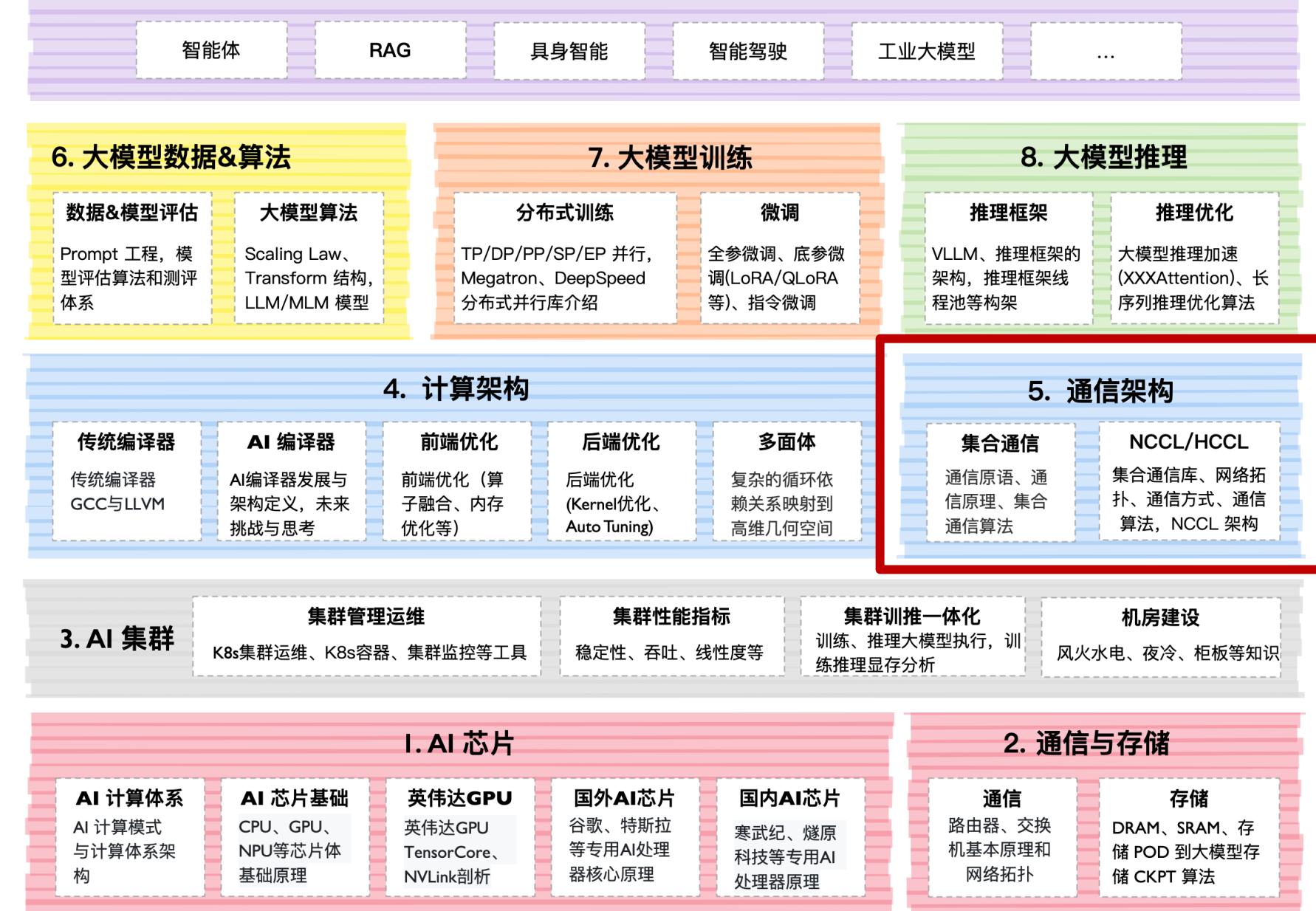


大模型系列 - 集合通信库

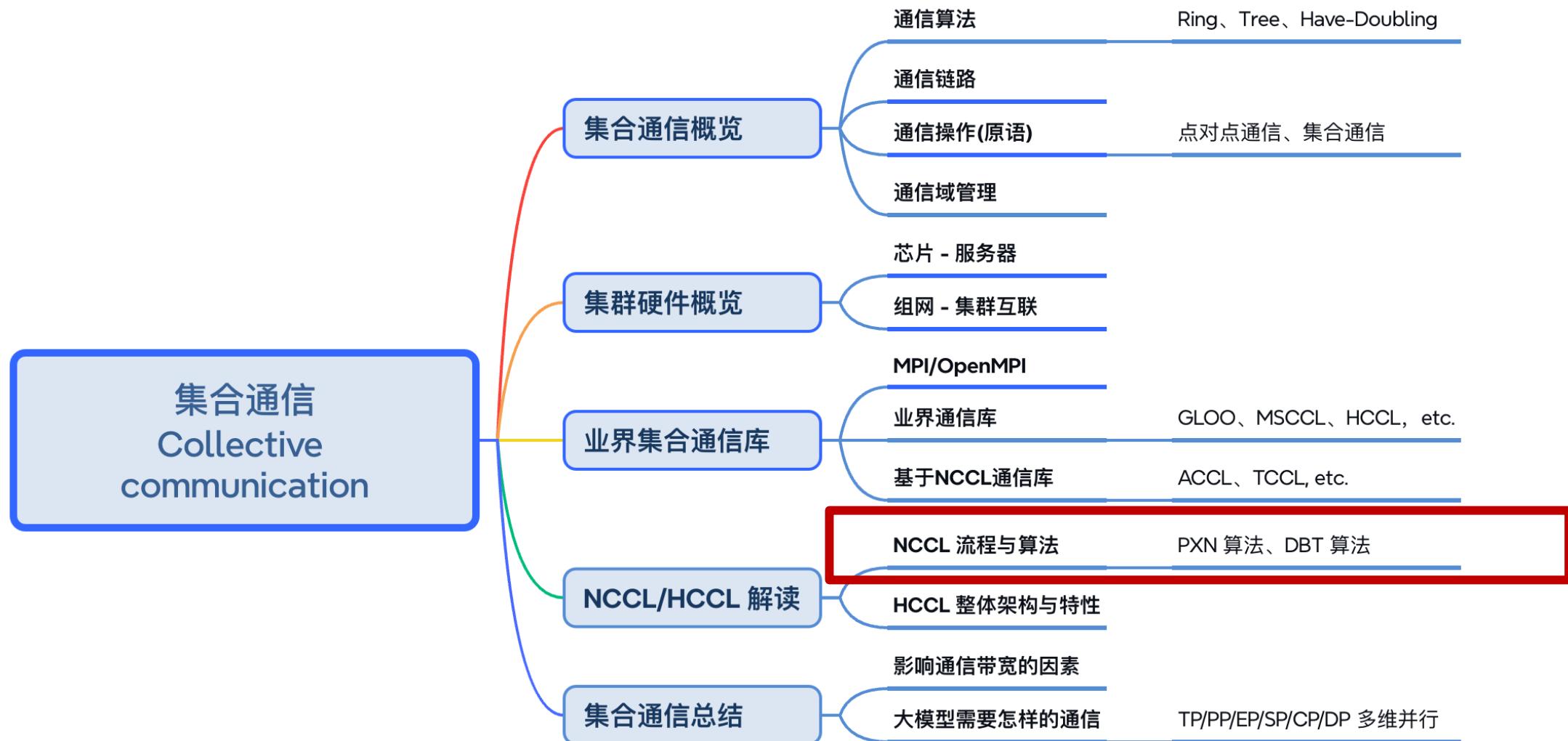
拓扑与算法结合



ZOMI



思维导图 XMind



Question

I. 如何在 GPU 集群上，网络拓扑和通信算法协同优化，提升集合通信的性能？



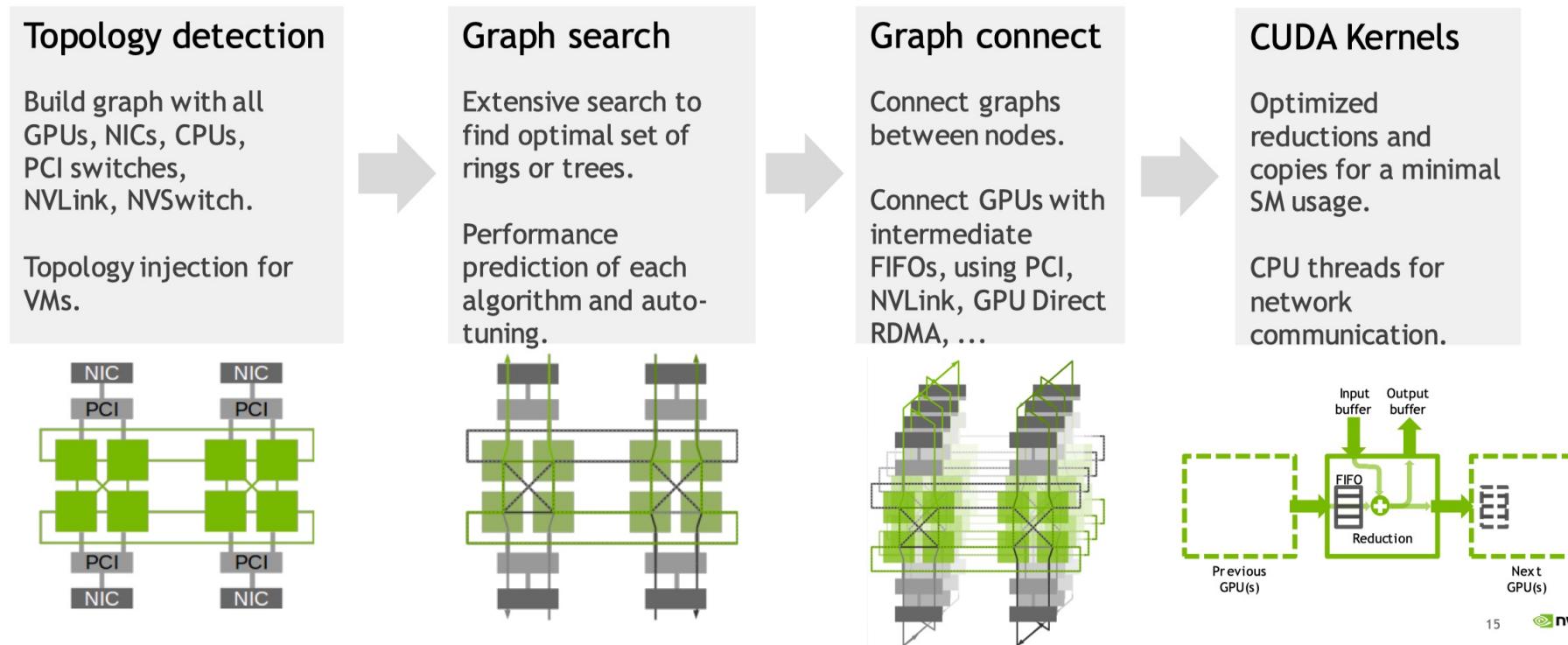
01. NCCL 简介

Basic Concept



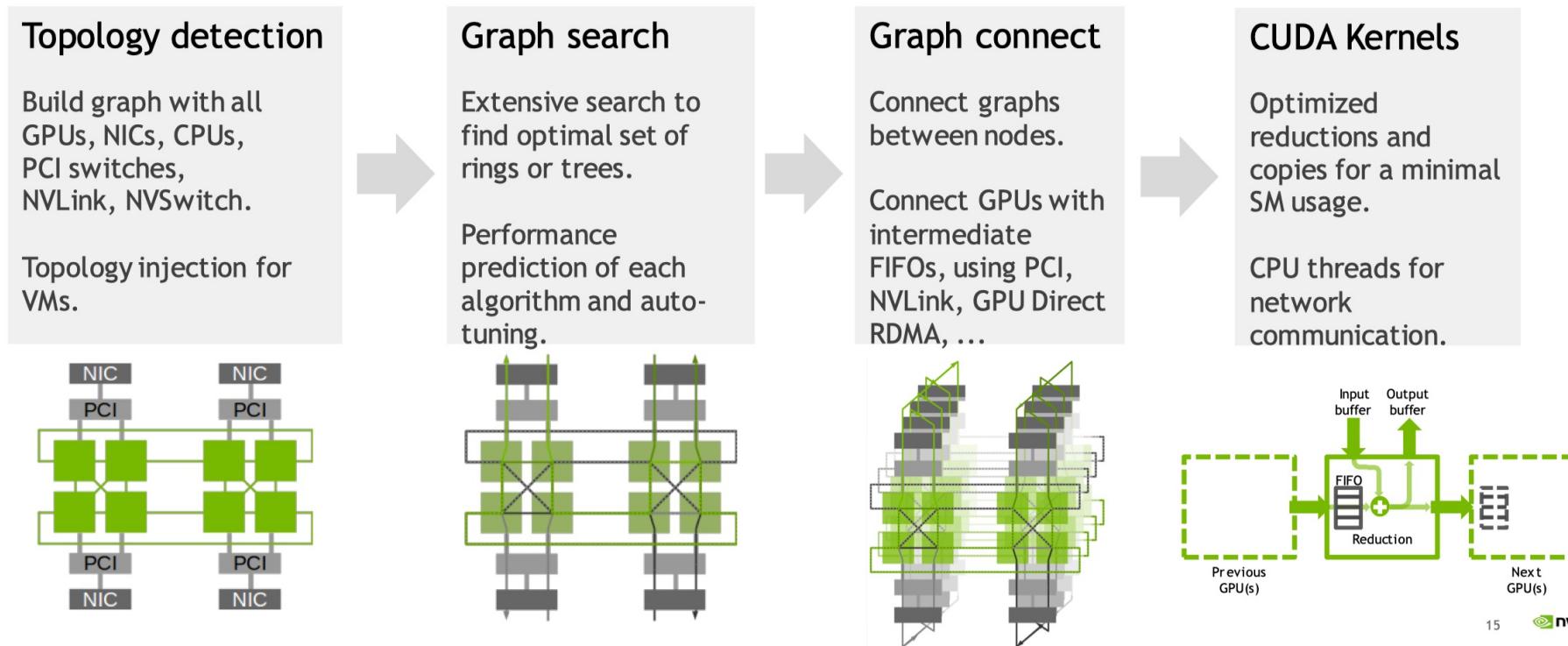
