

MPI Process Synchronization in Space and Time

J. Schuchart¹, S. Hunold², G. Bosilca¹

EuroMPI'23

September 11-13, 2023

Bristol, UK



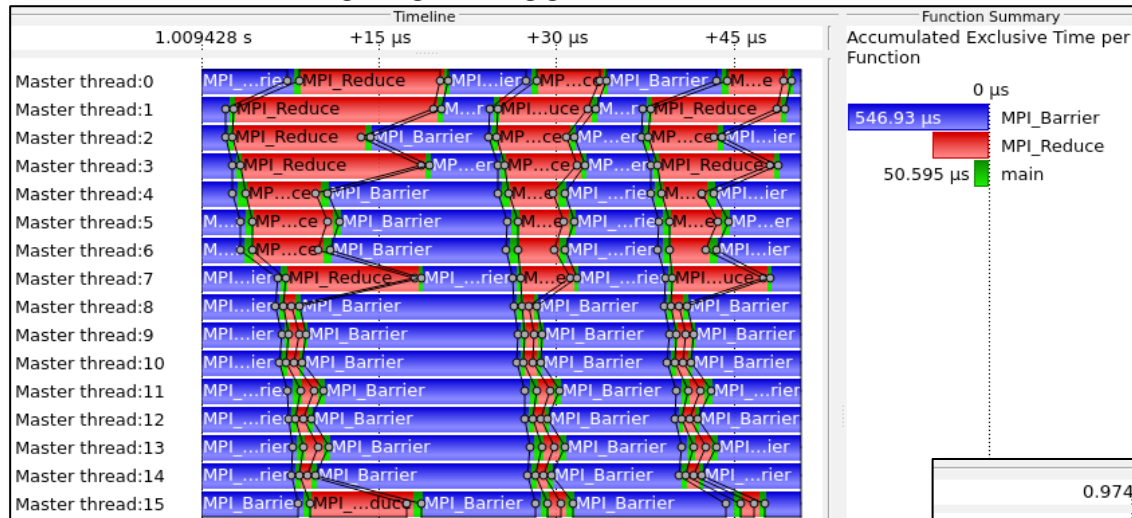
THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

