

# Distributed Queues: Faster Pools and Better Queues

Christoph Kirsch  
Universität Salzburg

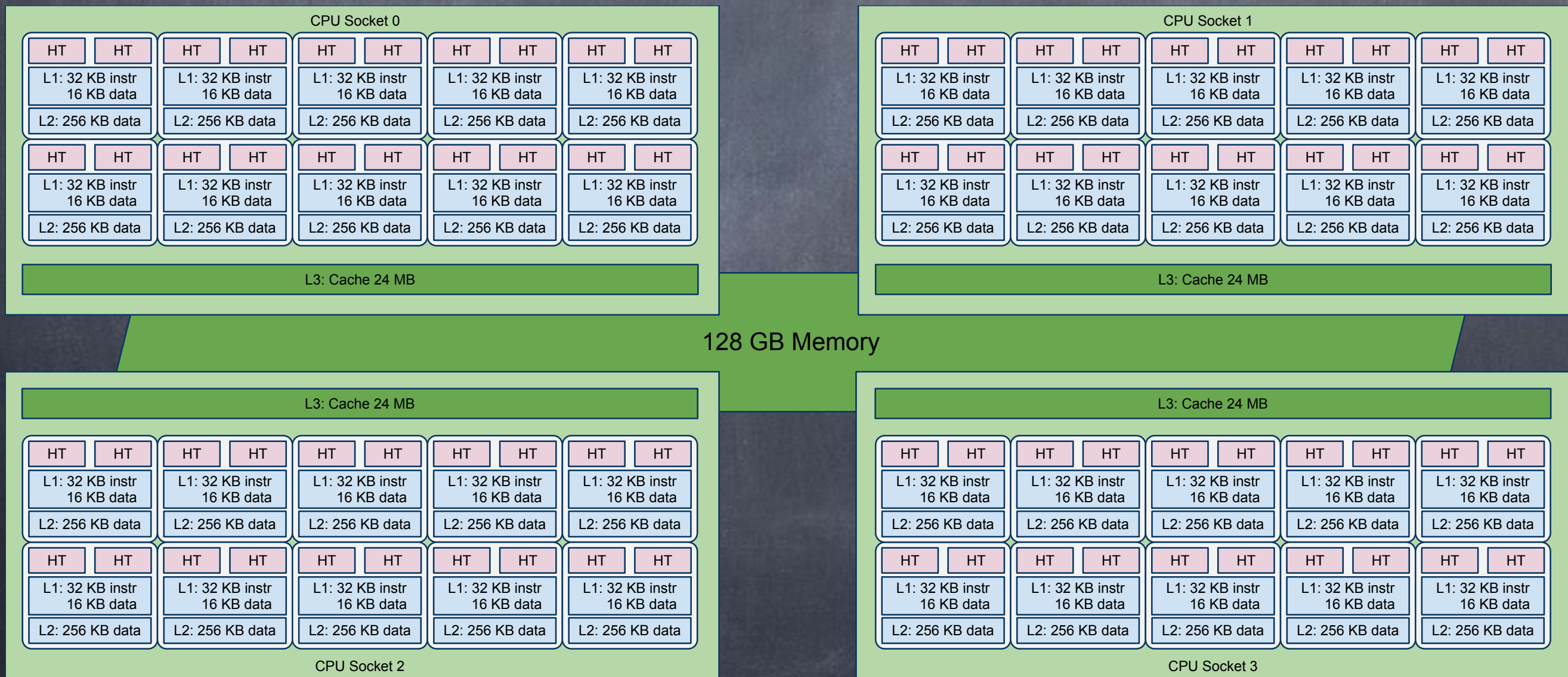


Oracle, Belmont, California, December 2012



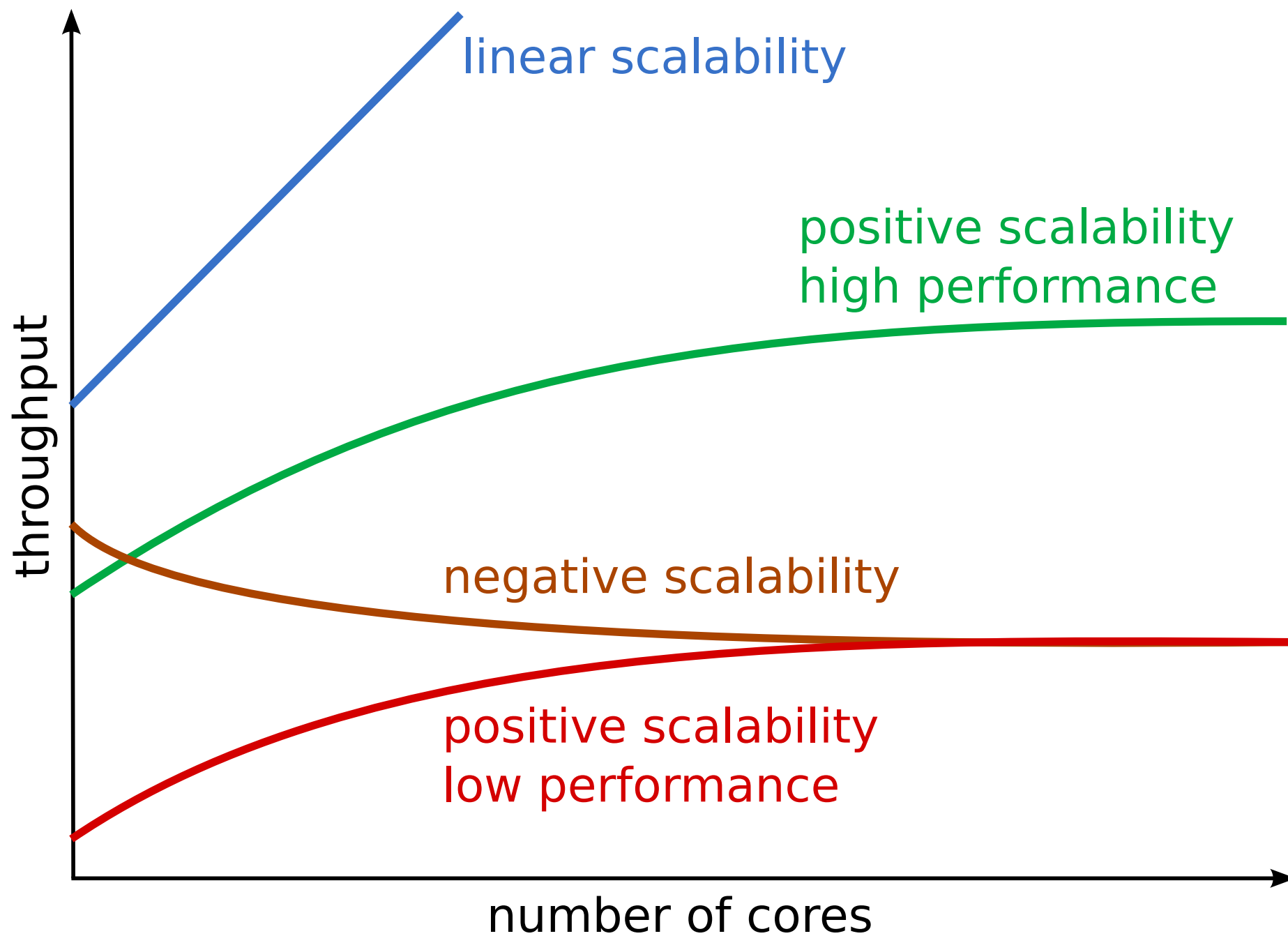
Joint work w/ A. Haas,  
M. Lippautz, H. Payer,  
A. Sokolova and our  
collaborators at IST Austria  
T. Henzinger, A. Sezgin

