



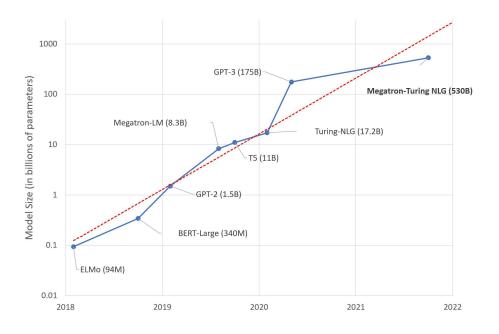
When Light Bends to the Collective Will A Theory and Vision for Adaptive Photonic Scale-UP Domains

Vamsi Addanki

https://stygianet.cs.purdue.edu

ACM HotNets 18 November 2025

Exponential Growth in Model Sizes







Chip-to-Chip Communication is Essential

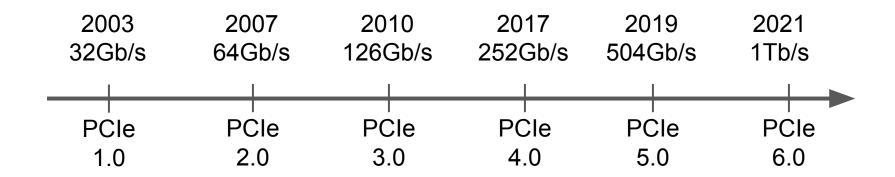


e.g., NVIDIA DGX server





Chip-to-Chip Communication is a Bottleneck!







Chip-to-Chip Communication is a Bottleneck!

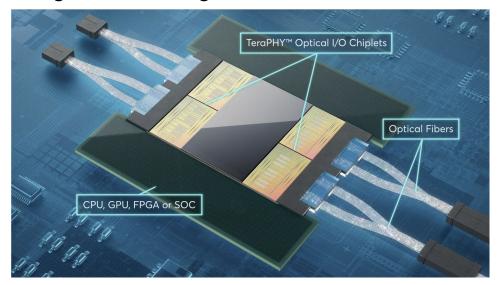
Memory bandwidth ~doubles every two years

But bandwidth requirement is growing much faster



Chip-to-Chip Photonic Interconnects are on the Rise

Ayar Labs: "Moving Data with Light!"



