



German Research School
for Simulation Sciences

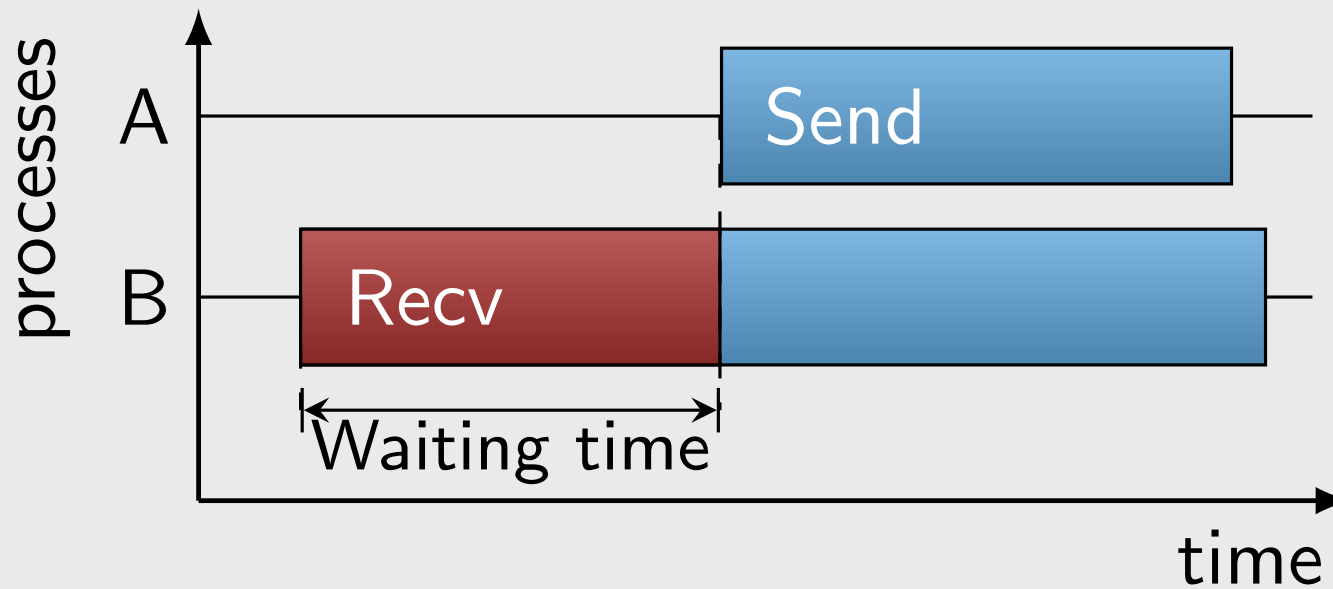
Catching Idlers with Ease: A Lightweight Wait-State Profiler for MPI Programs



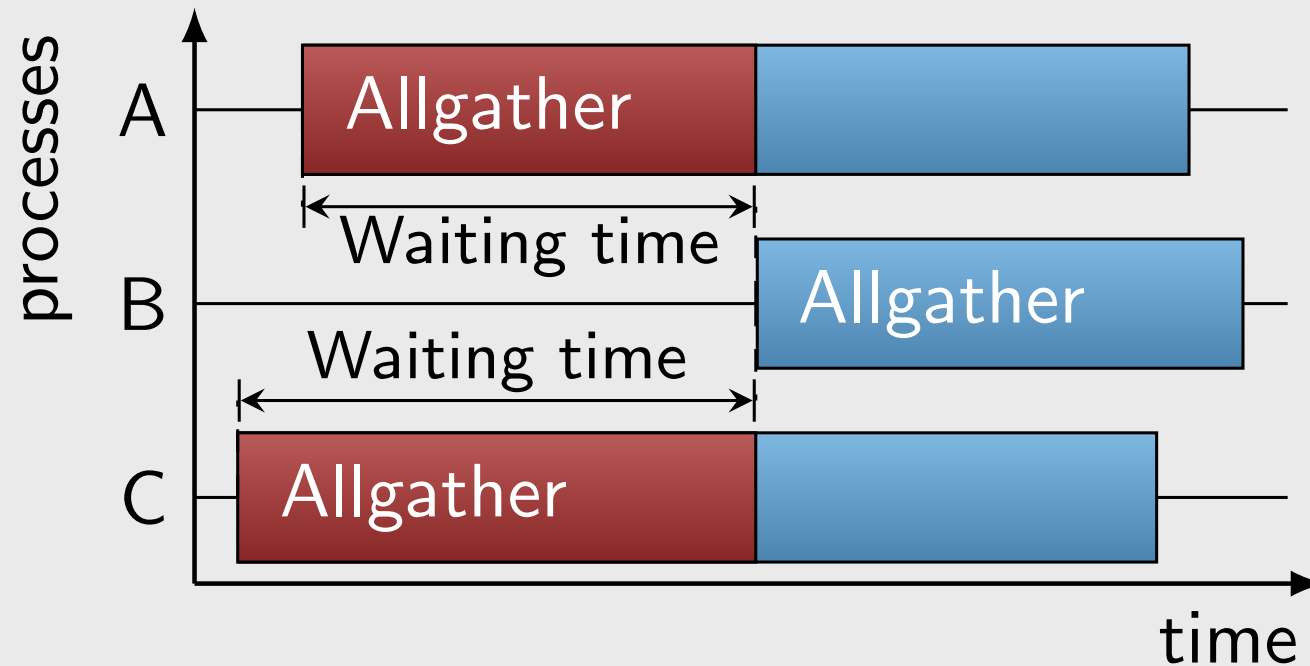
Guoyong Mao, David Böhme, Markus Geimer, Marc-André Hermanns,
Daniel Lorenz and Felix Wolf

