# MPI Collective Algorithm Selection in the Presence of Process Arrival Patterns

Majid Salimi Beni[1], Biagio Cosenza[1], Sascha Hunold[2]

[1] Department of Computer Science
University of Salerno, Salerno, Italy

[2] Faculty of Informatics, TU Wien
Vienna, Austria

# Outline

## Background

- MPI, Collective Operations and Algorithms
- MPI Collective Algorithm Selection

## Motivation

- Process Arrival Patterns and Algorithm Selection

## Methodology

- Micro-benchmarking technique

## Experimental results

- Simulation study
- Real-world experiments
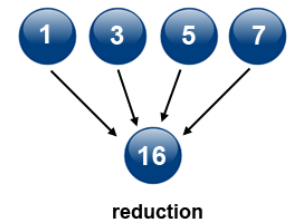- Arrival patterns in the applications
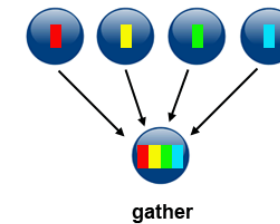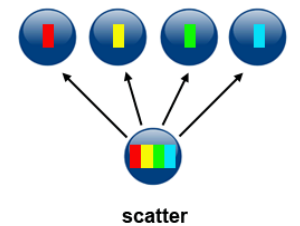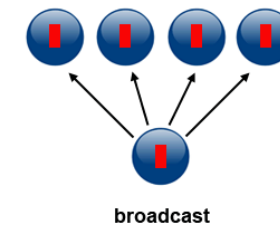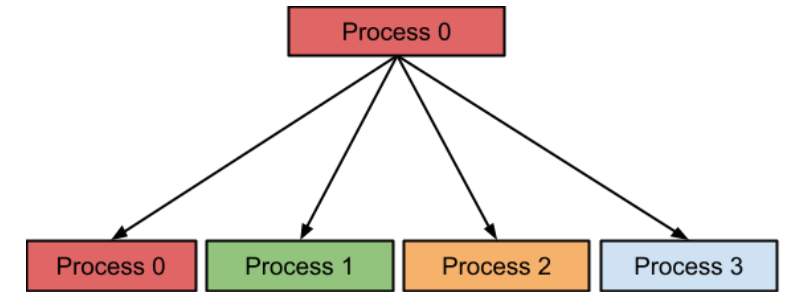
## Conclusion and future work

# Background: MPI, Collective Operations, and Algorithms

❑ MPI (Message Passing Interface)

❑ A standard message-passing library designed to function on parallel computing architectures

❑ MPI collectives

❑ Time-consuming: Big share of HPC applications' runtime is spent while performing collective communications

❑ Efficient implementation of collective operations

❑ Optimal performance

❑ Scalability

❑ Communication overhead

❑ Resource utilization

# Background: MPI Collective Algorithm Selection

❑ MPI standard defines the **semantics** of collective operations

❑ Leaves their **algorithmic implementations** to MPI libraries

❑ MPI libraries provide several algorithms for each collective operation

❑ A decision logic selects one of these algorithms

❑ Algorithm selection of MPI collectives

   ❑ Message size, process count, network topology, available hardware resources, network utilization

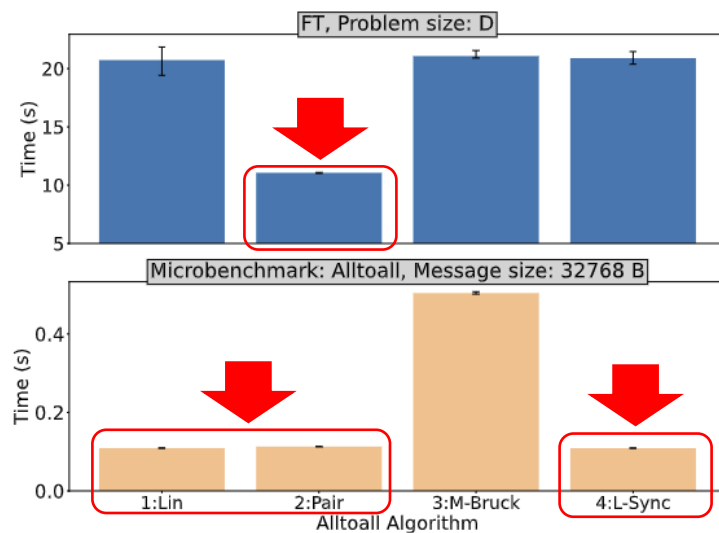❑ Based on the scenario, one algorithm may outperform the others