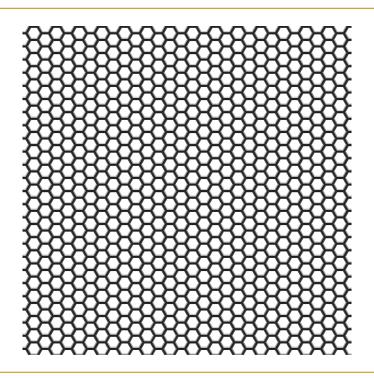














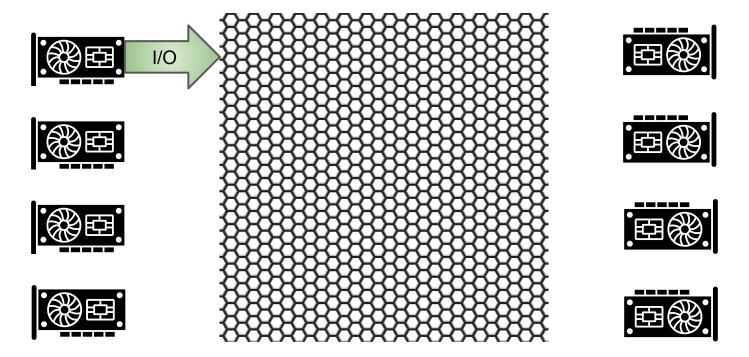






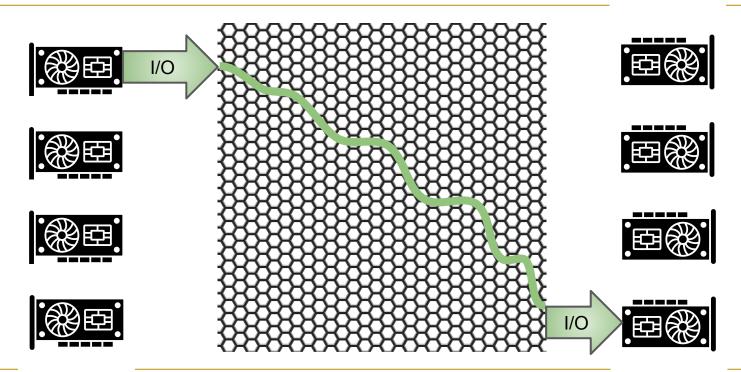






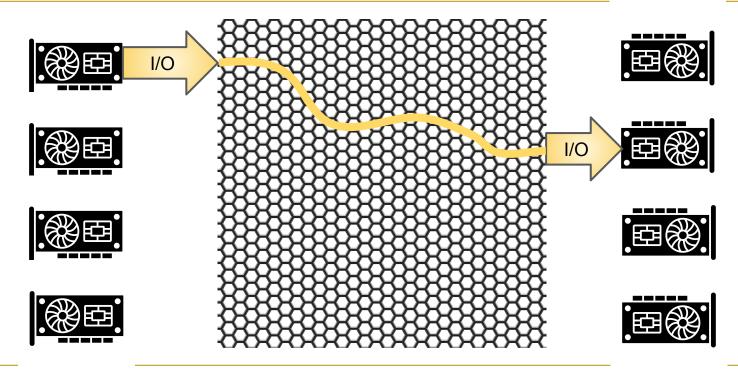






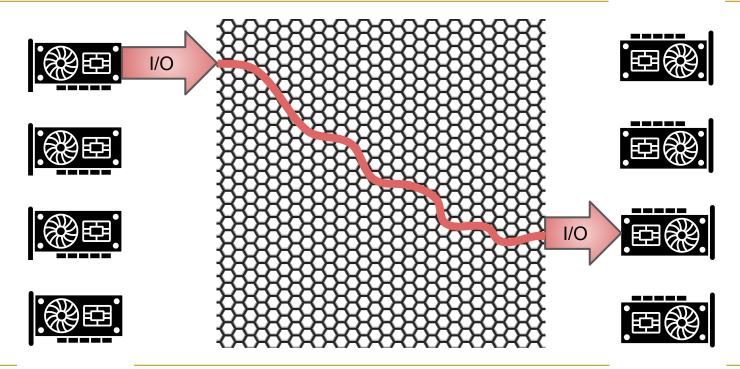
















- Circuit-switched network
- Bufferless
- Reconfiguration delays





- Reconfigurable networks are well studied in the recent past
 - Demand-aware reconfigurable networks
 - BvN decompositions
 - Greedy matchings
 - Oblivious reconfigurable networks
 - Periodic circuit switching





- Reconfigurable networks are well studied in the recent past
 - Demand-aware reconfigurable networks -
 - BvN decompositions
 - Greedy matchings
 - Oblivious reconfigurable networks
 - Periodic circuit switching

Assumption:

Reconfiguration delay is negligible



- Reconfigurable networks are well studied in the recent past
 - Demand-aware reconfigurable networks -
 - BvN decompositions
 - Greedy matchings
 - Oblivious reconfigurable networks
 - Periodic circuit switching
 - One-shot optimization
 - Failure mitigation



Reconfiguration delay is negligible

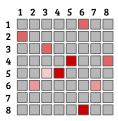
Assumption:

Reconfiguration delay is too high



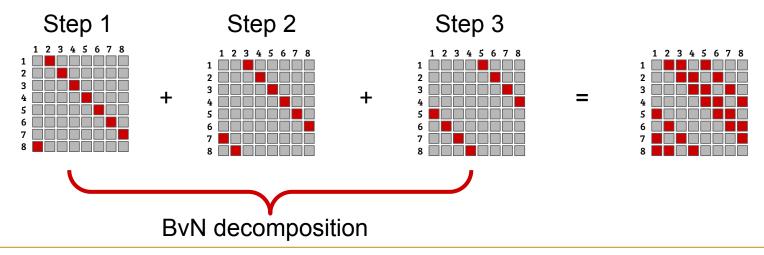


- Reality is much different!
 - Different range of reconfiguration delays across technologies, vendors...
 - Anywhere between 10s of nanoseconds to 100s of milliseconds
 - GPU communication has data dependencies
 - Aggregate demand matrix is not a good abstraction!





