

Studying the Influence of the InfiniBand Packet Size to Guarantee QoS*

Francisco J. Alfaro, José L. Sánchez
Dept. de Informática
Escuela Politécnica Superior
Universidad de Castilla-La Mancha
02071- Albacete, Spain
{falfaro, jsanchez}@info-ab.uclm.es

José Duato
Dept. de Informática de
Sistemas y Computadores
Universidad Politécnica de Valencia
46071- Valencia, Spain
jduato@gap.upv.es

Abstract

InfiniBand (IBA) has been proposed as an industry-standard architecture both for I/O server and interprocessor communication. IBA employs a switched point-to-point network, instead of using a shared bus. IBA is being developed by the InfiniBandSM Trade Association to provide present and future server systems with the required levels of reliability, availability, performance, scalability, and quality of service (QoS).

In previous papers we have proposed an effective strategy for configuring the IBA networks to provide users with the required levels of QoS. This strategy is based on the proper configuration of the mechanisms IBA carries to support QoS. Specifically, our methodology configures the InfiniBand Arbitration Tables and uses the different Service Levels and Virtual Lanes that are available, in order to segregate the different traffic flows. Thus, each flow receives the treatment it has previously requested. Moreover, by using our methodology, applications can be assured that their requirements will be satisfied.

In this paper, we review the basis of our methodology and we study the influence of the packet size on the QoS guaranteed to the applications.

1. Introduction

The InfiniBand Trade Association (IBTA) [8] was formed in 1999 to develop a new standard for high-speed I/O and interprocessor communication. InfiniBand defines a technology for interconnecting proces-

*This work was partly supported by the Spanish CICYT under Grant TIC2003-08154-C06 and Junta de Comunidades de Castilla-La Mancha under Grant PBC-02-008.

sor nodes (hosts) and I/O devices to form a system area network [7]. In a first stage, instead of directly replacing the PCI bus with a switch-based interconnection to access I/O devices, these devices are attached to a Host Channel Adapter (HCA), which is connected to the PCI bus. In this way, the communication is switched from HCA, affording the desired reliability, concurrency and security. Moreover, it is foreseen [1] that the PCI bus could be replaced in a near future by other advanced technologies like PCI Express Advanced Switching [2].

InfiniBand has been developed to provide present and future server systems with a cost-effective, high-performance solution with reliability, availability, and serviceability support. InfiniBand implements mechanisms to provide each kind of application with the required QoS. In previous works, [3] and [6], we have developed a methodology to configure such mechanisms. The proposed methodology successfully provides applications with both bandwidth and latency guarantees. In this paper, we study the influence of the packet size on the QoS guaranteed to the applications.

The structure of the paper is as follows: Section 2 presents a summary of the most important mechanisms included in IBA to support QoS; in Section 3, we review our proposal to give QoS guarantees; Section 4 presents the evaluation methodology used to study the influence of the packet size on the QoS guaranteed to the applications, as well as the obtained results; finally, some conclusions are given in Section 5.

2. InfiniBand

InfiniBand hardware provides highly reliable, fault-tolerant communication, improving the bandwidth, latency, and reliability of the system. The InfiniBand

architecture simplifies and speeds server-to-server connections and links to other server-related systems, such as remote storage and networking devices, through a message-based fabric network.

Specifically, IBA has three main mechanisms to support QoS: Service levels (SLs), virtual lanes (VLs), and a virtual lane arbitration for transmission over links. IBA defines a maximum of 16 SLs, although it does not specify which characteristics the traffic of each service level should have. Therefore, the distribution of the different existing traffic types among the SLs may be stated by the manufacturer or the network administrator. By allowing the traffic to be segregated by categories, we will be able to distinguish between packets from different SLs and to give them a different treatment according to their needs.

IBA ports support VLs as a mechanism for creating multiple virtual links within a single physical link. Each VL must be an independent resource for flow control purposes. A VL represents a set of transmission and reception buffers in a port. IBA ports can support a minimum of two and a maximum of 16 VLs. Since systems can be constructed with switches supporting a different number of VLs, the number of VLs used by a port is configured by the subnet manager. Moreover, packets are marked with a SL, and a relation between SL and VL is established at the input of each link by means of the table *SLtoVLMappingTable*.

