

Mitigating Network Noise on Dragonfly Networks through Application-Aware Routing

Daniele De Sensi
desensi@di.unipi.it
University of Pisa
Pisa, Italy
ETH Zurich
Zurich, Switzerland

Salvatore Di Girolamo
salvatore.digirolamo@inf.ethz.ch
ETH Zurich
Zurich, Switzerland

Torsten Hoefler
htor@inf.ethz.ch
ETH Zurich
Zurich, Switzerland

ABSTRACT

System noise can negatively impact the performance of HPC systems, and the interconnection network is one of the main factors contributing to this problem. To mitigate this effect, adaptive routing sends packets on non-minimal paths if they are less congested. However, while this may mitigate interference caused by congestion, it also generates more traffic since packets traverse additional hops, causing in turn congestion on other applications and on the application itself. In this paper, we first describe how to estimate network noise. By following these guidelines, we show how noise can be reduced by using routing algorithms which select minimal paths with a higher probability. We exploit this knowledge to design an algorithm which changes the probability of selecting minimal paths according to the application characteristics. We validate our solution on microbenchmarks and real-world applications on two systems relying on a Dragonfly interconnection network, showing noise reduction and performance improvement.

CCS CONCEPTS

• **Networks** → **Network performance evaluation**; • **Computer systems organization** → **Distributed architectures**; • **Computing methodologies** → **Distributed computing methodologies**.

KEYWORDS

network noise, dragonfly, routing

ACM Reference Format:

Daniele De Sensi, Salvatore Di Girolamo, and Torsten Hoefler. 2019. Mitigating Network Noise on Dragonfly Networks through Application-Aware Routing. In *The International Conference for High Performance Computing, Networking, Storage, and Analysis (SC '19)*, November 17–22, 2019, Denver, CO, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3295500.3356196>

1 INTRODUCTION

Interconnection networks are the backbone of large-scale supercomputers often connecting tens of thousands of servers. The cost,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SC '19, November 17–22, 2019, Denver, CO, USA

© 2019 Association for Computing Machinery.
ACM ISBN 978-1-4503-6229-0/19/11...\$15.00
<https://doi.org/10.1145/3295500.3356196>

performance, and maintainability depend on details of the networking technology and topology. Commonly deployed low diameter (e.g., Dragonfly) and hierarchical topologies (e.g., Fat Tree) usually share network resources between applications running in different allocations [5, 34]. This sharing can lead to interference between applications, e.g. if one application communicates heavily and fills shared links and switch buffers with large numbers of packets, another application that only sends small synchronization messages may be severely delayed by the resulting head-of-line blocking.

The caused communication delays can destroy the performance of large-scale applications, similarly to operating system (OS) noise. Indeed, it has been shown that OS noise limits the scalability of many applications to 15k processes on the 200k core Jaguar supercomputer [28]. Operating systems have been tuned to minimize interference through isolation, for example, by scheduling management tasks on separate cores [20]. However, network isolation and efficient use of computing resources cannot easily be achieved on today's supercomputers. The resulting *network noise* has been observed in various research groups [13, 43, 47]. Like operating system noise, we define network noise an *external effect* on application performance, caused by *sharing resources* with activities *outside of the control of the affected application*. Network noise can either be caused by the High-Performance Computing (HPC) system itself (e.g., control of distributed filesystems) or by other applications running simultaneously (e.g., cross traffic). Thus, in general, *the programmers of an application cannot avoid noise*, they can at best mitigate it, for example, using nonblocking collective operations [25, 26].

At the system level, network noise can be avoided by different application allocations to isolated partitions of the system. For example, to sub-trees in a Fat Tree topology [38] or groups of a Dragonfly interconnect. However, this is only possible if the allocation sizes exactly match the topology layout and resources are available—introducing such a strategy in any batch system will significantly reduce the system utilization. Thus, network noise is generally accepted in HPC systems and was of not much concern until recently. Yet, we argue that growing system sizes, as well as the introduction of adaptive routing technologies, aggravate the situation, causing up to 2X slowdowns, as we will show in Section 5. Adaptive routing has been introduced to increase the overall utilization of the network—it is in fact *required* in low-diameter topologies for most traffic patterns [9, 32, 34]. The downside of adaptive routing is that even two communicating nodes can congest resources on all paths that the adaptive routing utilizes, not just a single path in static routing. This is often referred to as *packet spraying* and

is feared by datacenter operators in so-called *incast* or *hot-spot* (many-to-one) patterns. In addition to this, adaptive routing may be affected by the so-called *phantom congestion* problem [46]. Namely, as congestion information is propagated with some delay, a node may react too late to congestion events. In this paper, we will be first to demonstrate and quantify the influence of different routing schemes on network noise in practice.

Analyzing network noise in detail is delicate because in practice, when observing application delays, it is hard to distinguish between network noise, operating system noise, and application imbalance. Other works have used network counters but may run into the fallacy that correlation is not causation by ignoring the aspects of unrelated traffic. We will clarify several potential problems in investigations of network noise and develop a set of general guidelines for our analysis. Using these guidelines, we study the relationship between different adaptive routing schemes, application performance, and network noise. Our findings on two large-scale Cray Aries systems, Piz Daint and NERSC’s Cori, are remarkable: not only will changing the adaptive routing mode reduce communication times and speed up applications up to twice, but it will also significantly reduce performance variation.

We find that the best routing mode that minimizes network noise and maximizes performance depends not only on the characteristics of the allocation but also on the communication load. For example, large-scale alltoall communications are best routed with the default mode whereas many other communication patterns benefit from mostly minimal routing. Because this all depends on the location of the communication peers, there is no simple static rule to select the best routing mode. To address this, we develop a simple but effective dynamic routing library that observes the network state through local network counters and adjusts the routing mode for each message based on application characteristics such as the message size and the observed network state. In essence, our routing library is application- and system-aware and acts as a coarse-grained guide that adjusts the packet-level adaptive routing based on application characteristics and dynamic network state.

In summary, the key contributions of this work are:

- We describe mechanisms to gather a detailed understanding of noise in real applications caused by the network using network counters.
- We provide a detailed analysis of two real-world systems with low-diameter topologies.
- We show that much of the application delay is due to network noise that stems from non-trivial interactions between routing strategies and application characteristics.
- We evaluate different adaptive routing strategies using bias on Cray Aries (Cascade) systems.
- We develop an application-aware routing library that routes different applications and application phases with different routing modes, leading to speedups of 2x on some microbenchmarks and real applications.

2 NETWORK PERFORMANCE COUNTERS ON CRAY ARIES

We describe in this section the topology of the Cray Aries Dragonfly network, the different routing algorithms available, and the network

counters provided. Eventually, we introduce a performance model which, by using selected network counters, can be used to estimate the transmission time of a message.

2.1 The Cray Aries Interconnect

The Cray Aries Network [6] is a high performance interconnect based on the Dragonfly topology [34]. The underlying idea of the Dragonfly topology is to partition compute nodes and routers in fully connected groups that act as very high radix virtual routers. The Aries interconnect is organized in three connectivity tiers: groups, chassis, and blades. Each group contains six chassis and within each chassis there are sixteen blades. Each blade contains the Aries router and four nodes. Figure 1 sketches two groups of an Aries interconnect and the internal organization of an Aries device. An Aries device is a system-on-chip comprising a 48-port router and 4 Network Interface Controllers (NICs). The router is organized in tiles: each tile provides a bidirectional link with a bandwidth ranging from 4.7 and 5.25 GB/s per direction (depending on whether the link is optical or electrical). Each router can have up to ten optical connections to routers in different groups: the total number of groups is constrained by the number of routers per group. Often, systems are configured to use more than one tile per inter-group connection, increasing the inter-group bandwidth. The router is connected to all the other routers on the same chassis using the 15 intra-chassis tiles and to 5 other routers sitting on different chassis of the same group using three intra-group tiles per connection.

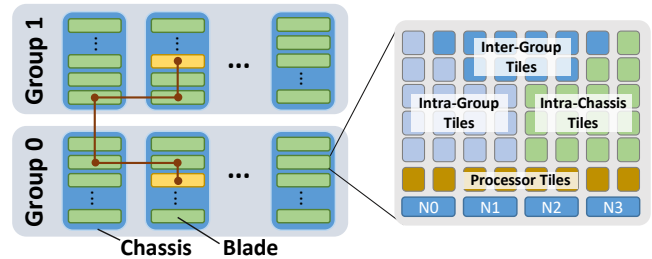


Figure 1: Cray Aries Topology and Aries Device.

As a consequence, routers inside a group are not fully connected. For example, in Figure 1, any router in Group 1 is directly connected to all the routers in the same “column” (i.e., in the same chassis), and to all the routers in the same “row” (i.e., all the routers in the same position but on the other chassis). For this reason, to reach its destination, a packet traverses between one and six tiles of the Aries interconnect. Figure 1 shows the example of a 5-hop minimal path between the two yellow blades, in the case where there is no direct link between their *inter-group* tiles. In this example, we assume that Group 0 is connected to Group 1 through the second blade in its first chassis. When a node on the yellow blade on Group 0 sends a packet to a node on the yellow blade on Group 1, the packet needs first to be routed to the second blade on its chassis (since it is the only blade directly connected to the second blade on the first chassis). Multiple paths are available between any pair of blades and the routing algorithm could take different decisions (e.g., routing first the packet from the yellow blade to the third

blade in the first chassis and then do an intra-chassis hop on the first chassis). The tiles have a finite input queue and use a credit flow control scheme to avoid overflowing the destination tile. As a consequence, a packet may stall in any of the tiles that are on its path to the destination.