# Motivation

❑ FT (problem size D) from NAS Parallel Benchmark
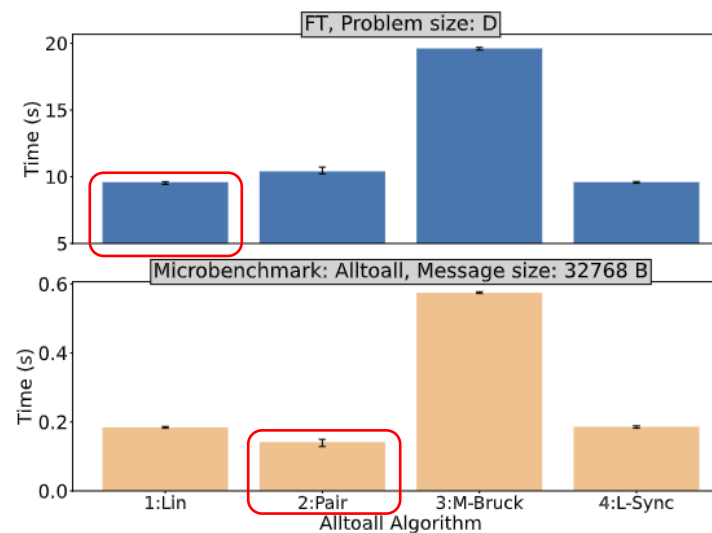
    ❑ Communication-intensive

    ❑ Profiling: MPI_Alltoall with a specific message size takes 50–70% of the total runtime

    ❑ Application vs Micro-benchmark (with the message size found in the application)
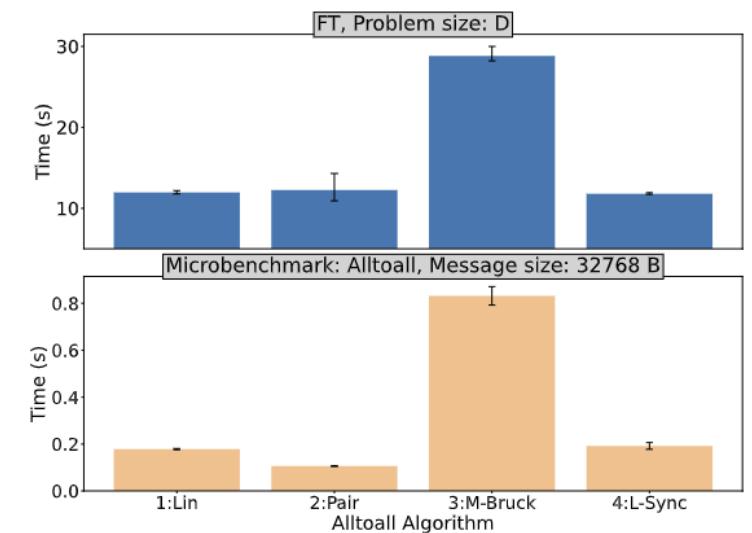
> **Observation:** Choosing the fastest algorithm in the micro-benchmark, doesn't lead to the best performance in the application



(a) *Hydra*  (b) *Galileo100*  (c) *Discoverer*

# Motivation: Process Arrival Patterns and Algorithm Selection

❑ In MPI applications, processes typically don't enter collective operations simultaneously

 ❑ System noise, performance variability, etc.

