



ZOMI

高性能计算HPC 定义

Content



Content

AI 系统 + 大模型全栈架构图



时事热点



大模型训推



编译计算架构



硬件体系结构



新基建，中国数字经济引擎的顶层规划

“数字经济不能建设在沙滩上”

习总书记2018年关于自主创新的重要讲话

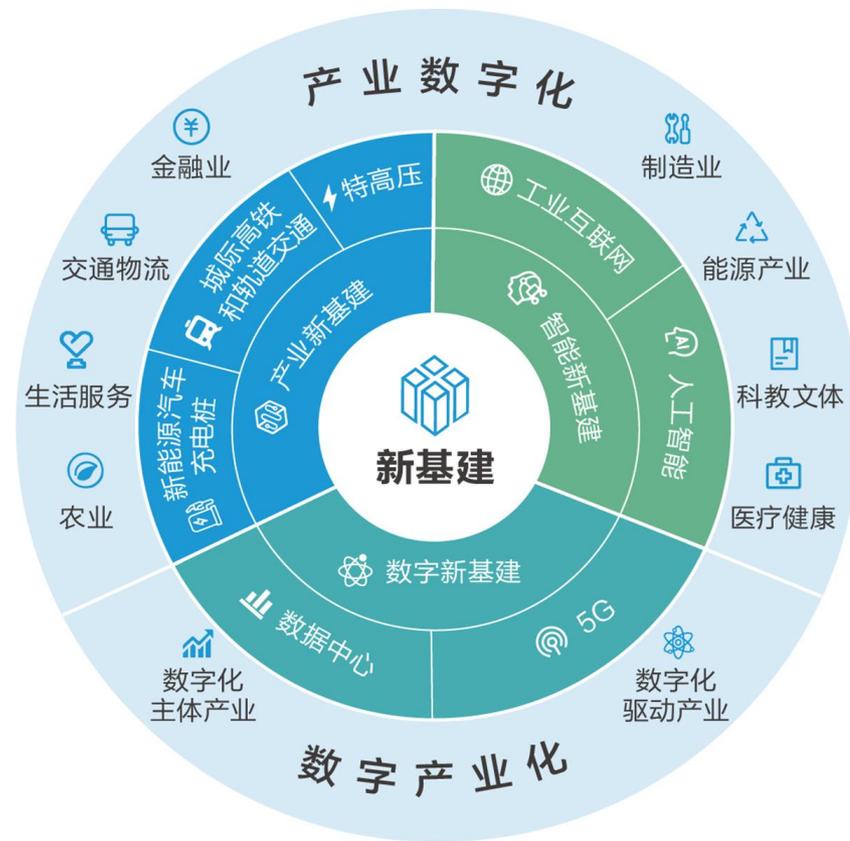
掌握创新和发展主动权，促进中国产业迈向全球价值链高端。

掌握核心技术，加快推进国产自主可控替代

“要适度超前，布局有利于引领产业发展和维护国家安全的基础设施”

习总书记2022年4月26日主持召开中央财经委员会第十一次会议

要加强信息、科技、物流等产业升级基础设施建设，布局建设新一代超算、云计算、人工智能平台、宽带基础网络等设施，推进重大科技基础设施布局建设



新基建和数字经济

数字新基建（**数据中心**）和智能新基建（**人工智能**）



Content

1. 什么是算力（算力单位、算力空间）
2. 高性能计算（数值分析、并行计算）
3. 集群计算（高性能计算中心、云数据中心、AI 智算中心）



01

算力

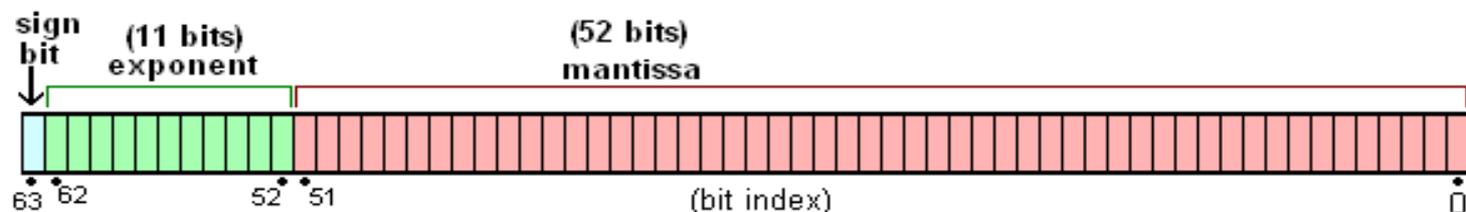
Computational Power



小知识：算力水平是如何衡量的？

- **FLOPS** = **F**loating point **O**perations **P**er **S**econd，每秒浮点计算次数
- HPC 一般数据类型 64位双字长精度（FP64），FLOPS一般采用FP64的算力水平来评估

