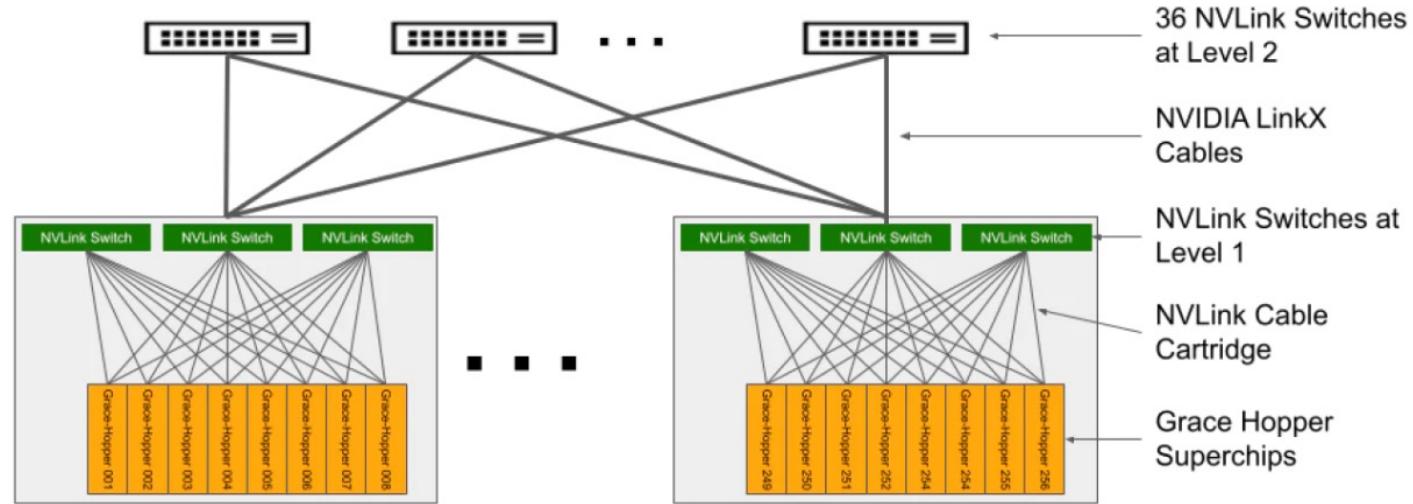
传统组网互联

- 传统方案中 (<2018) , GPU 互联采用 PCIe, 服务器节点间互联采用以太网 Ethernet。



大模型组网互联

- 大模型数据、参数量极大。服务器不同计算节点间，对超高带宽、超低延迟和超高可靠性的互联技术要求高。



AI 计算集群互联方式

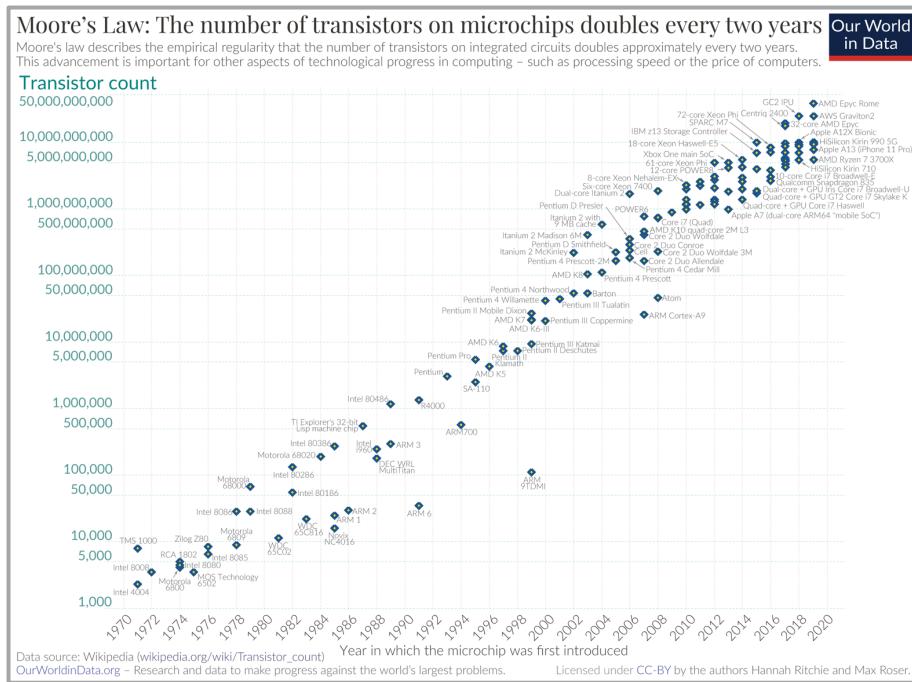
- **AI 计算集群互联演变，主要有三大趋势：**
 - Die 间，多芯粒互联技术和合封技术正加速崛起。