The router also has 8 processor tiles to connect to the 4 Aries NICs included on the Aries Device (N0, ..., N3 in the figure) and each pair of NICs shares four processors tiles. The NICs are connected to the respective compute nodes via independent x16 PCIe Gen3 host interfaces. The compute nodes can issue commands to the NIC via the host interface: in case of a data movement command, the NIC is in charge to packetize the data and issue up to 64 bytes request packet to the connected processor tiles. Each request packet is then acknowledged by a response packet.

Data can be carried either in request packets (in case of *RDMA PUT* calls), or in response packets (for *RDMA GET* calls). Each packet is composed by a number of *flits*. When a request is sent, the NIC splits the request packet in one *header* flit plus one to four payload flits¹, and transmits one flit per clock cycle. It is thus possible to estimate the number of packets and flits which will be sent by the NIC by knowing the type of *RDMA* call and the size of the application message (i.e., we will have 1 packet every 64 bytes, made of 1 request flit for *GETs* and 5 request flits for *PUTs*).

In this paper we will consider two different machines based on Cray Aries Network:

Piz Daint A Cray XC50 system hosted by CSCS, with compute nodes equipped with a 12-core Intel Xeon E5-2690 v3 CPU with 64 GiB RAM and with Hyper-Threading support.

Cori A Cray XC40 system hosted by NERSC, with compute nodes equipped with a 16-core Intel Xeon E5-2698 v3 CPU with 128 GiB RAM and with Hyper-Threading support.

2.2 Adaptive routing and bias

On the Aries network, each packet can be independently routed, and because two nodes can be connected by several (minimal and non-minimal) paths, adaptive routing is used so that packets are sent to less congested paths. The adaptive routing algorithm adopted is a variation of UGAL routing [31]. Packets sent on non-minimal paths will traverse an intermediate group connected to both source and destination groups, increasing the maximum number of hops up to 10 on the largest networks [6]. Every time a packet is sent, two minimal and two non-minimal paths are randomly selected, and the congestion of these paths is estimated by using both local information (e.g., queue occupancy) and estimation of *far-end* congestion, based on the current flow credits available.

However, due to the long inter-router latency, credit information may be delayed, resulting in inaccurate congestion information, leading the adaptive algorithm to select non-minimal paths even if that was not necessary anymore [46]. For this reason, known as *phantom congestion*, a *bias* value can be added to the congestion estimated for non-minimal paths, so that the higher is the bias, the higher is the probability that the packet will be routed on a minimal path. Although it is not possible to set an arbitrary value for the bias, for MPI applications, the bias can be selected by the

user among a restricted set of possibilities (which exact value is not public) by setting the `MPICH_GNI_ROUTING_MODE` environment variable before starting the application².

This variable can be set to one of the following values:

ADAPTIVE_0 No bias is added. We will refer to this algorithm as `ADAPTIVE`.

ADAPTIVE_1 Bias towards minimal routing increases as the packet approaches the destination [7]. It is the default routing algorithm used for `MPI_Alltoall` communications. We will refer to this algorithm as `INCREASINGLY MINIMAL BIAS`.

ADAPTIVE_2 A *low* bias is added.

ADAPTIVE_3 A *high* bias is added. We will refer to this algorithm as `ADAPTIVE WITH HIGH BIAS`.

Moreover, this variable can also be used to enforce deterministic routing rather than adaptive routing, by setting one of the following values:

MIN_HASH Packets are always routed minimally, and the path is selected based on a hash of some fields of the packet header.

NMIN_HASH Packets are always routed non-minimally, and the path is selected based on a hash of some fields of the packet header.

IN_ORDER Packets are always routed minimally, and the packets are received in the same order they were transmitted.

In this work, we will only focus on `ADAPTIVE_0`, `ADAPTIVE_1`, and `ADAPTIVE_3` routing algorithms. Indeed, the performance of `ADAPTIVE_2` lie between those of `ADAPTIVE_0` and `ADAPTIVE_3` because its bias lies between that of these two algorithms. Moreover, we will not consider `MIN_HASH`, `NMIN_HASH`, and `IN_ORDER` because they are not adaptive algorithms.

2.3 Counted events

Aries provides several network counters, which can either be accessed either by using the PAPI library [35] or the CrayPat tool, allowing the user to monitor the impact of the network on his application and vice-versa. Counters are present on both NICs and processor/network tiles. However, users can only access counters on the NICs and tiles associated with their jobs. Because adaptive routing is used, packets may traverse routers which are entirely allocated to other jobs, and by relying on network tiles counters, we would only have a partial view of the impact of our application traffic on the network. Moreover, each tile can be traversed by traffic coming for different jobs, and there is no way to isolate the contribution of each individual job.

For these reasons, we rely solely on the NIC network counters. Among the different provided counters, we will focus on the following ones, which provide enough information for our purposes:

Request Flits Number of request flits sent.

Request Flits Stalled Cycles This counter increments every clock cycle a ready-to-forward flit is not forwarded because of back-pressure. The ratio between this counter and the number of request flits represents the average number of cycles a flit must wait before being transmitted.

¹The size of the flit may be different in other points in the network (e.g., on network tiles), thus changing the number of flits which make up a packet.

²The routing algorithm for `MPI_Alltoall` calls can be separately selected through the `MPICH_GNI_A2A_ROUTING_MODE` environment variable.

Request Packets Number of request packets sent.

Request Packets Cumulative Latency Cumulative latency (microseconds) across all the request-response packets pairs. By dividing this counter for the previous one we get the average packet latency. This counter does not include the time a flit waits in NIC queues before being transmitted.

Detailed information about all the network counters available on Aries can be found on Cray’s technical documentation [1].

2.4 Performance model

We now show how to model how the significant network counters described in Section 2.3 influence application performance. Inspired by the well-known LogP model [14], we develop a model including the average stall cycles and the latency that we observe through the counters. We define with L the packet latency in NIC cycles, with RTT the round-trip-time of a flit, with s the average number of cycles a flit waits (due to stalls) before being transmitted, with k the number of flits per packet, and with f the number of flits which compose the application message.

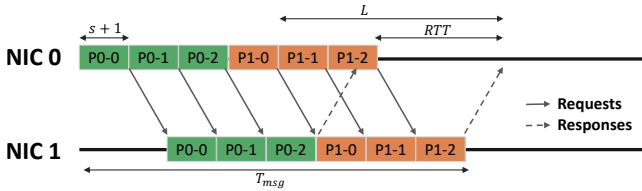


Figure 2: Relationship among NIC performance counters and transmission time of an application message.

We illustrate all parameters in Figure 2 for a *PUT* message and note that a similar model can be derived for *GET* messages. In the picture we show a scenario where the application message is decomposed in six flits (f), divided into two network packets (shown in green (P0) and orange (P1)). Each network packet comprises three request flits (k) traveling from the sender’s NIC to the receiver’s NIC, and one response flit acknowledging the reception of the request and traveling in the opposite direction. The way in which we defined s and L reflects the quantities measured by the counters available on Aries Networks (described in Section 2.3). For example, L measures the latency between the sending of the first request flit of the packet (excluding the waiting time $s + 1$) and reception of the last response flit³.

First of all, for reasons which will become clearer in Section 3.3, we are interested in estimating the impact of the network without including any host-side delays. This can be done by measuring T_{msg} , i.e., how many cycles elapse between the reception of a *PUT* call by the sender’s NIC to the moment when the last flit is delivered to the receiver’s NIC. In the absence of stalls, the NIC can transmit one flit per clock cycle. For this reason, if stalls are present, each flit is stalled on average for s cycles and the NIC transmits one flit every $s + 1$ cycles.

We can then define the transmission time of a packet as the time the first flit takes to reach the receiver’s NIC (i.e., $\frac{RTT}{2}$), plus the

³The counter on Aries NICs provides a measure of the latency in microseconds. This can be easily converted to NIC cycles given the clock frequency of the NIC.

time required to transmit f flits. The round-trip-time RTT can be obtained by removing from the latency L the time the sender’s NIC spent in transmitting $k - 1$ flits, i.e., $RTT = L - (k - 1) \cdot (s + 1)$. However, s is usually some order of magnitude smaller⁴ than L and we can approximate RTT as L , thus obtaining:

$$T_{msg} = \frac{L}{2} + f \cdot (s + 1) \tag{1}$$

However, Aries NICs can have at most 1024 outstanding packets. For this reason, if more than 1024 packets need to be sent, the NIC must wait for the reception of the responses for sent packets before transmitting additional packets. As a consequence, the transmission of some packets may be not fully overlapped and this would increase T_{msg} by a quantity proportional to the latency. In the best-case scenario, this would happen only once every 1024 packets. Defining p as the number of the packets, the transmission time of the message would then be:

$$\begin{aligned} T_{msg} &\cong \frac{p}{1024} \cdot L + \frac{L}{2} + f \cdot (s + 1) = \\ &= \frac{p + 512}{1024} \cdot L + f \cdot (s + 1) \end{aligned} \tag{2}$$

Whereas L and s can be obtained through network counters, f and p can be estimated from the message size and the type of *RDMA* request issued. To validate this model, we compared the estimations made by the model with the actual execution time of a ping-pong benchmark executed over 40 different allocations on the *Piz Daint* machine, obtaining an average 79% correlation on different message sizes, ranging from 128 bytes to 16MiB. We will leverage this performance model in Section 4 to analyze the impact of the routing algorithm on the network noise and to design our application-aware routing algorithm.

