

大模型：分布式训练

DeepSpeed



ZOMI

大模型业务全流程



大模型系列 – 分布式训练加速

- 具体内容

- 分布式加速库 :

- 业界常用分布式加速库 & 作用

- DeepSpeed 特性 :

- 基本概念 - 整体框架 – Zero-1/2/3 – ZeRO-Offload – ZeRO-Infinity

- Megatron 特性 :

- 总体介绍 – 整体流程 – 并行配置 – DP – TP – PP

2. DeepSpeed

<https://github.com/microsoft/DeepSpeed>

创新特性

- **Training**：为大模型训练供 ZeRO、3D-Parallelism、DeepSpeed-MoE、ZeRO-Infinity等特性；
- **Inference**：提供Tensor、Pipeline、Expert等并行特性与推理内核、通信优化和异构内存特性结合；
- **Compression**：提供 ZeroQuant 和 XTC 等 SoTA 在压缩方面的创新；
- **4Science**：遥远的你；

Training	Inference	Compression	Science
<ul style="list-style-type: none">• Speed Scale Cost• Democratization• MoE models• Long sequence• RLHF	<ul style="list-style-type: none">• Large models• Latency• Serving cost• Agility	<ul style="list-style-type: none">• Model size• Latency• Composability• Runnable on client devices	<ul style="list-style-type: none">• Speed• Scale• Capability• Diversity• Discovery

软件架构

1. APIs :

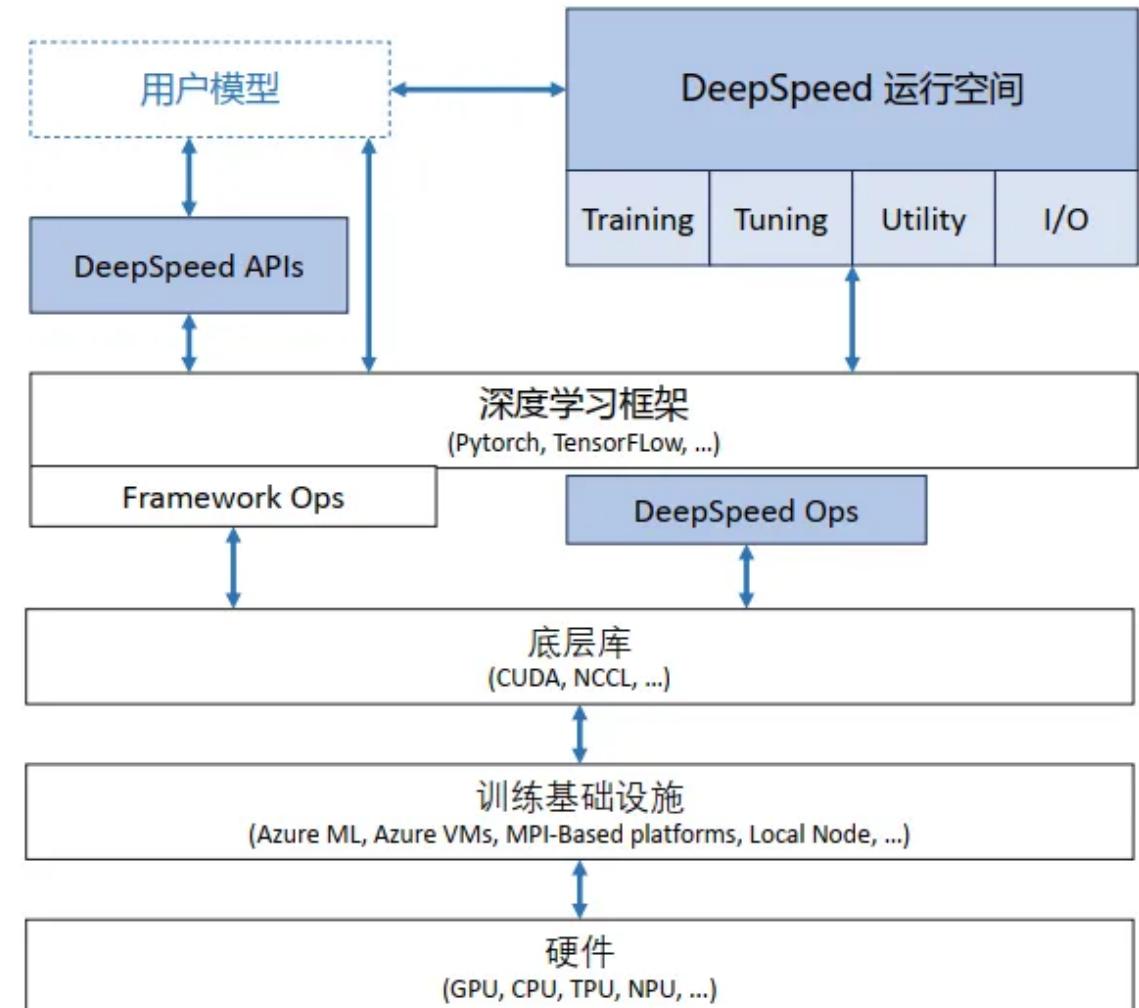
- 配置参数在 `ds_config.json` 中，通过 API 接口可以调用 DeepSpeed 训练/推理模型；