¹ Innovative Computing Laboratory, University of Tennessee, Knoxville, USA

² TU Wien, Vienna, Austria

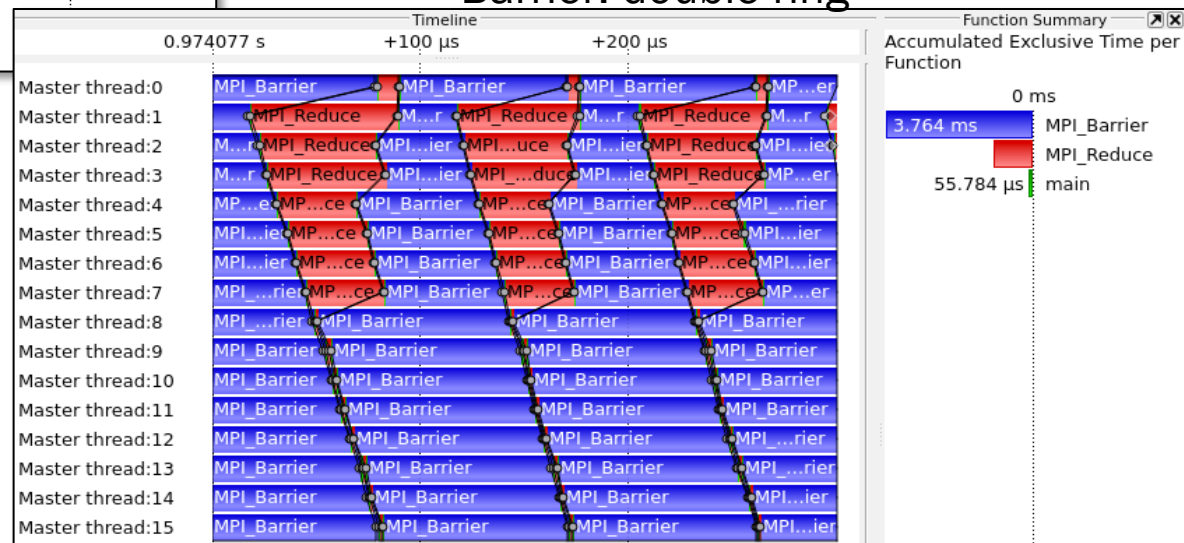
A Tale of Two Barriers

Barrier: linear



Reported Latency: 2.78us

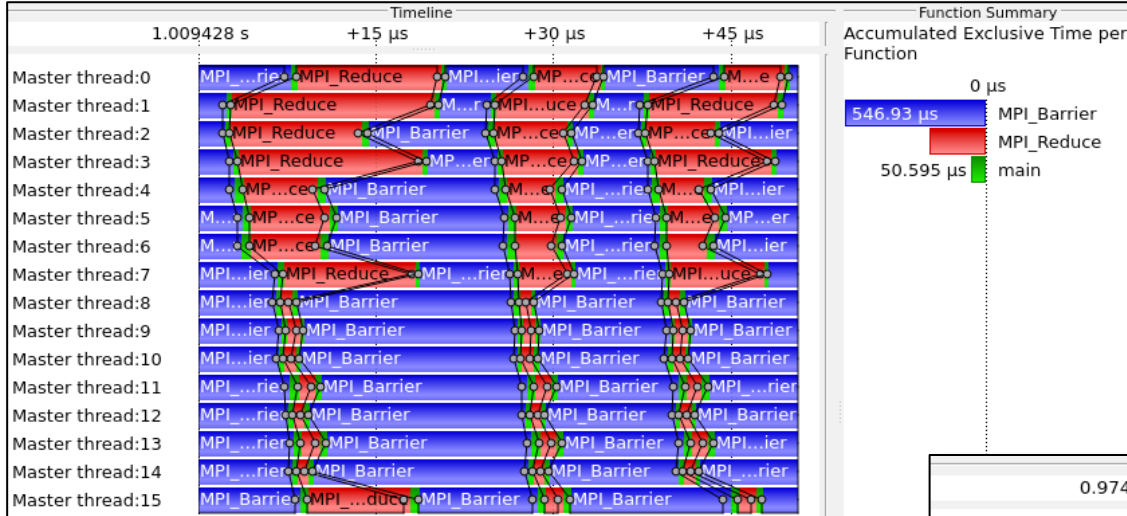
Barrier: double-ring



Reported Latency: 9.98us

A Tale of Two Barriers

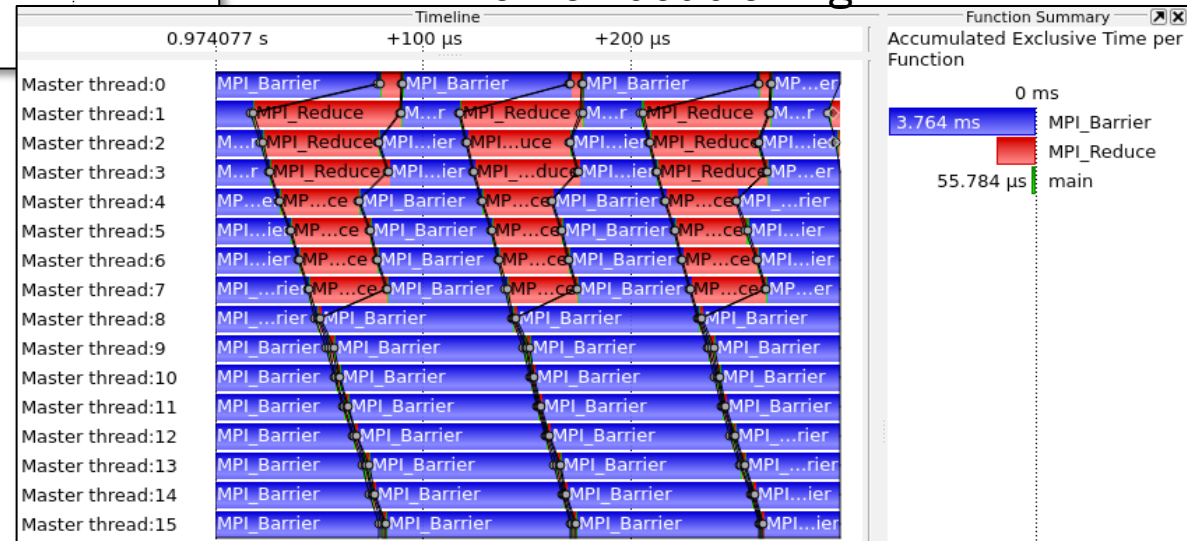
Barrier: linear



Reported Latency: 2.78us

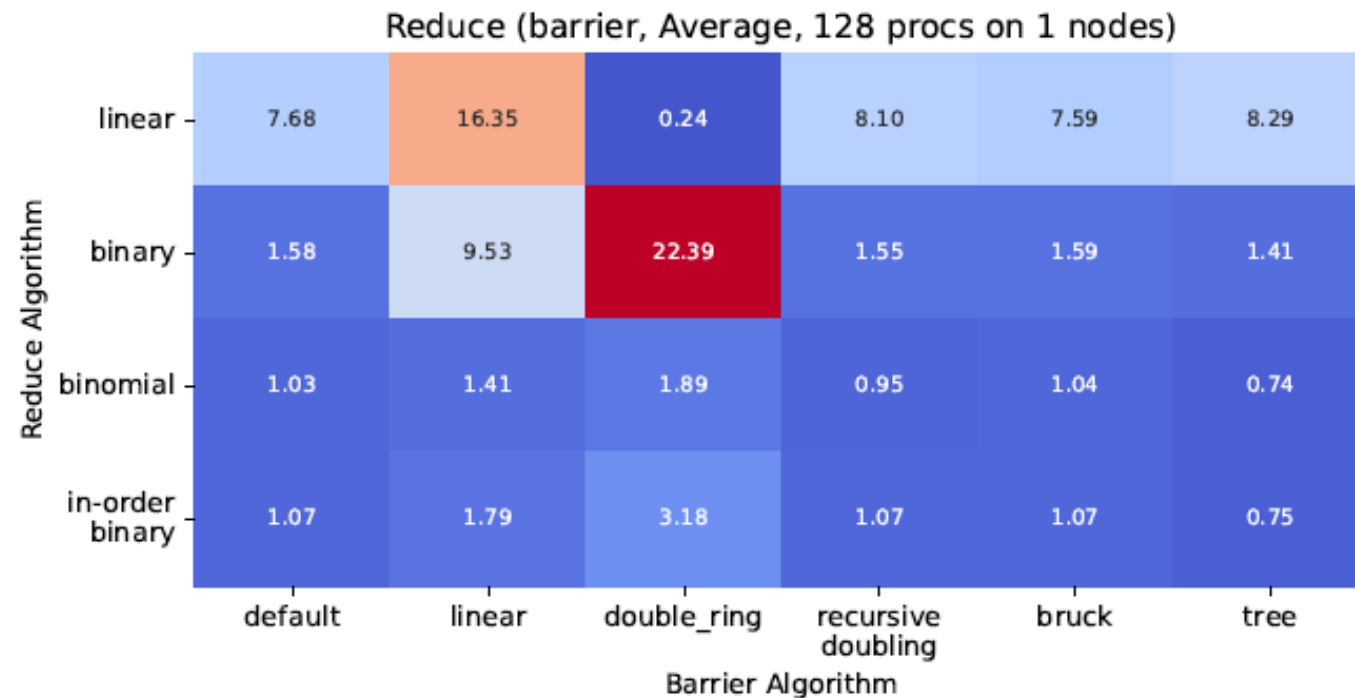
Arbitrary
Arrival
Pattern

Barrier: double-ring

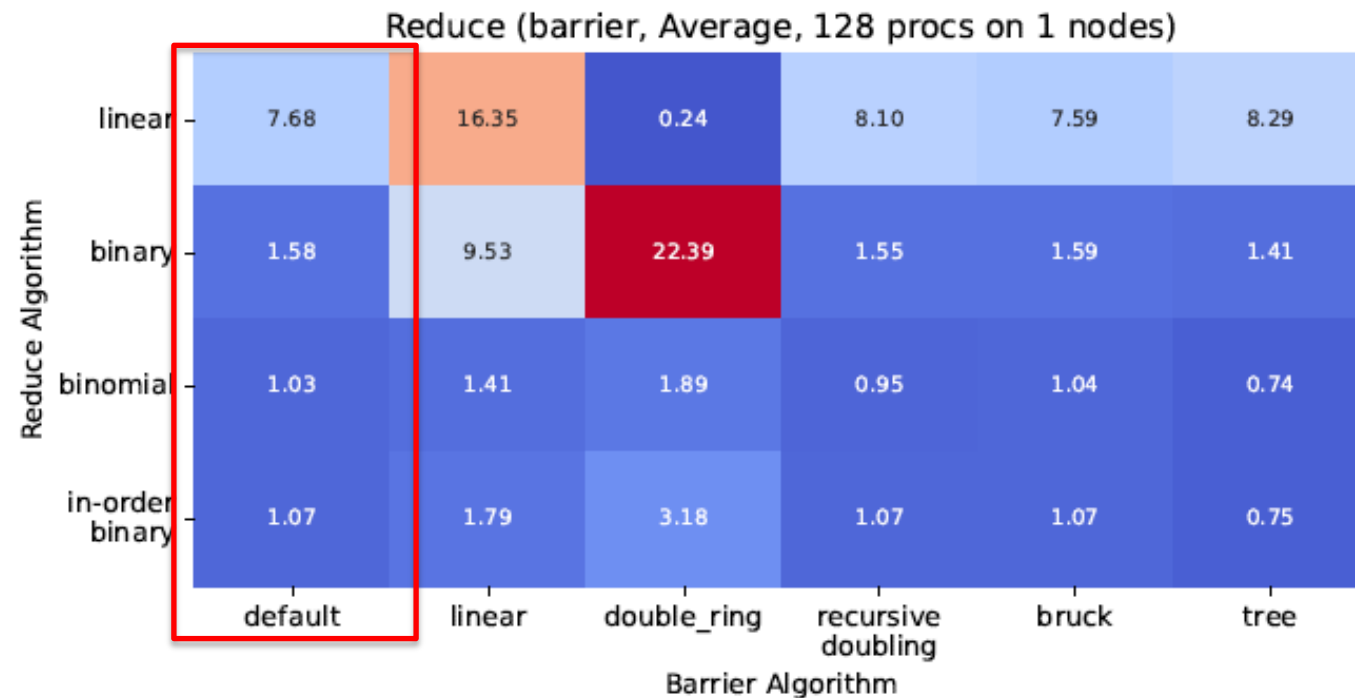


Reported Latency: 9.98us

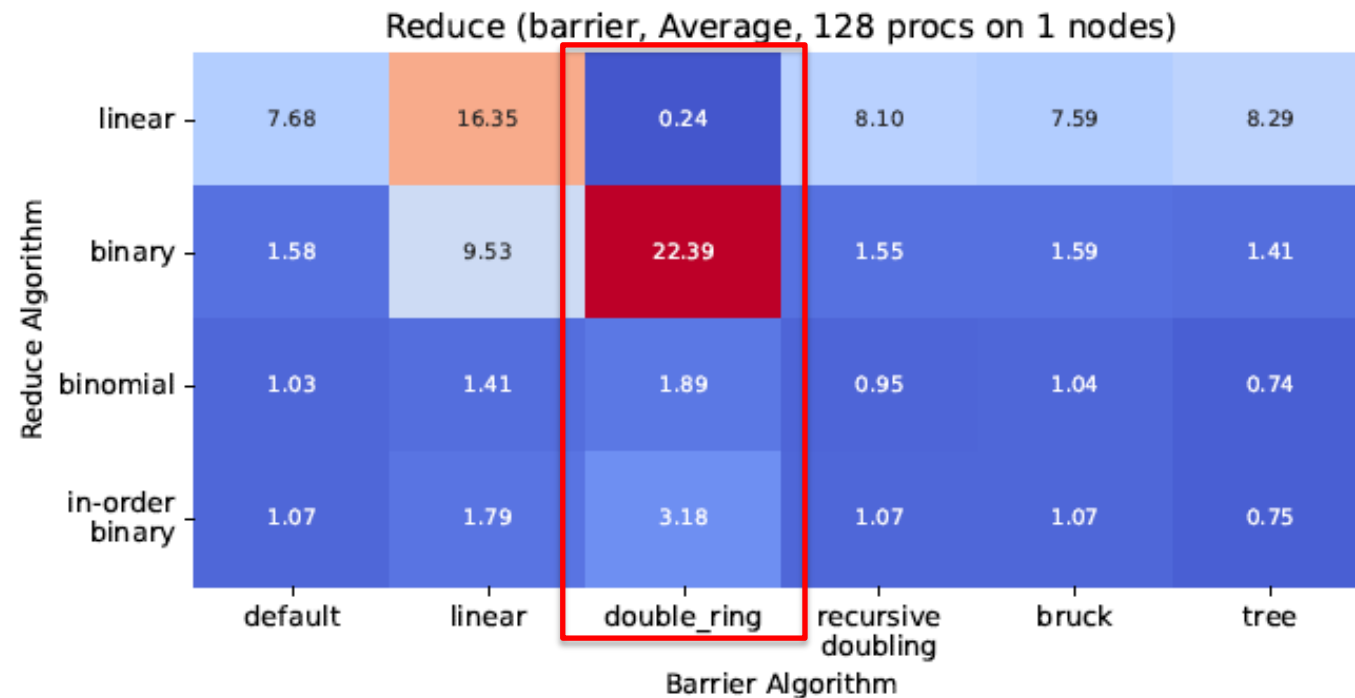
