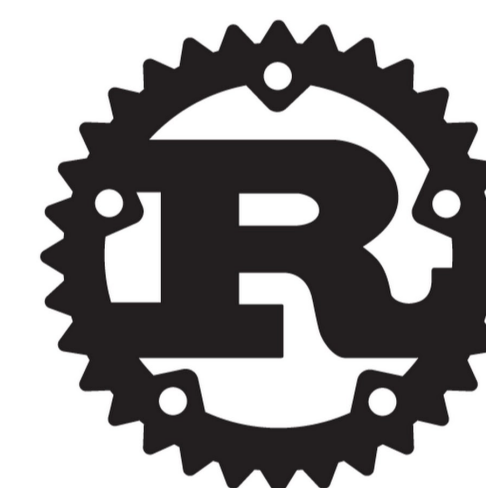# BBQ: A Block-based Bounded Queue
# for Exchanging Data and Profiling

Jiawei Wang,  **Diogo Behrens**,  Ming Fu,  Lilith Oberhauser,
Jonas Oberhauser,  Jitang Lei,  Geng Chen,  Hermann Härtig,  Haibo Chen

# Bounded queues (aka ring buffers) are everywhere...

# Why are they important to us?

Crucial for the
**performance and correctness**
of systems and applications!
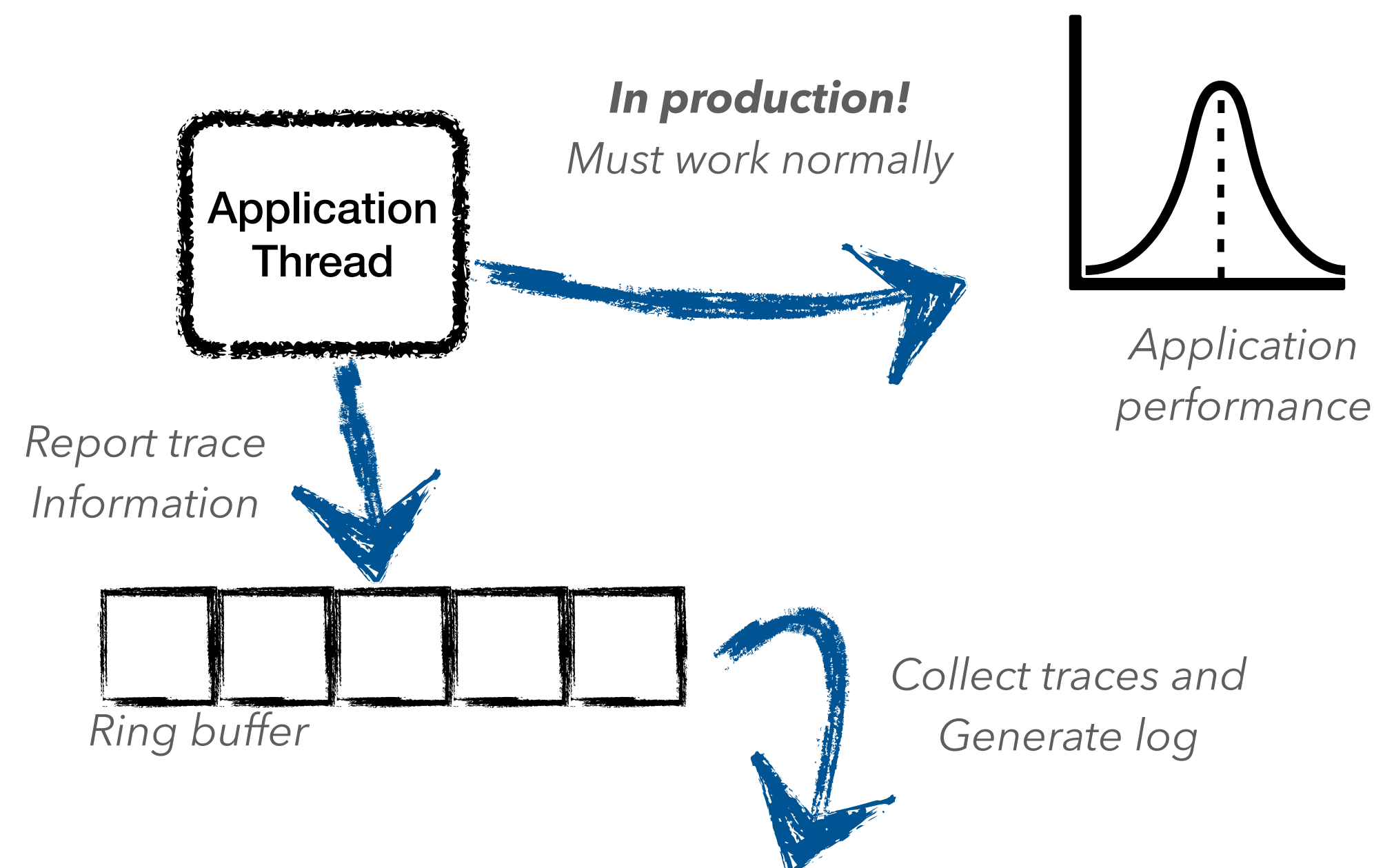
# Why are they important to us?

Crucial for the
**performance and correctness**
of systems and applications!

Next, **3 ring buffer stories**
from Huawei software development

# Story 1: Tracing overhead and operation interference

In-house OS with a **new tracing tool**

- Ring buffer used to collect traces

- Used to generate application profile

- Reporting **must be fast!**



*In production!*
*Must work normally*

Application
Thread

*Application
performance*

*Report trace
Information*

*Collect traces and
Generate log*

*Ring buffer*

```
openat(AT_FDCWD, "/usr/include/x86_64-linux-gnu/bits/stdint-intn.h", O_RDONLY|O_CLOEXEC) = 3
readlink("/proc/self/fd/3", "/usr/include/x86_64-linux-gnu/bi"..., 4096) = 48
fstat(3, {st_mode=S_IFREG|0644, st_size=1036, ...}) = 0
pread64(3, "/* Define intN_t types.\n    Copyr"..., 1036, 0) = 1036
rt_sigprocmask(SIG_SETMASK, ~[RTMIN RT_1], [], 8) = 0
close(3)                                = 0
rt_sigprocmask(SIG_SETMASK, [], NULL, 8) = 0
openat(AT_FDCWD, "/usr/local/include/genmc/bits/stdint-uintn.h", O_RDONLY|O_CLOEXEC) = -1 ENOENT
openat(AT_FDCWD, "/usr/include/x86_64-linux-gnu/bits/stdint-uintn.h", O_RDONLY|O_CLOEXEC) = 3
readlink("/proc/self/fd/3", "/usr/include/x86_64-linux-gnu/bi"..., 4096) = 49
fstat(3, {st_mode=S_IFREG|0644, st_size=1048, ...}) = 0
pread64(3, "/* Define uintN_t types.\n    Copy"..., 1048, 0) = 1048
rt_sigprocmask(SIG_SETMASK, ~[RTMIN RT_1], [], 8) = 0
close(3)                                = 0
rt_sigprocmask(SIG_SETMASK, [], NULL, 8) = 0
brk(0x55c3618b9000)                     = 0x55c3618b9000
brk(0x55c3618da000)                     = 0x55c3618da000
futex(0x7f185315087c, FUTEX_WAKE_PRIVATE, 2147483647) = 0
futex(0x7f185314b458, FUTEX_WAKE_PRIVATE, 2147483647) = 0
```
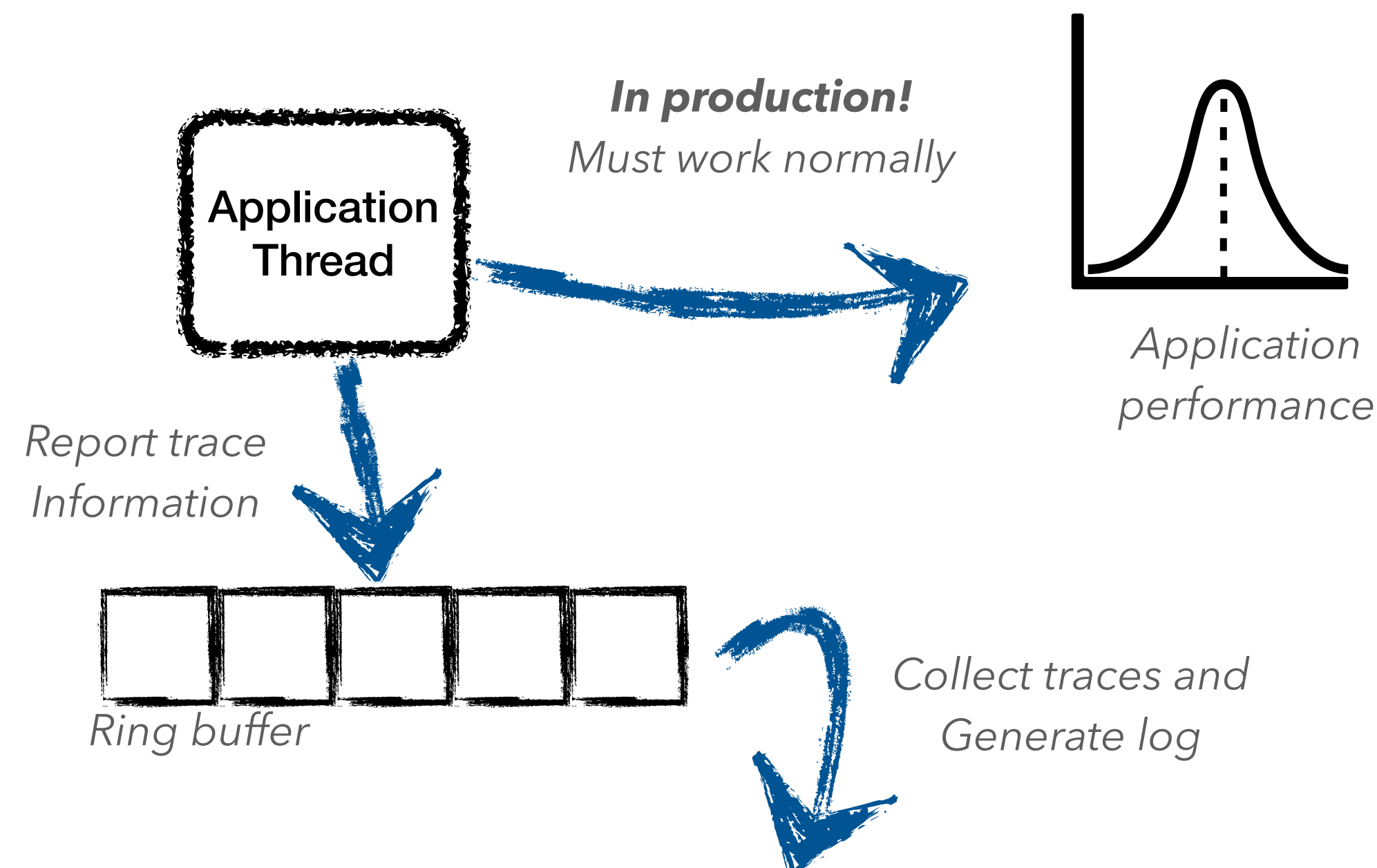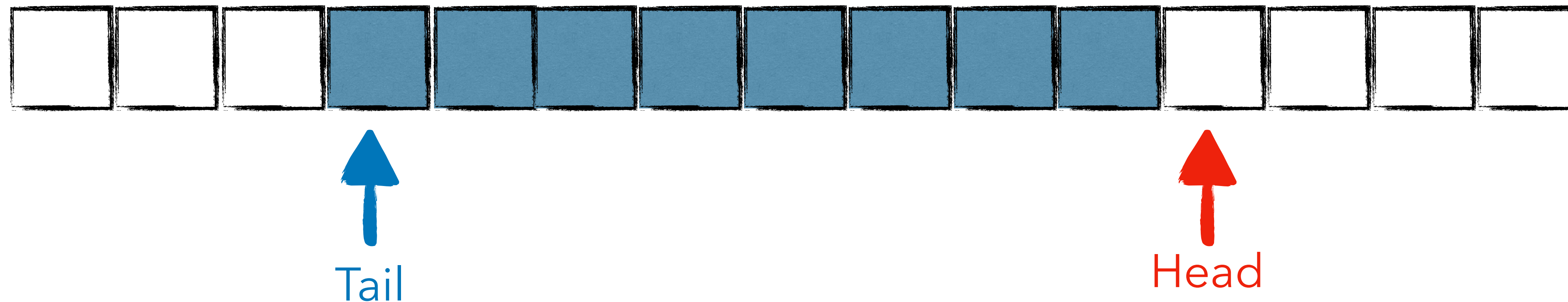
*Trace output*

# Story 1: Tracing overhead and operation interference

In-house OS with a **new tracing tool**

- Ring buffer used to collect traces

- Used to generate application profile

- Reporting **must be fast!**

Problem

- Consumer **slowdowns** producer!



*In production!*
*Must work normally*

Application
Thread

*Application
performance*

*Report trace
Information*

*Ring buffer*

*Collect traces and
Generate log*

```
openat(AT_FDCWD, "/usr/include/x86_64-linux-gnu/bits/stdint-intn.h", O_RDONLY|O_CLOEXEC) = 3
readlink("/proc/self/fd/3", "/usr/include/x86_64-linux-gnu/bi"..., 4096) = 48
fstat(3, {st_mode=S_IFREG|0644, st_size=1036, ...}) = 0
pread64(3, "/* Define intN_t types.\n   Copyr"..., 1036, 0) = 1036
rt_sigprocmask(SIG_SETMASK, ~[RTMIN RT_1], [], 8) = 0
close(3)                                = 0
rt_sigprocmask(SIG_SETMASK, [], NULL, 8) = 0
openat(AT_FDCWD, "/usr/local/include/genmc/bits/stdint-uintn.h", O_RDONLY|O_CLOEXEC) = -1 ENOENT
openat(AT_FDCWD, "/usr/include/x86_64-linux-gnu/bits/stdint-uintn.h", O_RDONLY|O_CLOEXEC) = 3
readlink("/proc/self/fd/3", "/usr/include/x86_64-linux-gnu/bi"..., 4096) = 49
fstat(3, {st_mode=S_IFREG|0644, st_size=1048, ...}) = 0
pread64(3, "/* Define uintN_t types.\n   Copy"..., 1048, 0) = 1048
rt_sigprocmask(SIG_SETMASK, ~[RTMIN RT_1], [], 8) = 0
close(3)                                = 0
rt_sigprocmask(SIG_SETMASK, [], NULL, 8) = 0
brk(0x55c3618b9000)                     = 0x55c3618b9000
brk(0x55c3618da000)                     = 0x55c3618da000
futex(0x7f185315087c, FUTEX_WAKE_PRIVATE, 2147483647) = 0
futex(0x7f185314b458, FUTEX_WAKE_PRIVATE, 2147483647) = 0
```
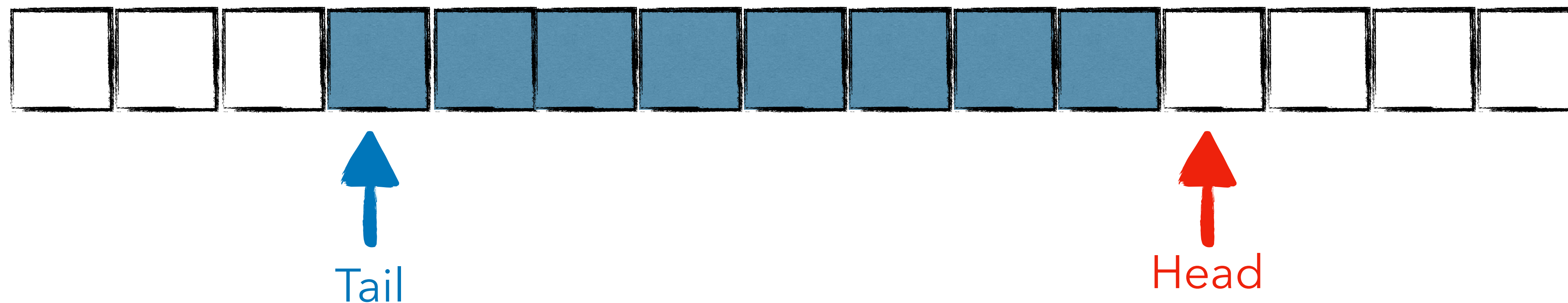
*Trace output*