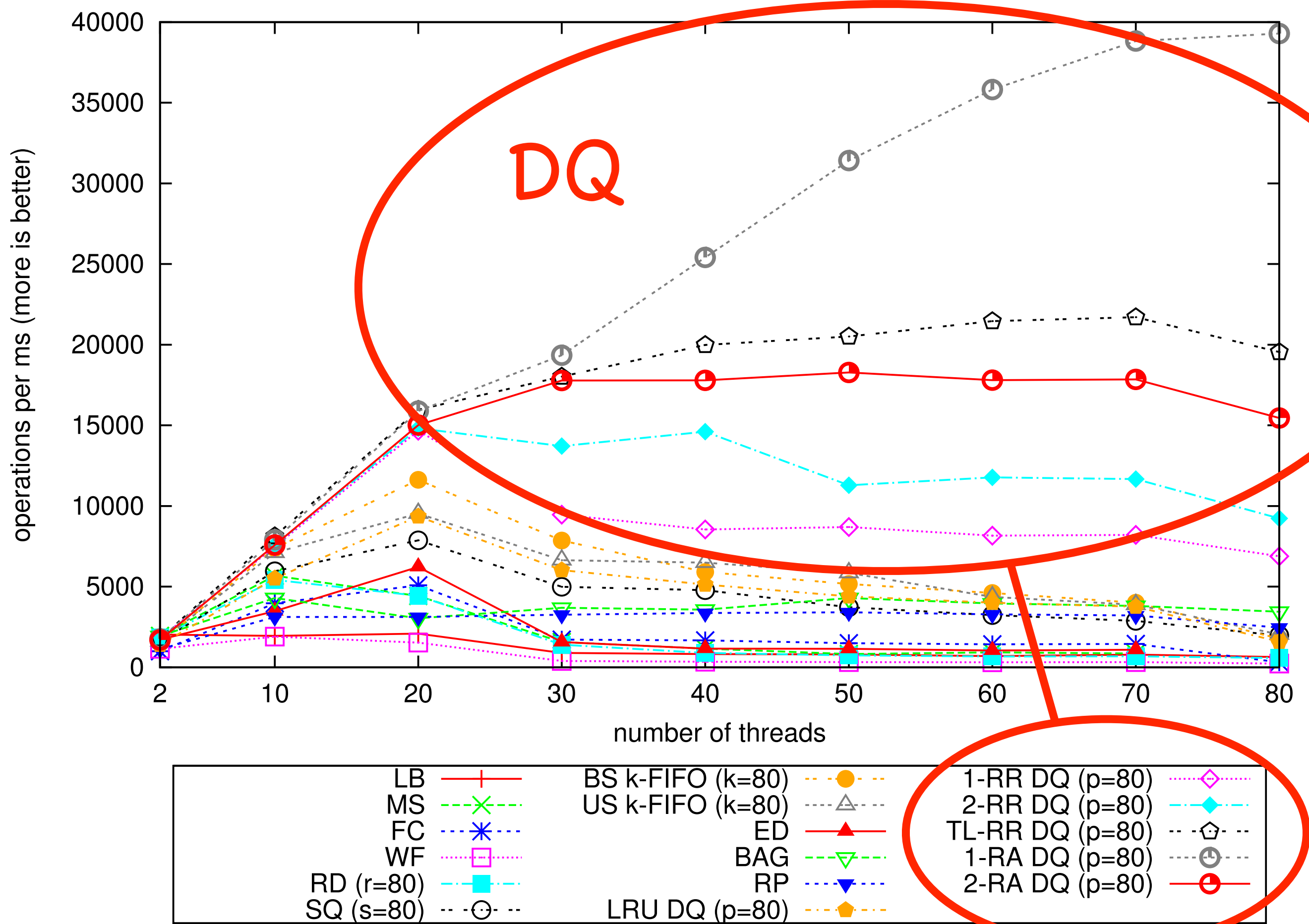
4 processors x 10 cores x  
2 hardware threads =  
80 hardware threads





# Performance & Scalability

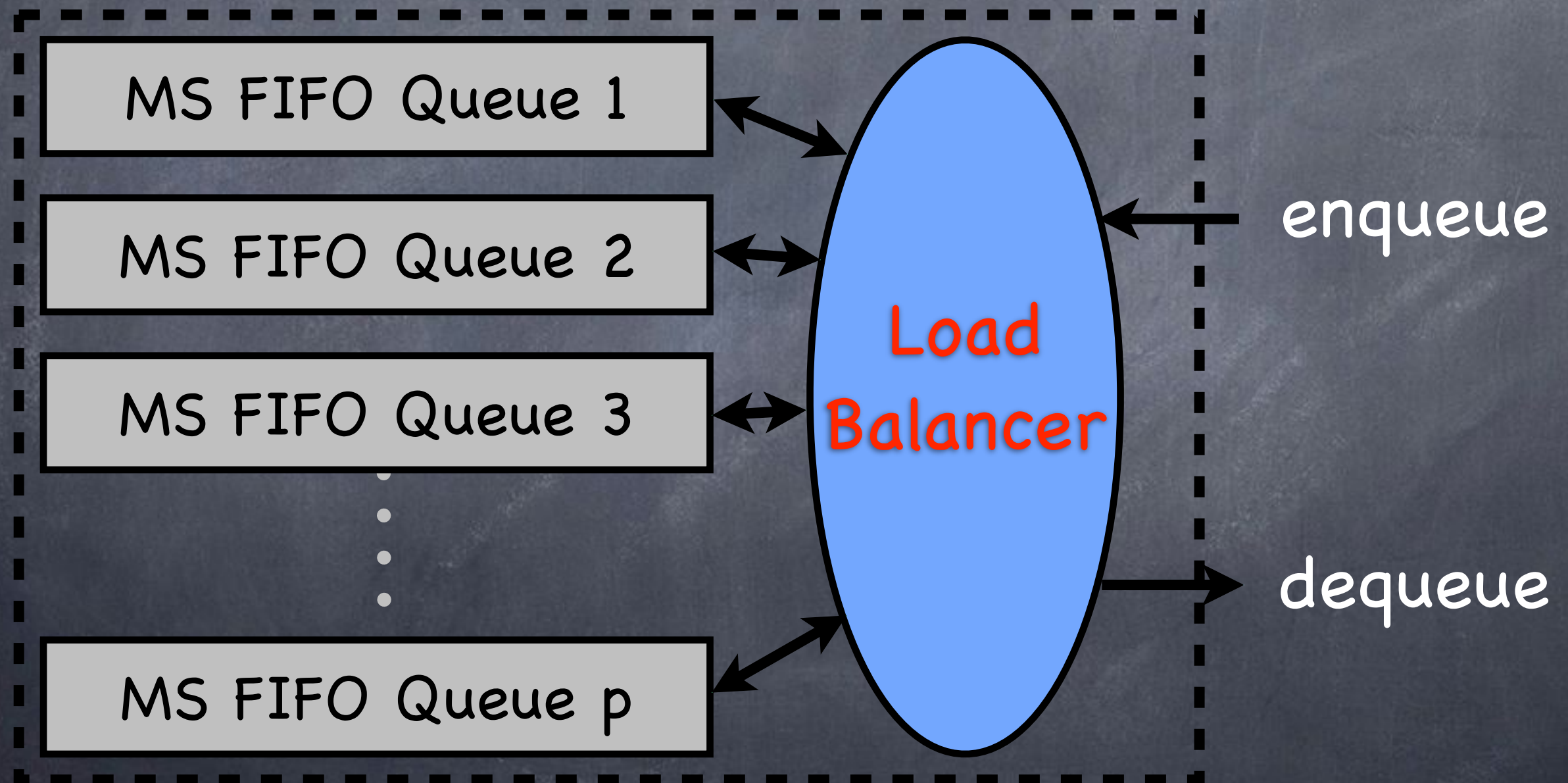




(a) High contention producer-consumer microbenchmark ( $c = 250$ )

