



Exploiting Software-Defined Networking Technology for Improving UGAL Routing in Dragonfly Networks

Ram Sharan Chaulagain, Tusher Chandra Mondol, Saptarshi Bhowmik, Xin Yuan

April 8, 2025

Accepted in **CCGRID'25**

Contents

① Background/Motivation

② Research Problem

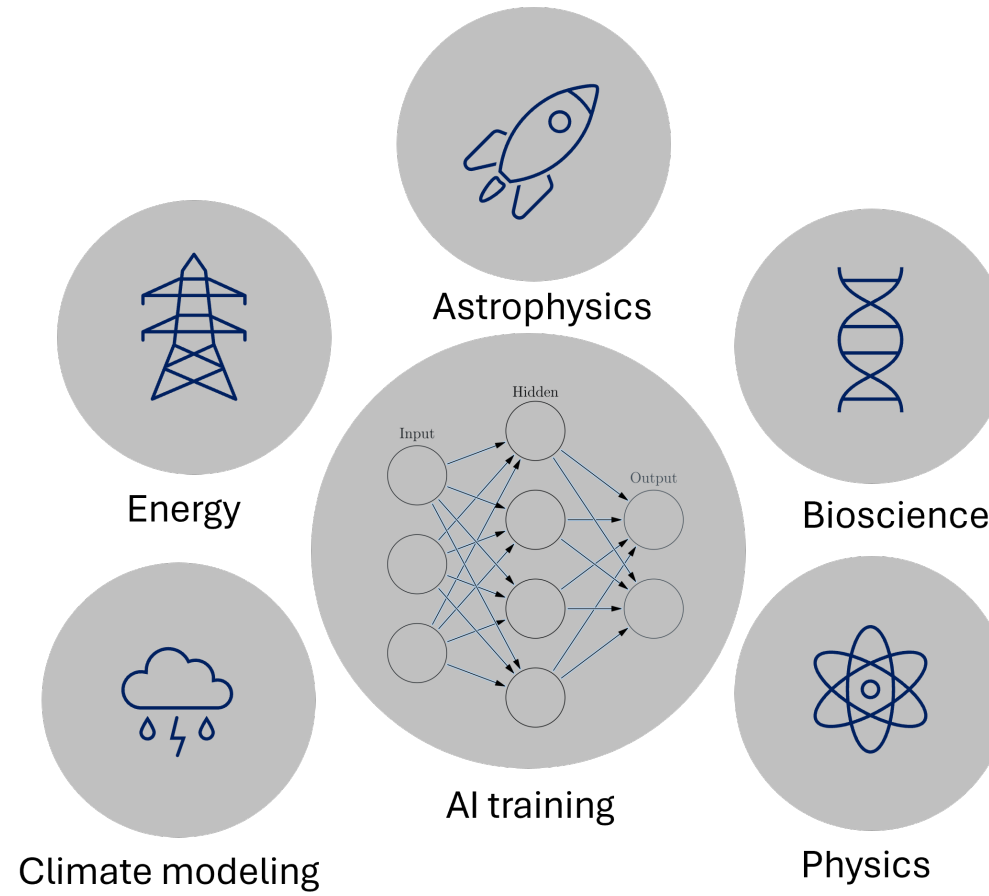
③ Proposed Methods

④ Evaluation and Results

⑤ Conclusion

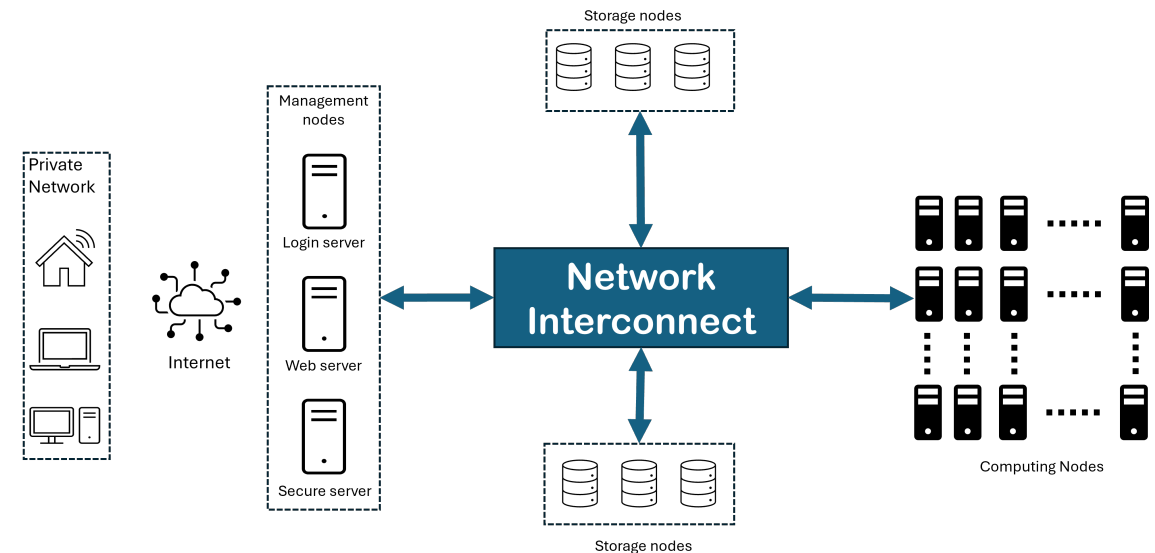
High Performance Computing System

HPC is about **solving complex problems** faster by combining powerful hardware and software.



The Backbone of HPC & AI Systems

- Multiple processing nodes are interconnected to achieve large-scale HPC.
- Performance relies heavily on **high-performance interconnects**.
- As compute scales up (thousands of GPUs/CPU), communication becomes a bottleneck.