 - 片间，由 PCIe 向多节点无损网络演进；
 - 集群间，互联方式从 TCP/IP 向 RDMA 架构转变；
- UCIe、CCITA 为代表联盟组织，在积极推进 Chiplet 标准化协议与生态建立与完善

02. DIE 间互联

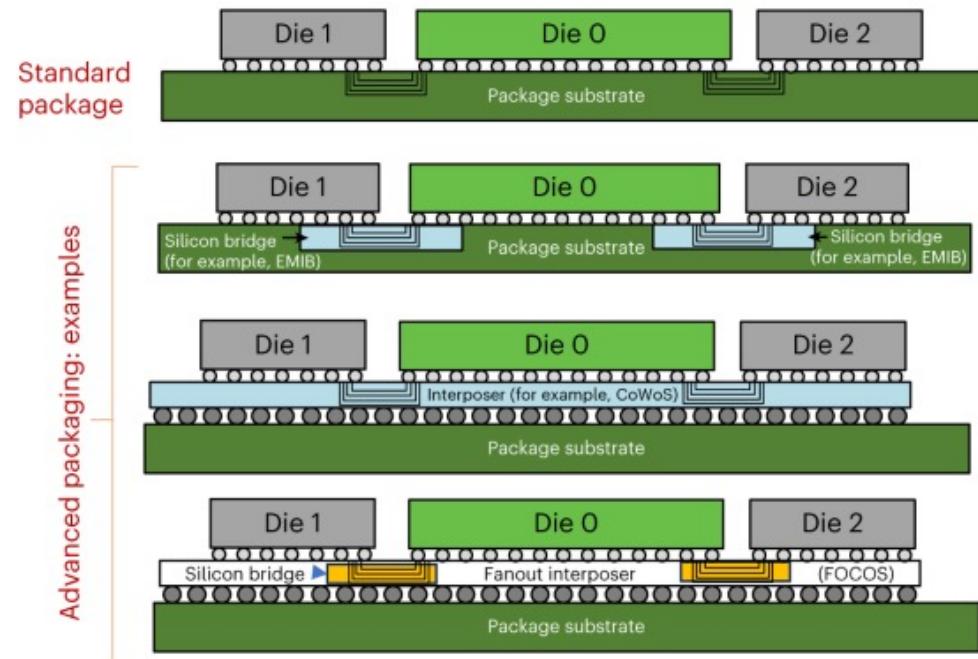
SoC 架构转向 Chiplet 异构

- 大模型对算力需求持续增长，在工艺发展较慢情况下，为继续提升算力，AI 芯片从传统 SoC 架构转向 Chiplet 异构。除芯粒数量不断增加，为有效发挥片内算力，也引发芯粒间互联挑战。