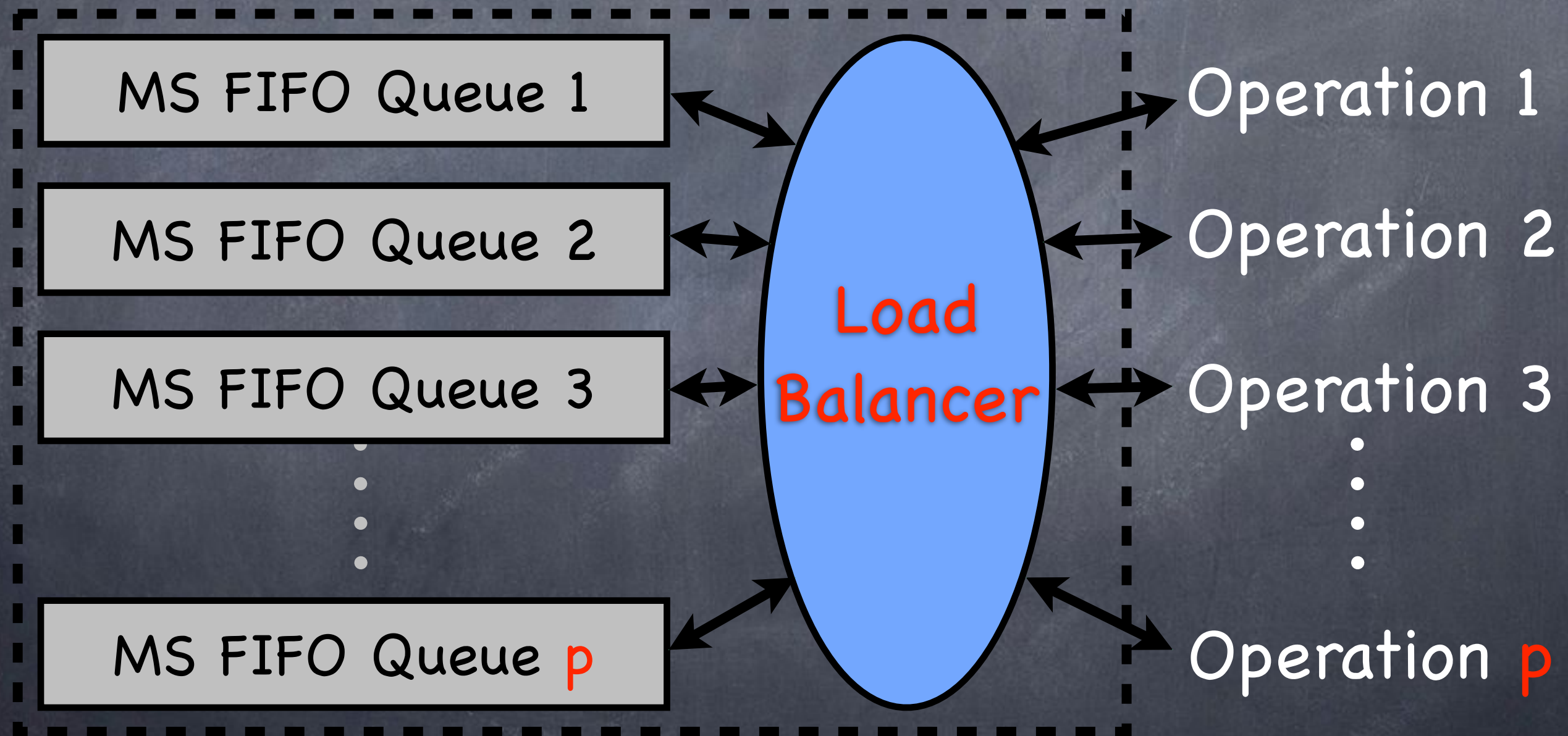
# Distributed Queues

[PODC BA 2011, ICA3PP 2012, Submitted 2013]

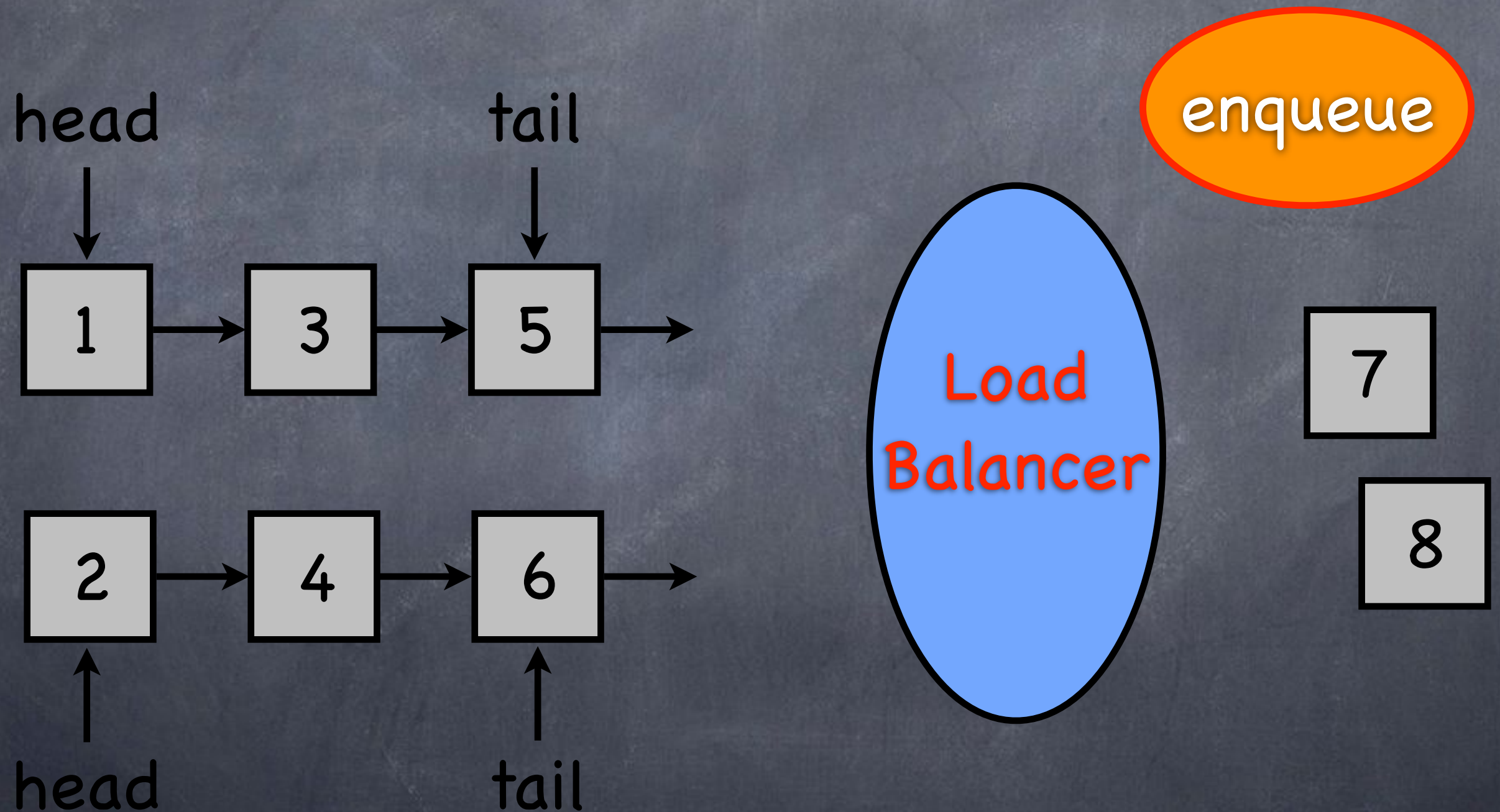




# Up to $p$ Parallel Enqueues and $p$ Parallel Dequeues



# Parallel Access





Emptiness Check?

→

Not Relaxed!

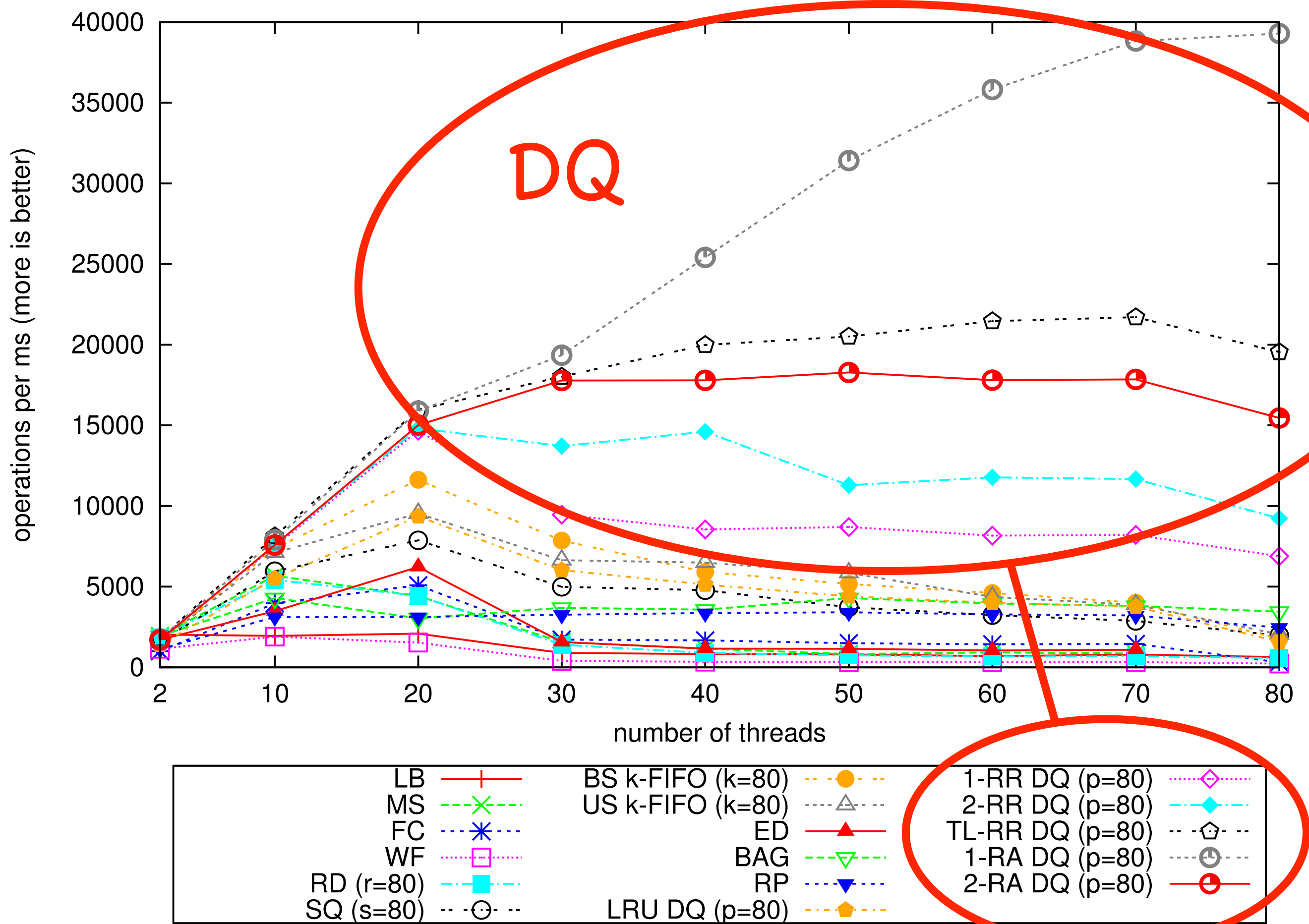
# Listing 1: Lock-free load-balanced distributed queue algorithm

```
1 enqueue(element):
2   index = load_balancer();
3   DQ[index].MS_enqueue(element);
4
5 dequeue():
6   start = load_balancer();
7   while true:
8     for i in 0 to p-1:
9       index = (start + i) % p;
10      element, current_tail = DQ[index].MS_dequeue();
11      1. if element != null:
12         return element;
13      else:
14         tail_old[index] = current_tail;
15      2. for i in 0 to p-1:
16         if get_tail(DQ[i]) != tail_old[i]:
17            start = i;
18            break;
19         if i == p-1:
20            return null;
```