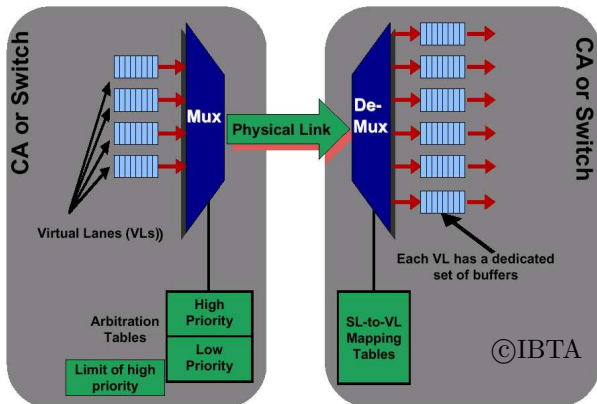


Figure 1. Operation of virtual lanes in a physical link.

When more than two VLs are implemented, the priorities of the data lanes are defined by the *VLArbitrationTable*. This arbitration affects only data VLs, because control traffic uses its own VL, which has greater priority than any other VL. The *VLArbitrationTable* consists of two tables, one for scheduling packets from high-priority VLs and another for low-priority VLs. However, IBA does not specify what is high and low-

priority. The arbitration tables implement weighted round-robin arbitration within each priority level. Up to 64 table entries are cycled through, each one specifying a VL and a weight, which is the number of units of 64 bytes to be transmitted from that VL. This weight must be in the range from 0 to 255, and is always rounded up as a whole packet.

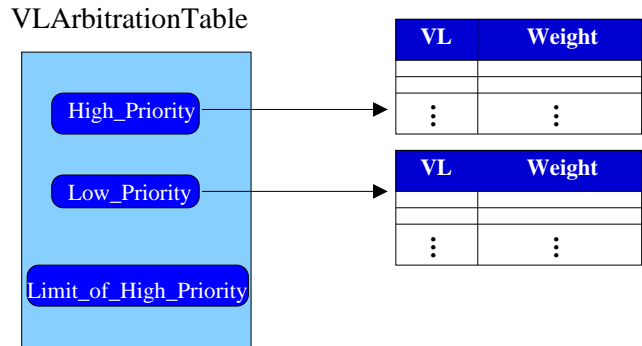


Figure 2. Structure of the VLArbitrationTable.

Moreover, a *LimitOfHighPriority* value specifies the maximum number of high-priority packets that can be sent before a low-priority packet is transmitted. Specifically, the VLs of the high-priority table can send *LimitOfHighPriority* × 4096 bytes before a packet from the low-priority table can be transmitted. If at a given time, no high-priority packets are ready for transmission, low-priority packets can also be transmitted.

3. Our proposal to give QoS guarantees

In previous works [3, 5] we have proposed a simple strategy to treat the requests of latency guarantee. Specifically, when an application requests latency guarantee, the maximum distance allowed between two consecutive entries in the high-priority table must be computed in order to allocate entries on that table to the application in question. Moreover, the application could also request a mean bandwidth that would result in a weight to be put in the entries of the arbitration table. Therefore, for a certain connection that requests a maximum delay that results in a distance d , and a mean bandwidth that results in a weight w , the number of entries needed is $\max\{\frac{64}{d}, \frac{w}{255}\}$.

Obviously, the maximum distance between two consecutive entries in the high-priority table requested by a connection ranges from 1 to 64. However, in order to optimize the filling up of the table, we only consider the following distances: 2, 4, 8, 16, 32, and 64 [4]. There-

fore, the applications' requests of a maximum distance between two consecutive entries in the high-priority table are turned into the closest lower power of 2 [6].

Traffic is grouped in SLs according to its maximum latency. Specifically, all connections using the same SL need the same maximum distance between two consecutive entries in the high-priority table, regardless of their mean bandwidth. For the most requested distances, we could distinguish between two or four different SLs considering the mean bandwidth. In this way, if we have enough available VLs, each kind of traffic could use a different VL.