- Each step of collective communication is a matching!
 - The sequence of matchings form a BvN decomposition of the aggregate demand matrix

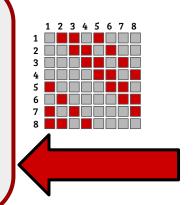






- Each step of collective communication is a matching!
 - The sequence of matchings form a BvN decomposition of the aggregate demand matrix

Decomposing the aggregate matrix may not align with the collective





a: Time to prepare the chunk, a setup latency

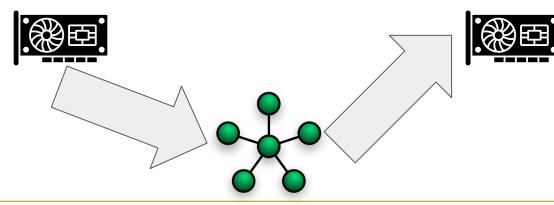
b: Link propagation delay

c: Time to transmit one bit

d: Congestion factor

Completion time

setup + (path length) propagation + (congestion) transmission delay





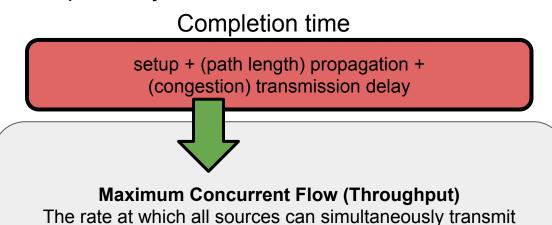


a: Time to prepare the chunk, a setup latency

b: Link propagation delay

c: Time to transmit one bit

d: Congestion factor



without exceeding link capacities in the network





 α : Time to prepare the chunk, a setup latency

δ: Link propagation delay

 β : Time to transmit one bit

θ: Congestion factor





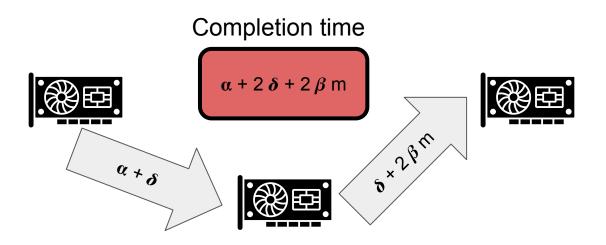


 α : Time to prepare the chunk, a setup latency

δ: Link propagation delay

 β : Time to transmit one bit

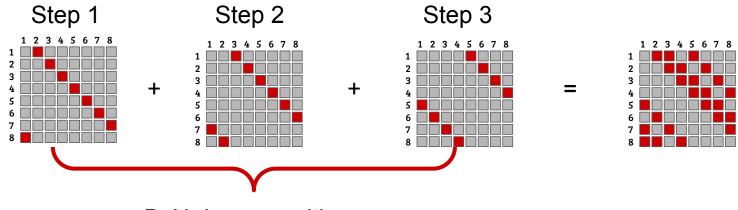
θ: Congestion factor







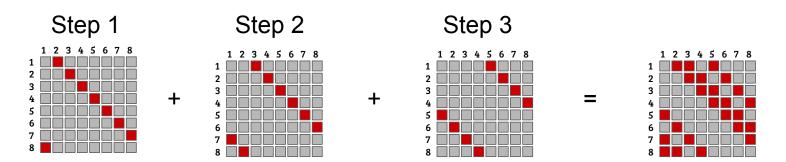
BvN, Concurrent Flow, and α - β Cost Model



BvN decomposition



BvN, Concurrent Flow, and α - β Cost Model



Completion time of each step of the collective

setup + (path length) propagation + (congestion) transmission delay

 $\alpha + 1 \delta + \beta \text{ m/}\theta$

 α - β is not merely a cost model.

