

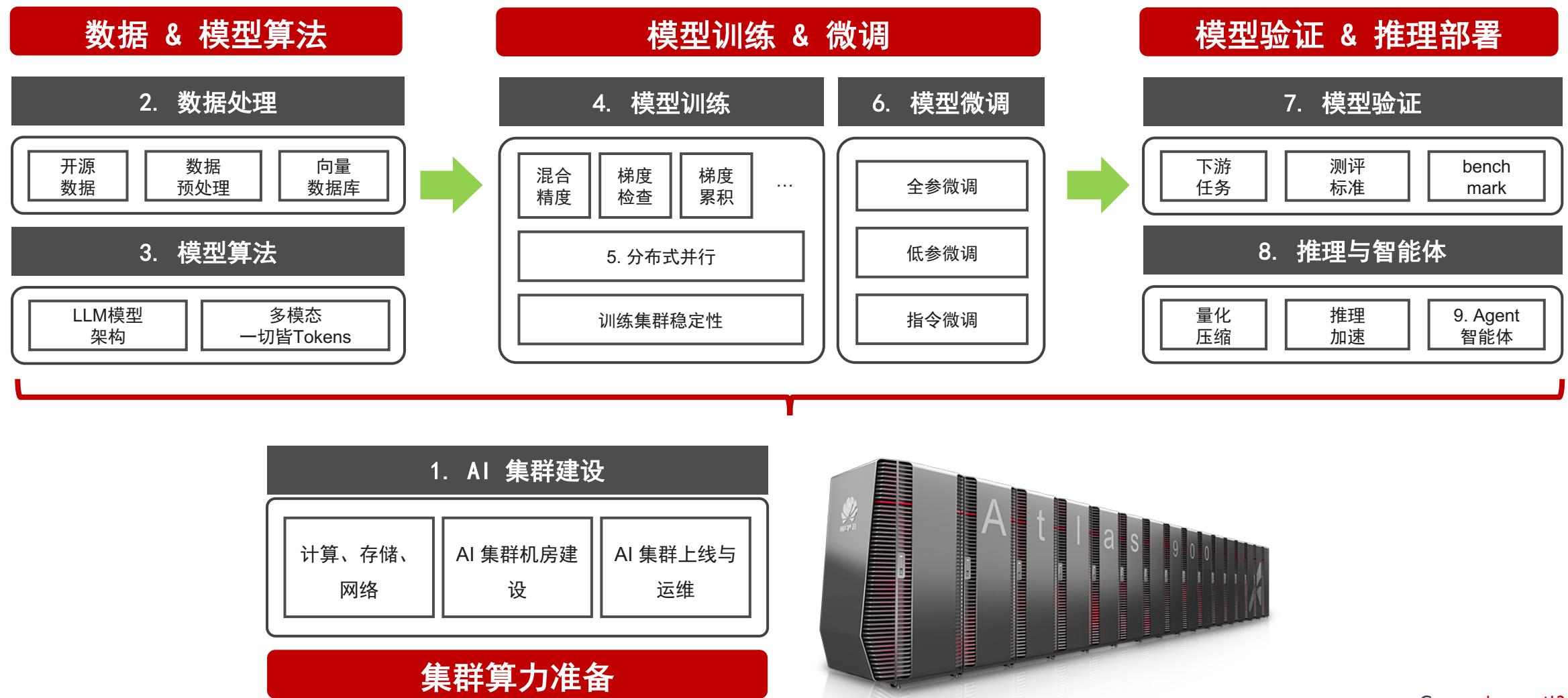
大模型-AI集群(存)

存算力架构的思考



ZOMI

# 大模型业务全流程



# 关于本内容

- **内容背景**

- AI 集群 + 大模型

- **具体内容**

- **数据存储现状和场景**：存储软件类型、存储硬件类型的发展
  - **大模型对存储的挑战**：存储性能指标、存储遇到大模型挑战与新机会点
  - **大模型训练CKPT优化**：大模型训练过程、CKPT过程分解、CKPT优化
  - **大模型时代对存储的思考**：什么样的存储架构才是AI大模型时代的选择？