摩尔定律发展



Chiplet 封装



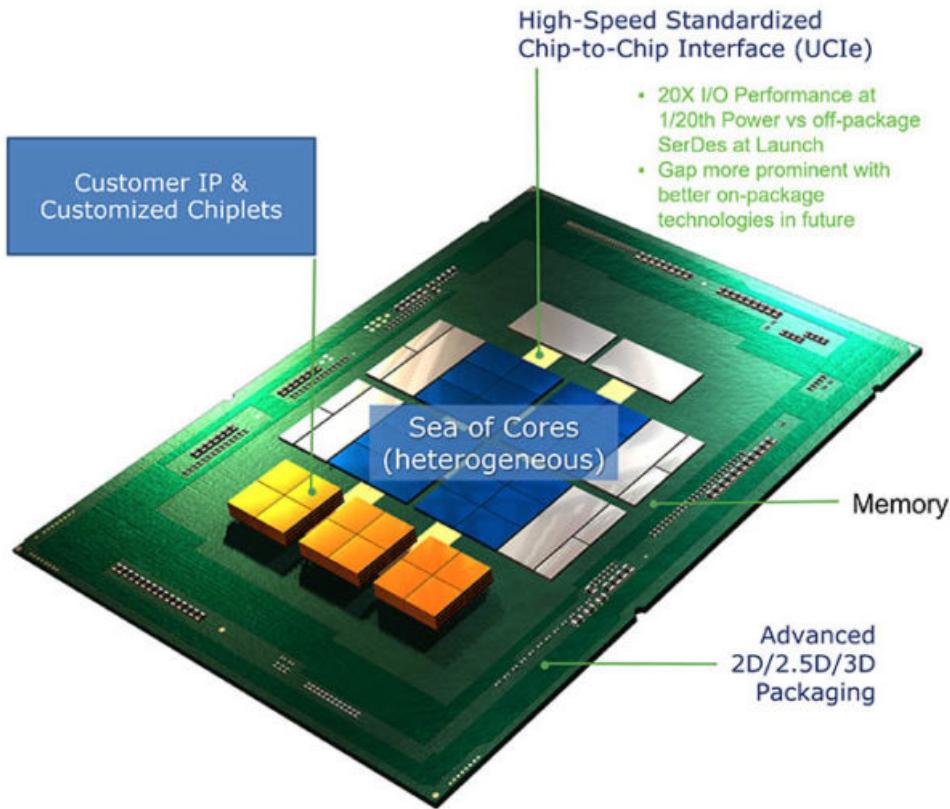
SoC 架构转向 Chiplet 异构

- 基于 Chiplet 架构，创新 Die 间互联技术正加速崛起。