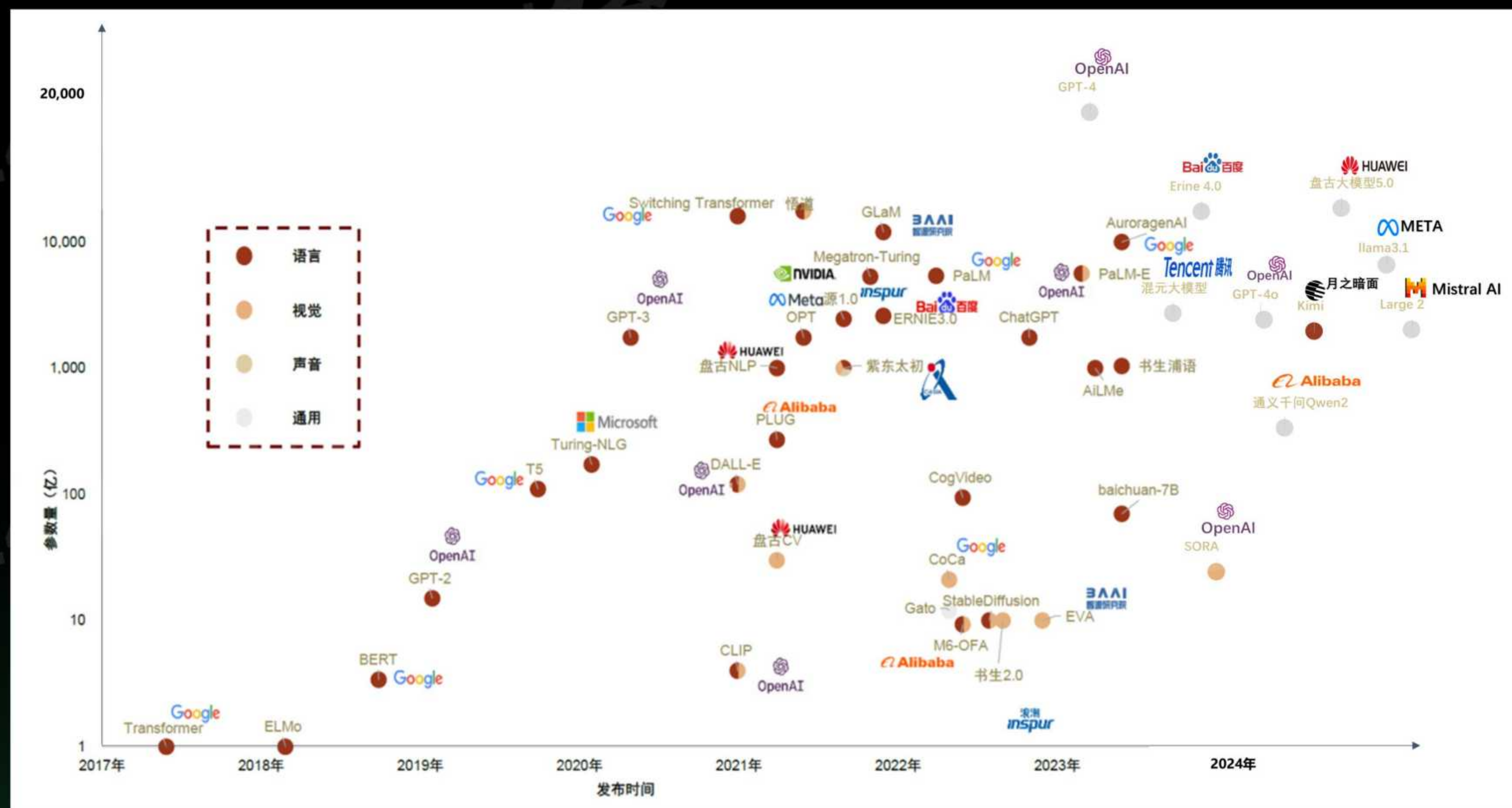


# 大模型重构AI基础设施

百度集团副总裁，OCTC轮值主席 侯震宇

# 大模型竞争白热化，模型厂商追求大参数+多模态

## 主流大模型参数量对比



资料来源：北京智源人工智能研究院、中金研究部

## 大模型更新迭代趋势

模型参数规模持续扩大：  
千亿->万亿

由单模态转向多模态：  
文字->图片音视频

# 2024年中国大模型应用步入落地期，AI原生应用有望爆发

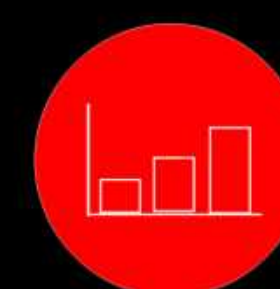
## 中国大模型市场未来5年预测增速突破50%



资料来源：上海社科院《全球数字经济竞争力发展报告（2023）》

## 以IDC为代表的多个机构预测

### 2024年中国大模型应用将进入落地期



垂直领域大模型的商业化应用正在加速



多模态大模型涌现，应用场景将更加丰富



大模型厂商打响“价格战”，降低AI应用研发门槛

## AI算力需求预计未来10年呈爆发式增长

	<u>2021</u>	<u>2030</u>	<u>21-30CAGR</u>
全球算力总规模	615 EFlops	56 ZFlops	65%
智能算力总规模	232 EFlops	53.5 ZFlops	80%
基础算力总规模	369 EFlops	3.3 ZFlops	27%