3 NETWORK NOISE ESTIMATION

Estimating the true impact of network noise on an application is a complex task. The main problem is to isolate the effects of network noise from other effects causing performance variability such as (operating) system noise or varying resource mapping strategies. In the following, we will describe and categorize such effects and derive simple but important rules for designing experiments with network noise. We note that these strategies for experimental design have not been described in previous works and, as we will show in this section, may lead to overestimation of network noise. We also quantify the potential influence that each strategy may have on the final outcome if it was ignored. Although in the following we will apply these rules to analyze noise on Dragonfly networks, they are general enough to be applied to other interconnection networks as well.

3.1 Fixing the allocation

Process-to-node allocation strategies attempt to solve a complex scheduling problem by trying to find the best balance between fairness, time to completion, topology mapping, throughput, and many

⁴For readability reasons, proportions in Figure 2 do not represent the true scales of L and s , i.e. s is in the order of a few cycles and L is in the order of thousands of cycles.

other metrics. Thus, it is not rare that the processes of a particular compute job are scattered throughout the network— making it vulnerable to network noise. Specifically in low-diameter networks such as Dragonfly [33, 34], paths between two arbitrary nodes often have widely different performance. Thus, changing the allocation will change the performance, even in the absence of any network noise.

Figure 3 quantifies the impact of different allocations on a simple ping-pong benchmark with a 16KiB message between two nodes on the *Piz Daint* system. We compare allocations with varying processor-to-node mappings for source and destination process: the same blade (*Inter-Nodes*), two different blades (*Inter-Blades*), two nodes on different chassis (*Inter-Chassis*), and two nodes on different groups (*Inter-Groups*). Each of these tests has been run for 5 hours recording one round-trip per second in the same allocation. It illustrates how different allocation strategies not only change the mean but also the variance and distribution of measured performance values. The 95% confidence interval for the median is represented as a notch around the median (not visible on this specific plot because its width is < 5% of the median).

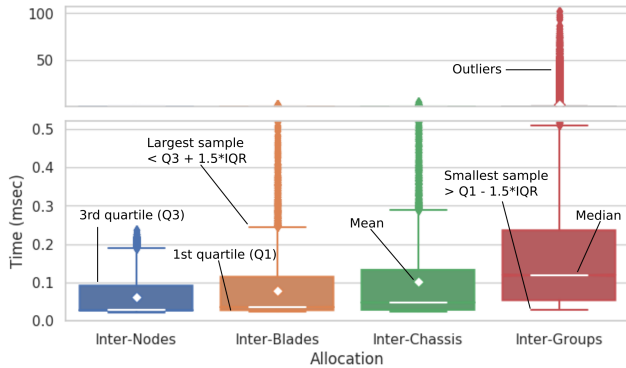


Figure 3: Performance of the ping-pong benchmark for different allocations on the *Piz Daint* system. IQR = Inter-Quartile Range (i.e. $Q3 - Q1$).

The figure clearly shows not only a growing median time but a massively increasing variance. Some outliers in the inter-group allocation are three orders of magnitude larger than the median, which even pulls the average into the regime of outliers for inter-group communication! Furthermore, whereas the inter-chassis median is barely higher than the inter-nodes median, the numerous outliers increase the average of the former to be nearly twice as high. We gathered similar results for more complex collective communication benchmarks. This demonstrates the huge potential influence (three orders of magnitude!) that different allocations could have when collecting data about network noise.

To isolate the effects of network noise, we must avoid effects from varying processor allocations. The simplest strategy is to only compare and analyze results that were collected within the same allocation/job execution.

3.2 Correlation is not (always) causation

A common approach in analyzing the impact of network noise is to correlate the execution time of network-intensive applications to the network traffic intensity, measured through network tiles counters [24]. For example, let us assume we observed an increase in both the execution time and the number of flits that traversed the routers used by the application. At a first sight, we may conclude that because there was more traffic, each packet had to wait for a longer time before traversing a link, increasing the message latency and slowing down the entire application. However, if the application was delayed for reasons not related to network (e.g., OS noise, imbalance between nodes, etc...), we would observe the network for a longer period and we would generally see a higher number of flits, due to other applications sending packets through the routers we are monitoring.

(Idle) Time (sec)	Incoming Flits	Stalled Cycles
1	110M	94M
2	255M	157M

Table 1: Relation of (Idle) Time, Flits, and Stalls

Table 1 demonstrates the effect. It shows the (idle) execution time, the number of flits and the number of stalled cycles for an application executed on 16 nodes, spanning over 5 blades. The application just waits 1 or 2 seconds, respectively, and then terminates. Although the execution time and the number of flits are correlated, it is clear that the longer execution time caused an increase in the number of observed flits, rather than the other way around.

This is also a relevant problem for those solutions that try to correlate network counters to execution time using machine learning approaches [43]. Indeed, such algorithms may conclude that the network intensity is the most relevant feature having an impact on the execution time even if there is no causal relationship. *This issue can be mitigated by normalizing the counters with respect to the observation interval or can be completely avoided by relying on NICs counters measuring latency and stalls, because they have a direct effect on the application performance.*

3.3 Communication time variation is not network noise

Another common way to estimate network noise is to analyze the variability in the execution time of the communication phases of the application, for example by focusing on the execution time of MPI routines [13, 24, 43]. However, especially for collective operations, this would also include other delays which do not depend on the network, such as OS noise [28], synchronization overheads due to application imbalance [17], or contention for shared resources. To provide evidence that not all the variations we observe on the execution time of the network routines are caused by network noise, we show in Figure 4 the performance of an `MPI_Alltoall` collective operation executed by 8 processes running on **the same node** of the *Piz Daint* machine, for different message sizes. Even though the network is not used at all, we can observe a significant performance variability. This demonstrates that varying execution

times of communication operations are not always a good indicator of network noise.

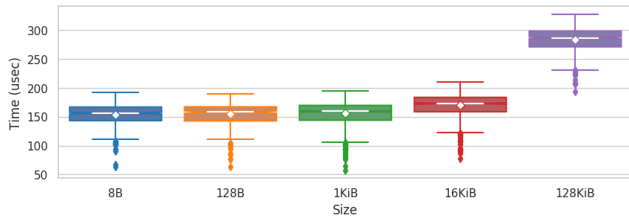


Figure 4: Execution time distribution of an MPI_Alltoall collective operation executed by 8 processes running on the same node of the Piz Daint machine, for different message sizes.

Moreover, this problem is not only relevant when multiple concurrent processes per node are used. Indeed, Figure 5 reports the variability of both the execution time and the network packet latency of a ping-pong benchmark between two nodes in two different groups on the Piz Daint machine, with only one process per node. Because in this specific case the output requests did not experience any stall, latency provides a good approximation of the variability in the time required to transmit the message.

We measure the variability by using the *Quartile Coefficient of Dispersion* (QCD), defined as:

$$QCD = \frac{Q3 - Q1}{Q3 + Q1}$$

where $Q3$ and $Q1$ are the third and the first quartile respectively. This would give us a measure of how much the data is concentrated around the median, i.e., the higher the value, the higher is the variability in the data.

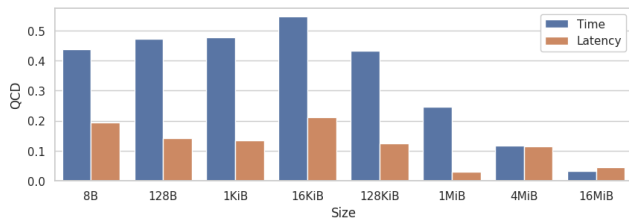


Figure 5: Quartile coefficient of dispersion of execution time and packet latencies, for a ping-pong benchmark between two different groups on the Piz Daint machine, for different message sizes.

As we can see from the figure, even when removing node-side contention, the variability in the execution time of communication routines is still an overestimation of the network noise. This is particularly true for small messages, whereas the impact of latency on execution time decreases for larger messages.

To avoid this problem, we must only consider the delays which are induced by the network, for example by using network counters that measure the actual packet latencies, and which do not include the host-side delays.

4 ROUTING IMPACT ON NETWORK NOISE

After establishing a baseline for measuring network noise in isolation, we now analyze the impact of the routing algorithm on noise. We will show how a significant share of the network noise on a Dragonfly network is caused by the adaptive routing algorithm. We analyze the causes of this behavior and we will exploit this information to design an algorithm that, at runtime, can detect and mitigate network noise and improve application performance. Because the algorithm does not make any assumption on the network topology, it could also be used to mitigate network noise on other networks relying on non-minimal adaptive routing.

4.1 Interactions between noise and routing

As described in Section 2.2, each time a packet is sent, the adaptive routing algorithm will estimate the congestion of two minimal and two non-minimal paths (chosen randomly), and will then route the packet on the path which is estimated to be the least congested one. This process introduces variability in the packets latencies for different reasons. First of all, each packet takes a different path, with a different number of hops and thus a different latency. Moreover, if on one side by taking non-minimal paths the packets avoid congestion, by traversing more routers there will be more traffic on the network, generating congestion on other applications but also on the application itself.

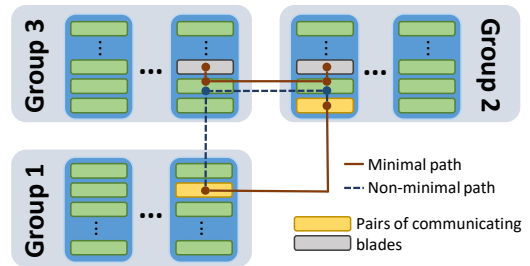


Figure 6: Example of network noise generated by adaptive routing.