UClie、CCITA 等组织积极的为 Chiplet 互连建立统一的接口标准。

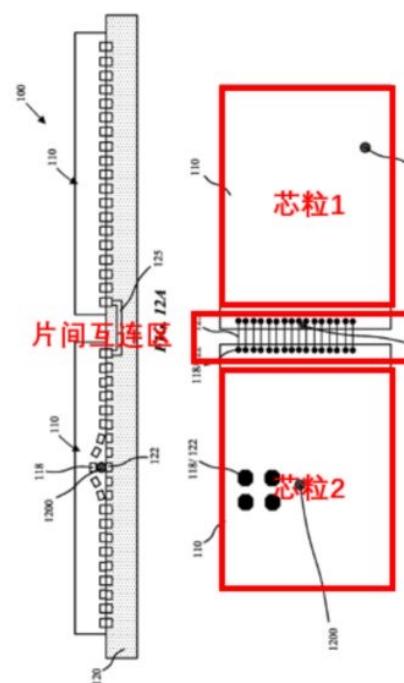


SoC 架构转向 Chiplet 异构

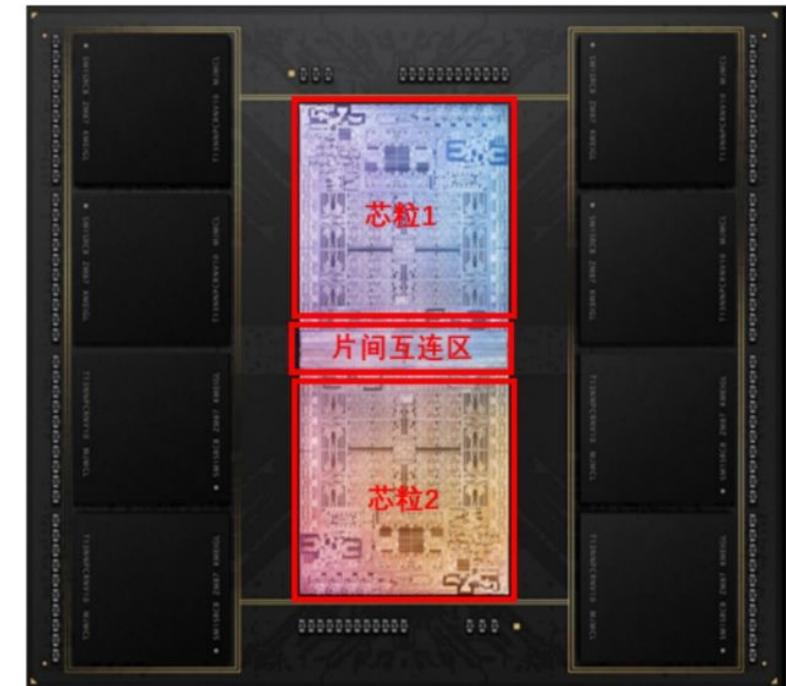