Moreover, in [6] and [4] we proposed an algorithm to select a free sequence of entries in the high-priority table to meet a new application's request. This algorithm successfully allocates a new sequence in the table if there are enough available entries. This is achieved because the available entries are always in the best situation to treat the most restrictive request. For a connection requesting $\frac{64}{d}$ entries with a maximum distance d between them, the algorithm looks for a previously established sequence for the corresponding VL, with enough available weight. If there is no available sequence, a new free sequence with the desired characteristics is looked for.

In a more formal way, in a table T , let the sequence $t_0, t_1, \dots, t_{62}, t_{63}$ represents the entries of the table. Each t_i has an associated weight w_i whose value can vary between 0 and 255. Thus, we say an entry t_i is *free* if and only if $w_i = 0$. For a table T and a request of distance $d = 2^i$, we define the sets $E_{i,j}$ with $i = \log_2 d$ and $0 \leq j < d$, as

$$E_{i,j} = \left\{ t_{j+n \times 2^i} \quad ; \quad n = 0, \dots, \frac{64}{2^i} - 1 \right\}$$

Each $E_{i,j}$ contains the entries of the table T spaced by an equal distance d which are able to meet a request of distance $d = 2^i$ starting with the entry t_j . We say a set $E_{i,j}$ is free if $\forall t_k \in E_{i,j}$, t_k is free. Other properties derived from this definition are available in [4].

In [6] we also presented a simple algorithm to maximize the number of requests allocated in the arbitration table. For a new request of distance $d = 2^i$, the algorithm studies all possible sets $E_{i,j}$ for this kind of request, in a certain order, and selects the first set that is free (so, all its entries are free). The order the sets are inspected in is based on the application of the bit-reversal permutation to the distance values in the interval $[0, d-1]$. Specifically, for a new request of maximum distance $d = 2^i$, the algorithm selects the first free $E_{i,j}$ in the sequence $E_{i,iR_0}, E_{i,iR_1}, \dots, E_{i,iR_{d-1}}$ where iR_j is the bit-reversal function applied to j codified with i bits. Note that this algorithm is only applied if

there is no previously allocated sequence for the same requested distance with available room in its entries.

For example, the order the sets are inspected for a request of distance $d = 2^3$ is $E_{3,0}, E_{3,4}, E_{3,2}, E_{3,6}, E_{3,1}, E_{3,5}, E_{3,3}$, and $E_{3,7}$. Note that this algorithm first fills in the even entries, and later the odd entries. In this way, if we have available entries, we can always meet a request of distance 2, which is the most restrictive. The same consideration can be made for longer distances.

In [4], we have also proved several theorems showing that the algorithm can always allocate a new request if there are enough available entries. This is achieved because the algorithm always selects the sequences in the optimal way for satisfying later the most restrictive possible request.

When a connection finishes, its bandwidth is deducted from the accumulated bandwidth in the entries it was occupying. When this accumulated bandwidth is zero those entries must be released. When some entries are released, a disfragmentation algorithm must be applied to leave the table in a correct state, such that the proposed filling in algorithm can be used. This disfragmentation algorithm and its properties are also described in [4]. Basically, it puts together small free sets to form a larger free set, moving the content of some entries.

Both algorithms together permit the allocation and release of sequences of entries in the arbitration table in a optimal and dynamical way [4]. This allows us to provide applications with QoS using the IBA mechanisms in an optimal way.

4. Performance evaluation

In [3] and [6] we have evaluated the behavior of our proposals using a fixed packet size. We have shown that our proposals are able to provide applications with QoS guarantee. In this section, we are going to determine if the packet size has influence on the QoS guaranteed to the applications. In the following points, we explain the network and the traffic models we have used.

4.1. Network model

We have used irregular networks randomly generated. All switches have 8 ports, 4 of them having a host attached, the other 4 being used for interconnection between switches. We have evaluated networks with sizes ranging from 8 to 64 switches (thus, with 32 to 256 hosts, respectively). We have also tested several packet sizes ranging from 256 to 4096 bytes, and the three link rates specified in IBA. Taking into account

the space limitation and the fact that all these variations present similar results, we have only included here results for the 16 switches-network and a link rate of 2.5 Gbps.