For example, let us consider the scenario depicted in Figure 6, where the node on the yellow blade on group 1 needs to send a message to the node on the yellow blade on group 2. After estimating the congestion on both the minimal and non-minimal paths, the node decides to send the packet on the non-minimal path traversing group 3. However, in the meanwhile, the two nodes on the gray blades start to communicate and the traffic generated by the yellow blade would introduce noise on the application running on the gray blades. If the two gray blades were allocated to the same job running on the yellow blades, this would introduce noise on the application itself. Clearly, a similar situation could happen also if only minimal paths are selected. However, due to the higher number of hops, the problem is more severe when using non-minimal paths.

To analyze the extent of this variability, we consider the performance of a ping-pong benchmark between two nodes exchanging a 4MiB message, considering some of the routing strategies described in Section 2.2. To avoid situations where transient OS noise

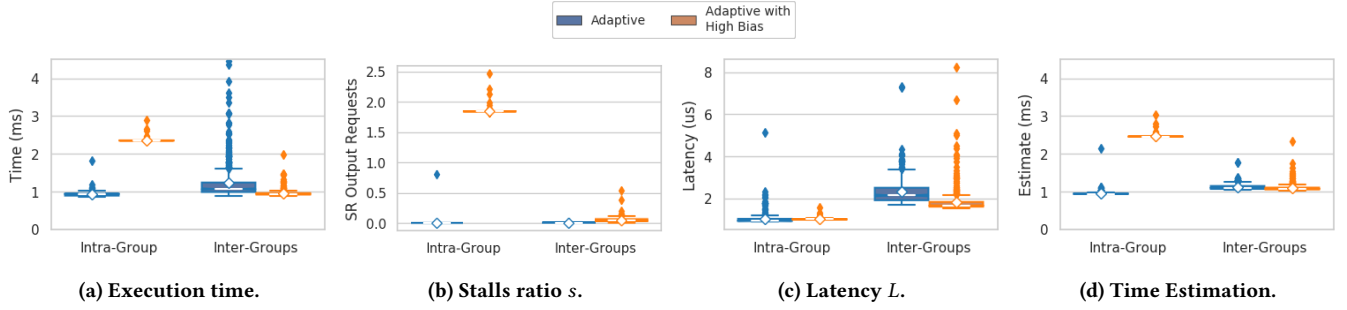


Figure 7: Performance comparison of the ping-pong benchmark, when using different routing algorithms. The test has been executed on the *Piz Daint* machine.

or network noise would affect a single routing strategy, we alternate the routing algorithm on successive ping-pongs.

We report in Figure 7 the performance of an iteration of the benchmark, when executed between two nodes in the same group (*Intra-Group*) and when executed between two nodes in different groups (*Inter-Groups*). We show both the execution time, the latency L , the stalls ratio s , and the time estimation according to the model shown in Equation 2. First of all, whereas for the *Intra-Group* case the time required to complete a ping-pong is lower when using ADAPTIVE routing, for the *Inter-Groups* case ADAPTIVE routing leads to an increase in network noise and to longer execution time.

Let us start analyzing the *Intra-Group* case. By analyzing the average stalled cycles s (Fig. 7b) we observe that packets sent with ADAPTIVE WITH HIGH BIAS routing are affected by more stalls than ADAPTIVE. Indeed, because the ADAPTIVE algorithm has a higher probability of selecting a non-minimal path, on average it will distribute the packets on a wider set of paths with respect to the ADAPTIVE WITH HIGH BIAS algorithm, thus decreasing the average stalls per flit. Because the two routing algorithms are characterized by a similar latency (Fig. 7c), the stalls determine the performance difference, as also estimated by our cost model (Fig. 7d).

On the other hand, in the *Inter-Groups* scenario, given the higher number of minimal paths connecting the two nodes, ADAPTIVE WITH HIGH BIAS routing algorithm can better distribute the packets, and the average number of stalled cycles decreases compared to the *Intra-Group* case (Fig. 7b). However, due to the higher number of hops between the two nodes, the latency of both routing algorithms increases (Fig. 7c). Moreover, ADAPTIVE routing is characterized by higher latency variations with respect to ADAPTIVE WITH HIGH BIAS because, due to *phantom congestion*, sometimes it may select a non-minimal path even if that was not necessary. Indeed, if the non-minimal path was selected to avoid actual congestion, we should see the effect of congestion on ADAPTIVE WITH HIGH BIAS in the form of higher average latency.

As a consequence, for the *Inter-Groups* case ADAPTIVE WITH HIGH BIAS performs better than the ADAPTIVE routing algorithm because it is characterized by a lower latency and a comparable number of stalled cycles with respect to ADAPTIVE. This evaluation clearly shows that a large part of the network noise can be attributed to the ADAPTIVE routing algorithm and that, under certain conditions, ADAPTIVE WITH HIGH BIAS may perform better

due to lower average latency. We will show in Section 5, when validating our application-aware routing algorithm, how several other microbenchmarks and real applications are also affected by the selection of the routing algorithm.

4.2 Noise-adaptive active routing

By leveraging these considerations, we can now devise an algorithm that uses information about latency and stalls to automatically change the routing algorithm according to the workload. As we have shown in Figure 7 the optimal choice does not only depend on the workload but also on its allocation, and for this reason we cannot derive any static solution to solve this problem. Moreover, by applying a static decision it would not be possible to react to transient changes in the network conditions, such as a temporary increase in the latency due to interfering jobs, or to intrinsic changes between application phases. For these reasons, we rely on a runtime approach which, after a message is sent, collects counters for latency and stalls. When sending a message, the algorithm will use the counters collected for the previous message to decide which routing algorithm should be used to send the current message.

To perform this decision, we assume that the application starts by using ADAPTIVE routing. Starting from Equation 2, we denote with L_{ad} and L_{bs} the latencies of the ADAPTIVE and ADAPTIVE WITH HIGH BIAS algorithms, respectively. Similarly, we denote the average stalls ratio with s_{ad} and s_{bs} . Then, when sending a message comprising f flits, the application would switch to ADAPTIVE WITH HIGH BIAS algorithm if:

$$\frac{p + 512}{1024} \cdot L_{bs} + f \cdot (s_{bs} + 1) < \frac{p + 512}{1024} \cdot L_{ad} + f \cdot (s_{ad} + 1) \quad (3)$$

i.e. if the message has a number of flits such that:

$$f < \frac{L_{ad} - L_{bs}}{s_{bs} - s_{ad}} \cdot \frac{p + 512}{1024} \quad (4)$$

Because we are assuming the application is already using the ADAPTIVE routing algorithm, we use the L_{ad} and s_{ad} monitored for the last message which was sent. L_{bs} and s_{bs} are instead estimated by multiplying L_{ad} and s_{ad} by appropriate scaling factors λ_{ad} and σ_{ad} , which we can derive by considering a median case over several runs of different microbenchmarks in different allocations. To correct mispredictions due to wrong choices of these scaling factors, the algorithm will store the last values observed for latency

and stalls for both routing algorithms. These values are discarded after a given number of samples, to avoid relying on data related to a different application phase. It is worth noting that this approach avoids the problems described in Sec. 3.2 and Sec. 3.3, because it only relies on NIC counters and is based on quantities which are independent from host-side delays.

Because reading the network counters for every message can introduce overhead on the application, we keep a cumulative counter of the message sizes, and we apply the algorithm when this counter is higher than a threshold (experimentally set to 4KiB). If the cumulative size is lower than the threshold, the message is sent with ADAPTIVE WITH HIGH BIAS routing. The reason behind that is that small messages are more affected by latency and ADAPTIVE WITH HIGH BIAS is usually characterized by a lower latency with respect to the ADAPTIVE algorithm. We decided to consider the cumulative size rather than only the size of the current message to avoid that applications which always send messages smaller than the threshold would never trigger the algorithm. To decide when to switch from ADAPTIVE WITH HIGH BIAS to ADAPTIVE, the dual equation of Equation 4 can be derived.

Algorithm 1: Application-Aware Routing

```

Function selectRouting msgSize
  if currentRouting == ADAPTIVE then
     $L_{ad} = L;$ 
     $s_{ad} = s;$ 
    if  $L_{bs}$  and  $s_{bs}$  too old then
       $L_{bs} = L_{ad} \cdot \lambda_{ad};$ 
       $s_{bs} = s_{ad} \cdot \sigma_{ad};$ 
    end
     $f = \text{getNumFlits}(\text{msgSize});$ 
    if  $f < \frac{L_{ad} - L_{bs}}{s_{bs} - s_{ad}} \cdot \frac{p+512}{1024}$  then
      currentRouting = ADAPTIVE WITH HIGH BIAS;
    else
      currentRouting = ADAPTIVE;
    end
  else
    // Similar to the other branch
  end
  return currentRouting;
end

```

The pseudocode of this decision process is shown in Algorithm 1. The function *selectRouting* is called before sending the message and selects the optimal routing algorithm. After sending the message, latency and stalls network counters are read and stored into the L and s variable. Counters are read after sending the message to do not introduce delays in the transmission. Moreover, for MPI_Alltoall communication, ADAPTIVE is replaced with INCREASINGLY MINIMAL BIAS routing because, as described in Section 2.2, this is also the current default behaviour.