- *Network interconnect impacts scalability, latency, cost, and energy efficiency.*



Network Interconnects for HPC systems

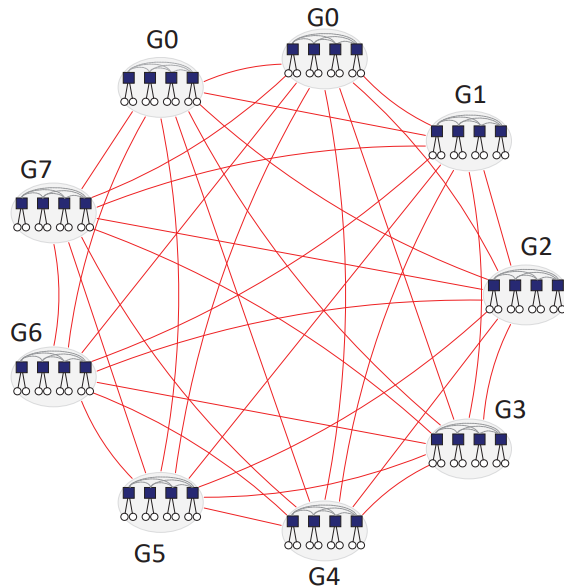


Figure 1: An example of **direct network**: Dragonfly

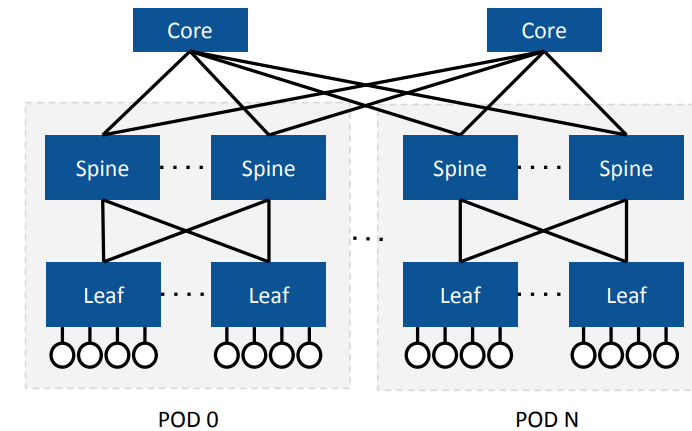


Figure 2: An example of **indirect network**: Fat-Tree

Topology	Max Servers	Switches Needed	Diameter	Real-world usage
Dragonfly	~262,656	~16,416	3	Modern HPC and AI systems
Fat-Tree	~65,536	~5,120	4	Conventional DC

Table 1: Comparison of Fat-Tree and Dragonfly topologies with K=64

Current Exascale Interconnect

- **El Capitan** is the **world's fastest supercomputer** as of November 2024.
- Three systems have achieved Exascale performance till now.
- 7/10 supercomputers from top 10 list of top 500 supercomputers uses dragonfly topology.