Choice of Barrier Algorithms Matters



Choice of Barrier Algorithms Matters



Choice of Barrier Algorithms Matters

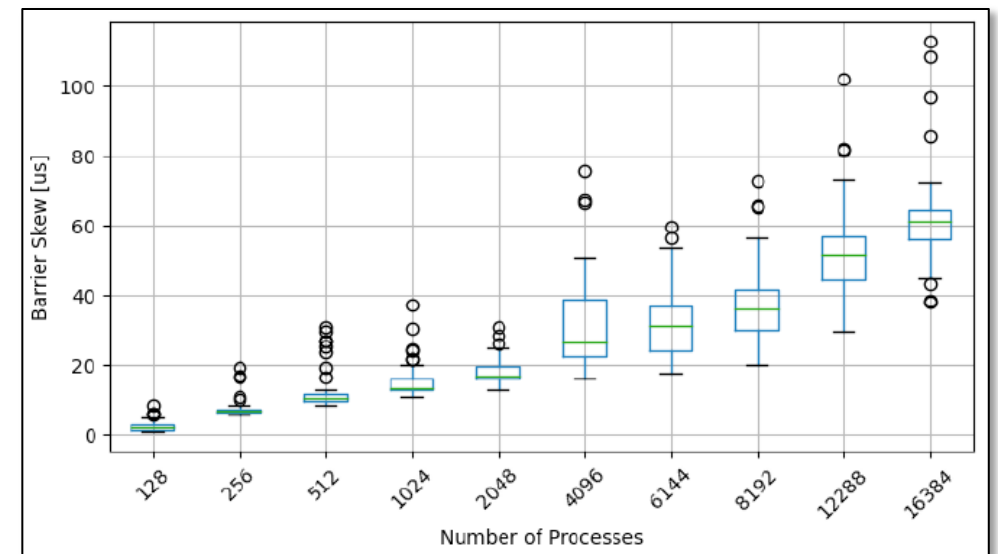
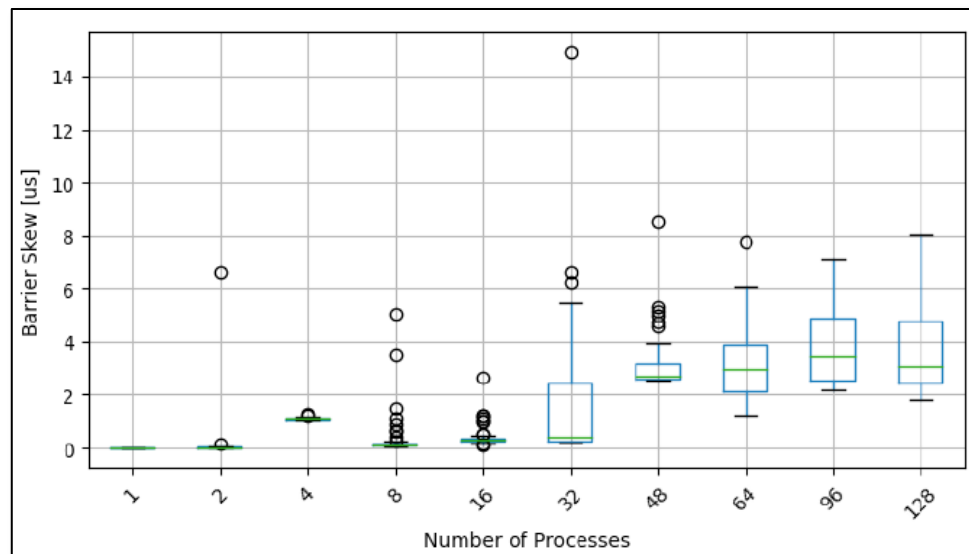


We have been
benchmarking collective
operations wrong for
decades.

Barriers only Synchronize in Space

- Barrier Skew:
difference between minimum and maximum barrier time
- Barriers do not guarantee any time synchronization
 - But that is what we want for benchmarks!

$$\sigma = \max_{t_0 \dots t_p} - \min_{t_0 \dots t_p}$$



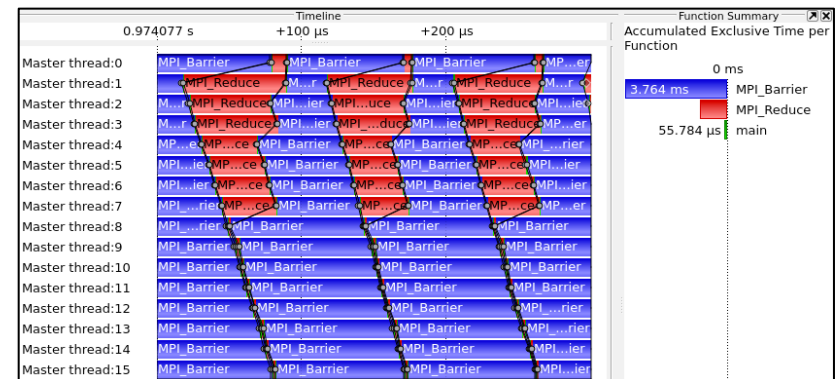
What we care about

What we care about:

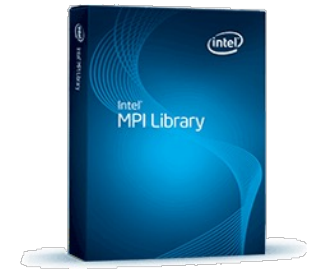
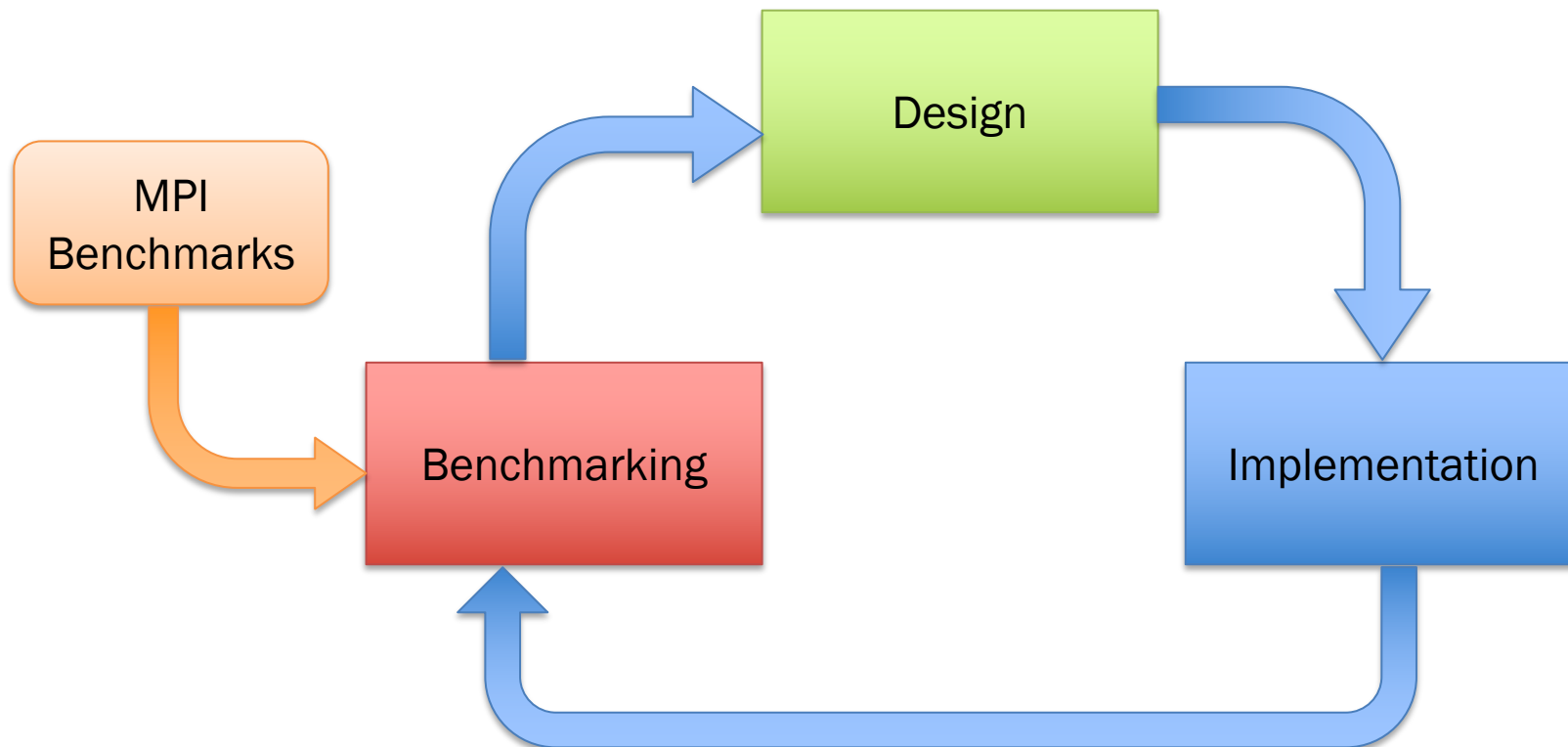
- Latency
- Latency hiding
- Bandwidth
- Concurrency
- Resource usage
- Hardware features
- Correctness

What we *don't* care about:

- Impact of barrier algorithms
(unless we're benchmarking barrier latency)



Why should we care?



MVAPICH

Why don't we care?

- We know this is bad, but all common benchmarks do it anyway
- Proper implementation is not trivial

LLNL mpiBench

```
__TIME_START__;  
for (i = 0; i < p->iter; i++) {  
    MPI_Reduce(sbuffer, rbuffer, p->count,  
              p->type, p->reduceop, p->root,  
              p->comm);  
    __BAR__(p->comm);  
}  
__TIME_END__;
```

Intel MPI Benchmarks

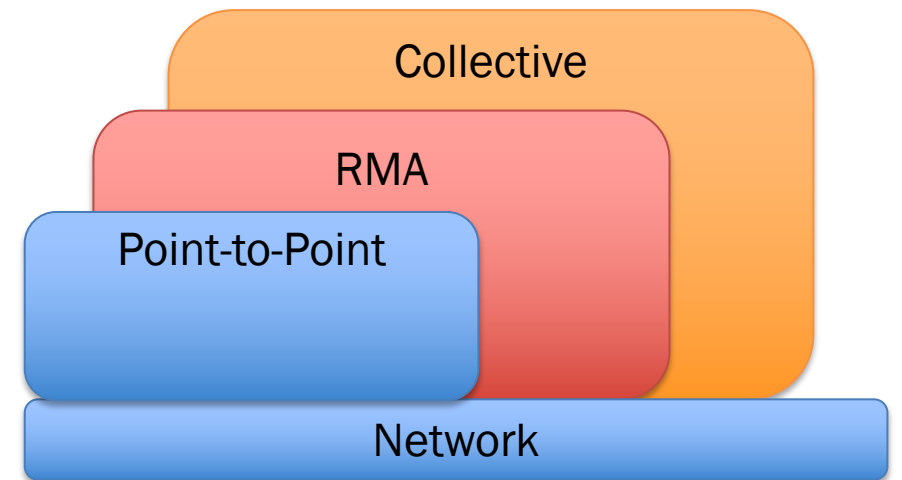
```
for (i = 0; i < ITERATIONS->n_sample; i++) {  
    t1 = MPI_Wtime();  
    MPI_ERRHAND(MPI_Reduce((char*)c_info->s_buffer + i % ITERATIONS->s_cache_iter * ITERATIONS->s_offs,  
                          (char*)c_info->r_buffer + i % ITERATIONS->r_cache_iter * ITERATIONS->r_offs,  
                          s_num,  
                          c_info->red_data_type, c_info->op_type,  
                          root,  
                          c_info->communicator));  
  
    t2 = MPI_Wtime();  
    *time += (t2 - t1);  
  
    /* CHANGE THE ROOT NODE */  
    root = (root + c_info->root_shift) % c_info->num_procs;  
    IMB_do_n_barriers(c_info->communicator, c_info->sync);  
}
```

OSU Micro-Benchmarks

```
MPI_CHECK(MPI_Barrier(MPI_COMM_WORLD));  
  
t_start = MPI_Wtime();  
  
MPI_CHECK(MPI_Reduce(sendbuf, recvbuf, num_elements,  
                    omb_curr_datatype, MPI_SUM, 0,  
                    MPI_COMM_WORLD));  
  
t_stop = MPI_Wtime();
```

How can we care?

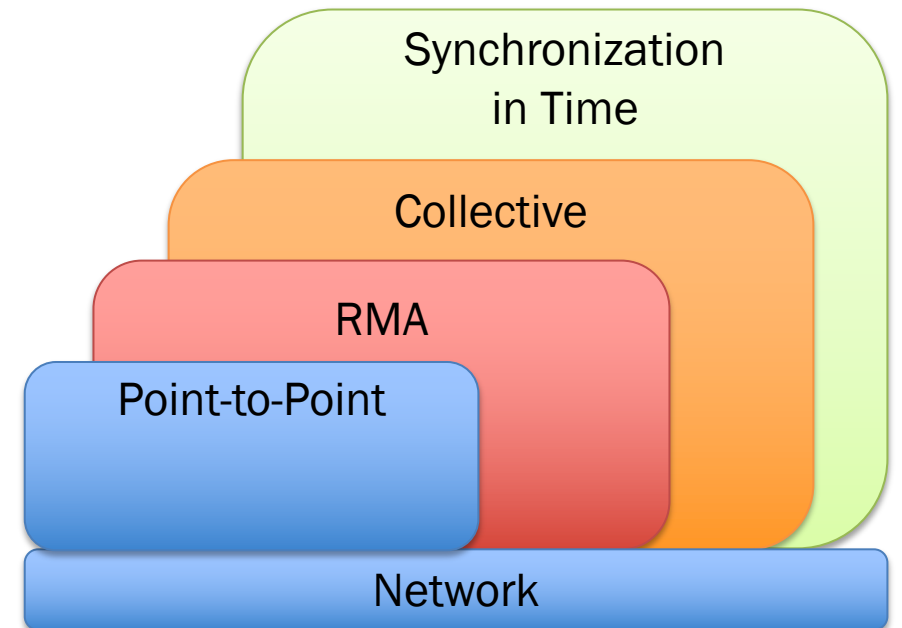
- MPI provides abstractions on many levels
- Abstract away the hard stuff
 - No user should have to implement reductions
 - Or broadcast
 - Or alltoall
 - Or message matching



How can we care?

