


RESEARCH ARTICLE

WILEY

Cache-aware design of general-purpose Single-Producer–Single-Consumer queues

Vincenzo Maffione  | Giuseppe Lettieri | Luigi Rizzo

Dipartimento di Ingegneria
dell'Informazione, Università di Pisa, Pisa,
Italy

Correspondence

Vincenzo Maffione, Lungarno Antonio
Pacioti 43, 56126 Pisa PI, Italy.
Email: v.maffione@gmail.com

Present Address

Vincenzo Maffione, Largo Lucio
Lazzarino 1, 56122 Pisa, Italy.

Luigi Rizzo is now with Google.

Summary

Data processing pipelines normally use lockless Single-Producer–Single-Consumer (SPSC) queues to efficiently decouple their processing threads and achieve high throughput, minimizing the cost of synchronization. SPSC queues have been widely studied, mostly for applications such as streaming data or network monitoring, where the main goal is maximizing throughput. There are now many applications, such as virtual-machine–virtual-machine communication, software-defined networking, and message-based kernels, where low latency is also important, and the tradeoffs between high-throughput and low-latency algorithms have not been studied equally well. Furthermore, at high or variable transaction rates, the effect of memory hierarchies and cache coherence subsystems may be dominant and yield surprising results. In this paper, we make two contributions. First, we provide a comprehensive study of the two main families of SPSC queues, namely, “Lamport” and “FastForward” queues, with a detailed analytical and experimental characterization of their behavior in terms of operating regimes, throughput, latency, and cache misses. Second, we propose two new queue variants, namely, improved FastForward and batched improved FastForward, which have better worst-case behavior than other variants in terms of cache misses, which is an important feature for a number of applications. Together, these two contributions provide practical guidelines to choose the best solution depending on the application requirements.

KEYWORDS

cache-aware design, FastForward, Lamport lock-free queue, SPSC

1 | INTRODUCTION

The large number of CPUs available on modern computing systems makes it possible to run complex parallel processing algorithms on a single machine. Beyond specialized High Performance Computing platforms, the maximum CPU count is approaching 100 units, even for commodity shared-memory machines in data centers and IT departments.^{1–3} With such a high degree of hardware parallelism, it is important to keep the cost of inter-thread synchronization under control to actually benefit from the increased CPU count and improve overall efficiency.

In particular, the use of locks for synchronization at high rates (ie, a million operations per second or more) is notoriously inefficient^{4,5}; locking operations cause the threads to repeatedly issue relatively expensive atomic instructions,

such as compare-and-swap (CAS),⁶ load-linked/store-conditional,⁷ fetch-and-add (FAA),⁶ or memory fences⁸ (barriers), and above all, the communicating threads suffer from continuous cache conflicts on the cache lines that store the lock variables. Locks do not scale well: the average cost of lock operations increases quickly with the number of conflicting threads, and they cause significant performance loss even when just two threads contend for the lock. As reported in our previous work,⁹ cache conflicts are particularly problematic on multisocket NUMA machines, where a single cache miss can cost up to 200 ns, severely limiting the maximum data rate.

To overcome the inherent limitations of locks and other traditional synchronization primitives (ie, semaphores, monitors, etc), several efficient lockless and lock-free^{6,10–12} algorithms have been designed. These algorithms are nonblocking and still rely on atomic operations on shared-memory variables; hence, they are still affected by cache conflict issues. However, every access to a shared variable is an essential part of a lock-free algorithm, and it is not hidden inside any synchronization primitive. As a consequence, a carefully designed lock-free algorithm can achieve higher efficiency by trying to minimize the time spent on synchronization and, in particular, minimize cache misses.

One of the most popular categories of lock-free data structures is the class of Single-Producer–Single-Consumer (SPSC) queues, where a (single) producer thread sends a stream of data items to a (single) consumer thread. SPSC queues are widely used in data processing pipelines, where multiple threads perform a complex task by the composition of simpler tasks. Each thread receives one or more streams of data items from other threads, performs some processing on the data, and sends them to other threads in the pipeline. These processing elements are interconnected using one-to-one unidirectional links to form a directed graph; each unidirectional link is implemented using an SPSC queue, to efficiently decouple the two threads (producer and consumer) attached to the ends of the link. Such processing pipelines are largely used in various forms of network packet processing, including traffic monitoring, software switching and routing, communication between virtual machines, and network function virtualization.¹³ All these applications have to deal with millions or tens of millions of packets per second, quite often also with strict latency requirements. Another interesting use case is software-based redundant multi-threading fault tolerance,¹⁴ which is a technique to detect transient memory errors caused by *bit flipping*. The same application runs both in a main thread and in a *checker* thread. The main thread also sends to the checker thread all the data read or written from memory, using an SPSC queue. The checker thread can compare the memory transactions coming from the queue with the ones produced locally and detect a fault if they differ.

Although many efficient SPSC queue algorithms are available in the literature, most of them are designed or analyzed for specific use cases. As an example, several proposals^{14–19} assume a continuous stream of messages from the producer and use fixed, large batch sizes to optimize throughput. This approach breaks or causes large delays if the producer temporarily slows down or stops before a batch is complete. Similarly, if data items cannot be embedded in the queue (as is the case for many network applications, where items are variable-size, possibly large, packets), algorithms need further memory barriers to make sure memory operations are properly serialized. This may impact the performance significantly and possibly require a redesign of the algorithm to adapt to the different requirements.

The goal of this work is to provide a comprehensive and general study on how to design practical and performant SPSC queues. Our contributions are a detailed analysis of the two main families of SPSC queues (Lamport and FastForward queues) in terms of their throughput, latency, and cache behavior, and the introduction of two queue variants (improved FastForward [IFFQ] and batched IFFQ [BIFFQ]) that have improved the worst-case behavior over existing proposals. The result of our analysis can be applied to network processing use cases, including network function virtualization deployments, which often pose challenging problems in the form of tradeoffs between throughput and latency.

We start from the basic implementations provided by Lamport²⁰ and Giacomoni et al,²¹ which are two complementary approaches that can be compared—together with some of their variants—against different metrics, namely, throughput and latency. We also study the impact of extending the queue application programming interface (API) to let the producer and the consumer amortize queue synchronization operations over a batch of messages. Depending on how the queue is implemented, this may or may not have a significant impact on performance. Moreover, we do not take for granted that a batching-capable API can be easily integrated in a given data processing software; this may be unfeasible in practice when the software itself is not designed to operate in batch (as an example, this is the case for many parts of the Linux kernel networking subsystem). Even if these software engineering constraints are extraneous to the design of the SPSC queue itself, they must be taken into account to make the optimal choice in spite of the given limitation.

In detail, Section 2 defines the problem addressed, the assumptions made, and the metrics considered; Sections 3 and 4 provide a detailed description of the SPSC algorithms under investigation, together with an analysis of their best-case and worst-case cache behaviors; in particular, Section 4.2 describes a new queue variant based on FastForward, contributed by our work; Section 5 reports the experiments carried out to validate the analysis in terms of cache misses, throughput,

and latency; Section 6 presents an example application (a virtual Ethernet switch) that makes use of many SPSC queues and evaluates the performance impact of choosing different SPSC queues; Section 7 discusses related work; and, finally, Section 8 contains our conclusions.

2 | PROBLEM STATEMENT

An SPSC queue allows two threads to exchange data items through a shared-memory first-in–first-out queue without using locks or other synchronization primitives. One thread—the producer P—only enqueues data items; the other thread—the consumer C—only dequeues items. We assume that the number of slots in the queue is fixed. This is common practice in high-performance processing systems: the queue size is chosen so that it can absorb short-term speed mismatches between the stages of the pipeline. The problem of dynamically growing an SPSC queue is substantially orthogonal to this study. Some general techniques (eg, as described by Aldinucci et al²²) are already available to efficiently chain fixed-size queues, which therefore act as basic building blocks.

A queue can be implemented without locks or *read-modify-write* atomic operations—such as CAS or test-and-set—only if there is a single producer and a single consumer.²⁰ Some state-of-the-art SPSC implementations^{14–19} can achieve extremely high data rates—up to a billion items per second or more—but only if very large, fixed-size batches of items are exchanged in every transaction with the queue. This approach works well when the producer generates a steady stream of data without significant idle periods or when latency is not important. However, it is completely impractical for network processing workloads, where the time between packets may vary by many orders of magnitude, and large or fixed-size batches would result in unacceptable latency. Moreover, these SPSC queues also assume that data items fit entirely in the slot of the queue (*embedded payload*), which simplifies the serialization of memory operations. For larger or variable-size items, such as network packets or disk operations, almost invariably, the queue can only store pointers to the actual data blocks (*indirect payload*), potentially requesting further memory barriers to make sure that updates to indices, slots, and data blocks are seen in the correct order by the data consumer. Adding those barriers has a performance impact that must be considered.

In this paper, we aim at giving effective design guidelines for general-purpose SPSC queues; hence, we discuss strategies to optimize latency and not just throughput, address the case of both streaming and nonstreaming producers, and discuss those of embedded and indirect payload. Finally, we analyze in detail the behavior of various algorithms in terms of memory accesses.

2.1 | The role of batching operations

The key strategy to improve throughput is to let P and C *synchronize in batch* as much as possible. Although items are always enqueued and dequeued one by one, some queue implementations allow P to *publish* many new items to C with $O(1)$ accesses to the shared synchronization variables. Similarly, C can report many freed slots to P with $O(1)$ accesses. In all the SPSC implementations we are aware of, the queue offers a single-item enqueue/dequeue API, so that P and C insert or extract entries one by one; any batching of synchronization operations is decided and hidden inside the enqueue and dequeue functions. However, when the time between two subsequent invocations of enqueue is not bounded, this hidden batching must be disabled, or it would produce unacceptable latency. In these cases, it is necessary to go beyond the single-item API and let the queue offer an API with batching capabilities, as shown in Sections 3 and 4. This approach is especially effective when the processing pipeline already operates in batches, and the extended API can be used to communicate batch boundaries to the queue. As an example, this is the case for frameworks like DPDK,²³ netmap,²⁴ and PF_RING,²⁵ which natively support batched input/output from/to the network interface (NIC).

We should note, though, that retrofitting an existing data processing pipeline to make use of batched input/output is not always possible or easy. An interesting example is the transmission path in the Linux kernel network stack. The PSPAT high-performance network scheduler⁹ uses the `dev_queue_xmit()` function to pass host-generated packets to the scheduler through an SPSC queue. This function is invoked by the Transmission Control Protocol/Internet Protocol stack to send packets through a network interface, with a separate function call for each packet; thus, the intercept code does not have the chance to enqueue more than one packet with a single operation. Delaying the enqueue to artificially create a batch would require using a timer, which would add artificial latency and add a cost comparable to that of synchronization. In this example, the only way to batch SPSC enqueue operations without adding unbounded artificial latency would be to completely refactor the Linux kernel network stack to expose a batching-capable API all the way up to the userspace

applications, whereas certainly feasible in theory, it would be a very intrusive and complex change, which is arguably not desirable for reasons unrelated to our analysis.

2.2 | SPSC queues

The most important difference among SPSC queue implementations is about how producer and consumer *synchronize* with each other, namely, how the producer learns which slots are empty, that is, ready to be used to enqueue new entries, and how the consumer learns which ones contain a produced item that can be dequeued. The final goal of our analysis is to find the best synchronization strategy that meets the requirements reported in Section 2 and the additional ones set by the user.

Synchronization happens through atomic control variables, eg, 32- or 64-bit integers, that are stored in shared memory. The algorithms presented in this paper do not need to use the relatively expensive *read-modify-update* atomic instructions (eg, CAS or FAA). As a result, the main source of overhead comes from cache conflicts and misses, which we need to minimize as much as possible to improve throughput and latency. Read cache misses are due to load instructions and are generally more expensive than write misses (resulting from store instructions), as we noted in our previous PSPAT work.⁹ However, to keep the analysis simpler, in Sections 3 and 4, we simply count the total number of cache misses, without differentiating between the two types. This choice also simplifies the experiment discussion in Section 5, as our machines do not have separate CPU counters for load and store miss events.

Control variables and queue slots should be laid out in memory in such a way as to minimize cache conflicts and, in particular, to avoid *false sharing*²⁶ problems. False sharing causes cache thrashing and can be avoided by making sure that variables belonging to two different groups among the following ones never reside in the same processor cache line.

- (A) Variables written only by P and read by both P and C.
- (B) Variables written only by C and read by both P and C.
- (C) Read-only variables read by both P and C.
- (D) Variables private to P (accessed only by P).
- (E) Variables private to C (accessed only by C).
- (F) The array of queue slots.

In practice, this is achieved by partitioning the queue internal variables into the groups above, possibly adding some padding between them to make sure that each group resides in a separate set of cache lines.

We study two families of SPSC queues: the first one comes from the work originally proposed by Lamport²⁰ and is analyzed in Section 3; the second one comes from FastForward²¹ and is studied in Section 4.

For simplicity and consistently with common implementation practices, in all algorithms, we assume that the number N of slots in the queue is some power of 2, $N = 2^W$, and indices are represented by unsigned integers on more than W bits. It follows that reductions by modulo N are only needed when accessing the slots and can be implemented with a bitwise AND operation. Index manipulations and comparisons can instead be performed without reductions, taking advantage of the native integer wraparound.

If N times the slot size is a multiple of the system page size, a simple virtual memory trick allows clients to access blocks of slots of arbitrary size without worrying about wraparounds. This feature can be supported by allocating the array of slots aligned to the page boundary and mapping (eg, with `mmap`) the same array a second time in the virtual pages that immediately follow the original array.

We denote with K the number of queue slots contained in each cache line, assuming that K divides the cache line size. A typical value for K is 8, as slots often contain 64-bit pointers, and the cache line size is often 64 bytes. For convenience, we also assume that the type `Entry`, which represents an item stored in a queue slot, contains a distinguished `NULL_ELEM` value that denotes invalid items. This can be the `NULL` pointer for indirect payloads and an agreed-upon special value for embedded payloads.

3 | LAMPORT QUEUES

It was Lamport²⁰ who proposed the first lockless SPSC queue algorithm. Assuming sequential consistency, he proved that locks are not required when the queue is accessed concurrently by a single producer and a single consumer.

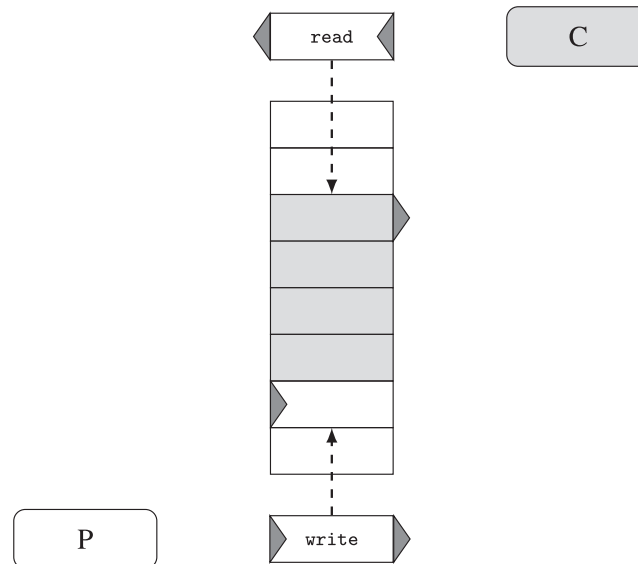


FIGURE 1 Lamport queue data structures. Each rectangle represents a cache line. In the middle, we have the `read` and `write` indices, with the `slots[]` array between them. We assume $K = 1$ for simplicity. The gray triangles on the left side of the cache lines denote possible cache misses for P: read misses if the triangle points to P and write misses otherwise. The gray triangles on the right have similar meaning for C

Synchronization is achieved by using two control variables that act as indices in the array of slots. The first index, `write`, belongs to group A (as defined in Section 2.2) and points to the next slot to be used by P. The other index, `read`, belongs to group B and points to the next slot to be used by C (Figure 1).