# 1. 大模型存储

思考

# 大模型全流程

- 大模型全流程：数据获取、数据处理、模型训练、模型微调、模型评估、模型部署
- 全流程每个阶段都涉及数据的存储与访问，各阶段的数据也需要协同

# 观点 I

I. 宏观上走向存算一体，走向近存加速；微观上会走向存算分离架构。



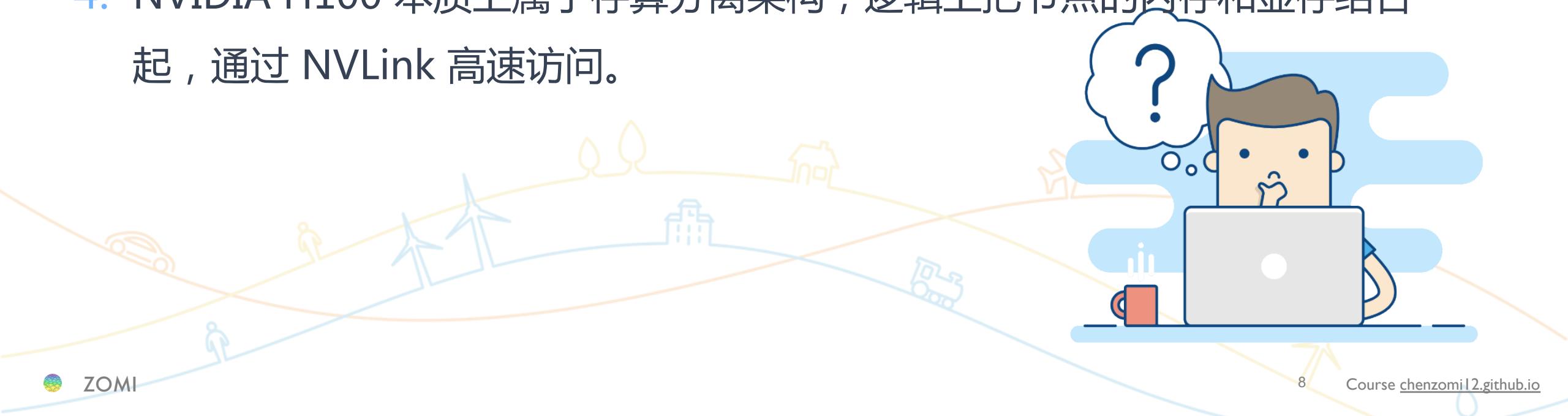
# 近存计算的好处

1. 减少数据在片内和片外的流动搬运，所带来的性能开销，减少时延；
2. 从而降低能耗，是系统架构发展的本质要求所驱动。

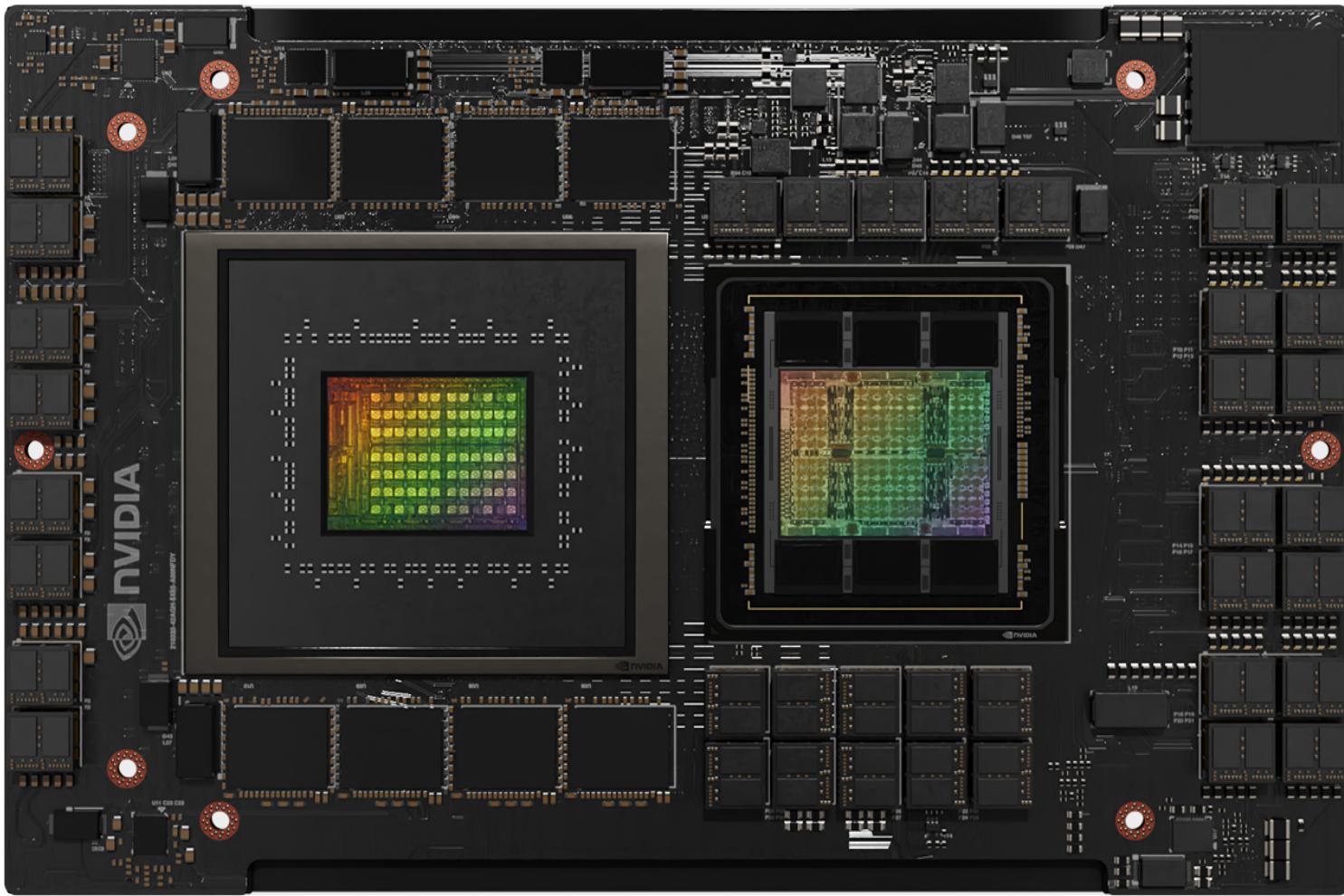


# 存算分离的好处

1. 计算算力和存储能够根据自身业务动态 Scale，更容易扩展；
2. 计算、存储和网络关键指标不同，计算重吞吐，存储重稳定性，网络重时延；
3. 系统复杂度不同，存储对于系统复杂度容错率低，高可用度要求严苛；
4. NVIDIA H100 本质上属于存算分离架构，逻辑上把节点的内存和显存结合一起，通过 NVLink 高速访问。