Figure 3: El Capitan Supercomputer

1.742 exaFLOPS

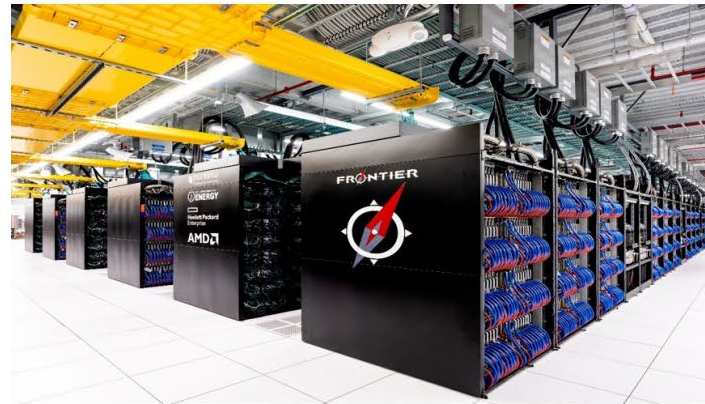


Figure 4: Frontier Supercomputer

1.135 exaFLOPS



Figure 5: Aurora Supercomputer

1.012 exaFLOPS

Dragonfly Network

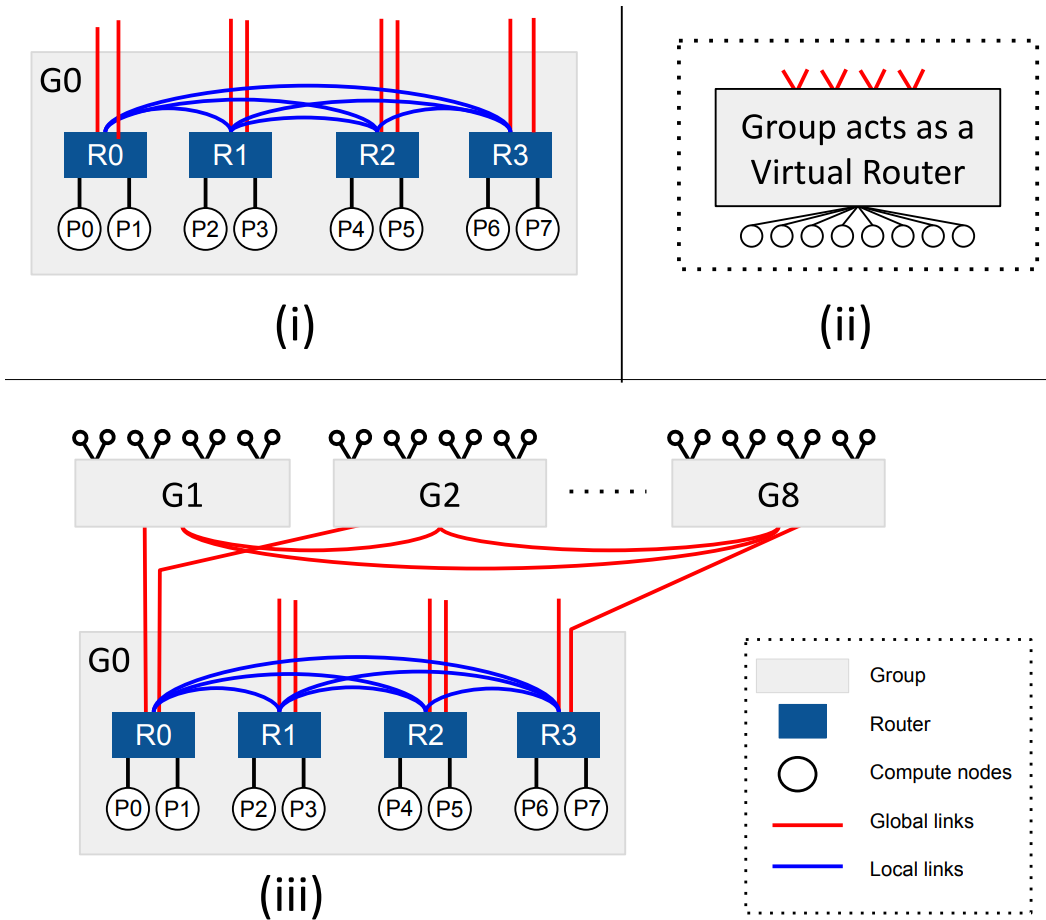


Figure 6: Dragonfly topology with $a=2$, $p=4$, $h=2$ and $g=9$ and denoted as $dfl y(a, p, h, g)$