UCle Chiplet 标准



苹果 M1 Chiplet 示例



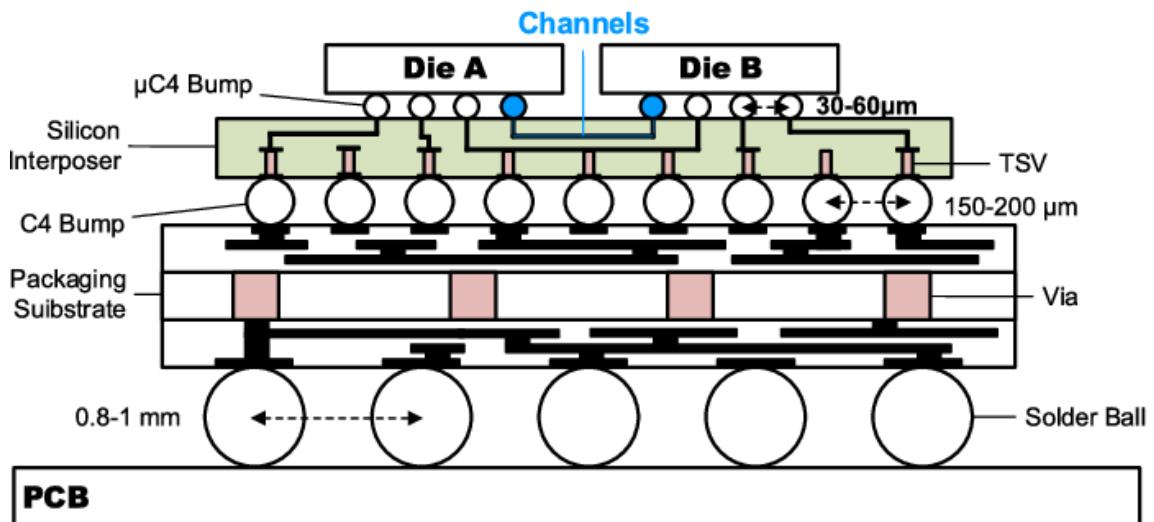
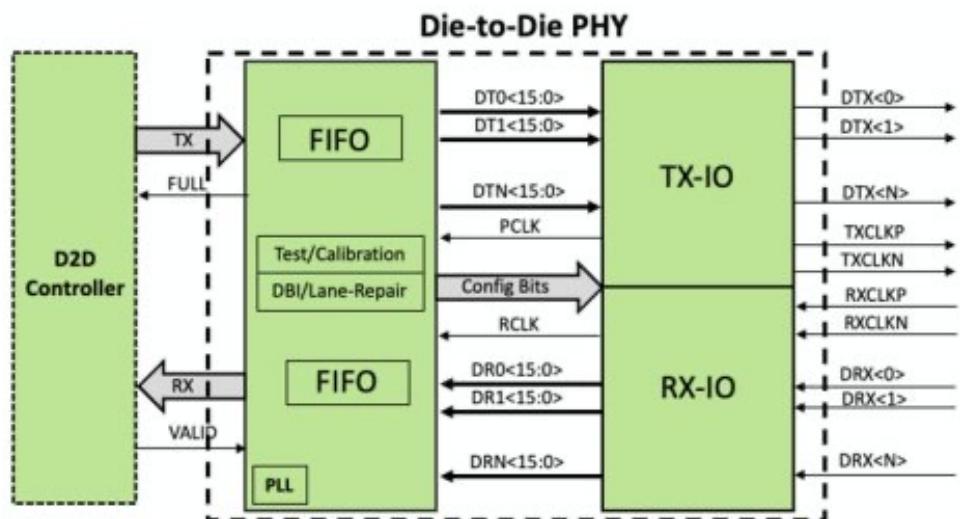
苹果Chiplet专利