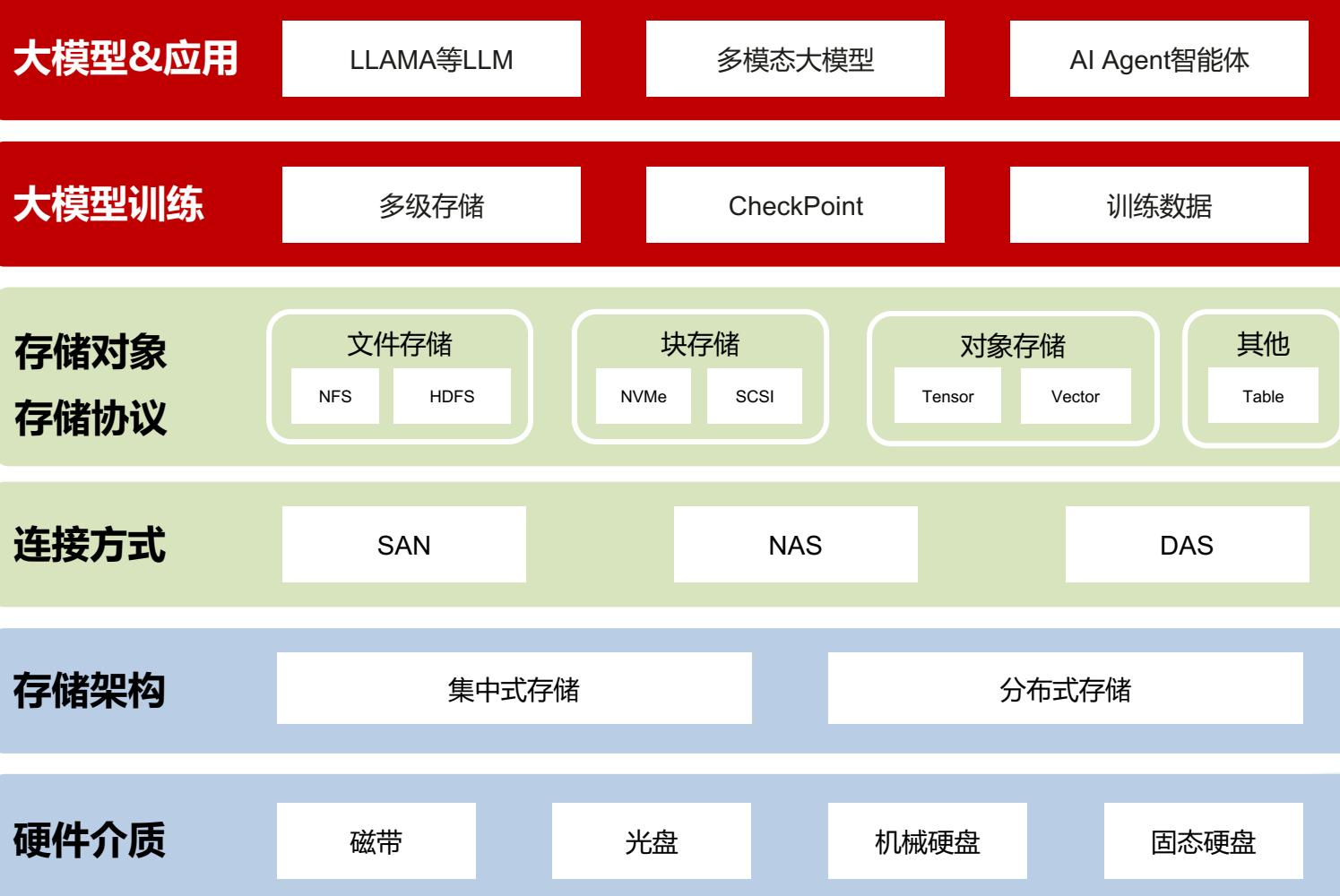
# NVIDIA GH200



# 2. 小结与回顾

# 大模型 & 存储技术架构

## 存储 技术 架构



# 存储硬件介质的组织层次

- 容量更小
- 速度更快
- 价格更高