❑ **Process Arrival Patterns**

❑ **Hypothesis:** Collective algorithms may perform differently when there is process arrival pattern

 ❑ Well-performing collective algorithm under a balanced process arrival pattern may show poor performance under an imbalanced process arrival pattern

❑ **Proposed solution**: Micro-benchmarking and exposing collective algorithms to different arrival patterns

 ❑ Simulation (SimGrid toolkit)

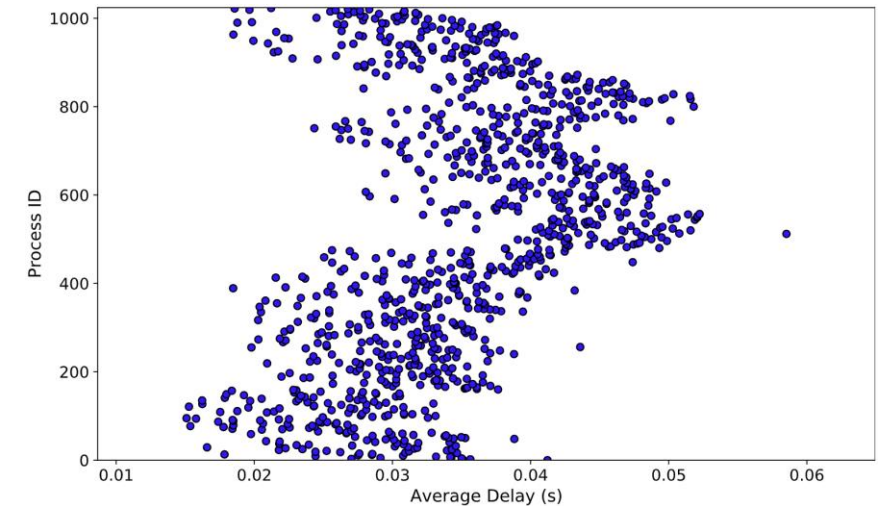 ❑ Real-world experiments on production machines (Hydra, Galileo100, Discoverer)



Fig: Avg. process delay (skew) across all MPI_Alltoall calls in FT (NAS parallel benchmarks) on Galileo100 with 32 × 32 processes.
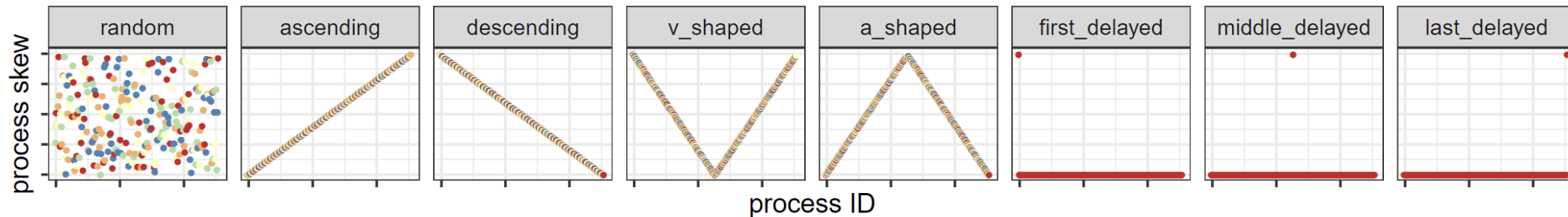
# Methodology

```
for (i=0; i<NREP; i++) {
#ifdef SIMULATOR
    double wait_time = get_arrival_pattern_delay();
    usleep(wait_time);
#else
    MPIX_Harmonize();
    double skew_time = MPI_Wtime() +
    get_arrival_pattern_delay();
    while( MPI_Wtime() < skew_time );
#endif
    double start_time = MPI_Wtime();
    MPI_COLLECTIVE(...);
    double end_time = MPI_Wtime();
}
```

Joseph Schuchart, Sascha Hunold, George Bosilca:
Synchronizing MPI Processes in Space and Time.
EuroMPI 2023: 7:1-7:11