3.1 | Baseline Lamport queue

The baseline Lamport queue (LQ) is a modern implementation of the original queue proposed by Lamport. The routines to enqueue and dequeue entries are shown in Figure 2. Modern compilers and CPUs may reorder store and load operations to speed up code execution. Such instruction reordering has no visible effect on the CPU executing those instructions (to preserve consistency), but it may have visible effects on *other* CPUs accessing the same memory. Since P and C are normally running on different CPUs, it is necessary to add *memory barrier* operations to prevent those reorderings that can lead to synchronization errors. In `lq_enqueue`, the store used to fill the queue slot must not be reordered after the

```

1  int lq_enqueue(LQ *q, Entry e) {
2      if (q->write - q->read == N)
3          return -1; /* no space */
4      q->slots[q->write & q->mask] = e;
5      /* store to q->slot must not be reordered after store to q->write */
6      store_release_barrier();
7      q->write++;
8      return 0;
9  }
10
11 Entry lq_dequeue(LQ *q) {
12     Entry e;
13     if (q->read == q->write)
14         return NULL_ELEM; /* queue empty */
15     /* load from q->slot must not be reordered before load from q->write */
16     load_acquire_barrier();
17     e = q->slots[q->read & q->mask];
18     q->read++;
19     return e;
20 }
21

```

FIGURE 2 Implementation of the basic Lamport queue (LQ). Synchronization happens through the `write` and `read` indices in the circular array of slots. Memory barriers are necessary to prevent the compiler and the CPU from reordering store/loads operations on slots and indices [Colour figure can be viewed at [wileyonlinelibrary.com](https://onlinelibrary.wiley.com/doi/10.1002/spe.2675)]

store that updates the `write` index; otherwise, C could observe a stale slot content. For the same reason, the load used in `lq_dequeue` to read from the queue slot must not be reordered before the load that reads the `write` index. On x86 CPUs, store-release and load-acquire barriers resolve to no-operations (NOPs) due to the strong memory ordering model of the x86 architecture. However, they can still impact performance by acting as barriers for compiler optimizations; this is especially true for SPSC enqueue and dequeue routines, which are normally inlined with the rest of producer/consumer code and, therefore, make more room for compiler optimizations.

Beside barriers, cache conflicts/misses are the main factor-limiting performance. Our data layout avoids cache conflicts due to false sharing, but the following conflicts are unavoidable because synchronization between P and C happens through read-write cache conflicts on each of the two indices, one for each direction (see also Figure 1).

- $P \rightarrow C$. When C loads the `write` variable to check if there are entries to read from the queue, a read cache miss happens if P has incremented the variable since the last time C loaded it. The cache coherence subsystem will fetch the cache line containing `write` to get the updated value. When P writes to the `write` variable to publish new slots, a write cache miss happens if C has loaded the variable since the previous time P updated it.
- $C \rightarrow P$. Specular cache misses happen on the `read` variable, with the role of P and C inverted. P loads the variable to check for more free queue slots, and C increments it to report more freed slots.

Additional cache misses are necessary to transfer the queue slots from P to C. When C learns that more slots are available to be read, it loads the next unread cache line from the slots array. If cache lines are 64 bytes wide (as in our 64-bit Intel CPUs), they may store 16 32-bit integers or eight 64-bit pointers each. Queue slots may even be larger than 64 bytes, so that more cache lines are necessary to store each item. In the common case where a cache line contains more than one slot, in the best case, P will incur only one write miss for a cache line worth of slots, and C will incur only one read miss for the same amount of slots.

The number of cache misses per item depends on the pattern of accesses to the queue. The best case happens (extremely unlikely) when P and C alternate at processing the whole queue and in such a way that they never access the queue (slots array or control variables) at the same time.

- While C is not active on the queue, P pays one read miss on `read` and finds that the queue is completely empty. It completely fills the queue, paying one write miss on the array every K items, because once a cache line is brought in, the L1 cache further writes to the same cache line do not cause more write misses. It also pays a single write miss on the `write` control variable, the first time it is incremented. On average, P pays only $\frac{1}{K} + \frac{1}{N}$ misses per item.
- While P is not active on the queue, C pays one read miss on `write` and finds that the queue is completely full. It then drains the queue paying one read miss on the array every K slots; it also pays a single write miss on the `read` control variable, the first time it is incremented. Moreover, C pays $\frac{1}{K} + \frac{1}{N}$ misses per item.

The worst case for the LQ algorithm causes three cache misses per item for both P and C.

- C pays a read miss to read the updated value of `write` and finds that it was incremented by just one unit. It then pays another read miss on the array to read the new slot and a write miss to increment `read` by one.
- P pays a read miss to read the updated value of `read` and finds that it was incremented by one. It then pays a write miss to fill a free slot in the array and one more write miss to increment `write` by one.

The worst case is actually very possible and occurs when the two parties are very aggressive in accessing the queue in parallel: P never has the chance to increment `write` more than once before C loads it, C never has the chance to increment `read` more than once before P loads it, and the queue occupancy oscillates between two consecutive values (often 0 and 1 or $N - 1$ and N). In practice, depending on the relative speed of P and C and on how often they access the queue, the average number of cache misses per item may vary between the best and the worst case.

3.2 | Lazy Lamport queue

It is possible to opportunistically reduce the number of read and write misses of LQ, using *lazy loading* techniques.^{14,15,18} The dequeue function can be improved by loading the `write` variable only when no more progress can be made. This optimization only requires an additional private index variable (`write_shadow` in group E) that tracks the latest known value of `write` and is mostly useful when C is slower than P (“Fast Producer” [FP], according to the terminology we introduced in our previous work²⁷), where each time `write` is loaded, it has advanced by many positions. A specular

optimization can be done in enqueue to lazy-load the read variable, using an additional read_shadow variable in group D. This helps when P is, on average, slower than C (“Fast Consumer” [FC]). Both optimizations require no changes to the queue API, as shown in Figure 4, and help in reducing read-write conflicts on the read and write control variables, by amortizing the read misses on more items (Figure 3).

In addition to lazy loading, it is also convenient to leave K entries unused in the queue (see function llq_enqueue in Figure 4): in this way, when the queue is full, P and C can never simultaneously read/write from slots belonging to the same cache line. This improves cache behavior whenever the queue is mostly full. A specular optimization is not possible in case the queue is mostly empty, as P and C need to work on the same cache line of the slots array, so that C can timely consume the new items.

The effect of these optimizations on cache misses is significant. While a single llq_enqueue or llq_dequeue can still cause three cache misses as for LQ, this cannot happen on a sustained basis. In the case of FC (the worst case for lazy Lamport queue [LLQ]), the queue is almost always empty. P will pay a write miss on write and on the slot for each item, but the read miss on read will occur only once every $N - K$ items, to find out that $N - K$ slots are available. The

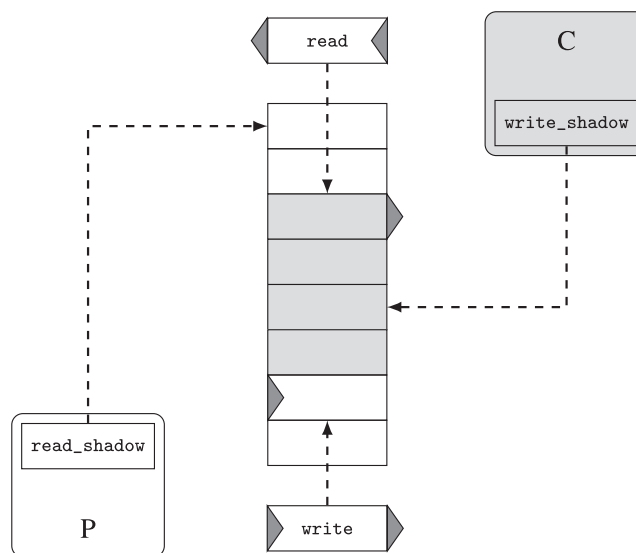


FIGURE 3 Lazy Lamport queue data structures. With respect to Figure 1, now, both P and C own a private cache line that stores a shadow copy of the queue index updated by the opposite party. The shadow copy is updated, reading from the shared variable only when necessary

```

1  int llq_enqueue(LLQ *q, Entry e) {
2      if (q->write - q->read_shadow == N - K)
3          q->read_shadow = q->read; /* lazy load */
4      if (q->write - q->read_shadow == N - K)
5          return -1; /* no space */
6      q->slots[q->write & q->mask] = e;
7      store_release_barrier();
8      q->write++;
9      return 0;
10 }
11
12 Entry llq_dequeue(LLQ *q) {
13     Entry e;
14     if (q->read == q->write_shadow) {
15         q->write_shadow = q->write; /* lazy load */
16         load_acquire_barrier();
17     }
18     if (q->read == q->write_shadow)
19         return NULL_ELEM; /* queue empty */
20     e = q->slots[q->read & q->mask];
21     q->read++;
22     return e;
23 }
24

```

FIGURE 4 Lazy Lamport queue (LLQ) improves the basic Lamport queue by loading the control indices write and read only when necessary and leaving at least a cache line worth of entries (K) empty to ensure P and C cannot work in the same cache line when the queue is full [Colour figure can be viewed at wileyonlinelibrary.com]

amortized cost is therefore $2 + 1/(N - K)$ misses per item. The cost for C is the same because every read (write) miss of P corresponds to a write (read) miss for C.

In the case of FP, the queue is almost always full: cache misses on read occur on every item, and misses on write occur every $N - K$ items. Misses on queue slots occur only once every K items, because P and C never work in parallel on the same cache line of the slots array. The amortized cost for FP is therefore $1 + 1/K + 1/(N - K)$ misses per item on both sides. The best case for LLQ is the same as for LQ, with the only difference being that the maximum queue size for LLQ is $N - K$.

3.3 | Batched Lamport queue

LLQ can amortize misses due to one peer reading the index incremented by the other peer. However, there is no way to reduce the frequency of a peer incrementing its index unless we introduce artificial latency^{14,15,17-19} or change the queue API. Since the first choice is not a viable option (see Section 2), we need to enhance the queue API with the ability to operate in batch; the goal is to allow P and C to advance the `write` and `read` indices by many units at a time. The resulting batched Lamport queue (BLQ) algorithm (see Figure 5 for the data structures used by the algorithm) is illustrated in Figure 6 for the producer and in Figure 7 for the consumer.

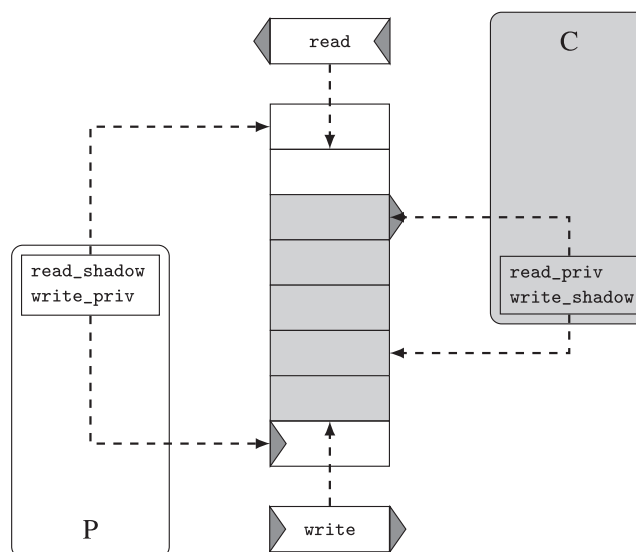


FIGURE 5 Batched Lamport queue data structures. Index updates are accumulated in private copies (`write_priv` and `read_priv`) before being made visible to the other party

```

1 unsigned blq_enq_space(BLQ *q, unsigned int needed) {
2     unsigned space = (N - K) - (q->write_priv - q->read_shadow);
3     if (space >= needed)
4         return space;
5     q->read_shadow = q->read;
6     return (N - K) - (q->write_priv - q->read_shadow);
7 }
8
9 /* No boundary checks, to be called after blq_enq_space(). */
10 void blq_enq_local(BLQ *q, Entry e) {
11     q->slots[q->write_priv & q->mask] = e;
12     q->write_priv++;
13 }
14
15 void blq_enq_publish(BLQ *q) {
16     store_release_barrier();
17     q->write = q->write_priv;
18 }
19

```

FIGURE 6 Producer routines to access a batched Lamport queue (BLQ). The producer can operate in batch, enqueueing one or many entries before advancing `write` [Colour figure can be viewed at wileyonlinelibrary.com]


```

1 unsigned blq_deq_space(BLQ *q, unsigned int needed) {
2     unsigned space = q->write_shadow - q->read_priv;
3     if (space >= needed)
4         return space;
5     q->write_shadow = q->write;
6     load_acquire_barrier();
7     return q->write_shadow - q->read_priv;
8 }
9
10 /* No boundary checks, to be called after blq_deq_space(). */
11 Entry blq_deq_local(BLQ *q) {
12     Entry e = q->slots[q->read_priv & q->qmask];
13     q->read_priv++;
14     return m;
15 }
16
17 void blq_deq_publish(BLQ *q) {
18     q->read = q->read_priv;
19 }
20

```

FIGURE 7 Consumer routines to access a batched Lamport queue (BLQ). The consumer can operate in batch, dequeuing one or many entries before advancing `read` [Colour figure can be viewed at wileyonlinelibrary.com]

The new workflow of P is as follows: (1) call `blq_enq_space` to get the number b of available slots; (2) enqueue up to b items without publishing them to C, by calling many times `blq_enq_local`; and, finally, (3) call `blq_enq_publish` to publish the new items to C. In this way, P can amortize both the read misses on the `read` index (note that lazy loading is used by `blq_enq_space`) and the write misses on the `write` index, without introducing unbounded latency: if just one or a few items are produced, they are published immediately, even if the queue has space for more items. Note that updating `write` less frequently is also beneficial for C, which will have fewer chances to suffer from read misses. To support batched updates, an additional variable `write_priv` is necessary within group D. This producer-local variable is used in place of `write` while filling unpublished slots, so that `write` is updated only once. The `needed` argument of `blq_enq_space` is used by the caller to control the lazy loading logic, ie, to force a (possible) read miss if the number of slots already available is lower than `needed`.

Similarly, the new workflow for C is as follows: (1) call `blq_deq_space` to learn the number b of items that are currently available to be dequeued; (2) call `blq_deq_local` up to b times, in order to dequeue that many items without returning any slot to P; and (3) call `blq_deq_publish` to return, ie, publish, the used items to P. As a result, C can amortize read misses on `write` by means of the lazy loading in `blq_deq_space` and can amortize write misses on `read` by publishing updates in batches. All of these are achieved without adding unbounded latency. Batch updates on `read` are enabled by a consumer-private variable `read_priv`, which thus belongs to group E.

To determine the cache miss count for FC and FP, let us assume that P can enqueue in batches of size B and that C is willing to read B items at a time if available. In the FC case, C observes batches of size B , and the queue size oscillates between 0 and B . In detail, we have the following.

- P pays a read miss on `read` once every $N - K$ items to learn that there are $N - K$ free slots, as K slots are left unused. It enqueues B items at a time, paying a write miss every K items and, finally, publishes the new slots with a single write miss on `write`. Since B consecutive slots cover, at most, $\lceil (B - 1)/K \rceil + 1$ cache lines,[†] the total cost for P in the worst case is $\frac{1}{N-K} + \frac{\lceil (B-1)/K \rceil + 1}{B} + \frac{1}{B}$ misses per item.
- C pays a read miss on `write` to learn that there are B items to dequeue and dequeues them, paying a read miss every K items. It then returns the used slots, paying a write miss on `read` only once every $N - K$ items. The total cost is the same as for C.

The FP case is symmetrical, with the queue size oscillating between $N - K - B$ and $N - K$. C(P) pays a read (write) miss on `write` once every $N - K$ items, to discover (publish) the new items ready to be processed, whereas misses on `read` happen once every B items. Differently from FC, misses on the queue slots happen exactly once every K items, because P and C always work on different cache lines. The total cost for both P and C is $1/(N - K) + 1/B + 1/K$.

If B is large enough (eg, 64 or more), the worst-case cache misses behavior of BLQ is dominated by the term $1/K$ and practically coincides with the best case for Lamport-based queues (Section 3.1), which is indeed the same for both LQ, ie,

[†]The +1 and -1 terms are needed because the batch of B slots may not be aligned to a cache line boundary.

LLQ and BLQ. Effectively, the use of large batches by P triggers the alternate processing of the queue, which results in the best-case situation for LQ. The FC operating regime is the worst case also for BLQ.