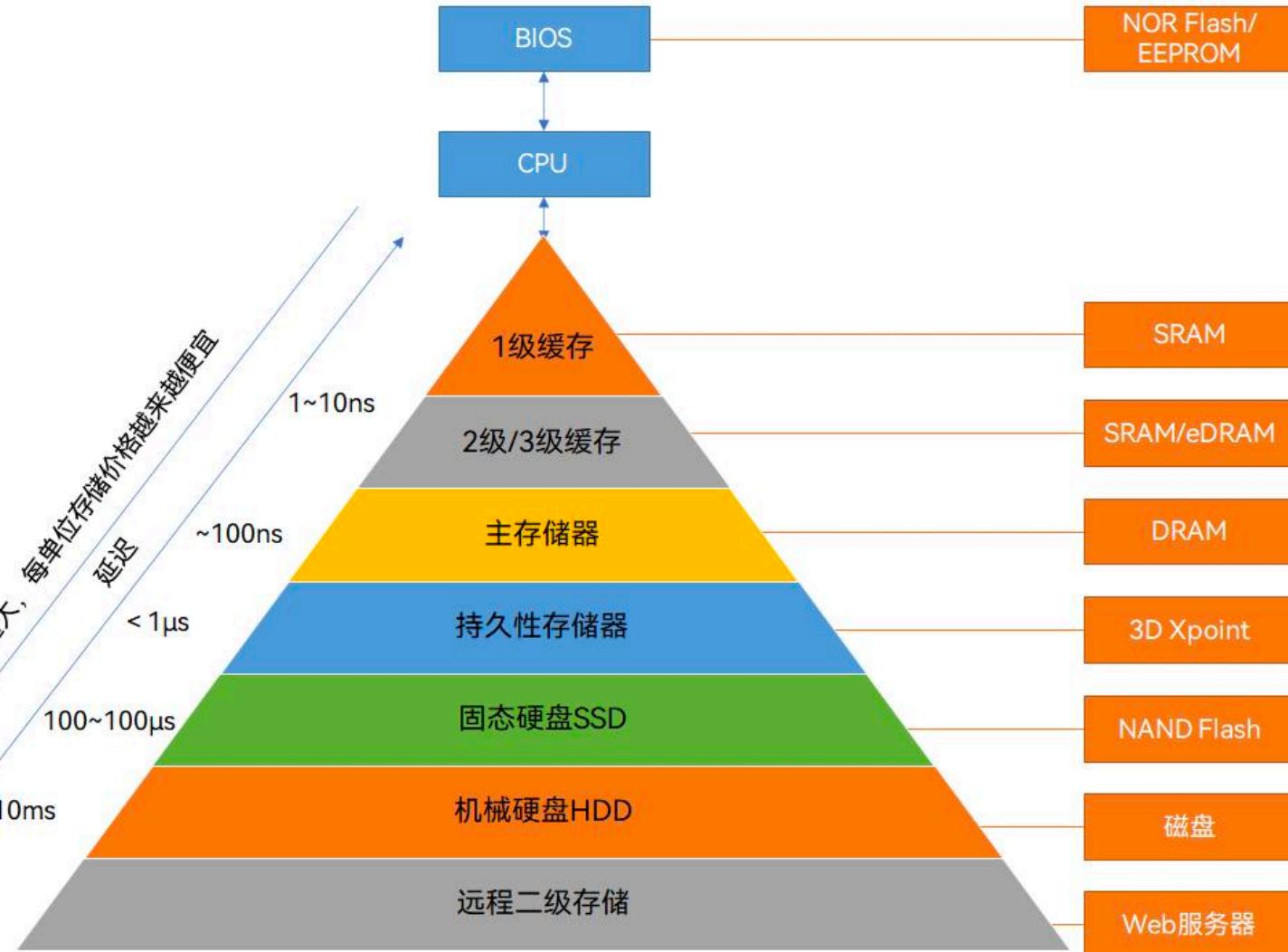


- 容量更大
- 速度更慢
- 价格更低

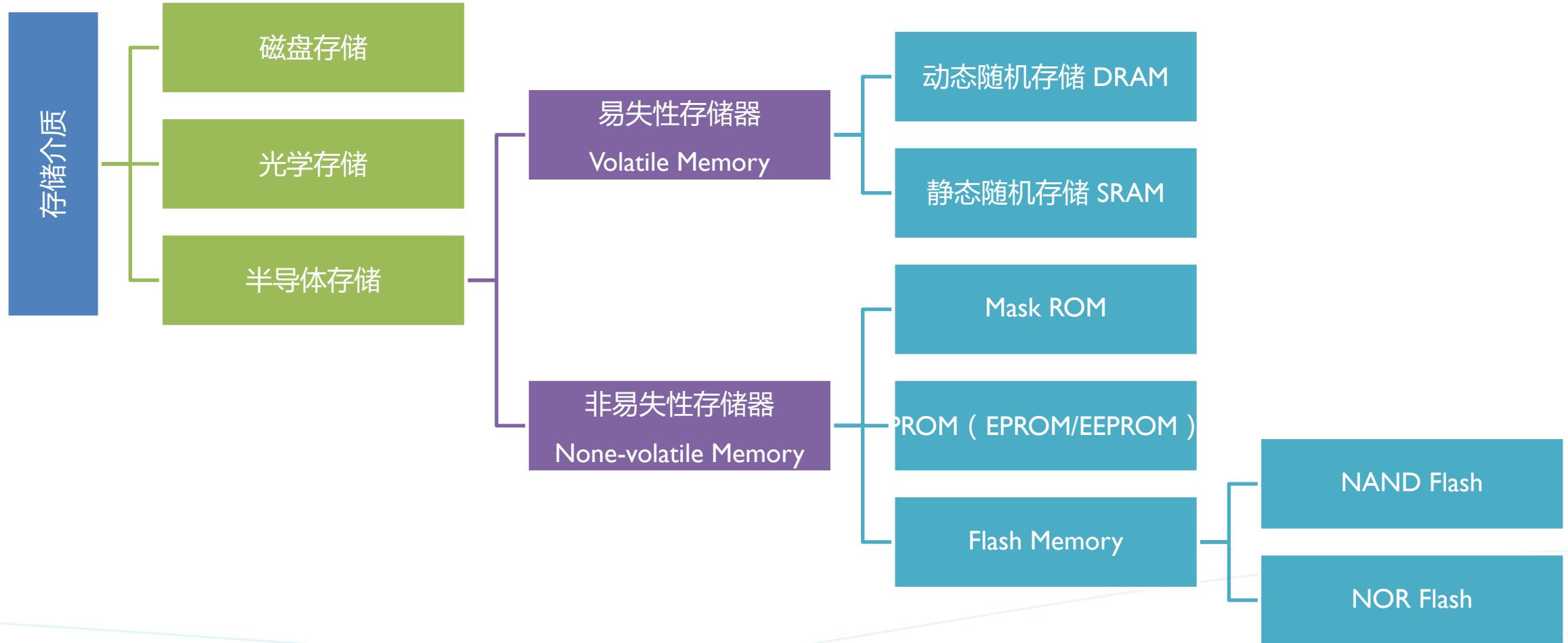


容量越来越大、每单位存储价格越来越便宜  
延迟