# Interference sources

Ring buffers are arrays with indices



Tail

Head

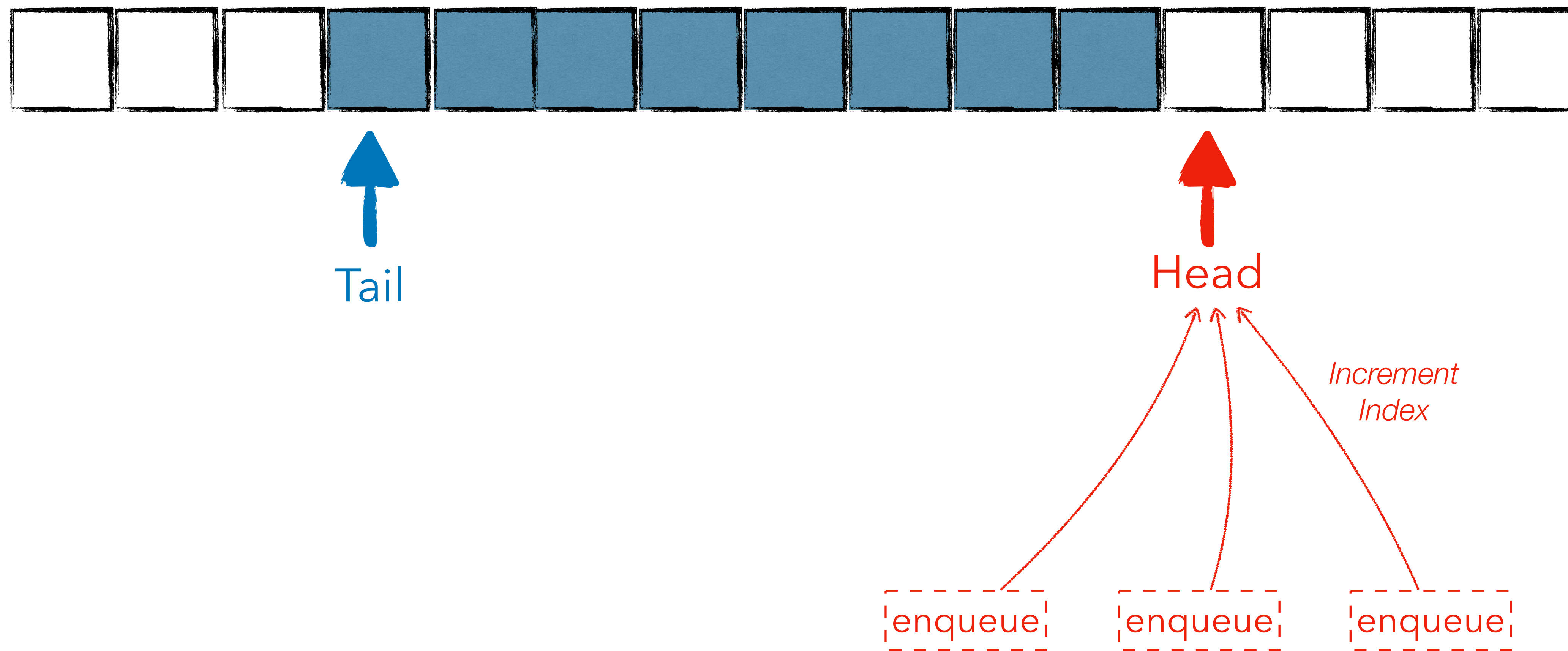# Interference sources

Ring buffers are arrays with indices



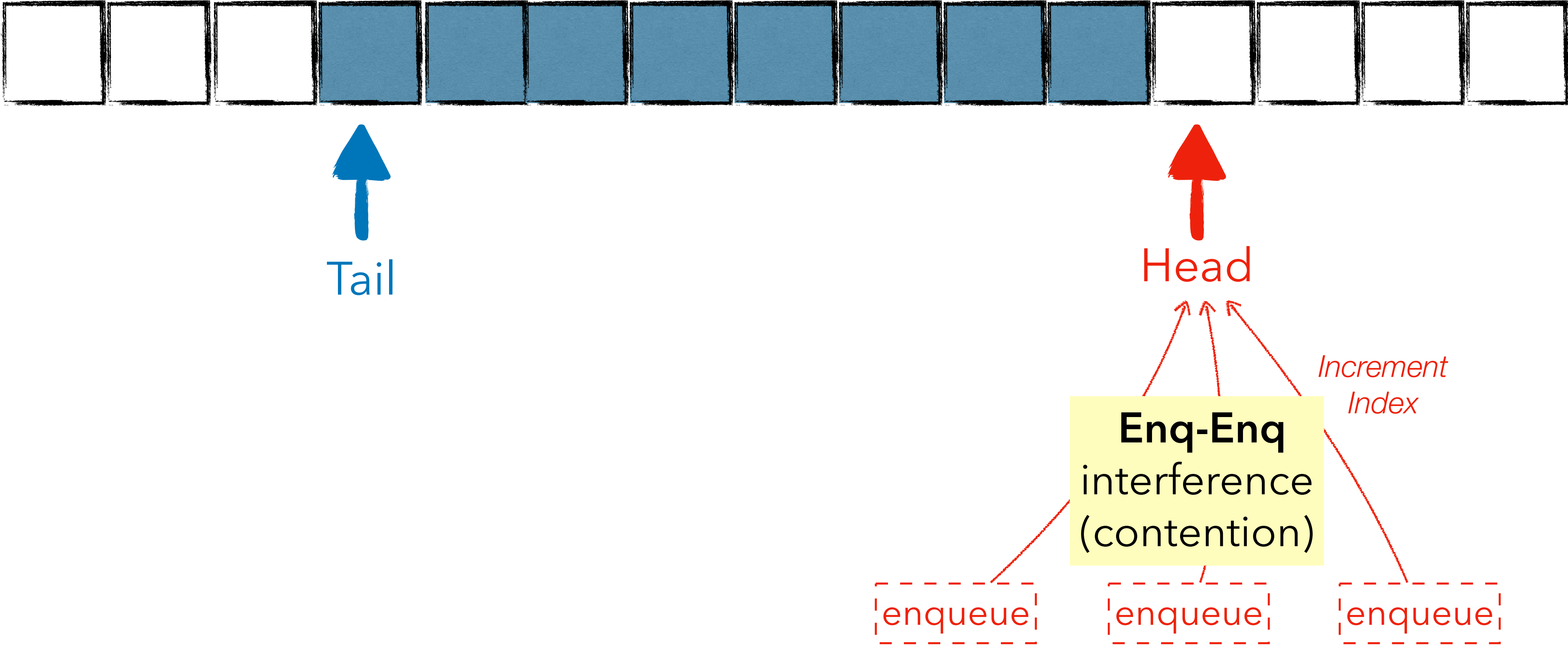Tail

Head

enqueue    enqueue    enqueue

# Interference sources

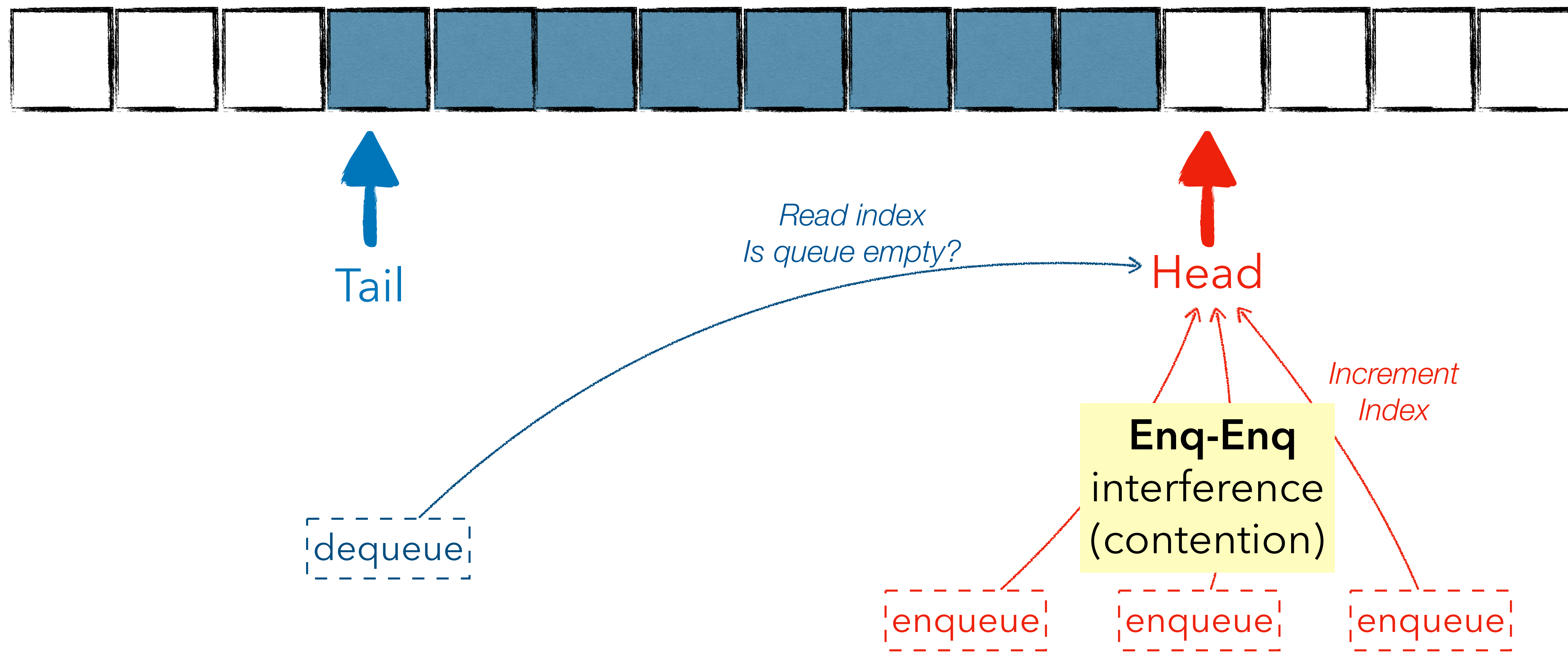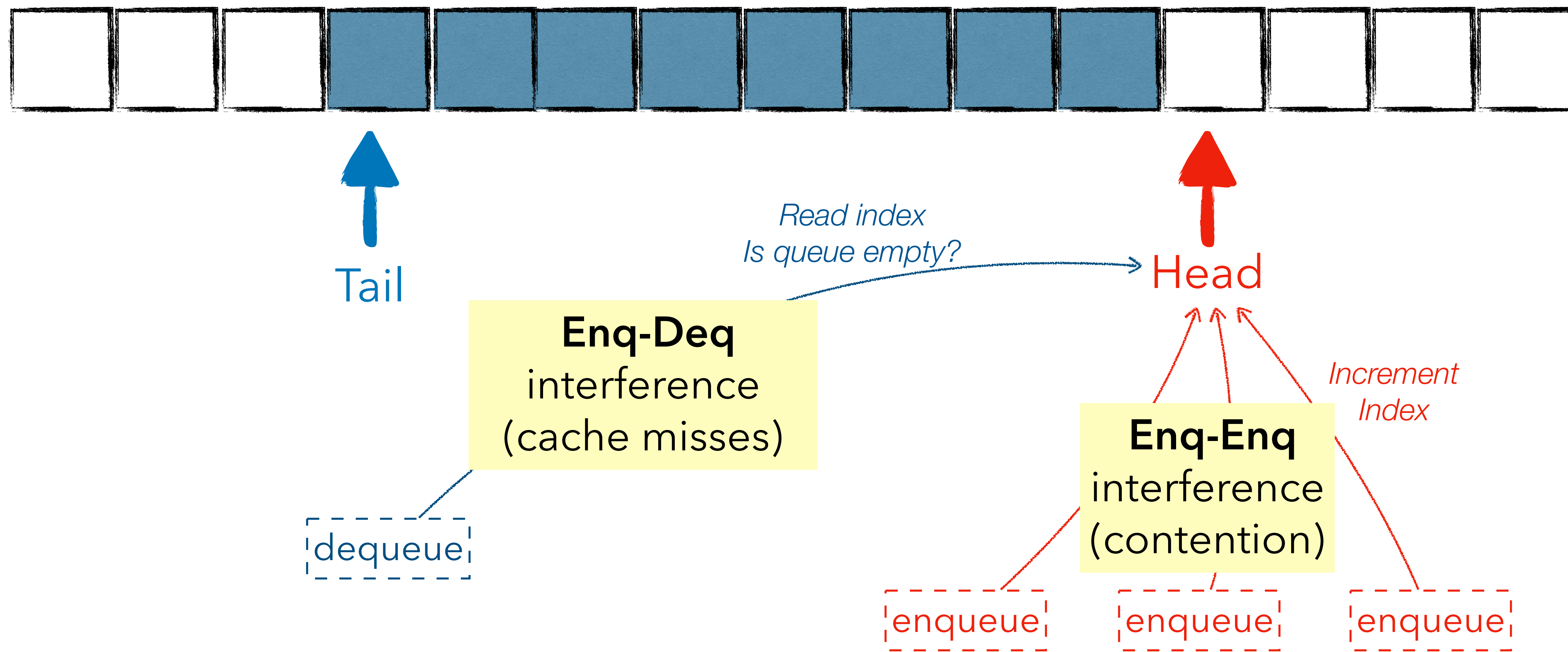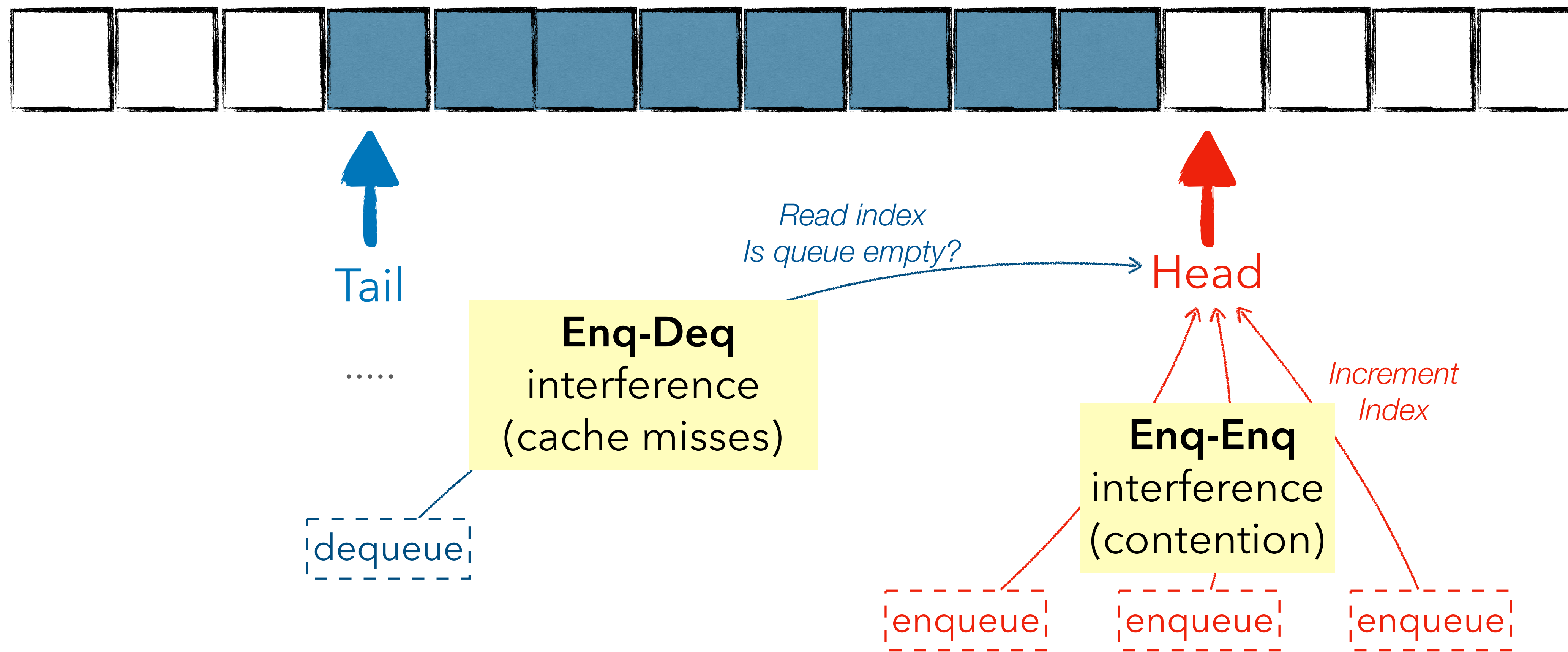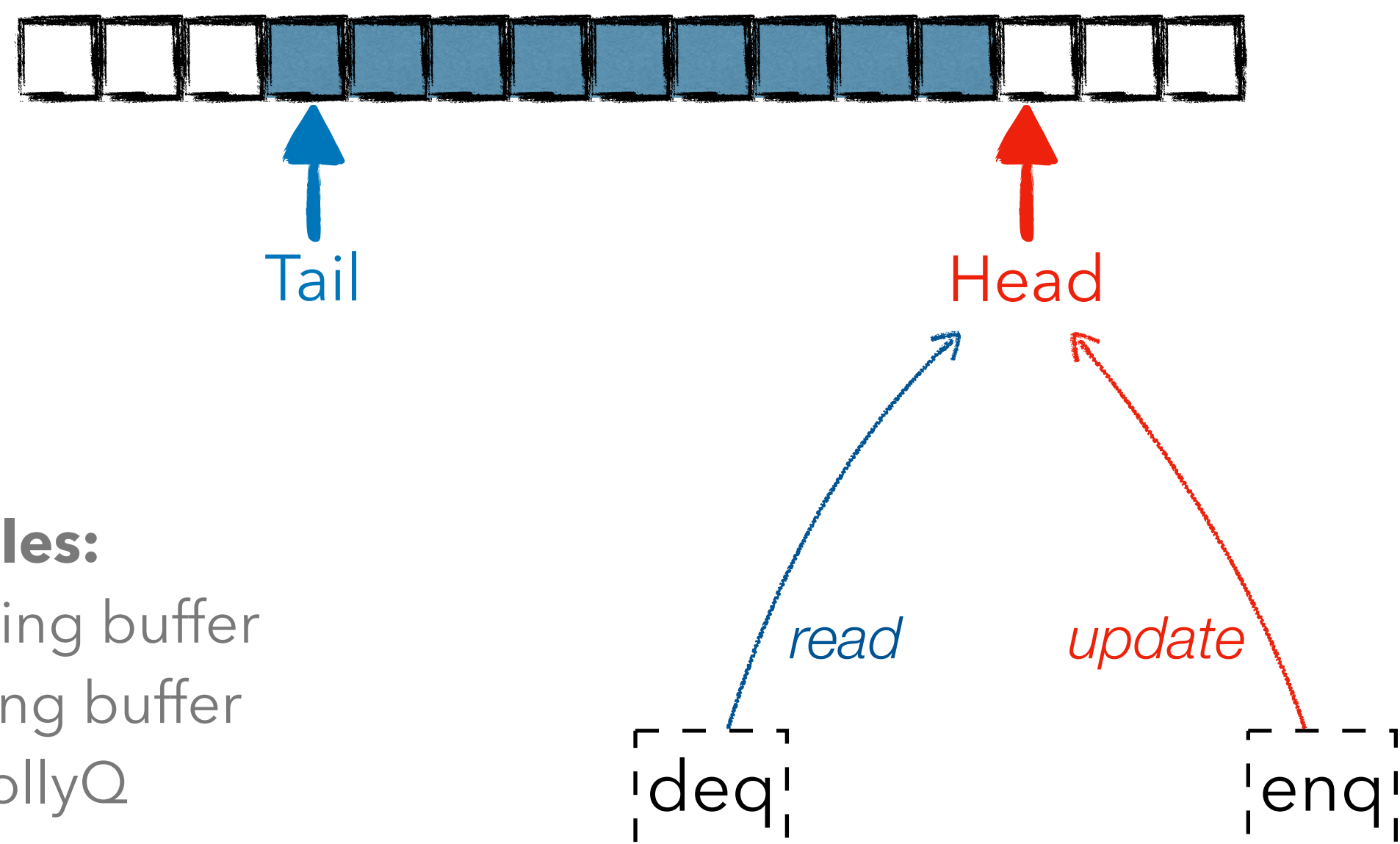Ring buffers are arrays with indices



Tail

Head

*Increment Index*

enqueue   enqueue   enqueue

# Interference sources



Ring buffers are arrays with indices

Tail

Head

*Increment Index*

**Enq-Enq** interference (contention)

enqueue    enqueue    enqueue

# Interference sources

Ring buffers are arrays with indices



Tail

*Read index*
*Is queue empty?*

Head

*Increment*
*Index*

dequeue

**Enq-Enq**
interference
(contention)

enqueue    enqueue    enqueue

# Interference sources

## Ring buffers are arrays with indices



Tail

*Read index*
*Is queue empty?*

Head

**Enq-Deq**
interference
(cache misses)

**Enq-Enq**
interference
(contention)

*Increment*
*Index*

dequeue

enqueue    enqueue    enqueue

# Interference sources

Ring buffers are arrays with indices

# Interference sources – Existing work

## Enq-Deq interference



Tail

Head

**Examples:**
DPDK ring buffer
Linux ring buffer
Meta FollyQ
SCQ

*read*

*update*

deq

enq

Mostly neglected

# Interference sources – Existing work
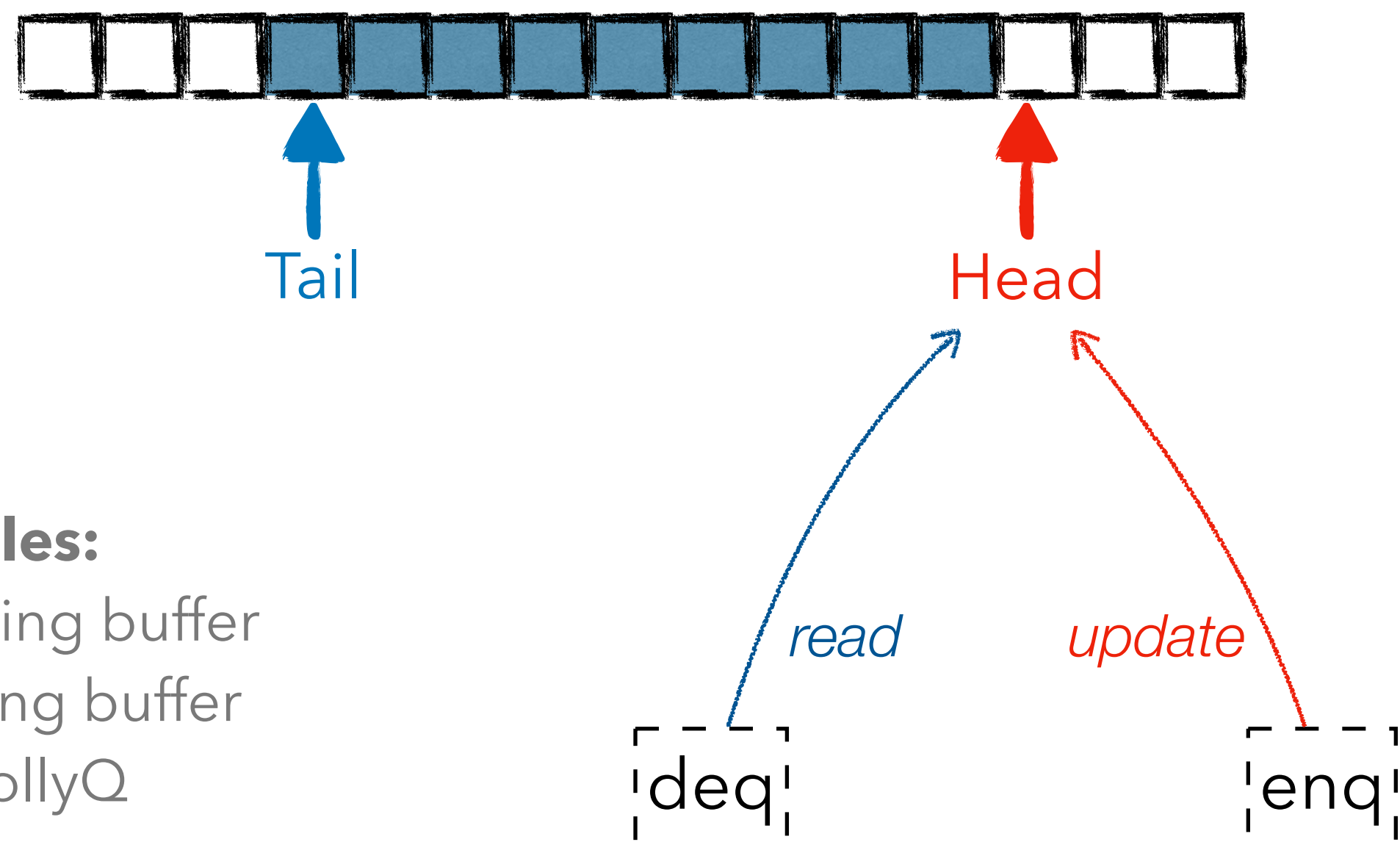
## Enq-Deq interference



Tail          Head

**Examples:**
DPDK ring buffer
Linux ring buffer
Meta FollyQ
SCQ

*read*     *update*

deq          enq

Mostly neglected

## Enq-Enq interference



Tail          Head

*update*                    …

**Example:**
FollyQ

enq    enq    enq    …

# Interference sources – Existing work

## Enq-Deq interference



**Examples:**
DPDK ring buffer
Linux ring buffer
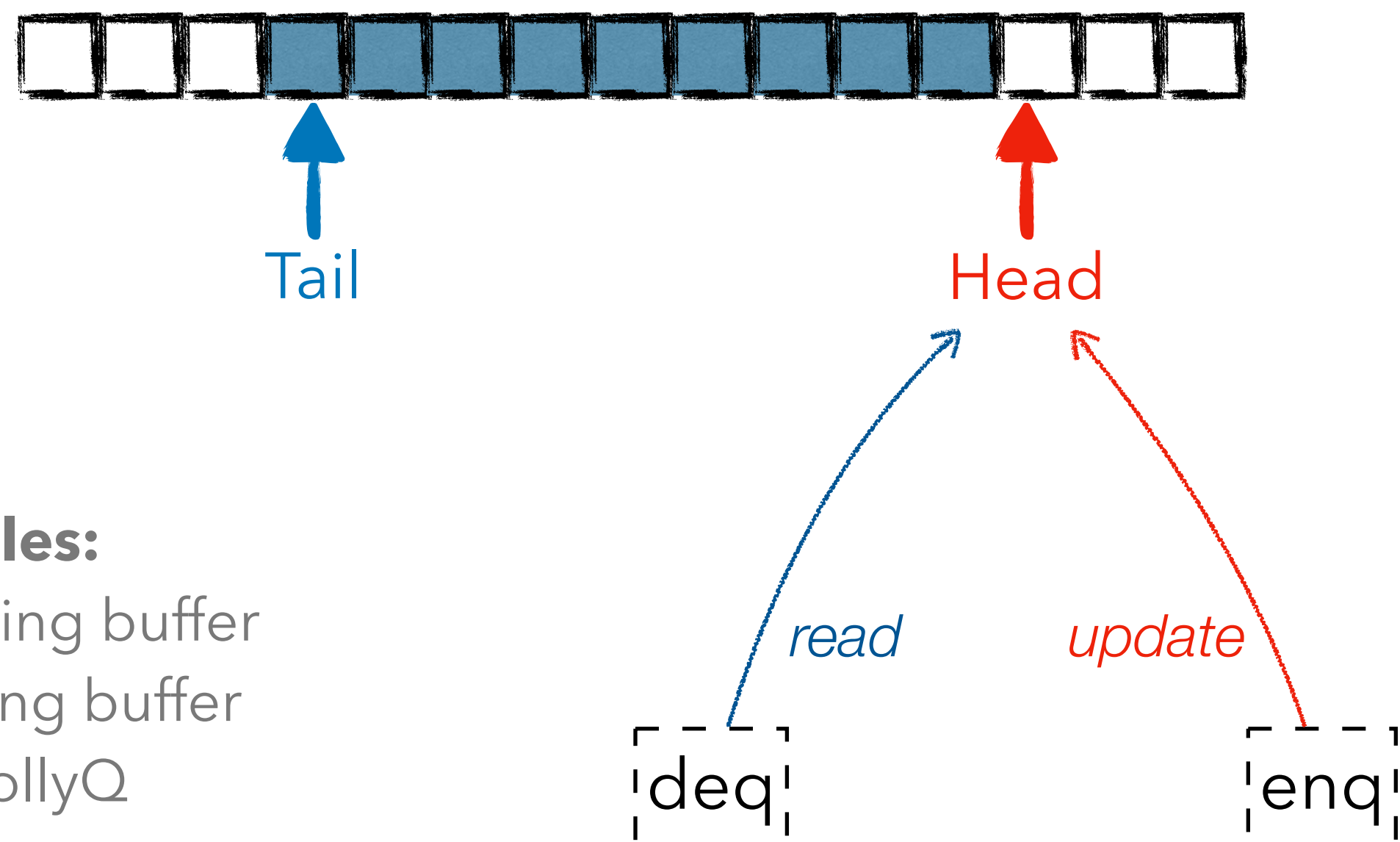Meta FollyQ
SCQ

Mostly neglected

## Enq-Enq interference



*FAA*

**Example:**
FollyQ

FAA typically faster than CAS

# Interference sources – Existing work
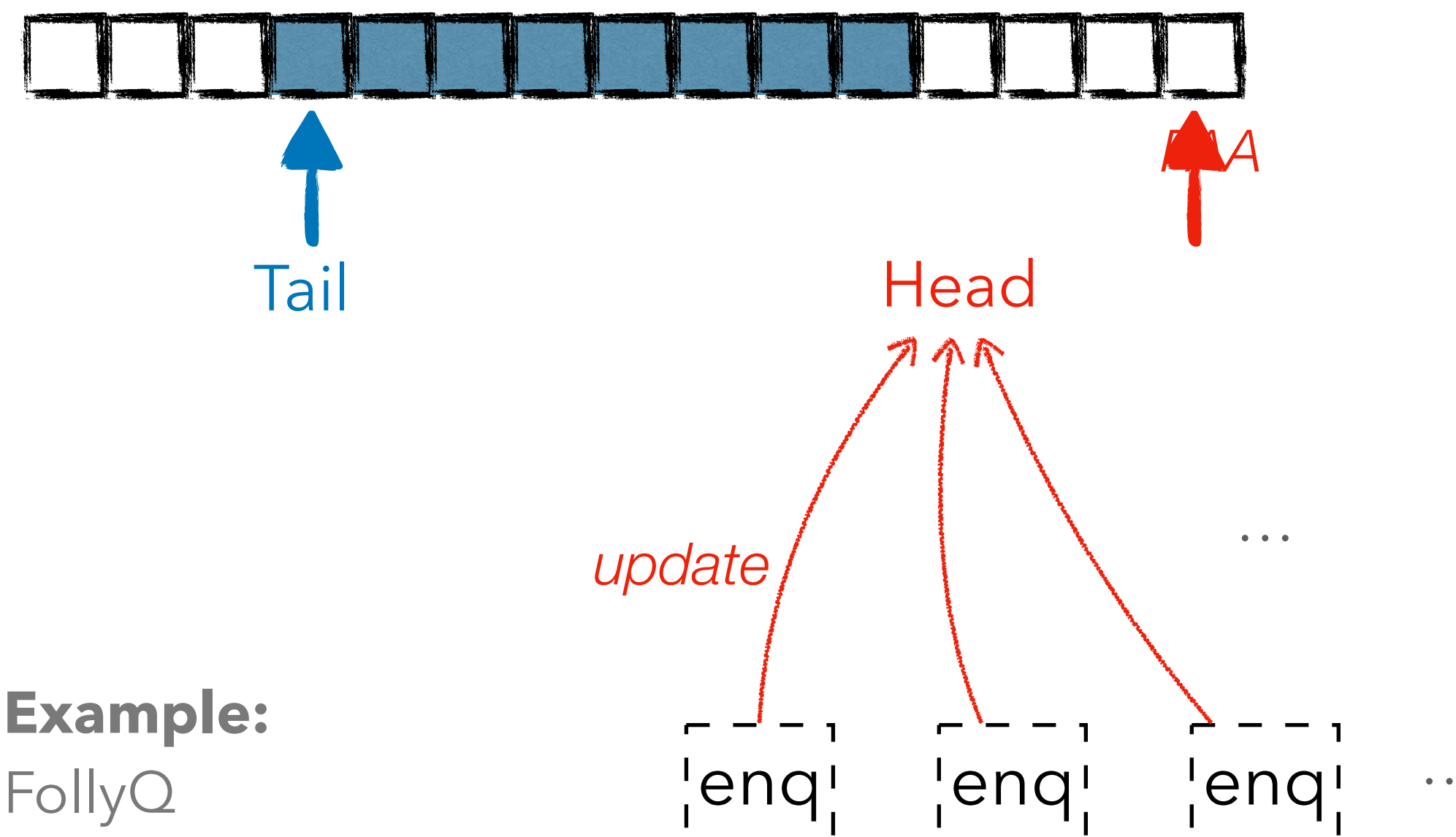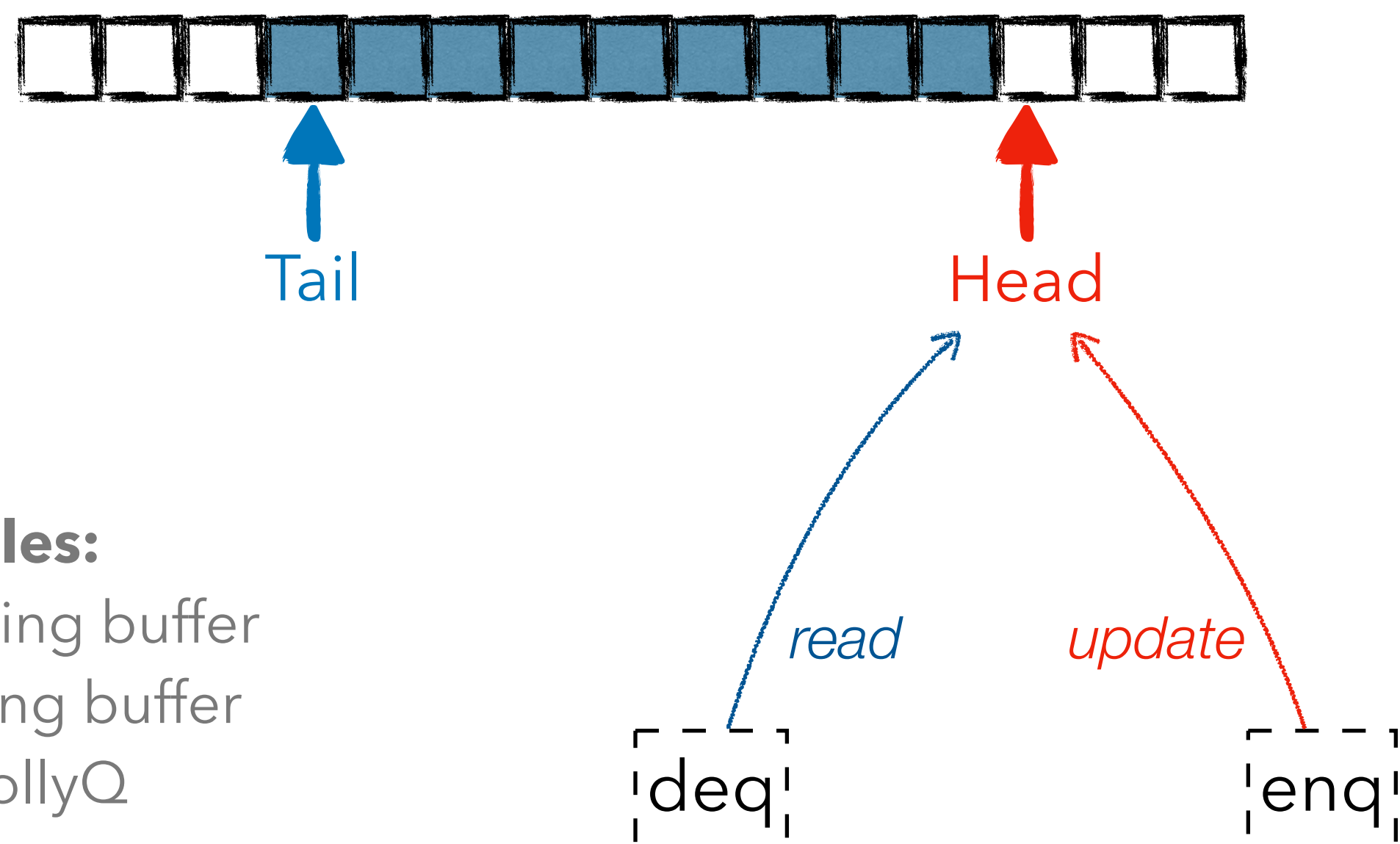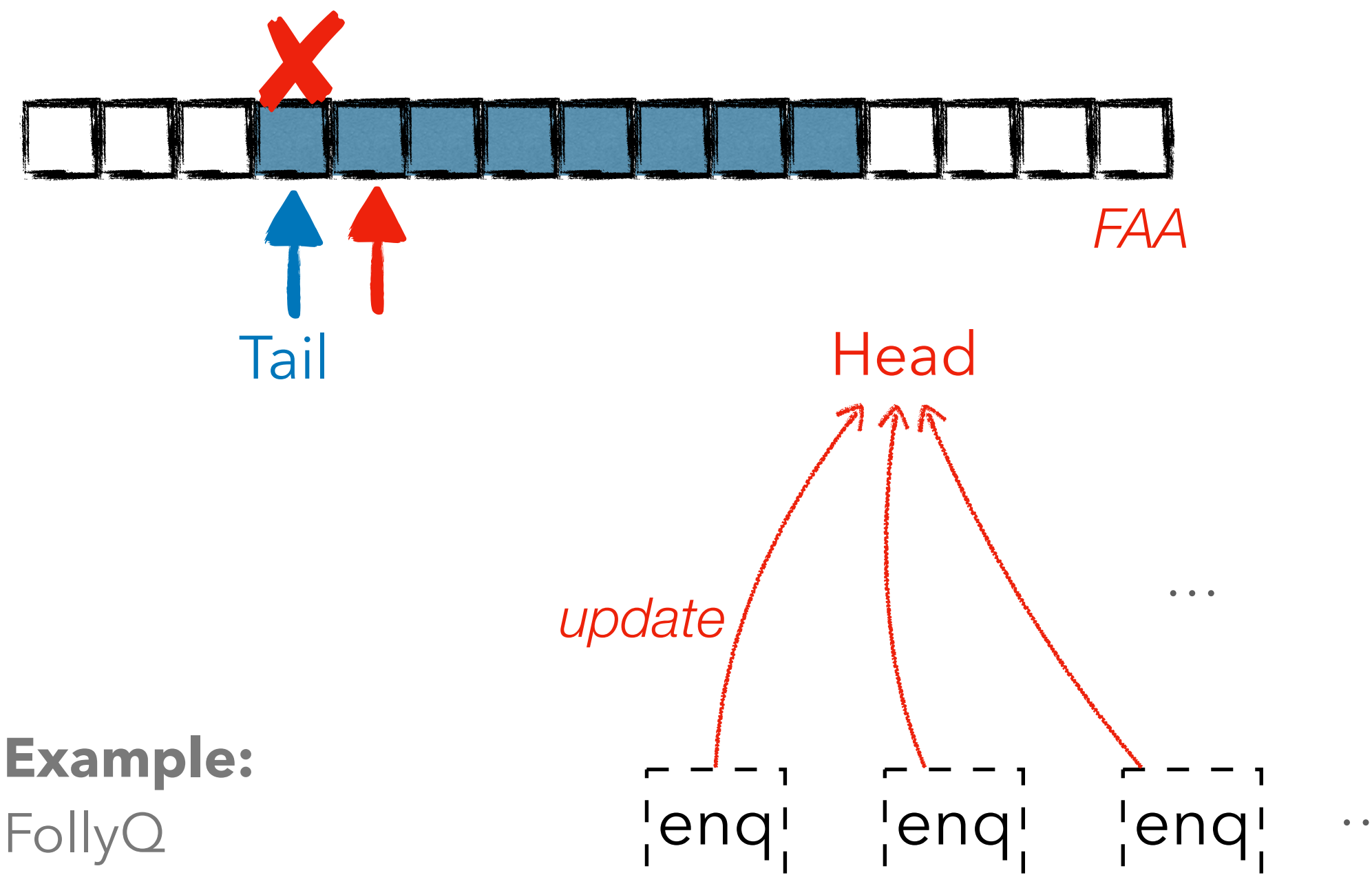
## Enq-Deq interference



Tail

Head

**Examples:**
DPDK ring buffer
Linux ring buffer
Meta FollyQ
SCQ

*read*    *update*

deq          enq

## Mostly neglected

## Enq-Enq interference



Tail

Head

FAA

*update*                    ...