2. RunTime :

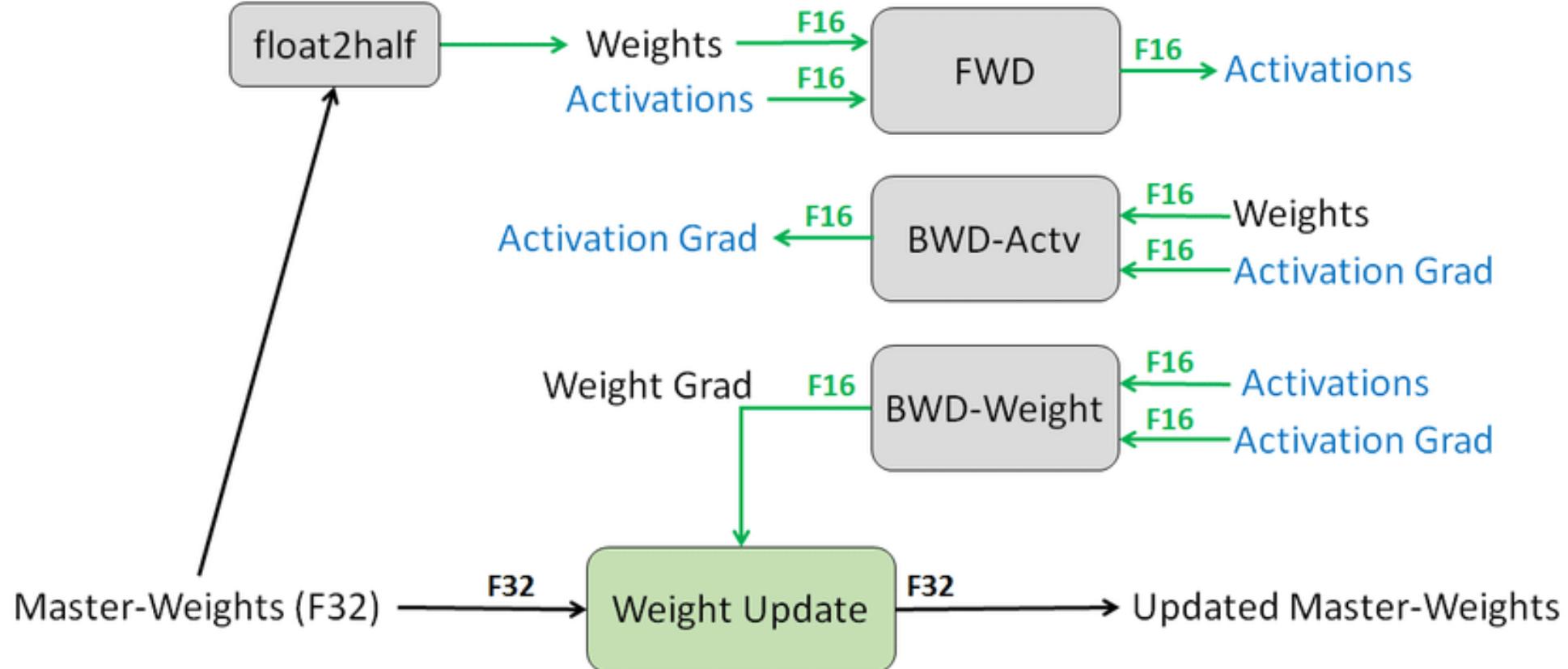
- 核心运行时组件，负责管理、执行和优化性能。包括数据、模型、并行优化、微调、故障检测以及 CheckPoint 保存和加载等任务。

