

Evaluating Quality of Service Traffic Classes on the Megafly Network

Misbah Mubarak*, Neil McGlohon[†], Malek Musleh[‡], Eric Borch[‡], Robert B. Ross*, Ram Huggahalli[‡],
Sudheer Chunduri[§], Scott Parker[§], Christopher D. Carothers[†], Kalyan Kumaran[§]

*Mathematics and Computer Science Division, Argonne National Laboratory, Lemont, IL, USA

[†]Department of Computer Science, Rensselaer Polytechnic Institute, Troy, NY, USA

[‡]Intel Corporation, Santa Clara, CA, USA

[§]Argonne Leadership Computing Facility (ALCF), Argonne National Laboratory, Lemont, IL, USA

Abstract—An emerging trend in High Performance Computing (HPC) systems that use hierarchical topologies (such as dragonfly) is that the applications are increasingly exhibiting high run-to-run performance variability. This poses a significant challenge for application developers, job schedulers, and system maintainers. One approach to address the performance variability is to use newly proposed network topologies such as megafly (or dragonfly+) that offer increased path diversity over a traditional fully connected dragonfly.

In this work, we select HPC application workloads that have exhibited performance variability on current 2-D dragonfly systems. Using the Scalable Workload Model (SWM) library and CODES HPC systems simulation framework, we evaluate the baseline performance expectations of these workloads on the megafly and 1-D dragonfly network models. Our results show that in the majority of the cases, a megafly network is more resistant to communication interference than a fully connected 1-D dragonfly network. However, in order to sufficiently mitigate the performance variability, additional quality of service (QoS) levels need to be explored. We use bandwidth capping and traffic differentiation to introduce multiple traffic classes in megafly networks. In some cases, our results show that QoS can completely mitigate application performance variability while causing minimal slowdown to the background network traffic.

Index Terms—Megafly Network, Quality of Service, Performance Variability, Dragonfly Network, Parallel Discrete-event Simulation

I. INTRODUCTION

With modern high-performance computing (HPC) systems shifting to hierarchical and low-diameter networks, dragonfly networks have become a popular choice. They have been deployed in multiple high-performance systems including Cori, Trinity, and Theta systems at NERSC, Los Alamos National Laboratory, and Argonne National Laboratory, respectively [1]–[3]. Dragonfly is a hierarchical topology that uses short electrical links to form groups of routers using a 1-D or 2-D all-to-all interconnect. These groups are then connected all-to-all via optical links. While this design offers low diameter and cost, it increases contention for the link bandwidth among multiple applications which introduces performance variability [4]. For next-generation exascale systems, HPC designers are considering variations of the dragonfly topology that offer increased path diversity, fairness, and

scalability [5]. One such topology that has been recently proposed is the dragonfly+, or megafly, which uses a two-level fat tree to form groups of routers. These groups are then connected all-to-all via optical links. Megafly networks use the path diversity of a two-level fat tree to alleviate the communication bottlenecks that can be introduced with standard dragonfly networks. Megafly networks also have the added advantage of using only four virtual channels (VCs) for deadlock prevention (including request and response traffic) as opposed to eight virtual channels used in a fully connected 1-D dragonfly network. Prior work [6] has shown that the design of megafly networks helps mitigate performance variability to some extent, it does not completely eliminate it.

Although quality of service (QoS) has been investigated and implemented on TCP/IP networks and data-centers [7], the mechanism remains largely unexplored in the context of HPC networks. In the past, HPC network topologies, such as the torus, had mitigated communication interference by dedicating isolated partitions to individual jobs. Introducing QoS traffic classes can be a useful way to mitigate interference on the now popular hierarchical networks. In this work, we use HPC application workloads that demonstrate performance variability on current dragonfly systems as shown by Chunduri et al. [4]. We replay them on CODES packet-level interconnect simulations [8], [9] to answer questions about dragonfly and megafly network topologies: How does the performance of a megafly network compare with a fully connected dragonfly network? How do traffic classes help with performance variability on a megafly network?

The contributions of this work are as follows. (1) We evaluate the performance variability of HPC application workloads on both a megafly and a 1-D dragonfly network using similar network configurations. We compare the performance of a megafly network with a 1-D dragonfly to determine whether megafly network are better resistant to perturbation. (2) We exploit the fact that megafly requires fewer virtual channels for deadlock prevention (as compared to conventional dragonfly), and we use the unused VCs to introduce QoS traffic classes. We evaluate two mechanisms through which QoS can be introduced in HPC networks. First, using bandwidth

capping and traffic prioritization, we quantify the impact of QoS when an entire high-priority traffic class is dedicated to an application or set of applications. Second, we dedicate the high-priority traffic class to latency-sensitive operations such as MPI collectives and observe the performance improvement. (3) We extend the CODES simulation framework to perform packet-level simulation of HPC networks in an online mode driven by the scalable workload models (SWM) [10] for use in the above-mentioned experiments.

II. EXPLORING QUALITY OF SERVICE ON HPC NETWORKS

In the past, HPC systems were often constructed with torus networks, and jobs were allocated onto partitions of the network that reduced resource sharing and communication interference. With hierarchical networks sharing resources such as switches and links, partitioning becomes more difficult and introducing traffic classes becomes an important step to mitigate communication interference. For example, Figure 1 shows the performance degradation of two HPC communication workloads on a model of Argonne’s Theta Cray XC system, which has a 2-D dragonfly architecture. The performance modeling is done using the CODES simulation framework, and the 2-D dragonfly network simulation in CODES has been validated against Argonne’s theta network using synthetic communication benchmarks [11]. The background traffic injection rates are at a percentage of the maximum link capacity and the application slowdown increases with more background injection. The slowdown in communication time can significantly impact the overall application performance as the typical range of communication time in communication intensive applications is in the range of 50-80% [12]. Performance variability along the same lines has been reported in [4], where the actual system is seeing a slowdown of up to 2x in overall performance.