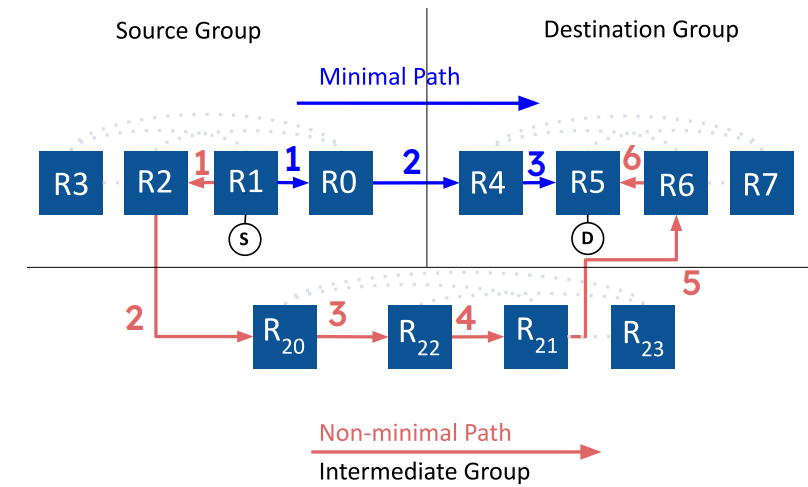


Figure 7: Minimal and Non-minimal paths

UGAL-L

UGAL with local information (UGAL-L):

- State-of-the-art routing scheme for dragonfly topology
- Uses local information of source router
- Queue occupancy estimates the link delay.
- Multiplies local link delay with path hop count to infer path delay.

$$Queue_{min} \times Hop_{min} \leq Queue_{nonmin} \times Hop_{nonmin} + Bias$$

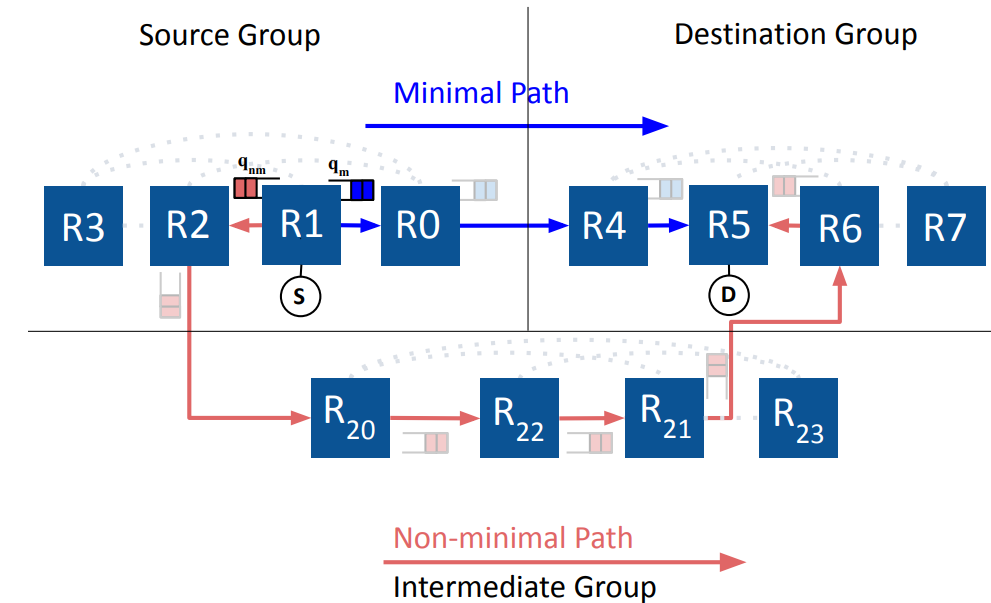


Figure 8: UGAL-L

Limitations in UGAL

Inaccurate Latency Estimation in UGAL-L:

- UGAL-L estimates path latency using local queue occupancy and hop count.
- Assumes uniform network traffic; becomes inaccurate under uneven load.
- Fails to account for congestion beyond the local router.
- Paths with congested links not directly connected to the source are under-estimated as low-latency.