Processors in Mate60 Pro	FP64 (双精度)	FP32 (单精度)
HiSilicon Kirin 9000	1,166GFLOPS	2,332 GFLOPS



- **IEEE-754 标准** 双精度浮点数 $D = (-1)^s \times 1.f \times 2^{e-1023}$ ，
- 其中，**s** 是**符号位**（0为正1为负），**f** 是（52位二进制）**尾数**，**e** 是（11位二进制）**指数**

小知识：算力水平是如何衡量的？

Name	Unit (单位)	Value	中文描述	算力水平约等于多少部Mate40 Pro
Kilo FLOPS	KFLOPS	10^3	每秒一千次	0.000000001部
Mega FLOPS	MFLOPS	10^6	每秒一百万次	0.000001部
Giga FLOPS	GFLOPS	10^9	每秒十亿次	0.001部
Tera FLOPS	TFLOPS	10^{12}	每秒一万亿次	一部Mate40Pro, 1/12部XBOX
Peta FLOPS	PFLOPS	10^{15}	每秒一千万亿次	一干部 (武汉超算, 算力200P)
Exa FLOPS	EFLOPS	10^{18}	每秒一百京次	一百万部 (深圳超算, 算力2.5E)
Zetta FLOPS	ZFLOPS	10^{21}	每秒十万京次	十亿部
Yotta FLOPS	YFLOPS	10^{24}	每秒十亿京次	一万亿部



数字经济时代，算力是国家间的核心竞争力

6 产经

2025年4月22日 星期二

经济日报

产业聚焦·AI新场景①

本报记者 李 晔

人工智能优化算力布局

中国算力目前主要集中在东部沿海地区，算力资源分布不均。算力是数字经济发展的核心要素，也是国家竞争力的重要体现。随着人工智能技术的快速发展，算力需求呈现出爆发式增长。如何优化算力布局，提高算力利用效率，成为当前亟待解决的问题。



算力规模不缩水

算力是数字经济发展的核心要素，也是国家竞争力的重要体现。随着人工智能技术的快速发展，算力需求呈现出爆发式增长。如何优化算力布局，提高算力利用效率，成为当前亟待解决的问题。

降低算力利用成本

算力是数字经济发展的核心要素，也是国家竞争力的重要体现。随着人工智能技术的快速发展，算力需求呈现出爆发式增长。如何优化算力布局，提高算力利用效率，成为当前亟待解决的问题。

车网互动开启规模化应用

车网互动开启规模化应用，是智能交通系统的重要组成部分。通过车网互动，可以实现车辆与电网、电网与电网之间的信息交互和能量流动，提高能源利用效率，降低碳排放。

重庆安全技术职业学院“三位一体”推动高职院校安全教育评价改革

重庆安全技术职业学院“三位一体”推动高职院校安全教育评价改革，旨在提高高职院校安全教育的质量和水平。通过构建“三位一体”的安全教育评价机制，实现安全教育评价的科学化、规范化和制度化。

数智赋能评价改革 提升技能人才培养质量

数智赋能评价改革 提升技能人才培养质量，是职业教育改革的重要方向。通过引入大数据、人工智能等先进技术，可以实现对技能人才培养过程的精准评价和动态调整，提高人才培养的针对性和实效性。

水库安全管理一刻不能松懈

水库安全管理一刻不能松懈，是保障人民群众生命财产安全的重要任务。随着经济社会的快速发展和人口规模的不断扩大，水库的数量和规模也在不断增加，水库安全管理面临着前所未有的挑战。

重庆安全技术职业学院“三位一体”推动高职院校安全教育评价改革

重庆安全技术职业学院“三位一体”推动高职院校安全教育评价改革，旨在提高高职院校安全教育的质量和水平。通过构建“三位一体”的安全教育评价机制，实现安全教育评价的科学化、规范化和制度化。

数智赋能评价改革 提升技能人才培养质量

数智赋能评价改革 提升技能人才培养质量，是职业教育改革的重要方向。通过引入大数据、人工智能等先进技术，可以实现对技能人才培养过程的精准评价和动态调整，提高人才培养的针对性和实效性。

● 中国人均算力水平当前仅为美国的1/5，提升空间巨大；预计到2026年，中国算力市场将持平甚至超越美国。

● 工信部数据显示，截至2024年三季度末，中国算力总规模达268EFLOPS（每秒百亿亿次浮点运算，以FP32单精度计算），算力应用项目超过1.3万个，在用算力中心机架总规模超过880万架标准机架，算力总规模居世界前列。