Exposing different
algorithms to different
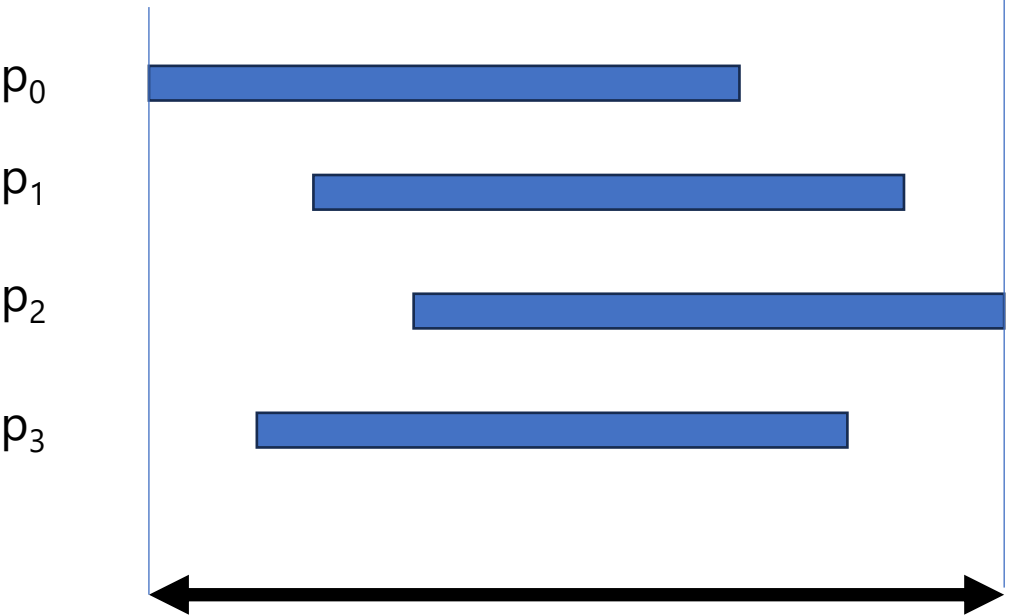(artificial) arrival patterns



Artificial process arrival patterns

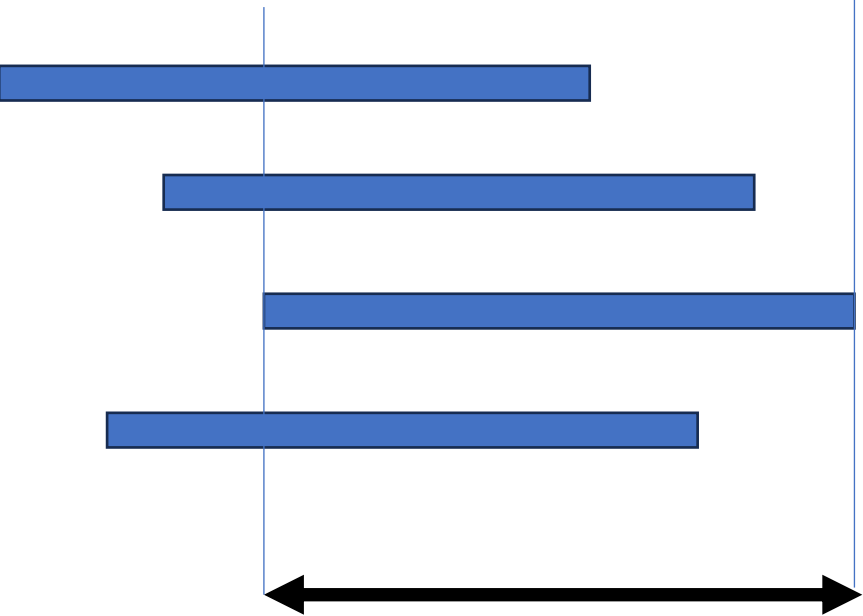# Last Delay Metric

**Total Delay**
(absolute makespan)

**Last Delay**
(more meaningful in case of load imbalance)

$p_0$

$p_1$

$p_2$

$p_3$

Since it's a collective call: it matters most how fast we can complete it when the last process has arrived!

# Simulation results

❑ 1024 processes (32 x 32)

