# *Experiences combining malleability and I/O control mechanisms*

David E. Singh and Jesús Carretero

University Carlos III of Madrid (Spain)

▶ I/O challenges

- ▶ Scientific applications (climate, genomics, high energy physics, astronomy etc.) ingest, generate, and process increasingly larger data sets

- ▶ Future high scale supercomputers need to deal efficiently with huge amounts of data

- ▶ Current I/O software stack needs to evolve in order to meet the oncoming scalability challenges

▶ CPU challenges

- ▶ Malleable applications can leverage unused computational resources

- Concurrent parallel data flows
  - Lack of data staging coordination
    - Among applications
    - Between applications and the system

- Lack of standards for dynamic monitoring of large scale infrastructures

- Need of coupled control and data mechanisms

- Lack of coordination with the job scheduler

- Integration of CLARISSE and FlexMPI into a framework
  - New coordination techniques between the applications and the scheduler
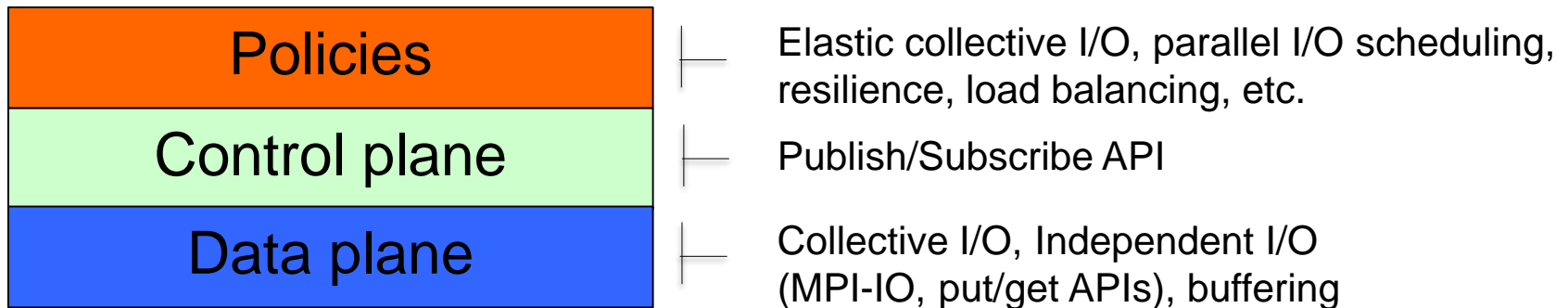  - Application monitoring

- Use of application malleability to enhance the I/O performance:

