

大模型时代 OS & AI 的思考与展望

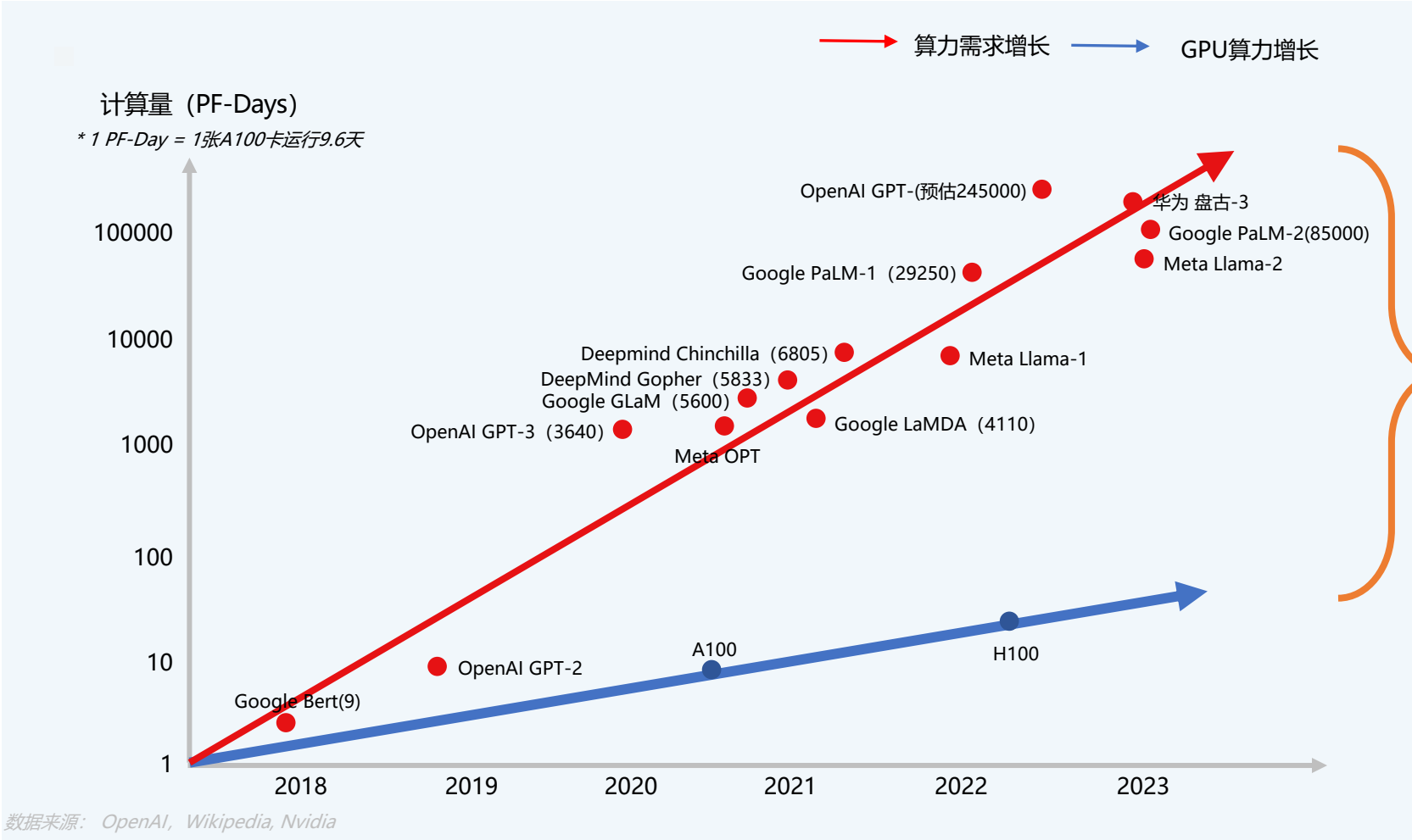
目录

1 为什么需要 OS for AI

2 OS for AI 优化

3 AI for OS 展望

大模型训练算力需求每年增长 **10** 倍，GPU算力年增 **1** 倍



软件优化

OS可为AI业务提供全生命周期优化



OS for AI: 开发场景优化

挑战：快速开发、快速测试、快速发布

模型管理

预置大模型

LLaMA

BLOOM

GLM

自有数据

+

精调

新模型

新模型

新模型

镜像管理

预置深度学习框架镜像

PyTorch

[M]^s
MindSpore

TensorFlow

K Keras

Caffe

训练管理

预置训练框架

DeepSpeed

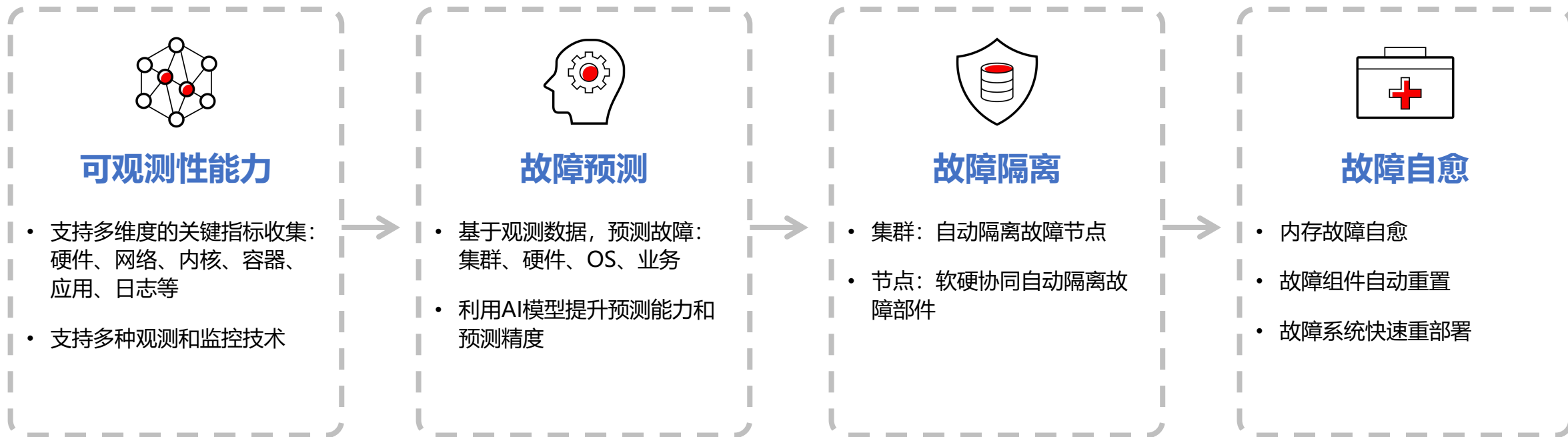