Petascale Tools Workshop, Madison, WI, USA, August 4, 2014

Late sender



Wait an NxN



What we want to know



What we measure



Execution time



Execution time

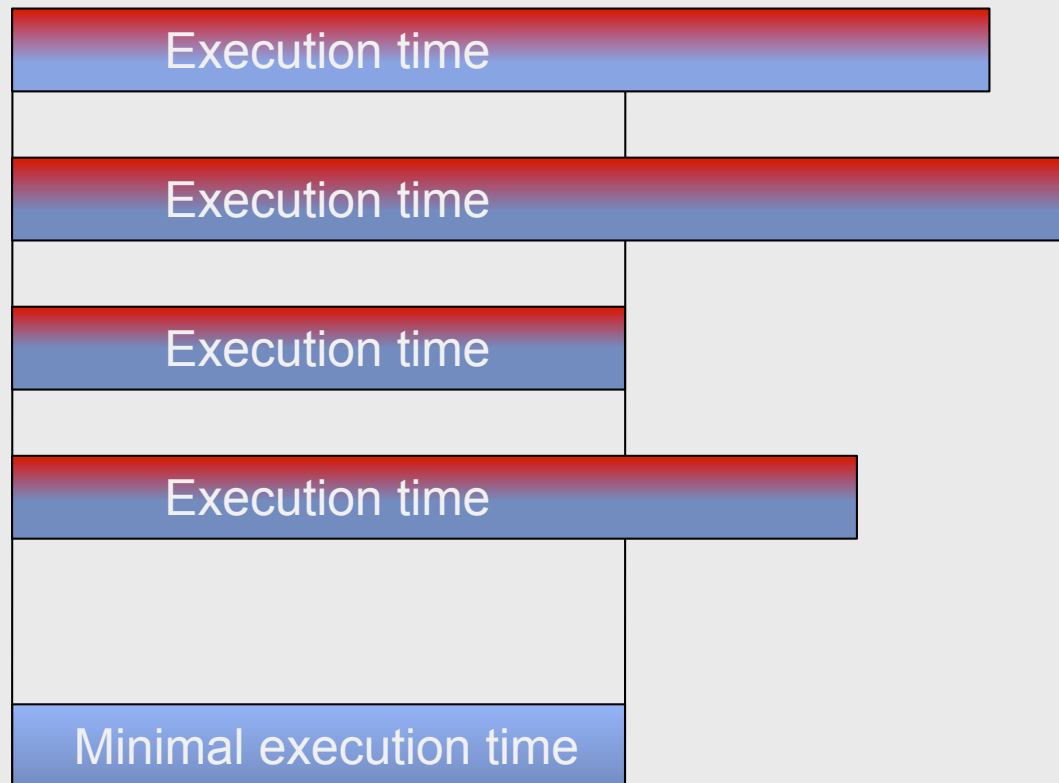


Execution time

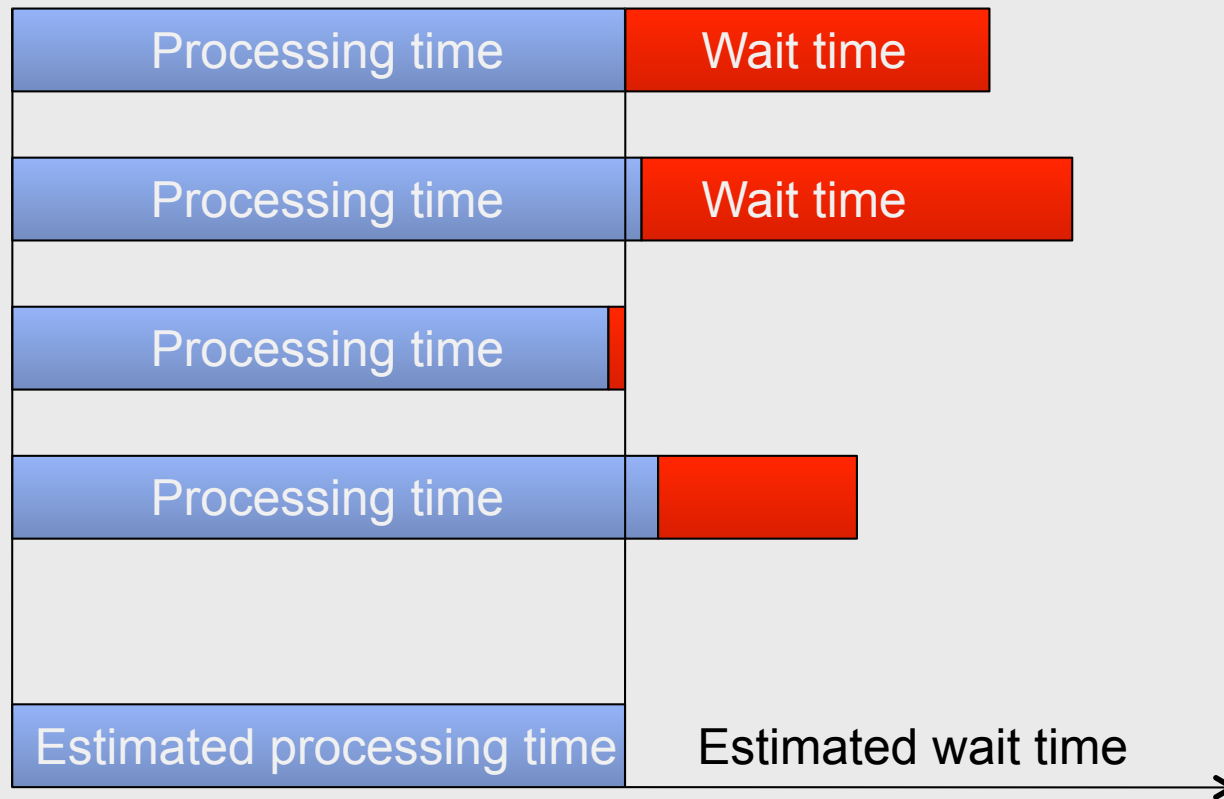


Execution time

The minimum idea



The minimum idea



Considered parameters

- We consider
 - MPI function
 - Message size
 - Receiver rank
- Other possible parameters
 - Sender rank
 - Data type
- Tradeoff between
 - Number of samples for a meaningful minimum and amount data
 - Parameters considered
- Need to find the relevant parameters.