国外AI巨头纷纷布局构建**10万卡集群**

- **Meta**: 24年3月公布拥有两个GPU大集群 (每个~2.4万卡GPU), 预计24年底拥有35万卡GPU构建AI基础设施
- **微软&OpenAI**: 24年4月透露投资千亿美元打造“星际之门”超算
- **xAI**: 24年7月宣布启动全球最大AI集群 (10万卡液冷GPU) 进行训练

# 并行计算是实现大模型黄金法则Scaling Laws的最优解, ...

**Scaling Laws** 由OpenAI在2020年提出, 为基于Transformer的AI大模型的训练提供了重要指导:

**模型规模要大:** 增加模型参数量、数据集和计算量, 就可以得到性能更优的模型效果

**智能涌现:** 当模型规模达到一个阈值时, 模型会出现涌现特质——未预期到的新能力, 推动模型性能提升

## Scaling choices for pre-training

Goal: maximize model performance

CONSTRAINT:  
Compute budget  
(GPUs, training time, cost)



**并行计算** 是一种将复杂的问题分解成多个较小的部分, 然后同时处理这些部分以加快计算速度的计算方法

并行计算关键特征:

- **多处理器:** 并行计算依赖多个处理器 (如CPU、GPU) 独立处理任务
- **任务分解:** 将大任务分解成小任务是并行计算的核心
- **并发执行:** 这些小任务被同时执行, 减少总体完成任务所需的时间
- **通信与协调:** 在并行计算中, 不同的处理器或节点需要有效地交换信息并协调其工作, 这通常通过网络或快速数据总线实现

# ..., 对AI基础设施带来重构要求: 极致高密、极致互联、极致规模

		移动互联网+			
基础设施	服务器	X86/ARM	风冷	整机柜	X-BOX
	网络	万兆/25G/100G	自研交换机		云网融合
	机房	自建/合建	预制化/模块化	东数西算/碳中和	
	调度	合池混布	在离线一体	云原生	存算分离
	运营	韧性供应		资源一体化交付	

**通用服务器为主**  
极致弹性、极致高效、极致性价比

人工智能+			
冷板AI服务器	多芯CPU	高密存储	
800G/T级互联	超大规模组网	光互联	
AIDC	源网荷储	风液融合	算电协同
异地异网异构调度	训推一体	存算分离	
柔性供应		灵活资源配置	

**AI服务器快速增长**  
极致高密、极致互联、极致规模

# AI新基建驱动产业链全面变革，多个赛道迎来发展新机遇

1 芯片&服务器



## 极致高密

向高算力、高性能、高密度演进

## 产业变革

- 服务器功耗大幅增加，液冷技术规模化应用
- 国际形势严峻，芯片国产化加速，多芯多元要求异构大规模组网能力

2 网络




## 极致互联

超高带宽、超低时延、可规模化扩展

- 柜内纳秒互联引领，大厂纷纷布局
- 柜外光互联成趋势，光模块及交换机产业迎来升级机遇
- 异构异网异地，跨AZ/Region RDMA网络

3 IDC



## 极致规模

资源布局、超高能耗、运管效率

- 政策引导八大节点聚集，绿电有明确比例
- IDC高效节能的供电及制冷技术加速迭代，源网荷储探索
- 智能运维，算电协同、AI技术加持，提高IDC运营效率及管理力

# 百度AIDC硬件产品及方案：面向大模型时代布局

**AI高密机柜**  
大模型时代新一代算力解决方案，  
支撑超大规模集群



高密算力、高密供电、高速互联

**“XMAN5.0” AI计算机**  
创新AI超级计算机，业内领先、  
引领行业架构设计



异构多芯、兼容液冷

**“乾坤”通用计算系统**  
新一代模块化服务器，兼具成  
本和弹性，实现高效交付



算力多元化、灵活高效

**“太行” DPU3.0**  
极致零损耗，高密虚拟机，超百万  
IOPS云盘，提供极致应用体验



400Gbps、Multihost、共池超发



# 百度智能云百舸平台：打造专业的AI基础设施，全面支持开源

## 大模型任务增强

### 大模型训推任务加速镜像

开源大模型定制优化

高性能算子

高效显存利用

高效并行策略

高性能训推框架

### 大模型IO加速方案

Flash checkpoint

大镜像预加载

大规模镜像P2P加速

30%

训练吞吐提升

98.8%

有效训练时长

60%

推理吞吐提升

95%

带宽有效性

支持万卡级别超大规模AI计算

超稳基础设施和自动化容错保障

丰富的运维和可观测工具

模块化AI组件设计 可灵活集成

## 百舸组件

### AI基础组件

高性能网络插件

高性能存储插件

异构资源调度

### AI编排调度

深度学习框架

AI任务编排

任务 workflow 管理

### 稳定性&容错

多维故障感知

自动任务容错

通信测试工具

### 可观测大盘

集群资源视图

任务稳定性大盘

性能监控&调优

## 百舸资源池

### CCE K8S集群

异构算力

高性能分布式存储PFS

万卡RDMA网络

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

**THANKS**

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会

2024开放计算中国峰会