**DQ[p]: array of MS queues**

**tail\_old[p]: array of MS tails**





(a) High contention producer-consumer microbenchmark ( $c = 250$ )

# Semantics

[Related Work]

Our Stuff

## Pools

### k-FIFO ( $k \geq 0$ )

1-RA DQ

2-RA DQ

TL-RR DQ

2-RR DQ

1-RR DQ

ED

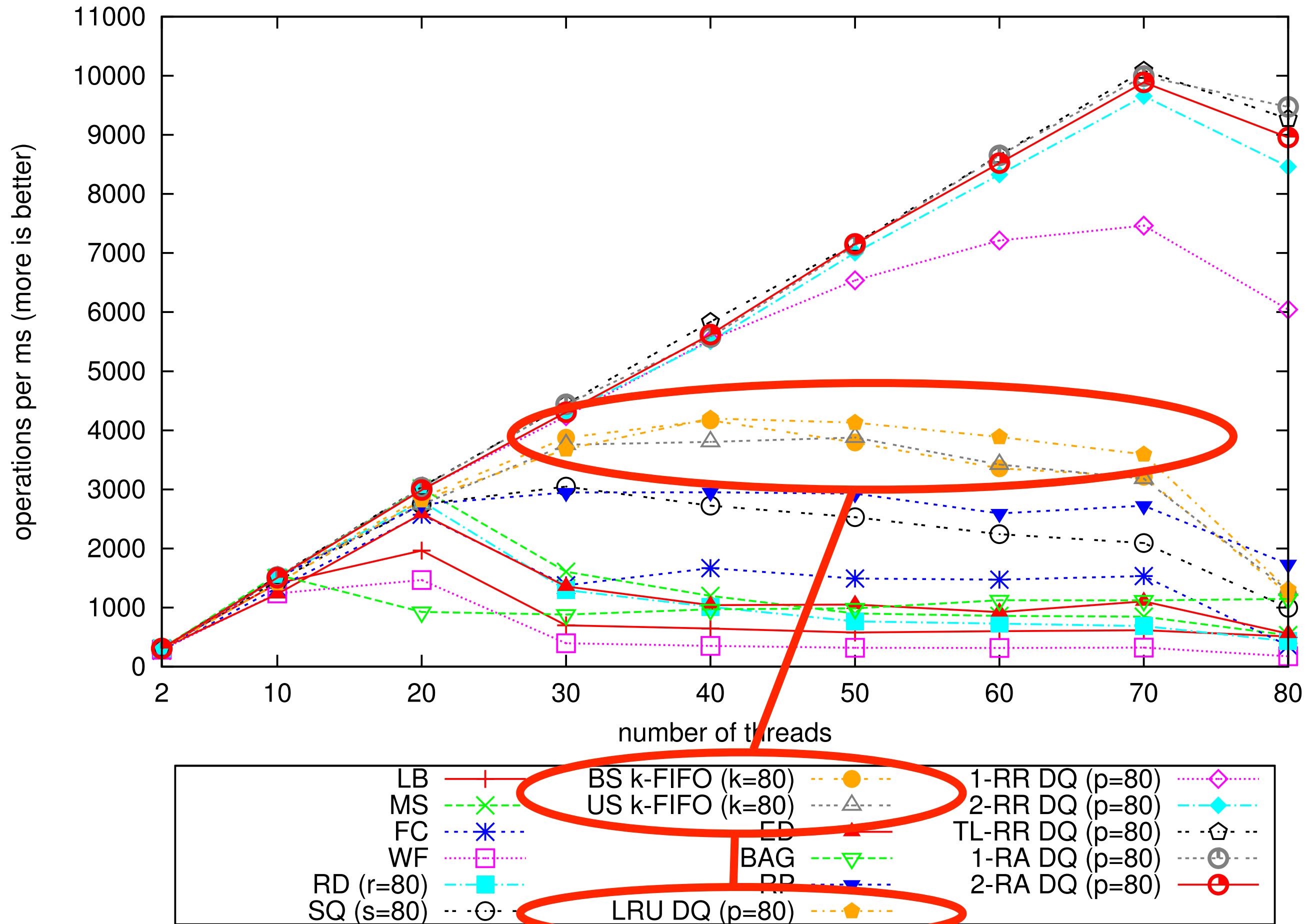
BAG

RP

[Sundell et al.'11]

[Afek et al.'11,'10]





(b) Low contention producer-consumer microbenchmark ( $c = 2000$ )

## Listing 2: Lock-free LRU distributed queue algorithm

```
1 enqueue(element):
2   start = random();
3   while true:
4     aba_index, aba_count = lowest_aba_tail(start);
5     for i in 0 to p-1:
6       index = (aba_index + i) % p;
7       current_tail = get_tail(DQ[index]);
8       if current_tail.aba == aba_count &&
9         DQ[index].try_MS_enqueue(element, current_tail):
10        return;
11
12 dequeue():
13   start = random();
14   while true:
15     aba_index, aba_count = lowest_aba_head(start);
16     check_emptiness = true;
17     clear(empty_queue);
18     for i in 0 to p-1:
19       index = (aba_index + i) % p;
20       current_head = get_head(DQ[index]);
21       if current_head.aba == aba_count:
22         element, current_tail =
23           DQ[index].try_MS_dequeue(current_head);
24         if element == FAILED:
25           check_emptiness = false;
26         else if element == null:
27           tail_old[index] = current_tail;
28           empty_queue[index] = true;
29         else:
30           return element;
31
32   if check_emptiness && there_is_any(empty_queue):
33     for i in 0 to p-1:
34       if empty_queue[i] &&
35         (get_tail(DQ[i]) != tail_old[i]):
36         start = i;
37       break;
38   if i == p-1:
39     return null;
```

LRU DQ:

max difference of  
tail/head

ABA counters  
is one!