An example of producer that can batch is mentioned in Section 2.1: a packet processing application using netmap²⁴ or DPDK²³ may read 256 packets at a time from a network interface on a busy link and push them all in the first SPSC queue of a processing pipeline, where the batch can be preserved across all the threads. Note that even if P cannot operate in batch, C can still do it, and hopefully, this will prevent P from having too many cache misses on the read control variable. This means that BLQ could still offer better performance than LLQ, although the worst case would be the same as LLQ.

4 | QUEUES BASED ON FAST FORWARD

Giacomoni et al²¹ proposed FastForward as an alternative approach to the original lock-free queue proposed by Lamport. One of the main characteristics of the class of LQs is that control variables used for synchronization, ie, `write` and `read`, are decoupled from the slots array. P and C can monitor the cache lines containing the control variables to know how much work can be done and access the queue slots only to actually read or write items. The main advantage of this approach is that operating in batch is easy and efficient: once C learns that *B* items can be read, it can work on those *B* items without worrying about what P is doing in the meantime and, in particular, without the need to keep looking at `write`. A similar reasoning applies to P. However, this decoupling has two drawbacks. First, both P and C need to look at up to three shared cache lines to process a single item, namely, a cache line in group A, one in group B, and one in group F, incurring in up to three cache misses per item, as detailed in Section 3.1. Second, potentially expensive memory barriers become necessary to guarantee that P and C see memory operations on these cache lines in the same order.

4.1 | FastForward queue

FastForward was designed to overcome the disadvantages of LQs by implicitly embedding the synchronization variables within the slots, as illustrated in Figure 8 and shown in Figure 9. The `write` and `read` indices are still present and have the same meaning as in LQs, but they are private to P and C, respectively, and are only used to keep track of the current slot to process. In particular, `write` is stored in group D, and `read` is stored in group E. The indices are thus not used for inter-thread synchronization, and accessing them does not normally cause any cache miss. P and C learn which slots are ready to be processed by looking at the contents of the slots themselves. A special `NULL_ELEM` value is used to mark empty slots; at initialization, all the slots contain `NULL_ELEM`, and both `write` and `read` are set to the same value (0 or any other valid index). P enqueues new items in empty slots starting at position `write`; C consumes items from nonempty slots starting at position `read`, overwriting each used slot with the `NULL_ELEM` marker, so that P can reuse it.

The main property of FastForward queue (FFQ) is that only a single shared cache line (containing *K* queue slots) is needed for P and C to exchange an item. However, the queue slots are written by both P and C, because notifications must be bidirectional: P notifies C about new items, and C notifies P about freed slots. Because updates and notifications are done with the same `write`, FFQ does not need memory barriers in the queue routines. The complete proof for this can be found in the work of Giacomoni et al,²¹ but a simple example is sufficient to understand the main idea. Assume that C

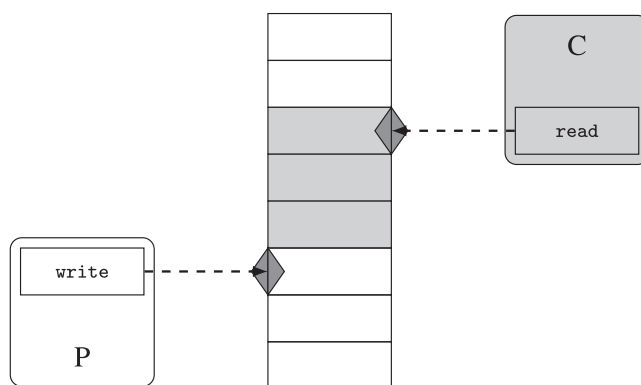


FIGURE 8 FastForward queue data structures. The queue indices are held in private cache lines because they are only needed locally. Both control and data information are exchanged through the queue slots, which are the only cache lines where the algorithm needs cache misses

```

1  int ffq_enqueue(FFQ *q, Entry e) {
2      unsigned wpos = q->write & q->mask;
3      if (q->slots[wpos] != NULL_ELEM)
4          return -1; /* no space */
5      q->slots[wpos] = e;
6      q->write++;
7
8      return 0;
9  }
10
11  Entry ffq_dequeue(FFQ *q) {
12      unsigned rpos = q->read & q->mask;
13      e = q->slots[rpos];
14      if (e != NULL_ELEM) {
15          q->slots[rpos] = NULL_ELEM; /* clear */
16          q->read++;
17      }
18
19      return e;
20  }
21

```

FIGURE 9 FastForward queue (FFQ) couples data transfer with synchronization to reduce worst-case cache conflicts. The producer and the consumer need to access only a cache line to process a single item [Colour figure can be viewed at [wileyonlinelibrary.com](https://onlinelibrary.wiley.com/doi/10.1002/spe.2675)]

is spinning on a `NULL_ELEM` in slot i and P issues a write to slot i and then a write to slot $i + 1$. Even if these writes get reordered, so that the write to $i + 1$ reaches C before the write on i , C will continue to spin on i until it receives the older write, and only then will it proceed to look at slot $i + 1$. Essentially, C is forced to see P writes in order. A symmetrical example can be constructed for P .

It is important to note that the producer may still need a memory barrier before calling the enqueue routine, eg, with indirect payload, where the producer issues stores to the payload before enqueueing. In this case, P must issue a store-release barrier between the store to the payload and the enqueue of the pointer, to make sure that payload is ready before C is notified of its presence. C does not need a corresponding load-acquire barrier, because the load from the payload has data dependency on the load from the slot pointer; no compiler or CPU reordering can happen between these two loads, to preserve sequential consistency.

```

1  void producer(FFQ *q) {
2      for (;;) {
3          Packet *pkt = ...;
4          // ...
5          pkt->x = y; /* store to the indirect payload */
6          store_release_barrier(); /* flush the store */
7          ret = ffq_enqueue(q, pkt); /* enqueue the pointer */
8          // ...
9      }
10 }

```

Another interesting property of the FFQ algorithm is that its worst-case cache behavior is better than that of LQ. FFQ needs, at most, two cache misses per item, whereas LQ may need up to three cache misses. The worst case for FFQ is similar to that described for LQ in Section 3.1; it happens when the queue size oscillates between 0 and 1 or between $N - 1$ and N . In detail, we have the following.

- P pays a read miss on the slots array cache line pointed by `write` and learns that the next slot to use is free (as it contains `NULL_ELEM`). P then pays a write miss on the same cache line when writing the new slot content, to upgrade the cache coherence protocol state from *shared* to *exclusive*.
- C pays a read miss on the slots array cache line pointed by `read` (whose state changes from *exclusive* to *shared*) and learns that the next slot to be read has a new valid content. After consuming the value, C pays a write miss on the same cache line to write back the `NULL_ELEM` marker, and the cache line state switches back to *exclusive*.

The best case for FFQ happens when P and C never work on the same cache line at the same time. This is true if the queue always contains at least K items and not more than $N - K$ items.

- P pays a read miss on the first slot of a cache line and learns that the slot is empty. P fills the slot, paying a write miss to upgrade the cache line from *shared* to *exclusive* state. The remaining $K - 1$ slots in the cache line are free and can be

read and filled without paying any miss because the cache line is already present in the L1 cache. On average, P pays $2/K$ misses per item.

- C pays a read miss on the first slot of a cache line and learns that the slot has a valid content. After processing the item, C writes `NULL_ELEM` back to the slot, paying a write miss to bring the cache line in exclusive state. The remaining $K - 1$ slots are ready and can be read and written-back without further misses, because the cache line is already present in the cache. The average cost for C is, thus, also $2/K$ misses per item.

The best case is extremely difficult to achieve in practice, as it would require P and C to operate at the same average rate while never approaching the boundaries of the queue, so that both P and C can operate in parallel on different cache lines. In contrast to LQ and LLQ, the best case for FFQ is desirable and more realistic to obtain, as the queue length can still vary freely in the range $[K, N - K]$. The original FastForward paper²¹ indeed proposes a control algorithm to maintain this property, possibly adding artificial delay in the processing of the producer or consumer; as this technique breaks our requirements on bounded latency (Section 2), we do not consider it here.

4.2 | Improved FastForward queue

Similarly to what has been done for LQ, there is room to reduce FFQ misses by limiting the situations where P and C may work in the same cache line. In this section, we present and discuss IFFQ, an improved version of FFQ based on the mailbox data structure used by the PSPAT⁹ high-performance network scheduler.

The FFQ producer pays one read miss to check whether a slot is available and one write miss to fill it. Read misses for P can be reduced by a factor of H by considering a slot available only if the first slot of the cache line H slots *ahead* of write is `NULL_ELEM`, as illustrated in Figure 10 and in function `iffq_enqueue` in Figure 11. The H parameter is a small multiple of K (eg, $H = 4K$) and a divisor of N , so that we can consider the queue as partitioned in groups of H/K consecutive cache lines. A private variable `limit` (group D) tracks the first slot in a cache line following the `write` position, which is the first slot that P cannot use. P needs to check the slot `limit + H` once every H insertions, causing $\frac{1}{H}$ read misses per item.

The FFQ consumer needs one read miss to fetch a new item and one write miss to clear the slot. Write misses can be reduced by a factor of K if C postpones the clear operation and frees at least K items at a time. To implement such a *lazy clear* strategy, C keeps an additional consumer-private index variable (group E) called `clear`, which points to the next slot to be cleared, as illustrated in Figure 10. The invariant to guarantee is that both `write` and `read` are always in a different cache line than `clear`, so that P can monitor the first slot of the cache line pointed by `clear` without interfering with C reading from `read`.

To maintain the invariant, C returns all the H slots in the i th partition only when both the i th and $(i + 1)$ th partitions have been completely processed, and C is already reading from the $(i + 2)$ th partition. As soon as C clears the first slot of partition i , P can immediately start to use the partition $i - 1$, so that C keeps clearing the rest of the slots in i concurrently with P filling the slots in $i - 1$. The drawback of this approach is that it is necessary to leave $2H$ slots unused, so that the effective maximum queue occupation becomes $N - 2H$.

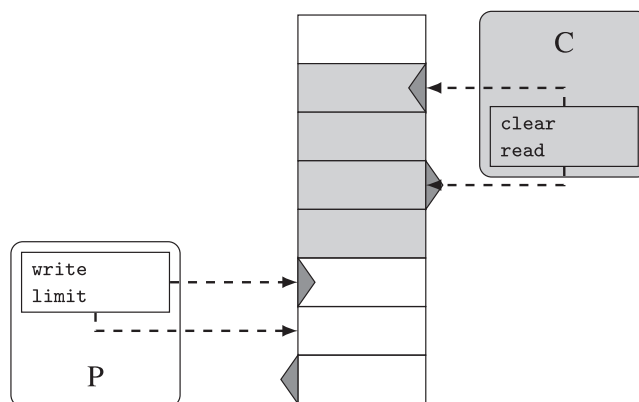


FIGURE 10 Improved FastForward data structures. With respect to Figure 8, the `limit` and `clear` private indices are introduced to reduce the cases where P and C access the same cache line at the same time; this happens when `write` and `read` (or `clear`) point to the same cache line

```

1  int iffq_enqueue(IFFQ *q, Entry e) {
2      if (q->write == q->limit) {
3          /* Check if the next queue partition is free. */
4          unsigned next_limit = q->limit + H;
5          if (q->slots[next_limit & q->mask] != NULL_ELEM)
6              return -1; /* no space */
7          q->limit = next_limit; /* Free partition, advance producer limit. */
8      }
9      q->slots[q->write & q->mask] = e;
10     q->write++;
11     return 0; /* OK */
12 }
13
14 Entry iffq_trydeq_local(IFFQ *q) {
15     Entry e = q->slots[q->read & q->mask];
16     if (e != NULL_ELEM)
17         q->read++;
18     return e;
19 }
20
21 /* Opportunistically clear slots, making sure that q->read is at least H slots
22  * ahead of q->clear, and that q->clear stops at the beginning of a queue partition. */
23 void iffq_deq_publish(IFFQ *q) {
24     unsigned next_clear = (q->read & ~(H-1)) - H;
25
26     while (q->clear != next_clear) {
27         q->slots[q->clear & q->mask] = NULL_ELEM;
28         q->clear++;
29     }
30 }
31

```

FIGURE 11 The improved FastForward queue (IFFQ) ensures that the producer and the consumer never access the same cache line in parallel when the queue is almost full, achieving optimal performance in this operating regime [Colour figure can be viewed at wileyonlinelibrary.com]

For the algorithm to maintain its invariants (and work properly), `clear` is initialized to 0, `write` and `read` are initialized to H , and `limit` is initialized to $2H$. Slots are initialized to `NULL_ELEM`, except for those between `clear` (included) and `read` (excluded), which must be initialized with some different value, to prevent P from quickly filling the queue and using them.

Similarly to BLQ, IFFQ supports batch operation on the consumer side, providing a consumer API different from a simple dequeue. As illustrated in Figure 11, C operates as follows.

- Call `iffq_trydeq_local` many times to read the available items, until there are no more items and the function returns `NULL_ELEM`.
- Call `iffq_deq_publish`, which opportunistically clears slots and advances `clear` until it reaches the first slot of the partition immediately preceding `read`.

Note that C delays the clear operation with no impact on latency. If the queue is almost full, the latency experienced by an item is in only due to the time needed by C to process all the items that were enqueued before; postponing the enqueue of further items (because of C clearing slots less frequently) does not worsen latency. If the queue is almost empty, P can go ahead without noticing the delayed clear.

The worst case for IFFQ is the same as FFQ, and it happens when the queue size oscillates between 0 and 1, with `write` and `read` always in the same cache line and `limit+H` and `clear` always in different cache lines. In this FC scenario, P and C pay a miss on writing/reading each item. Once every H items, P pays an additional read miss on the slot pointed by `limit+H`, learning that H new slots are free to be used, for a total amortized cost of $1 + 1/H$. C pays an additional write miss every K items when clearing the first slot of a previously read cache line. The amortized cost for C is therefore $1 + 1/K$.

In the FP case, the queue is almost full, `write` and `read` are guaranteed to point at different cache lines, and `limit+H` coincides with `clear` most of the time. The system evolves as follows.

- P pays a read miss on the slot pointed by `limit+H` to learn that H new items are available and a write miss when filling the first slot of a new cache line. Writing to the other $K - 1$ slots does not cause further misses because the cache line is already in exclusive state. The total amortized cost for P is $1/H + 1/K$ per item.

- C pays a read miss every K slots to read new items. When clearing H slots, a write miss occurs on the first slot of each returned cache line, which transitions to the exclusive state. However, P normally reads back the first slot of the first cache line (while checking for more space on $\text{limit}+H$) before C has the chance to clear the remaining $K - 1$ slots. This cache line becomes shared, and C pays an additional cache miss to bring it back to the exclusive state. The total amortized cost for C is therefore $\frac{1}{K} + \frac{2+H/K-1}{H} = \frac{2}{K} + \frac{1}{H}$ misses per item.

The best case for IFFQ happens when $\text{limit}+H$, `clear`, `read`, and `write` always point at four different cache lines. The amortized cost in this case is $1/K + 1/H$ for P and $2/K$ for C. Although its best case is unlikely, a major advantage of IFFQ is that its cache miss behavior is very close to optimal in the FP scenario, which is quite common.

4.3 | Batched IFFQ

IFFQ needs very few cache misses in most regimes except for FC, when the queue is almost empty most of the time. This regime is problematic for all the queues under study. For Lamport's queues, adding a batching API on the producer side allowed us to atomically publish multiple entries at once and amortize cache misses when reading new entries (Section 3.3). It is therefore desirable to extend IFFQ in such a way as to provide an API to enqueue in batch (batching is already supported on the consumer side).

Unfortunately, in FFQs, *it is impossible* to atomically publish an update involving multiple slots, because the synchronization information is implicitly embedded in the slots array, and P can only submit slots one by one. However, once P fills the first slot of a cache line, this transitions to the exclusive state, and a quick burst of back-to-back writes to the same cache line has a high probability to proceed without further write misses (or just an additional one if C was actively monitoring the same cache line). The core idea for a BIFFQ is therefore to accumulate new items in a temporary buffer, as illustrated in Figure 12, and copy them to the shared queue in a tight loop only at the end of the batch. The resulting producer-side routines for BIFFQ are shown in Figure 13. Consumer-side routines are not shown, as they are the same as illustrated in Figure 11.

The new workflow for a batching producer is similar to that described for BLQ (Section 3.3).