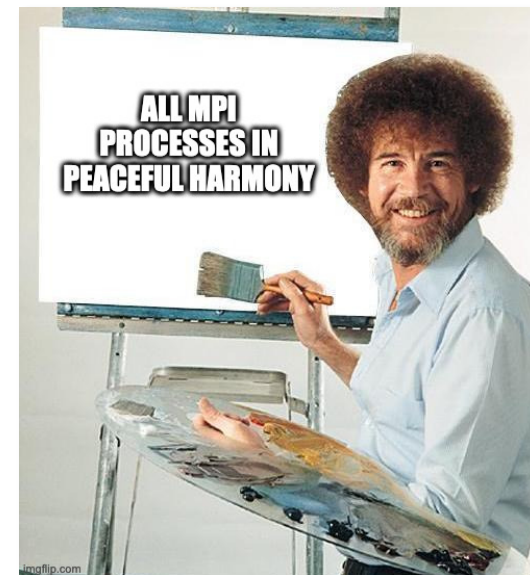
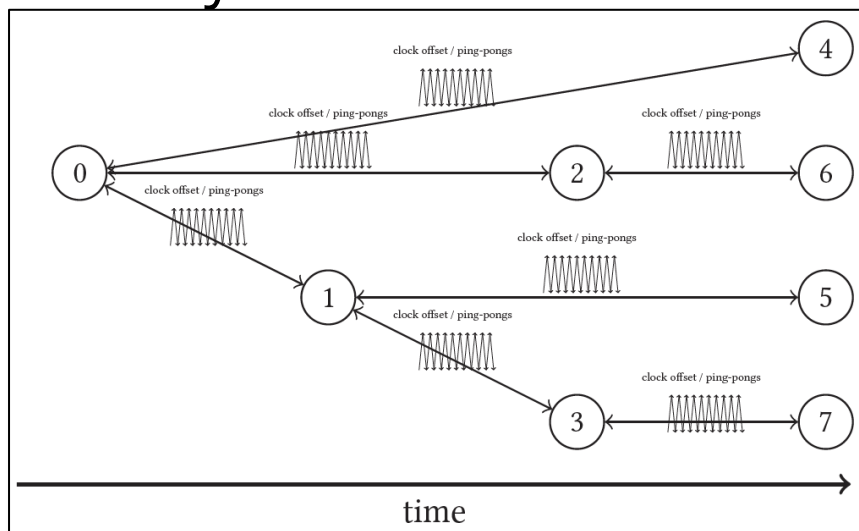
- MPI provides abstractions on many levels
- Abstract away the hard stuff
 - No user should have to implement reductions
 - Or broadcast
 - Or alltoall
 - Or message matching

MPI should provide
process time
synchronization!



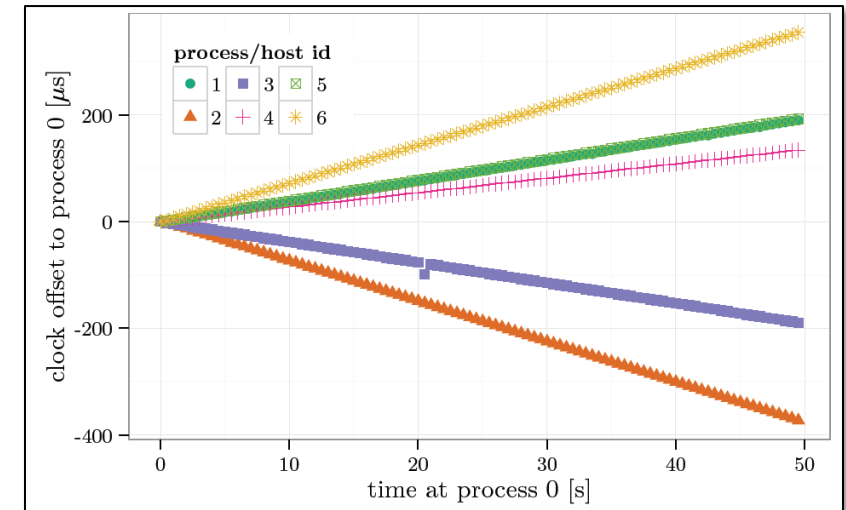
Process Synchronization in Time

- Few machines provide global clocks (MPI_WTIME_IS_GLOBAL)
 - Only have the work
- All others: Synchronization in 2 steps:
 1. Ensure synchronized virtual clocks
 2. Ensure synchronized execution

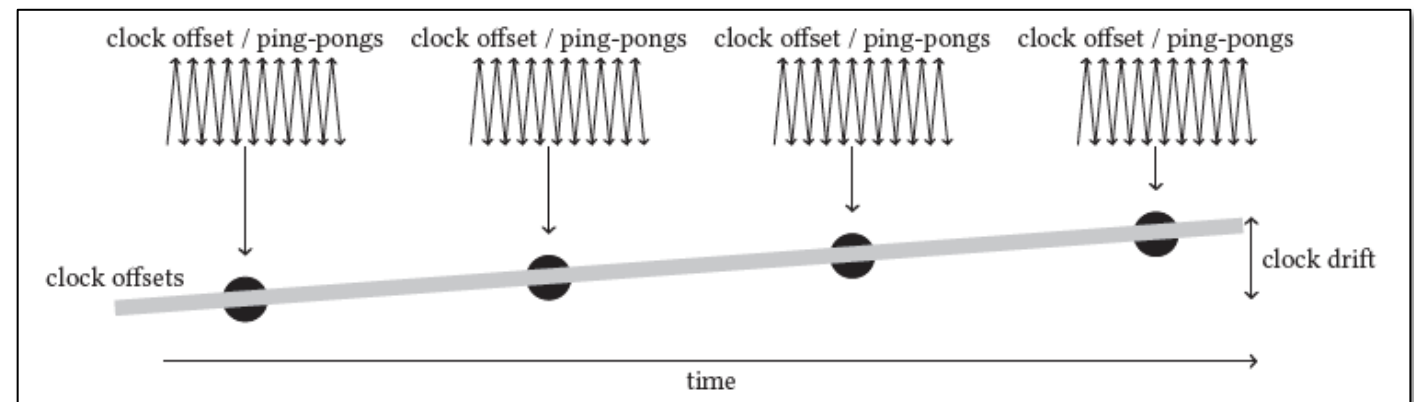


Challenge: Clock Drift

- Clocks run at different speeds (1-10ppm)
- Impacted by temperature and manufacturing differences
- Drift correction:
 - Include drift in local clock synchronization
 - Periodically re-synchronize



S. Hunold and A. Carpen-Amarie. 2015. On the Impact of Synchronizing Clocks and Processes on Benchmarking MPI Collectives. EuroMPI '15. <https://doi.org/10.1145/2802658.2802662>



Introducing: MPIX_Harmonize

- Synchronizes internal clocks in regular intervals
 - To correct for clock drift
 - Every few seconds (automatically adjustable?)
- Attempts to block processes until a common point in time
 - flag == true if calling process reached deadline
 - flag == false if calling process missed deadline
 - Application deals with synchronization misses
- Processes resume execution in harmony