重庆安全技术职业学院“三位一体”推动高职院校安全教育评价改革

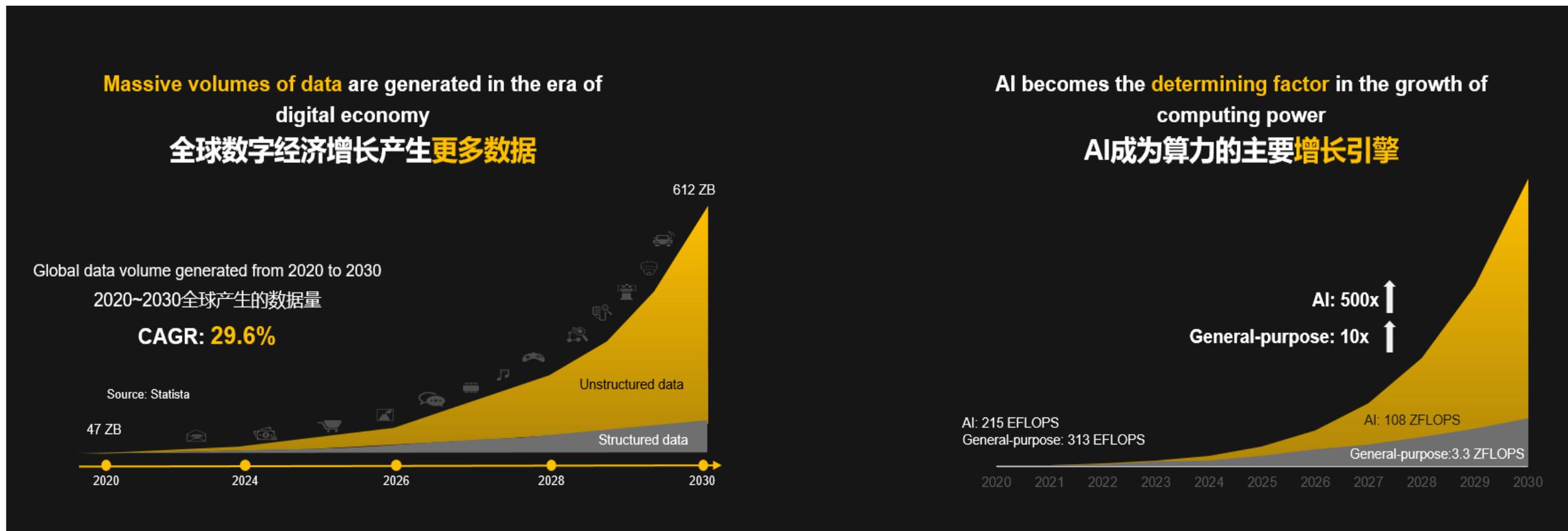
“三位一体”是指：安全教育评价、安全教育实施、安全教育保障。通过构建“三位一体”的安全教育评价机制，实现安全教育评价的科学化、规范化和制度化。

数智赋能评价改革 提升技能人才培养质量

数智赋能评价改革 提升技能人才培养质量，是职业教育改革的重要方向。通过引入大数据、人工智能等先进技术，可以实现对技能人才培养过程的精准评价和动态调整，提高人才培养的针对性和实效性。

算力规模

- 1961年，AI之父 John McCarthy 提出：“算力服务将成为未来的公共基础设施”，同电话网络一样重要
- 2030年，全球产生数据量年均复合增长29.6%，通用算力需求将增长10倍，AI算力需求将增长500倍