1~10ns  
~100ns  
< 1μs  
100~100μs  
~10ms

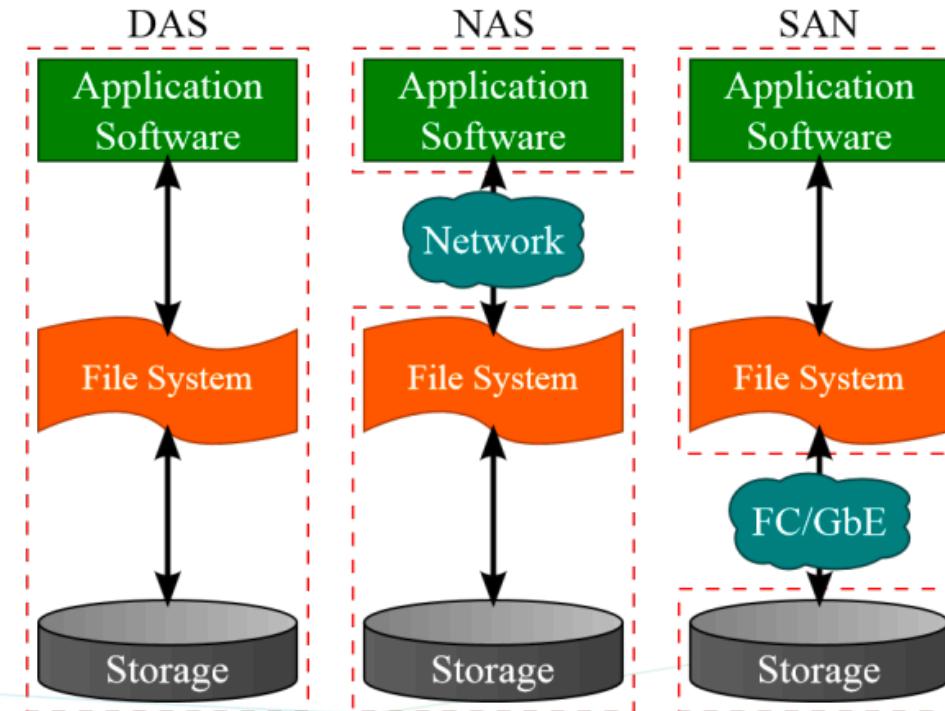


# 存储硬件介质的分类

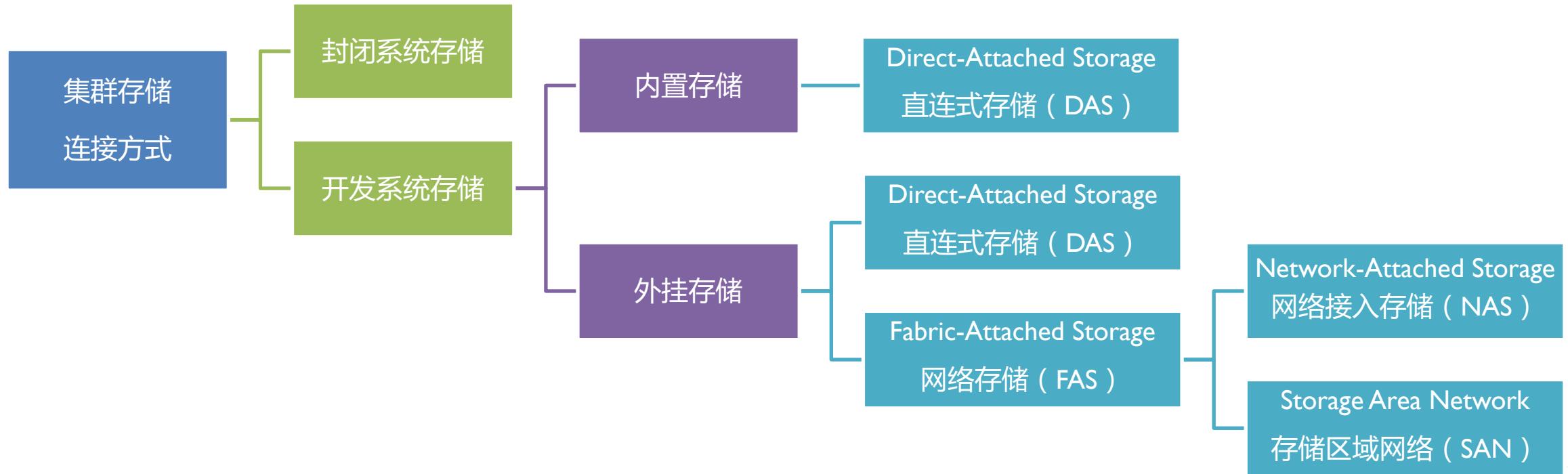


# 存储连接方式的DAS、SAN、NAS 关系

- DAS 存储一般应用在中小企业，与计算机采用直连方式，NAS 存储则通过以太网添加到计算机上，SAN 存储则使用 FC 接口，提供性能更加的存储。NAS 与 DAS 主要区别体现在操作系统所在位置。