The proposed function MPIX_Harmonize has the following signature:

```
int MPIX_Harmonize(MPI_Comm comm, int *flag);
```



MPIX_Harmonize Algorithm

- Periodically resynchronize clocks
 - Or if a process previously missed a deadline
- Process 0 broadcasts new deadline
- Processes block until deadline, return success or failure

Algorithm 1 Algorithm for MPiX_Harmonize.

```
Require: comm                                ▷ Input communicator
Require: outflag                               ▷ Output flag
1: slack ← sync_slack(comm)
2: flag ← 0                                       ▷ Local and global state
3: root ← 0
4: if last_sync_failed(comm) then
5:   flag ← SYNC_FAILED                           ▷ Locally failed
6: end if
7: if elapsed_since_last_sync(comm) > 1.0s then
8:   flag ← flag | SYNC_EXPIRED                 ▷ Locally expired
9: end if
10: flag ← MPI_Reduce(flag, MPI_BOR, root, comm)
11: if root = comm_rank(comm) then
12:   if flag then                                ▷ Global State
13:     sync_time ← -1.0                          ▷ Trigger clock sync
14:     if flag | SYNC_FAILED then
15:       slack ← slack × 1.5                    ▷ Increase slack
16:       sync_slack(comm) ← slack
17:     end if
18:   else
19:     sync_time ← global_time() + slack
20:   end if
21: end if
22: sync_time ← MPI_Bcast(sync_time, root, comm)
23: if sync_time < 0.0 then                        ▷ Negative time triggers sync
24:   sync_clocks();
25:   MPI_Reduce(0, root, comm)                 ▷ Reduce clock-sync jitter
26:   last_sync_time(comm) ← global_time()
27:   sync_time ← global_time() + slack        ▷ New sync time
28:   sync_time ← MPI_Bcast(sync_time, root, comm)
29: end if
30: if sync_time < global_time() then
31:   outflag ← 0                                  ▷ Missed deadline
32:   last_sync_failed(comm) ← TRUE             ▷ Store for next call
33: else
34:   outflag ← 1                                  ▷ Success, wait for sync time
35:   while sync_time > global_time() do
36:     end while
37: end if
```

Usage in Collective Benchmarks

1. Harmonize processes
2. Perform operation under test
3. Discard invalid measurements

```
5 for (int i = 0; i < NITER; ++j) {
6   int flag, check_idx;
7   /* synchronize in space and time */
8   MPIX_Harmonize(comm, &flag);
9   /* take measurement */
10  double begin = MPI_Wtime();
11  MPI_Reduce(..., comm);
12  check_idx = i%CHECK_EVERY;
13  check_time[check_idx] = MPI_Wtime() - begin;
14  check_flags[check_idx] = flag;
15  if (i == NITER-1 ||
16      check_idx == CHECK_EVERY-1) {
17    /* discard invalid timings */
18    MPI_Allreduce(MPI_IN_PLACE, check_flags,
19                 check_idx,
20                 MPI_INT, MPI_LAND, comm);
21    for (int k = 0; k < check_idx; ++k) {
22      if (check_flags[k]) { /* valid */
23        wtime += check_time[k];
24        valid_iterations++;
25      }
26    }
27  }
```

Evaluation

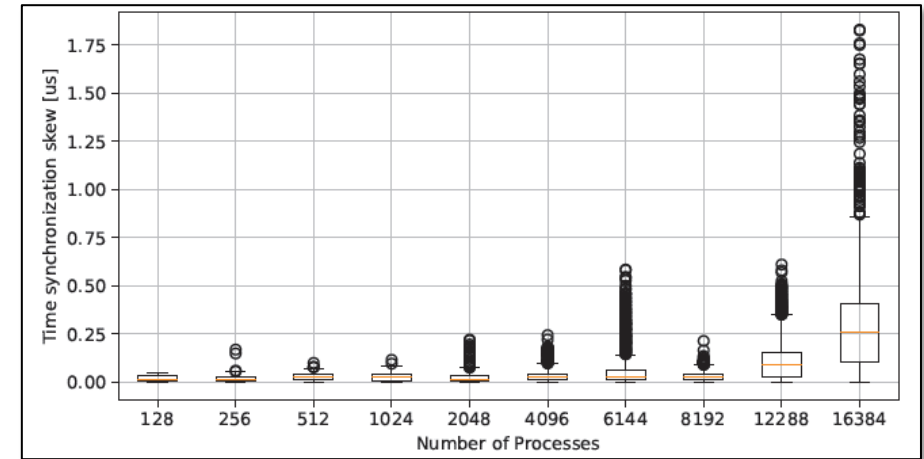
- All experiments performed on *Hawk* installed at HLRS, Stuttgart, Germany
- 2x64 core AMD EPYC Rome
- 200Gbit/s ConnectX-6
- Up to 64 nodes (64k processes)



Cost of Process Synchronization

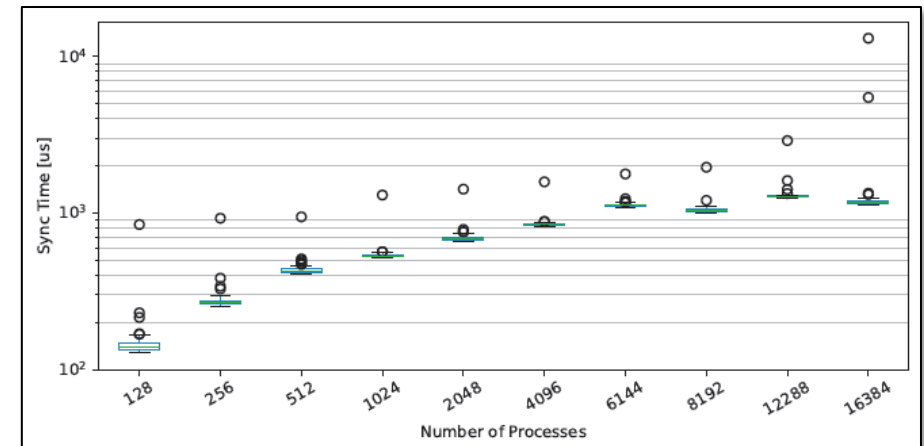
- Synchronization Skew:

- Majority: $< 0.5\mu\text{s}$
- Outliers: $< 2\mu\text{s}$
- $\sim 1/100$ of barrier



- Clock Synchronization $< 1\text{ms}$

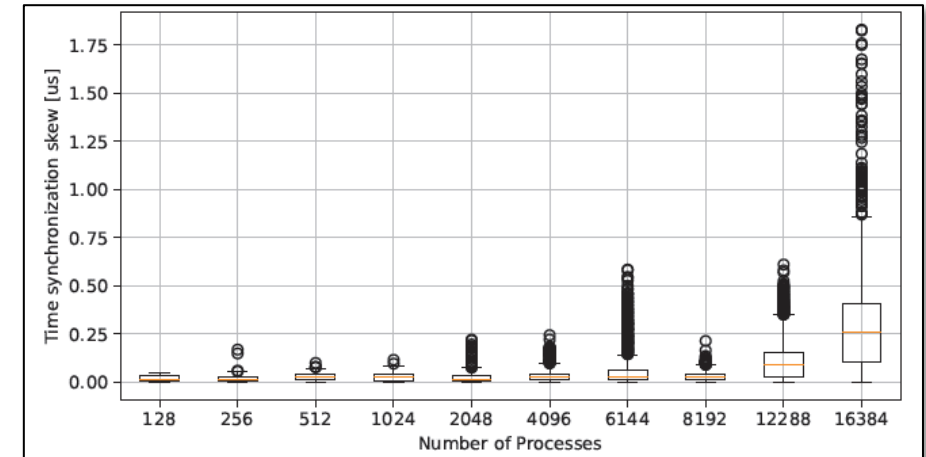
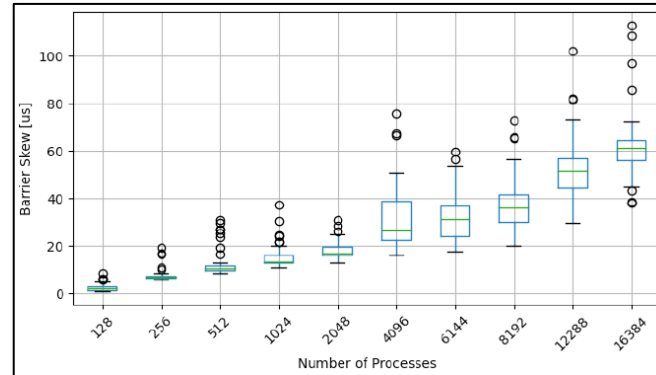
- First synchronization more costly due to connection setup
- 0.1 – 1% overhead if done every second