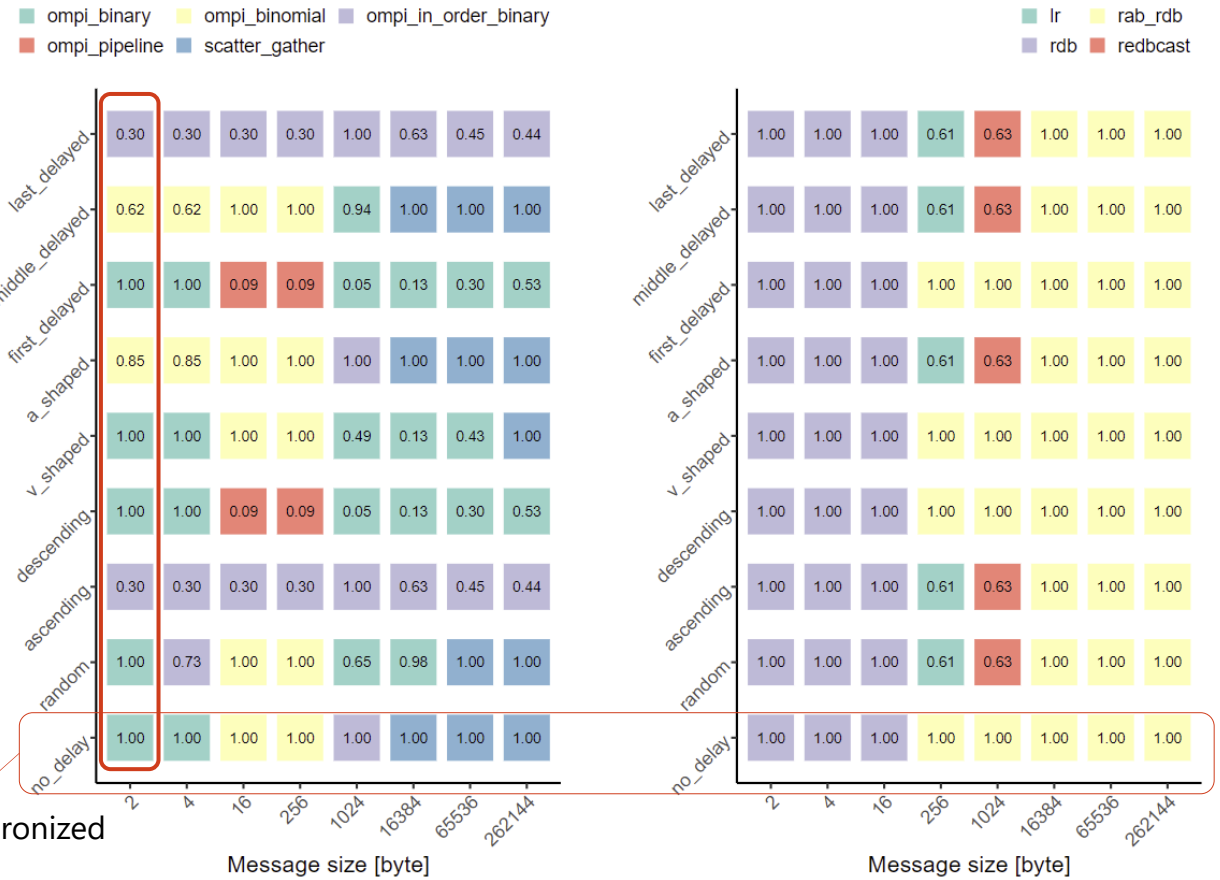The **color:** indicates the best algorithm found for a specific message size

❑ The **value**: denotes the relative performance of this algorithm compared to the best algorithm from the no_delay case

❑ **MPI_Reduce**

    ❑ The optimal algorithm for MPI_Reduce varies with different message sizes and process arrival patterns

❑ **MPI_Allreduce**

    ❑ The reduction step in an Allreduce is a strongly synchronizing sub-task

All processes are perfectly synchronized



(a) MPI_Reduce        (b) MPI_Allreduce

**Arrival patterns impact the collective algorithms**

# Hardware Overview

| Hydra (TU Wien) | Galileo 100 (Cineca) | Discover (Sofia Tech Park) |
|---|---|---|
| 36 nodes, 2x16-core Intel Xeon 2.1GHz | 512 nodes, 2x24-core Intel CascadeLake 8260 | 1128 nodes, 2x64-core AMD Epyc 7H12 |
| Dual-rail Intel **Omni-Path** (100 Gbit/s) | Mellanox **Infiniband** HDR100 | **Infiniband** HDR (Dragonfly+) |
| Open MPI 4.1.5 | Open MPI 4.1.1 | Open MPI 4.1.4 |

# Real-world Experiments

- 1024 processes (32 × 32) processes

- For each arrival pattern, algorithms within 5% of the fastest are in blue

- Knowing the arrival patterns, we can accurately select the best algorithm

- Detecting arrival patterns is time-consuming /infeasible in real-world

**Key Idea:**

Selecting a **robust** algorithm for MPI collectives, capable of performing well when facing various arrival patterns