**Example:**
FollyQ

enq    enq    enq    ...

## FAA typically faster than CAS

# Interference sources – Existing work



**Enq-Deq interference**

Tail     Head

**Examples:**
DPDK ring buffer
Linux ring buffer
Meta FollyQ
SCQ

*read*     *update*

deq     enq

Mostly neglected

**Enq-Enq interference**

Tail     Head     *FAA*

**Example:**
FollyQ

*update*     ...

enq     enq     enq     ...
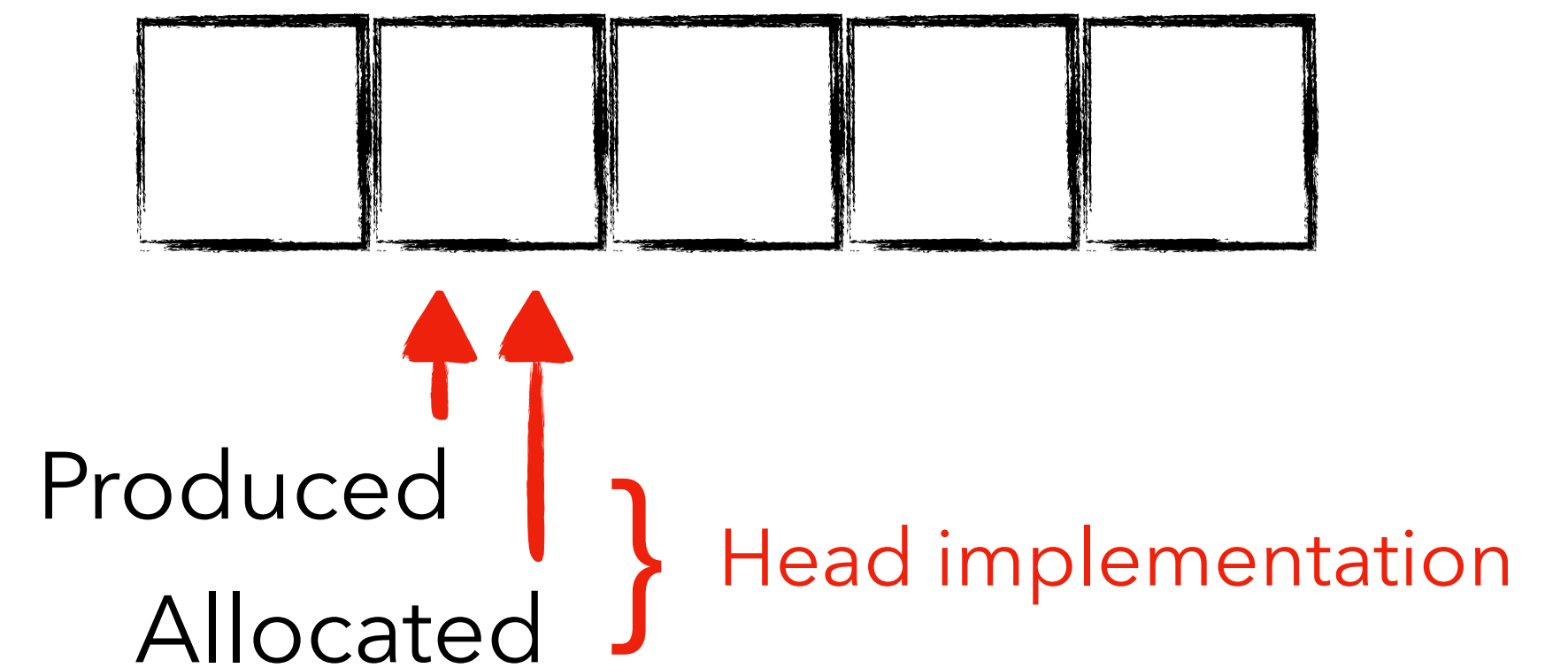
FAA typically faster than CAS

# Story 2: Oversubscription and out-of-order operations
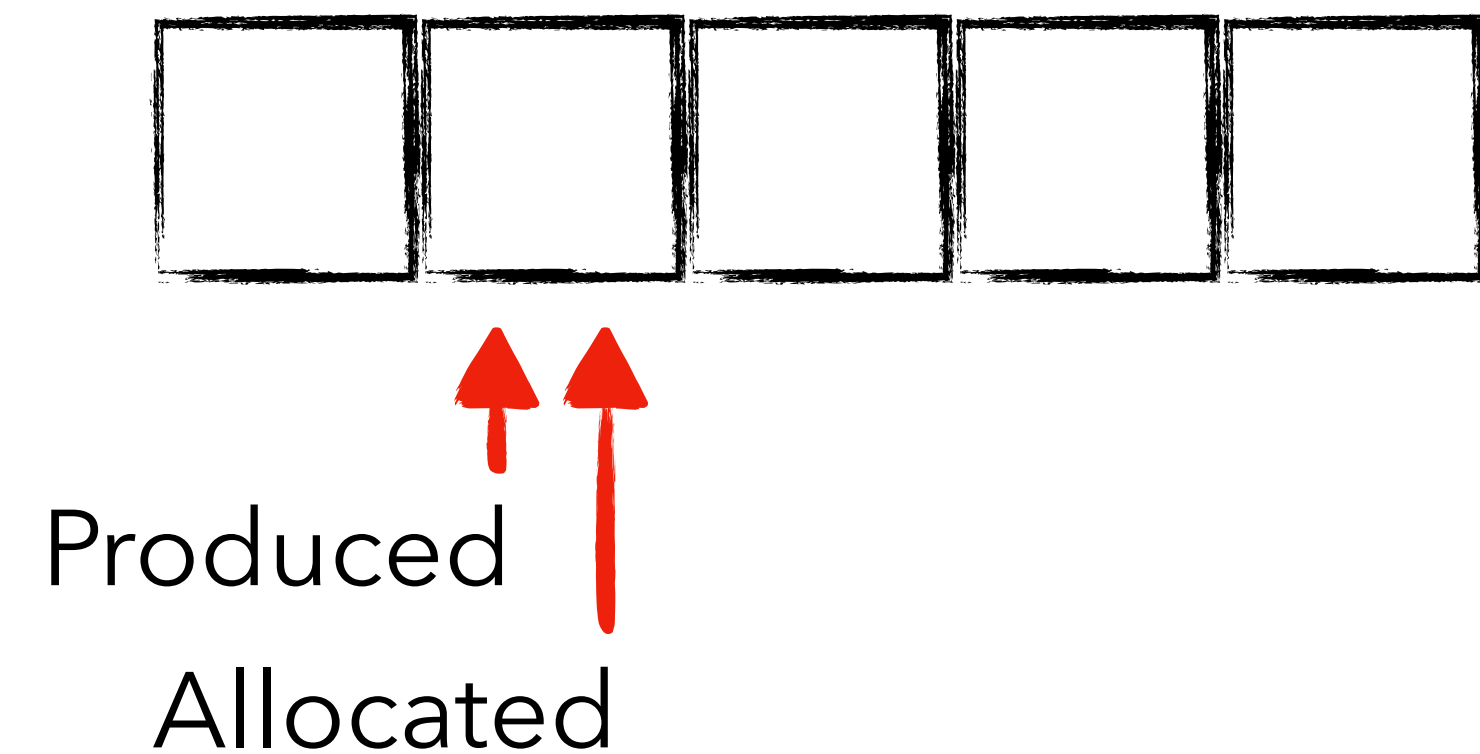
- New framework for mobile devices

  - **Few cores and many threads**

  - Communication via ring buffers

- Problem with initial implementation:

  - In-order operation **limits** performance



Produced
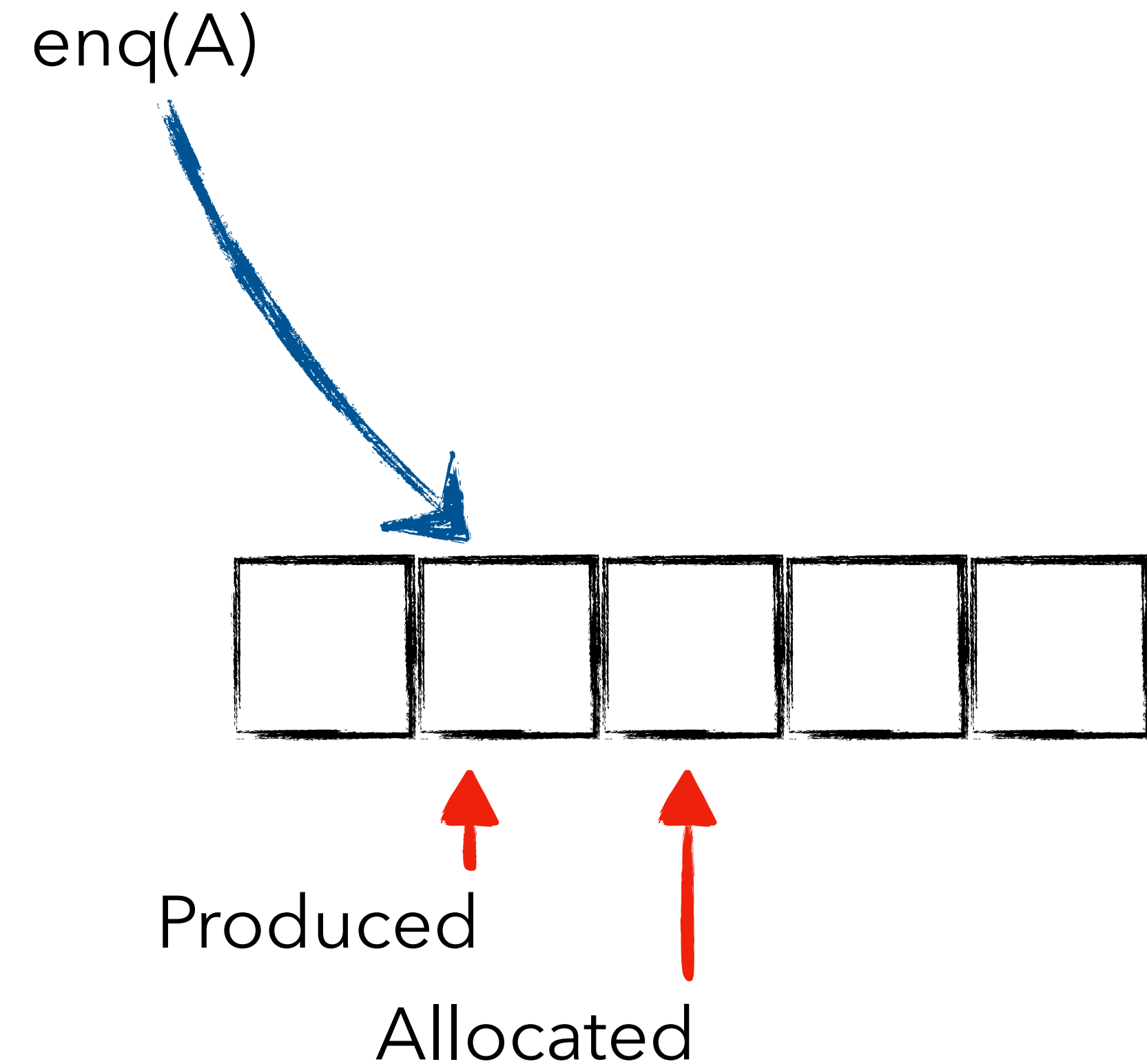
Allocated

} Head implementation

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices

  - **Few cores and many threads**

  - Communication via ring buffers

- Problem with initial implementation:

  - In-order operation **limits** performance
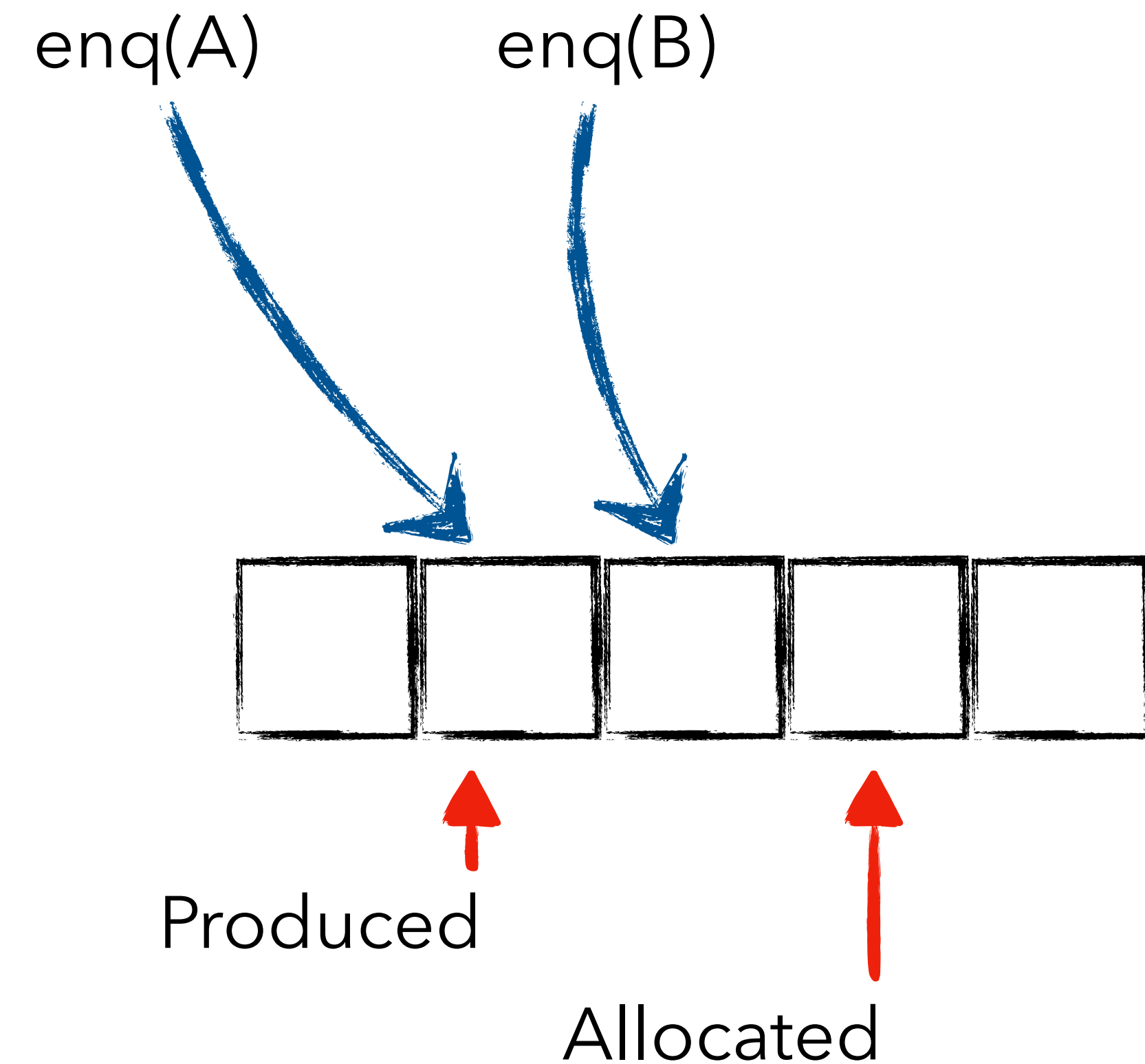
Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices

  - **Few cores and many threads**

  - Communication via ring buffers

- Problem with initial implementation:

  - In-order operation **limits** performance
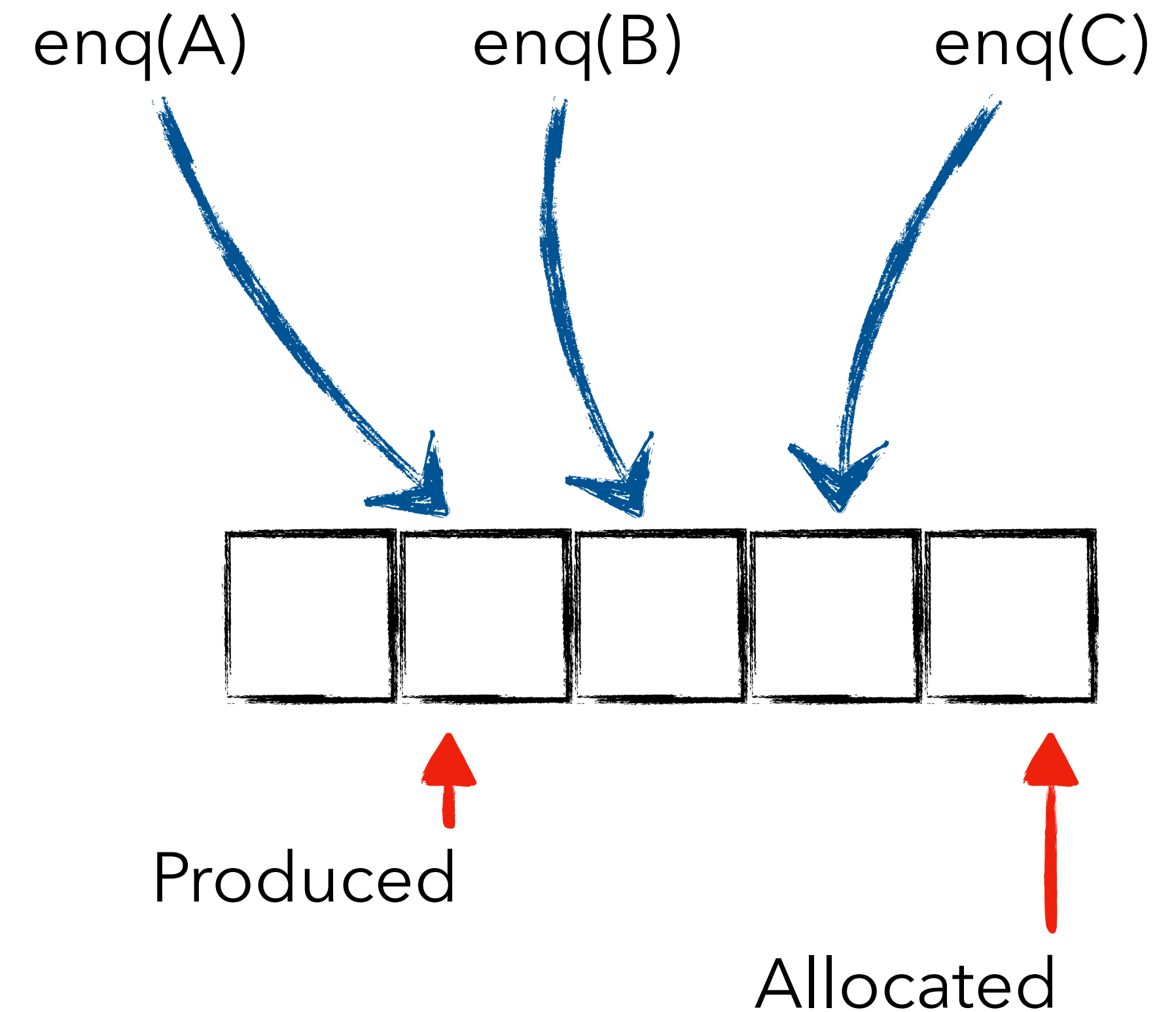
enq(A)

Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

enq(A)     enq(B)
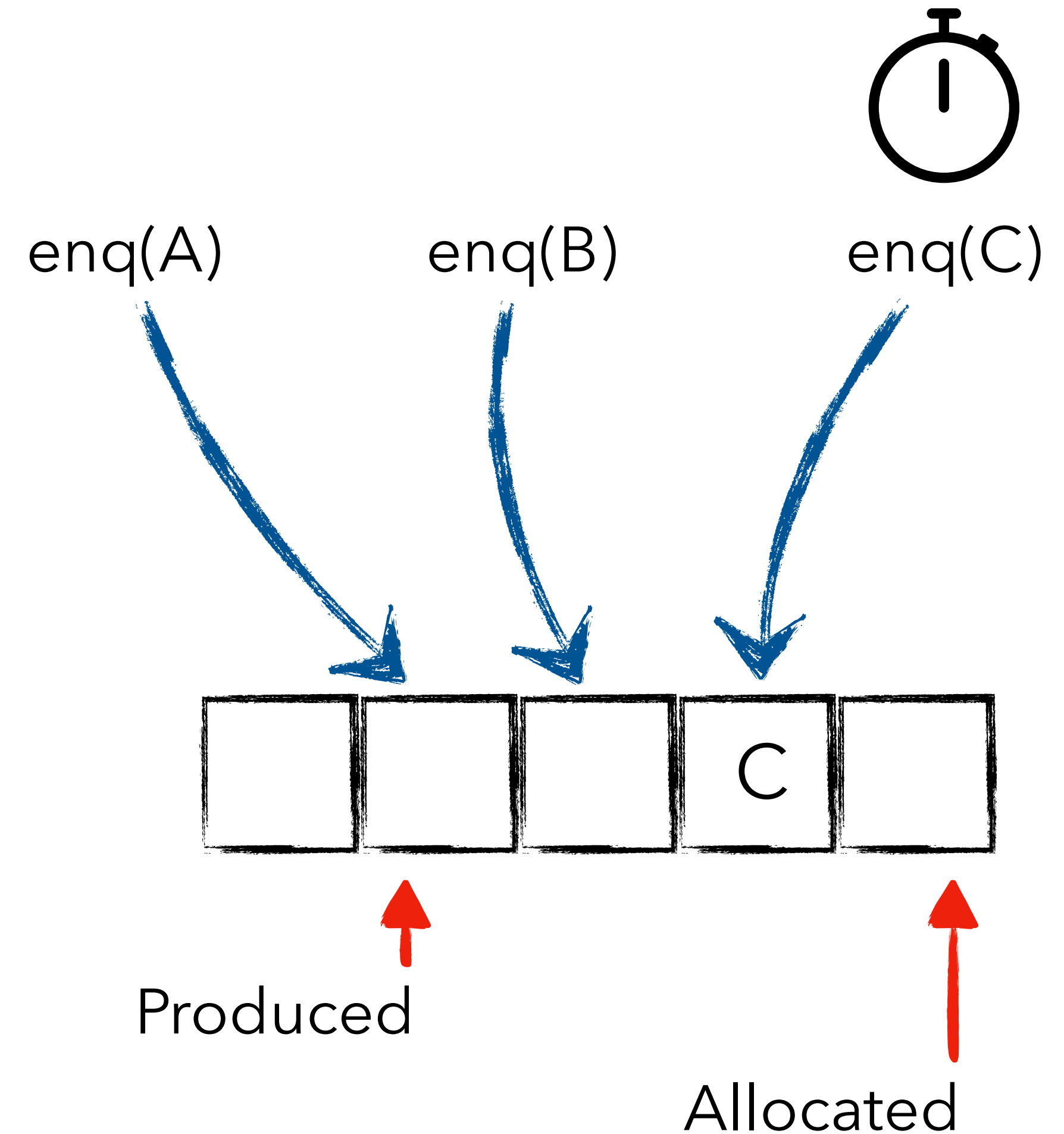
Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

enq(A)    enq(B)    enq(C)