Cost of Process Synchronization

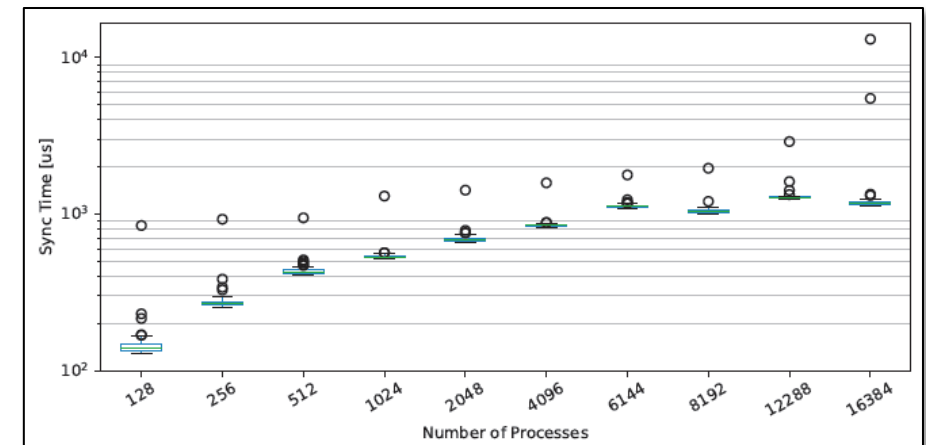
- Synchronization Skew:

- Majority: $< 0.5\mu\text{s}$
- Outliers: $< 2\mu\text{s}$
- $\sim 1/100$ of barrier



- Clock Synchronization $< 1\text{ms}$

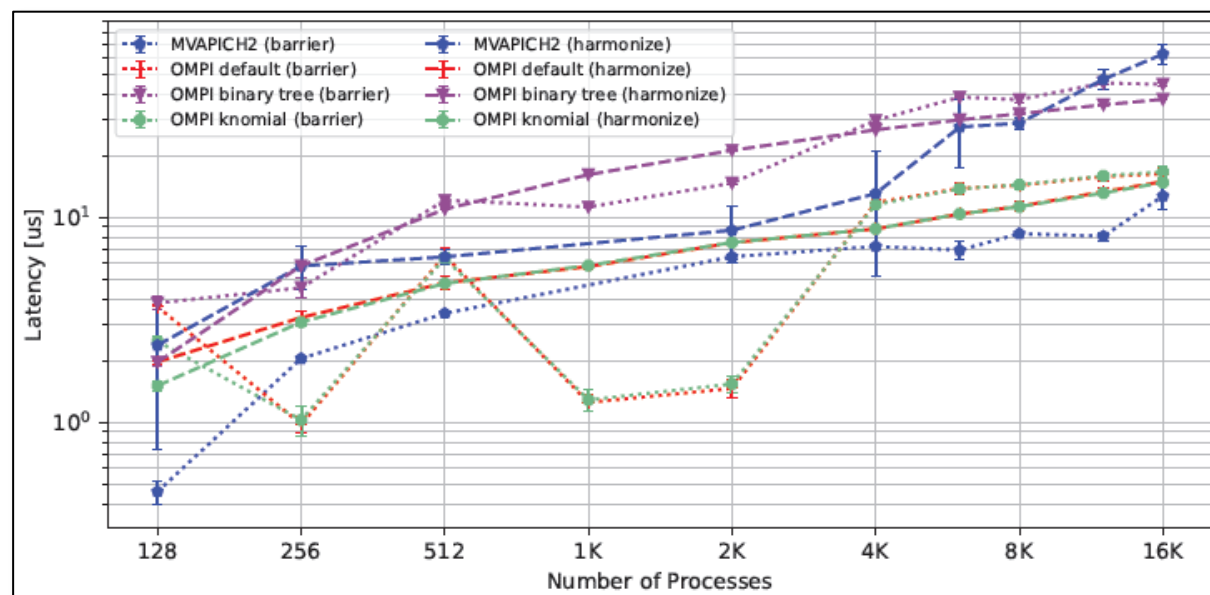
- First synchronization more costly due to connection setup
- 0.1 – 1% overhead if done every second



Impact on MPI_Bcast

- 4B messages, process scaling
- Open MPI: knomial impacted by barrier algorithm
- MVAPICH: lower reported latency with barrier-synchronization

Synchronization does not impact algorithm performance, but latency reported by OSU benchmarks.

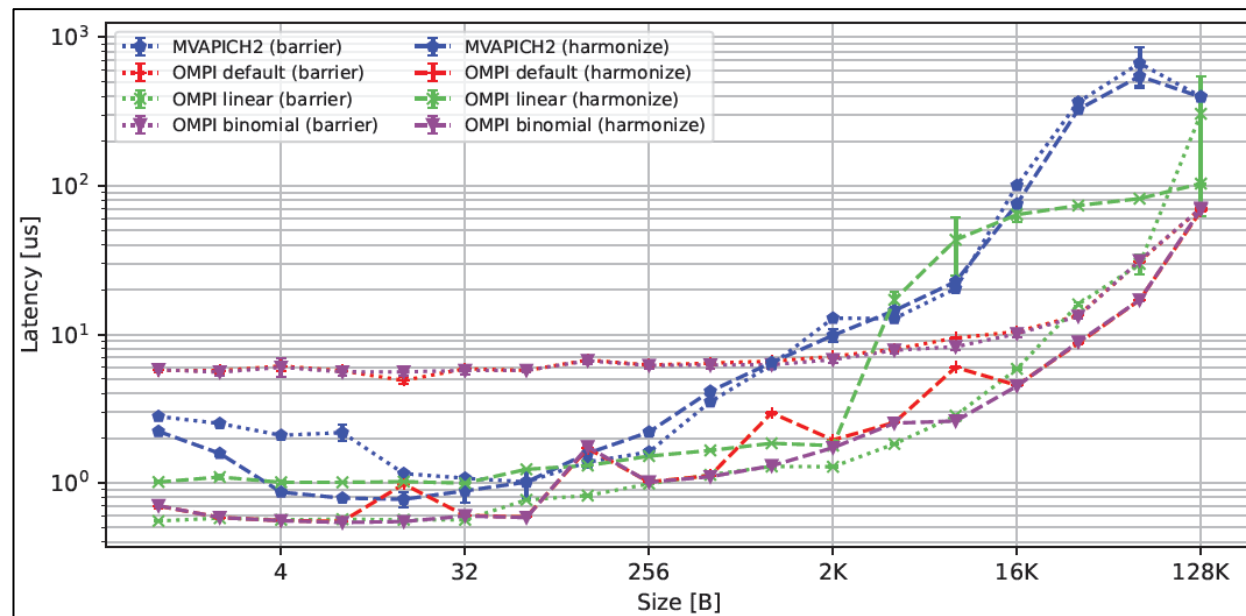


Impact on MPI_Reduce Benchmarks

- Some algorithms in Open MPI show large spread between barrier-synchronization and MPIX_Harmonize (binomial)
- Open MPI default heuristic needs to be re-evaluated

Synchronization does not impact algorithm performance, but latency reported by OSU benchmarks.