02

高性能计算

HPC



数值分析和并行计算

$$\frac{dy}{dx}$$



- 很多人知道韦神 但很多人不知道的是他的研究方向是N-S方程







并发计算 vs 并行计算

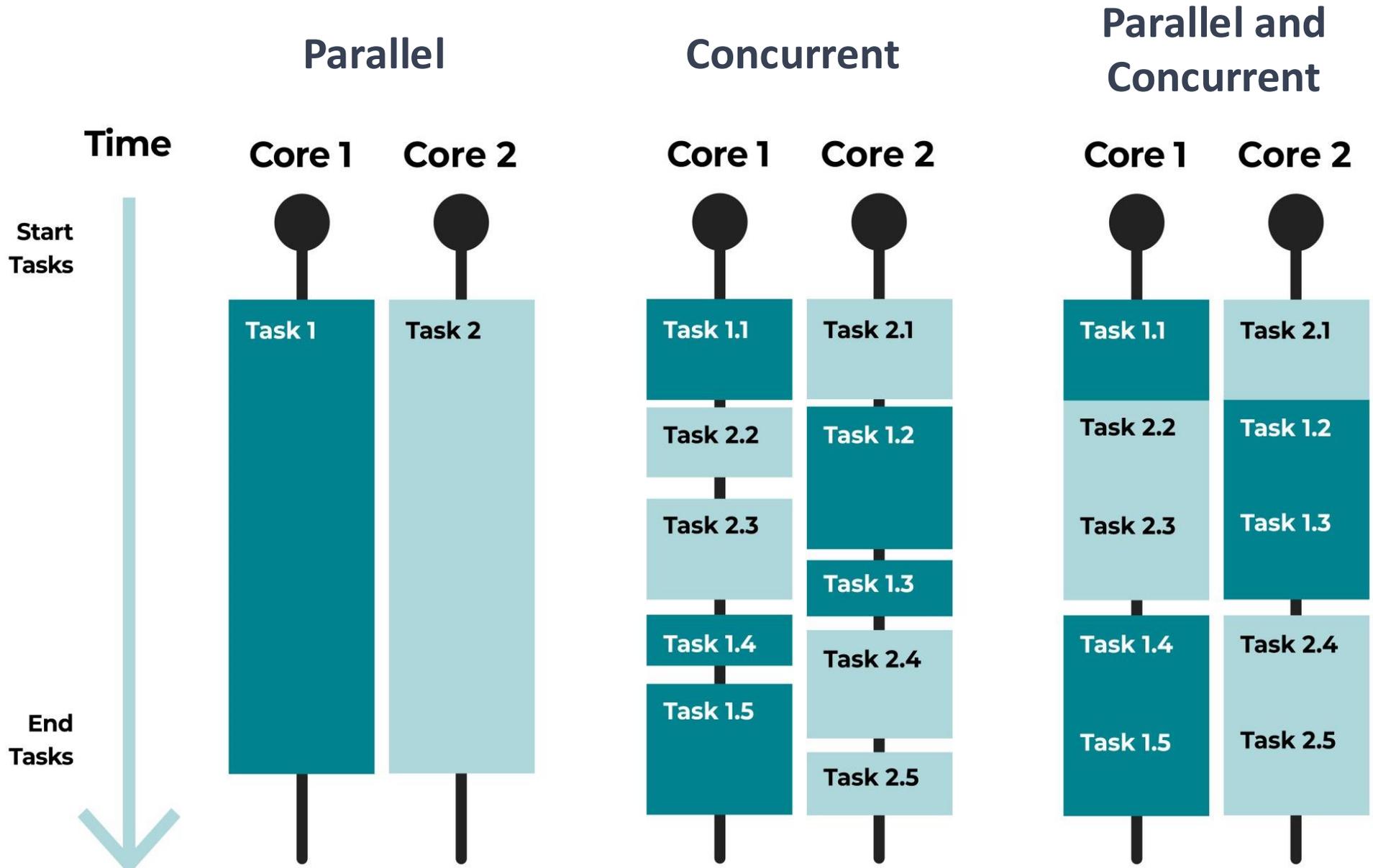
- **并发计算 Concurrent Computing:**

- 是一种程序计算形式，在系统中，至少有两个以上的计算任务在同时运作，计算结果可能同时发生。并发概念强调的是单个处理器在单位时间内完成多个任务，类似“一个人同时吃三个馒头”

- **并行计算 Parallel Computing:**

- 一般是指许多指令得以同时进行的计算模式。在同时进行的前提下，可以将计算的过程分解成小部分，之后以并发方式来加以解决。并行概念强调的是多个处理器同时完成多个任务。