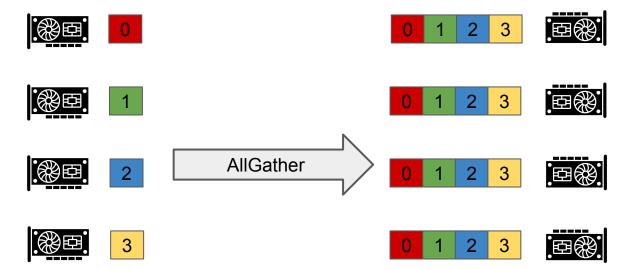
It should be seen as a *completion time* model.





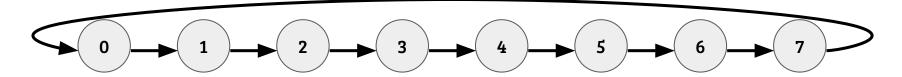
AllGather Operation

Every GPU transmits distinct chunks of its data to all other GPUs



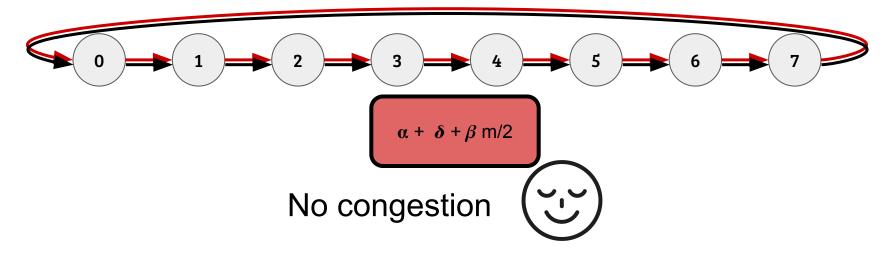


• 8 GPUs connected in a ring topology

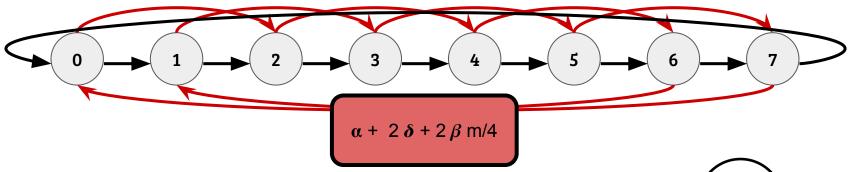




• Step 3



• Step 2

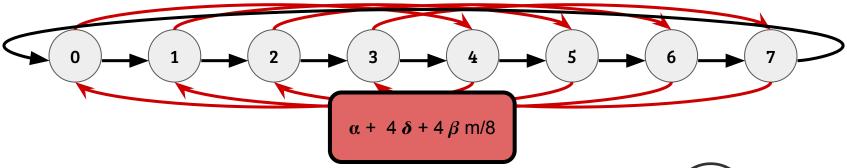


Moderate congestion→2 overlapping transfers





Step 1



Severe congestion→4 overlapping transfers



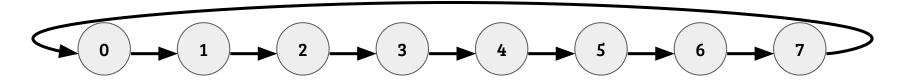


What if we can change the topology?





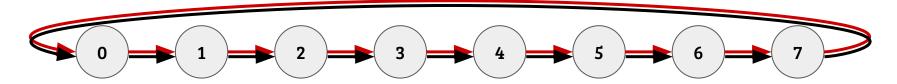
• 8 GPUs connected via *reconfigurable* photonic interconnect

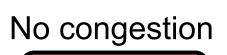


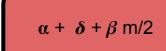
Topology can be changed



• Step 3





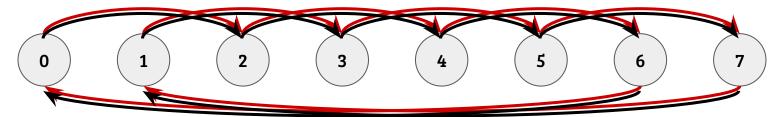






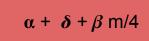


• Step 2



Reconfigure!
No congestion

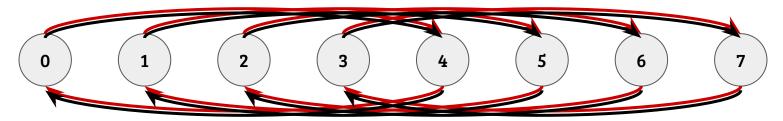




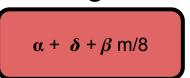




• Step 1



Reconfigure!
No congestion









Reconfigurable Photonic Interconnect



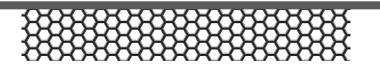




It is not free lunch unfortunately, there is a catch!