苹果M1 Ultra

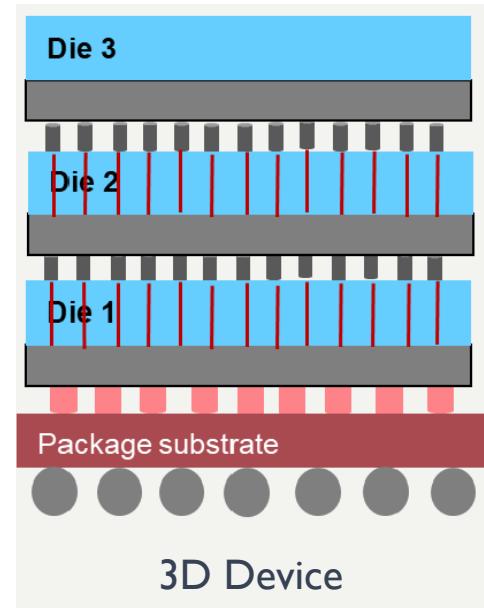
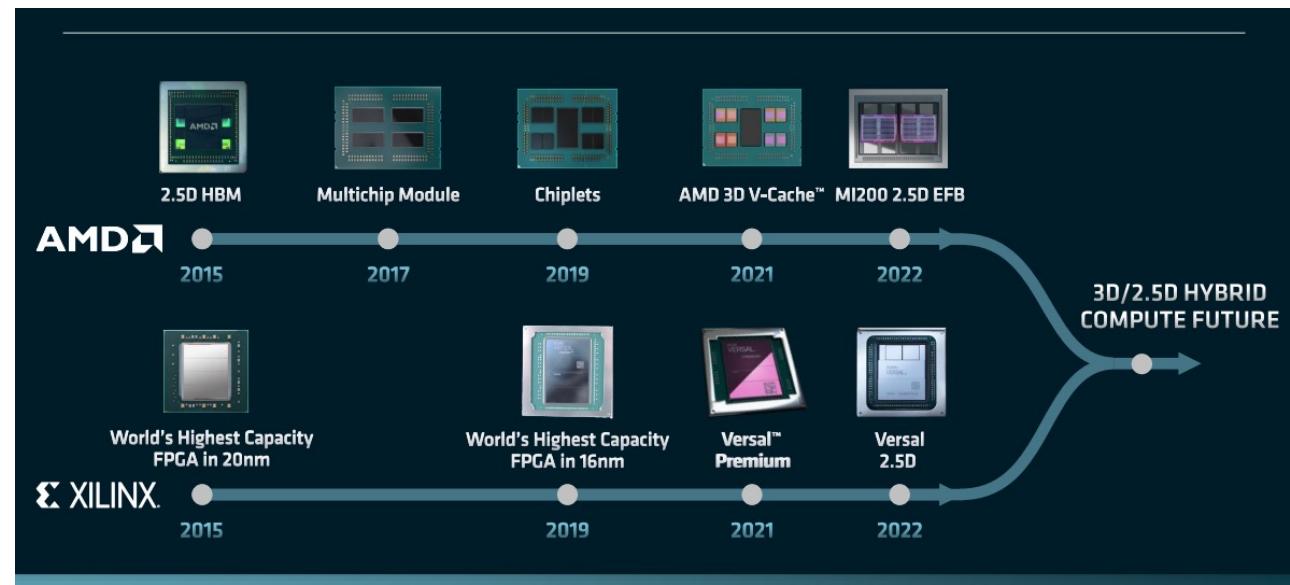
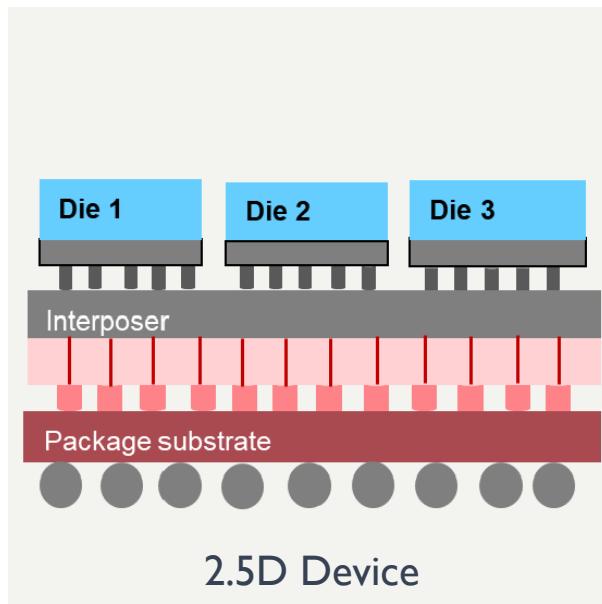
D2D Interface