- Call `biffq_wspace` to learn the number b of available slots. The `needed` argument is used by the caller to force the increment of `limit` by H if the number of slots already available is lower than `needed`.
- Call `biffq_enq_local` to enqueue up to b items without publishing them. The items are stored in a producer-private buffer. The buffer (`buf`) and the corresponding counter variable (`buffered`) are stored in group D.
- Call `biffq_enq_publish` to publish the buffered items. A tight loop is used to perform the burst of writes to the slots array, with the goal of minimizing the likelihood of cache conflicts with the consumer.

Essentially, BIFFQ relies on a race condition to improve over IFFQ in FC regimes, but its theoretical worst case is not different from IFFQ. The FP and best-case analyses are also the same as for IFFQ. Experiments in Sections 5 and 6 show how this best-effort optimization can be effective in practice, confirming that the race condition happens frequently enough to improve performance.

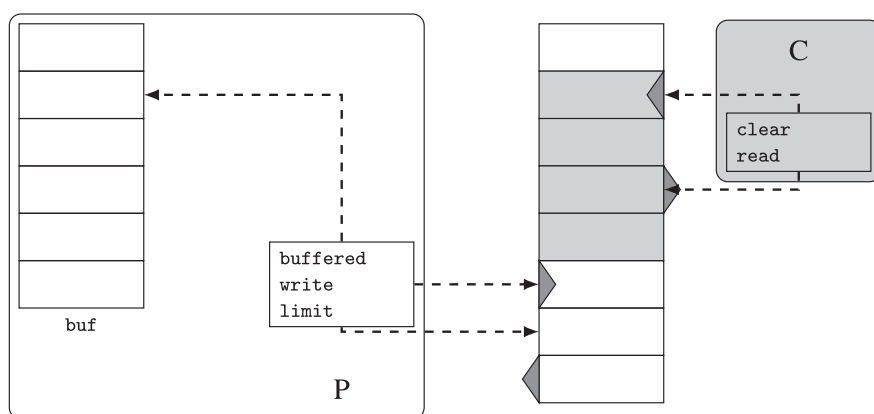


FIGURE 12 Batched improved FastForward queue. The producer accumulates new items in a private buffer and then tries to write them into the shared queue slots as fast as possible, to reduce the chance of conflict with the consumer when the queue is almost empty


```

1 void biffq_wspace(BIFFQ *q, unsigned int needed) {
2     unsigned int space = q->limit - q->write;
3     unsigned int next_limit;
4
5     if (space >= needed) {
6         return space;
7     }
8
9     next_limit = q->limit + H;
10    if (q->slots[next_limit & q->qmask] != NULL_ELEM)
11        return space;
12    q->limit = next_limit;
13
14    return next_limit - q->write;
15 }
16
17 /* Store items in a producer-local buffer. */
18 void biffq_enq_local(BIFFQ *q, Entry e) {
19     q->buf[q->buffered++] = e;
20 }
21
22 /* Copy the buffer into the queue slots as fast as possible
23  * to minimize possible interference with the consumer. */
24 void biffq_enq_publish(BIFFQ *q) {
25     for (unsigned i = 0; i < q->buffered; i++, q->write++)
26         q->slots[q->write & q->mask] = q->buf[i];
27     q->buffered = 0;
28 }
29

```

FIGURE 13 Producer routines for the batched version of the improved FastForward queue (BIFFQ), enabling batched operation for P. Publication of many items with a single atomic write is not possible with FastForward queues: BIFFQ optimistically relies on a race condition to try to achieve the same effect [Colour figure can be viewed at [wileyonlinelibrary.com](https://onlinelibrary.wiley.com/doi/10.1002/spe.2675)]

To analyze the probabilistic cache miss behavior for FC, let us assume that P produces B items at a time. On the first write of a burst, P pays a write miss to bring the cache line pointed by `write` in exclusive state. Since C was monitoring the same slot, C incurs a read miss and brings the cache line in shared state. If the written slot was not the last of the cache line (worst case), the following write in the burst causes an additional write miss for P and, later, an additional read miss for C. The remaining writes only cause a miss every K items for both P and C (with very high likelihood), because C is busy with the first item in the batch and does not interfere immediately. Since B items can cover, at most, $\lceil (B-1)/K \rceil + 1$ cache lines, it follows that the total cost for P is $\frac{2+\lceil (B-1)/K \rceil}{B} + \frac{1}{H}$, taking into account the read misses on `limit` + H . Similarly, considering the write misses on `clear`, the amortized cost for C is $\frac{2+\lceil (B-1)/K \rceil}{B} + \frac{1}{K}$ per item.

5 | EXPERIMENTAL VALIDATION

The analysis in Sections 3 and 4 identifies strengths and drawbacks of some promising SPSC queue implementations, with a focus on the cache miss behavior under various circumstances. Table 1 summarizes the result of this analysis, reporting for each queue the average number of cache misses per item in the best and worst cases, together with the common situations in which P is, on average, faster than C (FP) or the other way around (FC). In this section, we validate the analysis and perform some throughput and latency experiments, to verify our findings and study how the queues behave in practice on modern machines. With the goal of isolating the effects of those cache misses that are due to queue synchronization, it is extremely important to carefully design the validation experiment in such a way as to minimize or remove all the possible sources of noise that could hide or affect the phenomenon under observation. The precautions taken for this purpose, such as avoiding cache misses not pertaining to queue synchronization and disabling or deceiving cache line prefetching, are detailed in the following sections.

For the experiments, we use two different machines, called I7 and XEON40. I7 is a single-socket machine with an i7 processor, whereas XEON40 is a dual-socket machine with Xeon processors. Their specifications are reported in Table 2. Dual-socket machines are particularly interesting for SPSC measurements, because cache interactions between two threads running on different sockets are more expensive with respect to the case where threads run on the same socket. The cache line is 64 bytes in both machines, and the size of a pointer is 8 bytes. As a result, each cache line can contain up to eight pointers; hence, we have $K = 8$. The need for running P and C on separate sockets may arise in real

TABLE 1 Average number of cache misses per item that P and C must pay in the worst case. N is the queue size (in items), B is the number of items in the producer and consumer batches, K is the number of items that fit in a single cache line, and H is the partition size for the improved FastForward queue (IFFQ) and the batched improved FastForward queue (BIFFQ). It is common to have $N \gg K$; thus, $\frac{1}{N-K}$ terms have been replaced with $\frac{1}{N}$ to improve readability. Moreover, components with denominator N can usually be neglected. The ideal optimal number of cache misses per item is $\frac{1}{K}$

Average Number of Cache Misses Per Item Per Side (P, C) for Different Access Patterns				
Queue	Worst Case	Fast Consumer Worst Case	Fast Producer Worst Case	Best Case
LQ	3	3	3	$\frac{1}{K} + \frac{2}{N}$
LLQ	$2 + \frac{1}{N}$	$2 + \frac{1}{N}$	$1 + \frac{1}{K} + \frac{1}{N}$	$\frac{1}{K} + \frac{2}{N}$
BLQ	$2 + \frac{1}{N}$	$\frac{2 + \lceil \frac{B-1}{K} \rceil}{B} + \frac{1}{N}$	$\frac{1}{K} + \frac{1}{B} + \frac{1}{N}$	$\frac{1}{K} + \frac{2}{N}$
FFQ	2	2	2	$\frac{2}{K}$
IFFQ, BIFFQ (P)	$1 + \frac{1}{H}$	$1 + \frac{1}{H}$	$\frac{1}{K} + \frac{1}{H}$	$\frac{1}{K} + \frac{1}{H}$
IFFQ, BIFFQ (C)	$1 + \frac{1}{K}$	$1 + \frac{1}{K}$	$\frac{2}{K} + \frac{1}{H}$	$\frac{2}{K}$

Abbreviations: BLQ, batched Lamport queue; FFQ, FastForward queue; LQ, Lamport queue; LLQ, lazy Lamport queue.

TABLE 2 Specifications of the machines used in the experiments. I7 and XEON40 have similar sizes and configuration for both data caches and data translation lookaside buffer (DTLB). The L1 data cache, L2 cache, and DTLB are private to each core, whereas the L3 cache is shared by all the cores on the same socket

Test Machine Specifications		
	I7	XEON40
CPU model name	i7-3770K	XEON E5-2640
CPU frequency	3.5 GHz	2.4 GHz
Number of sockets	1	2
Number of CPUs	4 cores, 8 threads	20 cores, 40 threads
Memory speed and type	1.33 GHz DDR3	2.133 GHz DDR4
Kernel	Linux 4.15	Linux 3.10
L1 data cache	32 KB private per core	32 KB private per core
L2 cache	256 KB private per core	256 KB private per core
L3 cache	8 MB shared	25 MB shared per socket
DTLB for 2-MB pages	32 entries per core	32 entries per core
DTLB for 4-KB pages	64 entries per core	64 entries per core

applications like PSPAT,⁹ where the arbiter thread consumes packets generated by several producers (eg, virtual machines) running on potentially every single free core on the machine, both on the local and remote sockets. More generally, any High Performance Computing application that uses all or most of the cores of a multisocket machine may need two cores on different sockets to communicate through an SPSC queue.

Table 3 gives a summary of the experiments we have run, with pointers to the sections and figures where they are discussed and a brief summary of the main results.

5.1 | Validation methodology

To reduce measurement noise and increase experiment reproducibility, some general precautions have been adopted.

- Each producer or consumer thread used in any experiment is pinned to a separate dedicated physical core to avoid interference due to hyperthreading.
- On the XEON40 machine, the producer and consumer threads are pinned to two cores belonging to different CPU sockets, in order to highlight how cache misses affect performance in the worst case.

TABLE 3 We performed three kinds of measurements, namely, number of cache misses, throughput, and latency, with either embedded or indirect payload. The number of cache misses is reported only for I7, as that for XEON40 is similar

Summary of Validation Experiments		
Experiment	Sections/Figures	Main Results
No. of cache misses	Embedded payload: Section 5.2.1, Figure 14;	Highly dependent on regime (Fast Producer/Consumer); FF-based queues are generally better; batched queues are equally optimal, but BIFFQ degrades better
	indirect payload: Section 5.2.2, Figure 17	
Throughput	Embedded payload: Section 5.2.1, Figures 15 (I7) and 16 (XEON40); indirect payload: Section 5.2.2, Figures 18 (I7) and 19 (XEON40)	No clear winner and FFQ can be worse than LQ; IFFQ is a reasonable choice
Latency	Embedded and indirect payloads: Section 5.3, Figures 20 (I7) and 21 (XEON40)	IFFQ wins; batched queues can be detrimental

Abbreviations: BIFFQ, batched FastForward queue; FF, FastForward; FFQ, FastForward queue; IFFQ, improved FastForward queue; LQ, Lamport queue.

- CPU dynamic frequency scaling is disabled, with frequency pinned to the maximum supported one.
- The test machine is not used for additional jobs during the tests; it only runs producer and consumer threads and basic operating system services.
- When measuring cache misses, hardware data prefetching is kept disabled through the machine firmware, to remove the associated noise.
- Cache misses on data structures that are not part of the queue synchronization algorithm are avoided or minimized as much as possible, in order to better isolate the cache behavior to be attributed to the algorithm itself.
- Memory for SPSC queues and buffers is allocated using Linux *hugepages*, in order to reduce the pressure on the translation lookaside buffer (TLB). On our test machines, each hugepage is 2 MB in size. We are able to perform all our experiments using less than 4 hugepages, which can be always resident in the data TLB.

For throughput measurements, we use a single SPSC queue with $N = 256$ slots, a producer thread (P), and a consumer thread (C). P produces items as fast as possible, and C greedily consumes them, using either embedded or indirect payload. With embedded payload, items are 64-bit integers generated by P and read by C. With indirect payload, items are pointers to preallocated buffers, and each buffer contains a 32-bit integer field accessed by both P and C, whereas the other bytes in the buffer are never accessed. P and C are able to work with a configurable (maximum) batch limit B , which is meaningful only if the SPSC queue exposes batching capabilities (like BLQ and BIFFQ). In that case, P publishes, at most, B_P items at once, and C frees, at most, B_C slots at once. P (C) may publish (free) less than its batch limit when the queue has not enough free slots (available items). For queues without batching capabilities, it can be assumed that $B_P = B_C = 1$. Finally, IFFQ and BIFFQ use $H = 4K = 32$.

Buffer pool implementation

The pool of preallocated buffers is used for experiments with indirect payload. It is a circular array of $2N$ buffers—with N being the size of the SPSC queue—so that P can never run out of buffers. Storing the buffers in contiguous virtual memory helps in reducing the pressure on the TLB, since the whole array can be contained in a single 2-MB page and, thus, occupy a single entry in the TLB. P keeps an index to track the buffer in the pool to be used next. The reason why the pool contains more than N entries is subtle. Since we want to validate the cache behavior due to queue synchronization, we do not want P and C to conflict on the pool and on the buffers, as this would cause many additional cache misses that are not part of our model. This issue is avoided because the pool is not a real buffer allocator: the pool is accessed only by P, buffers are *lent* to C through the SPSC queue, but there is no way for C to *return* a buffer to the pool—the buffer actually never moves from the pool. C simply stops using a buffer when it is done consuming it. However, a very fast P keeps the queue constantly full; C may extract the first out of N pending buffers from the queue (thus immediately releasing a slot) and start accessing it, whereas in parallel, P gets and accesses the same buffer to fill in the slot just freed. This situation causes undesired conflicts between P and C, and it is easily solved if the pool contains more than N buffers; we chose $2N$ because N is a power of 2, to simplify modulo operations. Note that the problem here described does not affect the algorithms that leave empty cache lines to reduce conflicts.

To produce an item, P takes the next buffer from the pool, writes a sequence number in the integer field inside the payload, and enqueues the pointer, operating in batch if the queue offers such a capability. As explained in Section 4.1, with FastForward variants, P must issue a store-release memory barrier between the write access to the payload and the enqueue; such a memory barrier is not needed with LQs, because it is already included inside the enqueue routine. To consume items, C reads the sequence numbers from the pointed buffers and sums them up.

It is important to notice that with indirect payload, both P and C always pay a cache miss per item on the pointed buffer, which adds to the cache misses due to SPSC synchronization reported in Table 1. This normally causes a substantial throughput drop with respect to the embedded payload case but is also a necessary cost to pay, for instance, in packet processing systems. Each buffer in the pool is 4096 bytes large (the size of a physical page on our test platforms), even if P and C only access 4 bytes of it. As the CPU prefetcher does not cross page boundaries (according to the processor manual²⁸), it cannot guess the (regular) access pattern within the circular array and possibly speed up P and C. This is intentional, because we want to avoid those cache misses that are due to the prefetcher; moreover, we want to emulate a real system where buffers are scattered across memory and the prefetcher cannot be very effective.

Load emulation

Both P and C are instrumented to emulate additional per-item computations. Emulation is achieved by wasting CPU time with a tight loop that spins for a configurable amount of time; time sampling is quite efficient as it is based on the

timestamp counter register available on our x86 platforms. The amount of wasted time can be chosen separately for P and C. Load emulation for P happens before getting the next embedded value or indirect buffer pointer (and writing into the indirect buffer); this models a situation in which P performs some work to prepare the data and then writes the resulting data into the queue. Symmetrically, load emulation for C happens after reading the embedded value or the indirect buffer pointer from the queue (and after reading from the indirect buffer); this models a situation in which C performs some work to consume the data read from the queue. It is important to make sure there is data dependency between the emulated work and the value written/read to/from the queue slot, as it would happen for a real application. If this were not the case, the emulated work could hide the latency caused by a cache miss in the enqueue/dequeue routines, due to the out-of-order execution capabilities of the CPU; this would invalidate the meaningfulness of our measurements. To achieve data dependency, each iteration in the emulation loop increments a `trash` variable that has data dependency on the following enqueue or the preceding dequeue operation. A suitable compiler barrier is added to make sure the compiler cannot optimize the increment by aggregation. C increments the `trash` variable right after the dequeue operation, before the load emulation loop, and the value of the increment depends on the dequeued item. In the case of indirect payload, P uses the value of the `trash` variable after the emulation loop to fill the 32-bit field in the buffer right before enqueueing. With embedded payload, P uses the value of the `trash` variable to fill the payload itself. Emulating load for P and C is useful to bring any SPSC queue in a given FC or FP operating regime and observe how the queue behaves in that case. In the following, emulated load is represented by the L_P variable for the producer and by the L_C variable for the consumer.