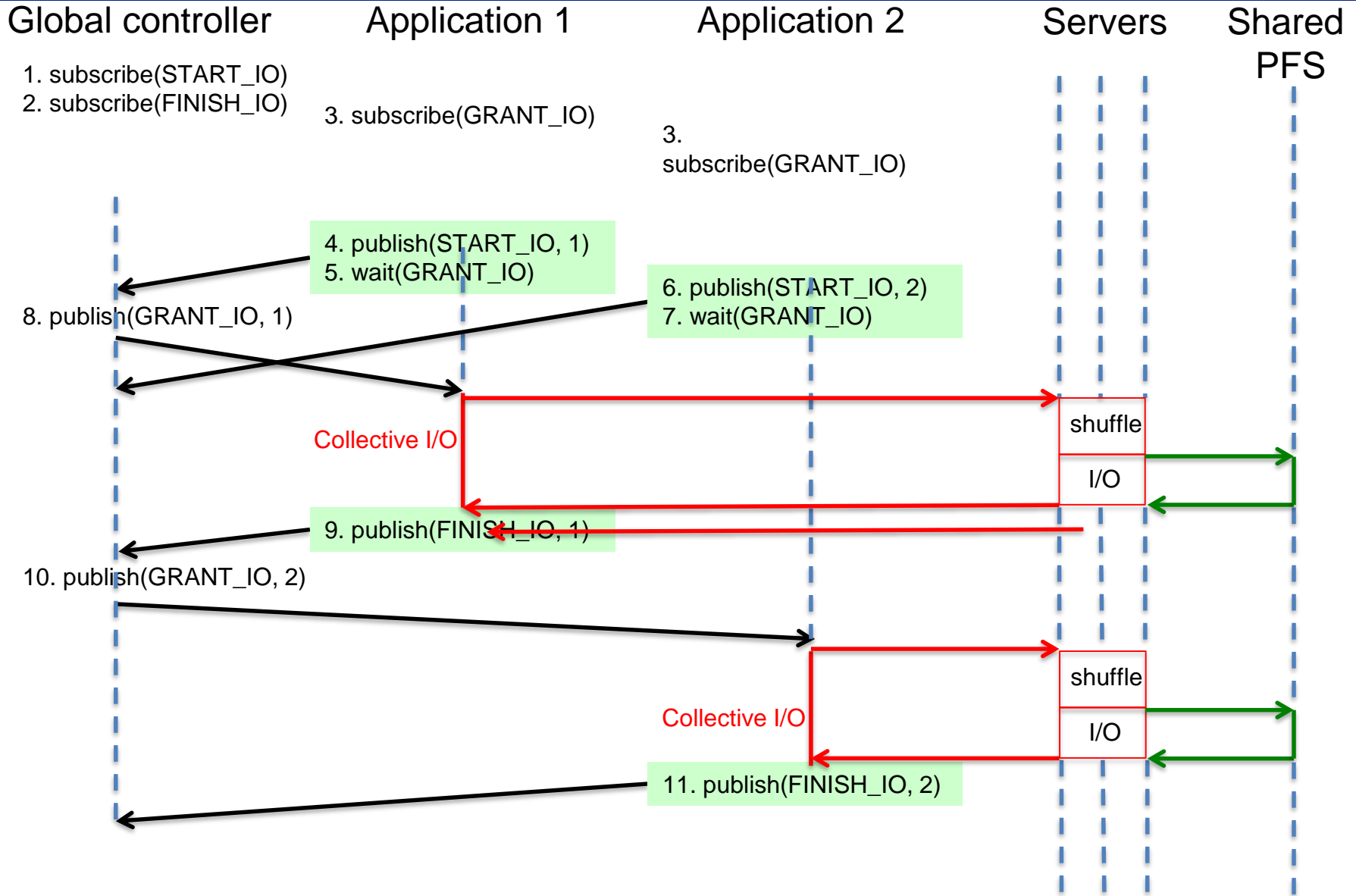  - Coordinated use of parallel I/O scheduling and malleability for reducing number of I/O interferences

    **I/O interference**: *two or more I/O operations that occur partially or totally at the same time competing for the I/O resources*