03

集群计算

cluster computing



集群计算

- **计算机集群**简称**集群**，是一种计算机系统，它通过一组松散集成的计算机软件或硬件连接起来高度紧密地协作完成计算工作。



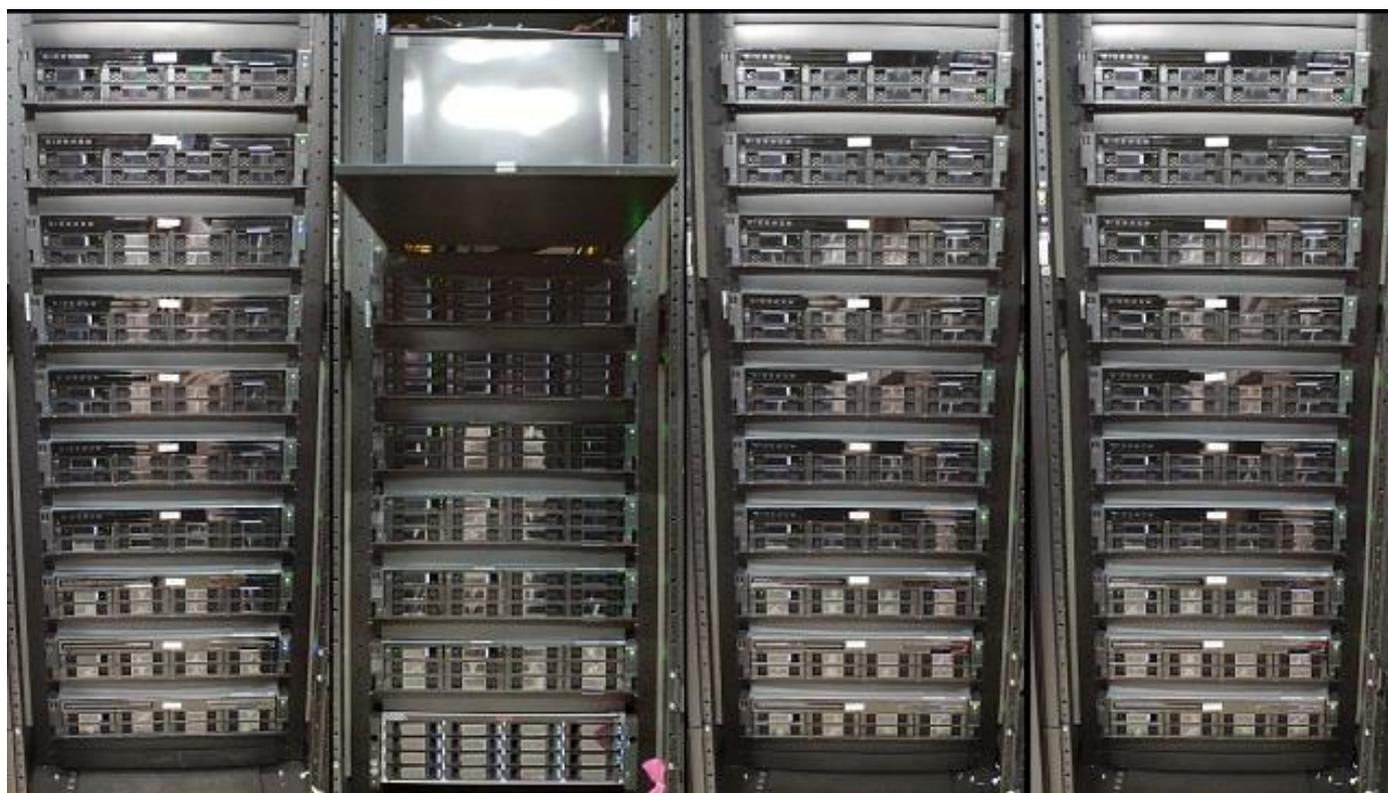
集群计算

- 集群是一个完整系统，而非多个计算机系统，可以被看作是一**整台计算机**（类似于“蚁群” / “蜂群”社会系统），是主要的算力基础设施建设方式。



算力基础设施建设

- 从分层堆叠走向集群全栈优化，达成规模系统最优。



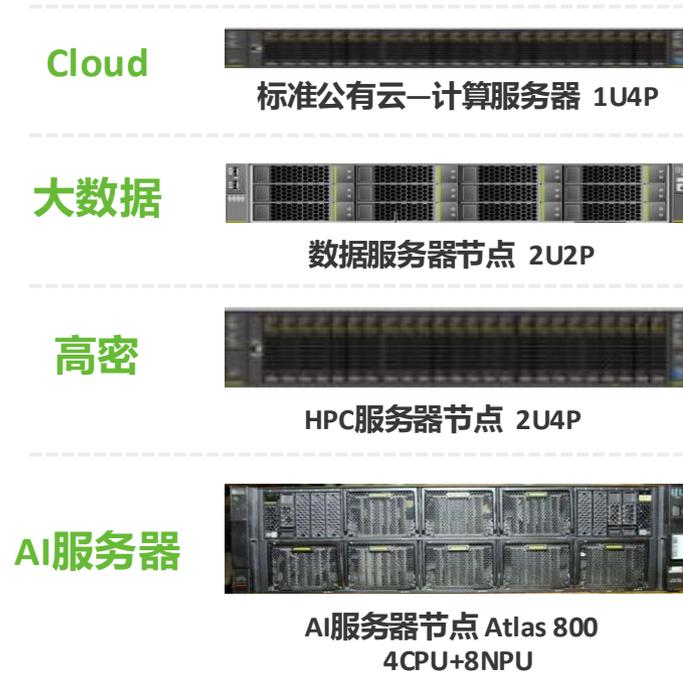
传统DC中服务器部署形态

1. 人工连线耗时长，故障率高，施工维护复杂；风冷为主，系统能效差
2. 计算、散热、电源网络松耦合集成，算力利用率低，算力密度低

服务器主板



服务器节点



标准服务器机柜内堆叠



集群计算规模应用的场景

高性能计算中心



人工智能计算中心



云数据中心



集群计算规模应用的场景： 计算中心

- **专用**：面向特定计算需求的高性能并行计算。面向国防、科研等重算力超算场景，以及人工智能大模型训练场景，算力和功耗惊人，体系架构专为特定应用算力摸高定制优化。