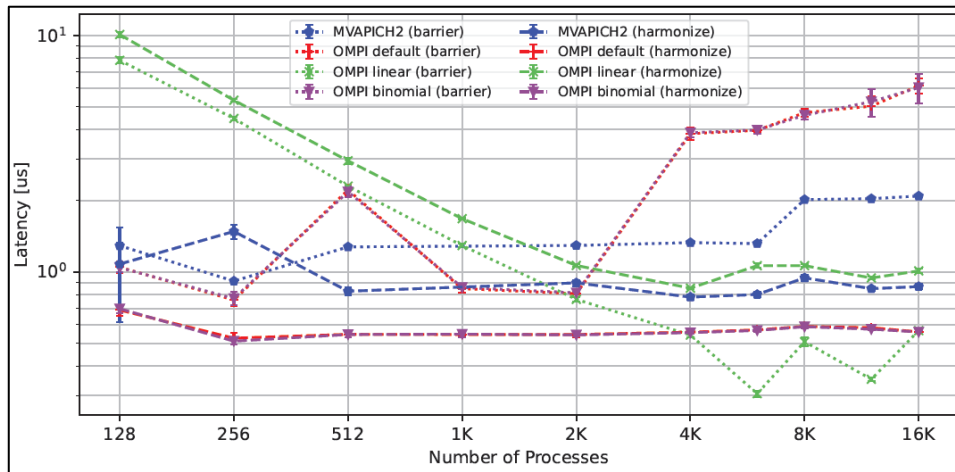
MPI_Reduce on 16k processes



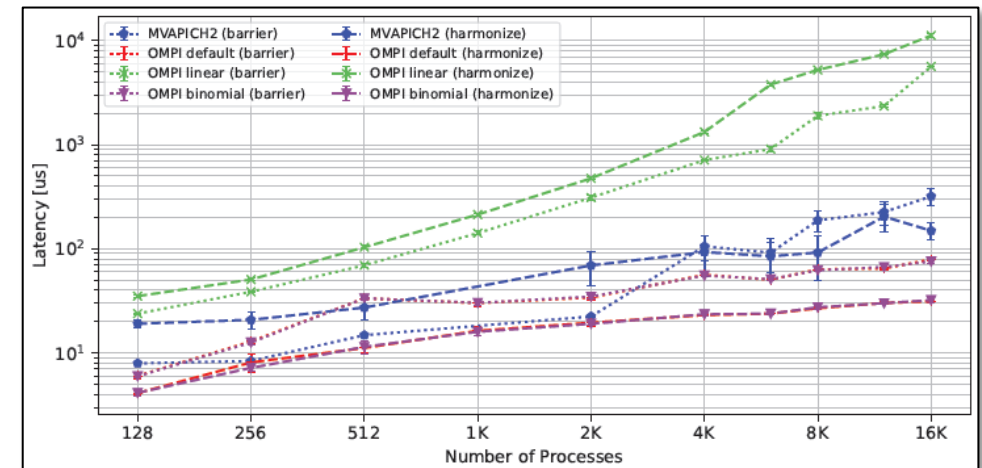
Detour: What do we actually measure?

- The mean of means is not a good metric
 - Mean over all processes
 - Is that really representative?
- Tightly coupled applications sensitive to imbalances
- Better: max of means?

$$t_{avg} = \frac{1}{p} \sum_p \left(\frac{1}{N} \sum_{i=0}^N t_i \right)$$



Process scaling, 4B messages, MEAN

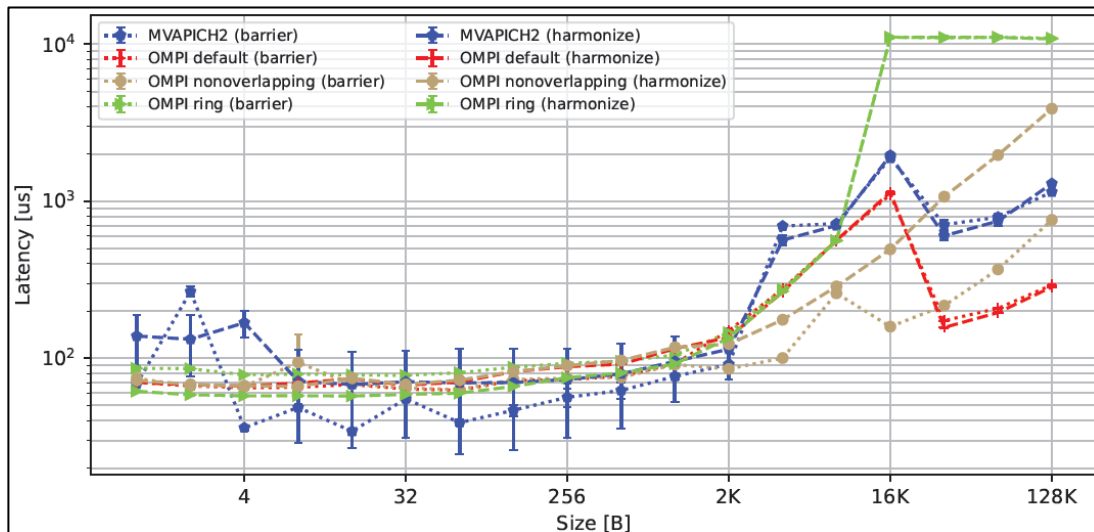


Process scaling, 4B messages, MAX

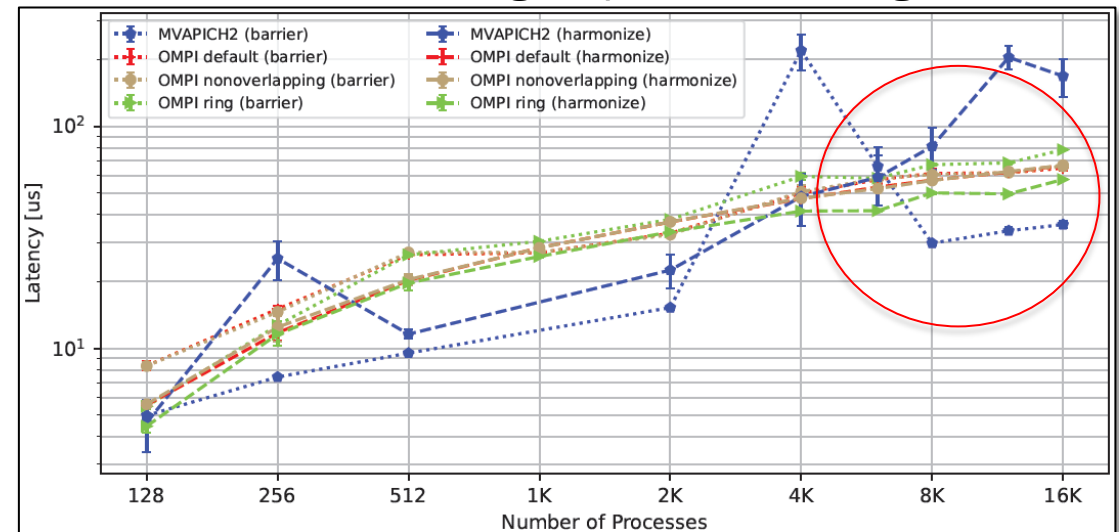
Impact on MPI_Allreduce

- Impact less pronounced on non-rooted collectives
- Open MPI & MVAPICH2 heuristics should be re-evaluated

Message scaling, 16k processes



4B messages, process scaling



Summary

- Proper time synchronization is hard but important
- Barriers introduce arbitrary arrival patterns
- MPI should provide synchronization infrastructure

MPIX_Harmonize

Process Synchronization in Space and Time
with local synchronization check

Acknowledgements

- HLRS: access to Hawk
- Exascale Computing Project (17-SC-20-SC)
- Austrian Science Fund (FWF): project P 33884-N