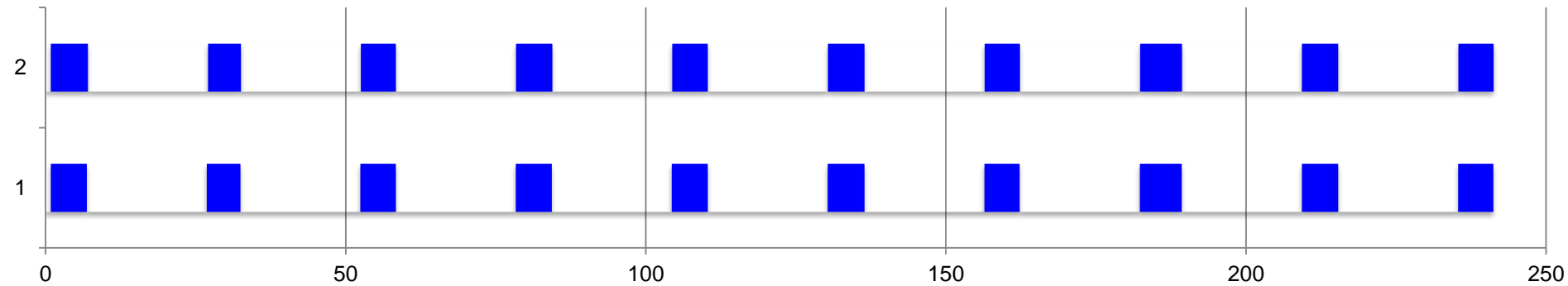  - I/O-aware scheduling policies

- ▸ Novel mechanisms for global data staging coordination to improve:
    - ▸ Load balance, resilience, parallel I/O scheduling, locality exploitation

- ▸ Decouple the data and control planes
    - ▸ Data plane
    - ▸ Control plane
    - ▸ Policy

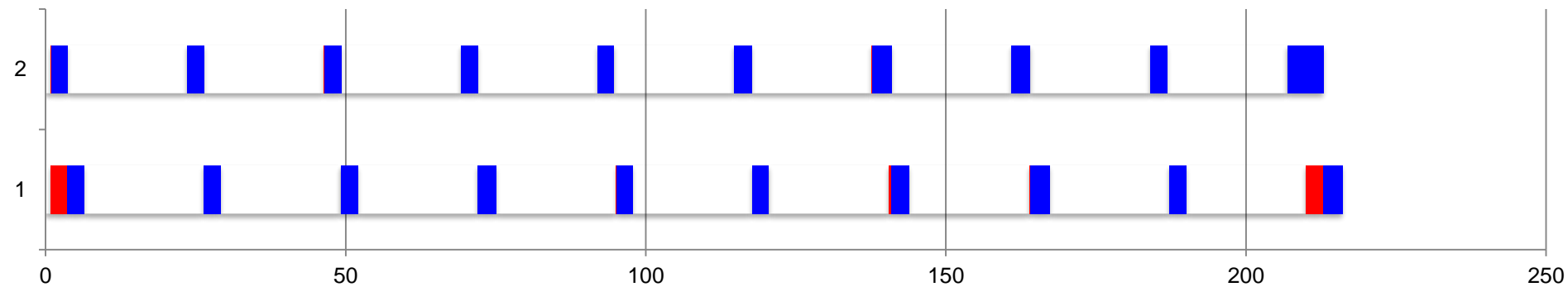- ▸ Facilitate the flow of control and data across the I/O stack

| Policies | Elastic collective I/O, parallel I/O scheduling, resilience, load balancing, etc. |
| Control plane | Publish/Subscribe API |
| Data plane | Collective I/O, Independent I/O (MPI-IO, put/get APIs), buffering |

Universidad Carlos III de Madrid
www.uc3m.es

ARCOS

| Global controller | Application 1 | Application 2 | Servers | Shared PFS |
|---|---|---|---|---|

1. subscribe(START_IO)
2. subscribe(FINISH_IO)

3. subscribe(GRANT_IO)

3. subscribe(GRANT_IO)

4. publish(START_IO, 1)
5. wait(GRANT_IO)

6. publish(START_IO, 2)
7. wait(GRANT_IO)

8. publish(GRANT_IO, 1)

Collective I/O

shuffle

I/O

9. publish(FINISH_IO, 1)

10. publish(GRANT_IO, 2)

shuffle

I/O

Collective I/O

11. publish(FINISH_IO, 2)