In the switches both at input and output ports, there are 16 VLs so that a different VL may be assigned to each SL. Each switch has a multiplexed crossbar. As IBA uses virtual cut-through we have only considered buffer sizes that allow to store completely whole packets. Specifically, we have considered buffer sizes of 4 whole packets of capacity, which is a usual value for this parameter in this type of study.

InfiniBand specifies that the packet size must be between 256 bytes and 4096 bytes. Therefore, we have tested several packet sizes in this range. Specifically, we have considered packets of size 256 bytes, 1024 bytes, 2048 bytes, and 4096 bytes.

4.2. Traffic model

We have used 10 SLs for traffic needing QoS. Each SL presents a different maximum distance and bandwidth requirements. We have used CBR traffic, randomly generated among the bandwidth range of each SL. For the most requested distances several SLs have been considered using the mean bandwidth of the connections. Specifically, the SLs used and their features are shown in Table 1.

Table 1. Features of the SLs used.

SL	Maximum Distance	Bandwidth Range (Mbps)
0	2	0.064 - 1.55
1	4	0.064 - 1.55
2	8	0.064 - 1.55
3	16	0.064 - 1.55
4	32	0.064 - 1.55
5		1.55 - 64
6	64	0.008 - 0.064
7		0.064 - 1.55
8		1.55 - 64
9		64 - 255

The connections of each SL request a maximum distance between two consecutive entries in the high-priority table and a mean bandwidth in the range shown in Table 1. Note that this is similar to requesting a maximum deadline and computing the maximum distance between two consecutive entries in the virtual lane arbitration tables.

Each request is considered by each node along its path and is accepted only if there are available resources. Connections of the same SL are grouped into the same sequence of entries in the arbitration table.

The total weight of the sequence is computed according to the accumulated bandwidth of the connections sharing that sequence. When the connection cannot be settled in a previously established sequence (or there is no previous one), our algorithm looks in the high-priority arbitration table for a new free sequence of entries with the correct distance between its entries.

When no more connections can be established, we start a transient period in order for the network to reach a stationary state. Once the transient period finishes, the steady state period begins, where we gather results to be shown. The steady state period continues until the connection with a smaller mean bandwidth has received 100 packets.

Although the results for best-effort traffic are not the main focus of this paper, we have reserved 20% of available bandwidth for this type of traffic, which would be served by the low-priority table. So, connections would only be established in up to 80% of the available bandwidth.

4.3. Simulation results

We can see in Table 2 several metrics measured for the different packet sizes considered. Specifically, we have computed the injected and delivered traffic (in bytes/cycle/node), the average utilization (in %), the average bandwidth reserved (in Mbps) in host interfaces and switch ports, and the number of established connections. Note that the maximum utilization reachable is 80%, because the other 20% is reserved for BE and CH traffic. We are therefore close to the maximum utilization achievable. Obviously, we could achieve a higher utilization establishing more connections, but we have already made many attempts for each SL. We could establish other connections, but these connections would be of SLs with a small mean bandwidth on account of the heavy network load, and it is unlikely that these new connections would provide us with more information. So, it seems reasonable to assume that with this load we can study the network behavior in a quasi-fully loaded scenario.

Note that achieved utilization grows for larger packet size. This is because with a larger packet size more connections are established. When the packet size increases then the overhead that the packet header means decreases. For small packet size this overhead is more important than for larger packet sizes. Therefore, for larger packet sizes more connections can be established, and so the network achieves a little higher utilization and load. The same happens for the other evaluated indexes. Note also that in all cases the network transmits all the packets it receives, and therefore

Table 2. Traffic and utilization for different packet sizes.

	Packet size (in bytes)			
	256	1024	2048	4096
Injected traffic (Bytes/Cycle/Node)	0,7258	0,7432	0,7507	0,7607
Delivered traffic (Bytes/Cycle/Node)	0,7258	0,7432	0,7507	0,7607
Av. utilization for host interfaces (%)	72,58	74,32	75,07	76,07
Av. utilization for switch ports (%)	73,48	75,50	75,04	75,13
Av. reservation for host interfaces (Mbps)	1848,67	1867,80	1882,44	1905,32
Av. reservation for switch ports (Mbps)	1871,75	1897,58	1881,69	1881,63
Established connections	111813	113233	115779	118938

delivered traffic is equal to injected traffic.