Reconfiguration adds latency











An Optimization Opportunity

- Tradeoff between congestion, propagation delay and reconfiguration delay
- Central question: When and how to reconfigure?



• Decision variable: Reconfigure or not

 $\alpha + \alpha_r + \Box + m \beta$: Reconfigure and align topology with communication

 $\alpha + \mathbf{I} \Box + \mathbf{m} \beta/\mathbf{\theta}$: Mismatch

min
$$s \cdot \alpha + \delta \cdot \sum_{i=1}^{s} \left(\begin{array}{c} \text{propagation delay} & \text{direct} \\ w/\text{o reconf.} & \text{with reconf.} \end{array} \right) + \sum_{i=1}^{s} \left(1 - z_i \right) \cdot \alpha_r$$

$$+ \beta \cdot \sum_{i=1}^{s} m_i \cdot \left(x_i \cdot \frac{1}{\theta(G, \mathcal{M}_i)} + \underbrace{(1 - x_i) \cdot 1}_{\text{no congestion}} \right)$$

$$\text{congestion} \quad \text{with reconf.}$$

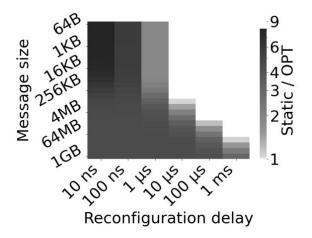
$$\text{w/o reconf.}$$



 Turns out the problem admits efficient dynamic programming solution for some collective communication algorithms, while retaining optimality



- Turns out the problem admits efficient dynamic programming solution for some collective communication algorithms, while retaining optimality
- The payoff is significant!







Open Questions

- What is the optimal schedule for All-to-All, Tree algorithms..?
- How should routing algorithms adapt to topology changes?
- How to implement custom routing algorithms in hardware?
- Can reconfigurations be overlapped with computations?
 - Potentially zero reconfiguration overhead!
- The reconfiguration delay is not just optical rise/fall e.g., 10ns-10us
 - GPU cuda kernels, PCIe latency, error correction, negotiation, thermal, electrical.... all contribute to the latency profile!
- Many more....