Measurement methodology

A single test run is executed by running the `spscq` program, which is freely available online.²⁹ The program spawns a thread for P and a thread for C, and both start working on an SPSC queue as greedily as possible. A third control thread is used to measure execution time and stop P and C when the test time duration has elapsed. To minimize the interference caused by the control thread, P and C stop when the control thread sets a global variable that P and C periodically read; this only causes an additional cache miss for P and C at the end of the experiment, which does not affect the tests. For each configuration, each test run lasts for 10 seconds and is repeated 10 times in order to allow the computation of average values and standard deviations.

Both P and C count the number of items processed and the number of *perceived batches of items*. A batch of items is defined as the number of items processed without stopping because the queue is full (for P) or empty (for C). To count the number of batches, we simply count the number of times the enqueue or dequeue fails (without counting repeated failures). If, for instance, P is significantly faster than C, the batch count for P will be very low, whereas the C count will be high and close to the item count. The ratio between the item count and the batch count is the perceived average batch. The average batch of P and C is an interesting indicator of the operating regime, as explained in our previous works.^{27,30} If P's batch is very low (eg, close to B_P) and C's batch is much higher than B_C , we can conclude that the system is operating in an FP regime. Symmetrically, the system is FC if C's average batch is very short and P's average batch is large. Ideally, if P and C have exactly the same speed, their perceived batch is infinite; more realistically, P and C may have similar speed and interact in such a way that they both perceive a large batch.

Finally, to validate the behavior of P and C in terms of cache interactions, the Linux `perf` tool is used to read from the CPU hardware counters that count L1 cache misses due to load and stores. Our machines do not have separate counters for store and load misses but have an aggregate metric. This is not an issue, as the analysis does not distinguish them. Cache miss counters are sampled separately for the two cores where P and C run. With these counters, it is possible to compute the average number of cache misses per item and validate the results reported in Table 1.

5.2 | Throughput experiments

To validate the analysis of Sections 3 and 4 under different workloads and constraints, we carry our throughput experiments using both embedded and indirect payloads. For both cases and for each of the six queues under test, we measure throughput, perceived batch, and cache misses per item, as described in Section 5.1. To explore different operating regimes, we use emulated loads on both P and C, with different combinations: for both P and C, we let load vary between 0 and 90 ns in steps of 10 ns, so that we have 100 different combinations for each queue. Loads larger than 90 ns are less interesting, as the impact of cache misses on throughput becomes smaller when P and/or C spend most of their time in the emulation tight loop; this aspect will be confirmed by the results presented in the following. For the queues that can operate in batch (BLQ and BIFFQ), first, we set the batch size to $B_P = B_C = 32$ items, and, second, we set $B_P = B_C = 1$ to measure how their performance degrades.

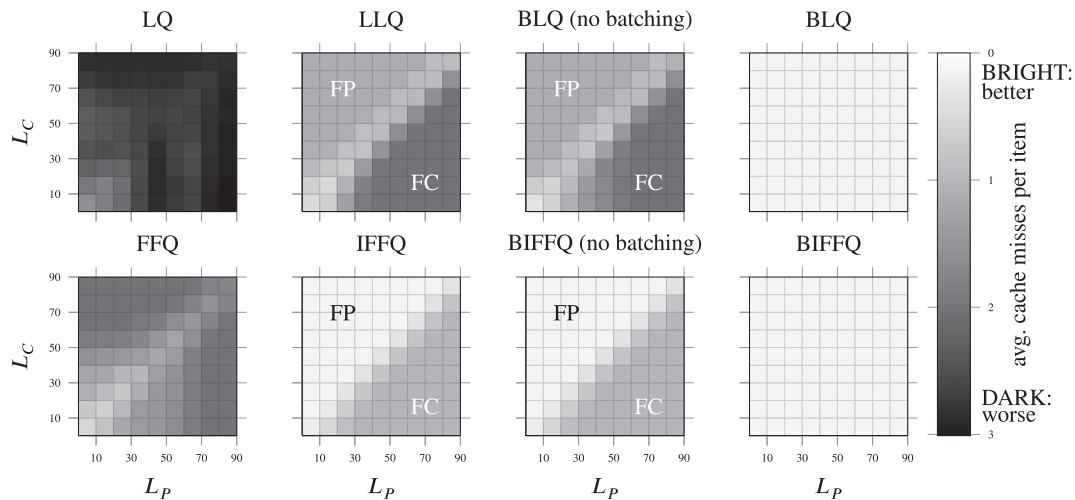


FIGURE 14 Average cache misses per item suffered by P during throughput experiments on I7 with embedded payload. Emulated load for both P and C ranges between 0 and 90 ns. Brighter colors correspond to a more favorable behavior, with less cache misses per item. BIFFQ, batched FastForward queue; BLQ, batched Lamport queue; FFQ, FastForward queue; FP, Fast Producer; IFFQ, improved FastForward queue; LLQ, lazy Lamport queue; LQ, Lamport queue

5.2.1 | Embedded payload

Cache behavior

Figure 14 shows the average number of cache misses per item on I7 in the case of embedded payload. Unless where explicitly noted, the values reported are only the ones measured for P; the values measured for C are practically identical most of the time.

Each matrix reports results for a different queue type, with each cell indicating average cache misses per item for a given combination of emulated loads of P and C (from 0 to 90 ns, as indicated on the axes). FP regimes are in the top-left region of each matrix, whereas FC regimes are in the bottom-right region. Brighter cells indicate fewer cache misses (better). Measured values range from approximately 1/8 to 3, as predicted in Table 1.

Results confirm that LQ is clearly worse than any of the other queues, for any combination of emulated loads. LQ suffers from many unnecessary synchronization operations, and this is particularly evident for FC or FP regimes, where both P and Q often end up paying three cache misses per item, which is the worst case. The small optimizations introduced by LLQ are already very effective at reducing cache misses, especially for FP regimes where we measured approximately 1.1 misses per item, whereas FC regimes report about two cache misses per item. These measurements perfectly agree with Table 1. For both FC and FP, the faster thread measured an average batch close to 1, whereas the slower thread measured a large batch. Along the main diagonal, P and C have very similar speed, and the number of cache misses is lower; this happens because both P and C perceive a large batch and, therefore, manage to amortize the cost of cache misses over many items. This phenomenon is not visible in LQ, because LQ has no logic to suppress synchronization operations; hence, a larger perceived batch does not necessarily lead to better performance.

Results for BLQ are significantly better because of the large batch size in use ($B_P = 32$). As a consequence, any synchronization cost due to accessing the `write` or `read` variable is always amortized on at least 32 items. Depending on the combination of emulated loads, the operating regime can still be FP or FC, so that the slower thread will perceive a much larger batch size. According to Table 1, we should expect between a theoretical minimum of $2/256 + 1/8 \approx 0.13$ and a maximum of $\frac{2 + [(32-1)/8]}{32} + \frac{1}{256} \approx 0.19$ cache misses per item. Measurements report values between 0.16 and 0.17, consistent with the observation that the best case is extremely unlikely in LQ variants.

The measurements also confirm that FFQ is clearly a better solution than LQ, as suggested by the original FastForward work.²¹ The FFQ worst case of two cache misses per item is observed in the FP and FC combinations, ie, the cells farther from the main diagonal. As described in Section 4.1, this is due to P and C always working in the same cache line. The situation improves for cells near the main diagonal, as P and C have less chances to work in the same cache line; the observed number of cache misses per item here is close to 1 or less. We do not go much below 1 or approach the theoretical minimum ($2/8$) because our FFQ implementation does not have a control algorithm to artificially keep P and C away from each other; thus, especially with a relatively short queue, P and C end up working on the same cache line quite often.

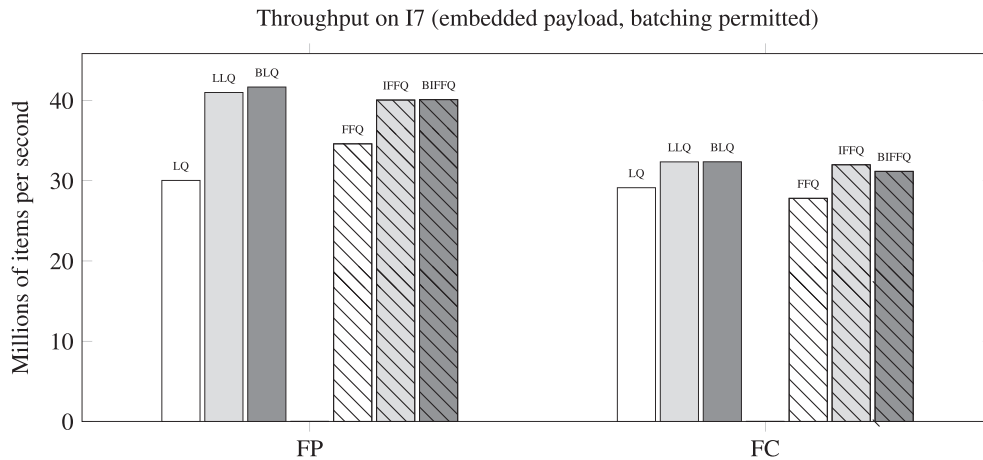


FIGURE 15 FP indicates a Fast Producer operating regime corresponding to the emulated loads $L_P = 0$ and $L_C = 10$ (in nanoseconds). FC indicates a Fast Consumer operating regime with $L_P = 20$ and $L_C = 10$. BIFFQ, batched improved FastForward queue; BLQ, batched Lamport queue; FFQ, FastForward queue; IFFQ, improved FastForward queue; LLQ, lazy Lamport queue; LQ, Lamport queue

IFFQ improves FFQ for all the combinations of loads, and this is expected because IFFQ has a better worst-case behavior. Measurements never show more than 1.1 cache misses per item, with a predicted worst-case value of $1 + \frac{1}{8} \approx 1.13$ (Table 1). FP regimes have very few cache misses, with the measured values matching the theory, which predicts, at most, $\frac{1}{8} + \frac{1}{32} \approx 0.16$ misses for P and $\frac{2}{8} + \frac{1}{32} \approx 0.28$ for C. This confirms that when P and C operate on slots more than one cache line apart, IFFQ shows optimal behavior irrespective of their distance. Finally, measurements confirm that BIFFQ has practically optimal and flat cache miss rates in all regimes (similarly to BLQ), because the batching on the producer side is effective and removes the extra cache miss per slot in FC regimes. The batch perceived by C is indeed at least 32, as P always manages to issue a quick burst of 32 writes without conflicting with C. The measured cache miss rate is about 0.26, which is slightly better than the predicted worst case for FC regimes, ie, $\frac{2 + \lceil (32-1)/8 \rceil}{32} + \frac{1}{8} \approx 0.31$, according to Section 4.3.

Effect on throughput

The matrices in Figure 14 show that average cache misses per item vary by a large factor, identifying three classes of behaviors, namely, FP, FC, and the mixed regime observed when P and C both manage to perceive large batches. The impact of cache misses on the performance of a queue type depends on their actual cost compared to the other parts of the producer and consumer loop. Hence, we selected two specific combinations of emulated loads that are representative of FP and FC regimes, where all the queues experience FP or FC, respectively.[‡] For these workloads, we measured the throughput (items per second) achieved on our two platforms, I7 and XEON40, with largely different cache miss cost.

Figure 15 shows how the incremental optimizations to LQ and FFQ affect throughput on I7. In this experiment, FP uses $L_P = 0$ and $L_C = 10$ ns, whereas in FC, we have $L_P = 20$ ns and $L_C = 10$ ns. Since P and C are on the same socket, L1 cache misses are served very quickly by L2 or L3. Hence, although fewer misses give higher throughput, the difference between the various queue types is relatively small (10%-30% across all cases). For processing times higher than 50 ns, the differences are almost negligible (not shown here). Note how, in the FC case, BIFFQ is slightly slower than IFFQ, because the perceived batch is already larger than 32 for both, and the extra data copy makes the BIFFQ producer slower.

Figure 16 shows the throughput measured in the same FP and FC configurations on XEON40 (same values of L_P and L_C as for I7). Here, P and C are on different sockets, resulting in much more expensive cache misses. The difference has a remarkable effect on throughput, making the improved versions of the algorithms (LLQ, IFFQ) run at 170%-290% the speed of their basic counterparts. Here too, we note that in the FC regime, BIFFQ is slightly worse than IFFQ, for the same reasons.

Disabling batched operation

The two matrices in the third column of Figure 14 show the average number of cache misses with embedded payload when batching operation is disabled. In contrast to the other six cases, here, neither P nor C can publish/return items

[‡]Because of the differences in how the queues are implemented, it may happen that the same combination of emulated loads leads to different regimes for different queues.

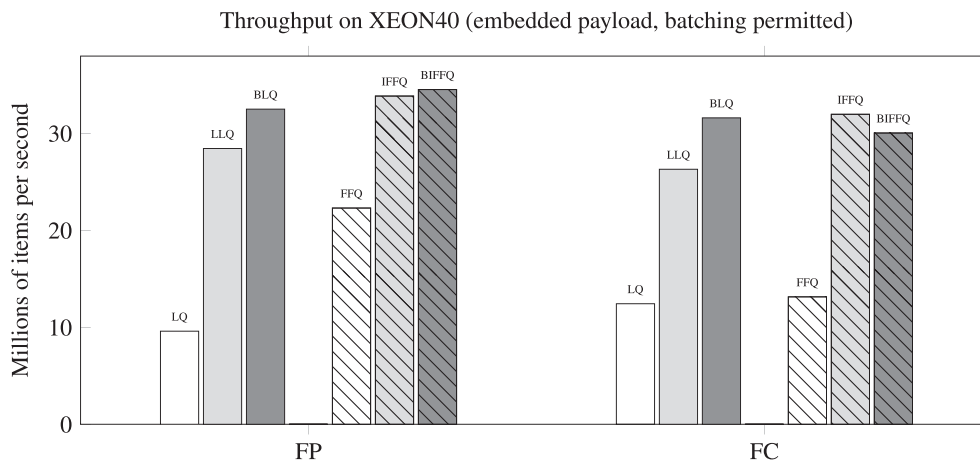


FIGURE 16 FP indicates a Fast Producer operating regime with emulated loads $L_P = 0$ and $L_C = 10$ (in nanoseconds), whereas FC is a Fast Consumer regime with $L_P = 20$ and $L_C = 10$. Reducing cache misses on XEON40 has more effect on throughput than on I7. BIFFQ, batched improved FastForward queue; BLQ, batched Lampport queue; FFQ, FastForward queue; IFFQ, improved FastForward queue; LLQ, lazy Lampport queue; LQ, Lampport queue

in batch; this constraint may arise because of the architectural limitations of the software that uses the SPSC queues, as explained in Section 2.1. This implies that the perceived batch for BLQ and BIFFQ, ie, the queues that expose batching capabilities, can be as small as 1 in the worst case. Note that the internal implementation of the queues can still batch their synchronization operations and, thus, improve performance. Both BLQ and BIFFQ show the same cache behavior as their corresponding nonbatching versions (LLQ and IFFQ, respectively), as expected. Regarding the effect on throughput, using an API capable of batching with a batch limit set to 1 is less efficient than using the simpler enqueue/dequeue; this effect is illustrated in Section 5.3. In any case, BLQ and BIFFQ should never be used when batching is not permitted.

5.2.2 | Indirect payload

Cache behavior

A queue with indirect payload requires extra memory access per item on both P and C, which results in an extra cache miss on both sides. This is perfectly reflected by our experiments, which result in a range of measured cache misses between 1.1 and 4 per item instead of 0.1 and 3 in the case of embedded payload. The average numbers of cache misses per item for indirect payload experiments on I7 are shown in Figure 17 with a color range [1, 4]. Apart from the extra cache miss, they resemble the results for embedded payload (Figure 14, color range [0, 3]). One interesting difference between the two figures is that the boundary between the FP and FC regions in Figure 17 is slightly concave, and the FP region is larger. This happens because on the consumer side, the extra read miss must wait for the slot read to complete, slowing down C more than the extra write miss slows down P. The effect is more visible for small values of L_C and L_P : for larger values, the emulated loads dominate the cache miss cost, and the boundary between FP and FC becomes again aligned to the main diagonal.

In any case, the interesting aspect confirmed by this experiment is that adding “external” cache misses unrelated to queue synchronization does not have an impact on how the SPSC queue under study behaves; such cost simply adds up to the rest of the workload (the emulated load in this case).