Algorithm

For every combination of

- MPI function
- Message size class
- Process

record the

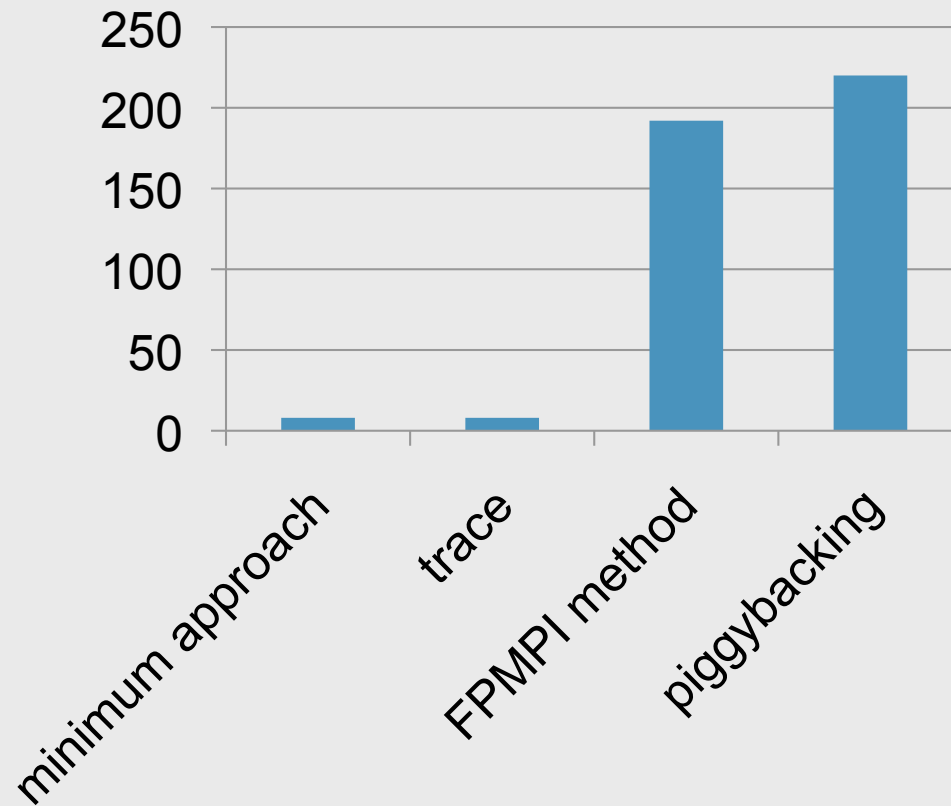
- Minimum execution time

For every combination of MPI call path and message size class record the

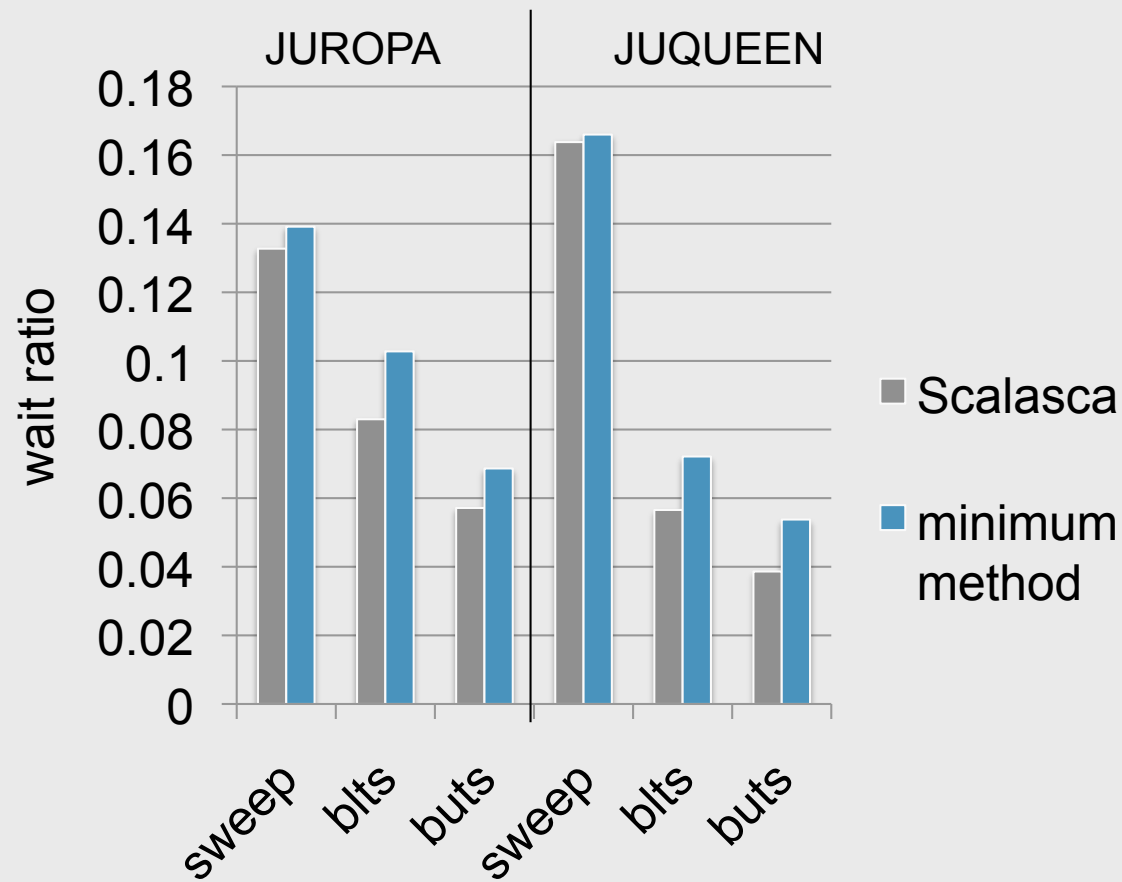
- Number of visits
- Total execution time

At the end of the profiling run, subtract the minimum from the execution time for every visit to calculate the wait time.

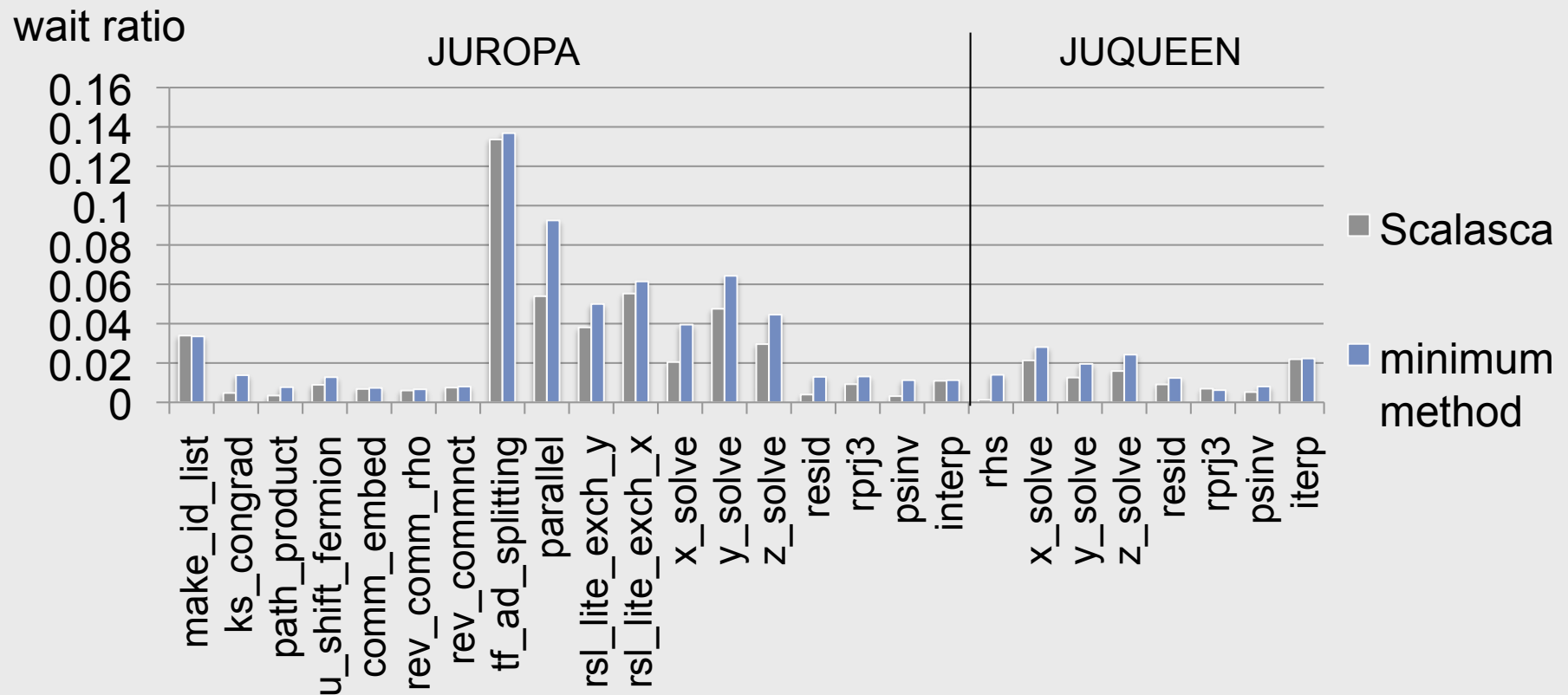
Per-call overhead increase compared to profiling overhead w/o wait state analysis (%)