NCCL 基本介绍

- NVIDIA 集合通信库 (NCCL) 是一个类 MPI 通信库，可实现 GPU 的集合通信算法相关操作：
- all-gather / all-reduce / broadcast / reduce / reduce-scatter / point-to-point send and receive



NCCL 基本介绍

- NCCL 具有拓扑感知能力，经过优化可通过 PCIe、NVLink、以太网和 InfiniBand 互连实现高带宽和低延迟。NCCL GCP 插件支持自定义网络连接的云环境中实现高性能 NCCL 操作。



02. 通信结合

网络拓扑



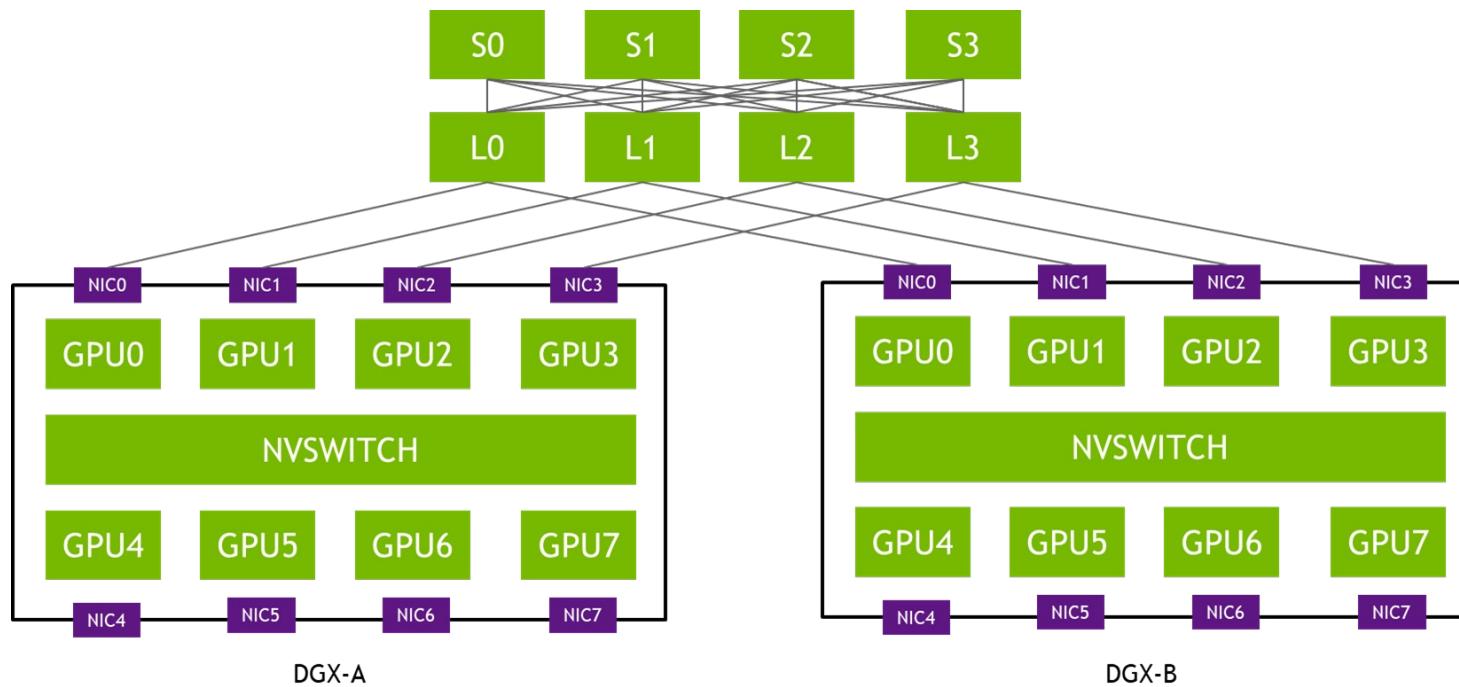
PXN: PCI × NVLink

- 2022 年 NCCL 2.12 引入新功能 PXN，即 PCI × NVLink，使 GPU 能够通过 NVLink 与节点上的 NIC 进行通信。
 - 不需要使用 QPI 或其他 CPU 协议。使得每个 GPU 仍然尝试尽可能多地使用本地 NIC 提升集合通信算法，当然也可以访问其他NIC。
-
- 快速通道互联 (Intel QuickPath Interconnect, QPI) ， Intel 开发并使用的点对点处理器互联架构，用来实现CPU之间的互联。