We have also computed the percentages of packets that meet a certain deadline threshold. These thresholds are different for each connection and are related to their requested maximum deadline. This maximum deadline is the maximum delay that has been guaranteed to each connection. In the figures, this deadline is referred to as D . The results for each SL are presented in Figure 3 for the considered packet sizes. In these figures, we can see that all packets of all SLs arrive at their destinations before their deadlines. However, for larger packet size, the packets take more time to arrive

at their targets. This is because the network is more loaded, as previously stated. There is more traffic in the network and so there is somewhat more contention. However, in all the cases, all packets of all SLs arrive at their destinations before their deadlines expire.

Moreover, as the thresholds are different for each connection (they are related to their maximum deadline), they are also different for each service level. Therefore, the packets belonging to stricter service levels arrive at their targets closer to their deadline expiry, as this is far smaller than for other service levels.

We have also measured the average packet jitter. We have computed the percentage of packets received in several intervals related to their inter-arrival time (IAT). Obviously, these intervals are different for each connection. The results for each SL are shown in Figure 4, for the considered packet sizes. We can see that for the strictest SLs (SLs 0, 1, 2, 3, and 4) all packets arrive at their targets in the central interval, except for packet size of 4096 bytes (Figure 4(g)). In this case we can observe that some packets have arrived at their targets in other intervals. As stated, this is because the network is a more loaded in this case and so some contention appears.

For the other SLs with bigger mean bandwidth (SLs 5, 6, 7, 8, and 9) the jitter has a Gaussian distribution. In these cases, the majority of the packets arrive at their targets in the central intervals, the number decreasing as we move from the central interval. This is because these SLs are lower priority, and also because some of them have larger bandwidth ranges, and so their IAT and intervals are smaller. For example, the connections belonging to SL 9 have bandwidth ranging from 64 Mbps and 255 Mbps, and so their IAT is far smaller than for other SLs with bandwidth around 1 Mbps.

Moreover, it should be noted that there are no appreciable differences in the behavior of the jitter for the different packet sizes considered. According to the obtained results there are other parameters (bandwidth,

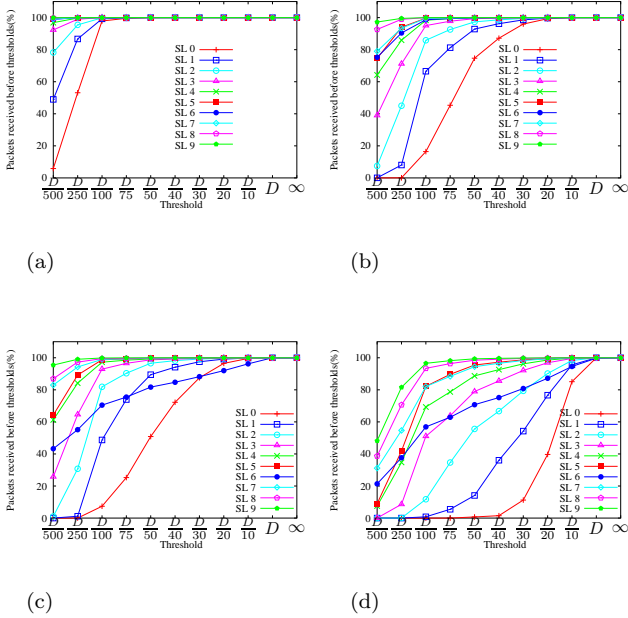


Figure 3. Packet delay for packet sizes of (a) 256 bytes, (b) 1024 bytes, (c) 2048 bytes, and (d) 4096 bytes.