Inefficiency with non-minimal path selection:

- UGAL chooses non-minimal paths randomly without considering congestion.
- May select congested non-minimal paths, worsening load imbalance.
- Misses the opportunity to utilize un-congested non-minimal paths effectively.

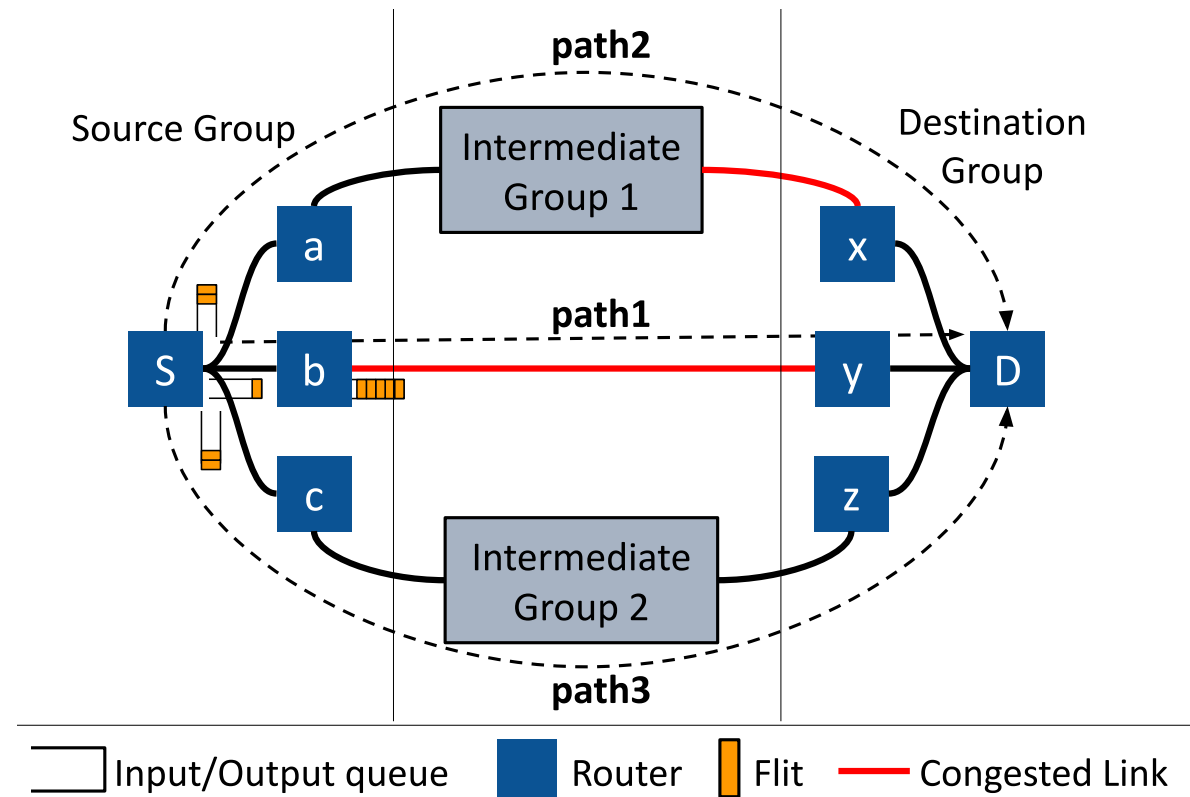


Figure 9: Limitations in UGAL-L

Contention factor

- We introduce a contention term that reflects the effect of link congestion with imbalanced traffic.
- The contention factor for a link is the number of active flows using the link, assuming minimal paths are used for the flows, shown in figure 10.
- *How can a router determine the contention factor of each link in the network when it only has a local view?*