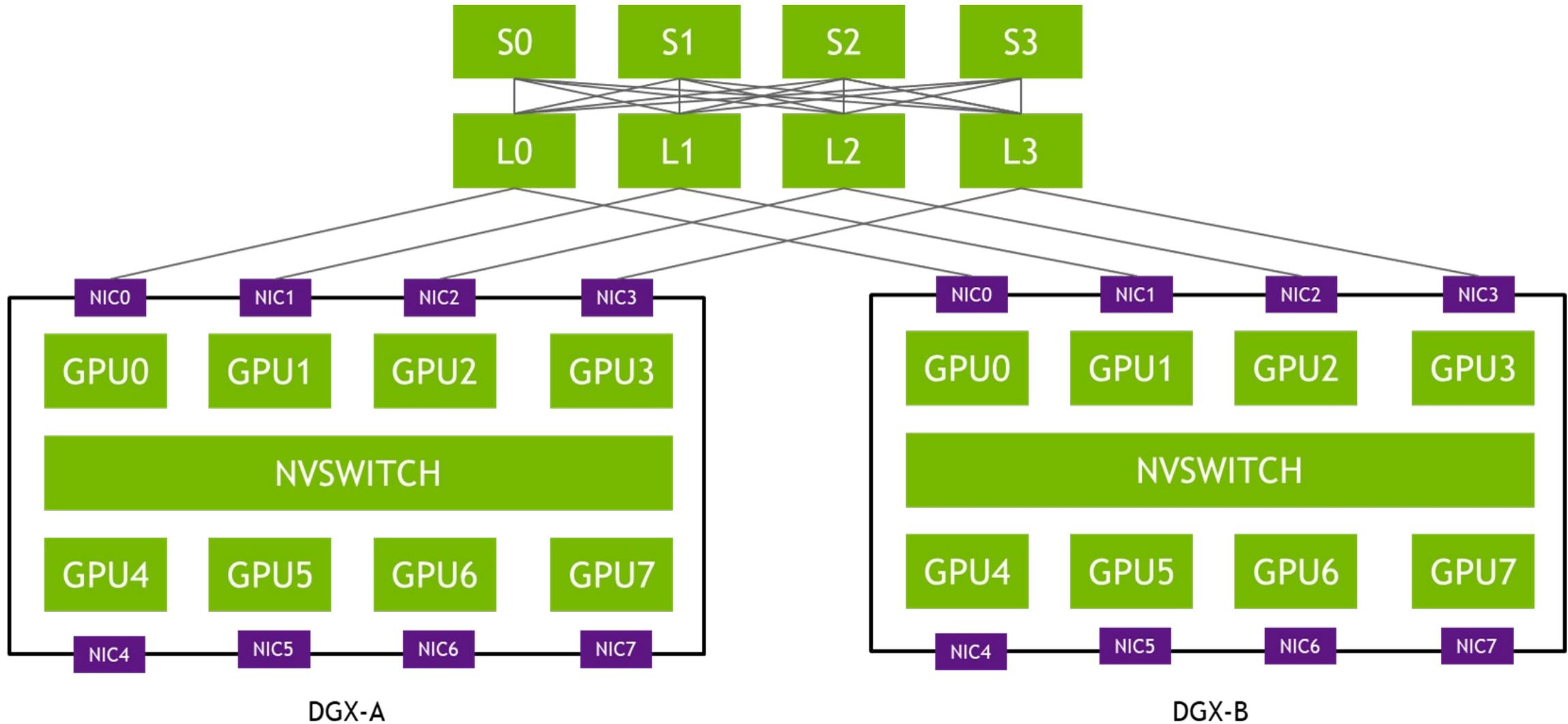


PXN: PCI × NVLink

- GPU 不是在本地内存上准备缓冲区供本地 NIC 发送，而在中间 GPU 上准备缓冲区，通过 NVLink 写入该缓冲区。接着通知管理该 NIC 的 CPU 代理数据已经准备好，而不是通知其自己 CPU 代理。