Accuracy MPI_Recv



Accuracy MPI_Wait



Accuracy MPI_Wait

wait ratio

0.16
0.14
0.12
0.1
0.08
0.06
0.04
0.02
0

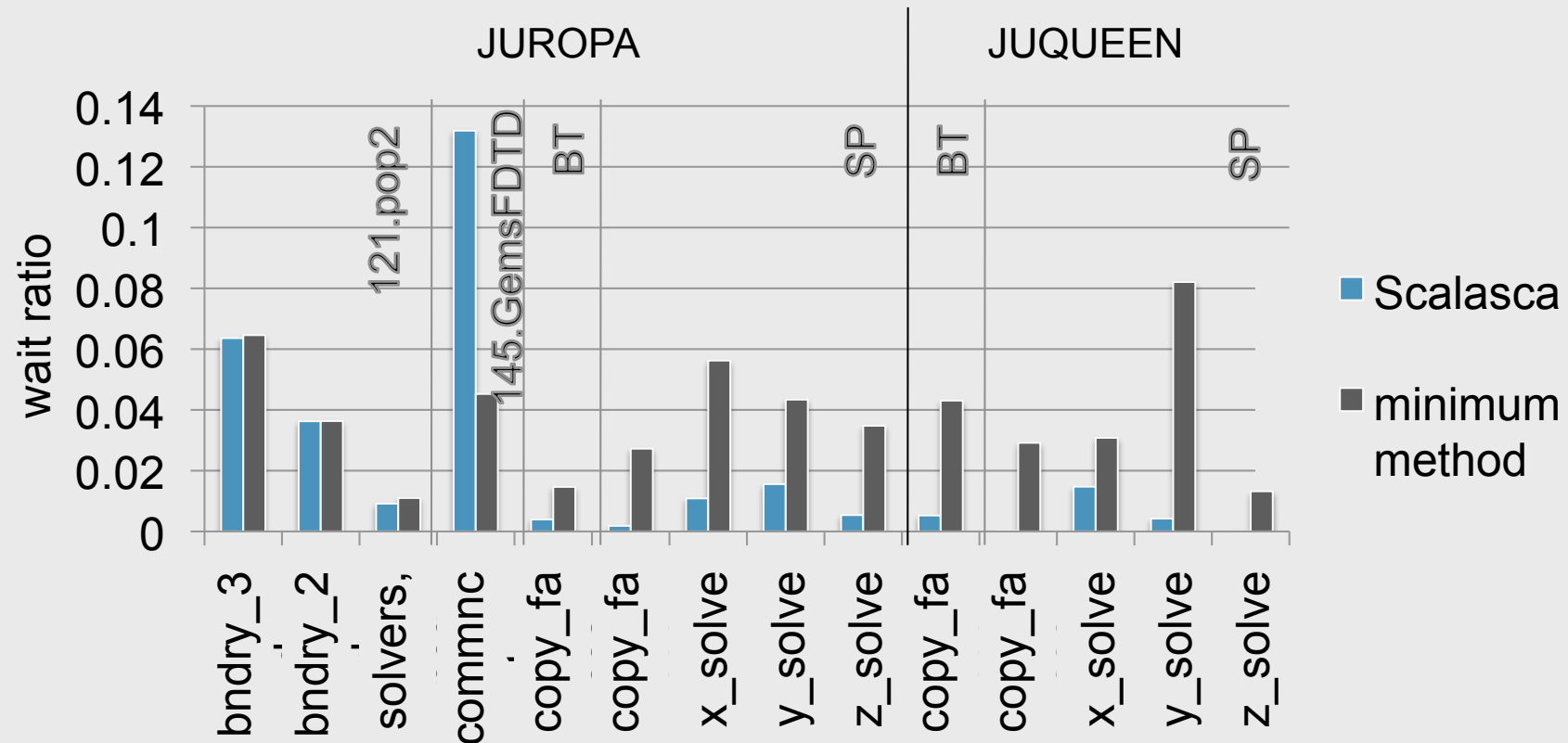
JUROPA

JUQUEEN

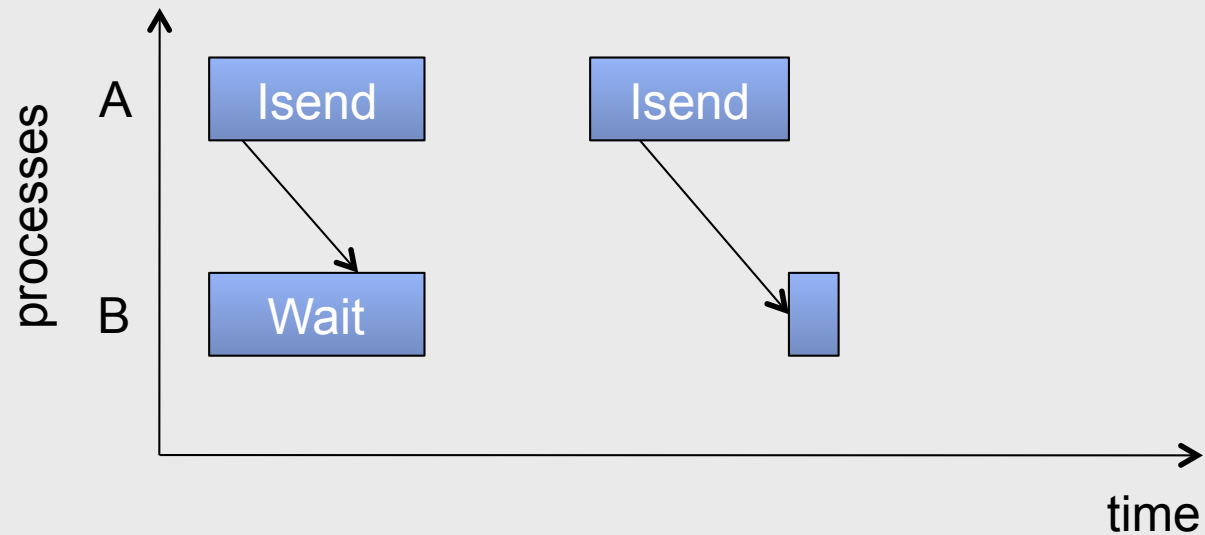
make_id_list
ks_congrad
path_product
u_shift_fermion
comm_embed
rev_comm_rho
rev_commct
tf_ad_splitting
parallel
rsl_lite_exch_y
rsl_lite_exch_x
x_solve
y_solve
z_solve
resid
rprj3
psinv
interp
rhs
x_solve
y_solve
z_solve
resid
rprj3
psinv
interp