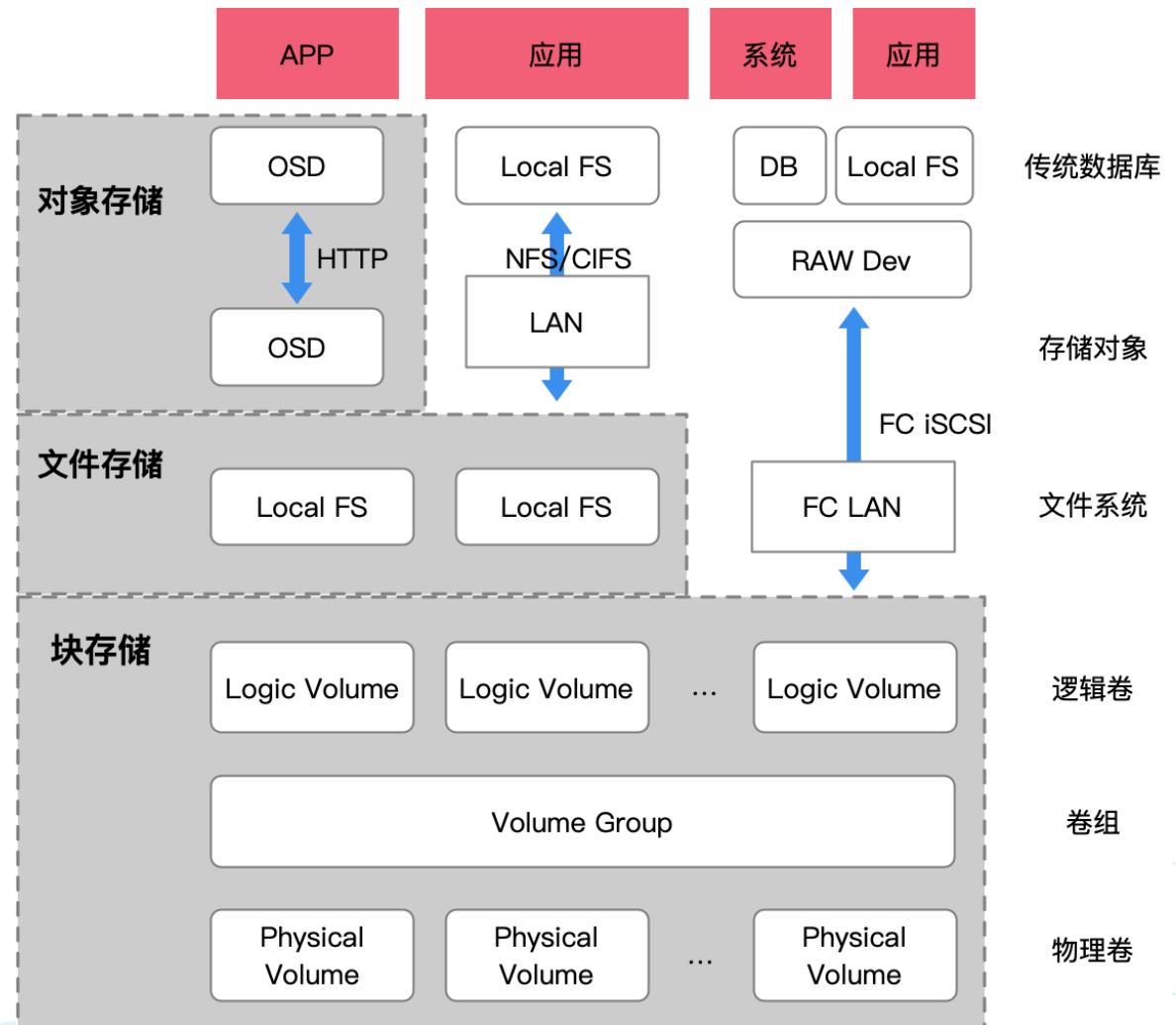


# 存储连接方式的分类



# 存储对象和存储协议

- 三者分别提供不同层级、不同粒度的存储方式，是一个逐层递增关系；
- 块存储是字节级存储，文件存储增加了文件层级语义，对象存储封装了对象语义，屏蔽底层细节。

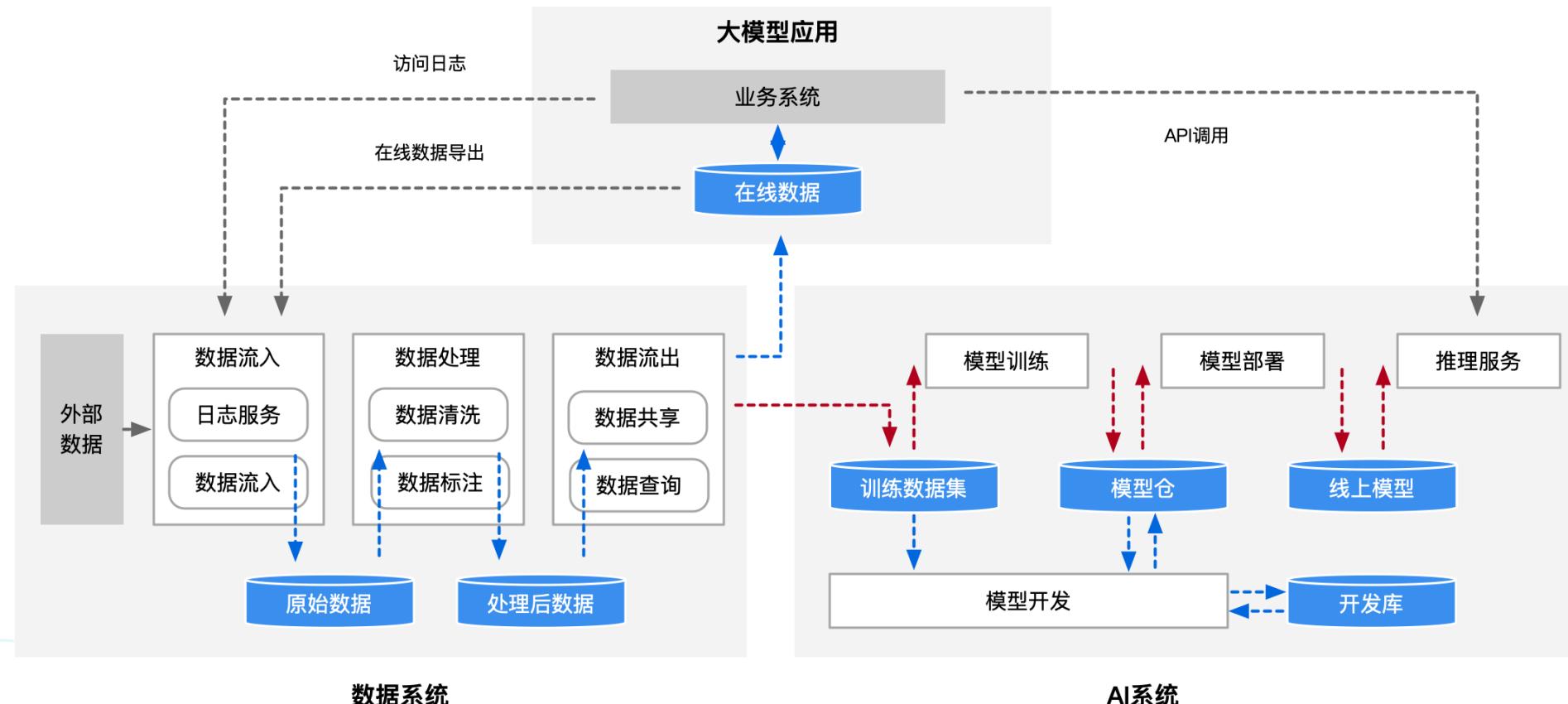


# 大模型对存储的需求和挑战



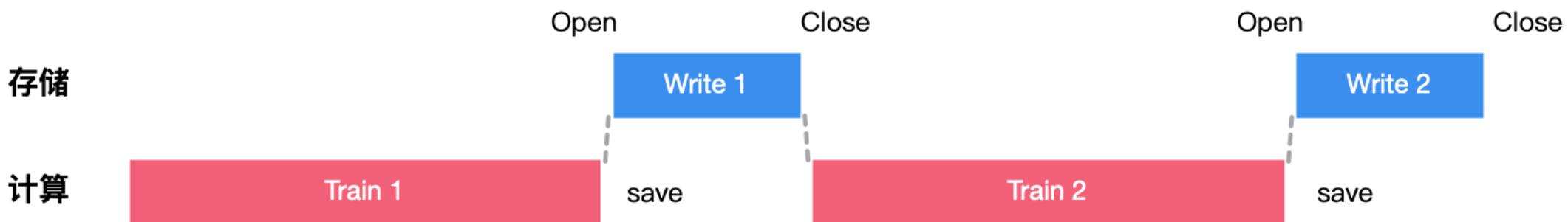
# 大模型在存储的优化手段

- 1) 大模型训练过程语料数据持续更新，2) 训练数据在不同数据预处理流程中频繁流动。有必要针对非结构化和半结构化数据提供专用存储系统，例如数据湖。



# 训练 Checkpoint 优化方案

优化思路	优化路径	成熟度
将CKPT模型数据存放在数据湖	save() and/or load()	*****
CKPT的save过程从同步到异步	save()	**
CKPT流式分块存储	save()	**
多文件加速聚合	save() and/or load()	***
本地内存缓存，同步写内存	save() and/or load()	***
数据拷贝过程使用零拷贝	save() and/or load()	





# 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](https://chenzomi12.github.io)

GitHub [github.com/chenzomi12/DeepLearningSystem](https://github.com/chenzomi12/DeepLearningSystem)