多轨通信



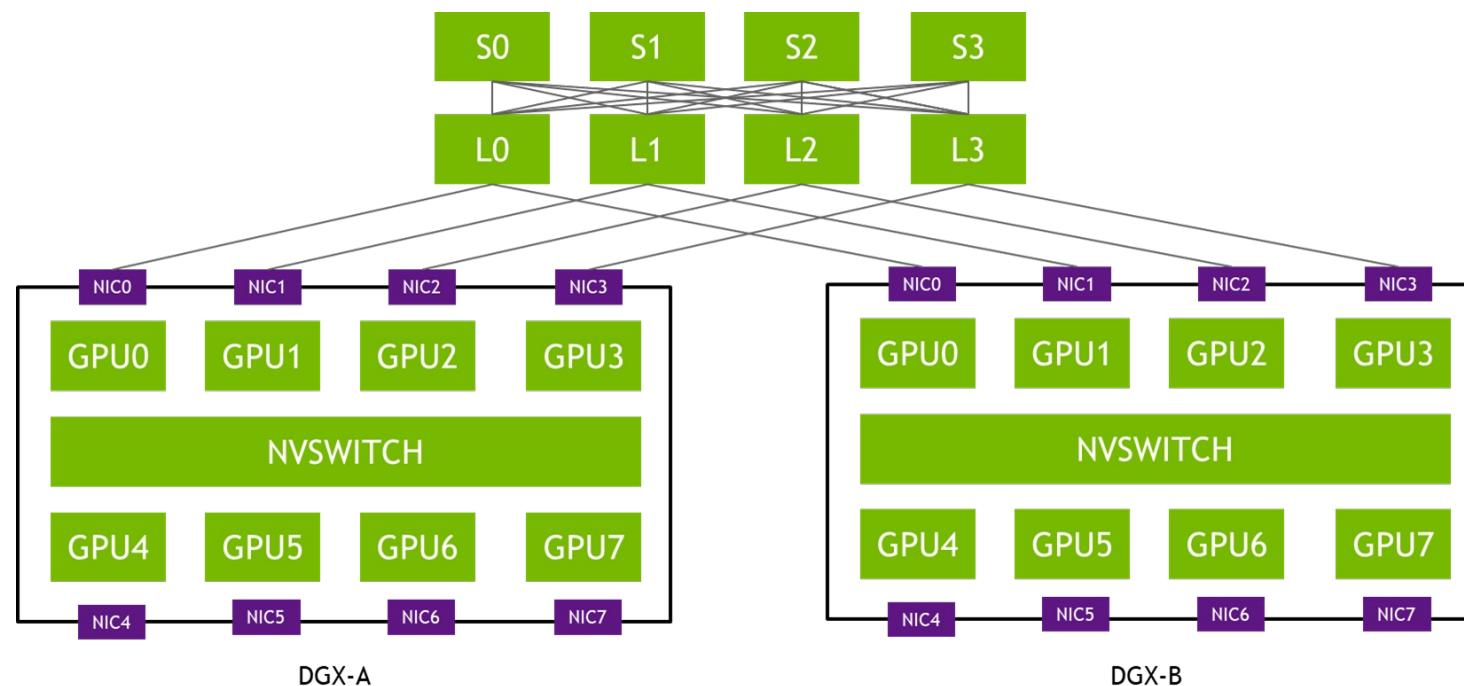
DGX-A

DGX-B

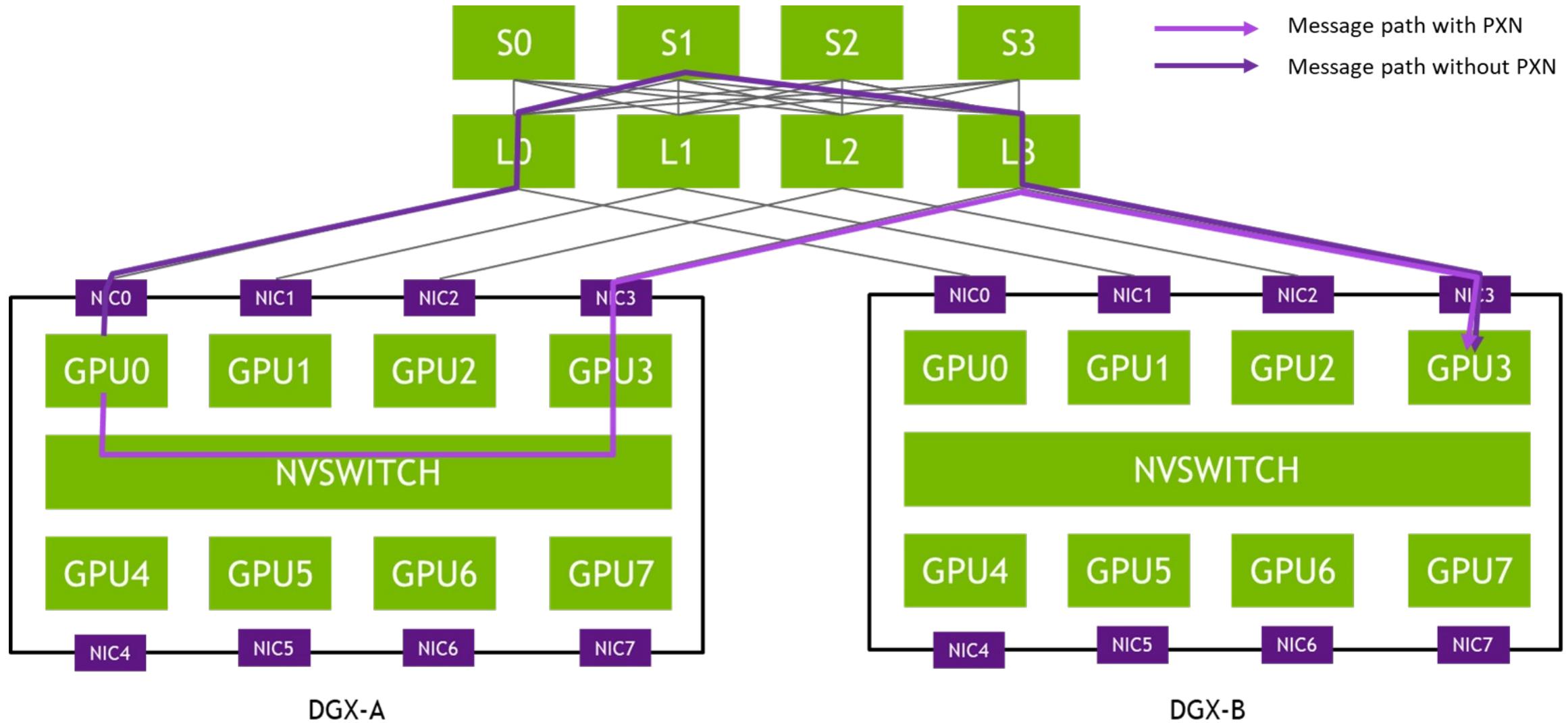


多轨通信

- PXN 利用节点内 GPU 间 NVSwitch 连接，首先将数据移动到与目的地位于同一轨道上 GPU，然后将其发送到目的地而无需跨轨道。这可以实现消息聚合和网络流量优化。