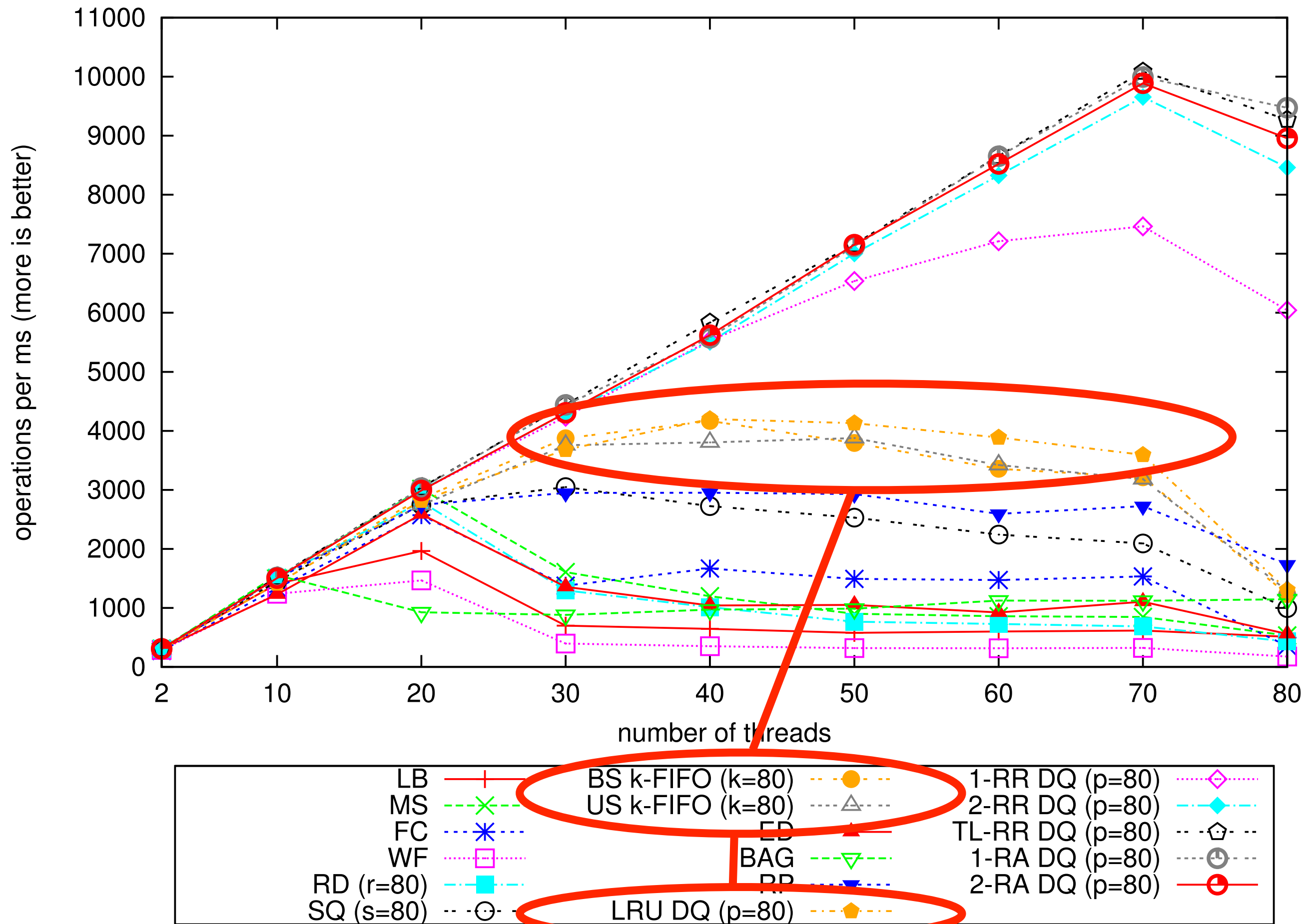
->

there are two  
partitions of MS queues  
with lowest/highest  
ABA counters

->

enqueue/dequeue  
@one\_of\_lowest

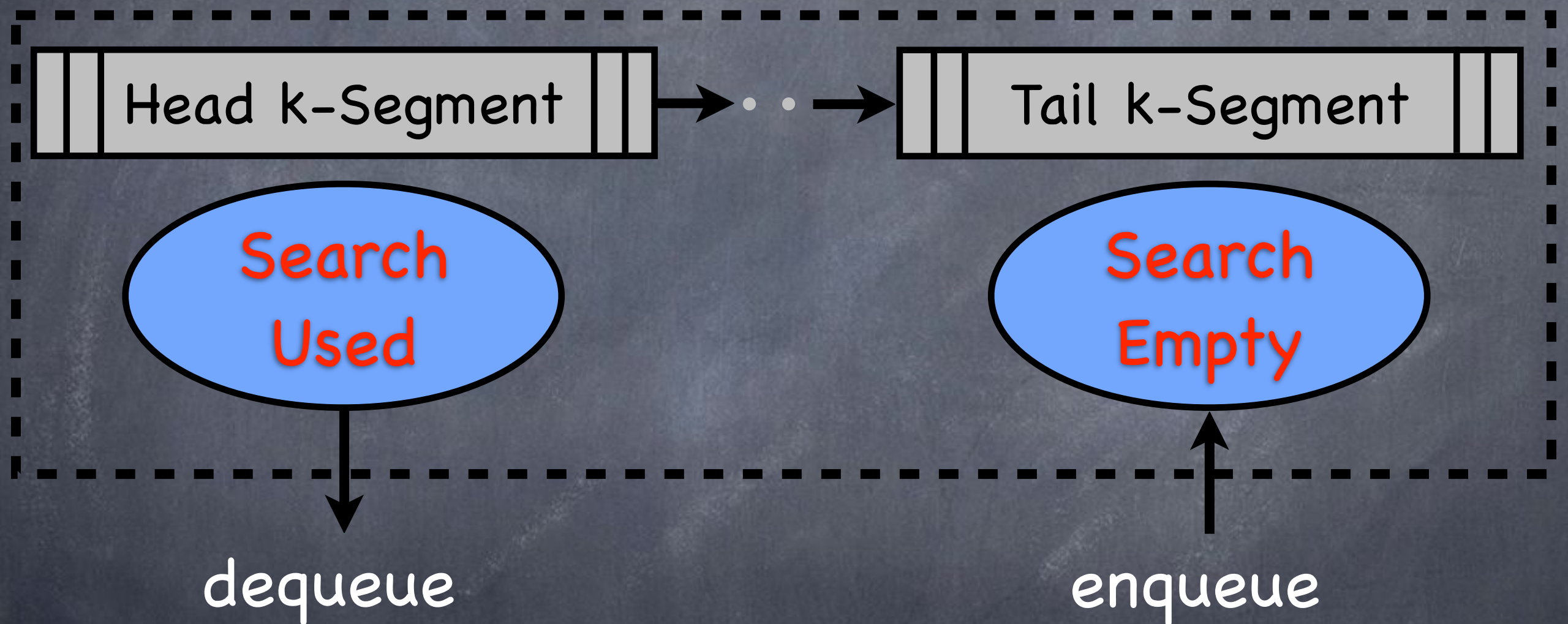




(b) Low contention producer-consumer microbenchmark ( $c = 2000$ )

# Segmented Queues (SQ)

[Afek, Korland, Yanovsky 2010]



→ BS, US k-FIFO Queues

[Lippautz, Payer 2012]



Emptiness Check?

→

Not Relaxed!

# enqueue

```
1 bool enqueue(item):
2   while true:
3       tail_old = get_tail();
4       head_old = get_head();
5       item_old, index = find_empty_slot(tail_old, k, TESTS);
6       if tail_old == get_tail():
7           if item_old.value == EMPTY:
8               item_new = atomic_value(item, item_old.counter + 1);
9               if CAS(&tail_old[index], item_old, item_new):
10                  if committed(tail_old, item_new, index):
11                      return true;
12           else:
13               if queue_full(head_old, tail_old):
14                   if segment_not_empty(head_old, k) && head == get_head():
15                       return false;
16                   advance_head(head_old, k);
17               advance_tail(tail_old, k);
18
19 bool committed(tail_old, item_new, index):
20   if tail_old[index] != item_new:
21       return true;
22   head_current = get_head();
23   tail_current = get_tail();
24   item_empty = atomic_value(EMPTY, item_new.counter + 1);
25   if in_queue_after_head(tail_old, tail_current, head_current):
26       return true;
27   else if not_in_queue(tail_old, tail_current, head_current):
28       if !CAS(&tail_old[index], item_new, item_empty):
29           return true;
30   else: //in queue at head
31       head_new = atomic_value(head_current.value, head_current.counter + 1);
32       if CAS(&head, head_current, head_new):
33           return true;
34       if !CAS(&tail_old[index], item_new, item_empty):
35           return true;
36   return false;
```



# dequeue

```
38 item dequeue():
39     while true:
40         tail_old = get_tail();
41         head_old = get_head();
42         item_old, index = find_item(head_old, k);
43         if head_old == head:
44             if item_old.value != EMPTY:
45                 if head_old.value == tail_old.value:
46                     advance_tail(tail_old, k);
47                     item_empty = atomic_value(EMPTY, item_old.counter + 1);
48                     if CAS(&head_old[index], item_old, item_empty):
49                         return item_old.value;
50             else:
51                 if head_old.value == tail_old.value && tail_old == get_tail():
52                     return null;
53                 advance_head(head_old, k);
```

# Semantics

[Related Work]

Our Stuff

[Afek et al.'10]

## Pools

### k-FIFO ( $k \geq 0$ )

1-RA DQ

2-RA DQ

TL-RR DQ

2-RR DQ

1-RR DQ

(SQ)