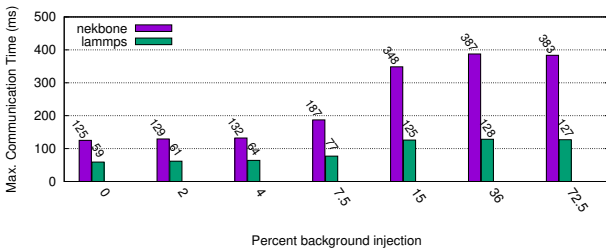


Fig. 1: Communication time slowdown of LAMMPS (1,024 ranks) and Nekbone (1,000 ranks) with background traffic on a Theta network model with 3,456 network nodes. Uniform random workload with large messages is used for background traffic. Background traffic is injected relative to maximum link bandwidth (x-axis shows percentage of link capacity consumed).

Current 2-D dragonfly networks use up to 8 virtual channels to prevent deadlocks, which is typically all the VCs available. Megafly networks use only 4 VCs for deadlock prevention (2 VCs for request traffic and 2 for response traffic), thus making

them a better candidate for enabling multiple traffic classes. In this paper, we explore quality of service on megafly networks based on traffic prioritization and bandwidth shaping [13]. Figure 2 shows one way to implement quality of service on HPC networks. In this implementation, a bandwidth monitor component in each switch tracks the bandwidth consumption of each traffic class for every port. The bandwidth monitoring is done over a static time window t_w . Each traffic class is assigned a certain fraction of maximum available link bandwidth, which serves as the upper bandwidth cap for that traffic class while the link is oversubscribed. If the bandwidth consumption of a traffic class reaches the cap and the link is oversubscribed, the traffic class is designated as inactive for the remaining duration of the static window t_w . An inactive traffic class has the lowest priority and it gets scheduled only if there are no packets in the remaining higher priority traffic classes. At the start of the window t_w , the bandwidth statistics for each traffic class are reset to zero, and the traffic class(es) marked as inactive are activated again. If all the traffic classes are violating their bandwidth cap, then a round-robin scheduling policy is used for arbitration.

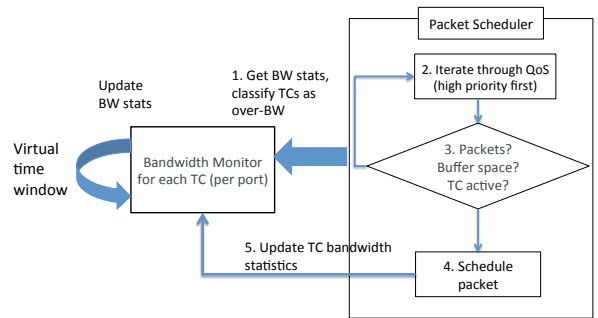


Fig. 2: Enabling quality of service on HPC networks (TC – traffic class, BW – bandwidth, QoS – quality of service)

The implementation of QoS can be beneficial in reducing communication interference on hierarchical networks. For instance, a common problem exhibited on such networks is that communication-intensive (or bandwidth-hungry) applications can “bully” less-communication-intensive applications [9]. With QoS-enabled networks, bandwidth-hungry applications can be prevented from exceeding their permissible bandwidth limits. This approach allows less-communication-intensive applications to have their fair share. Alternatively, one can assign a high-priority traffic class to latency-sensitive operations such as MPI collectives. We report on experiments with both of these QoS mechanisms in Section VI.

III. SIMULATION ENVIRONMENT

In this section, we describe the simulation environment that is used for evaluating interference and QoS on hierarchical networks with realistic HPC workloads. Prior to this work, the CODES simulation framework supported system simulations with post mortem communication traces. Although traces can illustrate realistic system behavior (for a given problem

size), their use inhibits flexibility and simulation scalability as compared to other workload representations. Therefore, we extended the CODES simulation suite to replay workloads in an online or in situ mode using the Scalable Workload Models (SWMs) presented in [10].

A. Scalable Workload Models

Scalable workload models are a workload representation approach that focuses on representing the communication patterns, dependencies, computation-communication overlap, and algorithms. The SWM code is decoupled from the original application code as well as from any particular simulator, enabling use across different simulation environments. The SWM runtime supports a set of low-level API communication primitives to support a number of MPI-based communication operations. The primitives used by the SWM closely resemble those of MPI and SHMEM, but they are not constrained to specific syntax or semantics. In this paper, we utilize several SWM representations for multiple HPC codes including Nekbone, LAMMPS and nearest neighbor [10].

B. CODES Simulation Framework

The CODES simulation framework provides high-fidelity, massively parallel simulations of prototypical next-generation HPC architectures. The framework has been extensively used for performance analysis of modern interconnect topologies (fat tree, torus, dragonfly, express mesh and slim fly) [8]. The network models have been validated against real architectures [11].

C. Integrating SWM with CODES

To avoid storing large communication traces, we enabled insitu workload replay with SWM through the Argobots library [14] to synchronize between CODES and SWM workloads. Argobots is a low level threading framework that is used to synchronize the execution of the SWM workload in parallel with the CODES simulation framework. More details on CODES and SWM integrations using Argobots are discussed in [15].

IV. EVALUATION METHODOLOGY

In this section, we discuss the network configurations, workloads, rank-to-node mapping policies, and routing algorithms used in the study.