4.3 Implementation details in Aries

Because the environment variables described in Section 2.2 can only be used to change the routing before running the application, we

had to adopt a different approach for changing the routing strategy message by message. High-performance communication libraries optimized for Cray Aries network, such as Cray’s MPICH, PGAS and others, rely on the uGNI and DMAPP APIs, which provide low-level communication services to user-space software. Before calling any of the uGNI and DMAPP functions to send messages over the Cray Aries network, it is possible to specify the routing algorithm to be used to send that message, usually by means of function parameters. However, this functionality is not exposed by higher-level APIs such as MPI. Accordingly, we implemented a dynamic library which defines the same functions used by uGNI to transmit on the network (with the same signatures). Inside each of these functions, the *selectRouting* function is called, the real uGNI function is invoked by specifying the chosen routing algorithm in the call and eventually, network counters are collected by using the PAPI library [35].

Every application which uses communication libraries which rely on uGNI (such as MPI), can then benefit from our application-aware routing by simply specifying our library in the LD_PRELOAD environment variable. A similar approach can be adopted for the DMAPP calls. This allows the algorithm to be applied transparently on almost any HPC application running on a Cray Aries platform. Moreover, in principle the algorithm could also be implemented on other networks relying on adaptive non-minimal routing, if appropriate mechanisms are provided for monitoring the network state and for changing the routing algorithm at runtime.

5 EXPERIMENTAL EVALUATION

In this section, we first analyze the performance of our application-aware adaptive routing on some microbenchmarks. Then, we will test it on some real applications.

5.1 Microbenchmarks

For the first part of our evaluation, we consider the following microbenchmarks.

ping-pong, allreduce, alltoall, barrier, broadcast These benchmarks use some common MPI calls. For ping-pong, alltoall and broadcast the size of the messages are expressed in bytes. Allreduce performs a sum reduction on an array of integers and the size of the messages is expressed as the number of elements of the array.

halo3d This benchmark performs a nearest neighbor communications, using a 3D stencil. Processes are logically arranged as a cube and each process communicates with the neighbors on the face of the cube formed around it as the center process. The input size corresponds to the size of the domain. For this benchmark, we used the implementation provided by the *ember* benchmark suite [3].

sweep3d This benchmark represents a wavefront communication pattern on a 3D grid. The pattern starts at a corner of the grid and “sweeps” out in a wavefront. The input size corresponds to the size of the domain. We use the implementation provided with the *ember* benchmark suite [3].

To avoid situations where transient OS noise or network noise would affect a single routing strategy, for each benchmark we alternate the routing algorithm on successive iterations. Because on

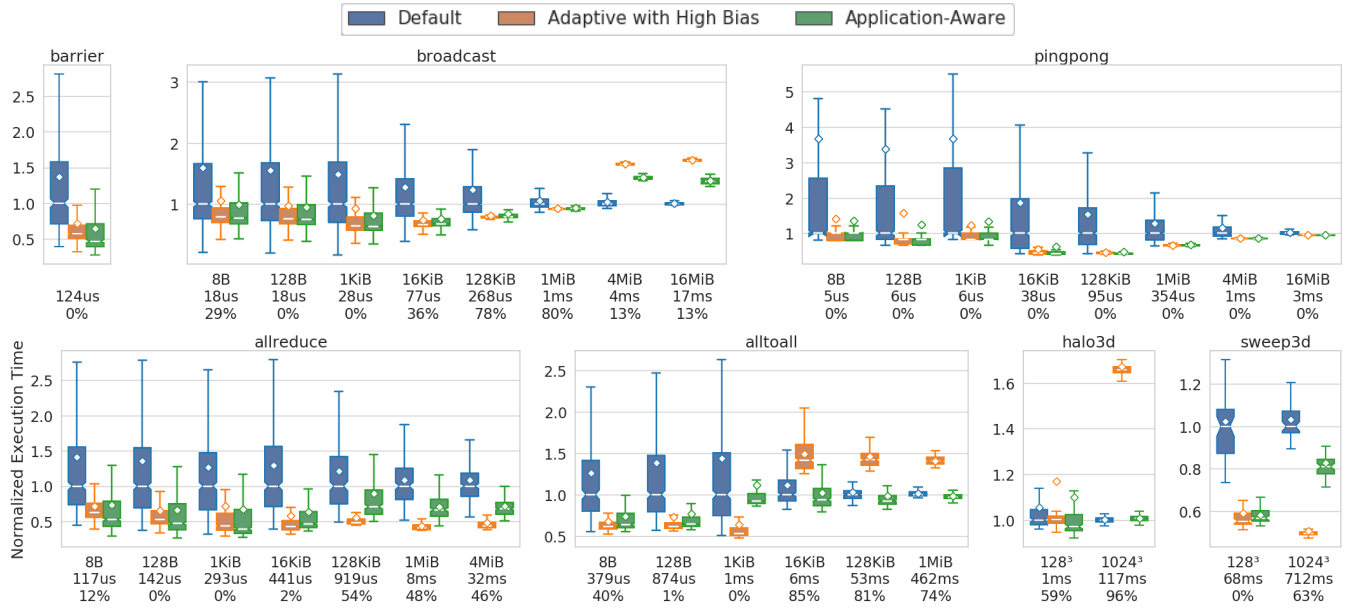


Figure 8: Execution time of microbenchmarks under DEFAULT, ADAPTIVE WITH HIGH BIAS and APPLICATION-AWARE routing algorithms, normalized with respect to the median of the ADAPTIVE algorithm, running on 1024 nodes on *Piz Daint*. For each test, we have on the x-axis: the input size, the median of the execution time for the DEFAULT routing, and the percentage of traffic which the APPLICATION-AWARE routing sends by using the DEFAULT routing.

Dragonfly networks the global available bandwidth depends on the number of nodes and links in the system, we verified that the system state (and thus the available bandwidth) did not change between different experiments. Moreover, to reduce the impact of resources contention we only execute one process per node.

We report in Figure 8 an overview of the performance of our algorithm across the microbenchmarks, for different input sizes, when executed on 1024 nodes on the *Piz Daint* machine. In this experiment, the job was allocated on 257 Aries routers spanning over 6 groups. We use a fixed allocation to avoid the problems described in Section 3.1. We show the execution time, normalized with respect to the median of the DEFAULT routing algorithm so that values lower than 1 represent a lower execution time compared to the DEFAULT algorithm. Because outliers can be some orders of magnitude higher than the median (as we saw in Section 3.1), although we still account them for computing means and medians, we do not show them on the plot to avoid shrinking the boxes too much. On the x-axis, we report for each test the input size, the median of the execution time for the DEFAULT routing strategy and the percentage of traffic that APPLICATION-AWARE routing sends by using DEFAULT routing. It is worth mentioning that in this specific case by reasoning in terms of execution time we are not affected by the problem described in Section 3.3. Indeed, if an application would be affected by OS noise or resources contention, this would affect all the routing algorithms in the same way, because the selection of the routing strategy does not have any host-side effect.

First, differences between DEFAULT and ADAPTIVE WITH HIGH BIAS are present across different communication patterns, and in some cases, by using a different routing algorithm we can reduce the

execution time by half. These differences depend from the communication pattern, the amount of data exchanged between the nodes and from their allocation. In general, whereas ADAPTIVE WITH HIGH BIAS performs better for benchmarks that do not generate much data (e.g., ping-pong or barrier), when the traffic intensity is higher (e.g., for alltoall, broadcast and halo3d), the DEFAULT routing algorithm outperforms ADAPTIVE WITH HIGH BIAS for large inputs. Besides, for some benchmarks the gap between the different routing algorithms is larger than for others, due to the different ability to *absorb* the noise. Namely, due to differences in communication and computation overlap, an increase of the average message latency may affect some benchmarks more than others.

Moreover, when using ADAPTIVE WITH HIGH BIAS we have a reduction in performance variability in almost all the cases, which clearly shows that a large part of the network noise we observe is due to the choices made by the routing algorithm rather than to actual congestion on the network. To further confirm this thesis, we can observe how for each benchmark, the performance variability of DEFAULT decreases when increasing the input size, due to the lower impact of latency on the transmission time of the messages. Our APPLICATION-AWARE routing selects most of the times the algorithm which provides the best performance among the two. For example, in the alltoall microbenchmark, it selects the ADAPTIVE WITH HIGH BIAS routing algorithm for small inputs whereas it selects the DEFAULT algorithm for larger message sizes.