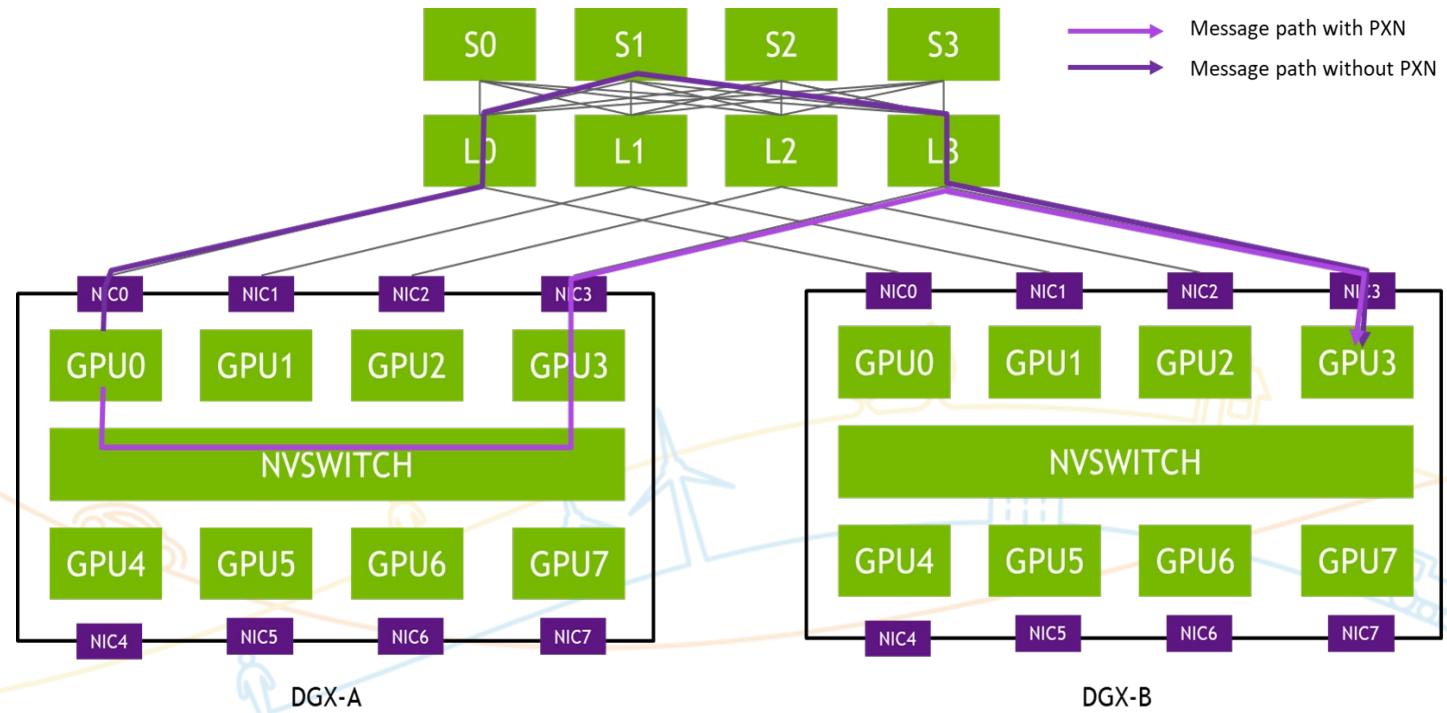
多轨通信



Question

I. 多轨通信的好处在哪里？

- 在同一对 NIC 之间传递的数据被聚合，最大限度地提高有效数据传输速率和网络带宽。



03 数据聚合

Message Aggregation



聚合操作

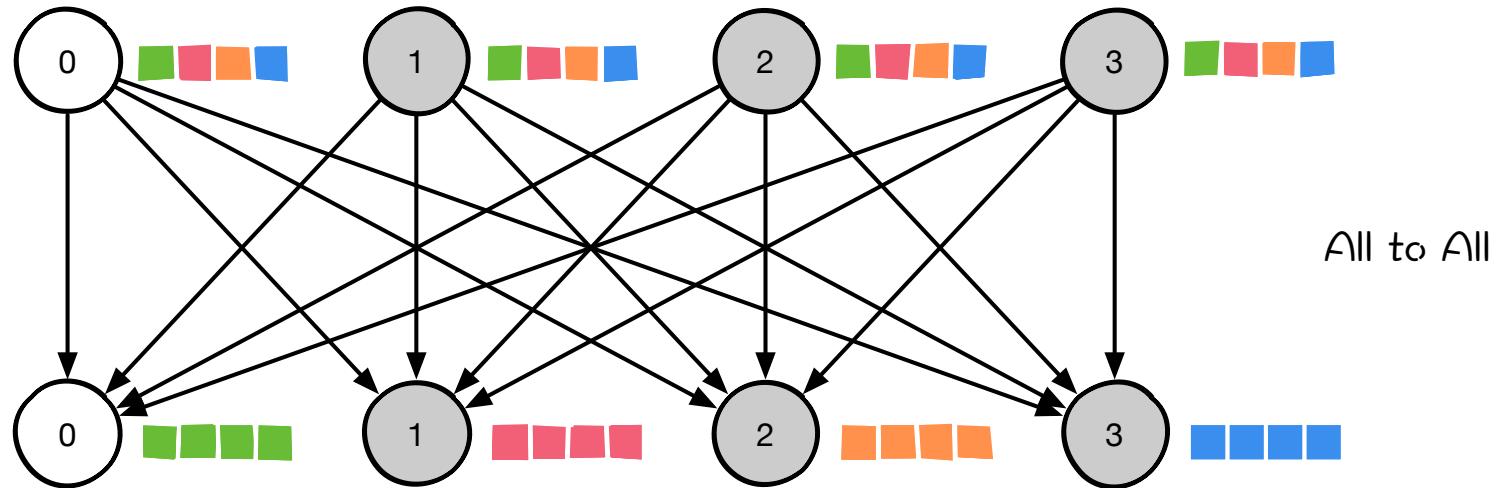
- Aggregation：给定节点上所有 GPU 数据，移动到目的单个 GPU。
- PXN 作用：使网络层能进行消息聚合，建立的网络连接更少，提升路由效率，减少 GPU 负载。