A. Topology and Routing Description

The dragonfly network topology, proposed by Kim et al. [16], consists of groups of routers that are connected to each other with one or more optical channels. Within each group, the routers are directly connected to each other in an all-to-all manner via electrical links. In this paper, we refer to this configuration as a 1-D dragonfly. A variation of a dragonfly topology, deployed in the Cray XC systems, uses a 2-D all-to-all within each group instead of all-to-all connections. We refer to this configuration as a 2-D dragonfly. A 2-D dragonfly traverses almost double the number of hops as a 1-D dragonfly. The hop count traversed by a 1-D dragonfly is close to a

megaflly network, therefore, we compare megaflly with a 1-D dragonfly to ensure a reasonable comparison. Various forms of adaptive routing have been proposed for a dragonfly, which detect congestion and determine whether the packet should take a minimal or non-minimal route.

Routing in a dragonfly network is done by taking either a minimal (typically direct) or nonminimal path. A minimal path uses the global channel connecting the two groups with the source and destination nodes. With a nonminimal path, Valiant’s algorithm [17] is used to route a packet on a minimal path to a randomly selected intermediate router and then route the packet minimally to the destination router. Nonminimal routes are taken if potential minimal routes are congested. Various forms of adaptive routing have been proposed that detect congestion and determine whether the packet should take a minimal or nonminimal route. We use the progressive adaptive algorithm (PAR) provided in [18]. The PAR algorithm in the simulation re-evaluates the minimal path until either the packet decides to take a nonminimal route or the packet reaches the destination group on a minimal path. A 1-D dragonfly network uses up to 6 virtual channels to avoid deadlock, 3 for request and 3 for response. However, prior work shows that 3 virtual channels for either request or response traffic can cause congestion in the intermediate group and suggests adding a fourth virtual channel, which is helpful in alleviating the congestion [18]. In this work, we use 8 virtual channels for progressive adaptive routing in a dragonfly network.

What separates various dragonfly topologies from each other is largely based on the interconnect within a group. Megaflly is a topology that belongs to the Dragonfly class of interconnection networks. At a high level, it is classified as having groups of routers which are, in turn, connected to each other with at least one global connection between any two groups. Megaflly is characterized by its connectivity in the form of a two-level Fat Tree network in each group. This locally defined network is also known as a complete bipartite graph; a graph with two sub-groups where all nodes within one subgroup are connected to all nodes in the other subgroup. There are no connections between the routers within the same subgroup. The two levels in each group have routers that will be referred to as *Leaf Routers*, those that have terminal/compute node connections but no global connections, and *Spine Routers*, those that have global connections to other groups but no terminal/compute node connections [5], [6]. In this paper, we use the progressive adaptive routing algorithm proposed in prior studies on Megaflly networks [6]. Megaflly requires only 4 virtual channels (VCs) to avoid deadlock and none to avoid congestion in the intermediate group.

B. Network Configurations

To perform a comparison of the megaflly network with a 1-D dragonfly, we maintain similar router radix and similar node counts. We used a router radix of 32 ports for both networks. Across the group, the routers are connected via

TABLE I: Configurations of megafly and 1-D dragonfly used for performance comparison

	Megafly	1-D Dragonfly
Router Radix	32	32
Groups	33	65
Nodes/Group	256	128
Node Count	8448	8320
Global Connections / Group	256	128
Link Bandwidth	25GiB/sec	25GiB/sec
Nodes Per Router	16 (Leaves only)	8

global channels. The configurations of the 1-D dragonfly and megafly are given in Table I.

C. Workloads

In order to quantify the slowdown of a particular job due to communication interference, multiple jobs need to be running in parallel to exhibit interference. We conduct two types of interference experiments: (i) replay HPC applications that serve as foreground communication traffic in conjunction with a job that generates synthetic background communication to understand the interference in a controlled manner, and (ii) replay multiple HPC applications in parallel to capture the dynamism of multi-phased communication and quantify the impact of perturbation.

1) *Foreground Traffic*: Previous work demonstrates the performance variability shown by LAMMPS, Nekbone, and MILC applications on the Cray XC40 system [4]. Thus, we use LAMMPS and Nekbone workloads as foreground workloads for our experimental analysis. We also use a 3-D nearest-neighbor communication pattern, which is a commonly used pattern in several HPC applications.

LAMMPS is a large-scale atomic and molecular dynamics code that uses MPI for communication. We use the SWM code that derives its communication pattern from the LAMMPS application. Figure 3(a) shows the message distribution of LAMMPS SWM per rank in a problem involving 2,048 ranks. The LAMMPS workload uses MPI_AllReduce with small messages as well as blocking sends and nonblocking receives for point-to-point communication with large messages.

Nekbone is a thermal hydraulics mini-app that captures the structure of the computational fluids dynamics code Nek5000. Nekbone’s SWM communication pattern is derived from Nekbone. Figure 3(b) shows the message distribution of the Nekbone SWM on a per rank basis in a problem with 2,197 ranks. Nekbone performs a large number of MPI collective operations with small 8-byte messages. It uses nonblocking sends and receives to transmit medium-sized messages.

Cartesian neighborhood communication is a pattern commonly used in multiple scientific applications including Hardware Accelerated Cosmology Code (HACC), fast Fourier transform solvers, and adaptive mesh refinement (AMR) codes. We use a 3-D nearest-neighbor SWM in this work that transmits large messages (64 KiB and 128 KiB) on a per rank basis with multiple iterations of MPI nonblocking sends and receives followed by MPI_Wait_All. A problem size of 4,096 ranks is used with the nearest neighbor SWM.