3. Ops :

- 底层内核组件，使用C++和CUDA实现。优化计算和通信，提供底层操作；



混合精度训练 && 显存占用分析



混合精度训练 && 显存占用分析

- **Model States** 模型本身相关且必须存储的参数：
 1. Parameters : 模型参数
 2. Gradients : 模型梯度
 3. Optimizer States : Adam中 momentum和variance
- **Residual States** 非模型必须，训练过程中产生的参数：
 1. Activation : 激活值
 2. Temporary Buffers : 临时存储
 3. Unusable Fragmented Memory : 碎片化存储空间

ZeRO 并行

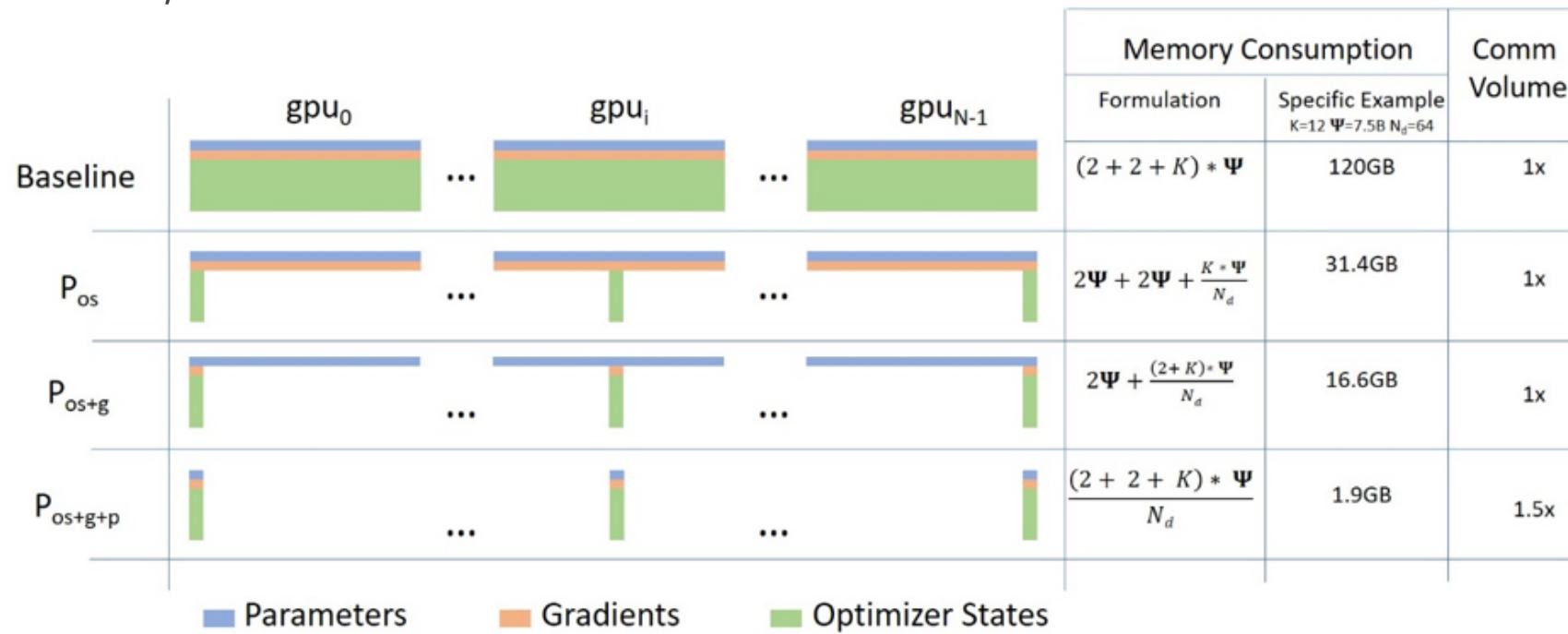
- ZeRO Zero Redundancy Optimizer , 一系列显存优化方法的统称：
 1. ZeRO-DP (Data Parallel) : ZeROI/2/3
 2. ZeRO-R (Reduce) : Activation Checkpointing、 Memory Defragmentation
 3. ZeRO-Offload : Offload Strategy && Offload Schedule
 4. ZeRO-Infinity : Breaking the GPU Memory Wall for Extreme Scale Deep Learning