If batched operation is disabled, we obtain the results shown in the two matrices in the third column of Figure 17. Similarly to what was explained in Section 5.2.1, the cache behavior for BLQ and BIFFQ regresses to their corresponding nonbatching versions (LLQ and IFFQ).

Effect on throughput

Figure 18 illustrates the throughput measured with indirect payload for a particular combination of emulated loads, which results in an FP operating regime for all the queues. Similarly to what was already shown in Figure 15, LLQ/BLQ improve LQ, and IFFQ/BIFFQ improve FFQ. However, BIFFQ cannot improve IFFQ that much because it is already optimal in FP regimes. Moreover, BLQ is only a small improvement over LLQ, because the batch perceived by P is already close to $B_P = 32$ for LLQ.

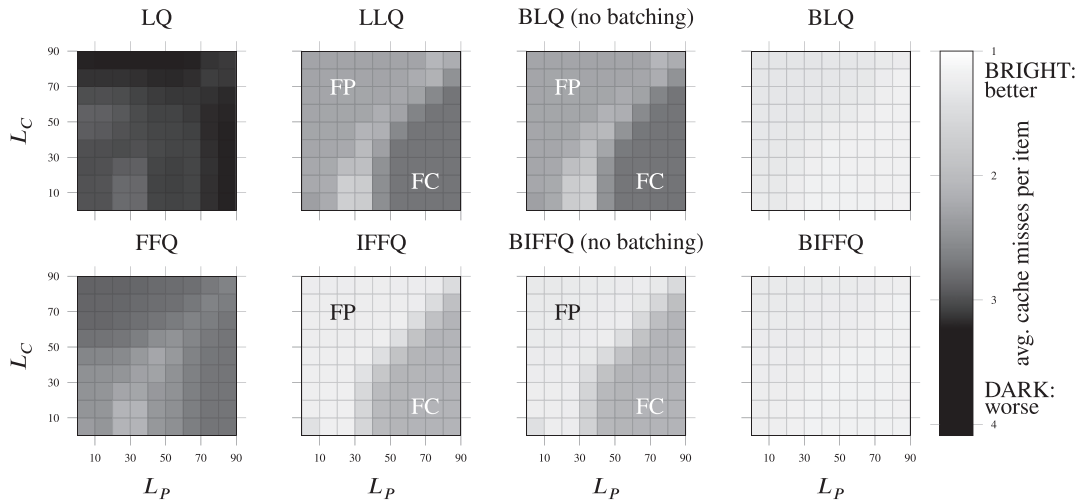


FIGURE 17 Average cache misses per item suffered by P during throughput experiments on I7 with indirect payload (brighter is better). Emulated load for both P and C ranges between 0 and 90 ns. In the plot for LLQ, we emphasize the two regions, Fast Producer (FP) and Fast Consumer (FC), which cause larger numbers of cache misses; details are explained in Sections 5.2.1 and 5.2.2. BIFFQ, batched improved FastForward queue; BLQ, batched Lamport queue; FFQ, FastForward queue; IFFQ, improved FastForward queue; LLQ, lazy Lamport queue; LQ, Lamport queue

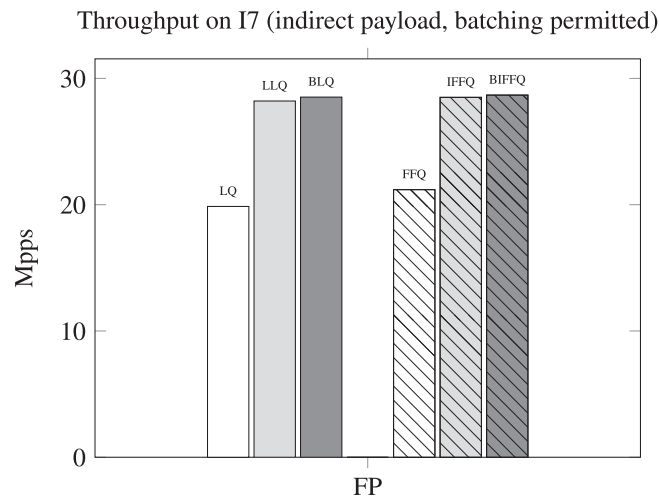


FIGURE 18 FP indicates a Fast Producer regime with $L_P = 0$ and $L_C = 10$ (in nanoseconds). Batched operation (batched Lamport queue [BLQ], batched improved FastForward queue [BIFFQ]) is not very useful as the improved FastForward queue (IFFQ) and the lazy Lamport queue (LLQ) are already able to generate a large perceived batch. FFQ, FastForward queue; LQ, Lamport queue

The same FP experiment on XEON40 is shown in Figure 19 and leads to similar observations, with a better relative throughput increment over the baseline LQ/FFQ queues (due to the higher cost of cache misses on XEON40). Figure 19 also shows an additional experiment, marked as NOLOAD because no emulated load is used. In such a configuration, all the queues experience a sharp FC regime, with C perceiving the minimum batch, ie, $B_P = 32$ for BLQ/BIFFQ and 1 for the others. In accordance with Table 1, the throughput for BLQ and BIFFQ is significantly larger than that for LLQ and IFFQ, respectively. LLQ and LQ show similar throughputs, because the LLQ producer is amortizing the load of the read variable on a small average batch (about five items). In this particular case, the IFFQ consumer experiences a slightly smaller perceived batch (three items) than FFQ (four items) and, thus also, a marginally worse throughput. On I7, we are not able to observe a sharp FC regime unless we use a large producer load (eg, $L_P = 70$ and $L_C = 0$). However, such a relatively large load dominates the cost of cache misses, so that the throughput differences between the queues become negligible; for this reason, we do not show an FC combination in Figure 18.

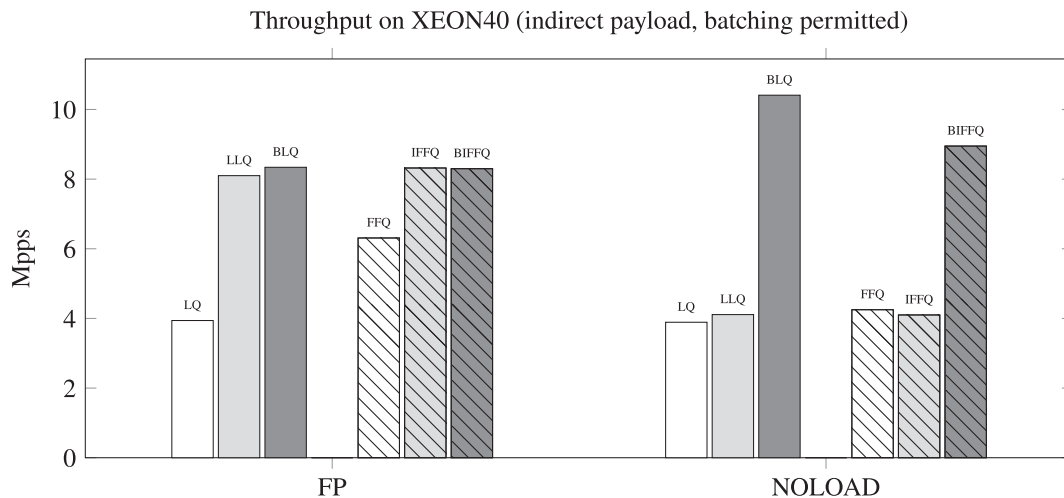


FIGURE 19 FP indicates a Fast Producer regime with $L_P = 0$ and $L_C = 10$ (in nanoseconds). The NOLOAD label refers to an experiment with no emulated loads, where all the queues happen to operate in a Fast Consumer regime. Batched operation for the batched Lamport queue (BLQ)/batched improved FastForward queue (BIFFQ) is effective whenever it guarantees a minimum batch larger than that perceived by the lazy Lamport queue (LLQ)/improved FastForward queue (IFFQ). FFQ, FastForward queue; LQ, Lamport queue

5.3 | Latency evaluation

The throughput experiments described in Section 5.2 are useful to validate the cache behavior of the SPSC queues under various configurations, resulting in different operating regimes. In all those configurations, P is always acting as a *streaming* producer, trying to achieve the maximum possible throughput. We now complete our experimental evaluation by measuring the worst-case latency introduced by the queue synchronization operations. We use two threads P1 and P2 connected by two SPSC queues in opposite directions. In each iteration, P1 writes an item (ie, a request) into the *request* queue, which P2 reads from. P2 reads the item and immediately enqueues it back to the *response* queue, where P1 reads it. P1 always waits for a pending response before enqueueing the next request; symmetrically, P2 always responds to a pending request before dequeuing the next request. In other words, P1 and P2 carry out a simple *request-response* (ping-pong) test in a loop, where P1 acts as a client and P2 as a server. This transaction-based workload is inspired to the TCP_RR test of the popular *netperf* benchmark tool,³¹ which is commonly used to measure network latency. The round-trip latency due to queue synchronization is the inverse of the transaction rate, since P1 and P2 do not perform significant operations other than enqueueing and dequeuing (except for accessing the packet payload once in the indirect payload case). Note that, with this workload, there is no possibility to batch, as both queues always store, at most, one item. P1 and P2 are forced to synchronize on each exchanged item on both the request and response queues. As a consequence, the regime experienced by any queue variant in these latency experiments must be exactly the same. This is in contrast with the throughput experiments of Section 5.2: even if two different queues are both experiencing an FP (or FC) operating regime, it is practically impossible that they observe the exact same evolution of queue occupation over time and, thus, the exact same sequence of perceived batches.

Although it is not possible for P1 and P2 to perceive a batch larger than one, we still include BLQ and BIFFQ in the experiments. This is useful to assess the overhead of using a batching API when batching is not possible. The same experiment infrastructure and methodology described in Section 5 (eg, 10 test runs of 10 seconds each) are also used here. Latency experiments use both embedded and indirect payloads and the latter mostly for the sake of completeness: because of how the experiment is designed, the additional cycles due to accessing the indirect payload simply add up to the baseline synchronization cost paid in the case of embedded payload. For similar reasons, emulated loads are not used. In more detail, P1 writes to the 32-bit integer field in the indirect payload before enqueueing to the request queue, and P2 reads the same 32-bit field after dequeuing from the request queue, but before enqueueing the item to the response queue. P1 does not perform more accesses to the indirect payload of the returned item. As a result, both P1 and P2 pay an additional cache miss per item with respect to the embedded payload scenario. Similarly to the *netperf* TCP_RR benchmark, rather than measuring latency directly, we count the number of request-response transactions that P1 and P2 manage to do with a given queue and compute an average rate. The transaction rate is inversely proportional to the queue latency; hence, the higher, the better, as in throughput experiments.

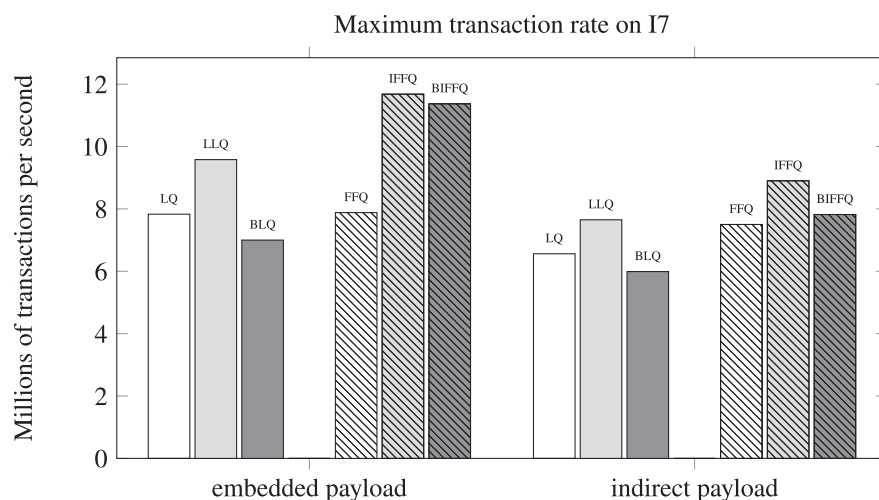


FIGURE 20 On I7, the improved FastForward queue (IFFQ) achieves a lower average latency as it pays less cache misses per item. Queues exposing a batching application programming interface (batched Lamport queue [BLQ] and batched improved FastForward queue [BIFFQ]) are inefficient because there are no batching opportunities. FFQ, FastForward queue; LLQ, lazy Lamport queue; LQ, Lamport queue

Transaction rates on I7 are shown in Figure 20. Values range from 6.5 million transactions per second (Mtps) to 11.5 Mtps, corresponding to latencies going from 87 to 154 ns per item. LQ has the worst performance overall, as this experiment deterministically triggers its worst case. Both P1 and P2 act as a producer and a consumer for each item, and thus, each of them must pay three cache misses on each enqueue or dequeue operation; this is confirmed by our measurements, which show an average of exactly six cache misses per item for both P1 and P2 with embedded payload and seven with indirect payload. LLQ reduces the cache misses to about two per operation, because the queue is always almost empty, and therefore, the test triggers the LLQ FC worst case (Table 1). LLQ measurements with embedded payload report 4.02 cache misses per item (both P1 and P2), with a transaction rate increase of 22% over LQ; indirect payload needs 5.09 cache misses per item. BLQ achieves a substantially lower throughput than LLQ (22% less with indirect payload), even if the measured cache miss rates are identical (as expected). This is an interesting outcome, and it can be explained with the higher number of instructions that both P1 and P2 must pay (twice) to use the BLQ batching API with a single element. Latency tests on FFQ report about 4.8 cache misses per item with embedded payload (5.9 with indirect payload) and, thus, a transaction rate lower than LLQ. Unfortunately, this result slightly disagrees with our analysis, which predicts two misses per item in the worst case and, at most, four misses for two operations. We believe this discrepancy is due to some quirk of the I7 machine, because the same does not happen on XEON40. Note that the FFQ cache behavior for the FC throughput experiments in Figure 14 perfectly agrees with the analysis, reporting an average of two misses or less; the discrepancy seems to affect only these latency tests. Nevertheless, IFFQ achieves the overall best performance, with an average of 2.16 misses per item with embedded payload, in line with the FC worst case predicted by the analysis. Similarly to BLQ, BIFFQ reports a lower throughput than IFFQ (12% less with indirect payload), because the additional overhead of the producer-side batching API is not amortized over a proper batch. The performance drop for BIFFQ is lower than it is for BLQ because the additional cost is lower, and it is paid on the producer side only. In summary, on I7, IFFQ is the best solution, with a latency that is 33% lower than LQ with embedded payload and 26% lower with indirect payload.

The corresponding latency numbers measured on XEON40 are illustrated in Figure 21. Cache misses are more expensive on the dual-socket machine, where P1 and P2 are pinned to different CPU sockets. As a consequence, transaction rates are lower, ranging from approximately 1.1 to 2.1 Mtps, corresponding to latencies between 476 and 909 ns. Nevertheless, the same discussion of I7 results is also mostly valid here. LQ has the worst performance in any case, whereas IFFQ performs best. Differently from I7, FFQ measurements here agree with the analysis, reporting about 3.9 cache misses per item with embedded payload and five with indirect payload. As a consequence, FFQ performs better than LLQ on XEON40. BLQ and BIFFQ have worse latency than LLQ and IFFQ, respectively. Overall, on XEON40, IFFQ latency is 43% lower than LQ latency with embedded payload and 31% lower with indirect payload. In conclusion, the results illustrated in this section meet the expectations of the analysis presented in Section 4.2: IFFQ is the optimal solution for latency-sensitive workloads.