MPI_Reduce

MPI_Alltoall

Fig: Runtimes of MPI collectives for various message sizes on Hydra

# Real-world Experiments – Robustness

- 1024 processes (32 × 32) processes

- Normalized runtimes to No-delay

- Green rectangles: at least 25% faster than No-delay; Red rectangles: at least 25% slower than No-delay

- For MPI_Reduce: most algorithms are sensitive to process arrival patterns

- **Selection strategy**: Algorithms with more green/grey areas can be good choices

### MPI_Reduce

**Message Size: 8 B**

| Pattern | 1:Lin | 2:Chain | 3:Pipe | 4:Bin | 5:Binom | 6:In-Bin | 7:Raben |
|---|---|---|---|---|---|---|---|
| Ascending | -0.516 | -0.070 | 0.013 | 0.015 | 0.054 | -0.224 | -0.514 |
| Descending | -0.319 | -0.307 | -0.992 | -0.263 | -0.330 | -0.020 | -0.325 |
| First-delayed | -0.018 | -0.959 | -0.992 | -0.824 | -0.489 | -0.112 | -0.019 |
| Mid-delayed | -0.801 | -0.060 | -0.490 | -0.198 | -0.452 | -0.211 | -0.800 |
| Last-delayed | -0.565 | -0.072 | 0.015 | -0.150 | -0.001 | -0.678 | -0.551 |
| Random | -0.564 | -0.100 | -0.036 | -0.113 | -0.111 | -0.062 | -0.130 |

**Message Size: 1024 B**

| Pattern | 1:Lin | 2:Chain | 3:Pipe | 4:Bin | 5:Binom | 6:In-Bin | 7:Raben |
|---|---|---|---|---|---|---|---|
| Ascending | -0.607 | 0.066 | 0.031 | -0.115 | 0.034 | -0.161 | -0.149 |
| Descending | -0.363 | -0.163 | -0.981 | -0.379 | -0.331 | -0.024 | -0.020 |
| First-delayed | -0.011 | -0.986 | -0.996 | -0.903 | -0.516 | -0.105 | -0.001 |
| Mid-delayed | -0.835 | 0.100 | -0.484 | -0.359 | -0.435 | -0.201 | -0.040 |
| Last-delayed | -0.631 | 0.078 | 0.045 | -0.297 | 0.010 | -0.734 | -0.183 |
| Random | -0.287 | 0.085 | 0.022 | -0.254 | -0.096 | -0.108 | -0.032 |

**Message Size: 1048576 B**

| Pattern | 1:Lin | 2:Chain | 3:Pipe | 4:Bin | 5:Binom | 6:In-Bin | 7:Raben |
|---|---|---|---|---|---|---|---|
| Ascending | 0.021 | -0.300 | -0.032 | -0.025 | -0.286 | -0.028 | 1.169 |
| Descending | 0.005 | -0.160 | -0.574 | 0.011 | -0.025 | -0.065 | 0.298 |
| First-delayed | 0.016 | -0.178 | -0.790 | -0.651 | -0.133 | -0.407 | -0.158 |
| Mid-delayed | -0.490 | -0.374 | -0.265 | -0.559 | -0.181 | -0.464 | 0.552 |
| Last-delayed | 0.008 | -0.299 | -0.027 | -0.530 | -0.470 | -0.676 | 0.455 |
| Random | -0.055 | -0.102 | -0.127 | 0.017 | -0.089 | -0.192 | 0.033 |

### MPI_Alltoall

**Message Size: 8 B**

| Pattern | 1:Lin | 2:Pair | 3:M-Bruck | 4:L-Sync |
|---|---|---|---|---|
| Ascending | -0.053 | -0.019 | -0.085 | -0.052 |
| Descending | -0.016 | -0.004 | -0.086 | -0.022 |
| First-delayed | -0.197 | -0.017 | -0.110 | -0.186 |
| Mid-delayed | -0.205 | -0.011 | -0.125 | -0.179 |
| Last-delayed | -0.226 | -0.005 | -0.133 | -0.206 |
| Random | -0.130 | -0.010 | -0.063 | -0.128 |