■ Scalasca
■ minimum method

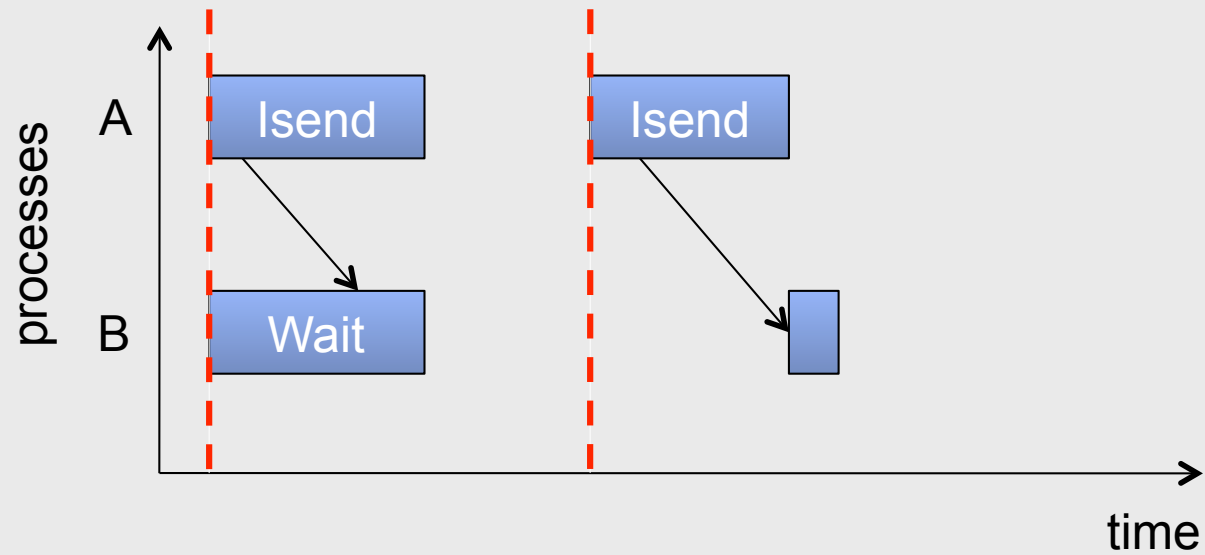
Accuracy MPI_Waitall



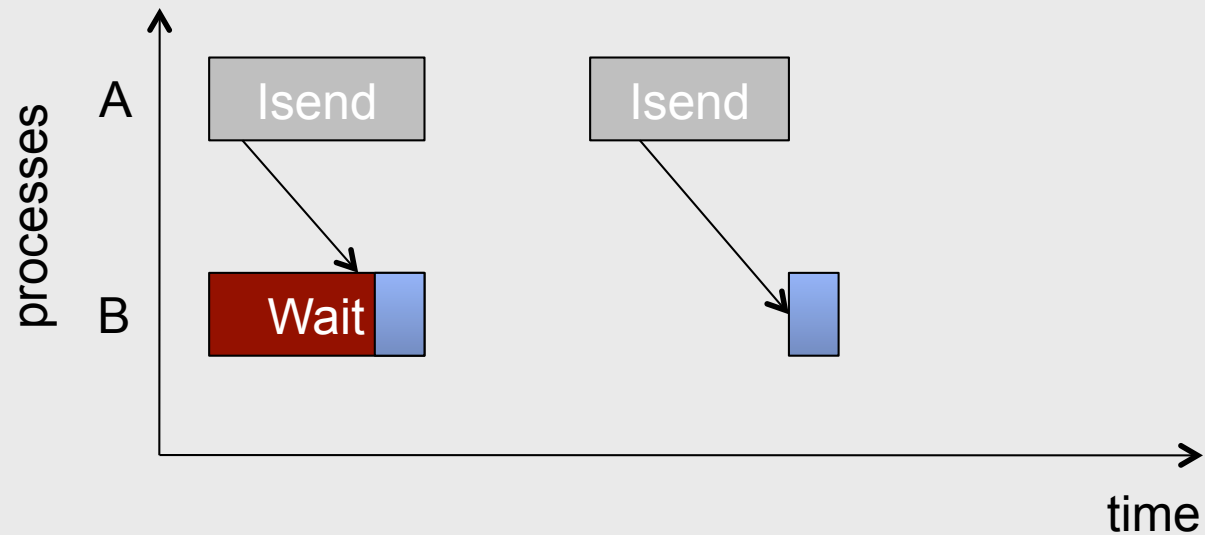
Non-blocking communication



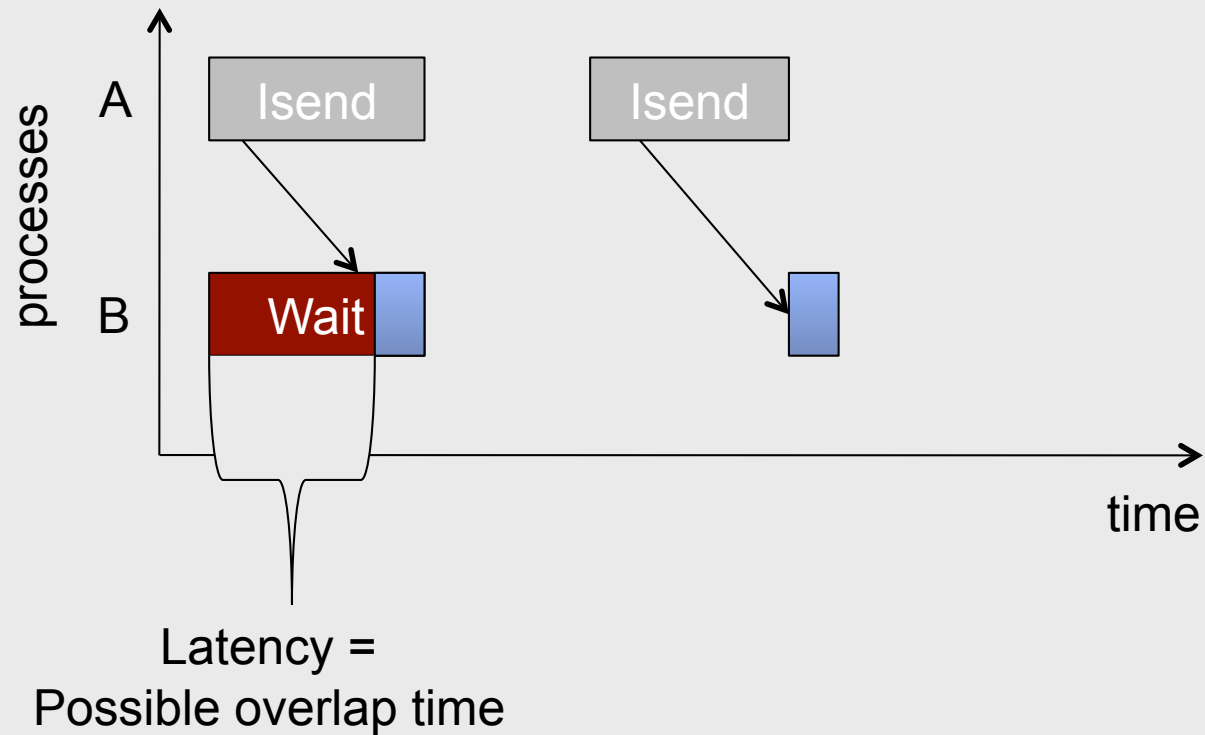