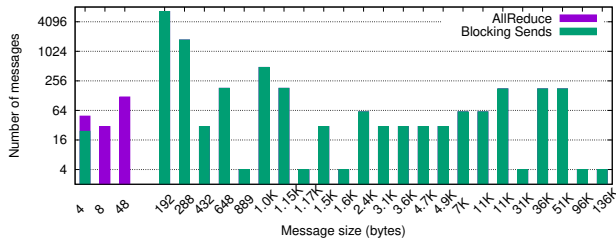
2) *Background Traffic*: The background communication traffic is needed to interfere with the foreground workloads. To ensure an even distribution of traffic that covers a significant fraction of the network, we use a uniform random communication pattern. This is generally considered a benign traffic pattern for dragonfly networks. However, with large messages randomly sent in the network, uniform random causes hotspots at multiple network locations and becomes a source of interference. We varied the amount of data transmitted via uniform random traffic and observed the effect of different data transmission rates on the foreground traffic. The background traffic generation is modeled as a separate job that runs in parallel with the foreground traffic and occupies at least 25% to 50% of the network. The background injection rates depend on the available compute node to router link bandwidth in the network. We inject traffic at a percentage of the available link bandwidth and vary the injection rates between 2% to 36.5% of the link bandwidth. At the 36.5% rate, each node is injecting 9GiB/sec of background traffic with an aggregate network background interference of 18TiB/sec. At this rate, we see significant slowdown (up to 4x for uniform random and up to 7x for random permutation) in application communication times for both networks. Therefore, we keep that as the maximum background injection rate.

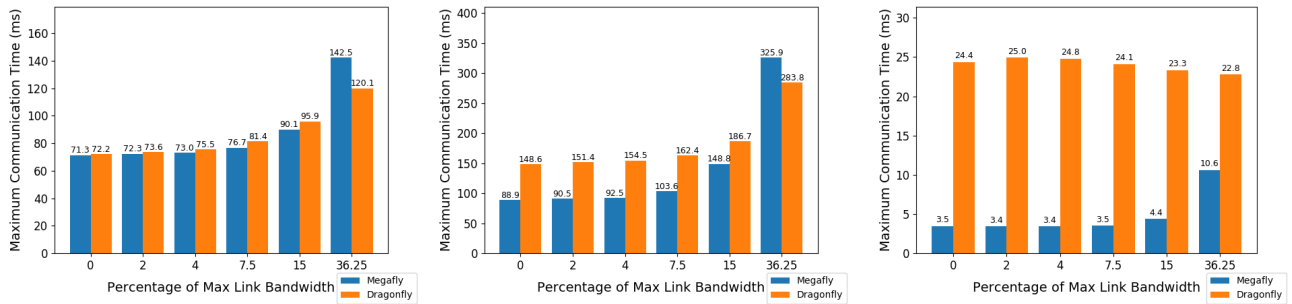
We experiment with two different background communication patterns: (i) a uniform random synthetic pattern where a rank randomly chooses a destination rank and transmits large messages and (ii) a random permutation traffic where a pair of ranks communicate and transmit data until a certain threshold is reached.

3) *Multiple Applications*: As a specific instance of representative HPC scenarios, we ran the three foreground workloads in parallel (Nekbone, nearest neighbor, and LAMMPS). We also ran each of these workloads in isolation on the network to determine the baseline performance and observed the slowdown introduced when the workloads are running in parallel.

D. Rank-to-Node Mappings

Ranks are placed on network nodes in a manner similar to that for production HPC systems, where clusters of available network nodes are assigned to a job. Therefore, we use a geometric job placement policy in which multiple clusters of network nodes are assigned to jobs. In the simulation, the clusters are formed by using the inverse transform sampling method for creating random samples from a given distribution.





(a) LAMMPS communication times comparison

(b) Nekbone communication times comparison

(c) Nearest-neighbor communication times comparison

Fig. 4: Performance of megafly vs. 1-D dragonfly with (a) geometrically allocated 2048/2048 ranks for LAMMPS/uniform random workloads and (b) geometrically allocated 2197/2197 ranks for Nekbone/uniform random workloads, and (c) geometrically allocated 4096/4096 ranks for nearest-neighbor/uniform random workloads. The intensity of the background traffic was scaled at a percentage of the maximum link capacity.

a linear slowdown in the performance of the foreground job as the number of bytes exchanged in the background traffic increases.

3) *Multiple Applications in Parallel*: In the third case, as a specific instance of representative HPC scenarios, we run the three workloads (LAMMPS, Nekbone, and nearest neighbor) in parallel without any synthetic communication traffic. This scenario clearly mimics a common system state, with multiple jobs completing for shared resources. Figure 6 shows the communication time of the three applications when running in parallel and in isolation on both dragonfly and megafly networks. We can see that with both LAMMPS and Nekbone, the applications are much less perturbed on a megafly network than on a 1-D dragonfly network.

Our analysis shows that in most cases, the megafly network is less prone to perturbation than is a dragonfly network. On both networks, however, HPC applications see a significant slowdown in communication ranging up to 700% in the presence of intense background communication traffic.

VI. EVALUATING QUALITY OF SERVICE ON MEGAFLY NETWORKS

Enabling quality of service on HPC networks requires that each traffic class have its own set of virtual channels. Megafly networks require a fewer number of virtual channels for deadlock prevention. When a fixed, limited number of VCs are available in the switch hardware, megafly needs half as many VCs as a dragonfly and has the opportunity to use the extra VCs for QoS. The mechanism for quality of service was introduced in Section II. In this section, we perform experiments to analyze the impact of QoS on traffic interference and application slowdown that we saw in Section V. We explore two configurations through which QoS can be introduced in megafly networks. Since there can be a large number of permutations for bandwidth caps, we performed a sensitivity analysis by sweeping different bandwidth values and picked the values that were most effective (for example, a 30% bandwidth cap