Produced

Allocated

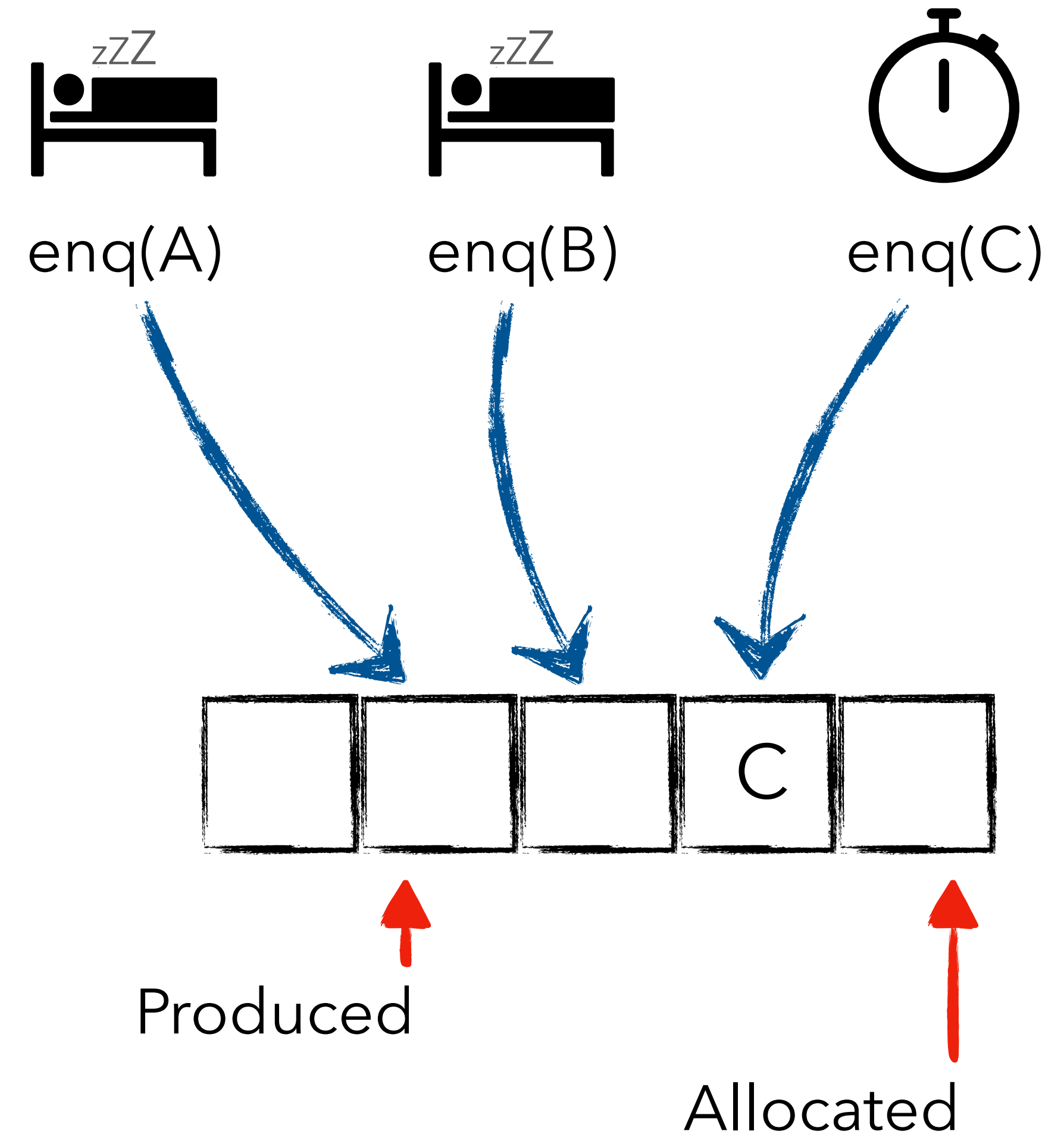# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

enq(A)        enq(B)        enq(C)

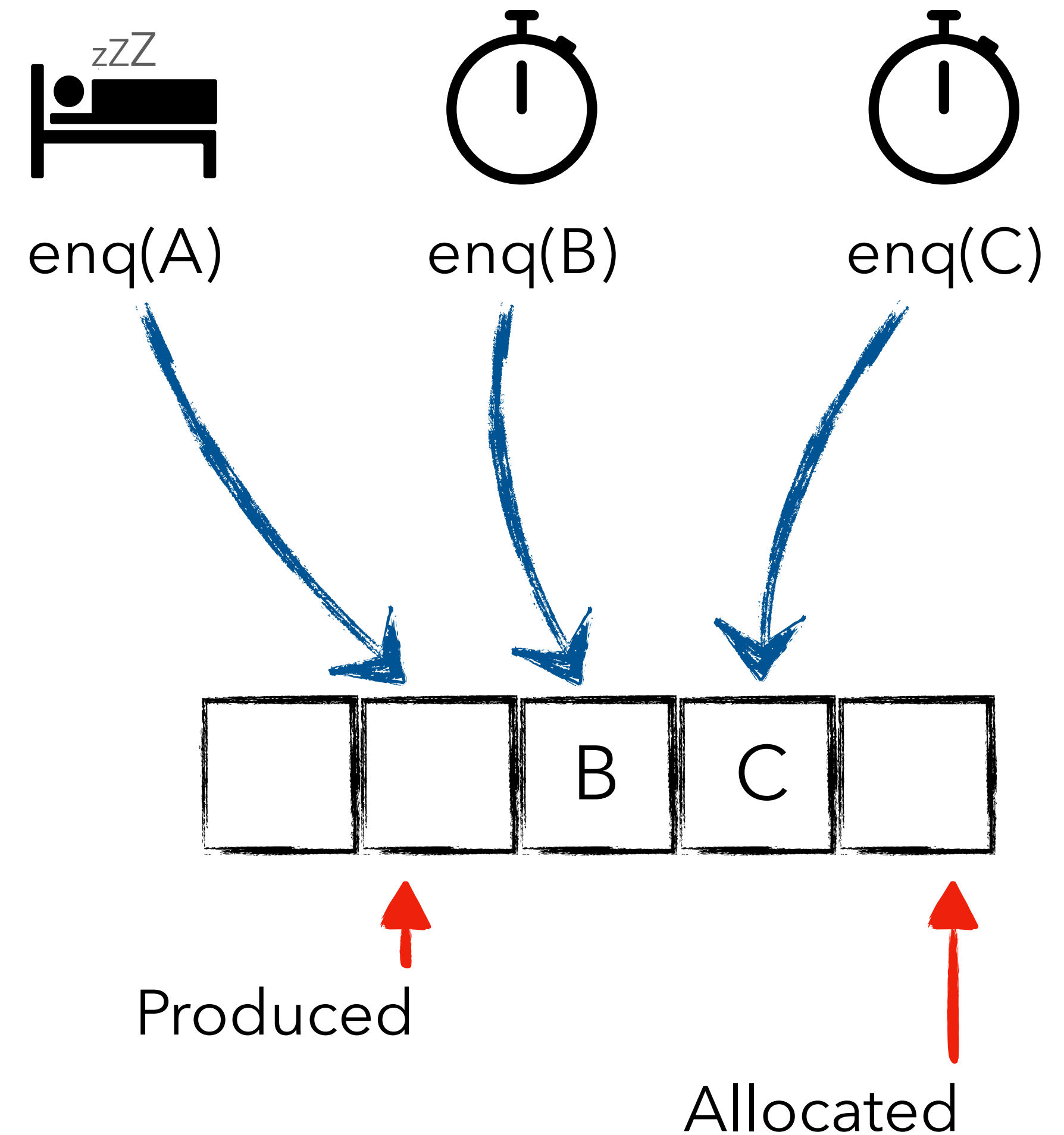| | | | C | |
|---|---|---|---|---|

Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

enq(A)        enq(B)        enq(C)

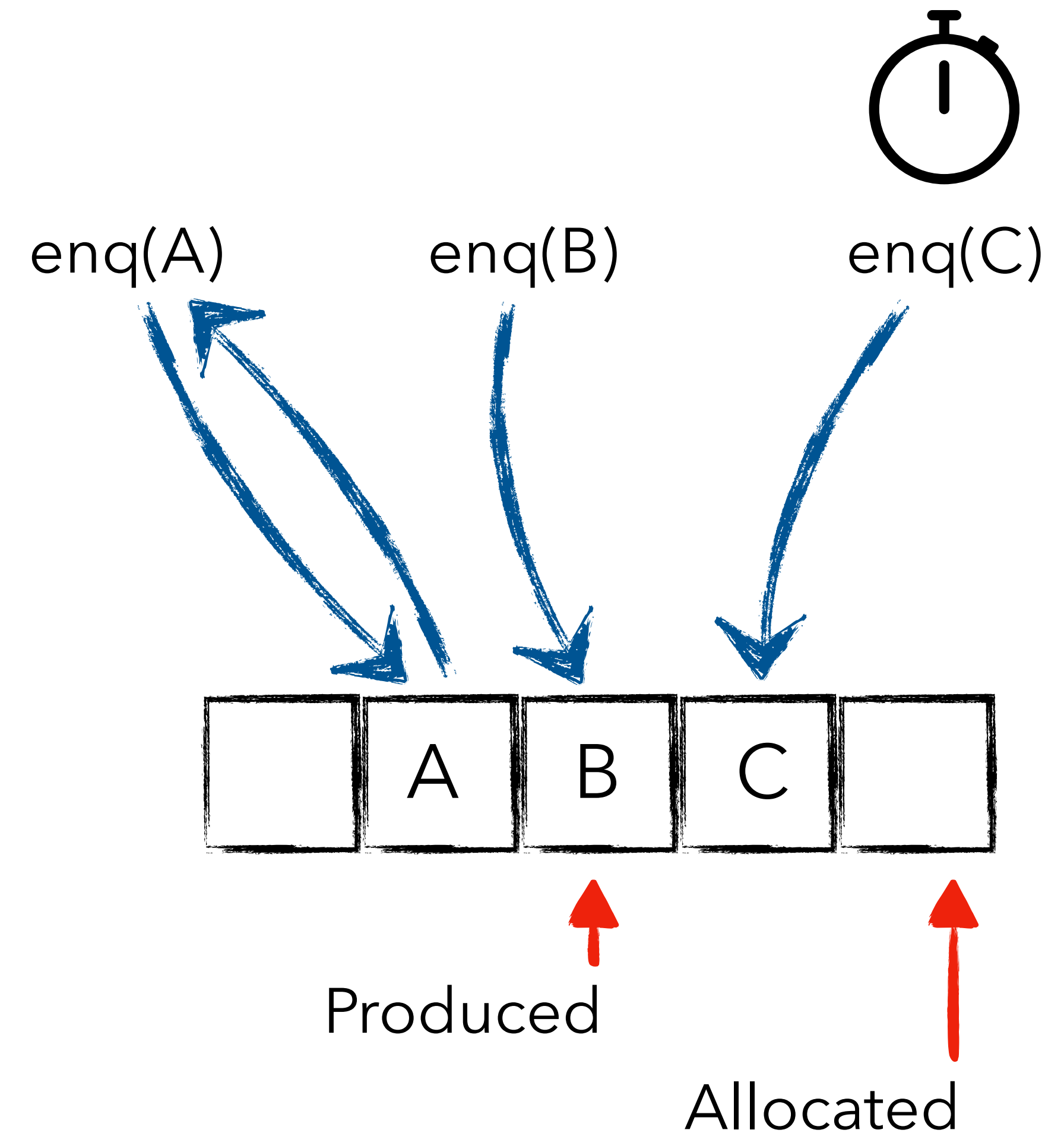| | | | C | |
|---|---|---|---|---|

Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

enq(A)    enq(B)    enq(C)

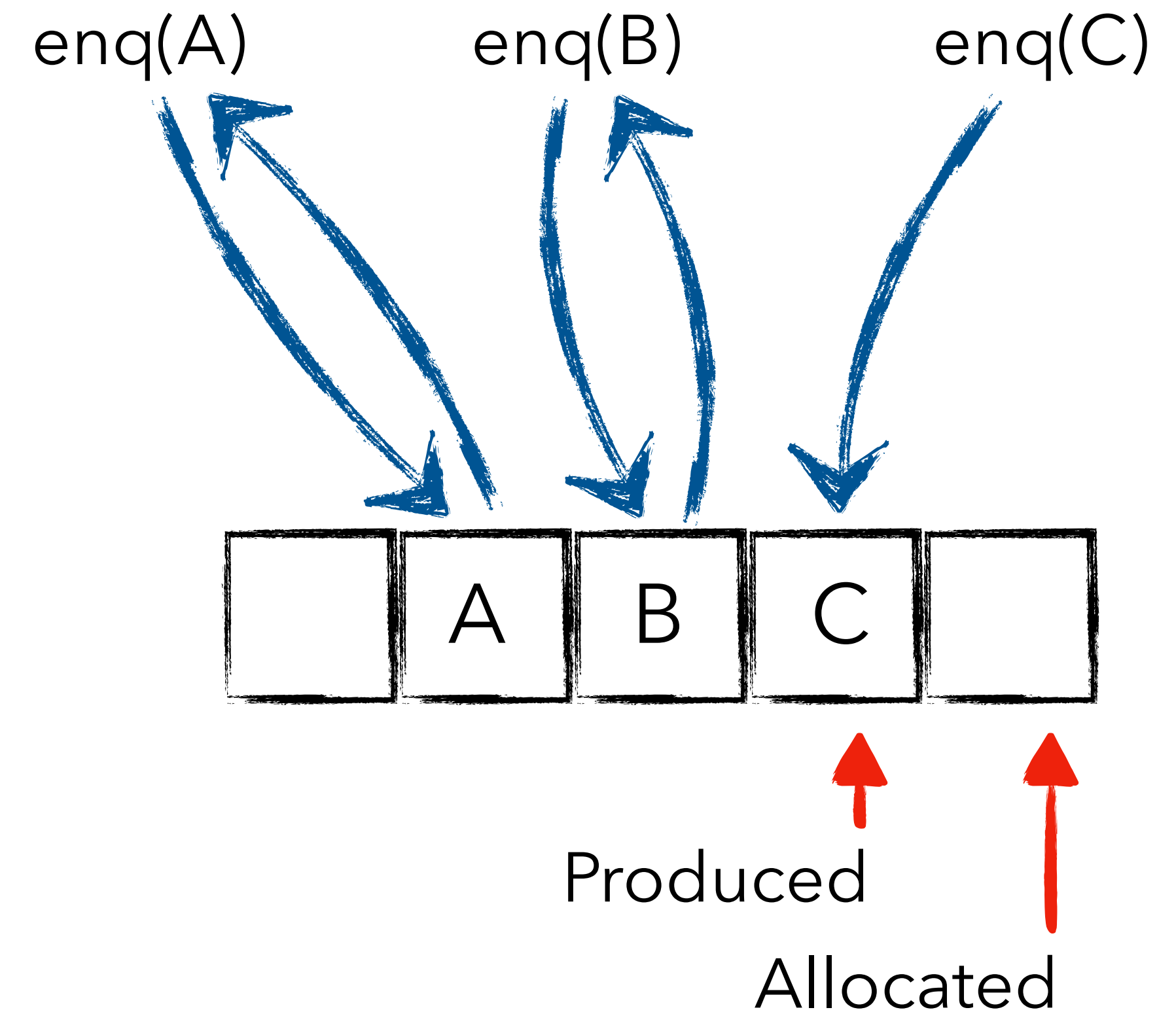| | | B | C | |

↑ Produced

↑ Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices

  - **Few cores and many threads**

  - Communication via ring buffers

- Problem with initial implementation:

  - In-order operation **limits** performance

enq(A)    enq(B)    enq(C)

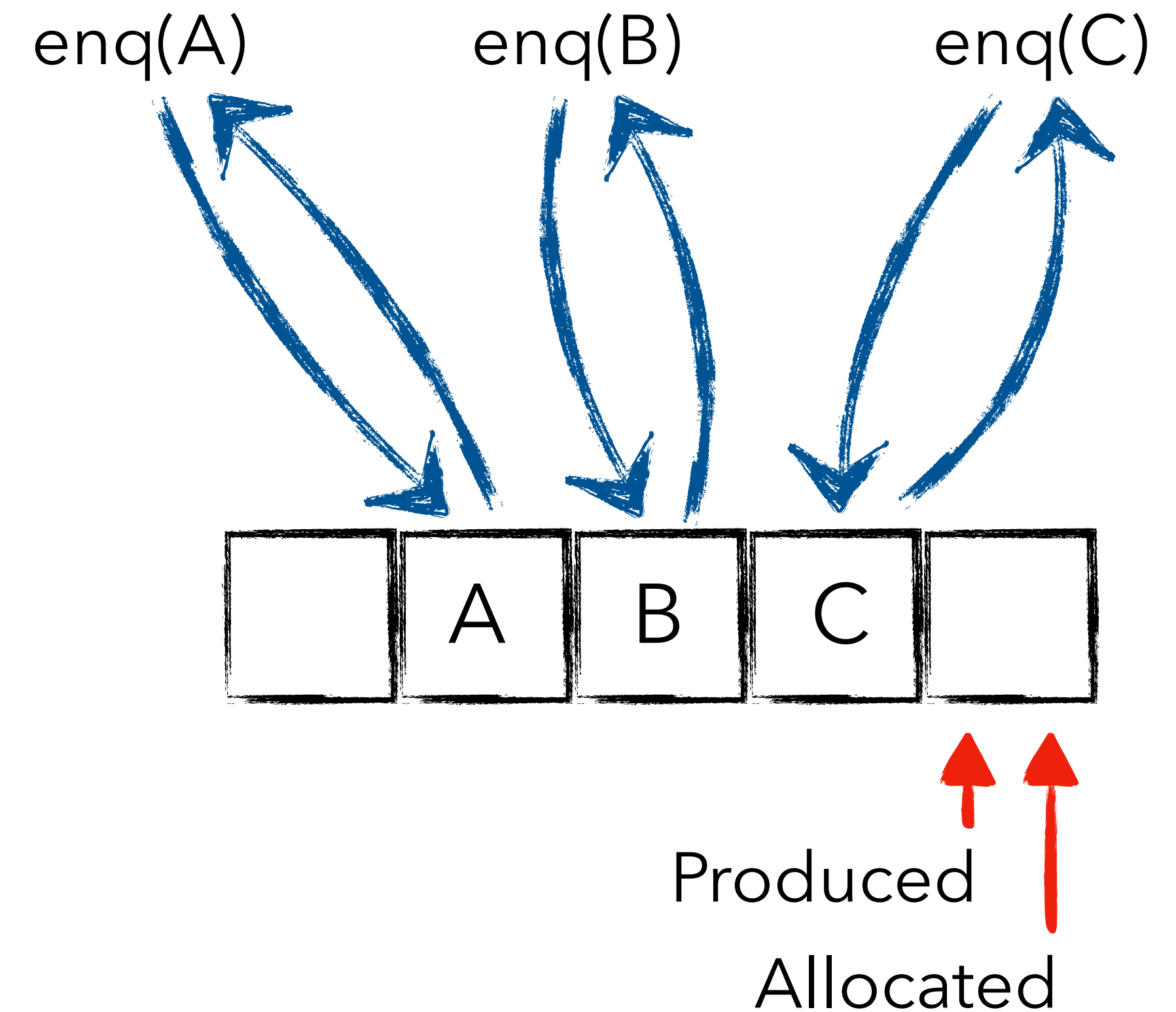|   | A | B | C |   |
|---|---|---|---|---|

Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**

  - Communication via ring buffers


- Problem with initial implementation:
  - In-order operation **limits** performance



enq(A)  enq(B)  enq(C)

| | A | B | C | |

Produced

Allocated

# Story 2: Oversubscription and out-of-order operations
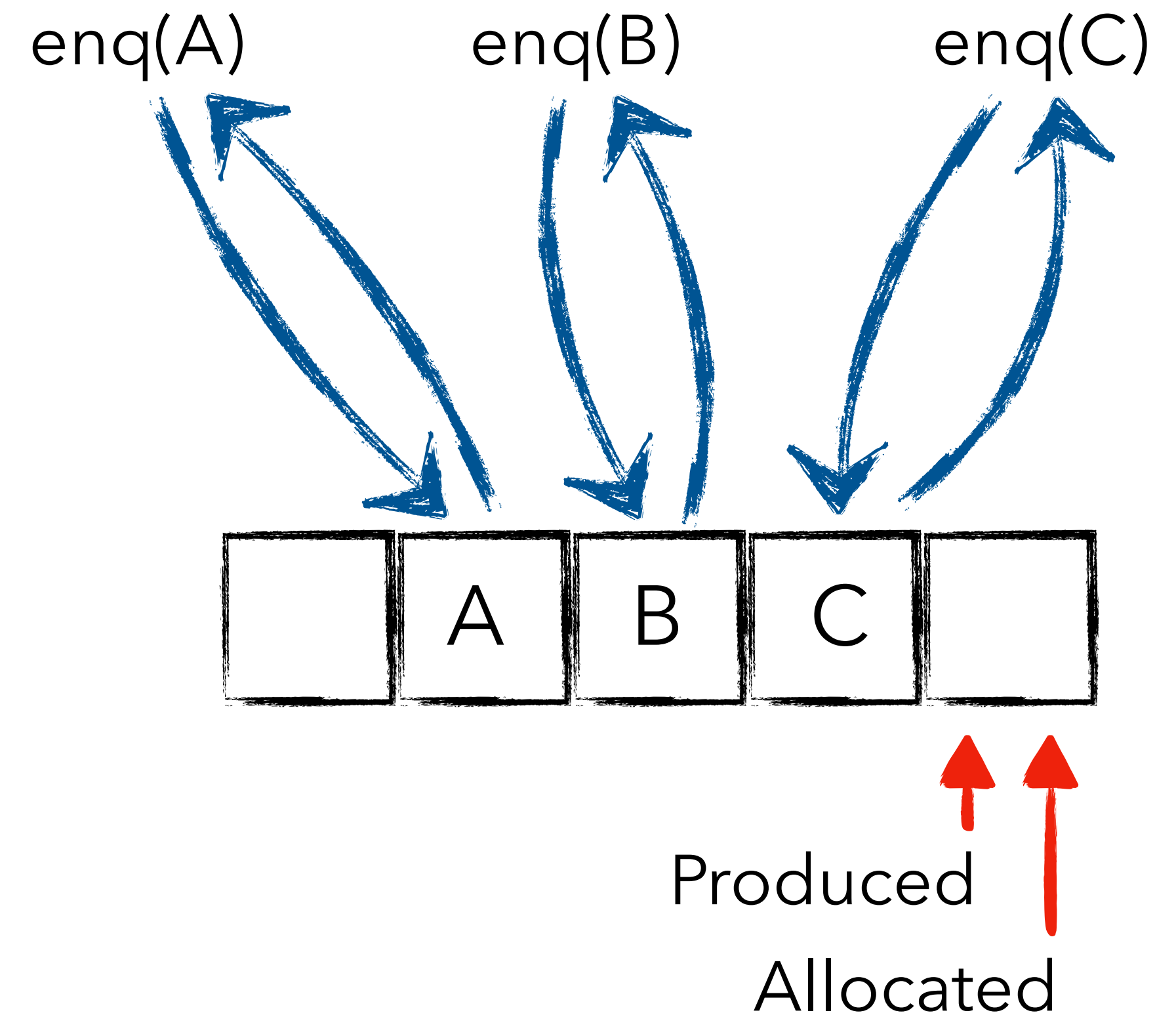
- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

enq(A)          enq(B)          enq(C)

| | A | B | C | |

Produced

Allocated

# Story 2: Oversubscription and out-of-order operations

- New framework for mobile devices
  - **Few cores and many threads**
  - Communication via ring buffers

- Problem with initial implementation:
  - In-order operation **limits** performance

- Out-of-order operations are **challenging**!
  - See paper for related work

enq(A)    enq(B)    enq(C)

| | A | B | C | |

Produced

Allocated

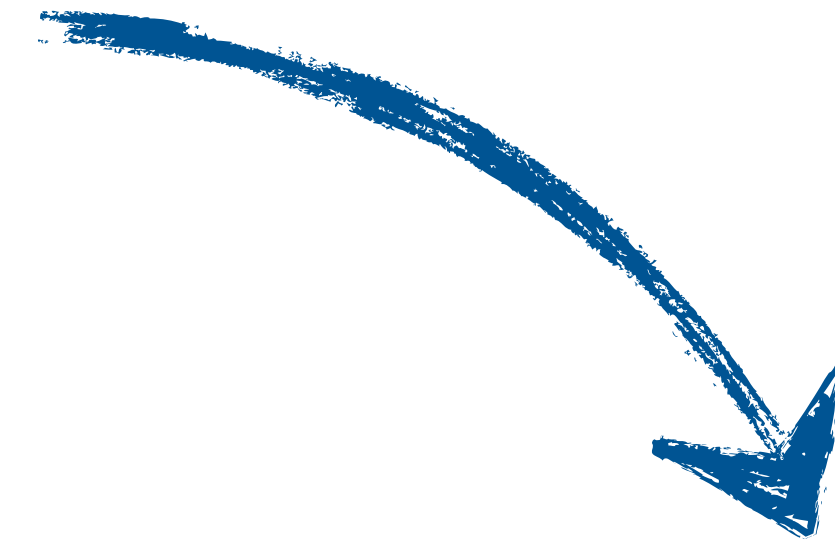# Story 3: Migrating from x86 to Arm

Product stable on x86 for years

- Must migrate to Arm (TaiShan servers)

- Internally uses **old DPDK ring buffer**

- Application grew **intertwined** with ring buffer

# Story 3: Migrating from x86 to Arm

Product stable on x86 for years

- Must migrate to Arm (TaiShan servers)
- Internally uses **old DPDK ring buffer**
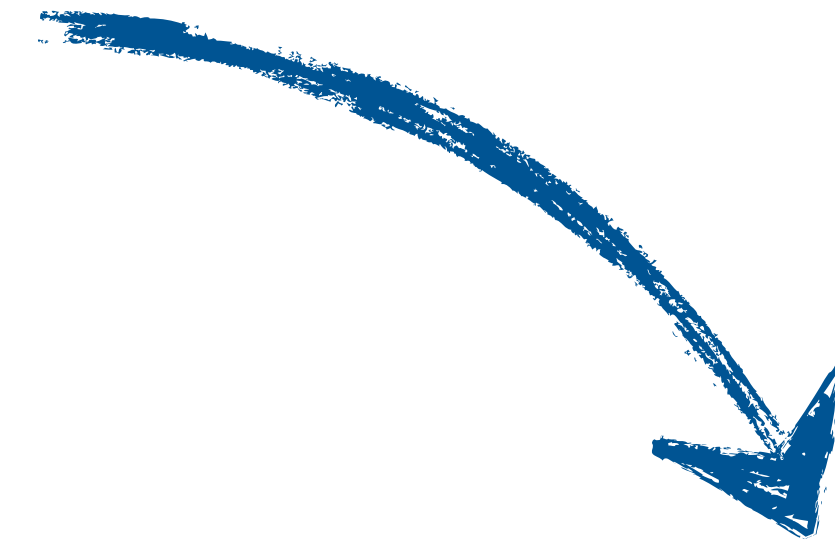- Application grew **intertwined** with ring buffer