Scalasca detects no wait state



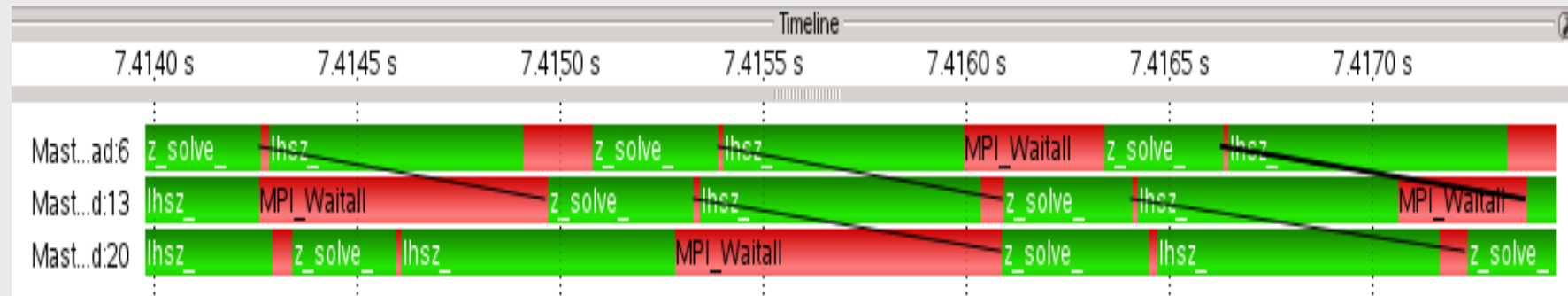
Minimum approach does calculate wait states



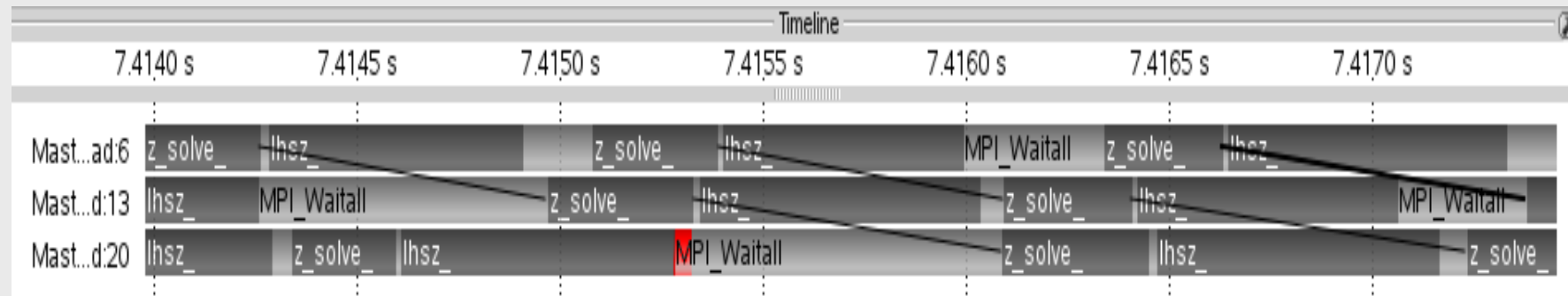
But is this wrong for performance analysis?