Write timeline for two parallel clients with 3840 processes each - No scheduling

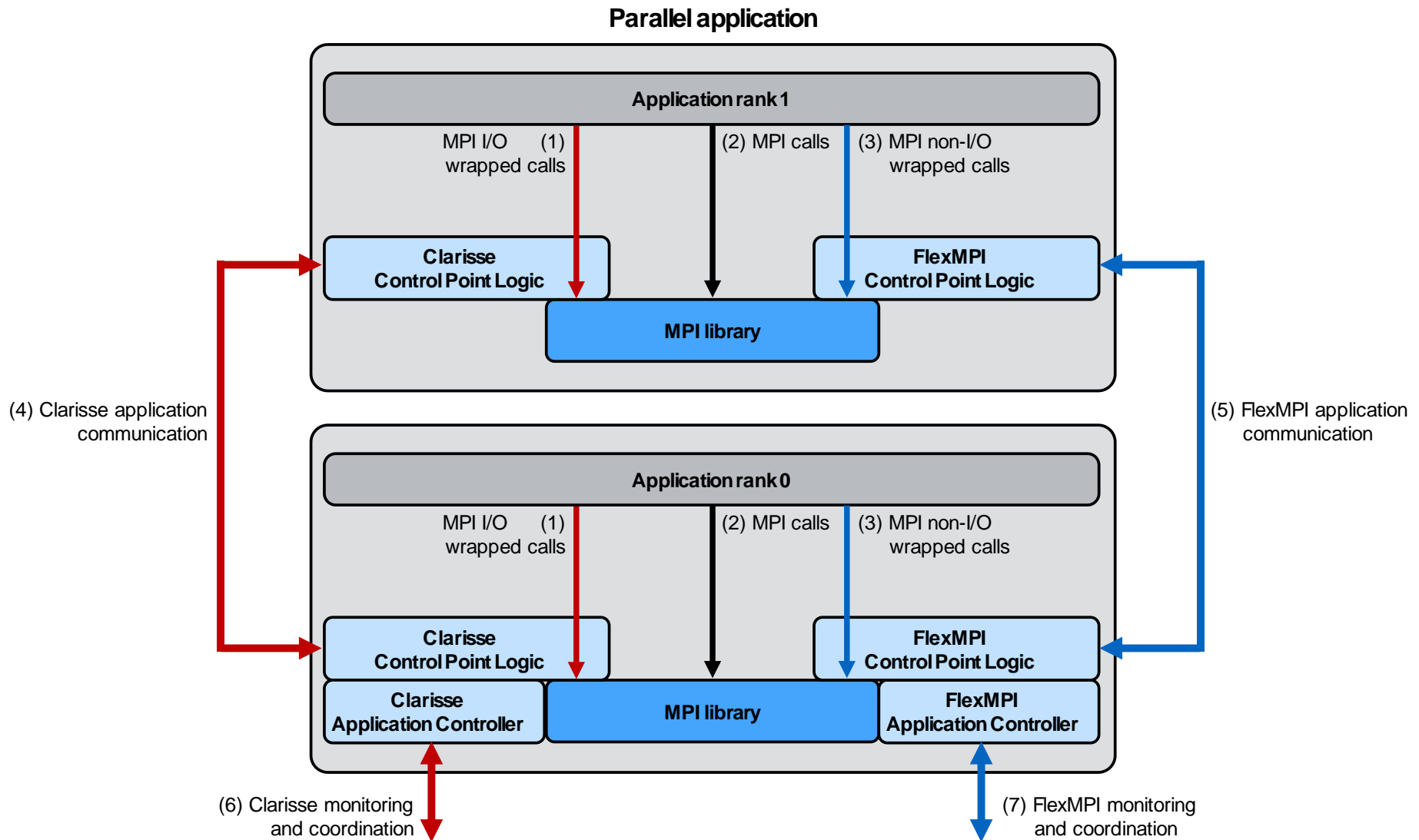Write timeline for two parallel clients with 3840 processes each - FCFS scheduling