Decision: **too high risk** of upgrading DPDK ring buffer

# Story 3: Migrating from x86 to Arm

Product stable on x86 for years

- Must migrate to Arm (TaiShan servers)
- Internally uses **old DPDK ring buffer**
- Application grew **intertwined** with ring buffer

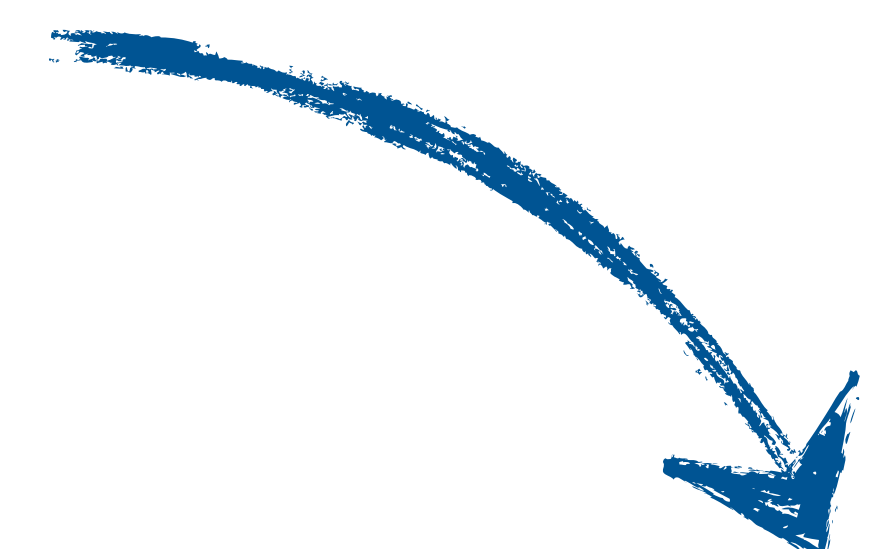Decision: **too high risk** of upgrading DPDK ring buffer

Consequence:

- **Annoying weak memory bug**
   due to a few missing fences
- More than **6 person-month to fix it**
- Decision wasn't the best

Product stable on x86 for years

- Must migrate to Arm (TaiShan servers)
- Internally uses **old DPDK ring buffer**
- Application grew **intertwined** with ring buffer

Decision: **too high risk** of upgrading DPDK ring buffer

Consequence:

- **Annoying weak memory bug** due to a few missing fences
- More than **6 person-month to fix it**
- Decision wasn't the best

```
                    Init
              data = ctrl = 0;

   Thread 1                    Thread 2
   data = 1;                   while(!ctrl) {}
   ctrl = 1;                   assert(data == 1); ✗
```
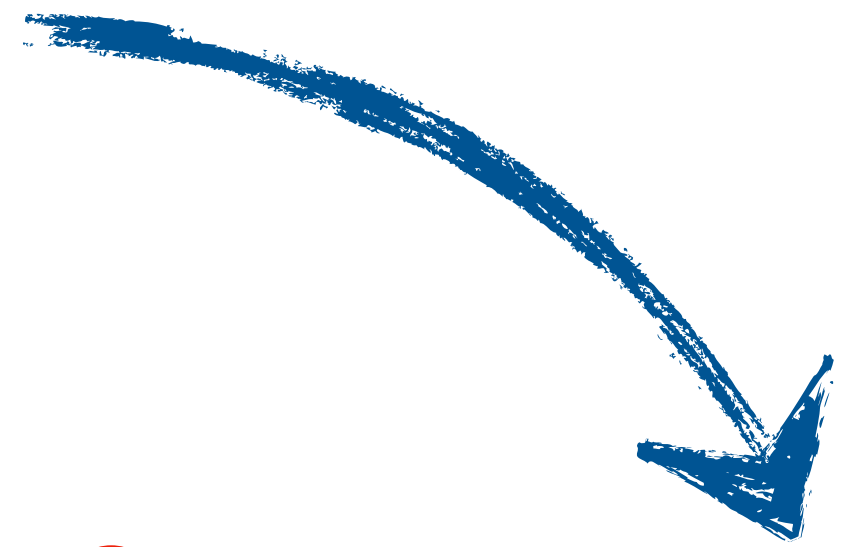
# Story 3: Migrating from x86 to Arm

Product stable on x86 for years

- Must migrate to Arm (TaiShan servers)
- Internally uses **old DPDK ring buffer**
- Application grew **intertwined** with ring buffer

Decision: **too high risk** of upgrading DPDK ring buffer

Consequence:

- **Annoying weak memory bug** due to a few missing fences
- More than **6 person-month to fix it**
- Decision wasn't the best

```
                    Init
        data = ctrl = 0;

    Thread 1                Thread 2
Fence→ data = 1;        while(!ctrl) {}
       ctrl = 1;        assert(data == 1); ✗
```

# How do people develop for WMM?

## Think hard and document

For example, printk_ringbuffer

```
/*
 * Guarantee the state is loaded before copying the descriptor
 * content. This avoids copying obsolete descriptor content that might
 * not apply to the descriptor state. This pairs with _prb_commit:B.
 *
 * Memory barrier involvement:
 *
 * If desc_read:A reads from _prb_commit:B, then desc_read:C reads
 * from _prb_commit:A.
 *
 * Relies on:
 *
 * WMB from _prb_commit:A to _prb_commit:B
 *    matching
 * RMB from desc_read:A to desc_read:C
 */
smp_rmb(); /* LMM(desc_read:B) */
```

# How do people develop for WMM?

## Think hard and document

For example, printk_ringbuffer

```
/*
 * Guarantee the state is loaded before copying the descriptor
 * content. This avoids copying obsolete descriptor content that might
 * not apply to the descriptor state. This pairs with _prb_commit:B.
 *
 * Memory barrier involvement:
 *
 * If desc_read:A reads from _prb_commit:B, then desc_read:C reads
 * from _prb_commit:A.
 *
 * Relies on:
 *
 * WMB from _prb_commit:A to _prb_commit:B
 *    matching
 * RMB from desc_read:A to desc_read:C
 */
smp_rmb(); /* LMM(desc_read:B) */
```

## Most just ignore topic

And wait to see what happens

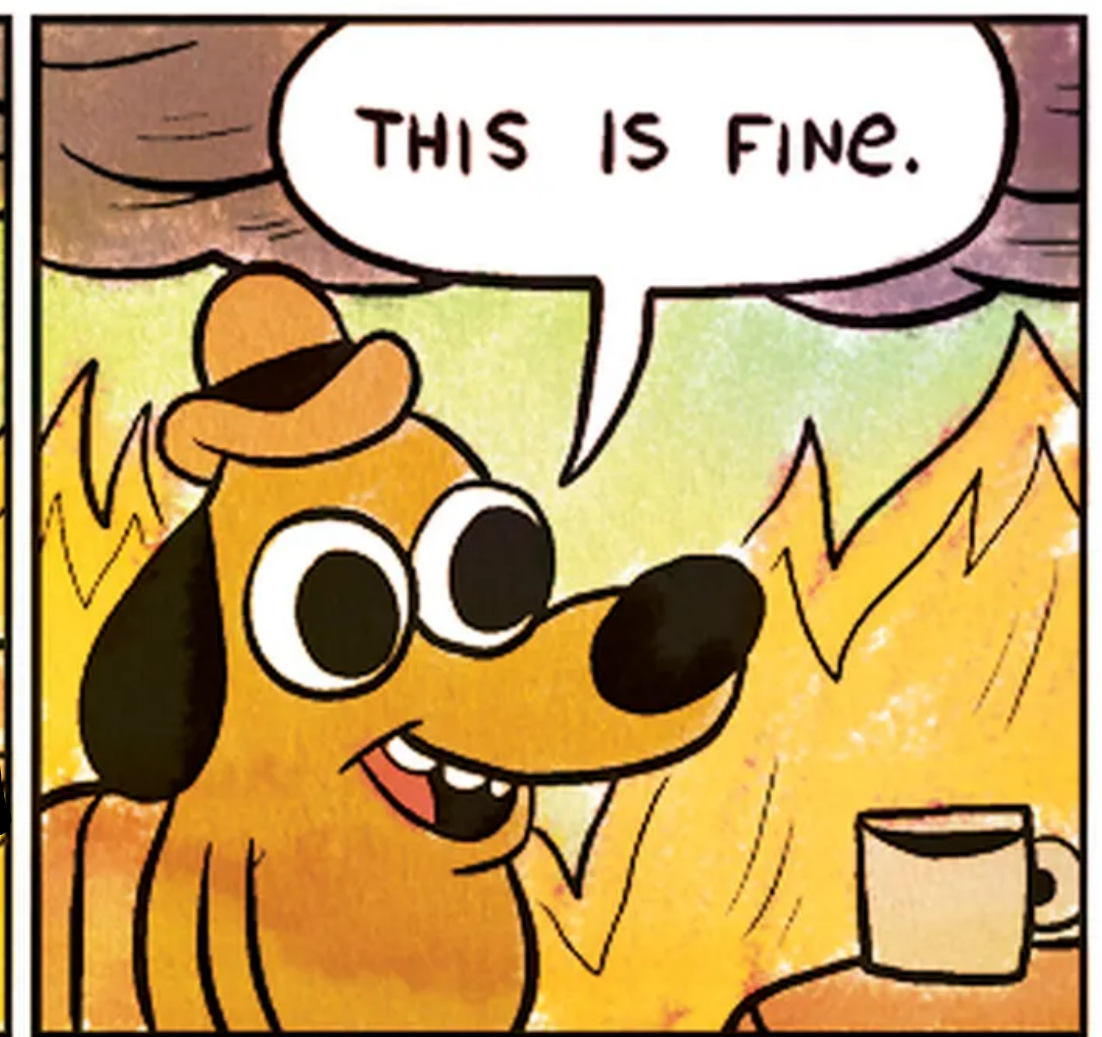# How do people develop for WMM?

## Think hard and document

For example, printk_ringbuffer

```
/*
 * Guarantee the state is loaded before copying the descriptor
 * content. This avoids copying obsolete descriptor
 * not apply to the descriptor state
 *
 * Memory barrier involvement:
 *
 * If desc_read:A reads from _prb
 * from _prb_commit:A.
 *
 * Relies on:
 *
 * WMB from _prb_commit:A to _prb_c
 *    matching
 * RMB from desc_read:A to desc_read
 */
smp_rmb(); /* LMM(desc_read:B) */
```

## Most just ignore topic

And wait to see what happens



THIS IS FINE.

**What about using tools?**

We are in 2022!

**There are scalable model checkers for WMM**!

e.g., GenMC, Dartagnan, VSync

# Our contributions

**BBQ: Block-based Bounded Queue**

- Novel block-based design

- Focus on enq-deq interference

- Support for out-of-order operations

- Verified for WMMs, pragmatically

# Our contributions

**BBQ: Block-based Bounded Queue**

- Novel block-based design

- Focus on enq-deq interference

- Support for out-of-order operations

- Verified for WMMs, pragmatically

**Bonus features**

- Single/multi producers/consumers

- Fixed- and variable-sized entries

- Retry-new and drop-old modes

- Use of efficient atomic operations
  - FAA and MAX (ARMv8.1 LSE)
  - No CAS at all if MAX available

# Agenda

☑ Motivation

☑ Stories and Challenges

*Interference, Out-of-order operations, Correctness on WMMs*

☐ BBQ – Block-based Bounded Queue

☐ Insights to Tackle the Challenges

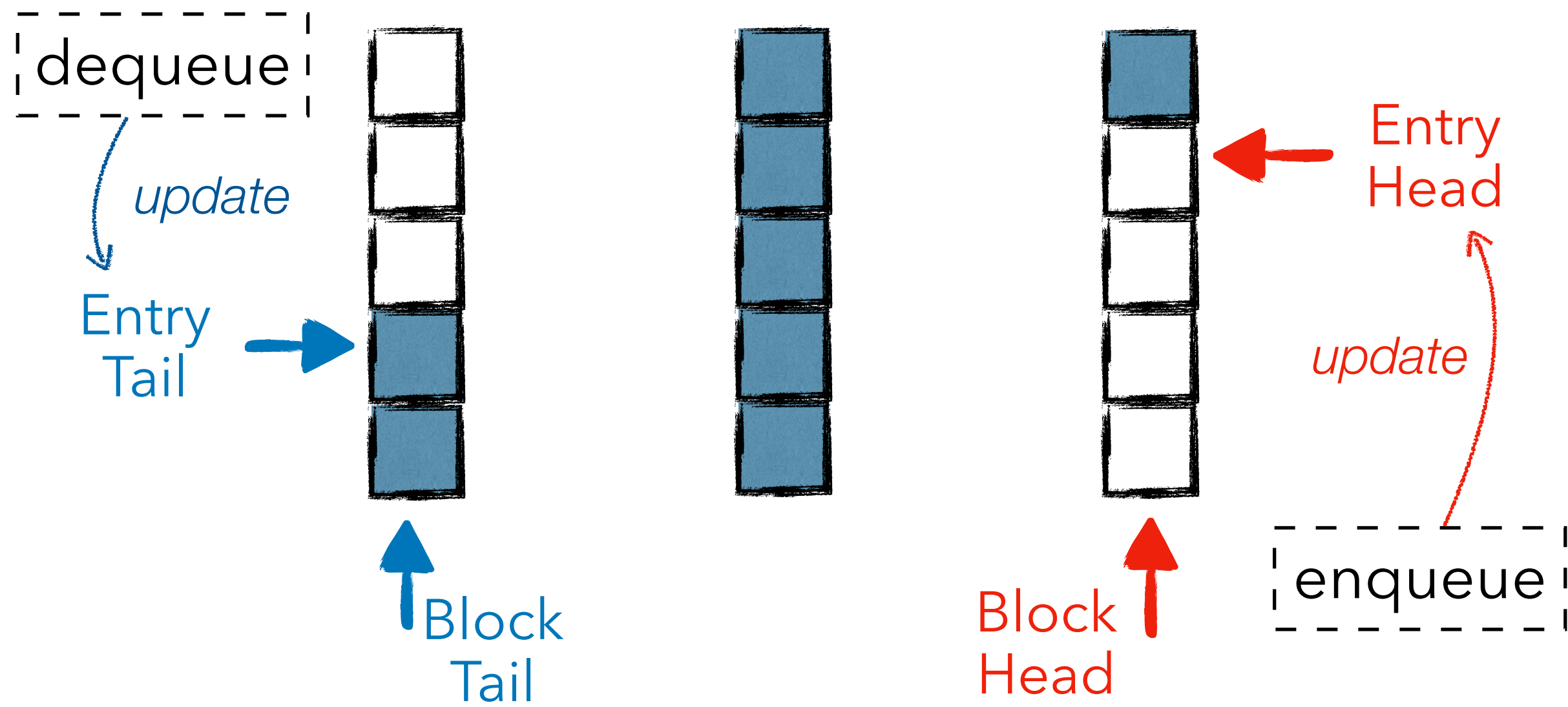☐ Selected Evaluation Results

# BBQ – Block-based Bounded Queue



- Ring buffer split into blocks

- Block Head points to current producer block

- Block Tail points to current consumer block

- In each block: Entry Head and Entry Tail

Tail

Head

Entry Head

Entry Tail

Block Tail

Block Head

12

## Enq-Deq interference

No interference when producer and consumer in different blocks

dequeue

*update*

Entry Tail

Block Tail

Entry Head

*update*

Block Head

enqueue

# Dealing with interferences

Enq-Deq interference

No interference when producer
and consumer in different blocks

dequeue

*update*

Entry
Tail

Block
Tail

Entry
Head

*update*

enqueue

Block
Head

Block head and tail only read
when moving to next block

13

# Dealing with interferences

## Enq-Deq interference

No interference when producer and consumer in different blocks

dequeue

*update*

Entry Tail

Entry Head

*update*

enqueue

Block Tail

Block Head

Block head and tail only read when moving to next block

## Enq-Enq interference

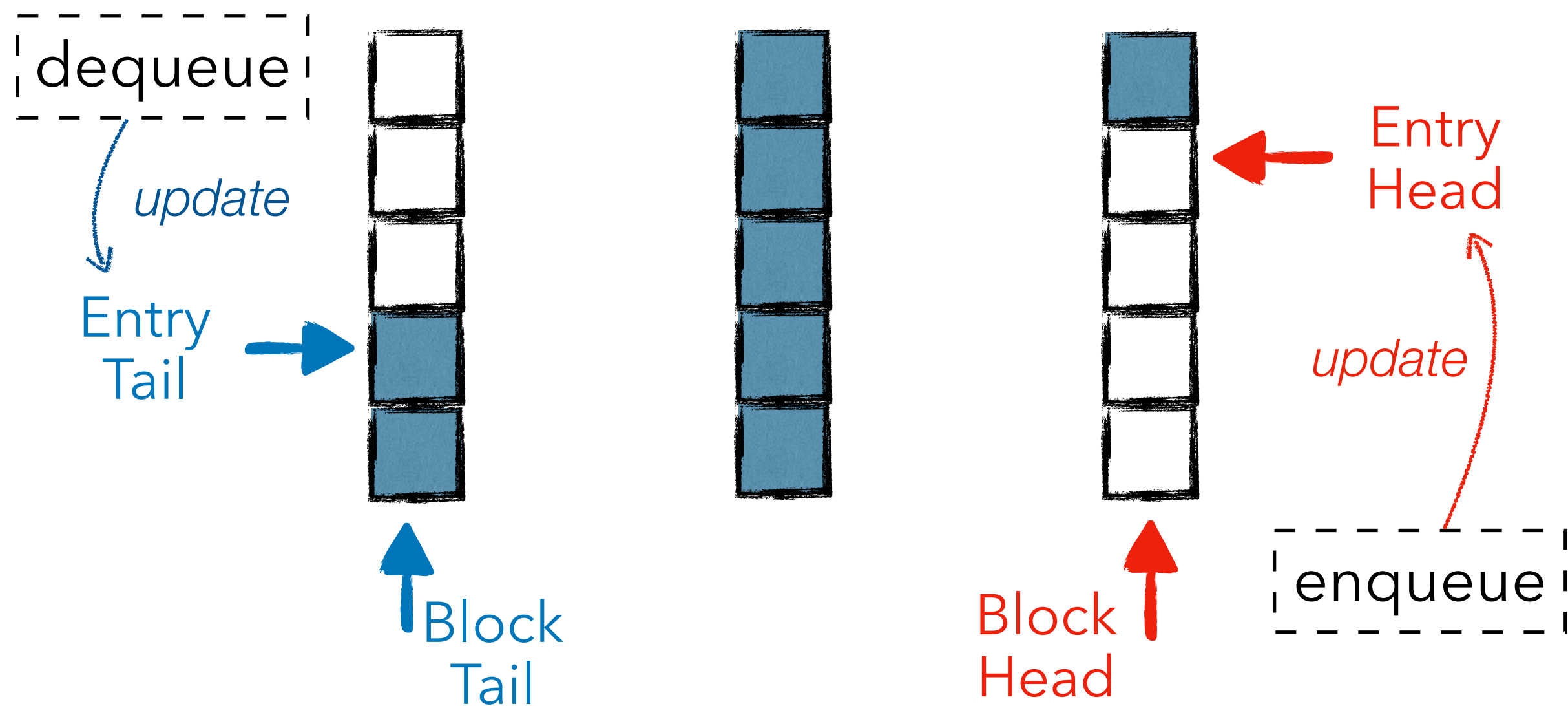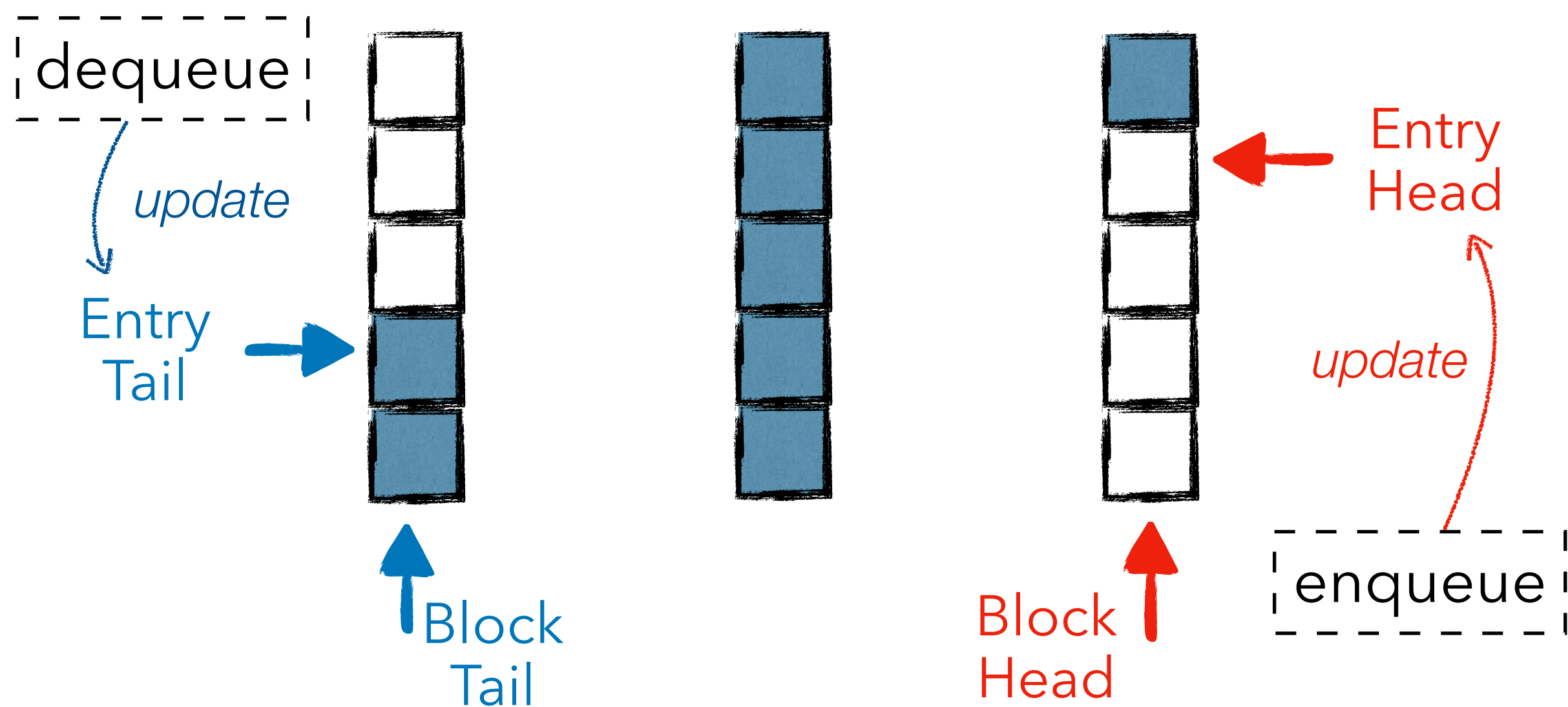Efficient use of FAA no side effects: neither rollback nor blocking

enqueue

*update*

enqueue

Entry Head

enqueue

Block Head

# Dealing with interferences



## Enq-Deq interference

No interference when producer and consumer in different blocks

dequeue

*update*

Entry Tail

Block Tail

Entry Head

*update*

enqueue

Block Head

Block head and tail only read when moving to next block

## Enq-Enq interference

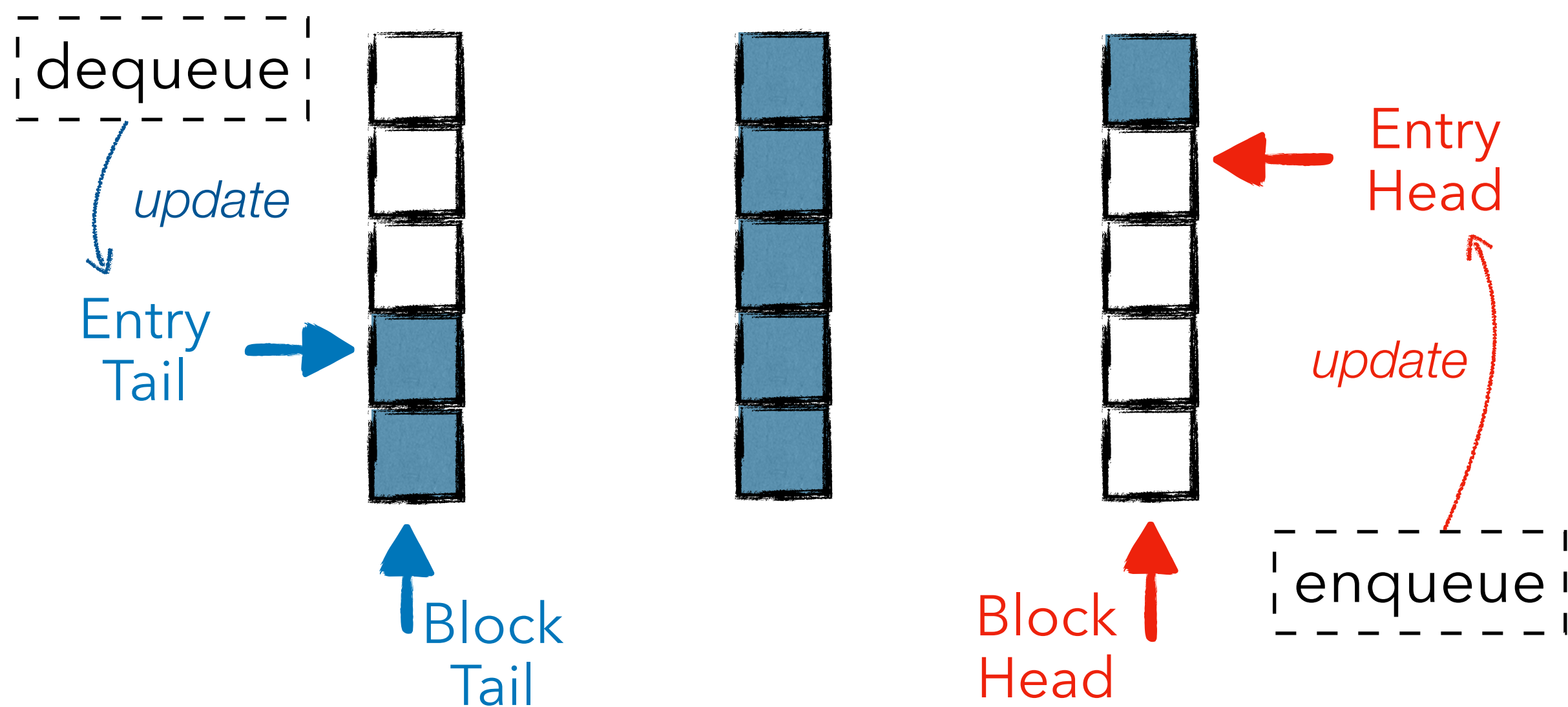Efficient use of FAA no side effects: neither rollback nor blocking

enqueue

*update*

Entry Head

enqueue

enqueue

Block Head

# Dealing with interferences

## Enq-Deq interference

No interference when producer and consumer in different blocks

dequeue

*update*

Entry Tail

Block Tail

Entry Head

*update*

enqueue

Block Head

Block head and tail only read when moving to next block

## Enq-Enq interference

Efficient use of FAA no side effects: neither rollback nor blocking

enqueue

*update*

Entry Head

enqueue

enqueue

Block Head

Head can move out-of-bounds, no consequence to following block.