- 随芯片系统变得越来越复杂，不同功能单元（芯粒）产生大量数据流需要专用的互联接口来实现数据的传输和调度。
- 这种专用互联接口简称为 Die2Die 接口，负责在不同芯粒间传输数据，协调调度数据流，确保芯片系统高效运行。



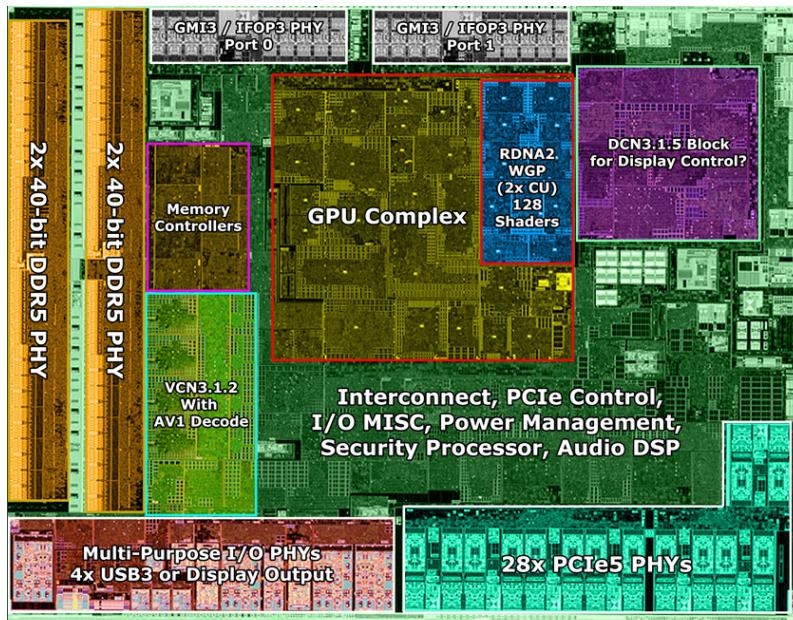
IO Die & Base Die

- IO Die：专用互联芯粒。IO Die 作为数据传输和调度核心，整合存储单元、Die2Die 接口和多种高速接口，通过自定义算法实现数据流和信息流的分发调度：
 - IO Die 通常适用于 2.5D Chiplet 芯片架构。
 - e.g. AMD 300I AI 芯片架构包含互联芯颗。

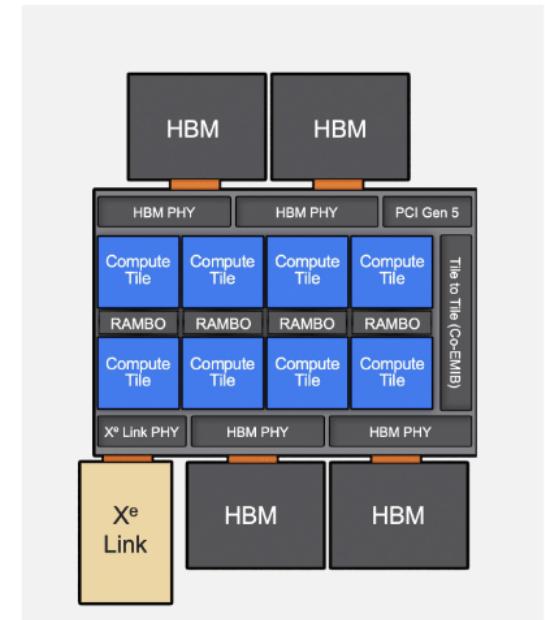
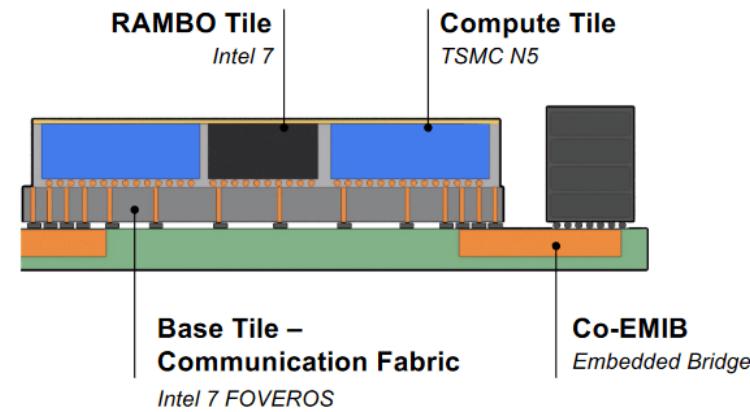


IO Die & Base Die

- Base Die: 当芯片性能继续增高, 平面维度也很难满足 Die 间互联需求。于是, 互联方式逐渐从 2D to 3D 垂直迭代:
 1. 芯片行业开始基于芯粒 3D 堆叠方式, 进一步提升芯片算力密度;
 2. 集成 die2die 3D 接口, Cache 等模块, 实现更快垂直互联, 减少片内存储延迟和功耗。



Ponte Vecchio Foveros
+ EMIB Construction





Thank you

把AI系统带入每个开发者、每个家庭、
每个组织，构建万物互联的智能世界

Bring AI System to every person, home and
organization for a fully connected,
intelligent world.

Copyright © 2023 XXX 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. XXX may change the information at any time without notice.



Course chenzomi12.github.io

GitHub github.com/chenzomi12/DeepLearningSystem