高性能计算中心



人工智能计算中心



集群计算规模应用的场景：数据中心

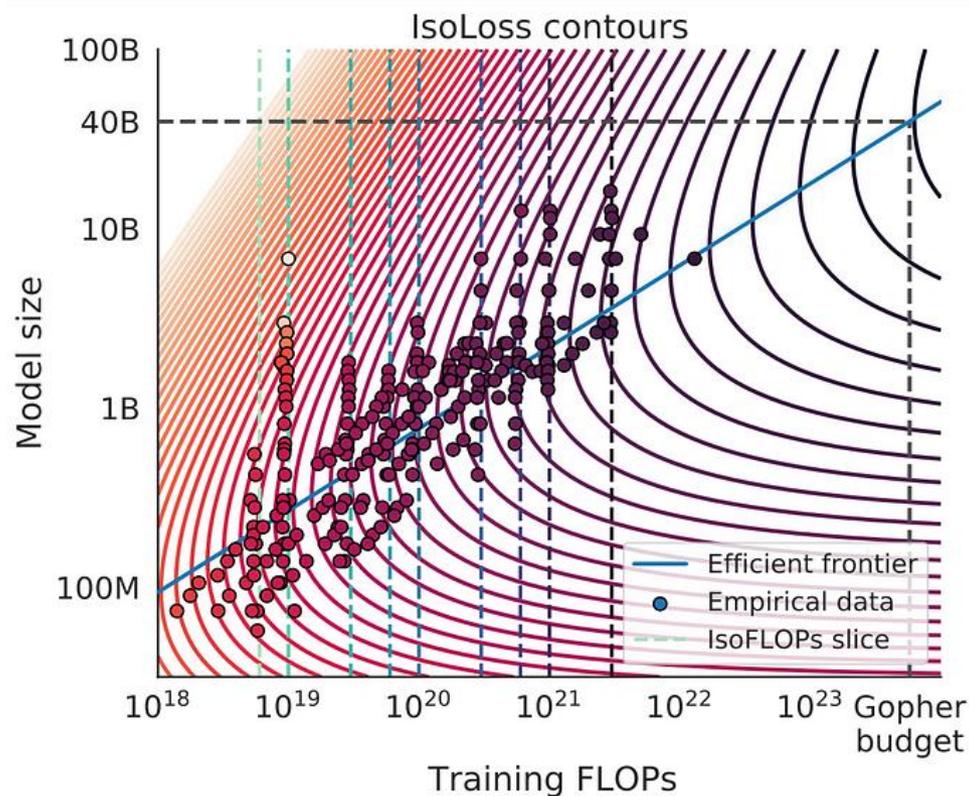
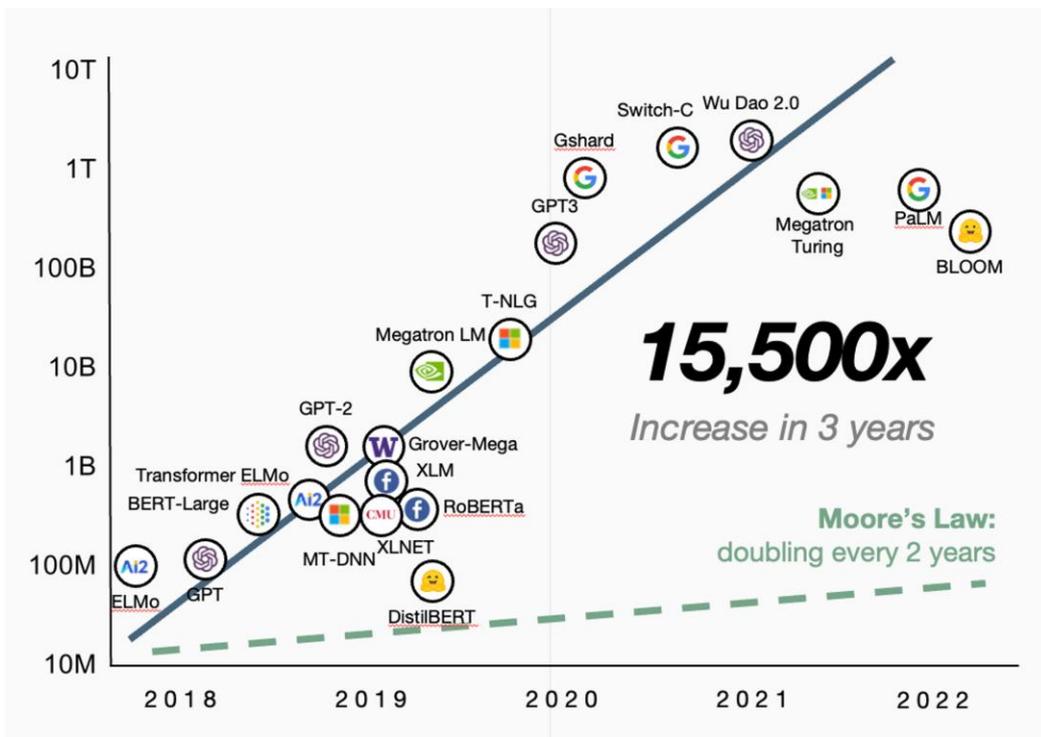
- **通用**：面向通用计算的**大规模并发处理**。面向大众商用，提供互联网IT服务，并行处理规模惊人，计算相对简单，强调高性价比和高兼容性，对可靠性、可用性、可服务性（RAS）要求高。