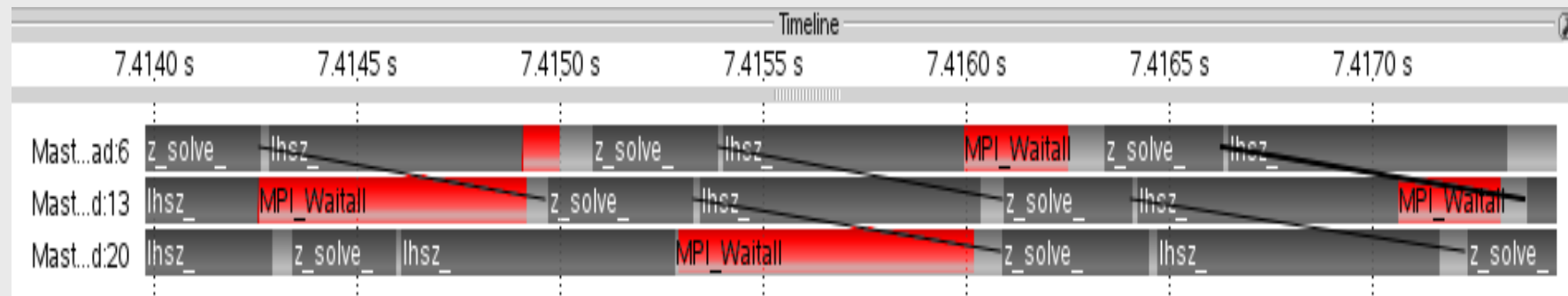
Detailed example from SP



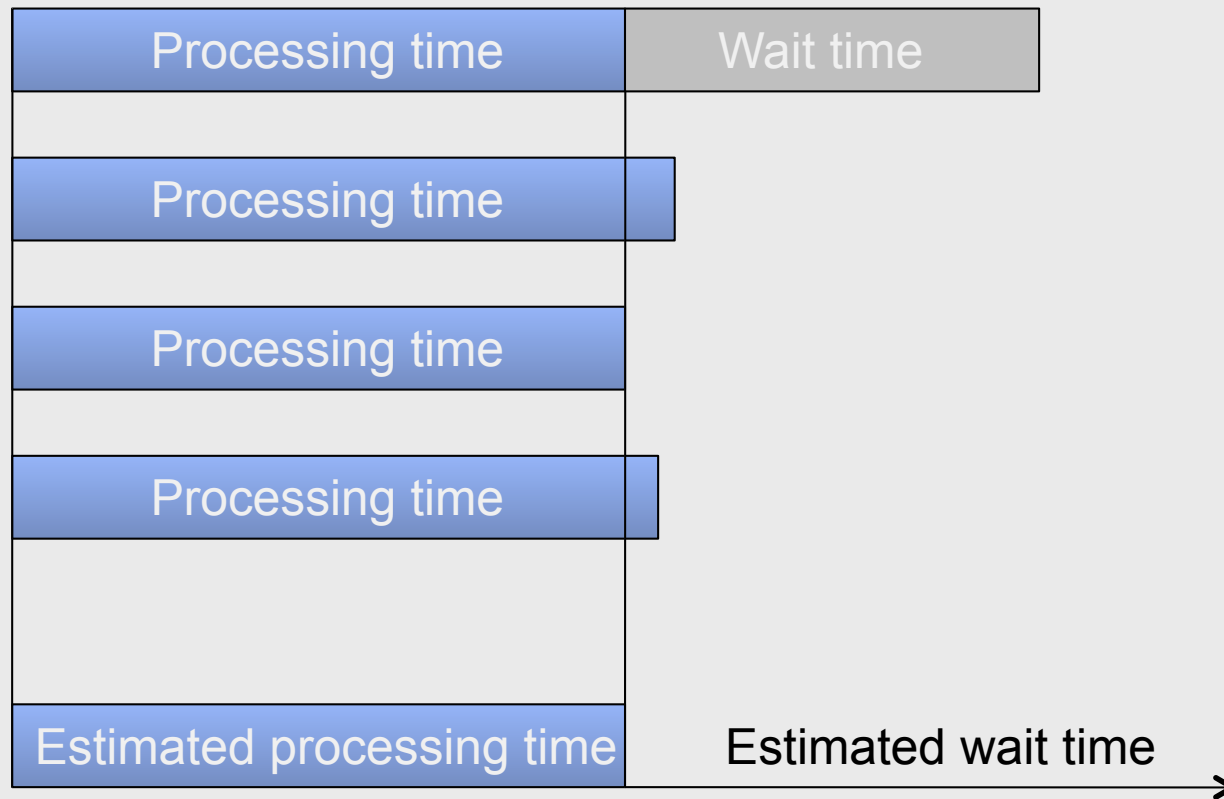
Wait time according to Scalasca



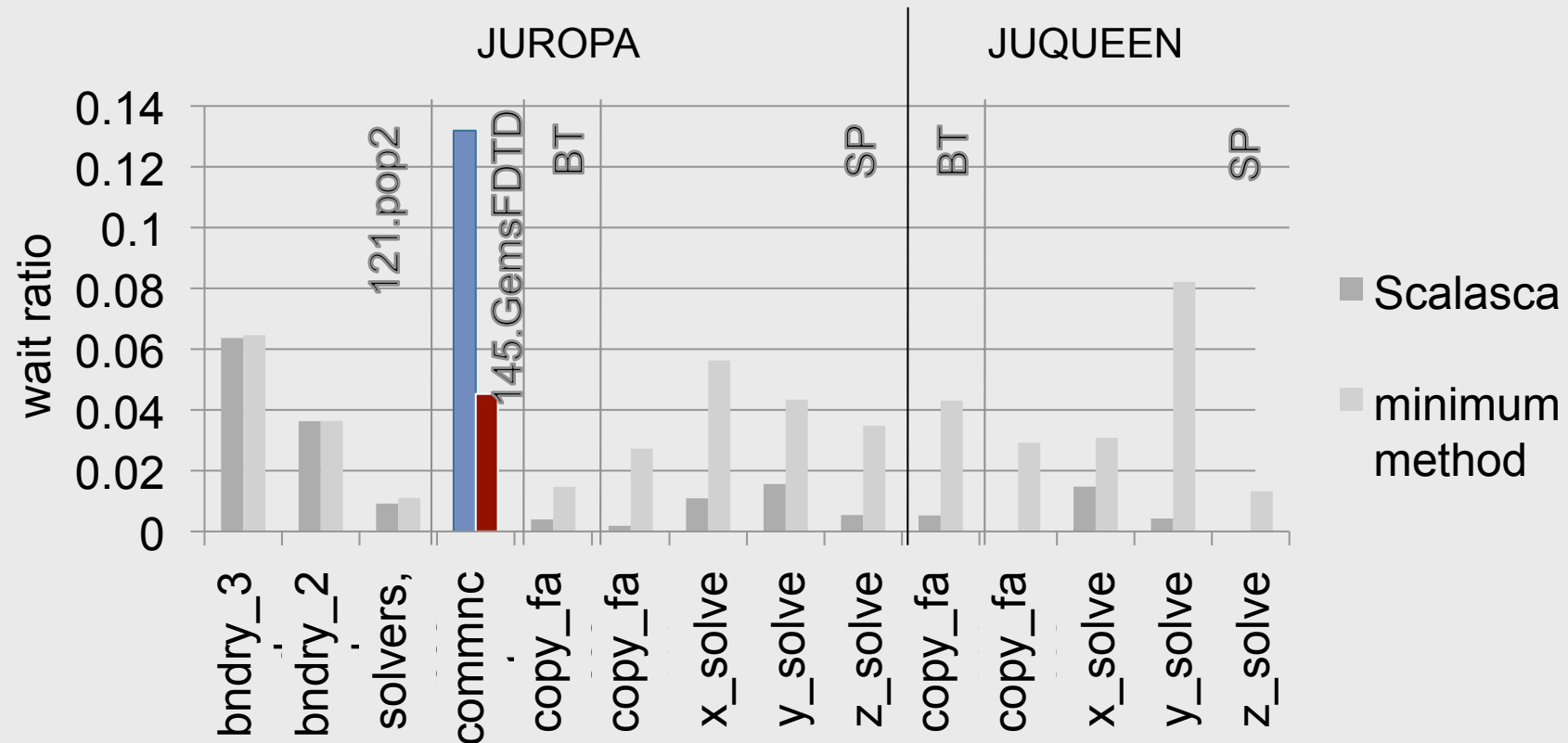
Wait time according to minimum method



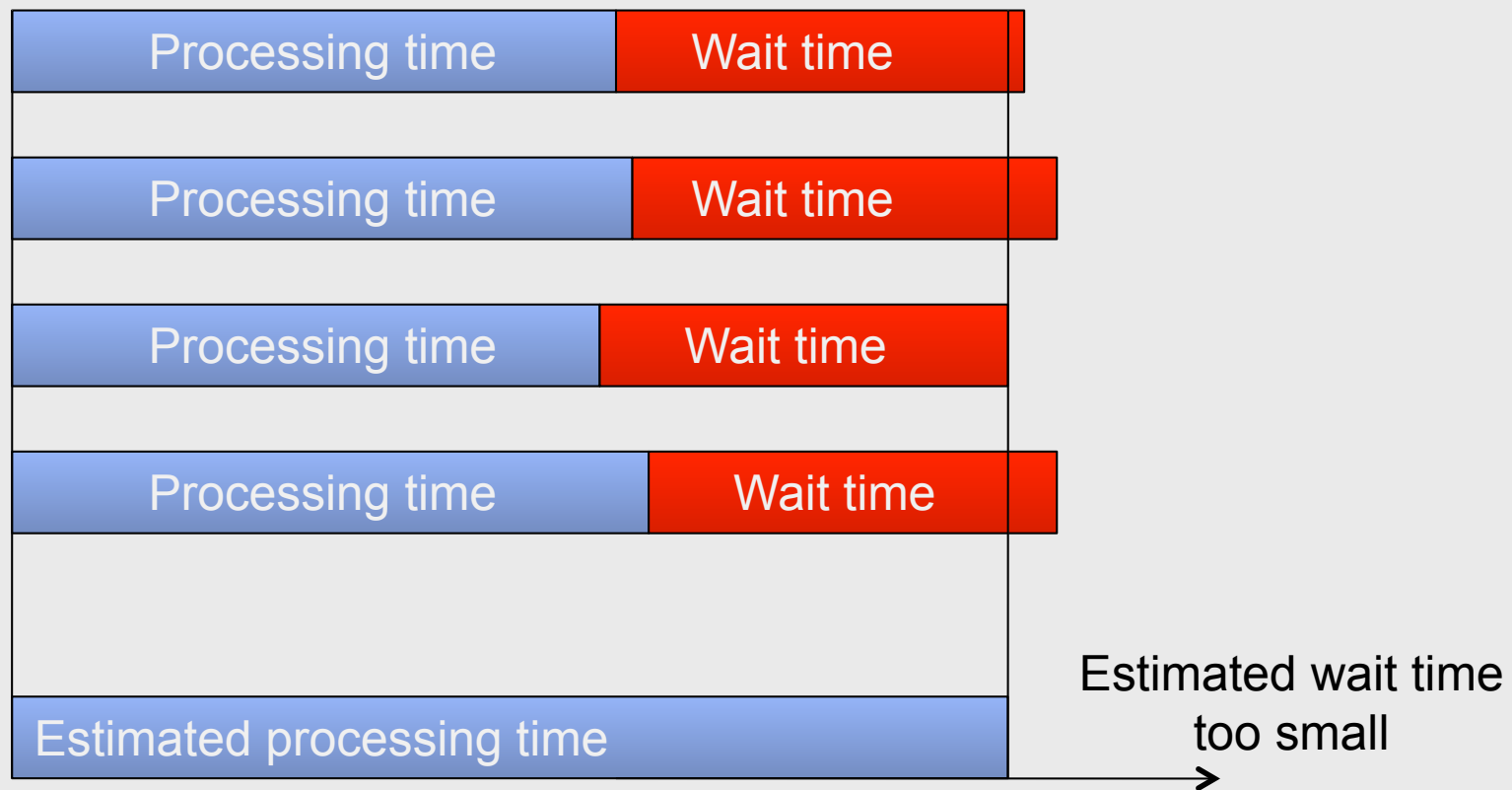
Jitter may cause a little higher wait time



Accuracy MPI_Waitall



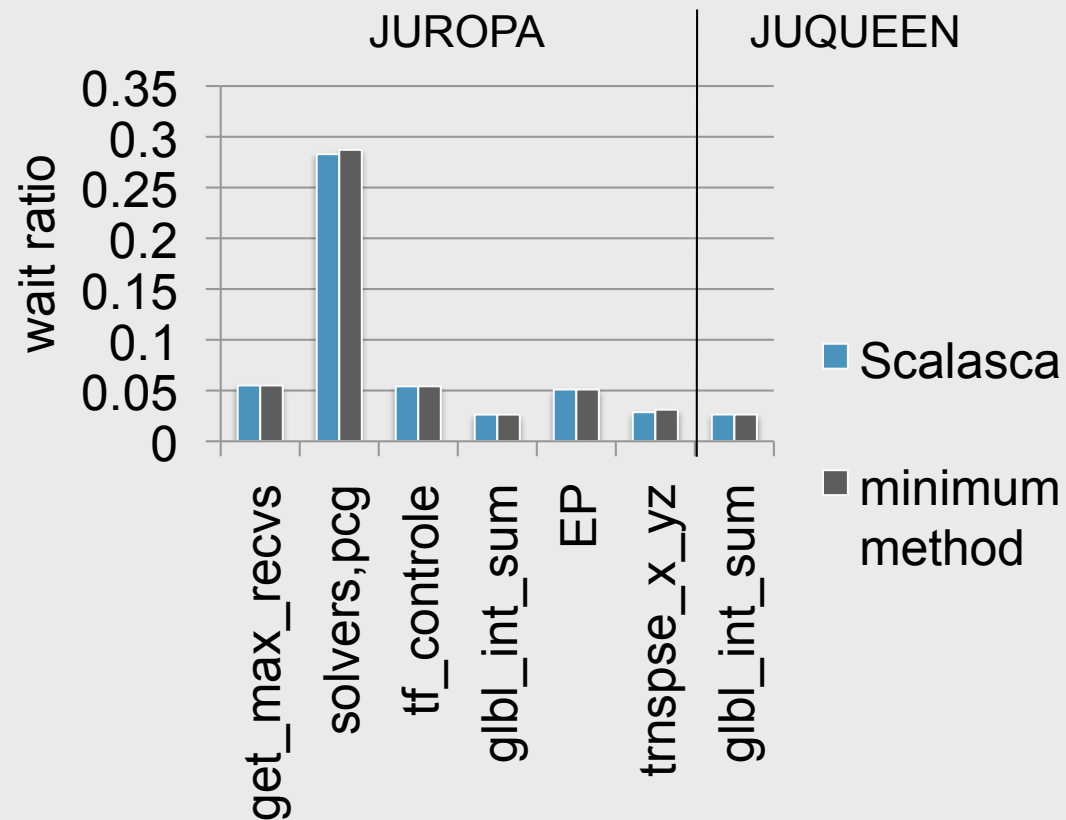
Static imbalance



Static imbalance

- Calculating global minima could resolve process local static imbalances
 - Reduction operation after measurement
 - No dilation at measurement time
- Loose sender/receiver parameterization of minima
- For collective operations, global minima were better

Accuracy for Wait at NxN



Conclusion (1)

- Minimum method works for the estimation of blocking and non-blocking communication
 - For blocking communication results similar to Scalasca
 - For non-blocking communication, in Waitall wait time do not match the Scalasca analysis.
- Low runtime overhead
- No trace recording or piggybacking
- May not produce 100% accurate numbers, but
 - Sufficient accuracy to locate performance problems
 - Point to places where we might want to investigate further with trace analysis

Conclusion (2)

- Detection of good minimum crucial
 - Static imbalance
 - Tradeoff between number of parameters and number of samples
- Jitter may lead to minor increase of measured wait state
- For non-blocking communication
 - Count possible overlap time
 - Might be larger than pure Late Sender time
 - Isn't this even more accurate to estimate the optimization potential?



Reference

Guoyong Mao, David Böhme, Marc-André Hermanns, Markus Geimer, Daniel Lorenz, Felix Wolf: *Catching Idlers With Ease: A Lightweight Wait-State Profiler for MPI Programs*. In: EuroMPI '14: Proc. Of the 21st European MPI Users' Group Meeting, Tokyo, Japan, Sep. 9-12, 2014