RD

### FIFO

LB MS

WF FC

LRU DQ

BS, US

ED  
BAG  
RP

configurable k

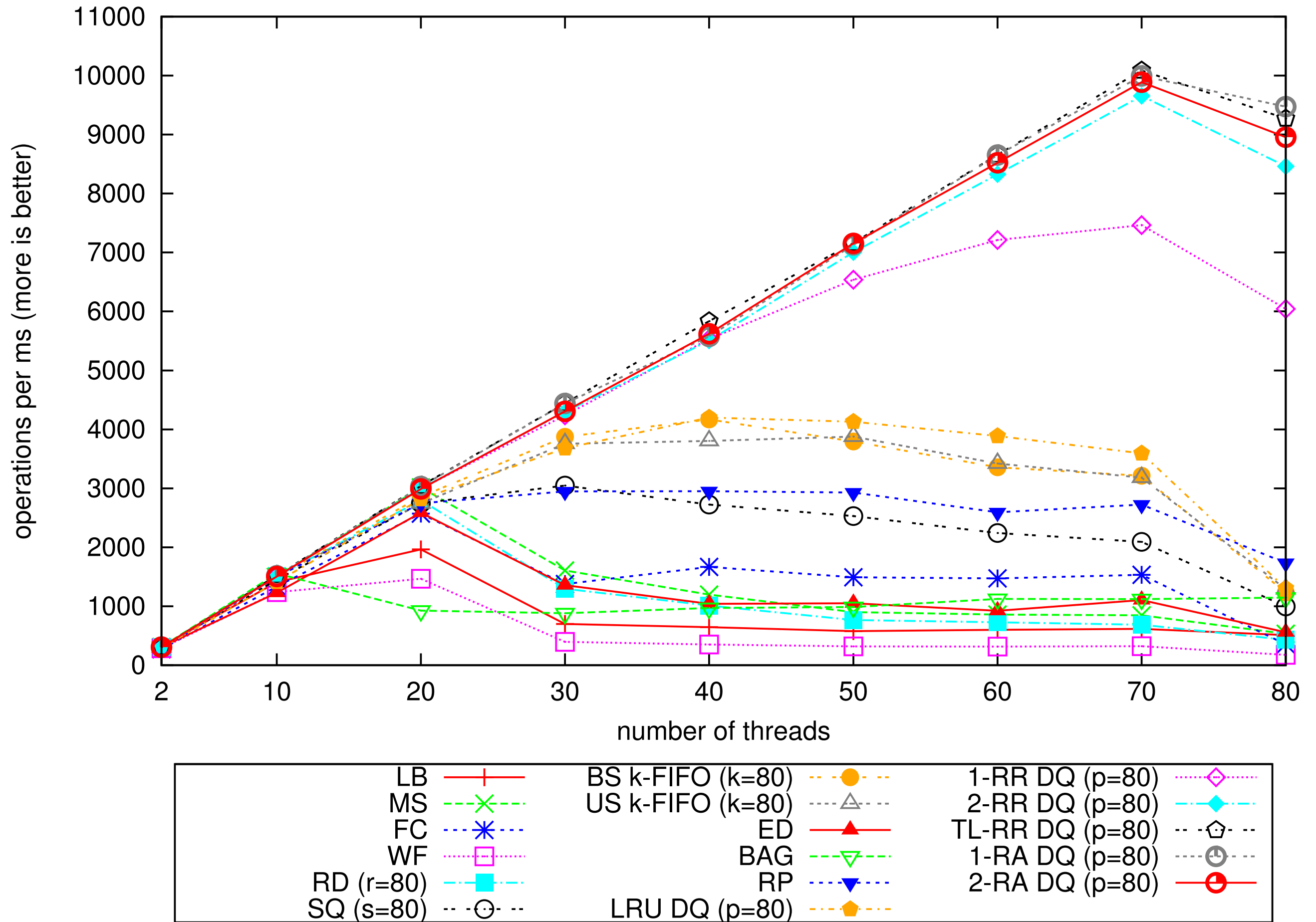
[Sundell et al.'11]

[Afek et al.'11, '10]

[Incze et al.'10]

[Kogan et al.'11]

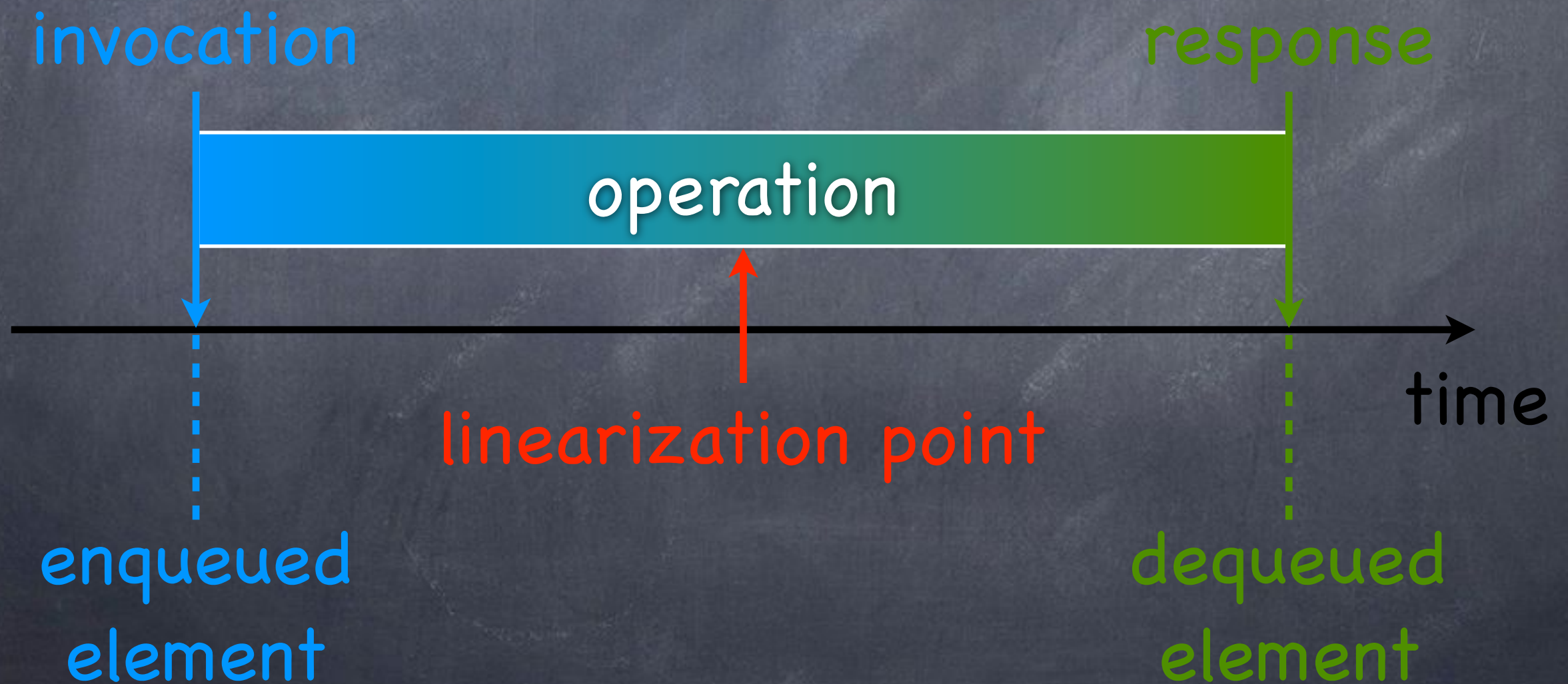




(b) Low contention producer-consumer microbenchmark ( $c = 2000$ )

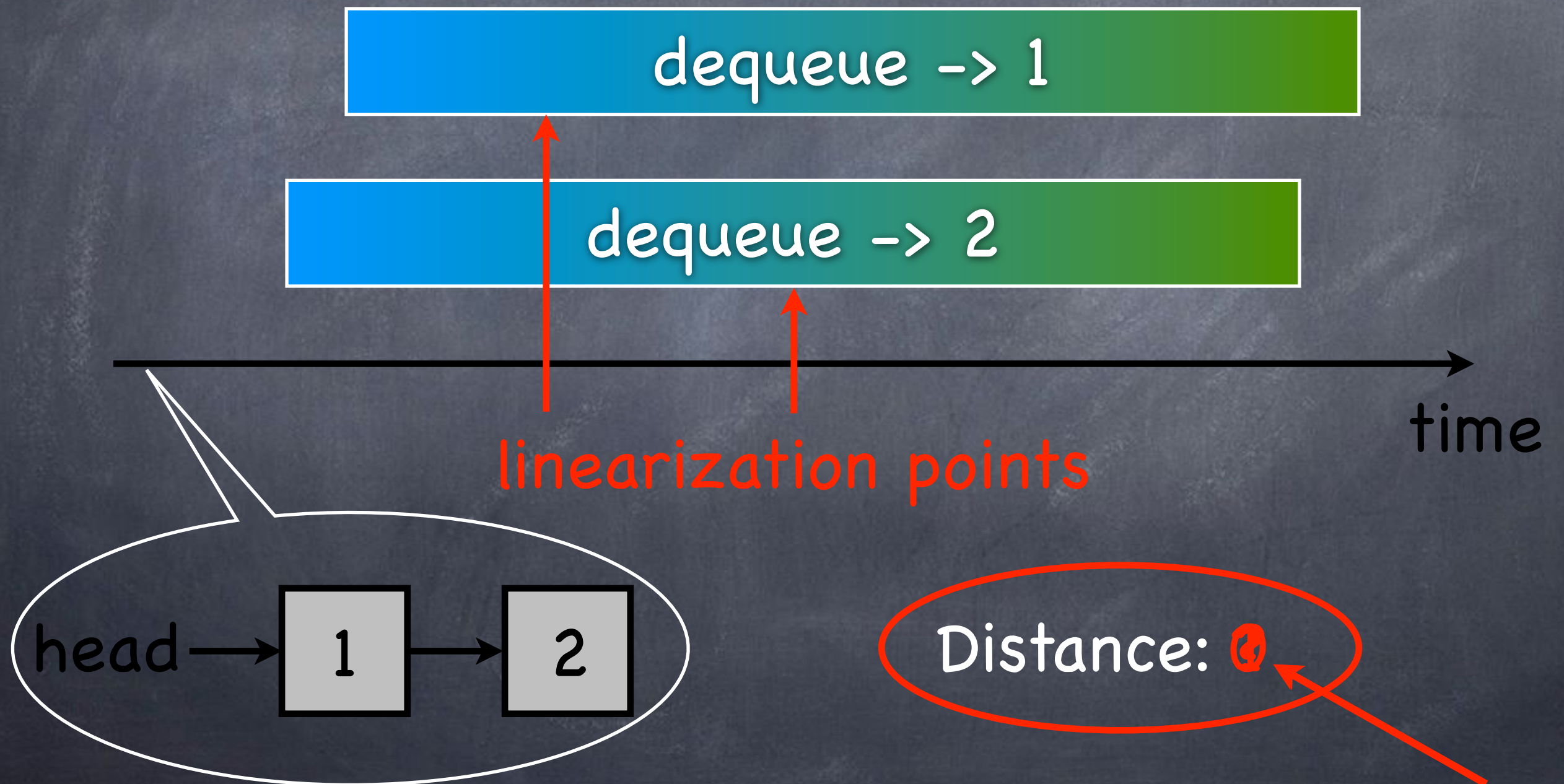
# (Enhanced) Concurrent History

Sequence of Time-stamped Invocation and Response Events as well as Time-stamped Linearization Points (Approximative)