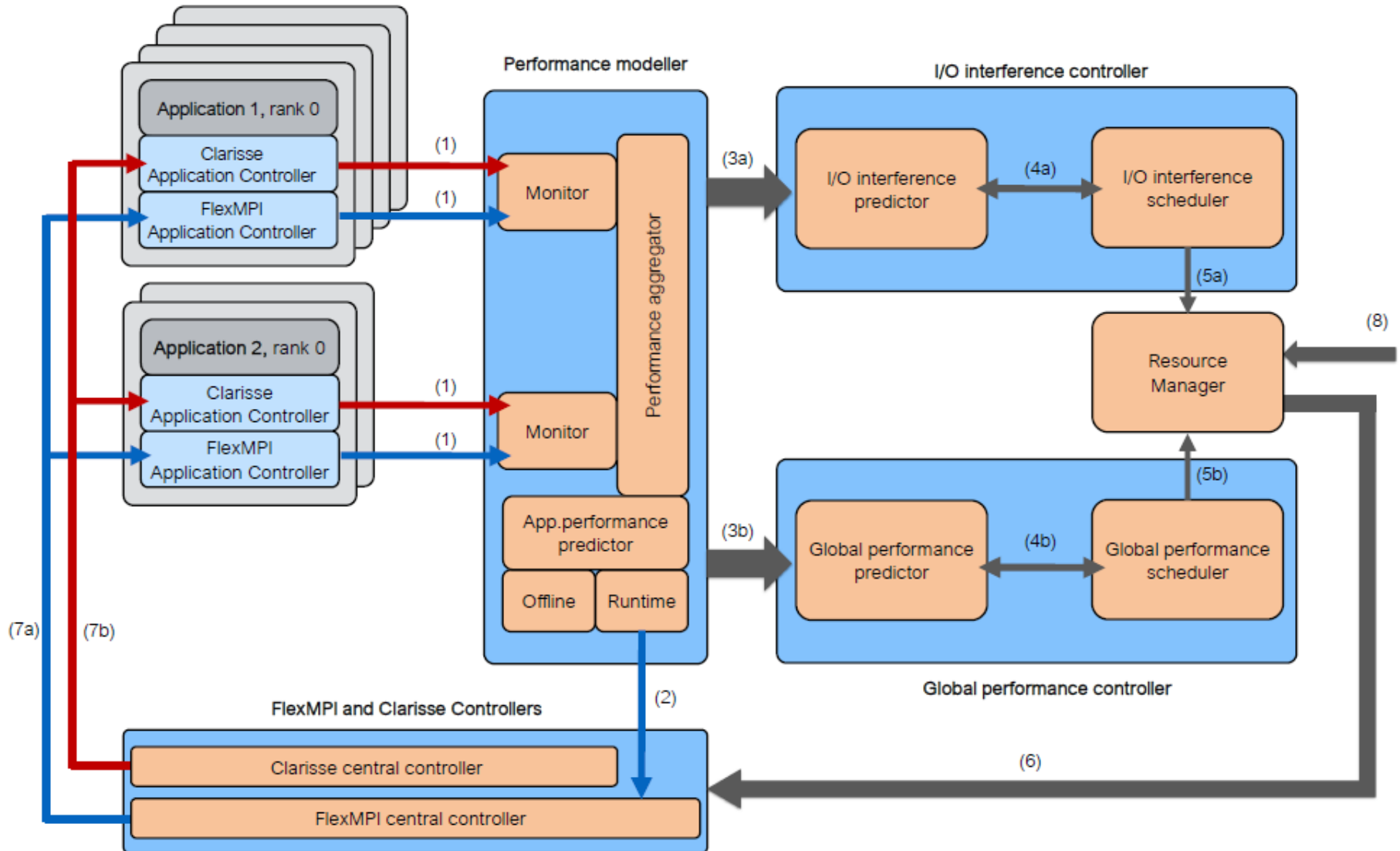▸ FLEX-MPI provides performance-aware malleability capabilities for MPI applications.