**Message Size: 1024 B**

| Pattern | 1:Lin | 2:Pair | 3:M-Bruck | 4:L-Sync |
|---|---|---|---|---|
| Ascending | -0.027 | -0.018 | -0.168 | -0.984 |
| Descending | -0.021 | -0.032 | -0.208 | -0.984 |
| First-delayed | 0.943 | -0.022 | -0.239 | -0.011 |
| Mid-delayed | 0.943 | -0.029 | -0.243 | -0.011 |
| Last-delayed | -0.004 | -0.015 | -0.258 | -0.011 |
| Random | -0.985 | -0.035 | -0.108 | -0.987 |

**Message Size: 1048576 B**

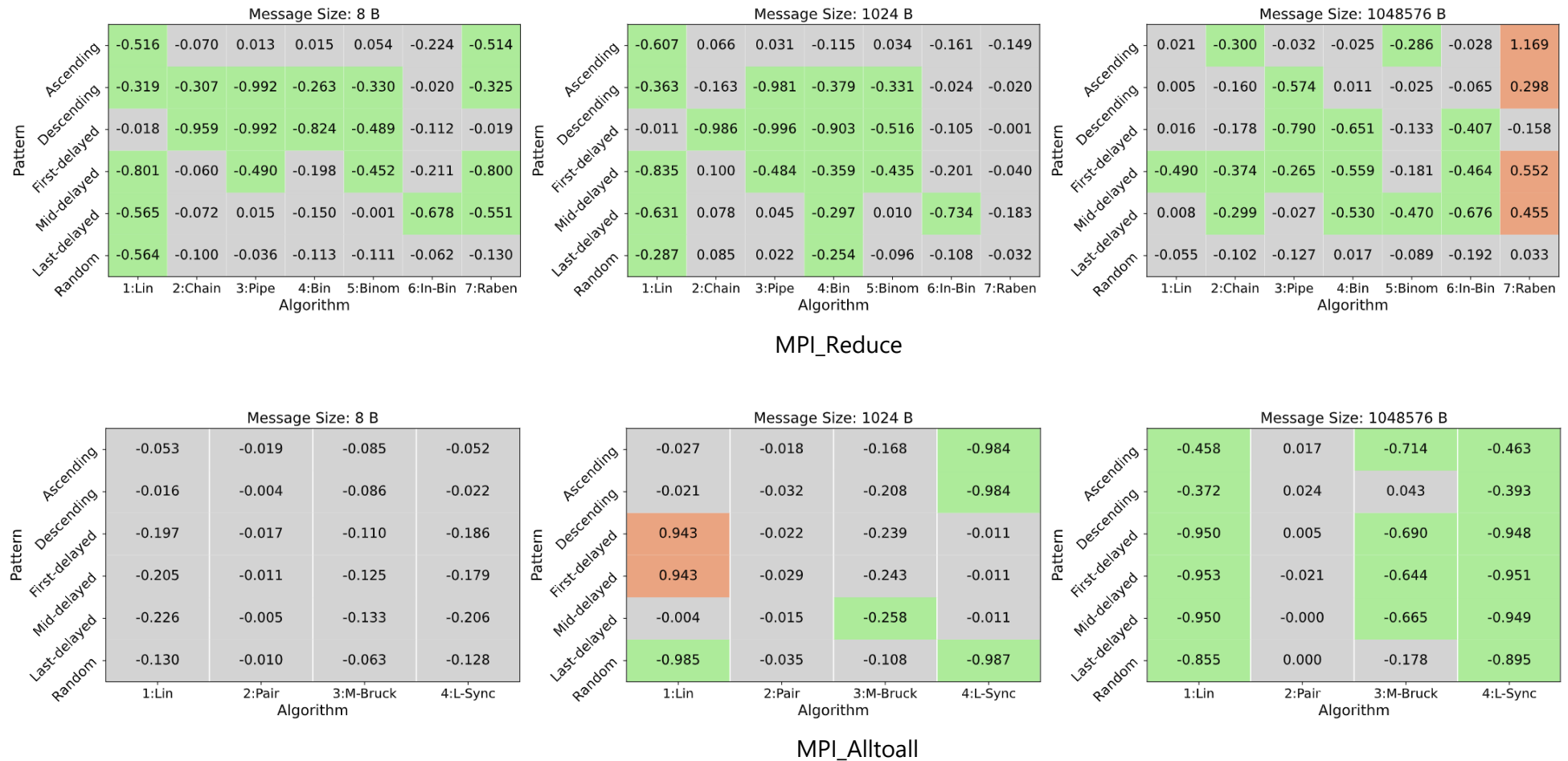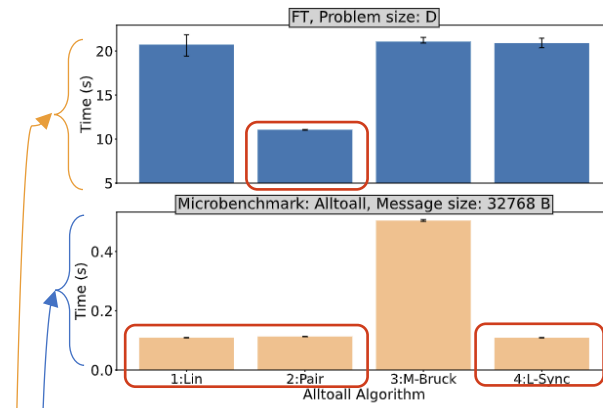| Pattern | 1:Lin | 2:Pair | 3:M-Bruck | 4:L-Sync |
|---|---|---|---|---|
| Ascending | -0.458 | 0.017 | -0.714 | -0.463 |
| Descending | -0.372 | 0.024 | 0.043 | -0.393 |
| First-delayed | -0.950 | 0.005 | -0.690 | -0.948 |
| Mid-delayed | -0.953 | -0.021 | -0.644 | -0.951 |
| Last-delayed | -0.950 | -0.000 | -0.665 | -0.949 |
| Random | -0.855 | 0.000 | -0.178 | -0.895 |