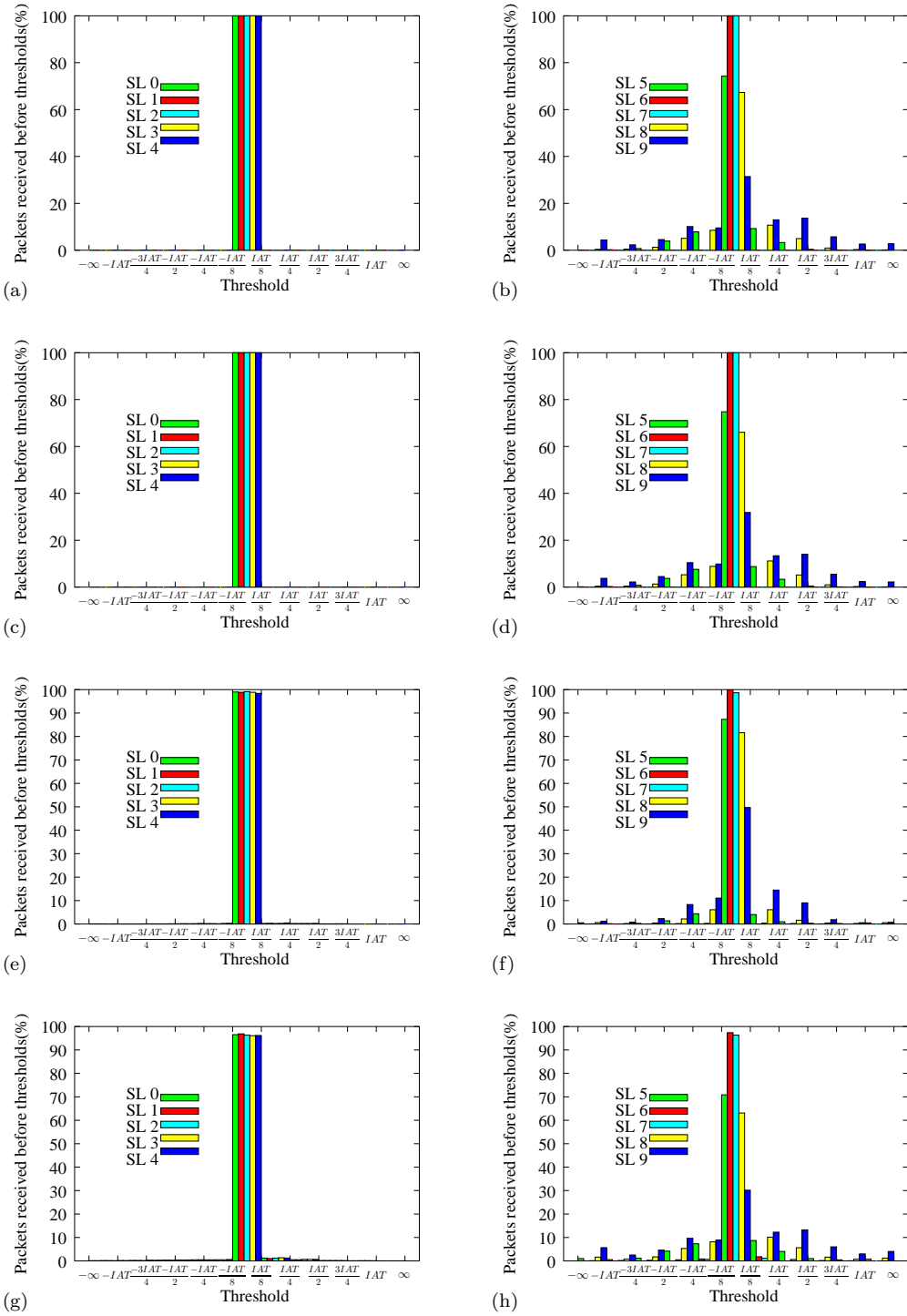


Figure 4. Packet jitter for packet sizes of (a) and (b) 256 bytes, (c) and (d) 1024 bytes, (e) and (f) 2048 bytes, (g) and (h) 4096 bytes.

priority of the SL, etc.) that have influence in the jitter, but all the packet sizes considered obtain similar

jitter results.