All2All

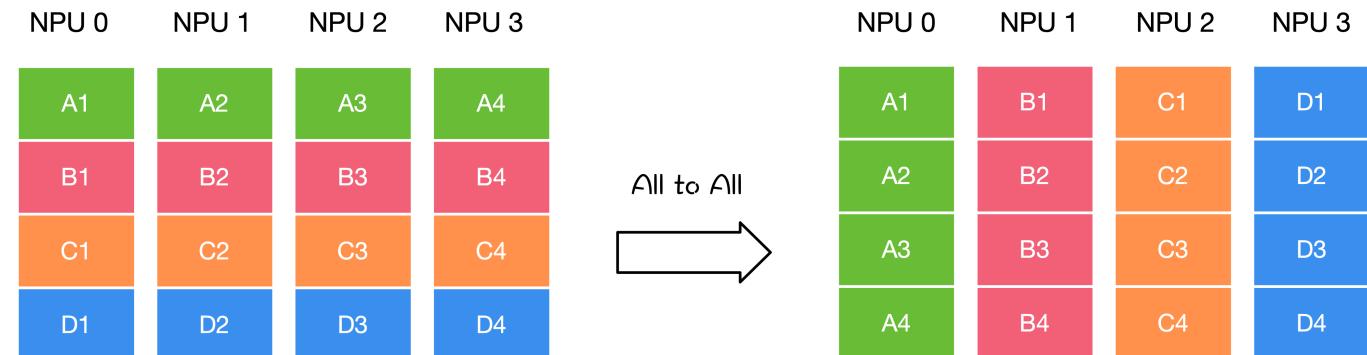
- 算法放方式：

- All-Gather 扩展，不同节点向某一节点收集到数据是不同的
- 每个进程与其他所有进程进行通信



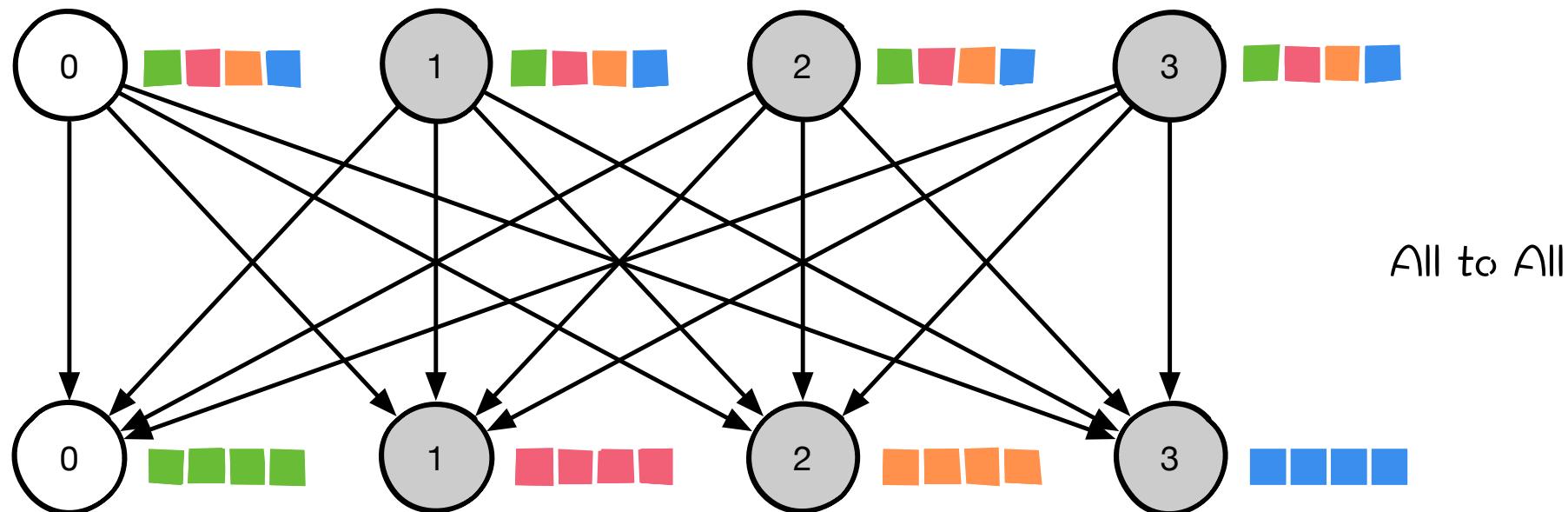
- 主要应用场景：

- 应用于模型并行(TP/SP/EP)
- 模型并行里的矩阵转置

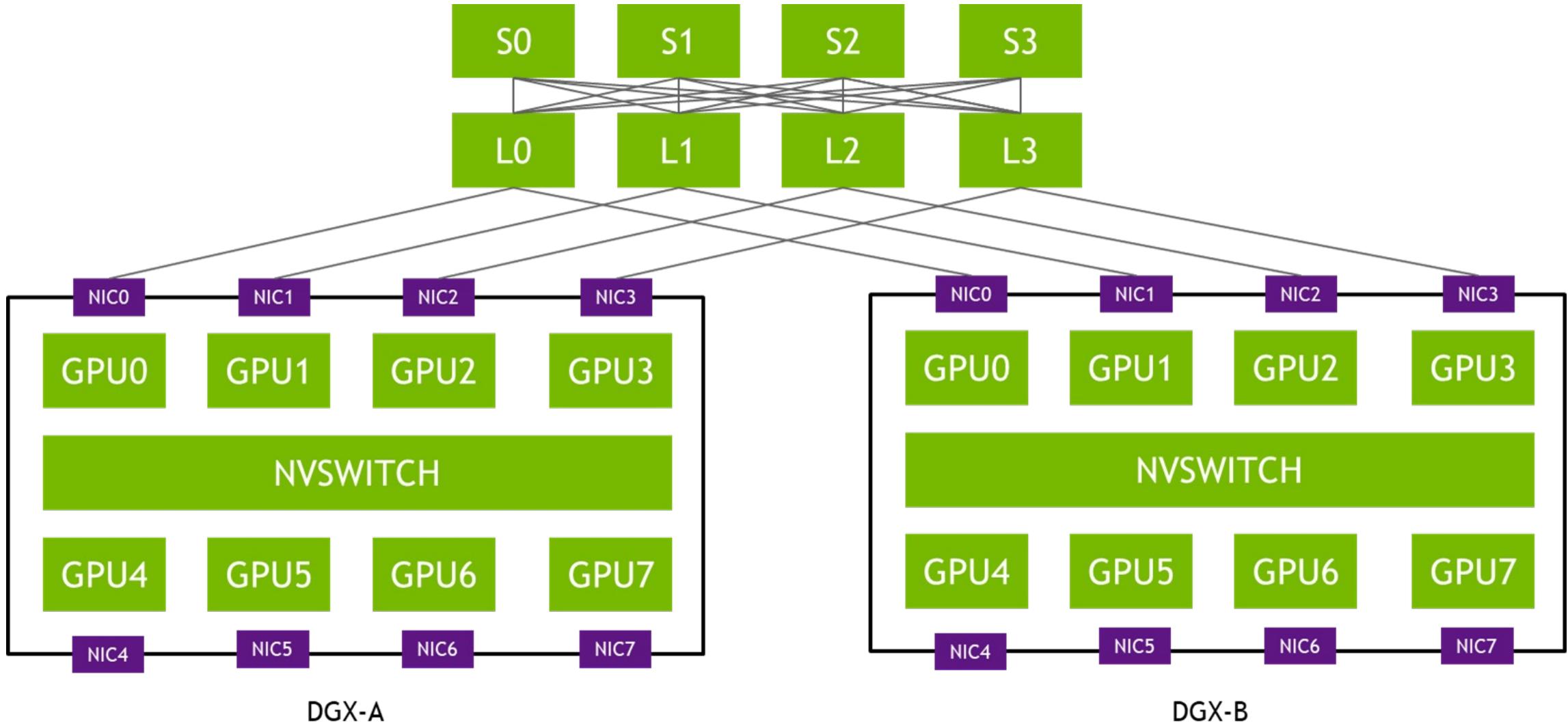


All2All 操作

- 如果节点上 GPU 执行 All2All 操作，要从远程节点的八个 GPU 接收数据：
 - 接收方：NCCL 会通过多个接收器，调用八个缓冲区的数据；