云数据中心



人工智能计算

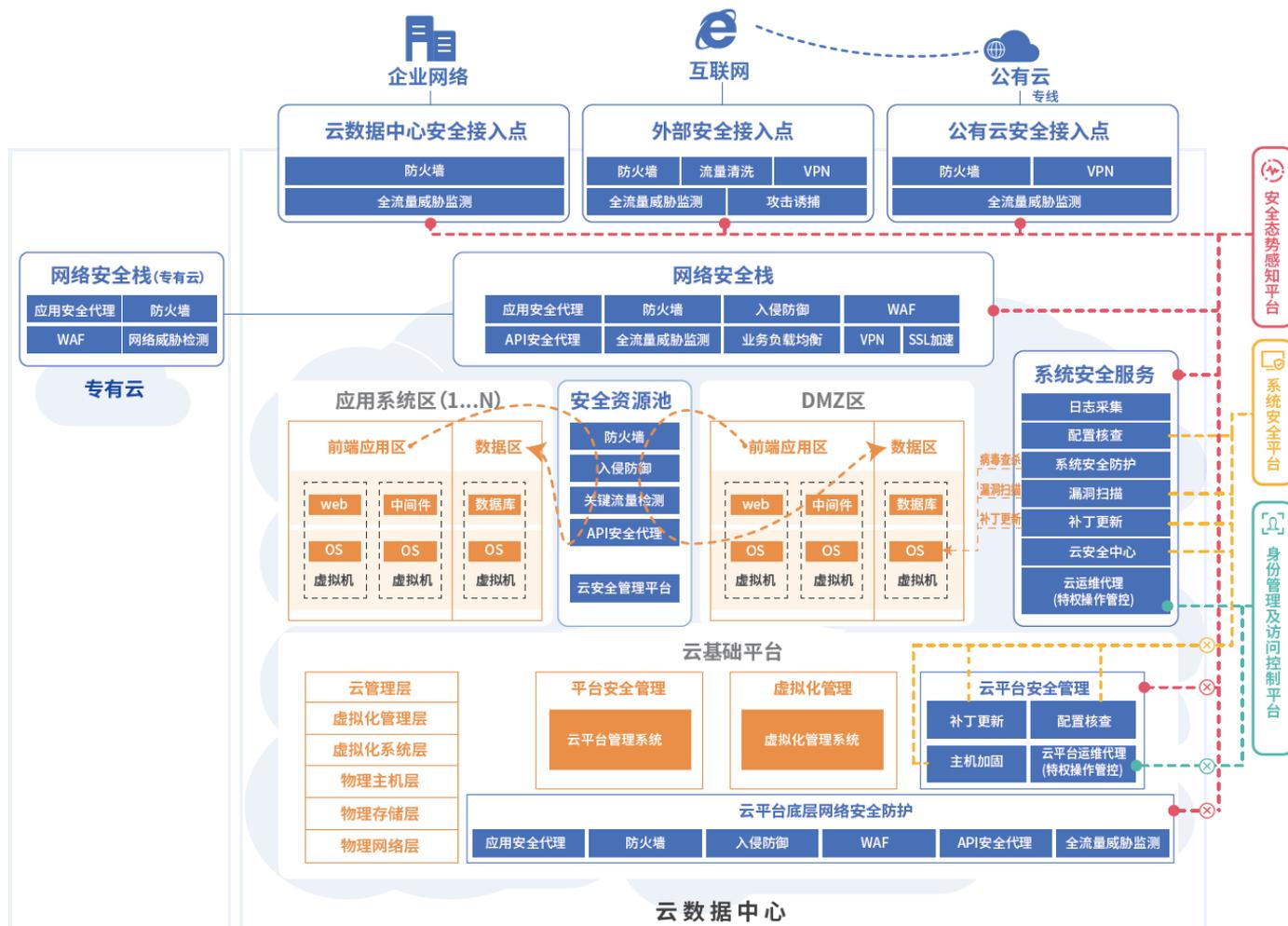
- AI成为新形势下中美科技竞争的新赛道，**大数据+大模型 驱动 训练&推理** 智能算力爆发增长
- 大模型算力需求每2年增长750倍 vs 硬件算力供给按摩尔定律每2年仅增长3倍
- 基于FFN前馈神经网络+Attention注意力机制的Transformer通用大模型架构



云数据中心

• 特点：高并发通用云计算

1. **资源按需分配**：资源全面虚拟化，不同粒度按需弹性分配
2. **服务无处不在**：服务通过网络公开，随时随地接入/控制
3. **绿色节能**：降低能耗、提高算效、降低OPEX
4. **集约化建设**：降低TTM和CAPEX



三类计算集群的主要区别

集群类型	应用目的	数据类型	计算特征	网络特征	存储特征
HPC高性能计算中心	国防科研	FP64双精度	重算力，密集计算，高并行	密集通信（每个应用流量特征不同，部分可隐藏）	密集&复杂IO
AI人工智能计算中心	AI训练推理	FP16半精度 BF16/FP8/FP4 FP32全精度		密集通信（部分可隐藏）	密集IO（按节奏迭代）
云数据中心	互联网云计算	INT32整型 FP32全精度	通用计算，高并发	分散通信	分散/密集IO