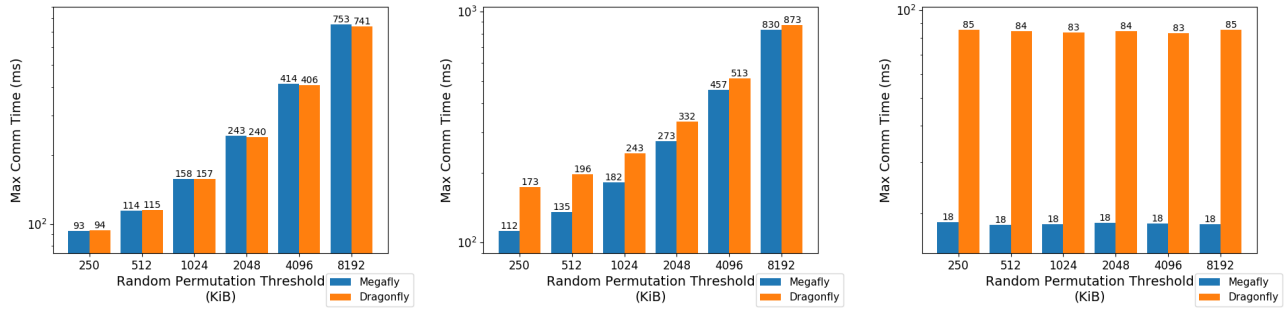
was effective for multiple applications in Section VI-C and a 10% cap for collective priority in Section VI-B). Discussion of the validation of the QoS mechanism is given in [15]. The static window over which the bandwidth statistics were monitored was kept to 5 ms throughout these experiments.

A. QoS Mechanism I: Prioritizing Entire Applications

With our first QoS mechanism, a higher priority and high bandwidth are assigned to the entire application (or set of applications) so that they face minimal slowdown relative to other traffic. We use uniform random background traffic and both LAMMPS and Nekbone foreground workloads that exhibited slowdown on megafly networks in Section V (Nearest neighbor was not getting perturbed). To understand the impact on background traffic, we measure the performance of both foreground workload and background traffic. The background traffic performance is measured by the maximum time to complete a message (all messages have the same size in the synthetic workload).

The benefit of using this QoS approach is that if the foreground application is not utilizing the full bandwidth allocated to it, then the background workload can consume the unutilized bandwidth. Figure 7 compares the performance of LAMMPS workload with and without QoS enabled on a megafly network. It also shows the slowdown to background communication traffic. The LAMMPS workload is not as communication intensive because it involves point-to-point messages along with a small number of MPI_AllReduce messages. Therefore, the perturbation to background traffic is not significant. Because of the high priority given to LAMMPS, it does not see any slowdown even though it is running in parallel with intense background traffic. Additionally, while we observe a significant speedup with LAMMPS, the background traffic observes only a small degree of slowdown as compared with the no-QoS case.

Nekbone SWM is a communication-intensive workload that transmits 4x more data than does LAMMPS SWM, and a



(a) LAMMPS communication times comparison

(b) Nekbone communication times comparison

(c) Nearest-neighbor communication times comparison

Fig. 5: Performance of megafly vs. 1-D dragonfly with (a) geometrically allocated 2048/2048 ranks for LAMMPS/random permutation workloads (b) geometrically allocated 2197/2197 ranks for Nekbone/random permutation and (c) geometrically allocated 2197/2197 ranks for nearest-neighbor/Random Permutation workloads. The amount of data exchanged between two nodes in a rotating random permutation was scaled from 250 KiB to 8 MiB.

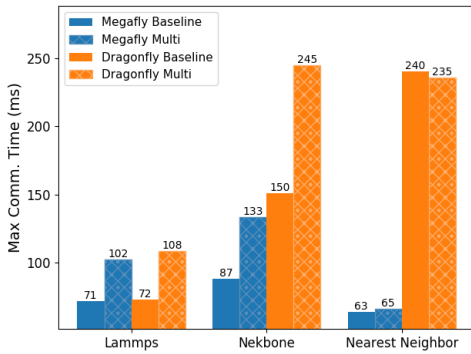


Fig. 6: Communication times of LAMMPS (2,048 ranks), Nekbone (2,197 ranks), and nearest neighbor (2,048 ranks) when running in parallel on 1-D dragonfly and dragonfly+ networks. Baseline indicates the application runs in isolation.

majority of the communication involves collectives. Figure 8 shows the performance of the Nekbone SWM when QoS is enabled. Once again we see Nekbone having minimal slowdown when QoS is enabled while causing minimal slowdown to background communication traffic. The primary reason for the improved performance is that both Nekbone and LAMMPS are given high priority and high bandwidth yet they do not consume all the bandwidth assigned to them. Therefore, the background traffic is able to get the required bandwidth that it needs while observing little slowdown. Both these results demonstrate that traffic differentiation with bandwidth shaping and prioritization can mitigate (or eliminate) communication interference to HPC workloads while causing minimal slowdown to the background traffic. Assigning a high priority to an application can eliminate the perturbation to that application while experiencing a reasonable slowdown in the remaining network traffic.

B. QoS Mechanism II: Prioritizing and Guaranteeing Bandwidth to Latency-Sensitive Operations

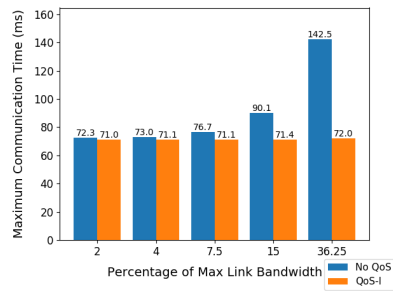
Several HPC applications rely on the performance of MPI collective operations. In a majority of the cases, collectives comprise small messages, and the application performance suffers when heavy background network traffic interferes with the transmission of these messages. An alternative application of QoS is to assign a high priority and guaranteed bandwidth to collective operations.