# Measuring “Out-of-Order Distance”





The **Actual-Time Linearization**  
of a concurrent history  
is  
the sequence of its operations  
ordered by  
their **linearization** points



The **Actual-Time Distance**  
of a concurrent history  
is

the **average**  
out-of-order distance  
of

its actual-time linearization



# Actual-Time Distance

measures

re-ordering

due to

semantical relaxation!



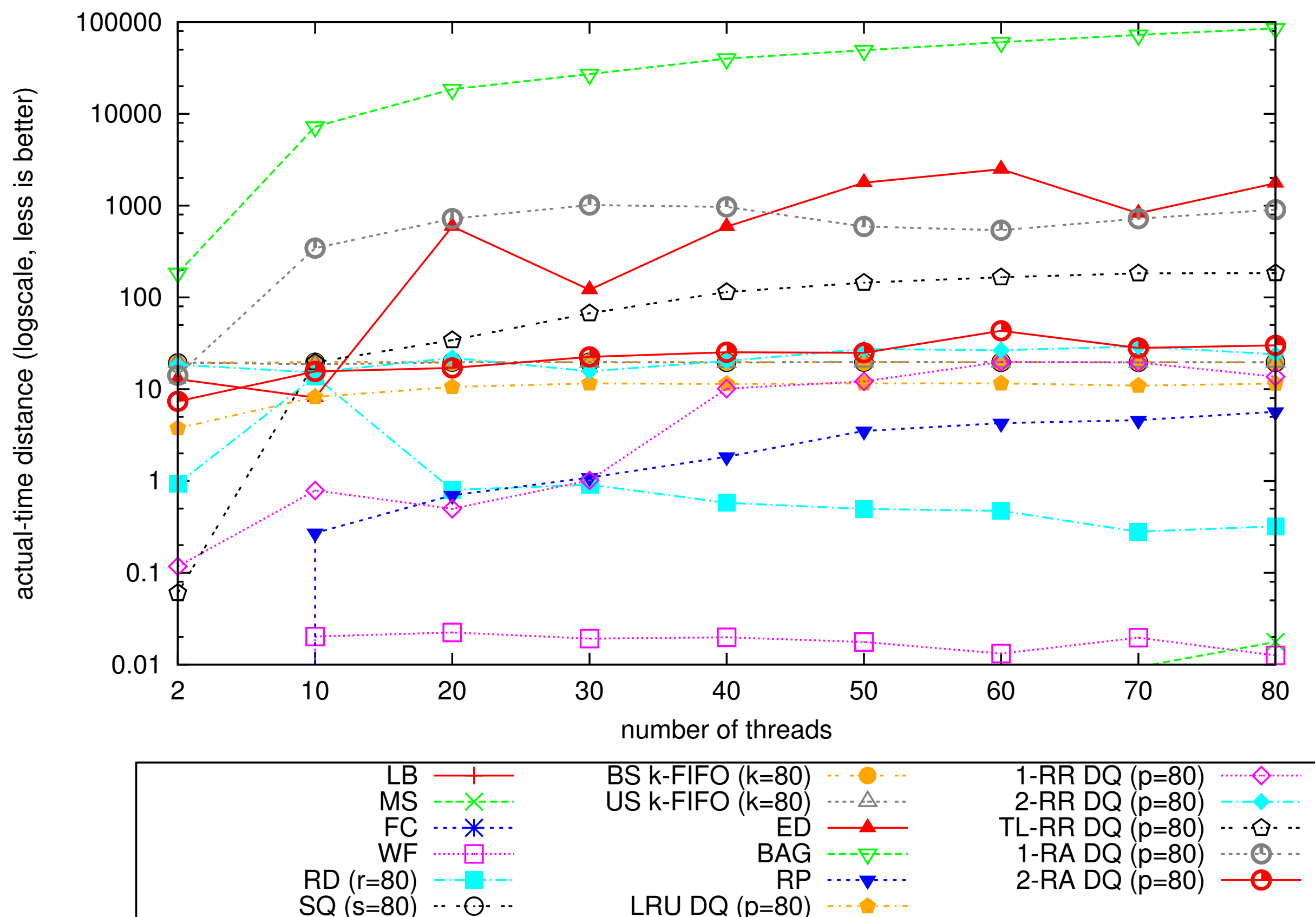
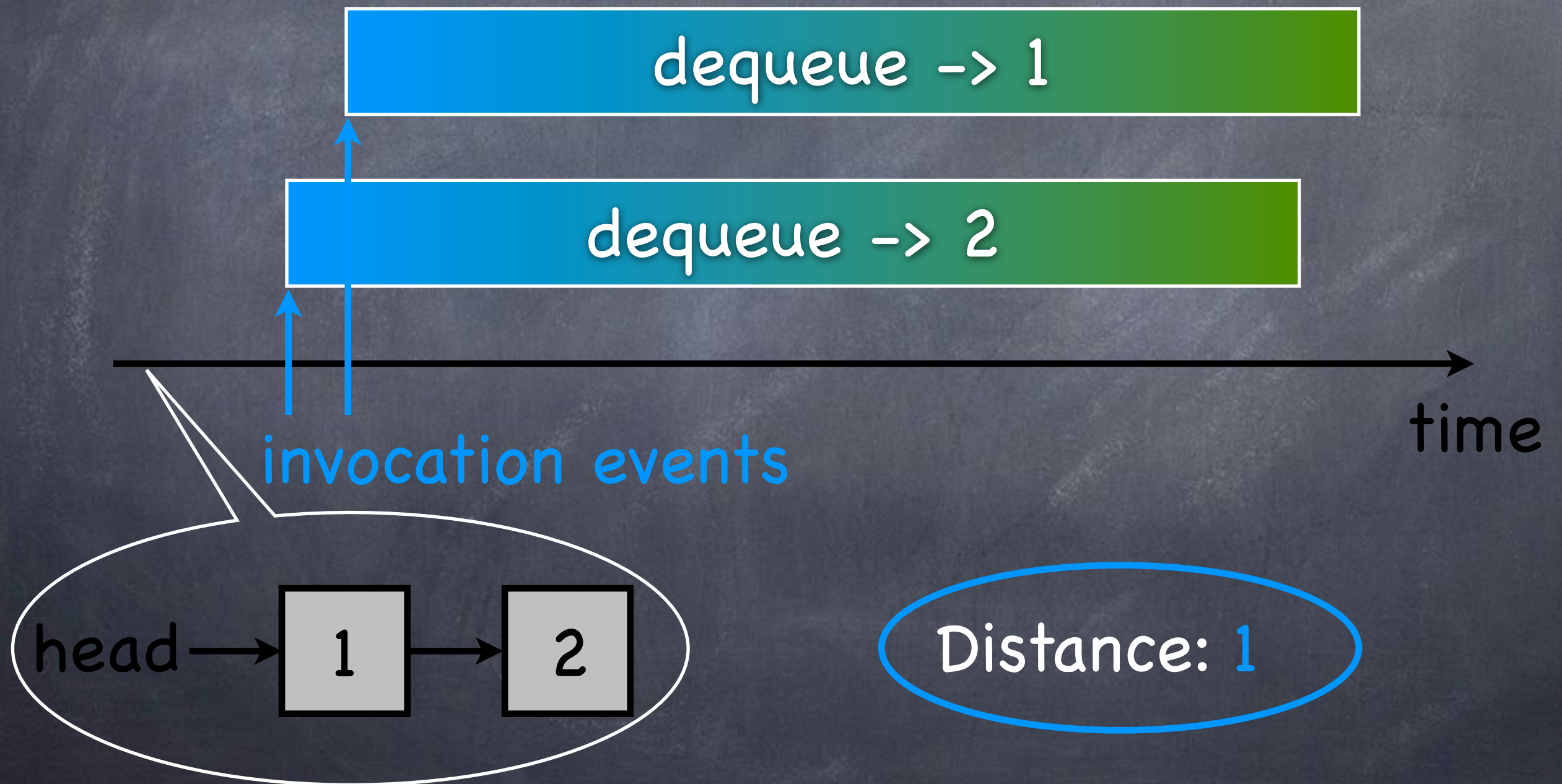


Figure 6: Actual-time distance of the high contention producer-consumer microbenchmark ( $c = 250$ )

# Invocation vs. Linearization





The Zero-Time Linearization  
of a concurrent history  
is  
the sequence of its operations  
ordered by  
their invocation events



The Zero-Time Distance  
of a concurrent history  
is  
the average  
out-of-order distance  
of  
its zero-time linearization



# Zero-Time Distance

measures  
re-ordering  
due to  
semantical relaxation  
and  
linearizability!

