

方缙

✉ fangjin98@mail.ustc.edu.cn

📞 (+86) 181-5566-1676

🌐 www.fangjin.site



教育背景

中国科学技术大学 硕博连读 计算机科学与技术 2020.9-2026.6 (预计)

- 研究方向: 分布式训练、在网计算
- 导师: 徐宏力、赵功名

湖南大学 学士 计算机科学与技术 2016.9-2020.6

- 湖南大学优秀毕业论文

项目经历

开源项目 **Triton-Distributed** 字节跳动-Seed (AI Infra), 北京
核心开发者

- 负责跨机计算通信融合算子开发，开发实现 AG+GEMM、AG+MoE 在 H800 上的融合算子
- 负责 GPU 跨机通信优化，基于多 QP 技术优化 All-to-All 算子在大规模 H20 集群上的性能
- 负责项目日常维护等工作

基于 GPU 主导通信的集合通信库设计与实现 字节跳动-Seed (AI Infra), 北京
主要开发者 2024.12-至今

- 分析 DeepEP 通信模式与性能瓶颈，复现实验并撰写分析报告
- 研究 GPU 与网卡性能调优机制，优化内部 AlltoAll 算子性能
- 设计通用 GPU 主导通信库，针对 AMD、寒武纪等平台实现 GPU 主导通信

针对自研硬件实现集合通信算子 字节跳动-Seed (AI Infra), 北京
主要开发者 2024.10-2025.1

- 基于 RDMA 编程为自研硬件实现集合通信库，包括 AllGather 和 AlltoAll 算子。
- 针对主机带宽受限的自研机型，形式化 ILP 问题并设计主机选路算法优化性能
- 在网卡 1 拖 2/4/8 场景下优化集合通信算法，优化后性能达到物理极限带宽

计算通信融合算子自动生成 字节跳动-Seed (AI Infra), 北京
主要开发者 2024.6-2025.3

- 分布式训练/推理任务中计算和通信算子交替执行，通过挖掘计算通信可重叠部分，编写融合算子，可以加速分布式训练/推理任务
- 基于 NVSHMEM 为 Triton 扩展分布式编程功能，能够直接借助 Triton 生态编写集合通信算子
- 实现跨卡和跨机的融合算子，分析 GPU 全局内存读写成本，优化算子性能
- 相比较非融合算子，实现性能提升 1.17 倍到 20.76 倍

大规模分布式训练任务在光网络中的节点部署优化 华为 2012 中央研究院, 合肥
主要开发者 2023.12-2024.5

- 大模型的分布式训练任务具有算力亲和性，然而算力由多层机间网络互联、带宽异构，导致跨节点网络成为训练瓶颈
- 调研现有大模型任务部署和算力调度优化方案，熟悉常见模型并行和数据并行方法
- 调研现有模型压缩工作，熟悉稀疏模型训练优化方法

- 对不同集合通信算法下的物理节点和逻辑节点通信模式建模，分析通信拓扑、链路对任务训练时间的影响
- 设计任务部署算法降低光网络下跨机架任务通信量

科研经历

基于可编程网络实现精准模拟网络故障

中科大苏高院，苏州

主要开发者

2022.12-2023.9

- 基于端主机的故障注入难以覆盖大量复杂网络故障场景，且无法针对应用流量精准注入故障
- 基于可编程控制面设计并实现用户友好的多后端故障注入系统，提供一系列参数供用户自定义流量协议
- 针对流依赖和流过滤，设计一个解析器生成算法，能够根据用户指示生成对应数据面程序
- 针对多租户和路由路径，形式化故障注入点选择问题
- 针对异构多后端网络设备，基于 P4 TNA 和 PSA 架构实现多种网络功能，系统资源消耗小于 10%
- 在 4 个流行的分布式系统任务 (Horovod, Redis, RDMA, Kafka) 中测试该系统并验证故障注入效果

使用可编程交换机加速分布式模型训练

之江实验室，杭州

主要开发者

2022.6-2022.9

- 大规模分布式模型训练具有通信瓶颈，该项目通过可编程交换机在网内聚合梯度以降低通信量，从而加速分布式模型训练
- 针对流量可变性，设计基于随机舍入算法解决网内聚合场景下的梯度路由问题
- 基于 Pytorch 实现包含 8 台服务器的 PS 架构分布式模型训练原型系统，主机间通过自定义协议进行通信（链接）
- 基于 Intel Tofino 可编程交换机实现网内聚合逻辑，并与主机端实现协同训练
- 相比较现有方案，降低分布式训练通信负载 81.2%

学术成果

1. S. Zheng, W. Bao, Q. Hou, X. Zheng, **J. Fang**, etc, *Triton-distributed: Programming Overlapping Kernels on Distributed AI Systems with the Triton Compiler*, (**Arxiv**)
2. S. Zheng, **J. Fang**, X. Zheng, Q. Hou, W. Bao, N. Zheng, Z. Jiang, D. Wang, J. Ye, H. Lin, L. Chang, X. Liu, *TileLink: Generating Efficient Compute-Communication Overlapping Kernels using Tile-Centric Primitives*, (**MLSys'25**)
3. **J. Fang**, G. Zhao, H. Xu, L. Luo, Z. Yao, A. Xie, *Non-Idle Machine-Aware Worker Placement for Efficient Distributed Training in GPU Clusters*, (**ICNP'24**), CCF B
4. **J. Fang**, G. Zhao, H. Xu, Z. Yu, B. Shen, L. Xie, *Accelerating Distributed Training with Collaborative In-network Aggregation*, (**ToN'24**), CCF A
5. J. Liu, Y. Zhai, G. Zhao, H. Xu, **J. Fang**, Z. Zeng, Y. Zhu, *InArt: In-Network Aggregation with Route Selection for Accelerating Distributed Training*, (**WWW'24**), CCF A
6. **J. Fang**, G. Zhao, H. Xu, C. Wu, Z. Yu, *GRID: Gradient Routing with In-network Aggregation for Distributed Training*, (**ToN'23**), CCF A
7. **J. Fang**, G. Zhao, H. Xu, Z. Yu, B. Shen, L. Xie, *GOAT: Gradient Scheduling with Collaborative In-Network Aggregation for Distributed Training*, (**IWQoS'23**), CCF B
8. **J. Fang**, G. Zhao, H. Xu, H. Tu, H. Wang, *Reveal: Robustness-Aware VNF Placement and Request Scheduling in Edge Clouds*, (**ComNet'23**), CCF B

奖项荣誉

- 中国电科十四所国睿奖学金 2023
- 英特尔 P4 中国黑客松优胜奖 2022
- 一等学业奖学金（博士）×2 2022, 2023
- 一等学业奖学金（硕士）×2 2020, 2021