However, despite it performs well in general, there are cases where it fails in selecting the best algorithm, such as for broadcast of large messages and for 1024³ sweep3d. By further investigating

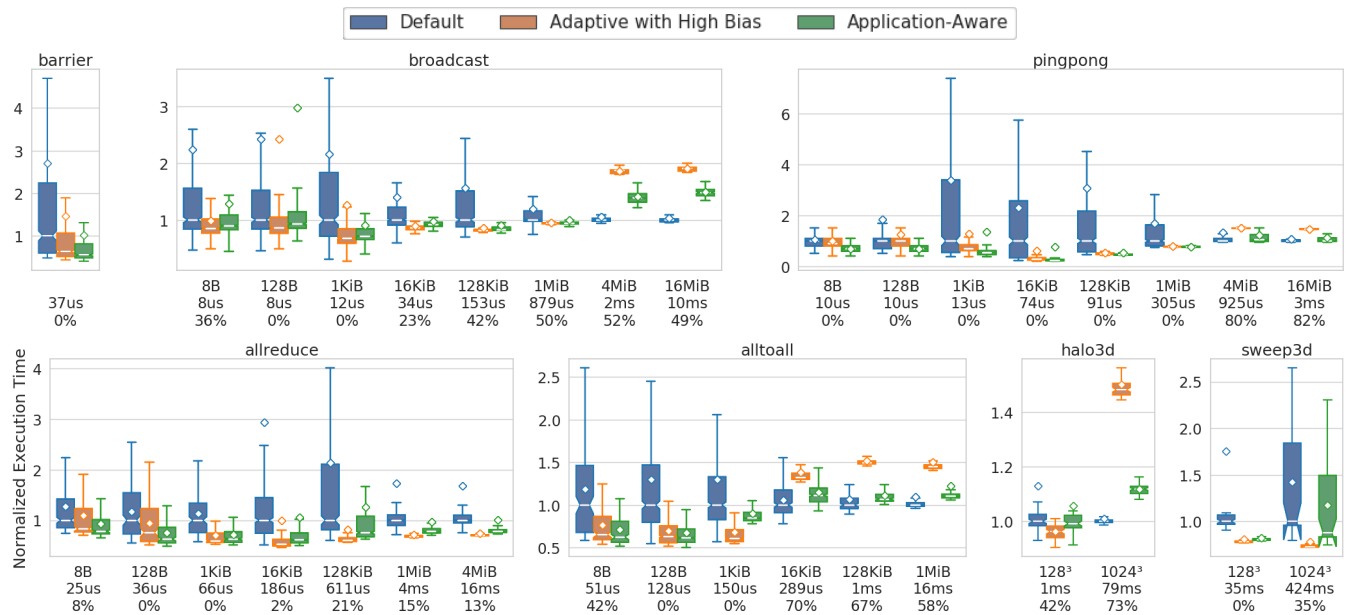


Figure 9: Execution time of microbenchmarks under DEFAULT, ADAPTIVE WITH HIGH BIAS and APPLICATION-AWARE routing algorithms, normalized with respect to the median of the ADAPTIVE algorithm, running on 64 nodes on Cori. For each test, we have on the x-axis: the input size, the median of the execution time for the DEFAULT routing, and the percentage of traffic which the APPLICATION-AWARE routing sends by using the DEFAULT routing.

the issue, we found out that this is caused by oscillations in the monitored stalls and latency, i.e., as soon as the algorithm detects that the DEFAULT algorithm should be used, the stalls start to decrease and it switches back to the ADAPTIVE WITH HIGH BIAS algorithm and the algorithm do not converge to the best routing algorithm. In other cases, even if APPLICATION-AWARE properly selects the optimal routing algorithm, we experience a performance drop with respect to the case where that routing strategy is statically set, such as for 1KiB alltoalls. This is due to the performance overhead introduced when reading the network counters and could be easily overcome by having more efficient access to network counters.

Eventually, in some cases, although DEFAULT routing provides better performance than ADAPTIVE WITH HIGH BIAS, APPLICATION-AWARE can achieve the same (or even better) performance than DEFAULT even by not sending the 100% of the traffic using such routing algorithm. By sending some data with ADAPTIVE WITH HIGH BIAS routing, less traffic is sent on non-minimal paths, reducing the network traffic. For example, this is the case of alltoalls larger than 16KiB. We performed the same set of experiments on 64 nodes on Cori, using 33 routers scattered on 5 Aries groups, obtaining similar results, as reported in Figure 9.

5.2 Applications

We now analyze our algorithm on the following applications:

CP2K Performs atomistic and molecular simulations of solid state, liquid, molecular, and biological systems [29].

WRF-B and WRF-T Simulations of a baroclinic wave and of a tropical cyclone performed with WRF, a weather prediction system designed for both atmospheric research and operational forecasting applications [41].

LAMMPS A molecular dynamics code that models an ensemble of particles in a liquid, solid, or gaseous state [37].

Quantum Espresso A suite for electronic-structure calculations and materials modeling at the nanoscale [21].

Nekbone Exposes the principal computation kernels of Nek5000, a fast and scalable high-order solver for computational fluid dynamics, used for many real-world applications [4, 44].

VPFET This application is an implementation of a mesoscale micromechanical materials model, which simulates the evolution of a material under deformation [2].

Amber A suite that allow users to carry out molecular dynamics simulations, particularly on biomolecules [40].

MILC/SU3_RMD The MILC benchmark represents part of a set of codes written to study quantum chromodynamics (QCD) by means of numerical simulations [8, 22].

HPCG Exercises computational patterns matching a wide set of applications, relying on operations like sparse triangular solvers and preconditioned conjugate gradient algorithms [15].

BFS and SSSP They perform, respectively, a breadth-first search and a single-source shortest path computation on a graph. We used the Graph500 reference implementation [36].

Fast Fourier Transform (FFT) It is a computation kernel which can be found in HPC applications across multiple domains. We used the benchmark provided by the `fftw` library [18].

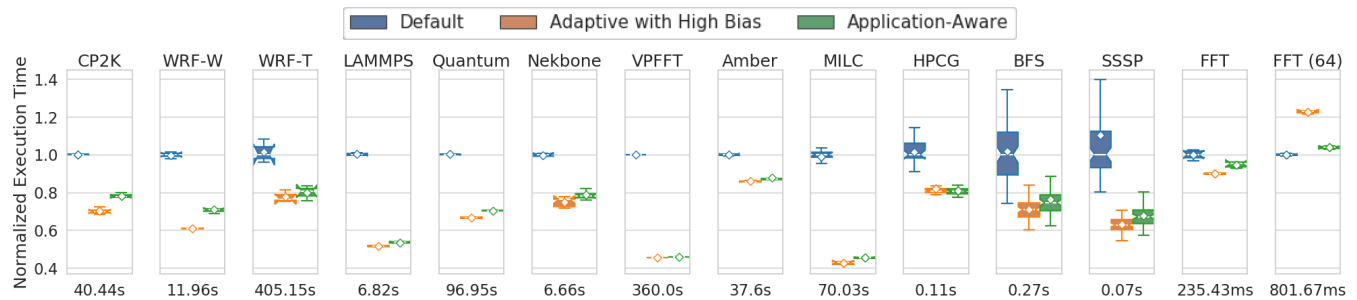


Figure 10: Execution time of applications under DEFAULT, ADAPTIVE WITH HIGH BIAS and APPLICATION-AWARE routing algorithms, normalized with respect to the median of the ADAPTIVE algorithm, running on 256 nodes on Piz Daint.

Because this work focuses on optimizing the communication phases, not all the applications can benefit from our approach, either due to their low communication intensity or to their good noise absorption capacity. For these reasons, on such applications we did not observe any difference between the different routing algorithms, and we excluded them from our analysis⁵. We report the results of our measurements in Fig. 10, when executed on 256 nodes on Piz Daint. We use in the plot the same notation used for the microbenchmarks. Even in this case, as for the microbenchmarks, a good selection of the routing algorithm can reduce the execution time of the application up to the 60%.

Moreover, as described in Section 4, the best algorithm to be used does not only depend on the workload but also on the allocation and the number of nodes used to run the application. In the rightmost part of Fig. 10 we report the execution time of the FFT application when using 64 nodes. While for the 256 nodes allocation ADAPTIVE WITH HIGH BIAS performs better than ADAPTIVE, for the 64 nodes allocation the opposite is true. As we can see from the results, our APPLICATION-AWARE algorithm selects the optimal routing strategy in both allocations, leading to performance similar to those of ADAPTIVE WITH HIGH BIAS for the 256 nodes allocation and comparable to those of DEFAULT for the 64 nodes allocation.. Interestingly, although both MILC and halo3d use a similar communication pattern, the optimal routing strategy is different in the two cases. Indeed, halo3d is a communication-oriented benchmark and does not perform any computation, generating much more traffic than MILC, thus taking advantage of the use of ADAPTIVE routing. This clearly shows how some decision which may seem optimal when stress-testing a communication network, could actually be sub-optimal for real applications. Overall, our APPLICATION-AWARE routing algorithm can select the optimal routing strategy across different applications and allocations, reducing the execution time by a factor of 2 on several applications.

6 DISCUSSION

Simulations. Although simulations would allow us to control different parameters and to analyze the impact and causation of noise in detail, we decided to do not rely on simulations for two main reasons. First, it is difficult to capture some aspects of the real network with a simulation, even because some information

⁵The full results are present in the artifact.

on how the network and the routing exactly work are not publicly available. Moreover, jobs are also affected by noise caused by other jobs in the system. To avoid introducing synthetically generated noise, which may not be representative of a realistic scenario, we decided to perform our analysis on a production system.

Static selection of the routing algorithm. Although the optimal routing algorithm depends on the size of the message to be sent and on its destination, performance differences also depend on the dynamic state of the network. For example, each rank in halo3d on a 1024^3 mesh generates, at each send, the same amount of data generated by a 64MiB ping-pong. However, even if each rank generates the same amount of data, in halo3d all the ranks are communicating simultaneously, generating more traffic than ping-pong (where there are only two ranks communicating). Indeed, whereas the former benefits from using the DEFAULT routing, the latter shows better performance when using ADAPTIVE WITH HIGH BIAS routing. Our algorithm captures the network state by observing the average latency and the stalls monitored through network counters.