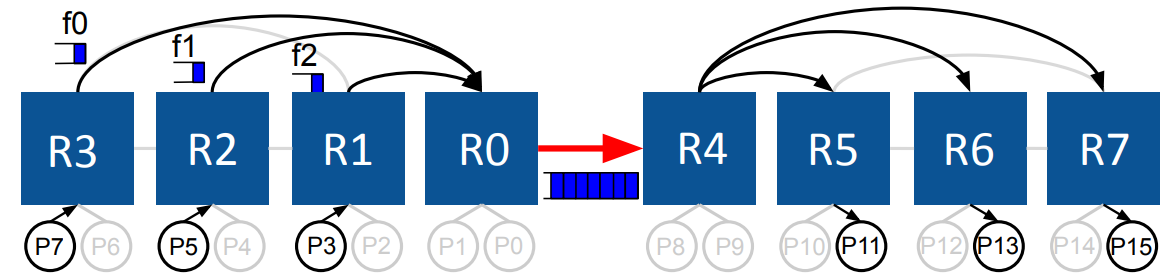


Figure 10: Limitations in UGAL-L

Software Defined Networking (SDN)

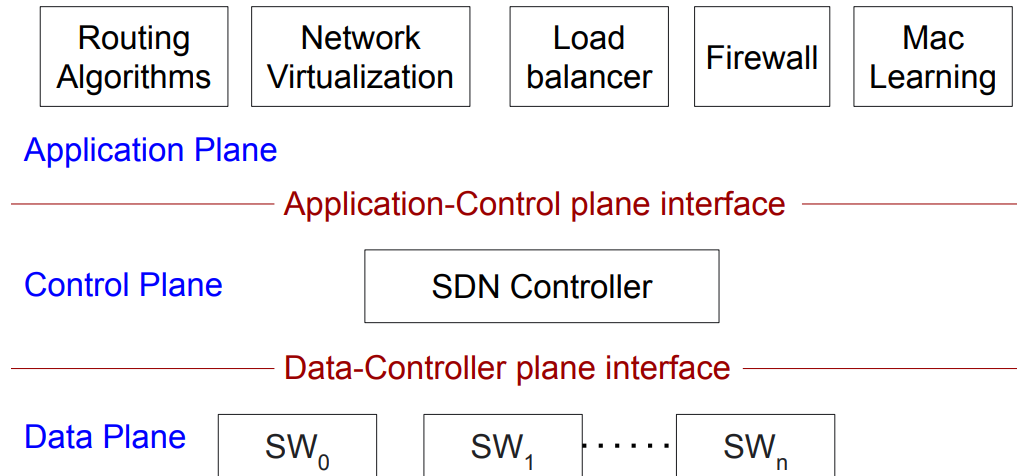


Figure 11: SDN Architecture

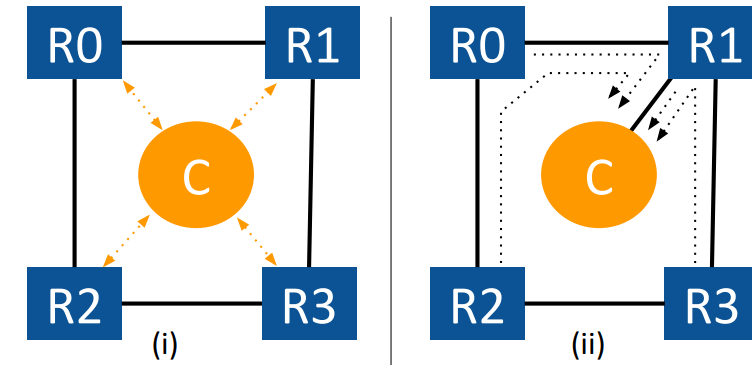


Figure 12: SDN Controller: (i) Out-band and (ii) In-band

Proposed Routing Scheme 1: SDN-UGAL-L

- We redesign the latency estimation formula to accurately predict latency for minimal and non-minimal paths within the UGAL scheme.
- Unlike UGAL-L, we consider the contention information during latency estimation.

Aspect	UGAL-L	SDN-UGAL-L
Queue Occupancy	q_m, q_{nm}	q_m, q_{nm}
Hop Count	H_m, H_{nm}	H_m, H_{nm}
Contention Awareness	None	$C_m[i], C_{nm}[i]$ (per-hop contention factors)
Adjustment Term	None	$A_m = q_m \cdot \sum_{i=1}^{H_m-1} \max(C_m[i] - 1, 0)$ $A_{nm} = q_{nm} \cdot \sum_{i=1}^{H_{nm}-1} \max(C_{nm}[i] - 1, 0)$
Minimal Path Latency	$q_m \cdot H_m + b$	$q_m \cdot H_m + A_m + b$
Non-minimal Path Latency	$q_{nm} \cdot H_{nm}$	$q_{nm} \cdot H_{nm} + A_{nm}$

Proposed Routing Scheme 2: SDN-UGAL-L+

- The contention factor provides a measure of relative congestion among global links.
- **SDN-UGAL-L+** uses this information to decide whether to **include or exclude non-minimal paths** through a particular intermediate group.
- If the contention factor of both global links is two or fewer, the group is considered as a candidate; otherwise, it is excluded.
- If the total number of candidate groups is below a threshold (e.g. 2), then all excluded groups are put back in the candidate groups ; to ensure sufficient number of path diversity.