Fig: Normalized runtimes to No-delay case on Hydra
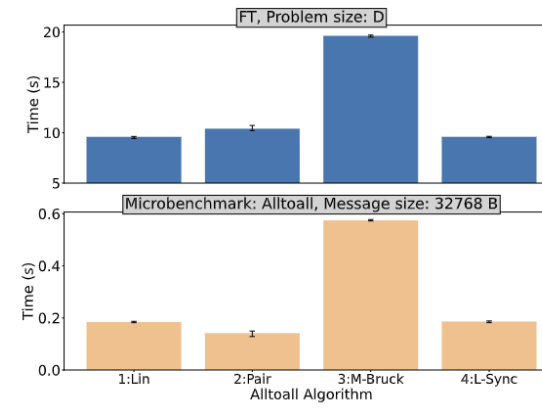
# Arrival Patterns in Applications

- ❑ FT (problem size D) from NAS Parallel Benchmark

- ❑ FT-Scenario: Real-world

  - ❑ Enables us to accurately predict the best performing algorithm

- ❑ Selection strategy: average is a good indicator

- ❑ An algorithm that **consistently performs well across multiple arrival patterns** will likely yield satisfactory results across various applications.



(a) *Hydra*  (b) *Galileo100*  (c) *Discoverer*

# Arrival Patterns in Applications (Cont'd)

❑ Expected FT Runtime, based on the **No-delay case**, does not align with the Actual FT Runtime.

❑ Expected FT Runtime, based on the **Average case**, aligns well with the Actual FT Runtime.

❑ The behavior of the collective algorithm in the application can be predicted!



(a) *Hydra*

Fig: The actual runtime of FT versus its projected runtimes (when processes enter collectives simultaneously, the No-delay case, and the average case) on Hydra with 32 × 32 processes.

# Conclusion and Future Work

- ❑ MPI collective algorithm selection problem

- ❑ Impact of arrival patterns on collective algorithms

- ❑ Micro-benchmarking strategy

  - ❑ Simulation study

  - ❑ 3 real-world production machines

- ❑ * Rooted collectives, such as MPI_Reduce, are more influenced

- ❑ * Algorithm selection without considering the process imbalance may lead to an inefficient choice

- ❑ * Considering robustness

- ❑ Future Directions

  - ❑ Studying more complicated applications

  - ❑ Studying arrival patterns on GPU clusters

THANK YOU

MPI Collective Algorithm Selection
in the Presence of Process Arrival Patterns

Majid Salimi Beni, Biagio Cosenza, Sascha Hunold

2024 IEEE International Conference on Cluster Computing
(CLUSTER)

Kobe, Japan
September 24-27, 2024

✉ Reach me at:
msalimibeni@unisa.it