System state. As mentioned in Section 2.2, traffic generated by other jobs may cross the routers used by the analyzed job. On Dragonfly networks it is not possible to isolate a job from the others because, even if we would allocate entire groups to our jobs, due to adaptive non-minimal routing packets may still traverse those groups. To mitigate these transient effects, each test has been run multiple times, and the 95% confidence interval from the median (which in most cases was lower than the 5% of the median) has been reported in the plots. Moreover, we alternated the routing algorithm on successive iterations to avoid having persistent noise affecting always the same routing algorithm. Furthermore, running the experiments on a production system without any kind of isolation allowed us to assess our algorithm in a realistic scenario where multiple and different jobs shared the network with our jobs.

Limitations. As shown in Section 5, due to its simplicity, the algorithm has some limitations and in some cases it does not select the absolute best routing strategy. Some of these problems could be mitigated by introducing some complexity in the algorithm. However, we wanted to keep the algorithm as simple as possible because, as we shown in our evaluation, the execution of the algorithm introduces some overhead. By striving for simplicity, we kept this overhead as low as possible, both in terms of time and memory.

7 RELATED WORK

Network noise. Different works recently investigated the impact of network noise for different network topologies, sometimes proposing solutions to mitigate this effect. Some studies quantified the effect of network noise on simple communication benchmarks based on MPI collective operations [13, 24, 27, 42]. However, besides being affected by some of the problems we described in Section 3, they do not propose any solution to mitigate network noise. The work by Chunduri et al. analyzes different sources of performance variability, including a brief analysis of the impact of the routing algorithm on MPI_Allreduce operations [13]. However, they do not analyze the reasons for such differences and do not exploit this information to mitigate network noise.

Most solutions optimize job allocation to minimize the contention on the links, either on dragonfly networks [39] or on other topologies [11, 38]. For example, Yang et al. [45, 47] show that whereas for communication intensive jobs a random allocation is more beneficial, for less communication-intensive jobs a contiguous allocation is better. Starting from this observation, they propose a hybrid allocation scheme, to allocate communication-intensive jobs randomly whereas the less communication-intensive jobs are allocated on contiguous nodes. However, due to adaptive non-minimal routing, it is not possible to fully isolate jobs on dragonfly networks.

Bhatele et al. [43] analyze the impact of network noise on both dragonfly and fat-tree networks, proposing an adaptive routing algorithm for fat-trees which, given the traffic matrix of the application, avoids hotspot by rerouting traffic on less loaded links. Eventually, some tools collect network counters across the entire network to provide visual information about congestion [10, 12, 23].

Adaptive routing. Other works analyze the limitations of adaptive routing and propose solutions to improve it. Won et al. [46] shown how the “far-end” congestion should be considered as *phantom* congestion, because it may be not properly represented with local information, such as the credit count. They propose a solution to overcome this problem and to avoid transient congestion, validating it by means of simulations. Other works address this problem [19], by proposing algorithms to improve congestion estimation, which are then simulated and compared to state of the art solutions. However, as also stated by the authors, approximations in the simulation are necessary because simulating the exact tiled structure of Dragonfly would be too costly. Similarly, Faizian et al. [16] propose a routing scheme which also considers traffic pattern information by using network counters. After analyzing the traffic pattern, the algorithm chooses a proper bias for the adaptive routing. However, this solution relies on counters which are not available on current networks. Eventually, Jain et al. [30] simulate different routing strategies and their interaction with job placement.

Our main difference with respect to these works is that whereas they are usually simulated and may require changes in the hardware infrastructure, in our case we are proposing a solution which is fully software-based and does not require any modification to existing hardware and software.

8 CONCLUSIONS

In this work we analyzed the impact of adaptive routing on network noise, proposing an application-aware routing algorithm to mitigate network noise and improve application performance. We first described how to measure variability caused by the interconnection network by isolating it from other sources of variability such as OS noise and resources contention. By following these guidelines, we shown that in some cases most of the network noise can be attributed to the adaptive routing algorithm and that noise can be reduced and performance improved by increasing the probability of selecting minimal paths.

By exploiting this knowledge we devised an application-aware routing algorithm which, before sending an application message, decides which algorithm should be used to route that message, based on information collected through network counters. Eventually, we validated this algorithm by comparing it with the ADAPTIVE, INCREASINGLY MINIMAL BIAS, and ADAPTIVE WITH HIGH BIAS algorithms provided in Cray Aries interconnection networks. We have shown how our algorithm is able to select, in most cases, the optimal routing strategy for different workloads on two different Cray machines, improving the performance up to a factor of 2 on both microbenchmarks and real applications.

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their insightful comments, CSCS for granting access to the Piz Daint machine, and NERSC for granting access to the Cori machine. We also thank Dr. Duncan Roweth for the useful discussions. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 programme (grant agreement DAPP, No. 678880). Daniele De Sensi is also supported by the EU H2020-ICT-2014-1 project *RePhrase* (No. 644235) and by the University of Pisa project PRA_2018_66 *DECLware: Declarative methodologies for designing and deploying applications*.

REFERENCES

- [1] [n. d.]. Aries Hardware Counters (4.0) S-0045. <https://pubs.cray.com/content/S-0045/4.0/aries-hardware-counters/about-ariestm-hardware-counters-s-0045>. ([n. d.]). Accessed: 17-05-2019.
- [2] [n. d.]. Crystal viscoplasticity proxy application. <https://github.com/exmatex/VPFFT>. ([n. d.]). Accessed: 17-05-2019.
- [3] [n. d.]. Ember Communication Pattern Library. <https://github.com/sstsimulator/ember>. ([n. d.]). Accessed: 10-04-2019.
- [4] [n. d.]. Nekbone. <https://github.com/Nek5000/Nekbone>. ([n. d.]). Accessed: 17-05-2019.
- [5] Mohammad Al-Fares, Alexander Loukissas, and Amin Vahdat. 2008. A Scalable, Commodity Data Center Network Architecture. *SIGCOMM Comput. Commun. Rev.* 38, 4 (Aug. 2008), 63–74. <https://doi.org/10.1145/1402946.1402967>
- [6] Bob Alverson, Edwin Froese, Larry Kaplan, and Duncan Roweth. 2012. Cray XC series network. *Cray Inc., White Paper WP-Aries01-1112* (2012).
- [7] Abdulla Bataineh, Thomas Court, and Duncan Roweth. 2017. Increasingly minimal bias routing. (2 2017).
- [8] G. Bauer, S. Gottlieb, and T. Hoefler. 2012. Performance Modeling and Comparative Analysis of the MILC Lattice QCD Application su3 rmd. In *Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012)*. IEEE Computer Society, 652–659.
- [9] M. Besta and T. Hoefler. 2014. Slim Fly: A Cost Effective Low-Diameter Network Topology. In *SC '14: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 348–359. <https://doi.org/10.1109/SC.2014.34>
- [10] A. Bhatele, N. Jain, Y. Livnat, V. Pascucci, and P. Bremer. 2016. Analyzing Network Health and Congestion in Dragonfly-Based Supercomputers. In *2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. 93–102. <https://doi.org/10.1109/IPDPS.2016.123>

