# Research summary – Junyeol Ryu
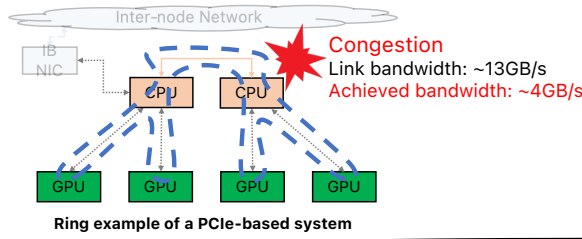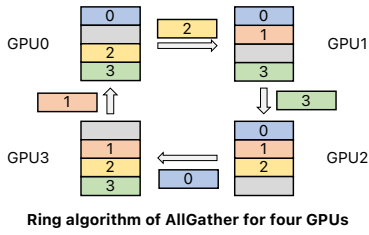
jyeol.ryu@gmail.com
https://junyeol.me/about.html

## Congestion avoidance for collective communication in PCIe-based systems

**Collective communication** is essential for parallelism in training AI models!

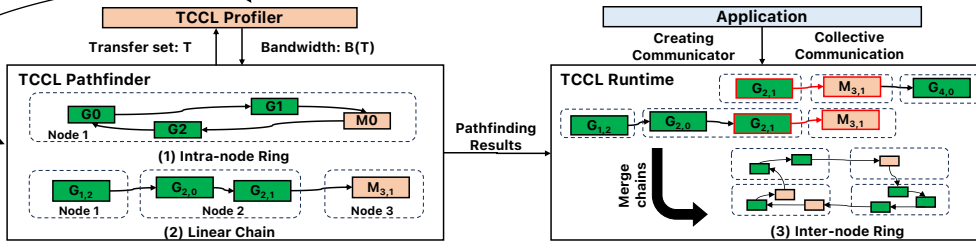GPU communication libraries exibit **low performance for PCIe-based systems**

Existing libraries can neither identify nor avoid congestion!


Ring algorithm of AllGather for four GPUs


Ring example of a PCIe-based system

Congestion
Link bandwidth: ~13GB/s
Achieved bandwidth: ~4GB/s

Key findings from my analysis of NCCL
- Existing libraries find paths based **solely on the bandwidths of individual links**
- However, multiple transfers are executed **simultaneously across the PCIe host bridge** during collective communication

Insight 1: Profiler specialized in measuring simultaneous multiple transfers



Insight 2: Enumerate all possible paths while minimizing the search time for performant path

Overview of TCCL

Result:
- Up to $2.07\times$ speedup for collective communication
- Up to $1.11\times$ speedup for training AI models

Further research directions
- Extending beyond ring algorithm (e.g., double binary tree, all-to-all)
- Overlapping dependent communication and computation by decomposition
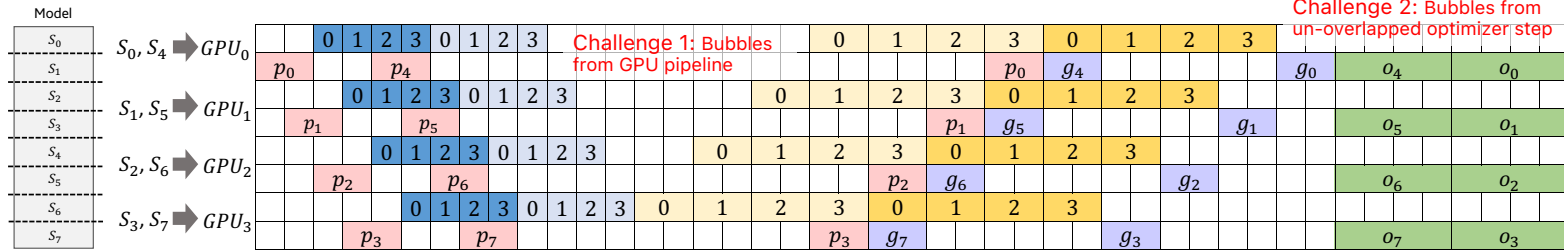- Utilizing multi-path opportunities

## Utilizing compute and memory of GPUs and CPUs for large-model training

**Parallelism + offloading** is essential for training under memory pressure!

$p_i$: Stage $S_i$'s prefetch of parameters
$g_i$: Stage $S_i$'s offload of gradients
$o_i$: Stage $S_i$'s CPU optimization steps

Legend: Bubble, Forward pass, Backward pass, Parameters transfer (CPU to GPU), Gradients transfer (GPU to CPU), Optimization step (CPU)



Challenge 1: Bubbles from GPU pipeline

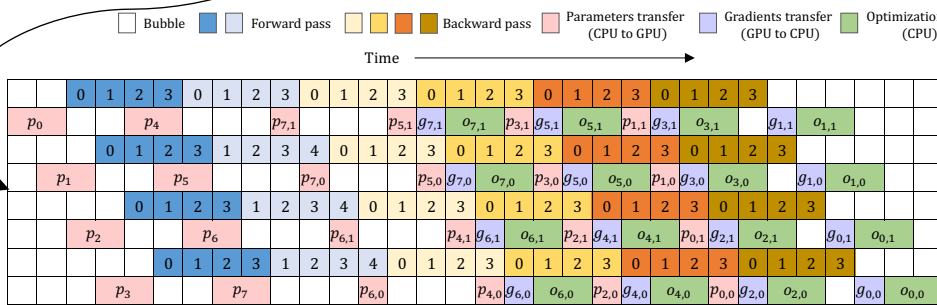Challenge 2: Bubbles from un-overlapped optimizer step

- Stage's memory should not exceed GPU memory
- Multiple stages are assigned to each GPU

Stages are fetched to GPU memory and gradients are offloaded to CPU memory in an overlapped manner

Optimizer updates the stages' parameters on the CPU

Insight 1: Shared parameters on CPU memory allows **decoupling** a stage's forward/backward pass onto different GPUs



Insight 2: CPU optimizer steps can execute **in parallel** with GPU's forward/backward pass

Result: Avg. $1.26\times$ speedup for training LLaMA2 models (~100B) using ~32 V100 32GB GPUs

Further research directions: Hybrid parallelism (e.g., PP+TP) + offloading technique

## Reinforcement learning-based resource management

**Mitigating network congestion** is essential when scheduling distributed jobs in GPU clusters!


Co-location of distributed jobs

Insight: Co-locating jobs yields **varying performance effects** due to model type, parallelism, placement

However, it is infeasible to try all co-location options on every new job request

Reinforcement learning (RL)
- Repetitive decisions leave abundant training data to RL algorithm
- Reward reflects complex objectives (e.g., min. congestion, max GPU utilization)
- Adapt to shifting or unseen circumstances by **explore-and-exploit**


Fixed state design as input to NN-based RL algorithm

- Penalize increase in congestion
- Incentivize increase in GPU utilization

⊕ Insight: **Simple heuristics** can effectively assist RL (e.g., selective multiplexing with greedy approach)

Result: Up to $18.2\%$ reduction for average job completion time

- Schedule
- Migration
- Preemption