ZeRO-DP (Data Parallel)

- ZeRO Zero Redundancy Optimizer , 一系列显存优化方法的统称：
 - I. ZeRO-DP (Data Parallel) : ZeRO I/2/3
- Optimizer state partitioning (ZeRO stage 1) :
 - 只对 optimizer 状态进行切分，占用内存原始1/4；
- Gradient partitioning (ZeRO stage 2) :
 - 对 optimizer 和 grad 进行切分，占用内存原始1/8；
- Parameter partitioning (ZeRO stage 3) :
 - 对 optimizer、grad 和模型参数进行切分，内存减少与数据并行度和复杂度成线性关系，同时通信容量是数据并行性的 1.5 倍；

ZeRO-DP (Data Parallel)

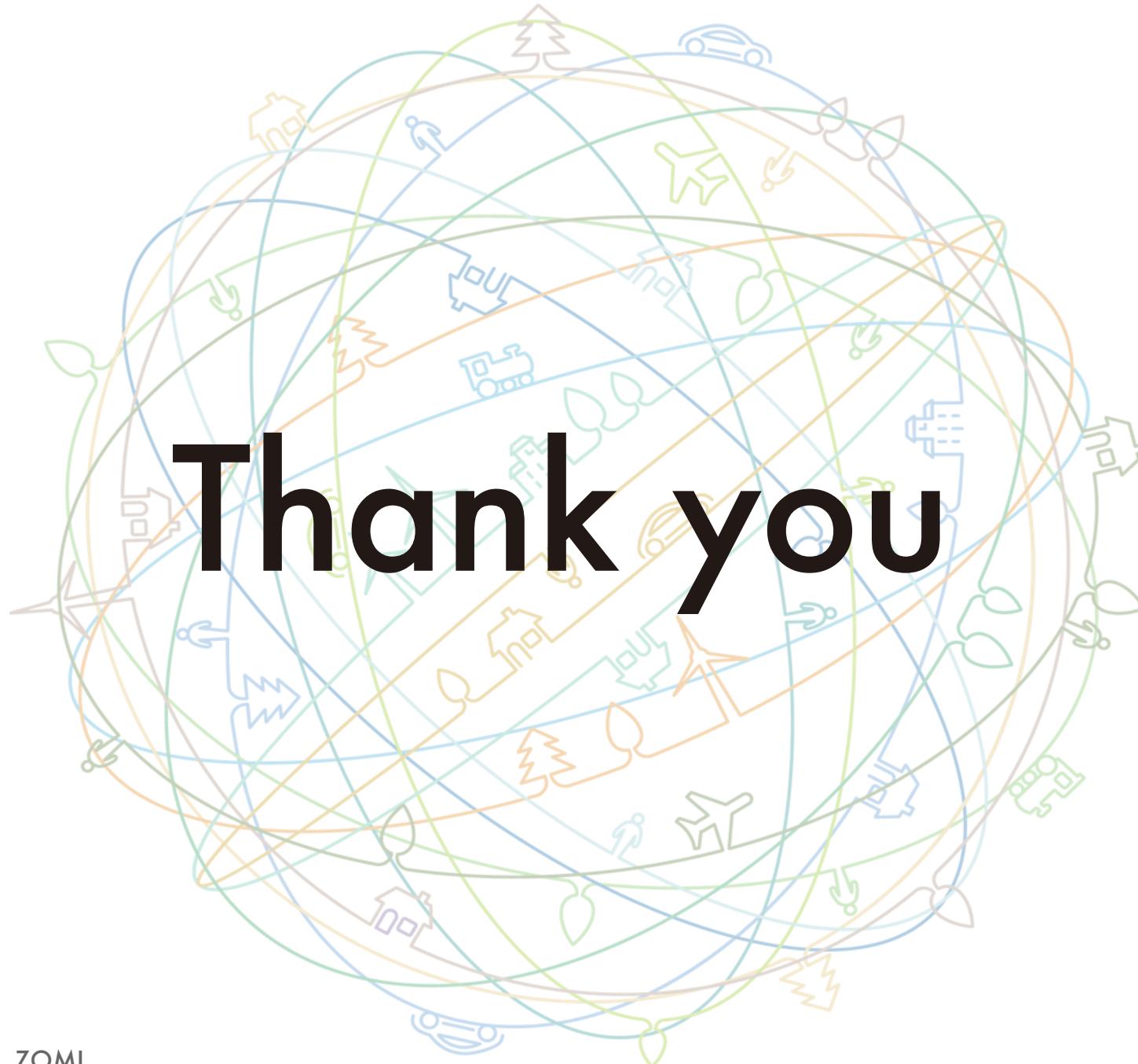
- Optimizer state partitioning (ZeRO stage 1) : 只对 optimizer 状态进行切分，占用内存原始 1/4；
- Gradient partitioning (ZeRO stage 2) : 对 optimizer 和 grad 进行切分，占用内存原始 1/8；
- Parameter partitioning (ZeRO stage 3) : 对 optimizer、grad 和模型参数进行切分，内存与数据并行线性关系，通信 DP 1.5X；



ZeRO: Where the memory goes?

- 模型本身相关且必须存储的参数
- Model States：
 1. Parameters (half) : 2 bytes
 2. Gradients (half) : 2 bytes
- Optimizer States：
 1. Master Weight (fp32) : 4 bytes
 2. Adam momentum (fp32) : 4 bytes
 3. Adam Variance (fp32) : 4 bytes

- 假设模型参数量 Ψ ，使用Adam优化器混合精度训练：
 - a. 模型参数和梯度 float16，显存消耗为 $2\Psi + 2\Psi$ ；
 - Adam 维护 float32 模型副本，消耗 4Ψ ；
 - Adam 辅助变量 fp32 momentum + fp32 variance，显存消耗 $4\Psi + 4\Psi$ ；
- 模型消耗 $2\Psi + 2\Psi = 4\Psi$ ，Adam 消耗 $4\Psi + 4\Psi + 4\Psi = 12\Psi$ ，总消耗 $\Psi + 12\Psi = 16\Psi$
- 优化器显存占用表示 $K\Psi$ （不同的优化器不同），则、混合精度训练显存占用 $4\Psi + K\Psi$



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