 - 发送方：网络层可以等到所有八个发送都准备就绪，然后一次发送 8 条消息。



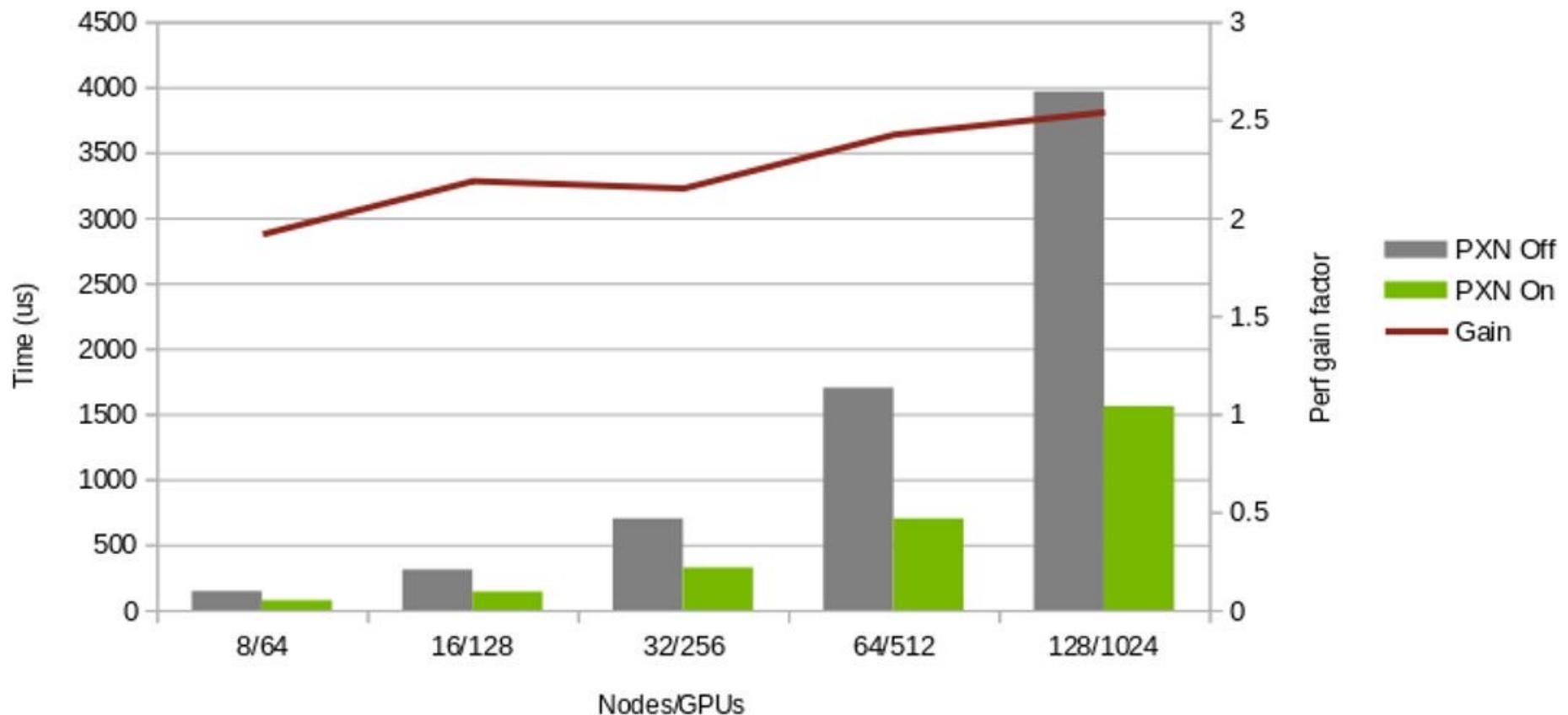
All2All 操作



PXN All2All

NCCL Alltoall latency

DGX A100, Infiniband, NCCL 2.12





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.



ZOMI

Course chenzomi12.github.io

GitHub github.com/chenzomi12/AIFoundation

参考资料

- http://www.njwdr.com/newsdetail_3252489.html