Finally, for a given deadline threshold, we have se-

lected the connections that deliver the lowest and the highest percentage of packets before this threshold. In the figures, these connections will be referred to as the worst and the best connections, respectively. We have selected a very tight threshold so that the percentage of packets meeting the deadline was lower than 100% in Figure 3. In particular, we have selected the threshold equal to $\frac{Deadline}{100}$. Note again that this threshold is different for each connection and depends on its own maximum deadline. Figure 5 shows the results for the considered packet sizes and for the SLs 0, 1, 2, and 3, which are the SLs with the highest deadline requirements. The results for other SLs are even better than these shown in the figures. It is noteworthy that, in all cases, even the packets of the worst connection arrive at their destination before their deadline.

In the same sense as for the average latency, when the packet size increases the considered best and worst connections deliver the majority of their packets later. However, in all cases, even the packets of the worst connection arrive at their destination before their deadline expires.

Note also that the difference between the connection considered the best and that considered the worst is very similar. In this sense, there is no important difference. This is very important since it means that the arbitration performed treats all the connections belonging to the same SL in a similar way.

5. Conclusions

In [3] and [6] we proposed a new methodology to provide each kind of application with the previously required QoS level. We also proposed an algorithm to select a free sequence of entries in the arbitration table. This algorithm successfully allocates a request in the arbitration table if there are enough available entries. It manages the requests in an optimal way, being able later to satisfy the most restrictive possible request. Some formal properties and theorems derived from this algorithm are shown in [4].

In this paper, we have studied the influence of the packet size on the QoS guaranteed to the applications. We have tested several packet sizes ranging from 256 bytes to 4096 bytes, which is the maximum packet size allowed in InfiniBand. The most important result may well be that for all packet sizes our proposals meet the QoS requirements.

According to the obtained results, packet size has a collateral influence in the QoS provided to the applications. The proposed methodology manages to provide applications with QoS guarantee according to the previously demanded requirements, regardless of the

packet size considered. However, for large packet sizes, where the packet header means a lower overload for the network, more connections can be established. This fact is very important for the obtained results because as we can establish more connections, a correspondingly higher utilization is achieved in the network, with the result that there is somewhat more contention. This affects the average latency that the applications present, which grows slightly when the packet size increases.

However, we can conclude that the packet size has no direct influence on the proposed methodology, because in all the cases the requirements demanded by the applications are achieved. Therefore, the resources are measured and scheduled in a correct way, and the performed arbitration is also correct.

References

- [1] *Advanced Bus and Interface Markets and Trends*. Electronic Trend Publications, 4th edition, Sept. 2004.
- [2] Advanced Switching Core Architecture Specification. Available at <http://www.asi-sig.org/specifications> for ASI SIG members only, 2003.
- [3] F. Alfaro, J. Sánchez, and J. Duato. QoS in InfiniBand Subnetworks. *IEEE Transactions on Parallel and Distributed Systems*, 15(9):194–205, Sept. 2004.
- [4] F. Alfaro, J. Sánchez, M. Mendiña, and J. Duato. Formalizing the Fill-In of the InfiniBand Arbitration Table. Technical Report DIAB-03-02-35, Dep. de Informática Universidad de Castilla-La Mancha, Mar. 2003.
- [5] F. J. Alfaro, J. L. Sánchez, and J. Duato. A Strategy to Manage Time Sensitive Traffic in InfiniBand. In *Proceedings of Workshop on Communication Architecture for Clusters (CAC'02)*, Apr. 2002. Held in conjunction with IPDPS'02, Fort Lauderdale, Florida.
- [6] F. J. Alfaro, J. L. Sánchez, and J. Duato. A New Proposal to Fill in the InfiniBand Arbitration Tables. In *Proceedings of IEEE International Conference on Parallel Computing (ICPP'03)*, pages 133 – 140, Oct. 2003.
- [7] InfiniBand Trade Association. *InfiniBand Architecture Specification Volume 1. Release 1.0*, Oct. 2000.
- [8] InfiniBandTM Trade Association. 1999. <http://infinibandta.com>.