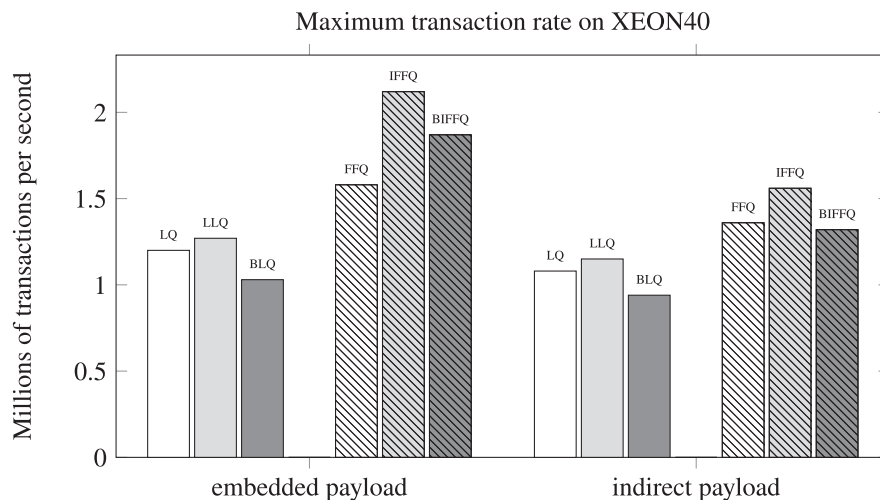


FIGURE 21 On XEON40, queues based on FastForward always achieve better average latencies. The batched Lamport queue (BLQ) and the batched improved FastForward queue (BIFFQ) perform poorly because of a lack of batching opportunities. FFQ, FastForward queue; IFFQ, improved FastForward queue; LLQ, lazy Lamport queue; LQ, Lamport queue

6 | AN EXAMPLE APPLICATION

The results presented in Section 5 are sufficient to validate our analysis (summarized in Table 1), but they do not evaluate the impact of choosing a particular SPSC queue in a real application. In particular, an application may need tens or hundreds of queues, and therefore, the benefit of using less cache misses is expected to grow with the number of queues. For this reason, we implemented an example application, *vswitch*, which uses a variable number of SPSC queues. The *vswitch* program is a *virtual switch* that performs software Ethernet switching among a set of N_C interconnected local client processes. Virtual switches are commonly used to interconnect virtual machines and implement data center networking. Each port of the virtual switch consists of two SPSC queues: a transmit queue for a process to transmit network packets and a receive queue to be used for packet reception.

The design of our switch is partially inspired by VALE,³² as it is able to operate in batch. The virtual switching logic runs in a dedicated thread, similarly to Open vSwitch³³ accelerated by DPDK²³ or to Snabb.³⁴ Using a single dedicated thread comes with some advantages; locking is not necessary to forward packets between ports, and it is feasible to busy-wait on the receive queues—as opposed to using a notification scheme—in order to minimize the worst-case forwarding latency. VALE follows a different approach, where the switching logic runs in the context of all the transmitting clients. This may enable better throughput (more than one CPU can be used for forwarding), but it requires locking and offers much less batching opportunities.

The *vswitch* application is implemented as an infinite loop that, in each iteration, executes a *pull* phase followed by a *push* phase, as depicted in Figure 22. In the *pull* phase, the switch scans the transmit queues of all the ports and dequeues up to $B_V > 0$ packets from each one. If less than B_V packets are available on a transmit queue, the switch dequeues them all and moves forward to the next queue. The destination port for each packet is computed by looking up the destination MAC address in the switch forwarding table. Packets are grouped by destination port, using a separate list for each port. These lists are private to the switch thread and are built using the *next* pointer field available in the packet metadata.[§] In the *push* phase, the per-port lists are drained by enqueueing the packets to the corresponding receive queues. Note that limiting the number of dequeue operations in the pull phase is useful to provide some degree of fairness to the interconnected clients and avoid aggressive clients to monopolize the *vswitch* thread. Moreover, the limitation is required to control the worst-case queue servicing latency, which would otherwise be unbounded. A smaller B_V corresponds to a smaller latency upper bound and better fairness. At the same time, it is convenient to choose B_V large enough, in such a way that the virtual switch can benefit from the batching capabilities of queues like BLQ, IFFQ, and BIFFQ. Batching operation is possible both when dequeuing from transmit queues and when enqueueing to receive queues, since packets to be enqueued to the

[§] Packet metadata are commonly used by network stack implementations to store information such as packet length, a pointer to the current network header, and a pointer to the next packet in a list. Examples are the FreeBSD `struct mbuf` and Linux `struct sk_buff`.

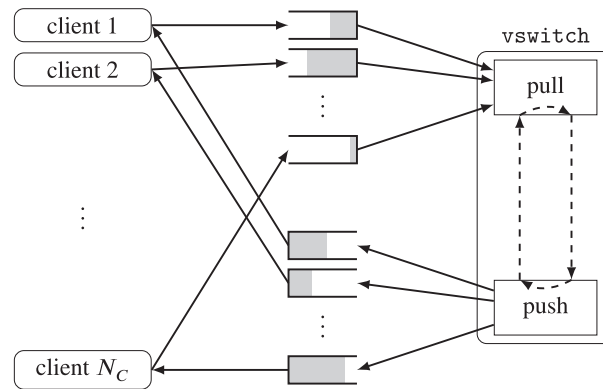


FIGURE 22 The `vswitch` example application. A dedicated thread *pulls* Ethernet packets from the clients' transmit queues and *pushes* them to the destination receive queues

same receive queues are collected in the same list. Note that VALE³² can only operate in batch with packets that come from the same source port, which means that batching will only happen if the client transmits several packets at once. In contrast, `vswitch` can build batches by aggregating packets originating from different source ports and directed toward the same destination port, which does not require the clients to transmit in batches. The switch does not perform packet copy when forwarding, which implies the use of indirect payload, as defined in Section 2. Network packets, including metadata, are dynamically allocated and deallocated by the clients. When a packet is enqueued to a transmit queue, packet ownership is passed to the switch, which will normally pass the ownership to the destination receive queue. Packets are dropped (deallocated) by the switch itself if the receive queues are full. Being a zero-copy switch, its performance is not dependent on packet size. The forwarding code only needs to access the packet metadata and the Ethernet header, which both fit in the first cache line (64 bytes) of the indirect payload. The experiments presented in Sections 6.2 and 6.3 use minimum-sized Ethernet packets (60 bytes) to enable faster packet initialization and, therefore, more aggressive clients, with the goal of stressing the switching thread as much as possible.

For the sake of simplicity, the forwarding table in the current `vswitch` prototype is prepared in advance, and it is never changed. In a real deployment, a protocol such as OpenFlow could be used to dynamically update the forwarding table. However, under the (common) assumption that updates are extremely infrequent with respect to look-ups, updates can be easily implemented in such a way as not to affect the forwarding performance in a measurable way. An optimized synchronization mechanism like RCU⁵—or even a custom strategy such as checking for updates once every second—would add negligible or nonmeasurable look-up overhead. The source code of the `vswitch` program is available online.²⁹

6.1 | Experiment methodology

Sections 6.2 and 6.3 report the measured forwarding rate of `vswitch` under two different workloads, different SPSC queues, and a variable number of clients. For these experiments, we use a less constrained execution environment than that described in Section 5.1 for validation experiments. In particular, we do not disable (nor try to deceive) the hardware prefetcher, and we reserve a physical core only for the `vswitch` thread (one of the two hyperthreads is left idle). Client threads run on both hyperthreads of the remaining available cores. These constraints have been removed because they are not appropriate for a real deployment. Nevertheless, the other measures described in Section 5.1 are desirable (and commonly used) in those deployments that need performance to be as deterministic as possible to meet service-level objectives. This is, for instance, the case of cloud computing providers that may want to provide guarantees on network latency and bandwidth to their customers. As a side effect, these precautions are beneficial for the reproducibility of our experiments.

The `vswitch` experiments are run only on XEON40, which has enough CPUs (40) to observe how the switch scales with an increasing number of clients. Each client is pinned to a different available CPU, and, at most, 38 clients are created. Preliminary experiments have shown that more clients sharing the same CPU result in lower aggregate throughput, because of the context switch overhead.

The measurement methodology is also very similar to that described in Section 5.1. A global variable is used to stop the virtual switch and the client threads, and each test run—for a given combination of N_C , SPSC queue, and workload—lasts

10 seconds, and it is repeated 10 times. The metric measured in these experiments is the total switch forwarding rate, which only requires the switch to keep an aggregate counter of packets pulled from the transmit queues. Note that it is important to use enough clients to saturate the CPU running the `vswitch` thread, so that more efficient SPSC queues can show an increased forwarding rate.

6.2 | Flooding experiments

In the first set of experiments, each client transmits as many packets as possible to all the other clients (*flooding*). A client is implemented as an infinite loop that, in each iteration, performs a *receive* phase followed by a *transmit* phase. During the *receive* phase, a client drains its receive queue, consuming (deallocating) all the available packets. In the *transmit* phase, a client allocates and enqueues $B_C > 0$ packets to its transmit queue, with B_C being the client batch. The destination MAC address for each allocated packet is selected in a round-robin fashion among all the possible clients. Packets allocated but not enqueued because the transmit queue is full are recycled for the next iteration in order to save CPU cycles.

Figure 23 (left) reports the aggregate forwarding rate of `vswitch` for $B_V = 8$ and $B_C = 1$. Since clients transmit a single packet in each iteration, they cannot enqueue in batch to the transmit queue. However, there are other batching opportunities: clients can dequeue in batch from the receive queues, and the switch can both dequeue and enqueue in batch from the transmit and receive queues, respectively. The switch throughput increases up to eight clients in a similar way for all the queues. For more clients, the throughput starts to saturate at different values for different queues. Fluctuations depend on how much the clients and the switch are able to batch for the given input workload, on the amount of packet drops, and on the combination of operating regimes of the many queues involved. The switch in saturation acts as a bottleneck: transmit queues tend to operate in an FP regime (queues generally full), and receive queues tend to operate in an FC regime (queues generally empty). The plot shows that IFFQ and BIFFQ have a consistently higher throughput than the others. This is expected because they are optimal in the FP case, and the efficiency of the dequeue operations is less affected by batching than it is for BLQ and LLQ. The effectiveness of the BIFFQ optimization for enqueue operations is also visible, and it explains the throughput improvement over IFFQ. The BLQ throughput is consistently higher than LQ and LLQ, which confirms that switch and clients are able to operate in batch. Figure 23 (right) shows the measured forwarding rate with $B_V = B_C = 8$. The plot is very similar to that on the left, with the values being generally higher for LLQ, BLQ, IFFQ, and BIFFQ because of the additional batching opportunities.

6.3 | Request-response experiments

The second set of experiments stresses `vswitch` with a request-response workload. A client is implemented as an infinite loop that, in each iteration, executes a *request* phase followed by a *response* phase. In the *request* phase, a client transmits B_C requests to the other clients, selected in a round-robin fashion. During the *response* phase, the client waits for the

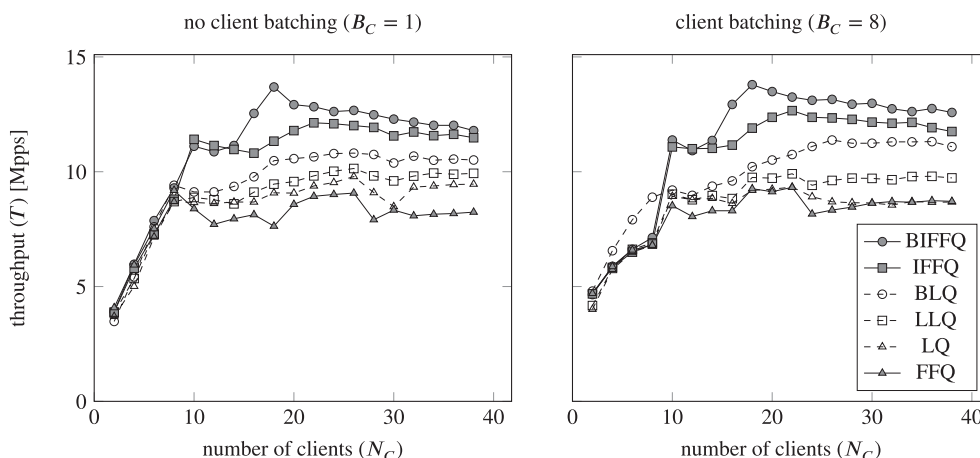


FIGURE 23 Throughput of the virtual switch when flooded by clients, with $B_V = 8$ and either nonbatching clients ($B_C = 1$, left plot) or batching clients ($B_C = 8$, right plot). The improved FastForward queue (IFFQ) and the batched improved FastForward queue (BIFFQ) provide the highest throughput when the switch is in saturation. BLQ, batched Lamport queue; FFQ, FastForward queue; LLQ, lazy Lamport queue; LQ, Lamport queue

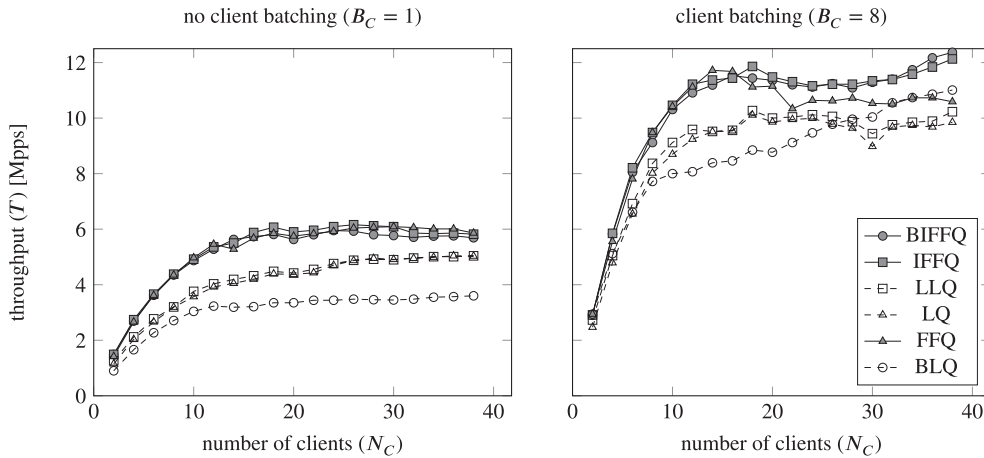


FIGURE 24 Throughput of the virtual switch under a request-response workload, with $B_V = 8$ and either nonbatching clients ($B_C = 1$, left plot) or batching clients ($B_C = 8$, right plot). On switch saturation, queues based on FastForward provide a higher throughput than the Lamport queue (LQ), lazy Lamport queue (LLQ), and batched Lamport queue (BLQ). The BLQ provides the poorest performance when there are no chances for batching. BIFFQ, batched improved FastForward queue; FFQ, FastForward queue; IFFQ, improved FastForward queue

corresponding B_C responses to arrive (in any order), by busy-waiting on its receive queue. Note that this workload is relevant to how scalable network services are designed: a frontend server may receive a request from a remote client machine and issue N subrequests (*outcast*) to backend servers. Once the corresponding subresponses arrive (*incast*), the frontend server can reply to the remote client. Batching opportunities are the same as those described in Section 6.2. However, note that, here, there is a form of flow control, because a client waits for the responses to come before proceeding with the next requests. In particular, the receive queue of any client will never contain more than $N_C \times \lceil B_C/N_C \rceil$ items, that is, at most, $\lceil B_C/N_C \rceil$ packets from each requestor. This is actually a worst case: receive queues are expected to contain, at most, B_C packets most of the time. In our experiments, $B_C \leq 8$ and $N_C \leq 38$, which means that receive queues will never overflow if the queue length is $N = 256$, and therefore, *vswitch* will never drop packets.

Figure 24 (left) illustrates the aggregate forwarding rate of the virtual switch for $B_V = 8$ and $B_C = 1$. The batching opportunities in this case are very limited, with all the queues containing zero or one packet most of the time. It follows that this particular experiment is very similar to that described in Section 5.3, although with many queues. Queues based on FastForward achieve significantly higher throughput (~ 6 Mpps in saturation) than the ones derived from Lamport's work (~ 4 Mpps), because of the lower amount of cache misses in the worst case. BIFFQ performs slightly worse than IFFQ because of the additional overhead of the producer-side buffering. It is not surprising that FFQ is very competitive in this particular experiment, especially with higher numbers of clients: FFQ enqueue and dequeue routines contain very few instructions, and neither IFFQ nor BIFFQ—both requiring more instructions per operation—can leverage batching. For similar reasons, BLQ performs worse than LLQ and LQ because of the higher overhead of its batching API. Note that the BLQ throughput gap increases as the number of clients increases. As expected, these results agree with the ones shown in Figures 21 and 20, confirming that the two experiments are of a similar nature. Figure 24 (right) shows the aggregated *vswitch* throughput for $B_V = 8$ and $B_C = 8$. In this case, the chances for batching are slightly higher, and therefore, the throughputs in saturation are higher, ranging from approximately 8 to 12 Mpps. IFFQ, BIFFQ, and FFQ provide higher forwarding rates than LQ, LLQ, and BLQ also in this case. As N_C increases, FFQ becomes worse than its variants because it cannot benefit from batching. Note that as N_C increases, BLQ becomes more efficient than LLQ and LQ because of the higher chances of batching operation.