Evaluation

- Evaluated using synthetic traffic patterns and application workloads.
- In synthetic traffic patterns all flows are considered as elephant flows.
- In application workloads, flows > 10MB within a 3s pooling period are classified as elephant flows.
- Compared with UGAL-L, UGAL-LE, UGAL-G and PAR.

Topology	Num of nodes	No of routers	No of groups	Links per group pair
<i>d fly</i> (4, 8, 4, 33)	1056	264	33	1
<i>d fly</i> (4, 8, 4, 17)	544	136	17	2
<i>d fly</i> (4, 8, 4, 9)	288	72	9	4

Table 2: Topologies Used in Experiments

Synthetic Traffic Patterns on dfly(4,8,4,33)

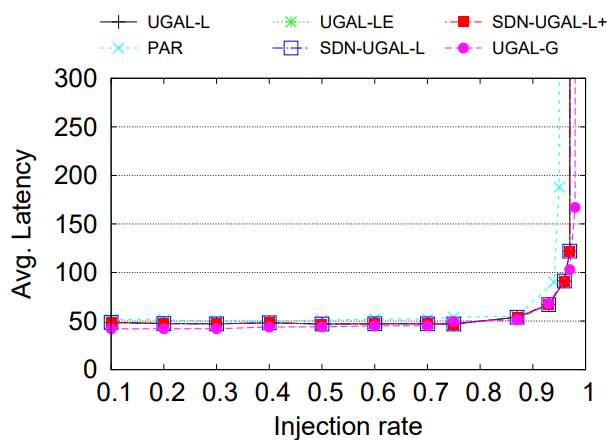


Figure 13: Uniform traffic

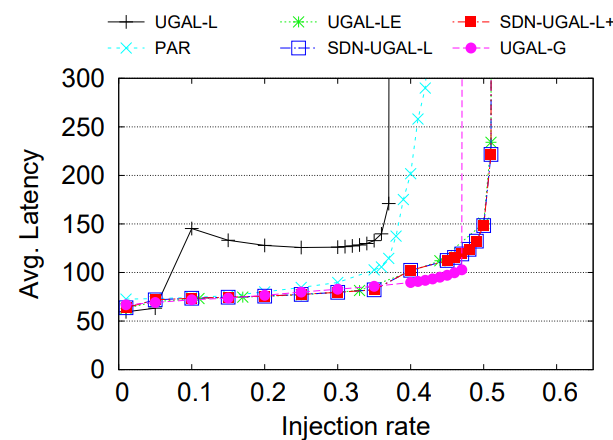


Figure 14: Adversarial shift

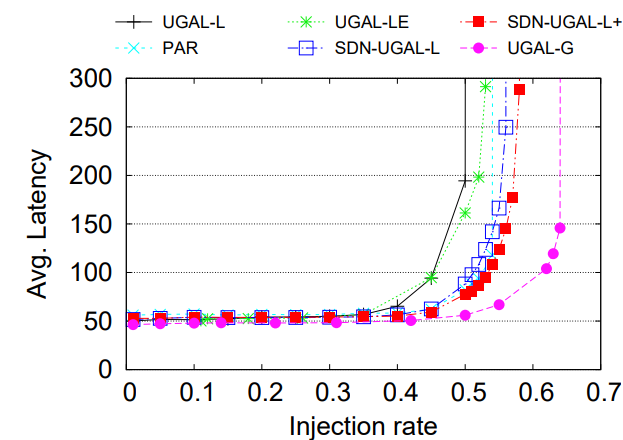


Figure 15: Random permutation

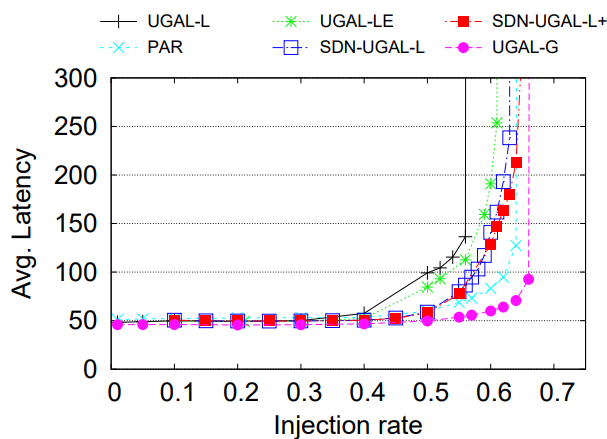


Figure 16: 50% mixed permutation

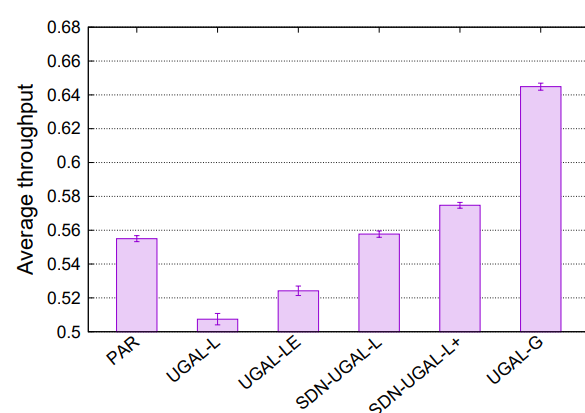


Figure 17: Throughput (36 permutations)

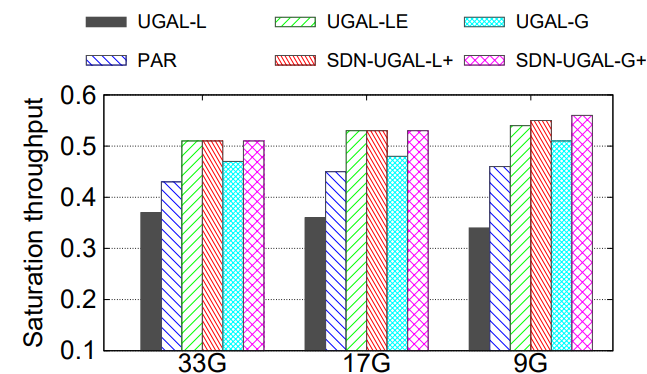


Figure 18: Shift traffic across topologies

Evaluation with HPC Application Workloads

Routing	MILC	LAMMPS	Nekbone	Stencil
PAR	64.49	51.15	112.26	19.73
UGAL-L	62.58	50.04	130.13	18.44
UGAL-LE	51.94	44.96	121.46	14.69