总结与思考



Summary

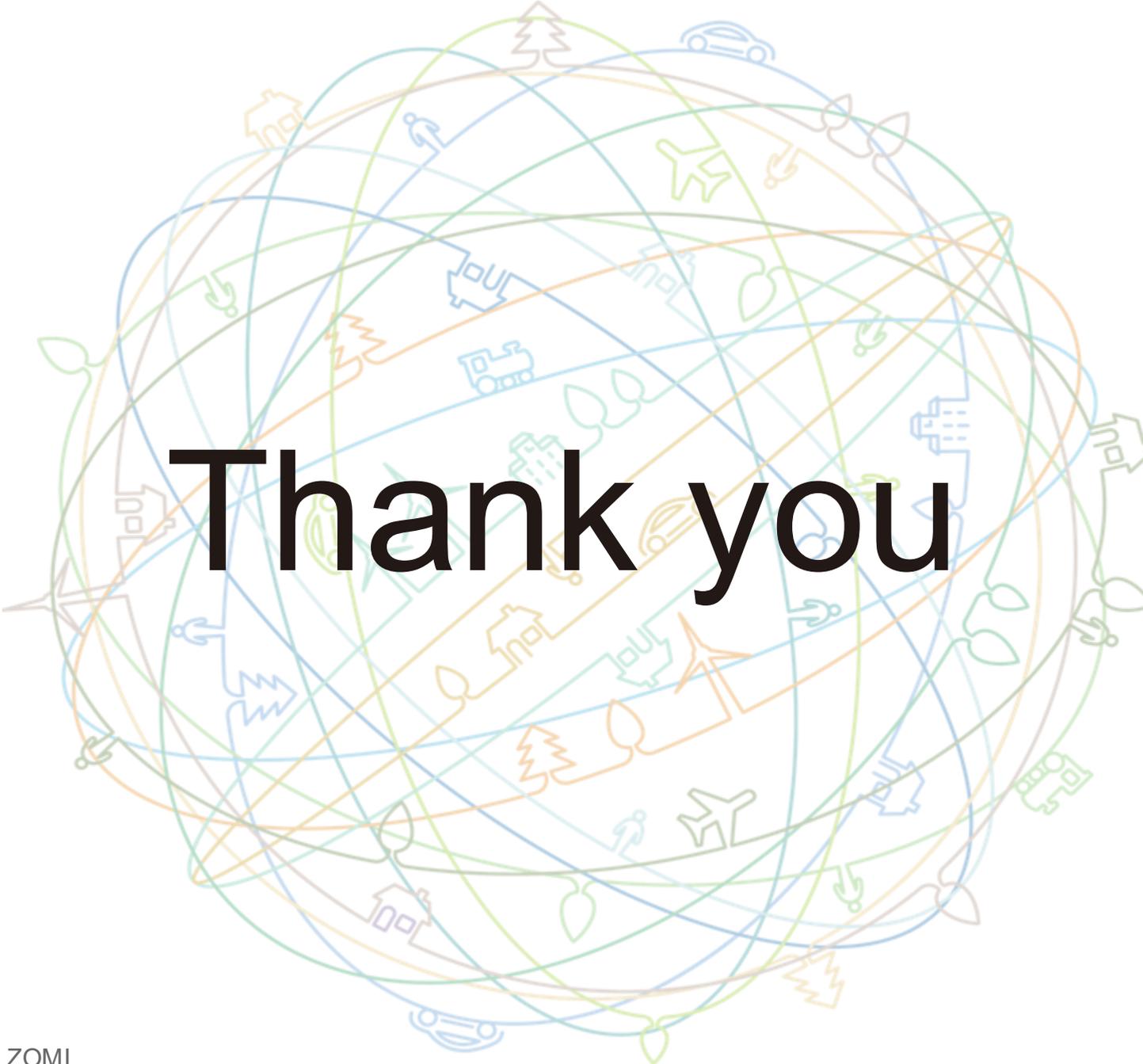
- 三类集群计算中心的区别：超算&智算中心为高并行计算，云数据中心为高并发计算
 - 超算中心：国之重器，主要进行FP64双精度的HPC高性能并行计算
 - 智算中心：智算利器，主要进行FP16半精度大模型并行训练和集中推理
 - 云数据中心：互联网云计算基座，主要进行整形和FP32全精度的高并发通用计算



Summary

- 集群计算是主要的算力集成设施建设方式，像水厂、电厂一样
- 算力基础设施建设，从分层堆叠走向集群全栈优化，达成规模系统最优





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 © 2024 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.



GitHub <https://github.com/chenzomi12/AllInfra>

引用与参考

- <https://zhuanlan.zhihu.com/p/683671511>
- PPT 开源在: <https://github.com/chenzomi12/AllInfra>