# Dealing with interferences



## Enq-Deq interference

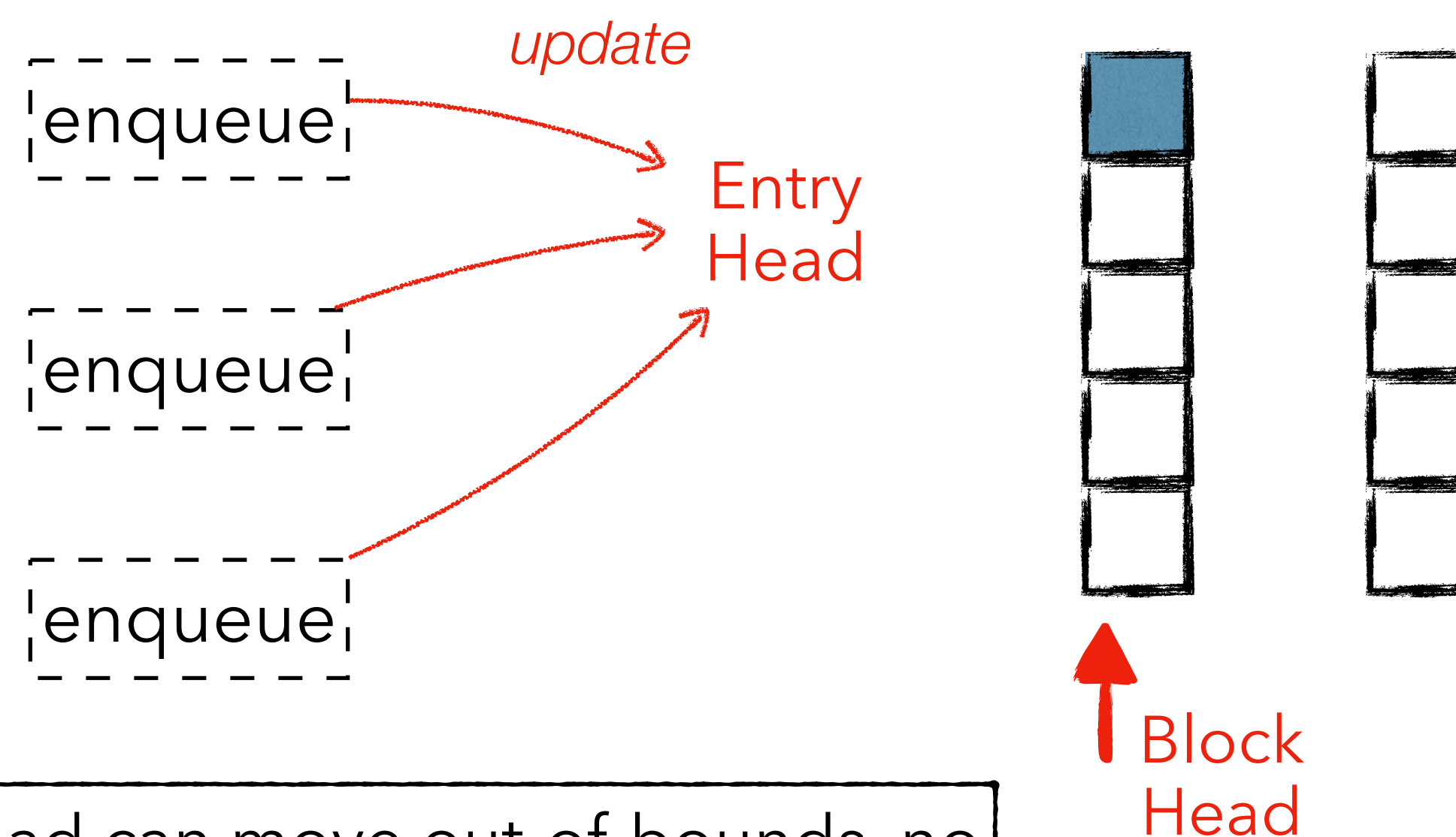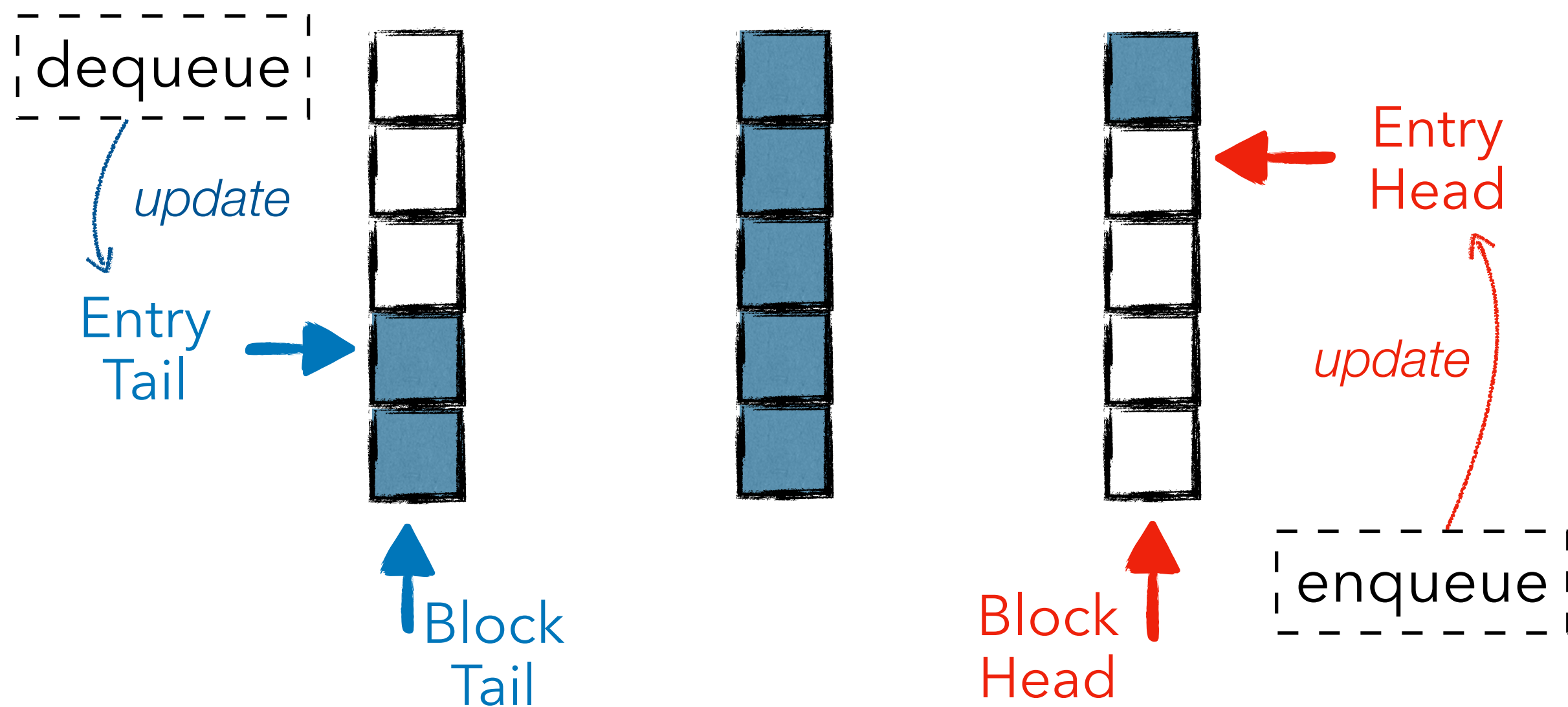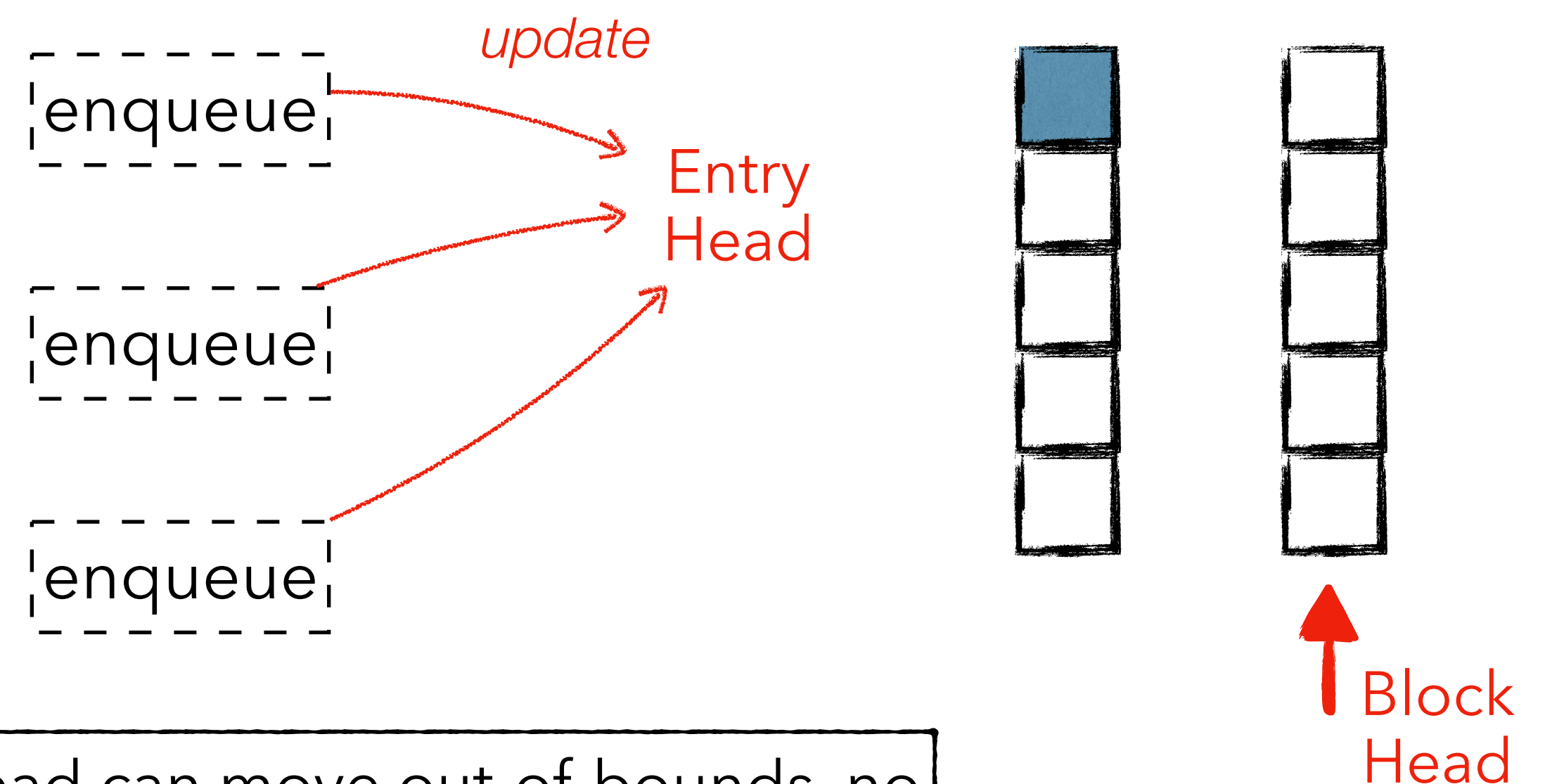No interference when producer and consumer in different blocks

dequeue

*update*

Entry Tail

Block Tail

Entry Head

*update*

enqueue

Block Head

Block head and tail only read when moving to next block

## Enq-Enq interference

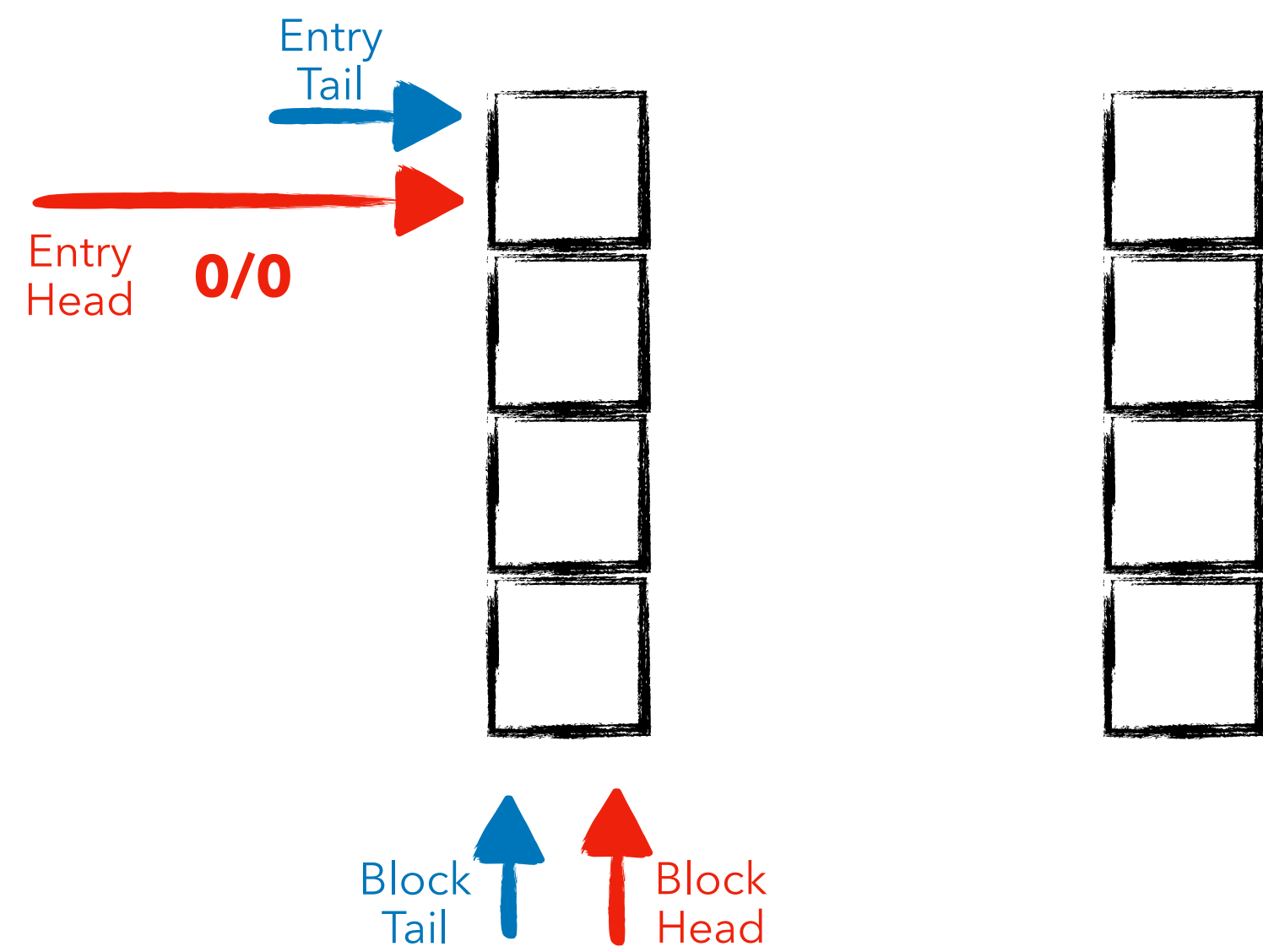Efficient use of FAA no side effects: neither rollback nor blocking

enqueue

*update*

Entry Head

enqueue

enqueue

Block Head

Head can move out-of-bounds, no consequence to following block.

OK

13

# Dealing with interferences
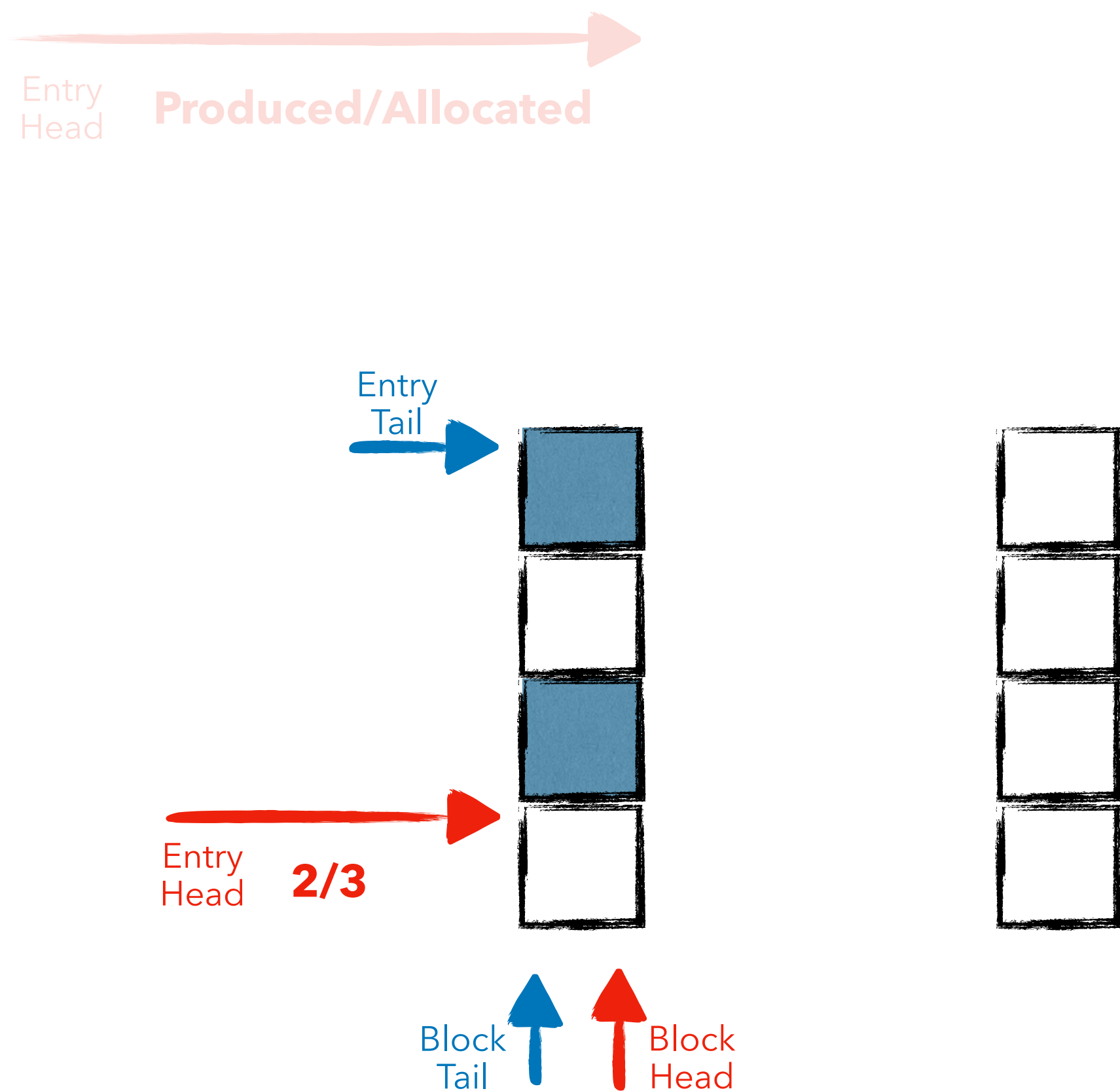
# Dealing with out-of-order operations

# Dealing with out-of-order operations

# Dealing with out-of-order operations

# Dealing with out-of-order operations



Enqueue calls:
- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

# Dealing with out-of-order operations

**Enqueue calls:**
- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues
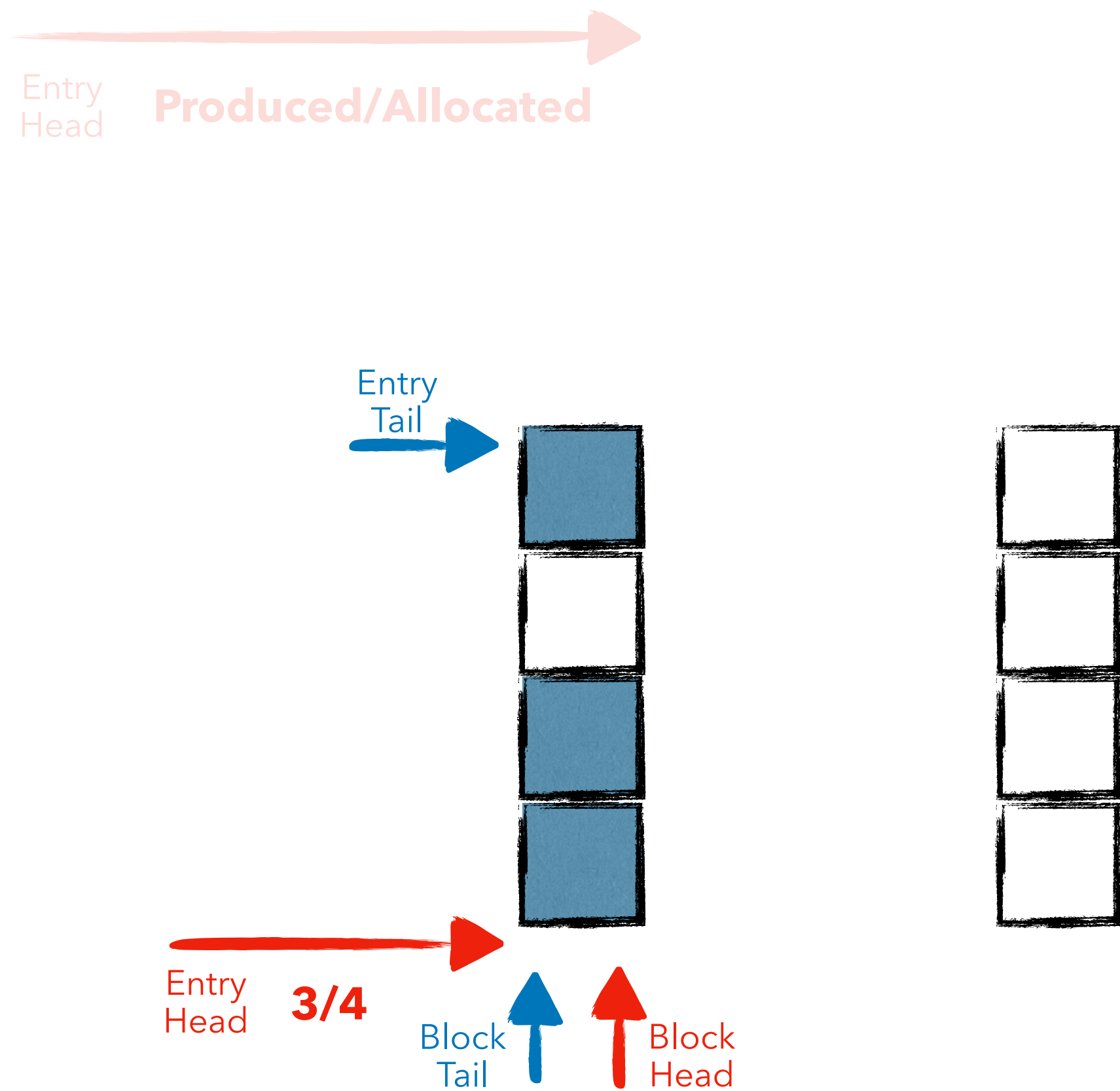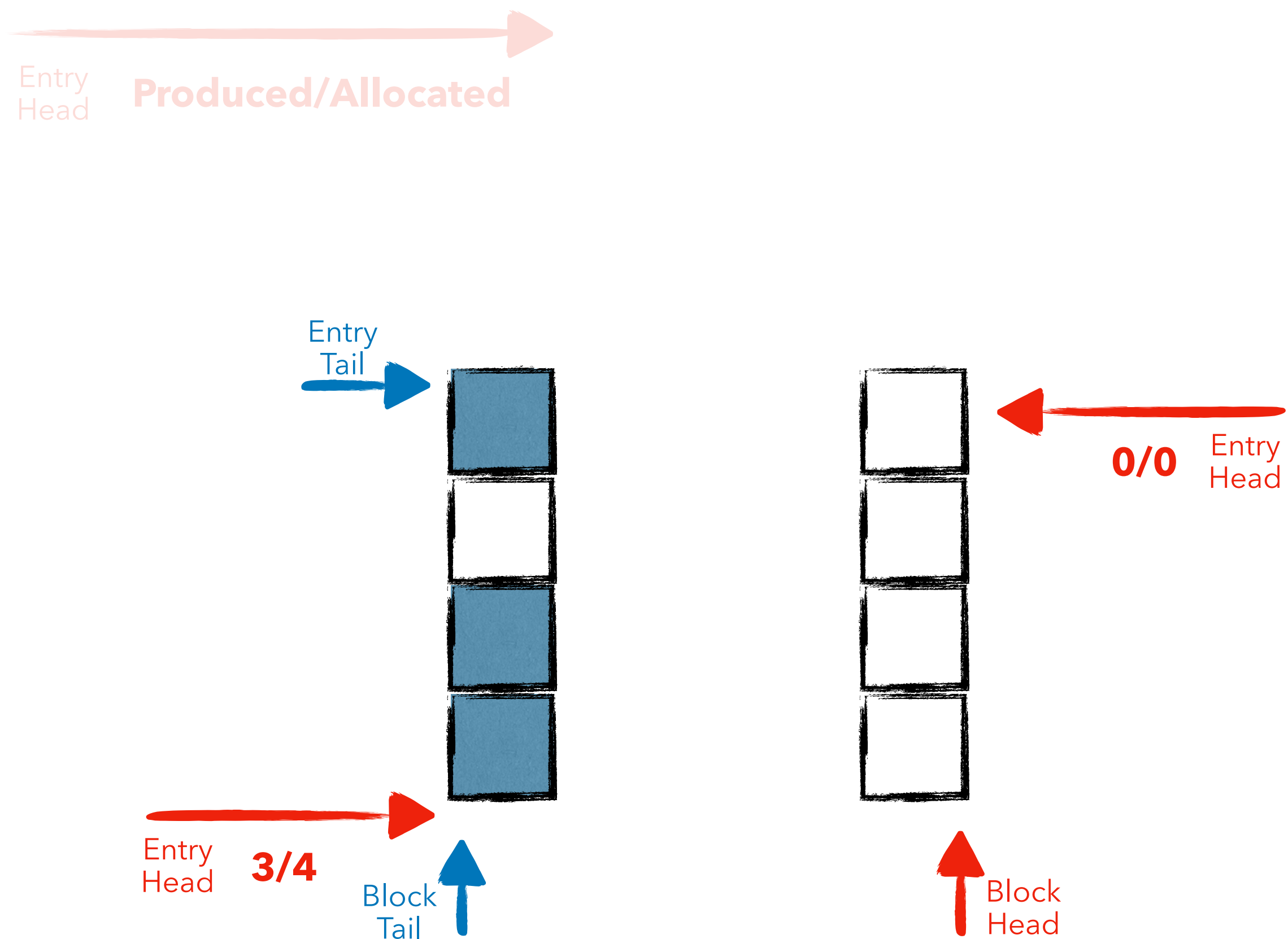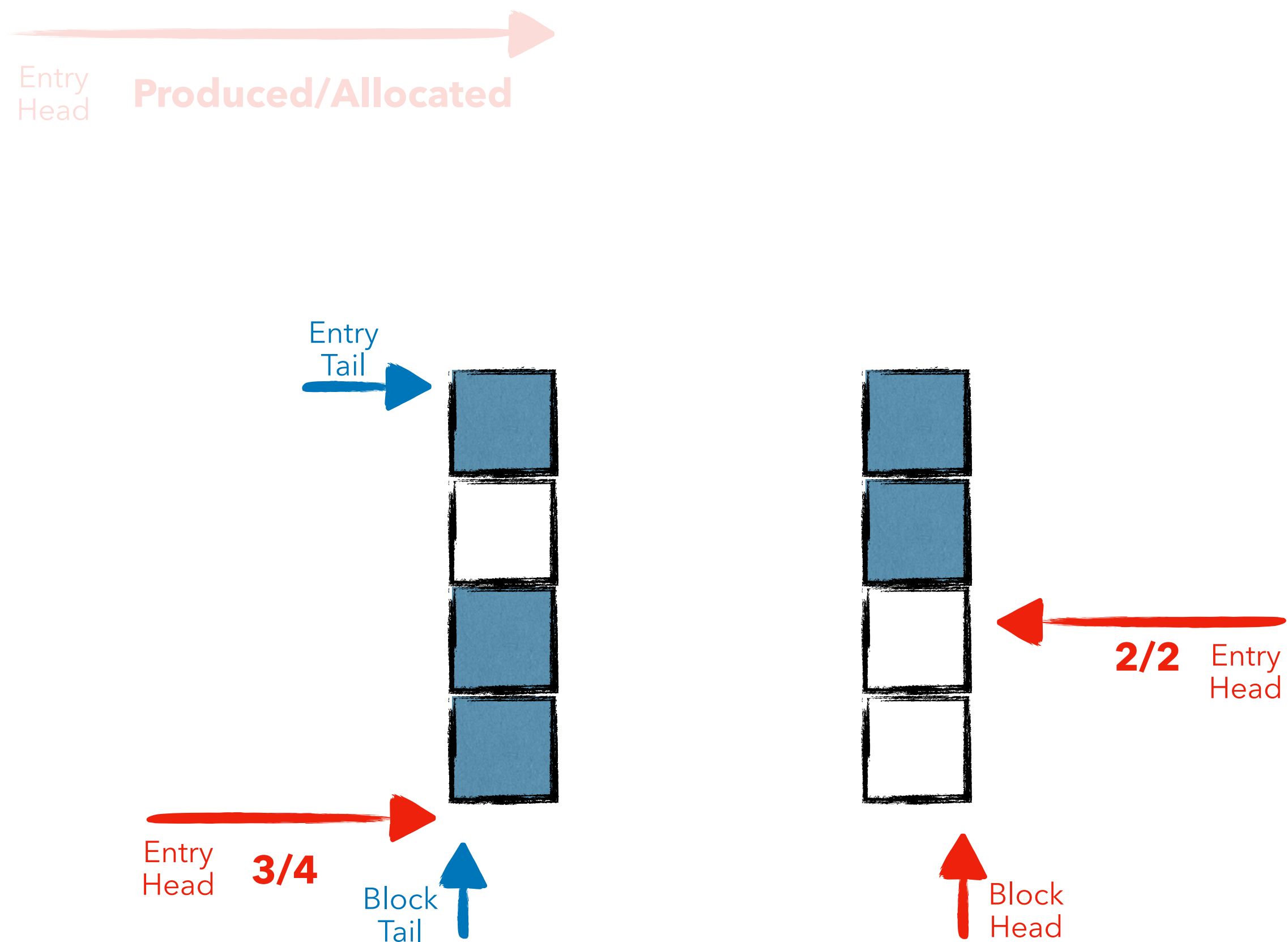
# Dealing with out-of-order operations