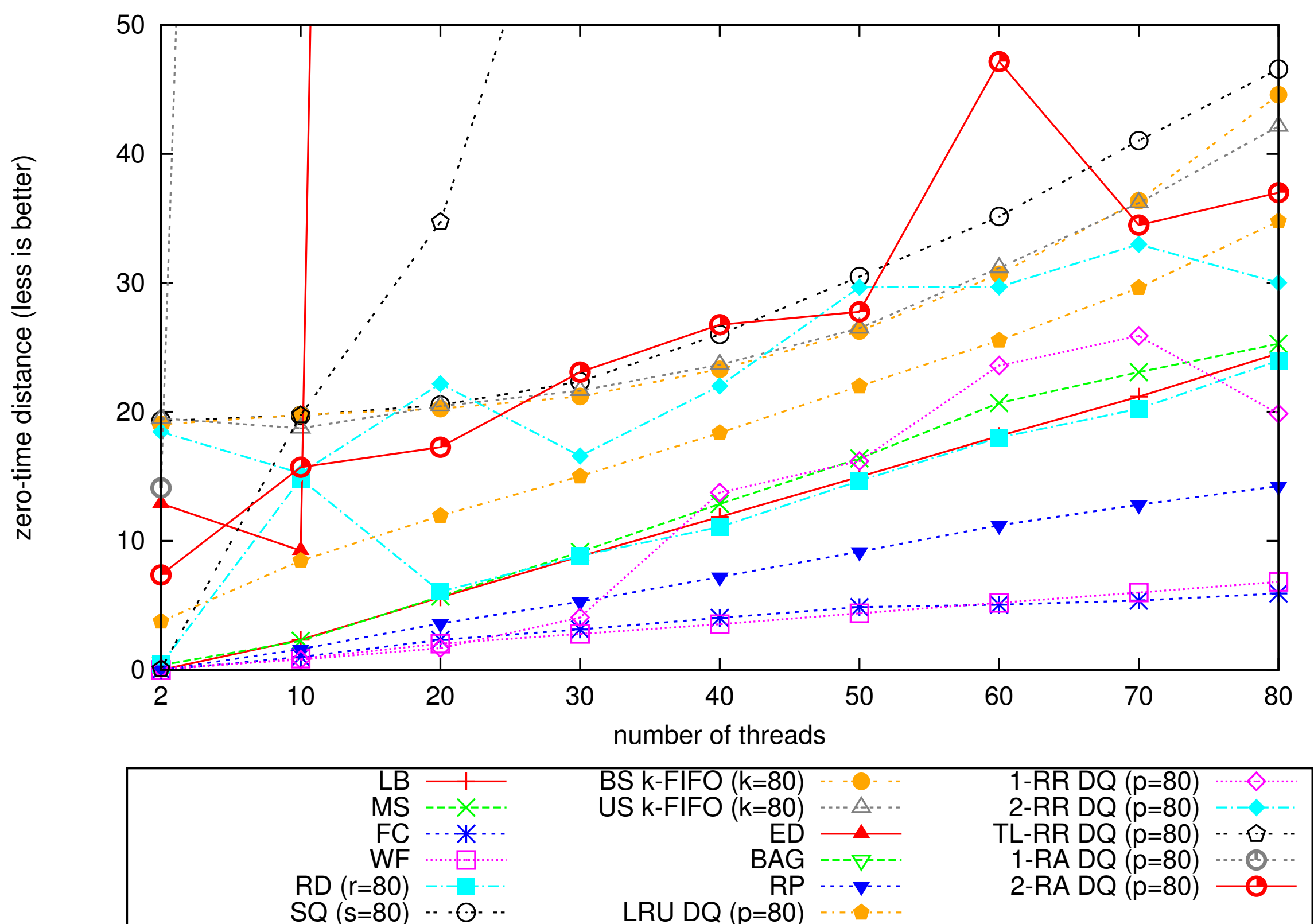


Figure 7: Zero-time distance of the high contention producer-consumer microbenchmark, zoomed-in for better resolution cutting off ED, BAG, TL-RR, and 1-RA ( $c = 250$ )



# Linearization Difference

(difference of zero- and  
actual-time distance)

measures

re-ordering

due to

linearizability!

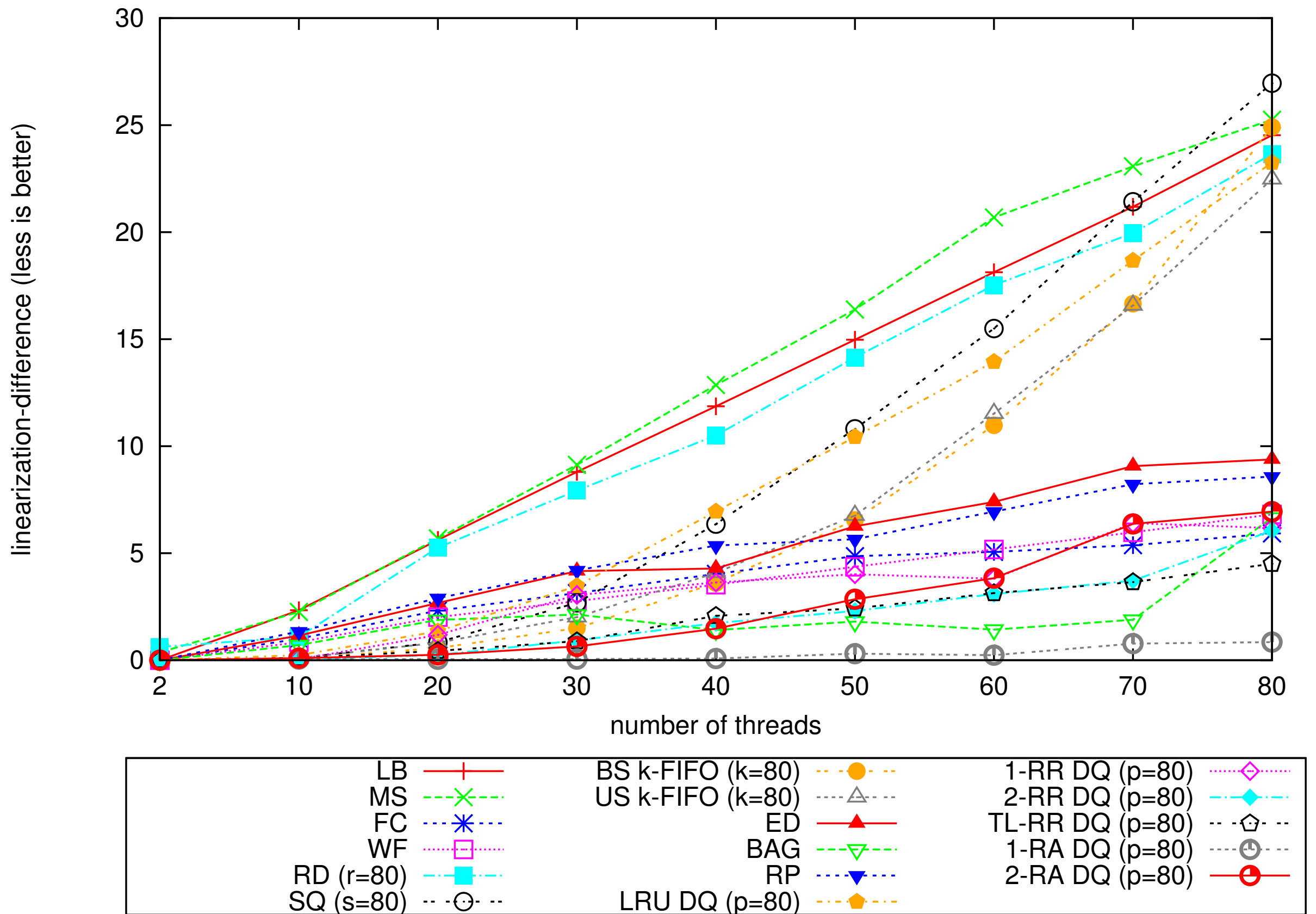


Figure 8: Linearization difference of the high contention producer-consumer microbenchmark ( $c = 250$ )



Thank you