Figure 9 shows the performance of LAMMPS when high priority is given to collectives and compares it with the case where no QoS is enabled. In this case, we are assigning a high priority but a small fraction of bandwidth to collective operations; the point-to-point operations and background traffic are given a lower priority and higher bandwidth cap (90%). We see that although there is some slowdown in foreground traffic when the background traffic becomes intense, the foreground workload is still 10% faster than the case where no QoS is enabled (specifically in the case of 15% background traffic injection). Additionally, LAMMPS uses more point-to-point operations and has fewer collective operations.

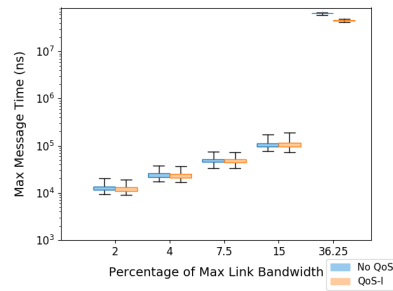
Figure 10 shows the performance of Nekbone when given high priority and a guaranteed bandwidth to collectives. Nekbone relies heavily on collective operations whereas LAMMPS uses fewer collectives. Therefore, we see a significant performance improvement of up to 60% speedup compared with the case where no QoS is enabled. The background communication traffic does not show a slowdown in message communication times; instead it shows a slight performance improvement compared to no-QoS options in one case.

C. Applying QoS Mechanisms to Multiple Application Workloads in Parallel

In this section, we examine both QoS mechanisms in the case where multiple applications are running in parallel, which is a specific instance of a representative HPC system. We compare the QoS-enabled performance with the case where

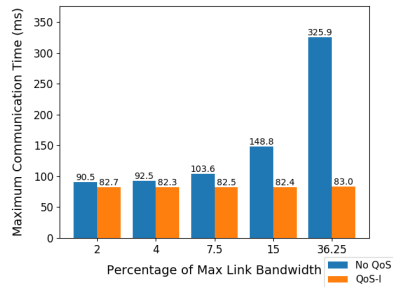


(a) LAMMPS communication times

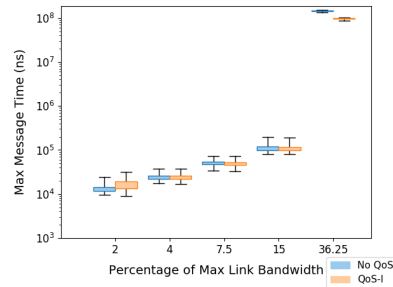


(b) Background traffic performance

Fig. 7: QoS Mechanism I (Application Priority): Performance of LAMMPS and background traffic on megafly network with QoS enabled and disabled. The entire LAMMPS application is given a high priority and high bandwidth (70%).

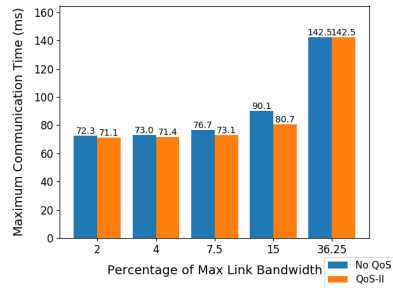


(a) Nekbone communication times

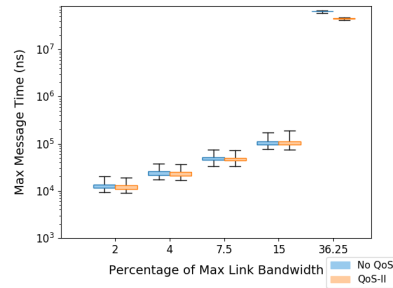


(b) Background traffic performance

Fig. 8: QoS Mechanism I (Application Priority): Performance of Nekbone and background traffic on megafly network with QoS enabled and disabled. The entire Nekbone application is given a high priority and high bandwidth (70%).



(a) LAMMPS communication times



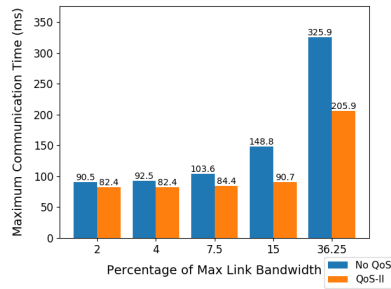
(b) Background traffic performance

Fig. 9: QoS Mechanism II (Collective Priority): Performance of LAMMPS and background traffic on megafly network with application-based QoS enabled and 10% bandwidth guaranteed to collectives.

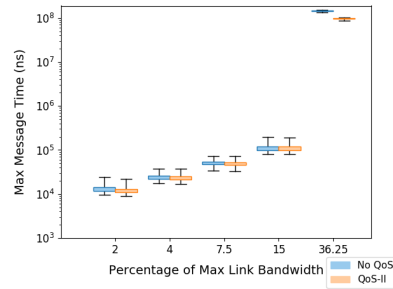
there are multiple applications running without any QoS. For the first QoS mechanism, since Nekbone is more communication intensive than LAMMPS and nearest neighbor (shown in Figure 3) we assign it a separate traffic class with a bandwidth cap of 30% and a high priority. The rest of the bandwidth (70%) is available to both LAMMPS and nearest neighbor. For the second QoS mechanism, we assign a higher priority to all collective communication in both LAMMPS and Nekbone and then see the impact on application performance.

Figure 11 shows the performance of multiple applications running in parallel with and without QoS enabled. In short,

both schemes are beneficial, and lead to reduced communication time. With the QoS Mechanism I, we give a high priority and guaranteed bandwidth to Nekbone (30%). Nekbone is communication intensive; and with a high priority and 30% of the bandwidth cap, it does not get any slowdown due to background communication traffic. In contrast, LAMMPS and nearest neighbor have a lower priority, and they still see a performance improvement compared with the case where there was no QoS enabled. Adding bandwidth caps on Nekbone (which is a bandwidth-intensive application) helps improve the performance of LAMMPS and nearest neighbor as well.