Enqueue calls:
- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

# Dealing with out-of-order operations



Entry Head    **Produced/Allocated**

Entry Tail

**2/2** Entry Head

Entry Head    **3/4**

Block Tail

Block Head

## Enqueue calls:
- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

## Dequeue calls:
- **return BUSY if an enqueue is ongoing** in same block
- **succeed when block full** or when Produced = Allocated

# Dealing with out-of-order operations

Entry Head
**Produced/Allocated**

Entry Tail

Entry Head **4/4**

Block Tail

**2/2** Entry Head

Block Head

## Enqueue calls:

- **do not wait for others** in same block to complete
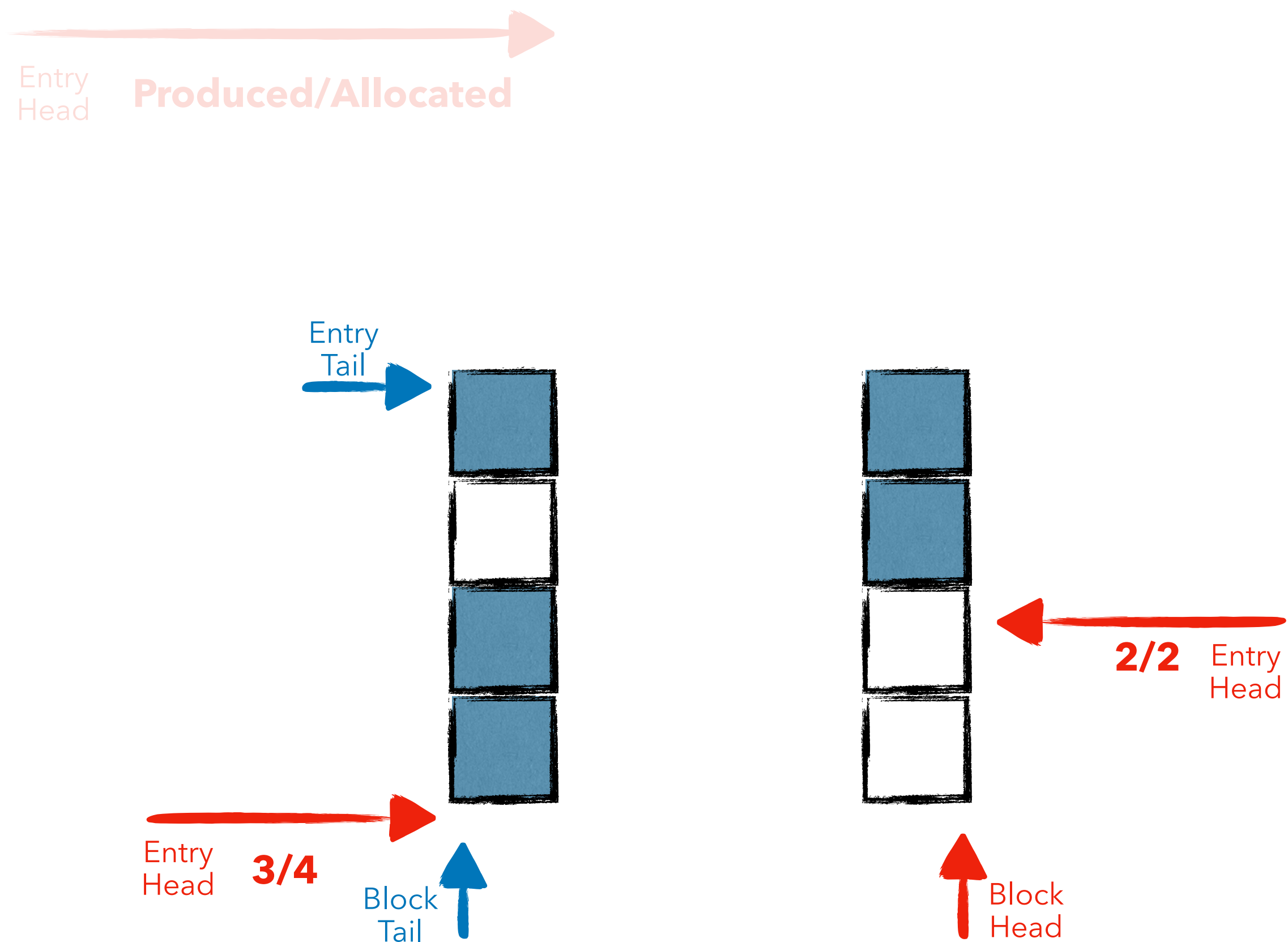- can **move to next block** even if current block has ongoing enqueues

## Dequeue calls:

- **return BUSY if an enqueue is ongoing** in same block
- **succeed when block full** or when Produced = Allocated

# Dealing with out-of-order operations



Entry Head **Produced/Allocated**

**2/2** Entry Head

Entry Tail

Entry Head **4/4**

Block Tail

Block Head
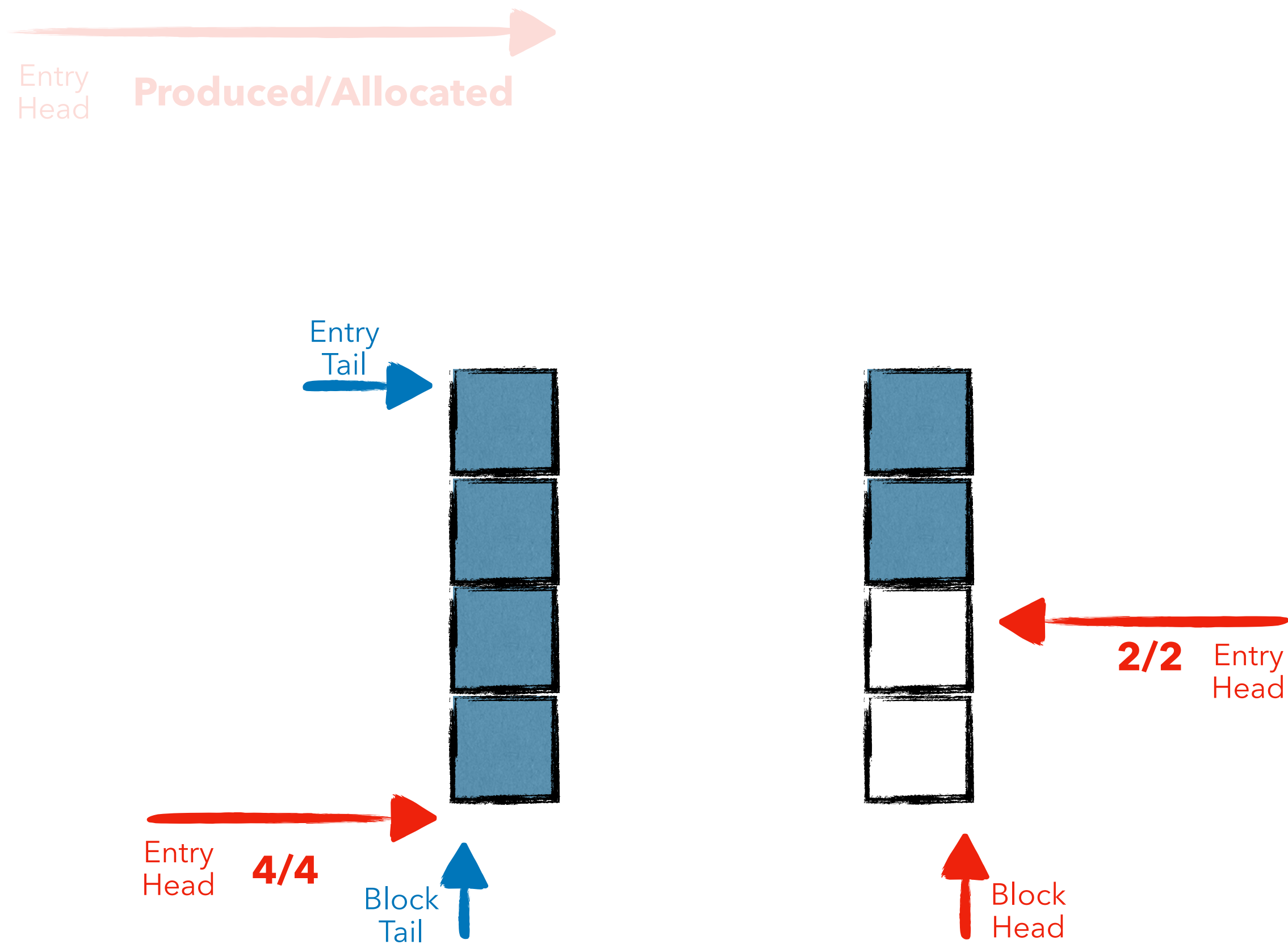
Enqueue calls:
- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

Dequeue calls:
- **return BUSY if an enqueue is ongoing** in same block
- **succeed when block full** or when Produced = Allocated
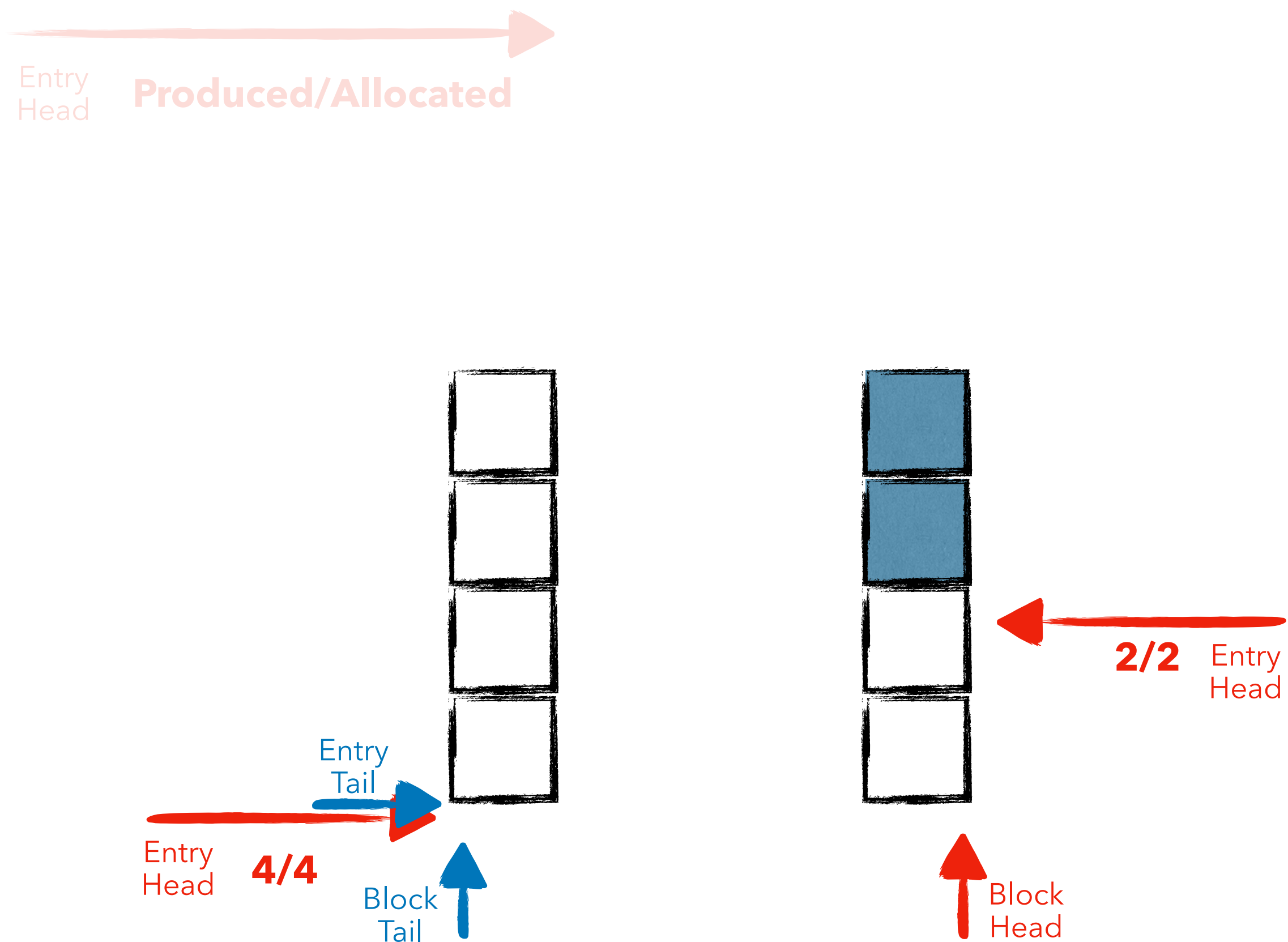
# Dealing with out-of-order operations

Entry
Head

**Produced/Allocated**

Entry
Tail

**2/2** Entry
Head

Entry
Tail

Entry
Head **4/4**

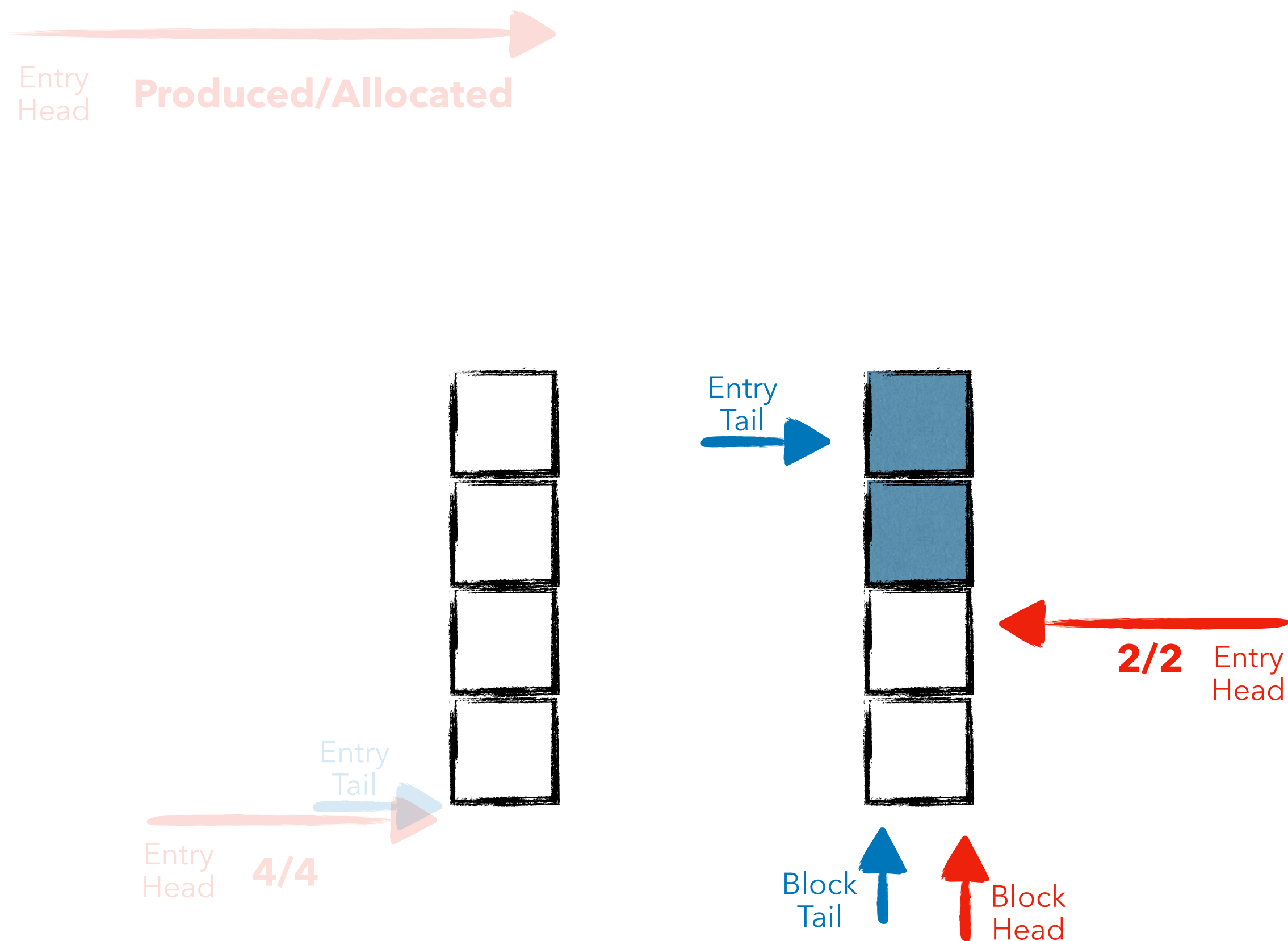Block
Tail

Block
Head

**Enqueue calls:**

- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

**Dequeue calls:**

- **return BUSY if an enqueue is ongoing** in same block
- **succeed when block full** or when Produced = Allocated

14

# Dealing with out-of-order operations



**Enqueue calls:**

- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

**Dequeue calls:**

- **return BUSY if an enqueue is ongoing** in same block
- **succeed when block full** or when Produced = Allocated
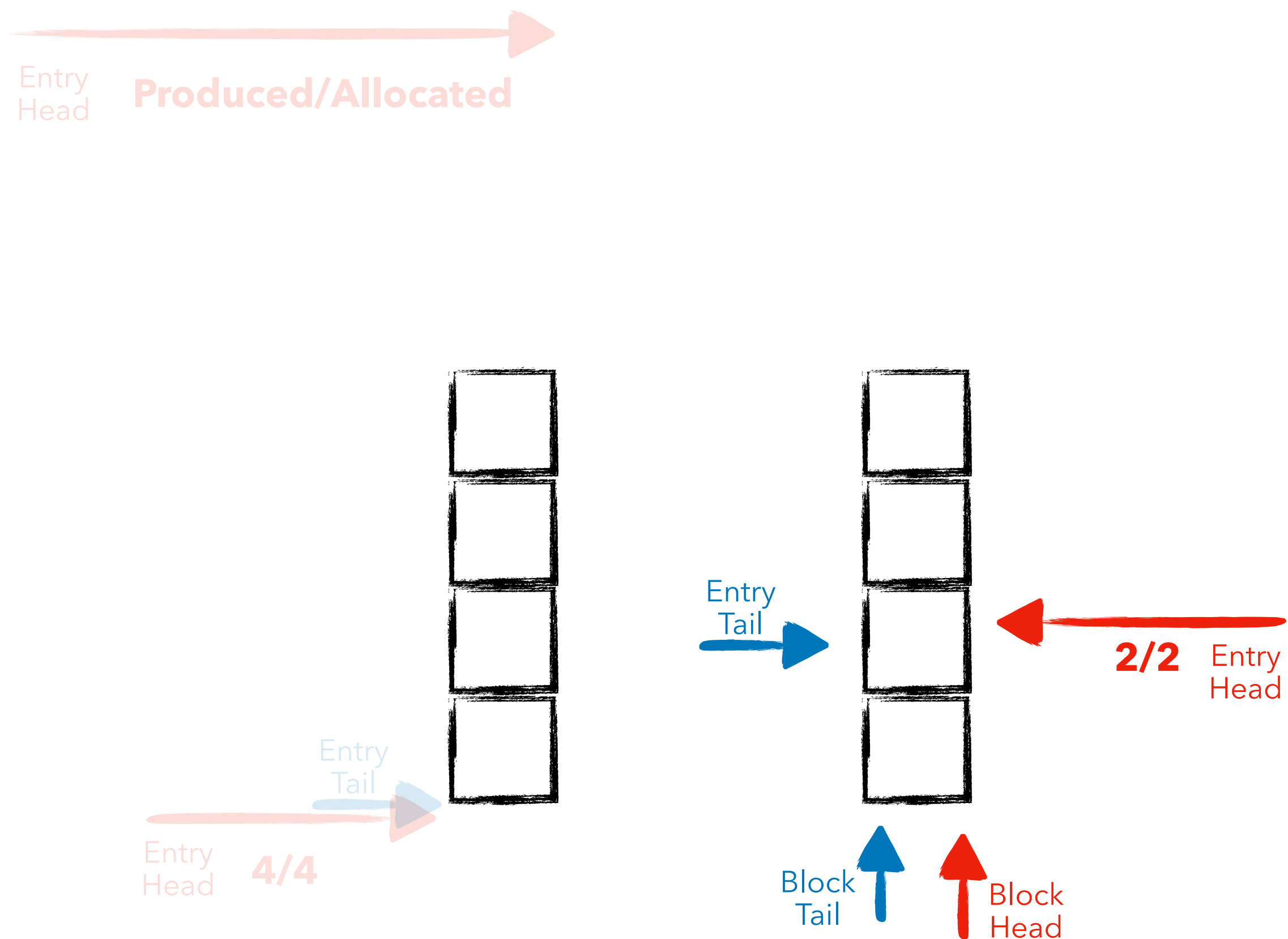
# Dealing with out-of-order operations

Entry
Head     **Produced/Allocated**

**Many cool tricks in the paper:**

- update block and entry indices at the same time without D-CAS

- Avoid ABA issues with versioning

- Cache block indices for speed

**Enqueue calls:**

- **do not wait for others** in same block to complete
- can **move to next block** even if current block has ongoing enqueues

**Dequeue calls:**

- **return BUSY if an enqueue is ongoing** in same block
- **succeed when block full** or when Produced = Allocated