Thank You

Vamsi Addanki vaddank@purdue.edu

DT

**Talina **Tali



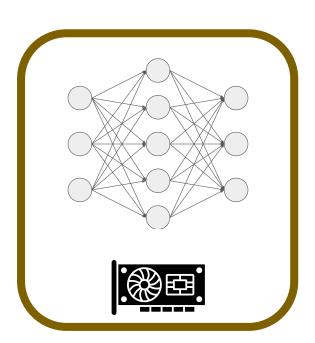


Backup Slides



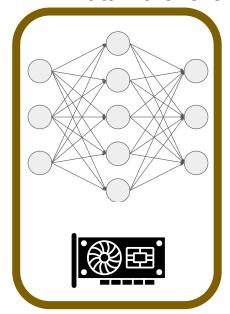


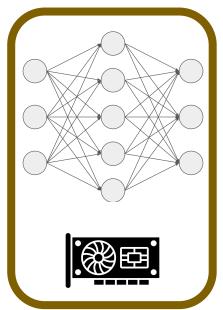
Single GPU

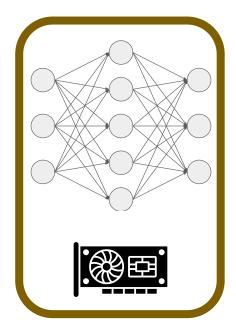


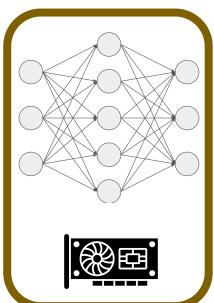


Data Parallelism







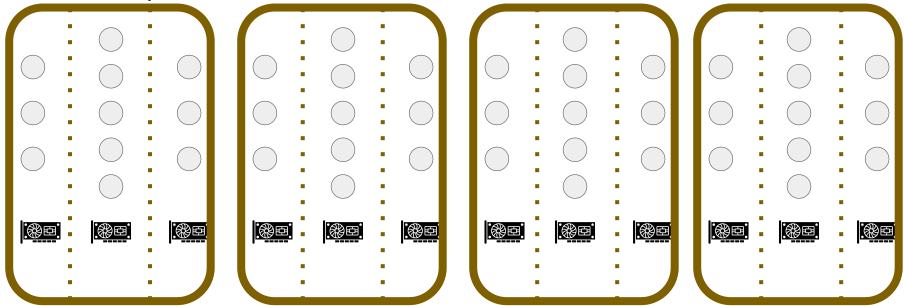


Data Parallelism **AllReduce**



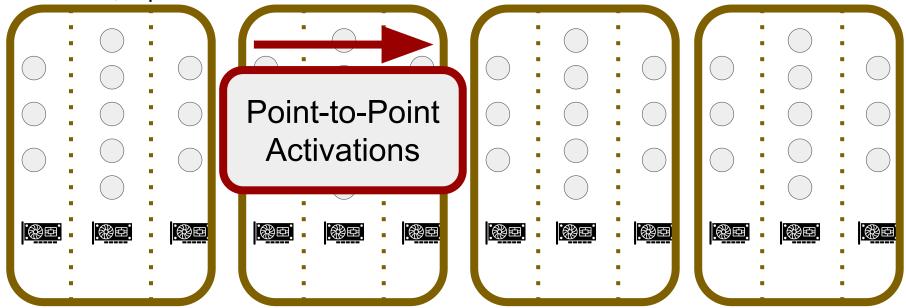


• Data, Pipeline Parallelism



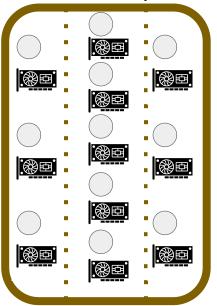


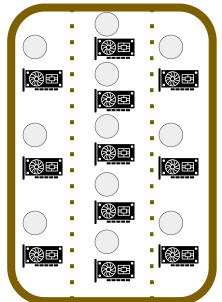
Data, Pipeline Parallelism

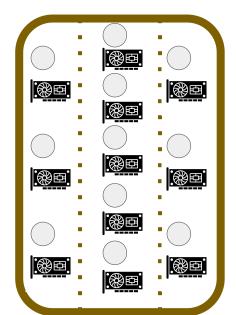


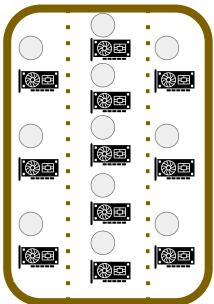


Data, Pipeline, Tensor Parallelism



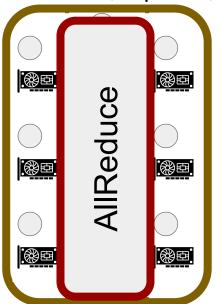


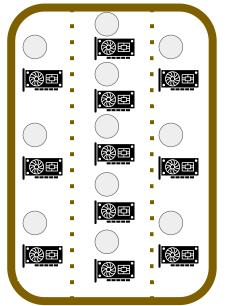


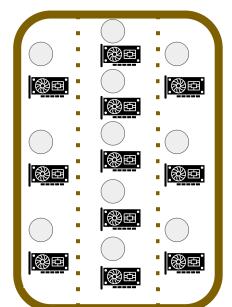


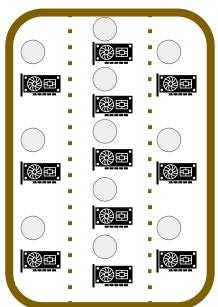


Data, Pipeline, Tensor Parallelism













Data, Pipeline, Tensor Parallelism **黎**国 **AIIReduce** Point-to-Point **Activations 89**E **AllReduce**