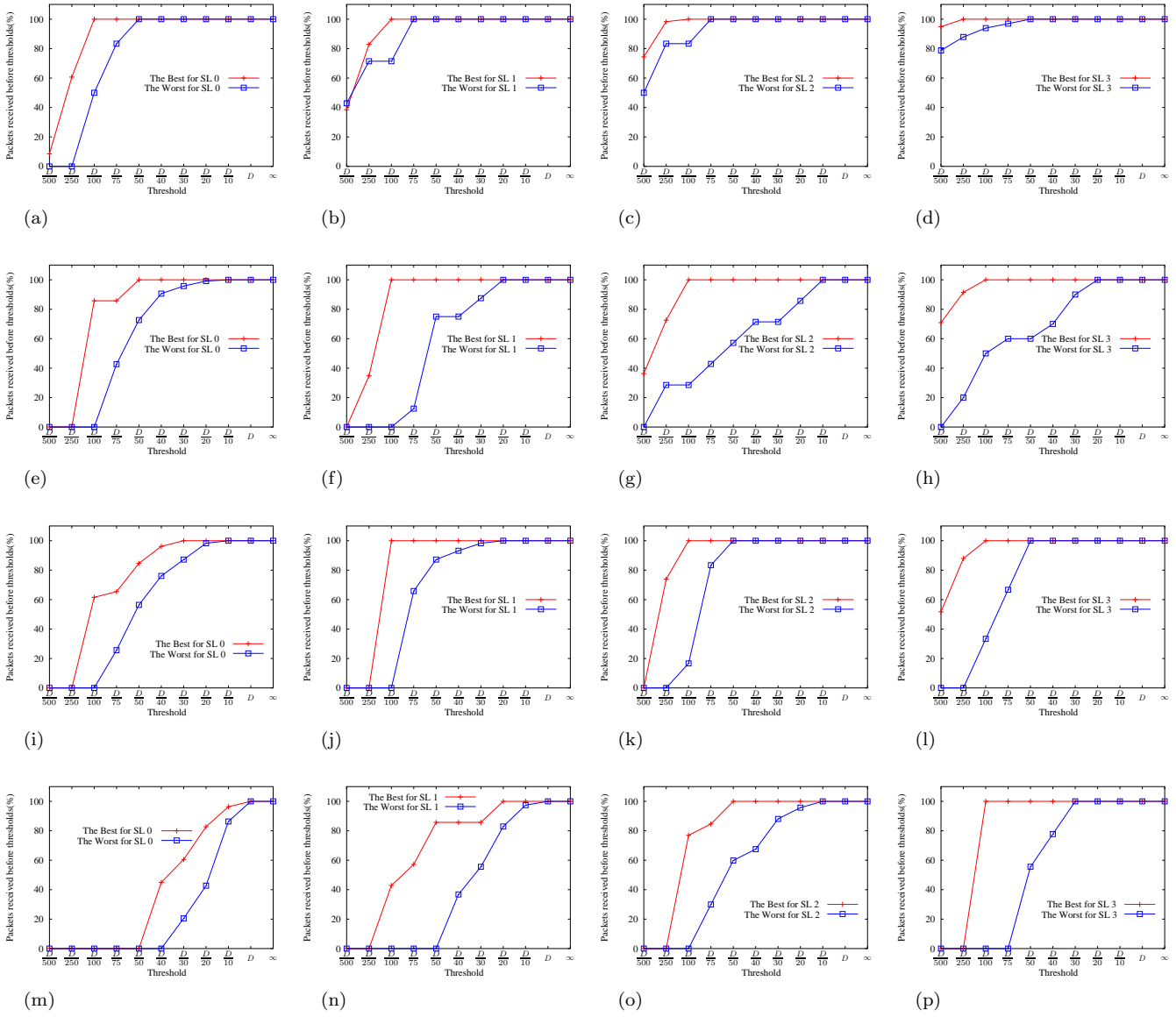


Figure 5. The best and the worst connection for SLs with the strictest delay requirements, for packet sizes of (a), (b), (c) and (d) 256 bytes, (e), (f), (g) and (h) 1024 bytes, (i), (j), (k) and (l) 2048 bytes, and (m), (n), (o) and (p) 4096 bytes.