14

# Correctness on WMMs with **practical** verification



DPDK-like algorithm

~10 atomics

Part of BBQ

More than 20 atomics

BBQ is not easy to digest

```
1  enqueue(data){
2  again:
3    ph = LOAD(P.head);
4    pn = ph + 1;
5    if (pn > LOAD(C.tail) + SZ)
6      return FULL;
7    if (!CAS(P.head, ph, pn))
8      goto again;
9    entry[pn % SZ] = data;
10   while(LOAD(P.tail) != ph);
11   STORE(P.tail, pn);
12   return OK;
13 }

14 dequeue(){
15 again:
16   ch = LOAD(C.head);
17   cn = ch + 1;
18   if (cn > LOAD(P.tail))
19     return EMPTY;
20   if (!CAS(C.head, ch, cn))
21     goto again;
22   data = entry[cn % SZ];
23   while(LOAD(C.tail) != ch);
24   STORE(C.tail, cn);
25   return data;
26 }
```

```
1  <Head, Block> BBQ<T>::get_phead_and_block(){
2    ph = LOAD(phead);
3    return (ph, blocks[ph.idx]);
4  }
5  state BBQ<T>::allocate_entry(Block blk){
6    if (LOAD(blk.allocated).off >= BLOCK_SIZE)
7      return BLOCK_DONE;
8    old = FAA(blk.allocated, 1).off;
9    if (old >= BLOCK_SIZE)
10     return BLOCK_DONE;
11   return ALLOCATED(EntryDesc{.block=blk, .offset=old});
12 }
13 void BBQ<T>::commit_entry(EntryDesc e, T data){
14   e.block.entries[e.offset] = data;
15   ADD(e.block.committed, 1);
16 }
17 state BBQ<T>::advance_phead(Head ph) {
18   nblk = blocks[(ph.idx + 1) % BLOCK_NUM];
19   cons = LOAD(nblk.consumed);
20   if (cons.vsn < ph.vsn ||
21      (cons.vsn == ph.vsn && cons.off != BLOCK_SIZE)) {
22     reserved = LOAD(nblk.reserved);
23     if (reserved.off == cons.off) return NO_ENTRY;
24     else return NOT_AVAILABLE;
25   }
26   cmtd = LOAD(nblk.committed);
27   if (cmtd.vsn == ph.vsn && cmtd.off != BLOCK_SIZE)
28     return NOT_AVAILABLE;
29   MAX(nblk.committed, Cursor{.vsn=ph.vsn + 1});
30   MAX(nblk.allocated, Cursor{.vsn=ph.vsn + 1});
31   MAX(phead, ph + 1);
32   return SUCCESS;
33 }
34 class BBQ<T> {
35   shared<Head> phead, chead;
36   Block<T>[] blocks;
37 }
38 class Block<T> {
39   shared<Cursor> allocated, committed;
40   shared<Cursor> reserved, consumed;
41   T[] entries;
42 }
43 class EntryDesc {
44   Block block; Offset offset; Version version; }
```

```
45 <Head, Block> BBQ<T>::get_chead_and_block(){
46   ch = LOAD(chead);
47   return (ch, blocks[ch.idx]);
48 }
49 state BBQ<T>::reserve_entry(Block blk){
50 again:
51   reserved = LOAD(blk.reserved);
52   if (reserved.off < BLOCK_SIZE) {
53     committed = LOAD(blk.committed);
54     if (reserved.off == committed.off)
55       return NO_ENTRY;
56     if (committed.off != BLOCK_SIZE){
57       allocated = LOAD(blk.allocated);
58       if (allocated.off != committed.off)
59         return NOT_AVAILABLE;
60     }
61     if (MAX(blk.reserved, reserved + 1) == reserved)
62       return RESERVED((EntryDesc){.block=blk,
63         .offset=reserved.off, .version=reserved.vsn});
64     else goto again;
65   }
66   return BLOCK_DONE(reserved.vsn);
67 }
68 T BBQ<T>::consume_entry(EntryDesc e){
69   data = e.block.entries[e.offset];
70   ADD(e.block.consumed, 1);
71   allocated = LOAD(e.block.allocated);
72   if (allocated.vsn != e.version) return NULL;
73   return data;
74 }
75 bool BBQ<T>::advance_chead(Head ch, Version vsn){
76   nblk = blocks[(ch.idx + 1) % BLOCK_NUM];
77   committed = LOAD(nblk.committed);
78   if (committed.vsn != ch.vsn + 1)
79     return false;
80   MAX(nblk.consumed, Cursor{.vsn=ch.vsn + 1});
81   MAX(nblk.reserved, Cursor{.vsn=ch.vsn + 1});
82   if (committed.vsn < vsn + (ch.idx == 0))
83     return false;
84   MAX(nblk.reserved, Cursor{.vsn=committed.vsn});
85   MAX(chead, ch + 1);
86   return true;
87 }
```

retry-new mode | drop-old mode

# Correctness on WMMs with **practical** verification

- **Long stress testing**
  *by engineers*

- **Identification of corner cases**
  *by WMM experts and engineers*

- **Only a few corner cases necessary**
  *queue full/empty, FIFO, wrap-around*

- **Model check corner cases on WMM**
  *by engineers*

- **3 bugs found model checking them**
  *Not found while stress testing*

- **Reproducible on real hardware**
  *Test cases were built in retrospect*

# Agenda

☑ Motivation

☑ Stories and Challenges

 *Interference, Out-of-order operations, Correctness on WMMs*

☑ Block-based Bounded Queue (BBQ)

☑ Insights to Tackle the Challenges

☐ Selected Evaluation Results

# Micro-benchmark Results – SPSC

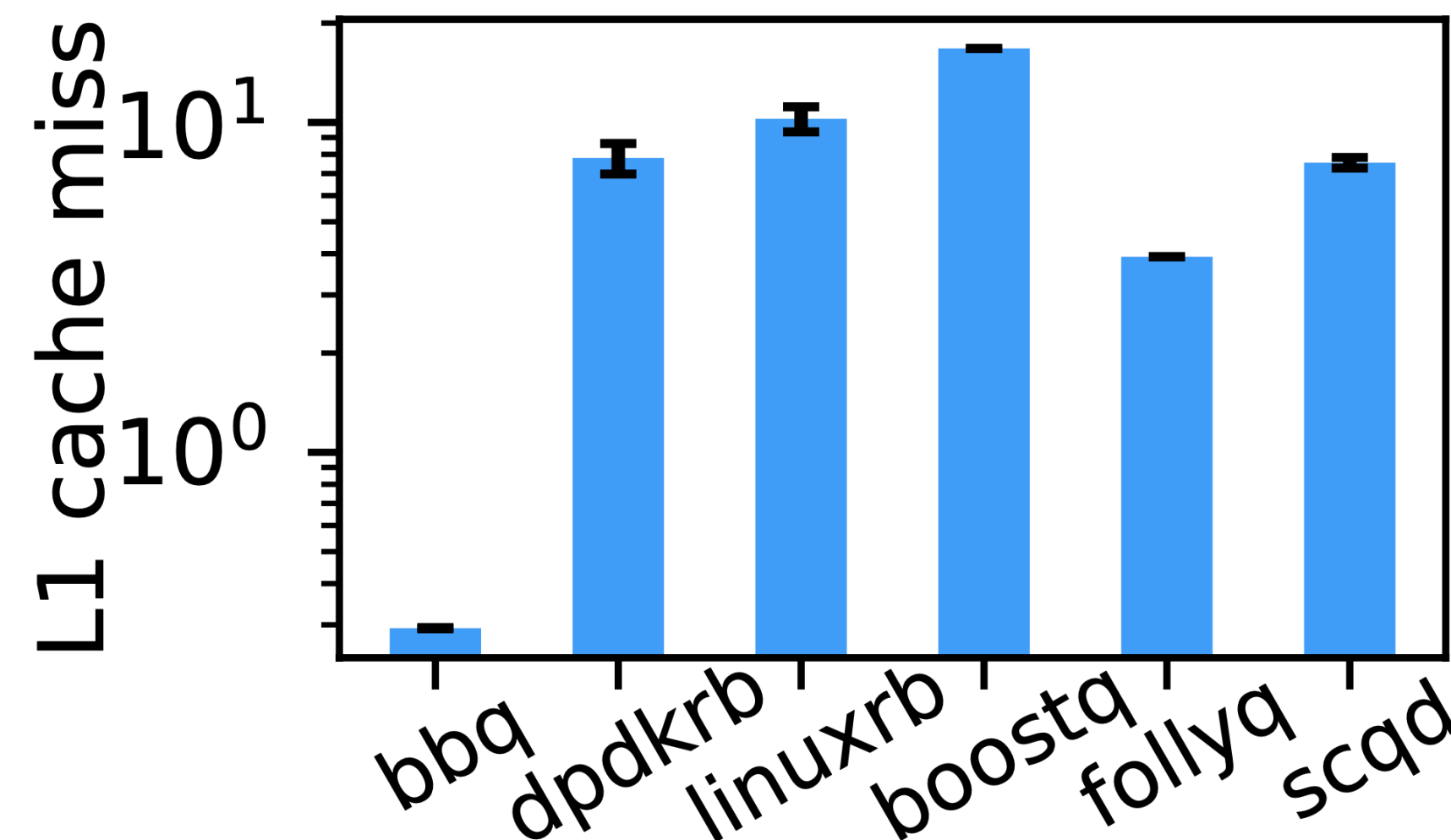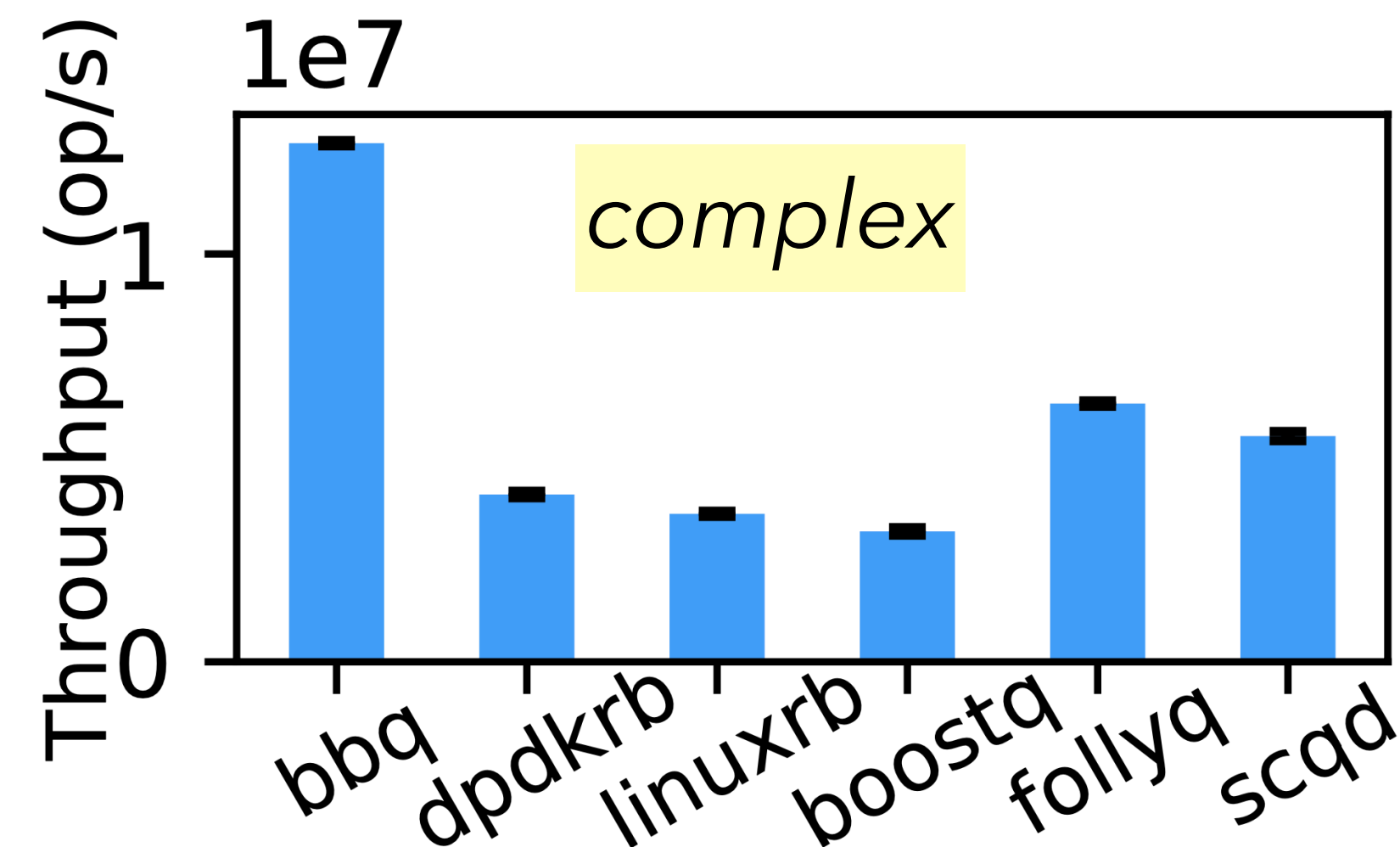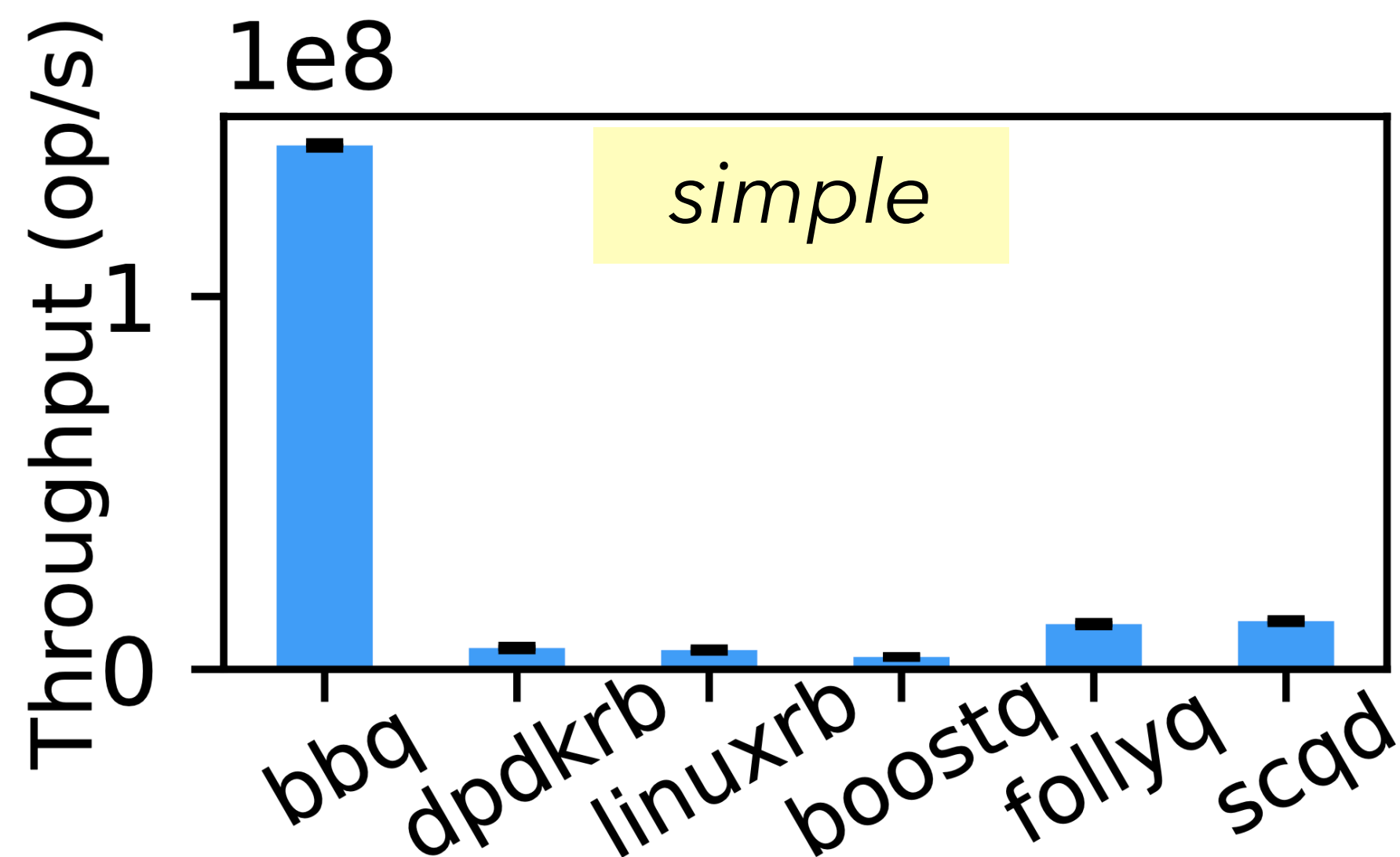*Compared against 5 state-of-the-art bounded queues*

- x86 machines with 88 hyper-threads

- 8 bytes data size, 32*k* bytes memory usage

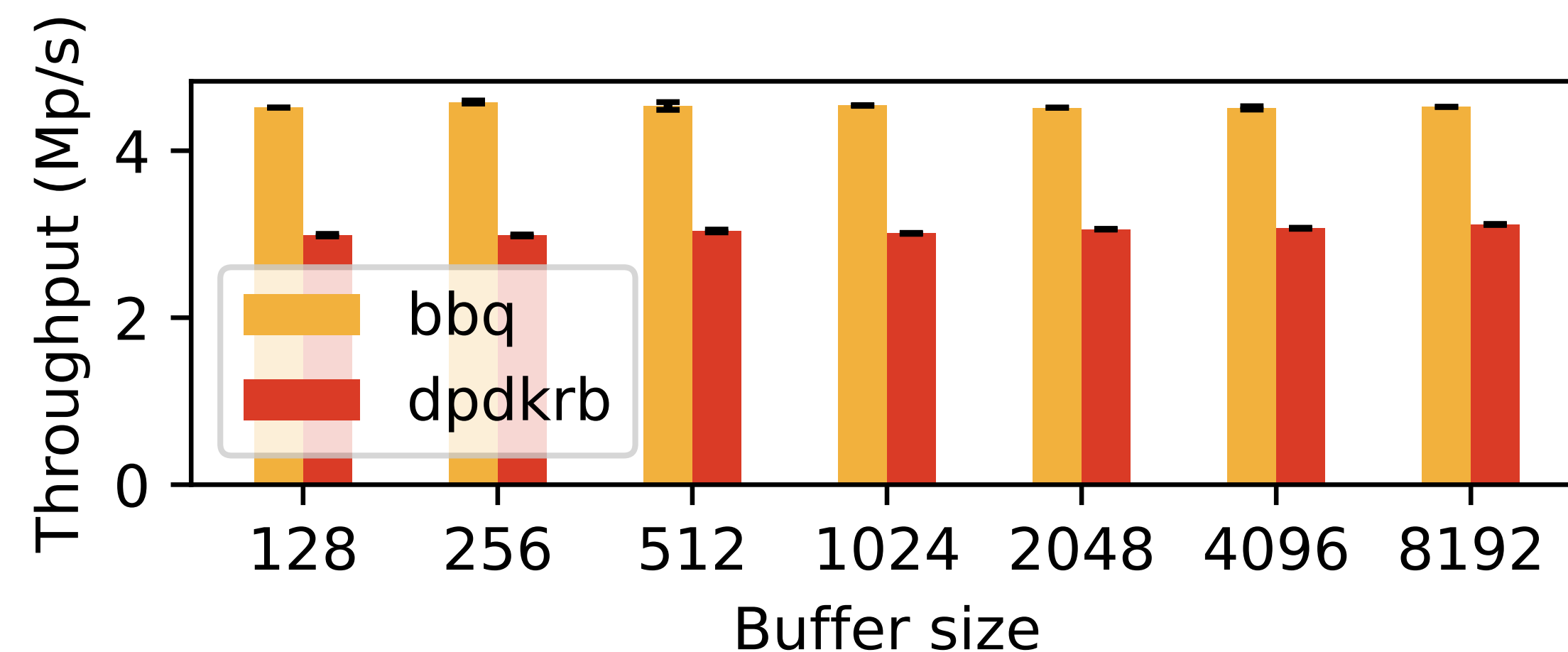- *simple:* 11.3x to 42.4x higher throughput
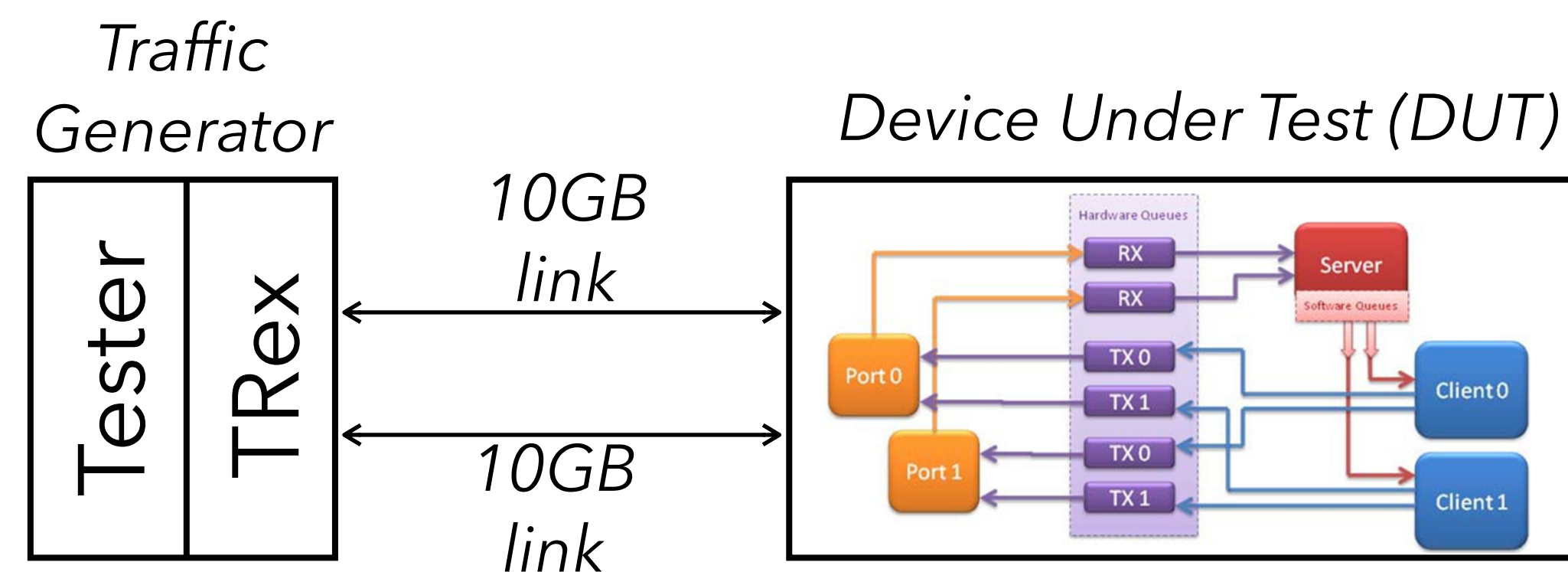
# Micro-benchmark Results – SPSC

*Compared against 5 state-of-the-art bounded queues*

- x86 machines with 88 hyper-threads

- 8 bytes data size, 32*k* bytes memory usage

- *simple:* 11.3x to 42.4x higher throughput

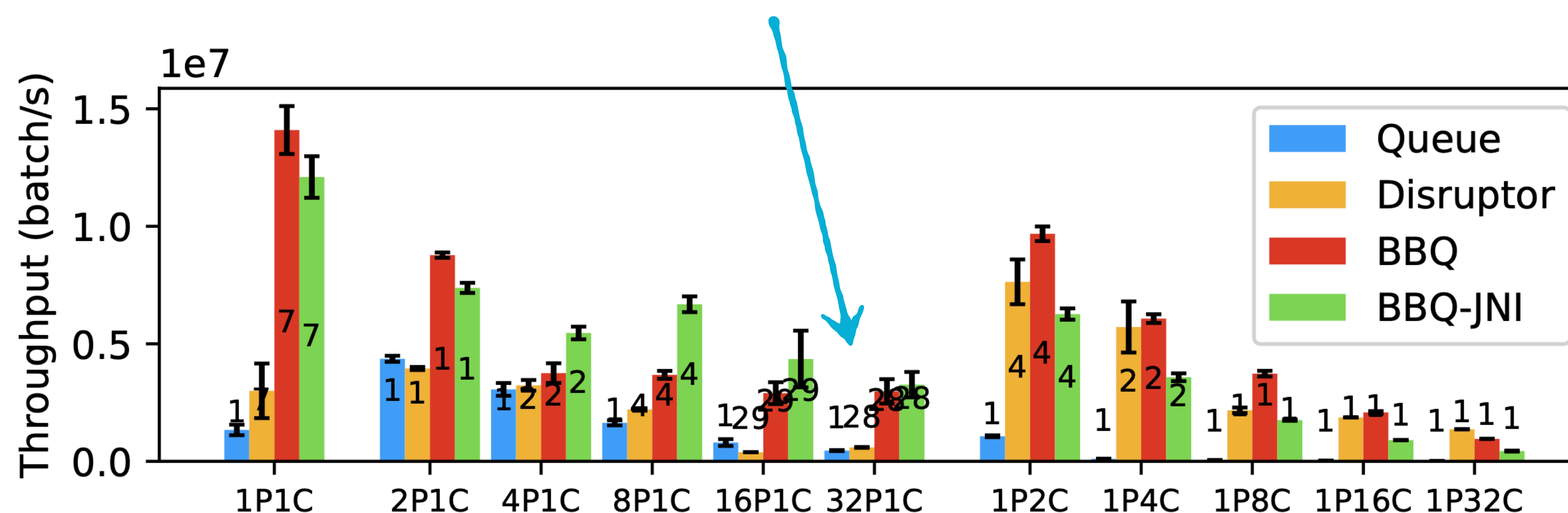- *complex*: **at least 2x higher than FollyQ**

# DPDK Test Suite (DTS) – Multiprocess benchmark

- **Device Under Test**
  - One server process
    *receiving and distributing packets*
  - Two client processes
    *performing level-2 packet forwarding*

- Tester and traffic generator run on another machine
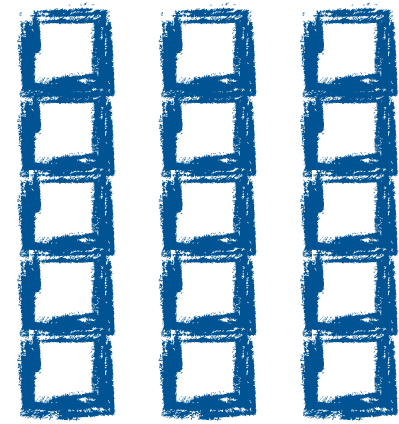
- BBQ yields **1.5x throughput** of DPDK

# Macro-benchmark Results – Disruptor

- LMAX Disruptor: bounded queue for high-performance trading

- Compared on three official Disruptor benchmarks
  Against Java queue, BBQ in Java, and BBQ in C via JNI

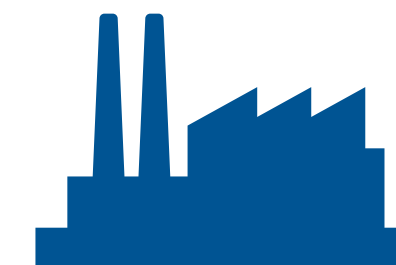- With 32 producers, **BBQ yields 3 Mop/s and Disruptor 0.6 Mop/s**

# Wrap up

**BBQ is a novel ring buffer design**
- Reduces enq-deq interference
- Supports out-of-order operations
- Model checked for WMMs

**Large spectrum of scenarios**
- Single/Multi Consumer/Producer
- Retry-new and Drop-old modes
- Etc

Greatly outperforms several industrial ring buffers

Please **look up the paper** for many more results

# Thank you! Questions?

(BTW, we are hiring in Dresden and Munich…)