SDN-UGAL-L	51.01	47.13	120.63	14.36
SDN-UGAL-L+	50.92	46.31	117.42	14.14

Table 3: Average communication time under **continuous allocation** with Out-band SDN Controller

Routing	MILC	LAMMPS	Nekbone	Stencil
PAR	63.99	51.92	118.55	21.05
UGAL-L	63.41	48.21	128.05	20.31
UGAL-LE	62.63	50.30	123.12	20.29
SDN-UGAL-L	60.12	47.52	122.03	19.68
SDN-UGAL-L+	59.22	47.23	120.32	19.24

Table 4: Average communication time under **random allocation** with Out-band SDN Controller

App	Alloc	Out-band	In-band	Ideal
MILC	Cont	50.92	51.10	49.99
	Rand	59.22	59.31	59.07
LAMMPS	Cont	46.31	47.41	46.52
	Rand	47.23	47.65	47.13
Nekbone	Cont	117.42	117.11	117.48
	Rand	120.32	120.08	120.35
Stencil	Cont	14.14	14.67	14.05
	Rand	19.24	19.53	19.14

Table 5: SDN-UGAL-L+ implementation variants under different allocation strategies and SDN-Controller types

Conclusion

- SDN offers significant potential for enhancing adaptive routing performance in Dragonfly systems.
- We developed an efficient adaptive routing scheme using SDN for Dragonfly interconnects.
- Evaluation results demonstrate that incorporating SDN information into UGAL-based routing can achieve substantial performance improvements in Dragonfly systems.

Thank you for your attention!

Q&A