(a) Nekbone communication times



(b) Background traffic performance

Fig. 10: QoS Mechanism II (Collective Priority): Performance of Nekbone and background traffic on megafly network with application-based QoS enabled and 10% bandwidth guaranteed to collectives.

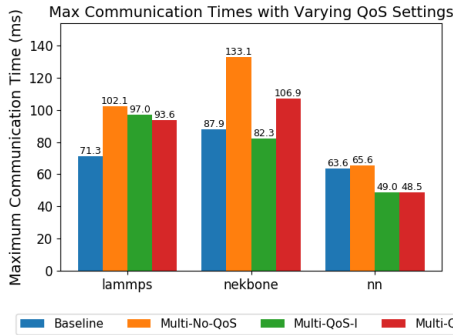


Fig. 11: Communication times of Nekbone, LAMMPS, and nearest-neighbor workloads when running in parallel. Both mechanisms of application-based QoS were enabled. The comparison is done with (i) the worst case when no QoS is enabled (Multi No QoS) and (ii) the best case when the workload is running in isolation with no interference (baseline).

With the QoS Mechanism II, where collective communication is given priority, both LAMMPS and Nekbone benefit by seeing a 10% and 20% speedup, respectively, compared with the case where QoS is not enabled. One interesting observation is the performance of the nearest-neighbor workload, which is much faster with QoS enabled than when the workload is running in isolation. Looking at the adaptive routing statistics, we see that the nearest-neighbor workload when running in isolation takes the maximum number of minimal routes because of the bias toward minimal routes. With QoS enabled, nearest-neighbor traffic has a lower priority with QoS mechanisms enabled, which causes it to take more nonminimal routes, coincidentally helping with the congestion points. Thus the workload sees improved performance. A similar phenomena is observed with Nekbone when it is running with QoS Mechanism I.

These experiments demonstrate the effectiveness of applying QoS to reduce or eliminate communication interference. With both mechanisms, Nekbone, being more bandwidth intensive, sees a 20% to 350% speedup in communication time

compared with the case where QoS is not enabled. LAMMPS sees a 10% to 200% improvement in communication time compared with the case where QoS is not enabled.

VII. RELATED WORK

There are different approaches to address run-to-run variability on HPC systems. One approach is based on partitioning the networks and providing an isolated partition for a job. While this approach has successfully worked for low-radix networks such as torus [20], it is a challenge to implement partitioning on networks such as dragonfly or megafly, due to their hierarchical nature. The other approach is QoS, which can be enforced through various mechanisms on data centers and HPC networks. Flow control [21] [22] is a high-level approach for avoiding interference in large-scale and datacenter-scale networks which takes a coarser-grained look at data within the network. Alizedah et al. [23] studied the impacts of sacrificing a portion of the total bandwidth while lowering the threshold for congestion sensing to provide a buffer zone within links in an attempt to reduce the overall latency of applications in a datacenter environment. On the algorithmic routing side of QoS implementation, many different approaches exist, from centralized global information methods to distributed routing algorithms with limited or incomplete network information and hierarchical algorithms that bridge the gap between globally and locally available information when making routing decisions. Chen and Nahrstedt [7] presented an overview of various routing algorithms solving different QoS problems for both unicast and multicast applications. Most of the literature available on quality of service is intended for data-centric and TCP/IP networks and does not explore HPC workloads, routing, and flow control mechanisms. Cheng et al. [13] provided high-level details about implementing quality of service on data-centric and HPC networks. The work is not specific to a topology and does not provide any performance results. Jakanovic et al. [24] provided an efficient QoS policy for HPC systems with InfiniBand network (fat tree topology).

VIII. DISCUSSION & CONCLUSION

With HPC applications showing performance variation on recent hierarchical interconnects, we analyze communication

interference for both megafly and dragonfly networks. We extend the CODES parallel simulation framework to replay the communication workloads of LAMMPS, Nekbone and nearest neighbor using the concept of Scalable Workload Models (SWM). We introduce moderate to intense background communication traffic during the execution of these communication workloads and compare the slowdown on megafly network with a 1-D dragonfly network. We demonstrate that performance variability is experienced in both topologies, while observing that in a majority of experiments the performance implication is less severe for megafly.

To further mitigate the variability, we introduce traffic differentiation and quality of service mechanisms in megafly networks, as megafly makes QoS implementation feasible by requiring fewer VCs for deadlock avoidance. We explore two different QoS mechanisms for HPC workloads (i) prioritizing and bandwidth capping entire HPC applications (ii) prioritizing and guaranteeing bandwidth to latency sensitive collective operations with small messages. With the first mechanism, performance results show that when a high priority and a bandwidth cap is given to entire HPC applications, it can eliminate performance variability while the rest of the background traffic also sees minimal impact. For the second mechanism, we show that when a small fraction of bandwidth is guaranteed to latency sensitive operations like the MPI collectives, it can mitigate the performance variability by 10% to 60% depending upon the intensity of collective communication in the application.

While this work is aimed to provide a proof of concept that QoS is effective in mitigating communication interference for realistic HPC workloads, there are a number of avenues that need to be further explored. First, real HPC systems have tens to hundreds of jobs running. Giving a high priority to more than one HPC application (as shown in QoS mechanism I) can introduce interference within the traffic class, which can slowdown high priority applications. Secondly, a statically allocated time window for bandwidth monitoring can lead to bandwidth imbalance with bursty and irregular workloads. Having a dynamic or sliding bandwidth monitoring window may be more effective for such cases. Finally, one would need to explore how to expose the traffic classes to the MPI interfaces and the scheduler.