7 | RELATED WORK

Lamport²⁰ was the first to propose a concurrent lock-free (CLF) queue, in order to efficiently decouple a single-producer thread from a single-consumer thread without resorting to locks. The LQ algorithm described in Section 3.1 is a modern implementation of the original CLF queue. Although Lamport proved the correctness of the queue only under the *sequential consistency* memory model,³⁵ the addition of proper memory barriers makes it correct for any memory model. Moreover, the queue slots and control variables can be laid out in separate cache lines to avoid cache thrashing, as already

suggested by many authors.^{15–19,21} LQ is therefore the natural baseline for the evaluation of the other PSC algorithms presented in Sections 3 and 4. The lazy loading optimization introduced in Section 3.2 is also used by DBLS,¹⁴ MCRingBuffer,¹⁵ and Liberty.¹⁸ FastForward²¹ showed that it is possible to tackle the same problem addressed by Lamport with reduced cache misses, by embedding the synchronization information within the queue slots. The FFQ algorithm presented in Section 4.1 is an implementation of the baseline FFQ. Moreover, FastForward addresses the indirect payload case and proposes an adaptive temporal slipping algorithm to make sure that the producer and the consumer never work on the same cache line. We do not include this technique in our FFQ as it would break the assumptions of Section 2.

More recent works try to go beyond FastForward and boost the queue throughput at the cost of unbounded latency. For this purpose, in addition to the lazy loading optimization, the LQ can be modified to publish or return items in batches of fixed size, regardless of the relative speed between the producer and the consumer. For instance, a producer could update the `write` control variable once every 64 enqueue operations or more. Adopting the terminology introduced by Lynx,¹⁹ these optimizations define the concept of *multisection* queue (MSQ), where the array of slots is partitioned into multiple sections of equal size, and synchronization between P and C only happens at section boundaries. These batching techniques are used, for instance, by MCRingBuffer,¹⁵ BatchQueue,¹⁶ Zhang's queue,¹⁷ DBLS,¹⁴ and Liberty.¹⁸ The fixed batch size gives strong guarantees on the maximum cache misses per item, and without the need to expose batch capabilities in the API, but it is acceptable only assuming a streaming producer. Our BLQ, presented in Section 3.3, does not make such an assumption in order to keep latency under control, but it is able to deliver comparable throughput at least while the producer is actually streaming and can operate in batch. Moreover, these MSQ algorithms assume to deal with embedded payloads and do not take into consideration the issues related to indirect payloads (eg, the need for memory barriers). Liberty¹⁸ also uses nontemporal writes to bypass the cache subsystem and, thus, avoid the overheads due to cache coherence. However, performing some preliminary experiments with nontemporal writes for our use cases, we have observed that they reduce both throughput and latency (as data need to go through central memory), and they are therefore not suited for the very high rates that we target. Lynx¹⁹ further specializes the MSQ approach by completely removing most of the CPU instructions from the critical path of the enqueue and dequeue implementation, namely, the mask operation (needed to rotate the `write` or `read` index), the comparison and jump instructions to check if the index reached the boundary of a section, together with the synchronization code itself. Leveraging CPU exceptions and operating system support, Lynx executes the removed operations in the context of a signal handler triggered by a CPU exception; the exception is triggered whenever the enqueue or dequeue code reaches a section boundary, which is marked nonaccessible in the process page tables. These extreme optimizations allow for huge rates, peaking to 2 billion items per second, according to the authors. Drawbacks include a higher worst-case latency—because of the CPU exception and signal handler overhead—and limitations on the queue size (eg, sections must be at least two memory pages).

Aldinucci et al²² proposed dynamically growing and shrinking *unbounded* SPSC queues, which are useful to avoid deadlocks with complex processing graphs. Their strategy is to combine multiple array-based queues in a linked list, in order to avoid a fixed-size queue while still preserving the optimal performance of array-based queues. These techniques are orthogonal to our study and can be applied to any queue described in this paper.

Several lock-free and wait-free algorithms have been designed to address the more general problem of Multiple-Producer–Multiple-Consumer (MPMC) queues. The seminal work of Michael and Scott³⁶ presents a simple concurrent and lock-free MPMC queue based on the CAS instruction. CAS atomic operations are used to atomically update the head and tail pointers of the queue, which is implemented as a linked list. Since CAS can fail, however, progress is not guaranteed for all the competing threads in a finite amount of steps (ie, the queue is not wait-free), and high contention on the head and tail pointers can cause severe performance degradation. Later proposals try to outperform Michael and Scott's basic queue by reducing or amortizing the CAS contention in several ways.^{17,37–41} As an alternative approach, the more recently available FAA instruction can also be used to handle contention on the queue insertion and extraction points, as showed by Morrison and Afek.⁴² Differently from CAS, FAA is guaranteed to progress for all the competing threads; this property leads to improved performance over CAS-based solutions. A recent work by Yang and Mellor-Crummey⁴¹ extends Morrison and Afek's lock-free queue to provide wait-free guarantees. However, the performance of these MPMC queues in case there is a single producer and a single consumer is significantly lower compared to the ones presented in our work. This is not surprising, as MPMC implementations need to be more sophisticated to address significantly harder problems, such as ABA,³⁶ scaling to hundreds of threads, or guaranteeing wait-freedom.⁴¹ As a consequence of their requirements, all of the abovementioned MPMC queues use dynamically allocated memory, linked list data structures (sometimes combined with arrays of slots), and/or *fast-path–slow-path* techniques.⁴³ When compared to SPSC queues, such a sparse memory layout and the need for memory reclamation schemes often lead to worse worst-case behavior in terms of cache misses and per-operation overhead.

Finally, an alternative approach to operate in batch is to use *futures*, as proposed by Kogan and Herlihy⁴⁴ and later by Milman et al,⁴⁵ in the context of MPMC queues. A thread accessing the queue can accumulate enqueue and dequeue

operations locally (in a private data structure) and force their actual *evaluation* on the queue only at a later point or when it is really necessary. By relaxing the operation ordering guarantees perceived by the interacting threads, some locally pending operations can be eliminated (eg, a pending dequeue can cancel a previous pending enqueue) or combined with each other to access the shared data structure less frequently, hence reducing contention. Although futures are useful to increase the degree of parallelism in the case of MPMC queues, the additional cost involved in maintaining the future objects makes them inconvenient for our SPSC workloads.

8 | CONCLUSION

In this study, we have described and analyzed six general-purpose high-speed SPSC lock-free queues, each one coming with a different degree of optimization. Three of them are based on the original Lamport's CLF queue. The other three queues are based on FastForward, with the improved versions—IFFQ and BIFFQ—detailed in this paper. Each queue has different properties in terms of cache misses, showing different behaviors under different conditions. The experiments carried out on both single- and dual-socket machines have validated the behaviors predicted by the theoretical analysis. We have also developed a virtual switch example to evaluate the efficiency of the different SPSC queues in a real application.

Minimizing the cache misses related to producer/consumer synchronization is important to increase throughput and reduce latency, especially when a queue needs to process millions of items per second or more. Our investigation has shown how the ability to operate in batch or at least to amortize synchronization-related cache misses over many items is key to achieving large throughput improvements. The BLQ achieves the lowest number of cache misses and the highest throughput in case the producer and consumer can operate in large batches, as it provides strict guarantees on the minimum batch perceived by producers and consumers. The BIFFQ may need slightly more cache misses and, therefore, provide a lower throughput than BLQ, but it behaves very well in applications with many queues and limited batching opportunities. On the other hand, the IFFQ delivers the lowest latency, because embedding the control information within the array of slots reduces the per-item cache misses needed in the worst case. Finally, if the producer or consumer cannot operate in batch, IFFQ provides the best worst-case cache behavior and achieves the highest throughput in our experiments.

ORCID

Vincenzo Maffione  <https://orcid.org/0000-0003-0020-8664>

REFERENCES

1. Amazon. Amazon EC2 instance types. <https://aws.amazon.com/ec2/instance-types/>
2. Machine types available on Google Cloud. <https://cloud.google.com/compute/docs/machine-types>
3. Virtual machines types available on Azure. <https://docs.microsoft.com/en-us/azure/virtual-machines/windows/sizes>
4. Gamsa B, Krieger O, Appavoo J, Stumm M. Tornado: maximizing locality and concurrency in a shared memory multiprocessor operating system. In: Proceedings of the 3rd USENIX Symposium on Operating Systems Design and Implementation (OSDI); 1999; New Orleans, LA.
5. McKenney PE, Slingwine JD. Read-copy update: using execution history to solve concurrency problems. In: Proceedings of the 10th IASTED International Conference on Parallel and Distributed Computing and Systems; 1998; Las Vegas, NV.
6. Herlihy M. Wait-free synchronization. *ACM Trans Program Lang Syst*. 1991;13(1):124-149.
7. Jensen EH, Hagensen GW, Broughton JM. *A New Approach to Exclusive Data Access in Shared Memory Multiprocessors*. Technical report UCRL-97663. Livermore, CA: Lawrence Livermore National Laboratory; 1987.
8. McKenney PE. Memory barriers: a hardware view for software hackers. Beaverton, OR: Linux Technology Center, IBM Beaverton; 2010.
9. Rizzo L, Valente P, Lettieri G, Maffione V. PSPAT: software packet scheduling at hardware speed. *Comput Commun*. 2018;120:32-45.
10. Fraser K. *Practical Lock-Freedom*. Technical report UCAM-CL-TR-639. Beaverton, OR: Cambridge, UK 2004.
11. Valois JD. Lock-free linked lists using compare-and-swap. In: Proceedings of the 14th Annual ACM Symposium on Principles of Distributed Computing (PODC); 1995; Ottawa, Canada.
12. Herlihy M. A methodology for implementing highly concurrent data objects. *ACM Trans Program Lang Syst*. 1993;15(5):745-770.
13. ETSI. Network functions virtualisation (NFV); architectural framework. GS NFV 002 (v. 1.1.1). 2013.
14. Wang C, Kim H-s, Wu Y, Ying V. Compiler-managed software-based redundant multi-threading for transient fault detection. In: Proceedings of the International Symposium on Code Generation and Optimization (CGO); 2007; San Jose, CA.
15. Lee PPC, Bu T, Chandranmenon G. Lock-free, cache-efficient multi-core synchronization mechanism for line-rate network traffic monitoring. Paper presented at: IEEE International Symposium on Parallel Distributed Processing (IPDPS); 2010; Atlanta, GA.

16. Preud'homme T, Sopena J, Thomas G, Folliot B. BatchQueue: fast and memory-thrifty core to core communication. Paper presented at: 22nd International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD); 2010; Petrópolis, Brazil.
17. Zhang Y, Ootsu K, Yokota T, Baba T. Clustered communication for efficient pipelined multithreading on commodity MCPs. *IAENG Int J Comput Sci*. 2009;36(4):275-283.
18. Jablin TB, Zhang Y, Jablin JA, Huang J, Kim H, August DI. Liberty queues for epic architectures. In: Proceedings of the 8th Workshop on Explicitly Parallel Instruction Computing Architectures and Compiler Technology (EPIC-8); 2010; Toronto, Canada.
19. Mitropoulou K, Porpodas V, Zhang X, Jones TM. Lynx: using OS and hardware support for fast fine-grained inter-core communication. In: Proceedings of the 2016 International Conference on Supercomputing (ICS); 2016; Istanbul, Turkey.
20. Lamport L. Specifying concurrent program modules. *ACM Trans Program Lang Syst (TOPLAS)*. 1983;5(2):190-222.
21. Giacomoni J, Moseley T, Vachharajani M. FastForward for efficient pipeline parallelism: a cache-optimized concurrent lock-free queue. In: Proceedings of the 13th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP); 2008; Salt Lake City, UT.
22. Aldinucci M, Danelutto M, Kilpatrick P, Meneghin M, Torquati M. An Efficient unbounded lock-free queue for multi-core systems. Paper presented at: European Conference on Parallel Processing (Euro-Par); 2012; Rhodes Island, Greece.
23. Intel. Data plane development kit. 2012. <http://edc.intel.com/Link.aspx?id=5378>
24. Rizzo L. netmap: a novel framework for fast packet I/O. In: Proceedings of the 2012 USENIX Annual Technical Conference (USENIX ATC); 2012; Boston, MA.
25. Deri L. PF_RING DNA. https://www.ntop.org/products/packet-capture/pf_ring/
26. Torrellas J, Lam JS, Hennessy JL. False sharing and spatial locality in multiprocessor caches. *IEEE Trans Comput*. 1994;43(6):651-663.
27. Rizzo L, Garzarella S, Lettieri G, Maffione V. A study of speed mismatches between communicating virtual machines. In: Proceedings of ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS); 2016; Santa Clara, CA.
28. Intel 64 and IA-32 architectures optimization reference manual. 2009. <https://www.intel.com/content/dam/www/public/us/en/documents/manuals/64-ia-32-architectures-optimization-manual.pdf>
29. Maffione V. SPSCQ library and example programs. <https://github.com/vmaffione/spscq>
30. Lettieri G, Maffione V, Rizzo L. A study of I/O performance of virtual machines. *The Comput J*. 2018;61(6):808-831.
31. The netperf source code. <https://github.com/HewlettPackard/netperf>
32. Rizzo L, Lettieri G. VALE, a switched ethernet for virtual machines. In: Proceedings of the 8th International Conference on Emerging Networking Experiments and Technologies (CoNEXT); 2012; Nice, France.
33. Pfaff B, Pettit J, Koponen T, et al. The design and implementation of open vSwitch. In: Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation (NSDI); 2015; Oakland, CA.
34. Paolino M, Nikolaev N, Fanguede J, Raho D. SnabbSwitch user space virtual switch benchmark and performance optimization for NFV. Paper presented at: 2015 IEEE Conference on Network Function Virtualization and Software Defined Network (NFV-SDN); 2015; San Francisco, CA.
35. Lamport L. How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Trans Comput*. 1979;C-28(9):690-691.
36. Michael MM, Scott ML. Nonblocking algorithms and preemption-safe locking on multiprogrammed shared memory multiprocessors. *J Parallel Distributed Comput*. 1998;51(1):1-26.
37. Tsigas P, Zhang Y. A simple, fast and scalable non-blocking concurrent FIFO queue for shared memory multiprocessor systems. In: Proceedings of the Thirteenth Annual ACM Symposium on Parallel Algorithms and Architectures (SPAA); 2001; Crete Island, Greece.
38. Ladan-Mozes E, Shavit N. An optimistic approach to lock-free FIFO queues. In: Proceedings of the 18th International Symposium on Distributed Computing (DISC); 2004; Amsterdam, The Netherlands.
39. Moir M, Nussbaum D, Shalev O, Shavit N. Using elimination to implement scalable and lock-free FIFO queues. In: Proceedings of the Seventeenth Annual ACM Symposium on Parallelism in Algorithms and Architectures (SPAA); 2005; Las Vegas, NV.
40. Scherer III WN, Lea D, Scott ML. Scalable synchronous queues. In: Proceedings of the Eleventh ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP); 2006; New York, NY.
41. Yang C, Mellor-Crummey J. A wait-free queue as fast as fetch-and-add. In: Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP); 2016; Barcelona, Spain.
42. Morrison A, Afek Y. Fast concurrent queues for x86 processors. *ACM SIGPLAN Not*. 2013;48(8):103-112.
43. Kogan A, Petrunk E. A methodology for creating fast wait-free data structures. *ACM SIGPLAN Not*. 2012;47(8):141-150.
44. Kogan A, Herlihy M. The future(s) of shared data structures. In: Proceedings of the 2014 ACM Symposium on Principles of Distributed Computing (PODC); 2014; Paris, France.
45. Milman G, Kogan A, Lev Y, Luchangco V, Petrunk E. BQ: a lock-free queue with batching. In: Proceedings of the 30th Symposium on Parallelism in Algorithms and Architectures (SPAA); 2018; Vienna, Austria.

How to cite this article: Maffione V, Lettieri G, Rizzo L. Cache-aware design of general-purpose Single-Producer-Single-Consumer queues. *Softw: Pract Exper*. 2019;49:748-779. <https://doi.org/10.1002/spe.2675>