- [11] Abhinav Bhatele and Laxmikant V. Kalé. 2009. Quantifying Network Contention on Large Parallel Machines. *Parallel Processing Letters* 19, 04 (2009), 553–572. <https://doi.org/10.1142/S0129626409000419> arXiv:<https://doi.org/10.1142/S0129626409000419>
- [12] James M. Brandt, Edwin Froese, Ann C. Gentile, Larry Kaplan, Benjamin A. Allan, and Edward J. Walsh. 2016. Network Performance Counter Monitoring and Analysis on the Cray XC Platform. (5 2016).
- [13] Sudheer Chunduri, Kevin Harms, Scott Parker, Vitali Morozov, Samuel Oshin, Naveen Cherukuri, and Kalyan Kumaran. 2017. Run-to-run Variability on Xeon Phi Based Cray XC Systems. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '17)*. ACM, New York, NY, USA, Article 52, 13 pages. <https://doi.org/10.1145/3126908.3126926>
- [14] David E. Culler, Richard M. Karp, David Patterson, Abhijit Sahay, Eunice E. Santos, Klaus Erik Schauer, Ramesh Subramonian, and Thorsten von Eicken. 1996. LogP: A Practical Model of Parallel Computation. *Commun. ACM* 39, 11 (Nov. 1996), 78–85. <https://doi.org/10.1145/240455.240477>
- [15] Jack Dongarra, Michael Heroux, and Piotr Luszczek. 2015. *HPCG Benchmark: a New Metric for Ranking High Performance Computing Systems*. Technical Report ut-eecs-15-736. <http://www.eecs.utk.edu/resources/library/file/1047/ut-eecs-15-736.pdf>
- [16] P. Faizian, M. S. Rahman, M. A. Mollah, X. Yuan, S. Pakin, and M. Lang. 2016. Traffic Pattern-Based Adaptive Routing for Intra-Group Communication in Dragonfly Networks. In *2016 IEEE 24th Annual Symposium on High-Performance Interconnects (HOTI)*. 19–26. <https://doi.org/10.1109/HOTI.2016.017>
- [17] Ahmad Faraj, Pitch Patarasuk, and Xin Yuan. 2008. A Study of Process Arrival Patterns for MPI Collective Operations. *Int. J. Parallel Program.* 36, 6 (Dec. 2008), 543–570. <https://doi.org/10.1007/s10766-008-0070-9>
- [18] M. Frigo and S. G. Johnson. 2005. The Design and Implementation of FFTW3. *Proc. IEEE* 93, 2 (Feb 2005), 216–231. <https://doi.org/10.1109/JPROC.2004.840301>
- [19] P. Fuentes, E. Vallejo, M. Garcia, R. Bevide, G. Rodríguez, C. Minkenberg, and M. Valero. 2015. Contention-Based Nonminimal Adaptive Routing in High-Radix Networks. In *2015 IEEE International Parallel and Distributed Processing Symposium*. 103–112. <https://doi.org/10.1109/IPDPS.2015.78>
- [20] Mark Giampapa, Thomas Gooding, Todd Inglett, and Robert W. Wisniewski. 2010. Experiences with a Lightweight Supercomputer Kernel: Lessons Learned from Blue Gene's CNK. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis (SC '10)*. IEEE Computer Society, Washington, DC, USA, 1–10. <https://doi.org/10.1109/SC.2010.22>
- [21] Paolo Giannozzi, Stefano Baroni, Nicola Bonini, Matteo Calandra, Roberto Car, Carlo Cavazzoni, Davide Ceresoli, Guido L. Chiarotti, Matteo Cococcioni, Ismaila Dabo, Andrea Dal Corso, Stefano de Gironcoli, Stefano Fabris, Guido Fratesi, Ralph Gebauer, Uwe Gerstmann, Christos Gougousis, Anton Kokalj, Michele Lazzeri, Layla Martin-Samos, Nicola Marzari, Francesco Mauri, Riccardo Mazzarello, Stefano Paolini, Alfredo Pasquarello, Lorenzo Paulatto, Carlo Sbraccia, Sandro Scandolo, Gabriele Sclauzero, Ari P. Seitsonen, Alexander Smogunov, Paolo Umari, and Renata M. Wentzcovitch. 2009. QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials. *Journal of Physics: Condensed Matter* 21, 39 (sep 2009), 395502. <https://doi.org/10.1088/0953-8984/21/39/395502>
- [22] Steven Gottlieb, W. Liu, William D Toussaint, R. L. Renken, and R. L. Sugar. 1987. Hybrid-molecular-dynamics algorithms for the numerical simulation of quantum chromodynamics. *Physical review D: Particles and fields* 35, 8 (1987), 2531–2542. <https://doi.org/10.1103/PhysRevD.35.2531>
- [23] Ryan E. Grant, Kevin T. Pedretti, and Ann Gentile. 2015. Overtime: A Tool for Analyzing Performance Variation Due to Network Interference. In *Proceedings of the 3rd Workshop on Exascale MPI (ExaMPI '15)*. ACM, New York, NY, USA, Article 4, 10 pages. <https://doi.org/10.1145/2831129.2831133>
- [24] T. Groves, Y. Gu, and N. J. Wright. 2017. Understanding Performance Variability on the Aries Dragonfly Network. In *2017 IEEE International Conference on Cluster Computing (CLUSTER)*. 809–813. <https://doi.org/10.1109/CLUSTER.2017.76>
- [25] T. Hoefler, P. Kambadur, R. L. Graham, G. Shipman, and A. Lumsdaine. 2007. A Case for Standard Non-Blocking Collective Operations. In *Recent Advances in Parallel Virtual Machine and Message Passing Interface, EuroPVM/MPI 2007*, Vol. 4757. Springer, 125–134.
- [26] T. Hoefler, A. Lumsdaine, and W. Rehm. 2007. Implementation and Performance Analysis of Non-Blocking Collective Operations for MPI. In *Proceedings of the 2007 International Conference on High Performance Computing, Networking, Storage and Analysis, SC07*. IEEE Computer Society/ACM.
- [27] T. Hoefler, T. Schneider, and A. Lumsdaine. 2009. The impact of network noise at large-scale communication performance. In *2009 IEEE International Symposium on Parallel Distributed Processing*. 1–8. <https://doi.org/10.1109/IPDPS.2009.5161095>
- [28] Torsten Hoefler, Timo Schneider, and Andrew Lumsdaine. 2010. Characterizing the Influence of System Noise on Large-Scale Applications by Simulation. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis (SC '10)*. IEEE Computer Society, Washington, DC, USA, 1–11. <https://doi.org/10.1109/SC.2010.12>
- [29] Jurg Hutter, Marcella Iannuzzi, Florian Schifmann, and Joost Van-deVondele. 2014. cp2k: atomistic simulations of condensed matter systems. *Wiley Interdisciplinary Reviews: Computational Molecular Science* 4, 1 (2014), 15–25. <https://doi.org/10.1002/wcms.1159> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/wcms.1159>
- [30] N. Jain, A. Bhatele, X. Ni, N. J. Wright, and L. V. Kale. 2014. Maximizing Throughput on a Dragonfly Network. In *SC '14: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 336–347. <https://doi.org/10.1109/SC.2014.33>
- [31] Nan Jiang, John Kim, and William J. Dally. 2009. Indirect Adaptive Routing on Large Scale Interconnection Networks. *SIGARCH Comput. Archit. News* 37, 3 (June 2009), 220–231. <https://doi.org/10.1145/1555815.1555783>
- [32] G. Kathareios, C. Minkenberg, B. Prisacari, G. Rodriguez, and T. Hoefler. 2015. Cost-Effective Diameter-Two Topologies: Analysis and Evaluation. ACM. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC15)*.
- [33] G. Kathareios, C. Minkenberg, B. Prisacari, G. Rodriguez, and T. Hoefler. 2015. Cost-effective diameter-two topologies: analysis and evaluation. In *SC '15: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 1–11. <https://doi.org/10.1145/2807591.2807652>
- [34] J. Kim, W. J. Dally, S. Scott, and D. Abts. 2008. Technology-Driven, Highly-Scalable Dragonfly Topology. In *2008 International Symposium on Computer Architecture*. 77–88. <https://doi.org/10.1109/ISCA.2008.19>
- [35] Philip J. Mucci, Shirley Browne, Christine Deane, and George Ho. 1999. PAPI: A Portable Interface to Hardware Performance Counters. In *In Proceedings of the Department of Defense HPCMP Users Group Conference*. 7–10.
- [36] Richard C Murphy, Kyle B Wheeler, Brian W Barrett, and James A Ang. 2010. Introducing the graph 500. *Cray Users Group (CUG) 19* (2010), 45–74.
- [37] Steve Plimpton. 1995. Fast Parallel Algorithms for Short-range Molecular Dynamics. *J. Comput. Phys.* 117, 1 (March 1995), 1–19. <https://doi.org/10.1006/jcph.1995.1039>
- [38] Samuel D. Pollard, Nikhil Jain, Stephen Herbein, and Abhinav Bhatele. 2018. Evaluation of an Interference-free Node Allocation Policy on Fat-tree Clusters. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis (SC '18)*. IEEE Press, Piscataway, NJ, USA, Article 26, 13 pages. <http://dl.acm.org/citation.cfm?id=3291656.3291691>
- [39] Bogdan Prisacari, German Rodriguez, Philip Heidelberger, Dong Chen, Cyriel Minkenberg, and Torsten Hoefler. 2014. Efficient Task Placement and Routing of Nearest Neighbor Exchanges in Dragonfly Networks. In *Proceedings of the 23rd International Symposium on High-performance Parallel and Distributed Computing (HPDC '14)*. ACM, New York, NY, USA, 129–140. <https://doi.org/10.1145/2600212.2600225>
- [40] Romelia Salomon-Ferrer, David A. Case, and Ross C. Walker. 2013. An overview of the Amber biomolecular simulation package. *Wiley Interdisciplinary Reviews: Computational Molecular Science* 3, 2 (2013), 198–210. <https://doi.org/10.1002/wcms.1121> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/wcms.1121>
- [41] William C. Skamarock, Joseph B. Klemp, Jimmy Dudhia, David O. Gill, Dale M. Barker, Wei Wang, and Jordan G. Powers. 2008. A description of the Advanced Research WRF version 3. NCAR Technical note -475+STR. (2008).
- [42] D. Skinner and W. Kramer. 2005. Understanding the causes of performance variability in HPC workloads. In *IEEE International. 2005 Proceedings of the IEEE Workload Characterization Symposium, 2005*. 137–149. <https://doi.org/10.1109/IISWC.2005.1526010>
- [43] Staci A. Smith, Clara E. Cromey, David K. Lowenthal, Jens Domke, Nikhil Jain, Jayaraman J. Thiagarajan, and Abhinav Bhatele. 2018. Mitigating Inter-Job Interference Using Adaptive Flow-Aware Routing. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '18)*.
- [44] H. M. Tufo and P. F. Fischer. 1999. Terascale Spectral Element Algorithms and Implementations. In *Proceedings of the 1999 ACM/IEEE Conference on Supercomputing (SC '99)*. ACM, New York, NY, USA, Article 68. <https://doi.org/10.1145/331532.331599>
- [45] X. Wang, M. Mubarak, X. Yang, R. B. Ross, and Z. Lan. 2018. Trade-Off Study of Localizing Communication and Balancing Network Traffic on a Dragonfly System. In *2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. 1113–1122. <https://doi.org/10.1109/IPDPS.2018.00120>
- [46] J. Won, G. Kim, J. Kim, T. Jiang, M. Parker, and S. Scott. 2015. Overcoming far-end congestion in large-scale networks. In *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*. 415–427. <https://doi.org/10.1109/HPCA.2015.7056051>
- [47] X. Yang, J. Jenkins, M. Mubarak, R. B. Ross, and Z. Lan. 2016. Watch Out for the Bully! Job Interference Study on Dragonfly Network. In *SC '16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 750–760. <https://doi.org/10.1109/SC.2016.63>