ACKNOWLEDGMENT

This research was supported by the Exascale Computing Project (17-SC-20-SC), a joint project of the U.S. Department of Energy's Office of Science and National Nuclear Security Administration, responsible for delivering a capable exascale ecosystem, including software, applications, and hardware technology, to support the nation's exascale computing imperative. The work has used resources from the Argonne's Leadership Computing Facility (ALCF), Rensselaer's CCI supercomputing center and Argonne's Laboratory Computing Resource Center (LCRC).

REFERENCES

- [1] NERSC, "Cori," <https://www.nersc.gov/users/computational-systems/cori/>.
- [2] Los Alamos National Laboratory, "Trinity Cray XC40 system," <http://www.lanl.gov/projects/trinity/>.
- [3] Argonne Leadership Computing Facility (ALCF), "Theta, Argonne's Cray XC System:" [Online]. Available: <https://www.alcf.anl.gov/theta>
- [4] S. Chunduri *et al.*, "Run-to-run variability on Xeon Phi based Cray XC systems," in *Supercomputing(SC)*. ACM, 2017, p. 52.
- [5] A. Shpiner, Z. Haramaty, S. Eliad, V. Zdornov, B. Gafni, and E. Zahavi, "Dragonfly+: low cost topology for scaling datacenters," in *High-Performance Interconnection Networks in the Exascale and Big-Data Era (HiPINEB)*, 2017 IEEE 3rd International Workshop on. IEEE, 2017, pp. 1–8.
- [6] M. Flajslik, E. Borch, and M. A. Parker, "Megafly: A topology for exascale systems," in *International Conference on High Performance Computing*. Springer, 2018, pp. 289–310.
- [7] S. Chen and K. Nahrstedt, "An overview of quality of service routing for next-generation high-speed networks: problems and solutions," *IEEE network*, vol. 12, no. 6, pp. 64–79, 1998.
- [8] M. Mubarak *et al.*, "Enabling parallel simulation of large-scale HPC network systems," in *IEEE Transactions on Parallel and Distributed Systems*. IEEE, 2017.
- [9] X. Yang *et al.*, "Watch out for the bully! job interference study on dragonfly networks," in *Supercomputing (SC)*, 2016.
- [10] Intel Corporation, "Scalable Workload Models for System Simulations," <http://hpc.pnl.gov/modsim/2014/Presentations/Thompson.pdf>.
- [11] Mubarak, Misbah and Ross, Robert B. and Carothers, Christopher D., "Validation of CODES dragonfly model with Theta Cray XC System," 2017. [Online]. Available: <https://publications.anl.gov/anlpubs/2017/05/135666.pdf>
- [12] S. Chunduri *et al.*, "Characterization of mpi usage on a production supercomputer," in *Supercomputing (SC)*. IEEE/ACM, 2018.
- [13] A. S. Cheng, T. D. Lovett, and M. A. Parker, "Traffic class arbitration based on priority and bandwidth allocation," Dec. 22 2016, uS Patent App. 15/120,038.
- [14] S. Seo, A. Amer, P. Balaji, C. Bordage, G. Bosilca, A. Brooks, P. Carns, A. Castelló, D. Genet, T. Herault *et al.*, "Argobots: a lightweight low-level threading and tasking framework," *IEEE Transactions on Parallel and Distributed Systems*, vol. 29, no. 3, pp. 512–526, 2018.
- [15] M. Mubarak, N. McGlohon, M. Musleh, E. Borch, R. Ross, R. Huggahalli, S. Chunduri, S. Parker, C. Carothers, and K. Kalyan, "Exascale Computing Project Hardware Evaluation Milestone Report: Evaluating Quality of Service on High-radix HPC Networks."
- [16] J. Kim, W. J. Dally, S. Scott, and D. Abts, "Technology-driven, highly-scalable dragonfly topology," *ACM SIGARCH Comput. Architecture News*, vol. 36, no. 3, pp. 77–88, Jun. 2008.
- [17] L. G. Valiant, "A scheme for fast parallel communication," *SIAM journal on computing*, vol. 11, no. 2, pp. 350–361, 1982.
- [18] J. Won, G. Kim, J. Kim, T. Jiang, M. Parker, and S. Scott, "Overcoming far-end congestion in large-scale networks," in *High Performance Computer Architecture (HPCA)*, 2015 IEEE 21st International Symposium on. IEEE, 2015, pp. 415–427.
- [19] M. Mubarak *et al.*, "Quantifying i/o and communication traffic interference on dragonfly networks equipped with burst buffers," in *IEEE Cluster Conference*. IEEE, 2017.
- [20] Z. Zhou, X. Yang, Z. Lan, P. Rich, W. Tang, V. Morozov, and N. Desai, "Improving batch scheduling on blue gene/q by relaxing 5d torus network allocation constraints," in *IPDPS*. IEEE, 2015, pp. 439–448.
- [21] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner, "OpenFlow: Enabling innovation in campus networks," *SIGCOMM Comput. Commun. Rev.*, vol. 38, no. 2, pp. 69–74, 2008.
- [22] A. R. Curtis *et al.*, "DevoFlow: Scaling flow management for high-performance networks," in *Proceedings of the ACM SIGCOMM 2011 Conference*, ser. SIGCOMM '11. New York, NY, USA: ACM, 2011, pp. 254–265.
- [23] M. Alizadeh *et al.*, "Less is more: Trading a little bandwidth for ultra-low latency in the data center," in *9th USENIX*, ser. NSDI'12. Berkeley, CA, USA: USENIX Association, 2012, pp. 19–19.
- [24] A. Jakanovic *et al.*, "Effective quality-of-service policy for capacity high-performance computing systems," in *HPCC-ICISS*. IEEE, 2012, pp. 598–607.