Megatron-LM

Horovod

AscendSpeed

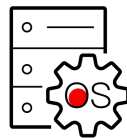
OS for AI: 训练场景优化

挑战：保障长时间、大集群的训练任务稳定执行，不被中断



OS for AI: 部署、运维优化

挑战：快速完成AI基础设施环境安装环境部署，并通过基本验证

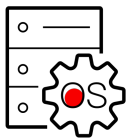


部署步骤繁琐

- OS安装
- NPU固件刷新
- GPU/NPU驱动
- IB驱动
- CUDA/CANN安装
- 容器镜像准备
- GPU/NPU直通配置
- 集群部署

OS for AI: 部署、运维优化

挑战：快速完成AI基础设施环境安装环境部署，并通过基本验证



部署步骤繁琐

- OS安装
- NPU固件刷新
- GPU/NPU驱动
- IB驱动
- CUDA/CANN安装
- 容器镜像准备
- GPU/NPU直通配置
- 集群部署

自动化运维

2
周

2
天

基础软件

算力调度

Volcano

集群K8S

容器镜像

NPU/GPU

AscendCL

CANN

CUDA

操作系统

FusionOS

固件/驱动
管理

系统配置

硬件

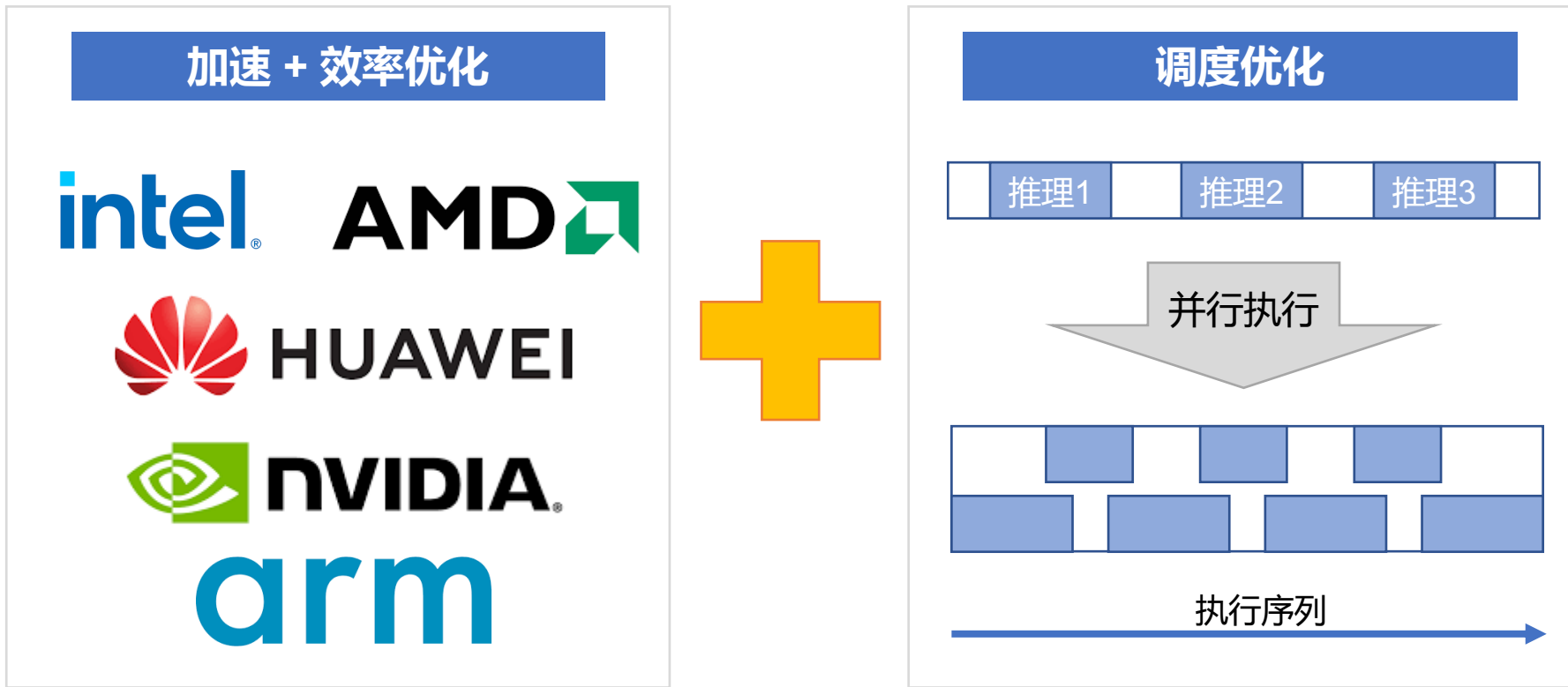
高性能AI服务器

低延迟RDMA网络

高吞吐存储

OS for AI: 推理场景优化

挑战: 推理设施利用率低



展望 AI for OS

挑战：AI可以为OS提供哪些智能化助力？OS需要怎样的AI能力？



展望 AI for OS

挑战：AI可以为OS提供哪些智能化助力？OS需要怎样的AI能力？



需要一个新的AI框架

- 为推理任务优化
- 高性能
- 依赖少
- 体积小
- 安全可靠

THANKS

THANKS

THANKS

THANKS