▸ Goals:

  ▸ Dynamic application reconfiguration
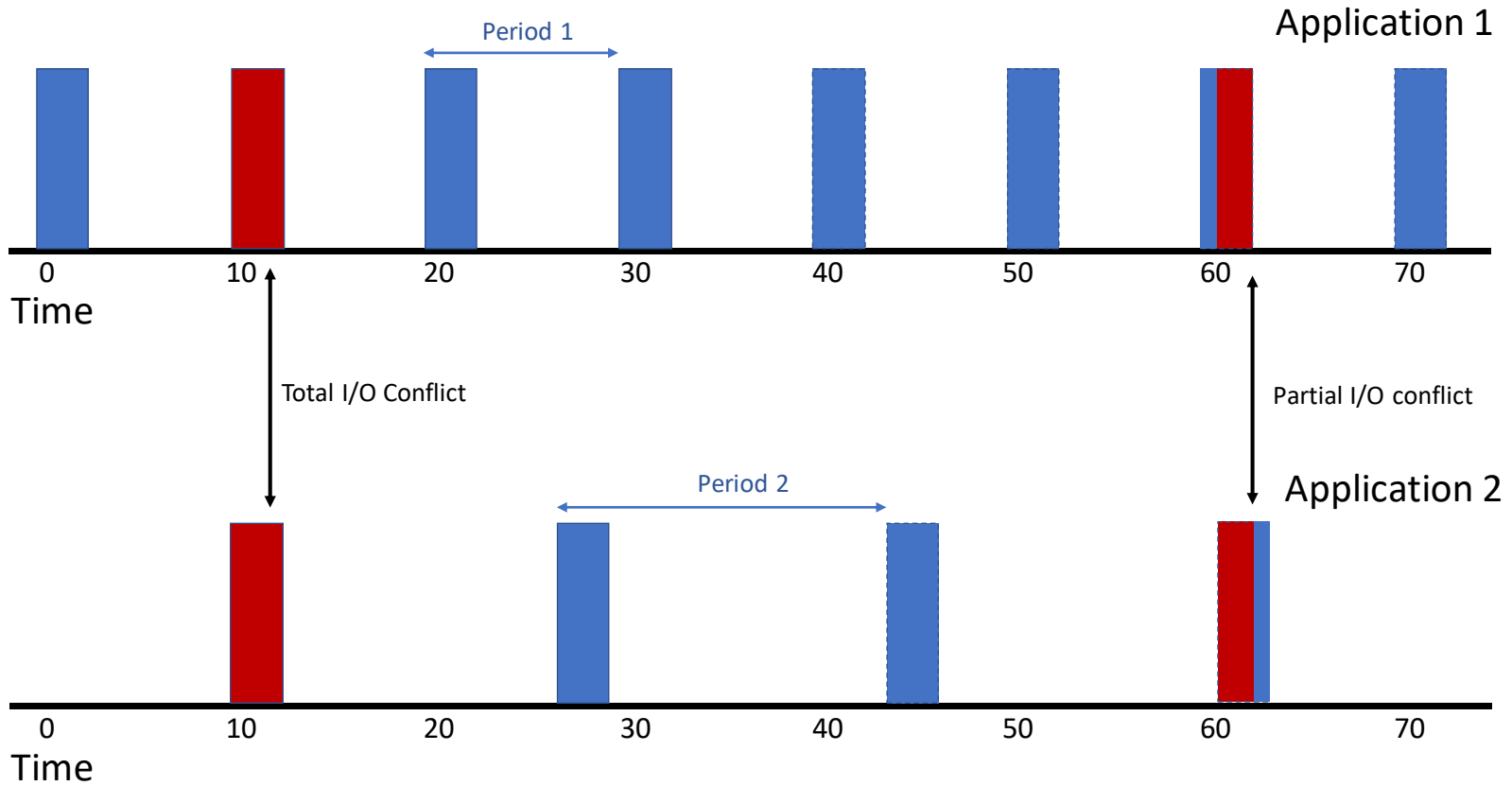
  ▸ Automatic load balancing

  ▸ Monitoring capabilities

# FlexMPI example

- ‣ CLARISSE and FlexMPI integrated at application-level

- ‣ CLARISSE and FlexMPI use separate external controllers

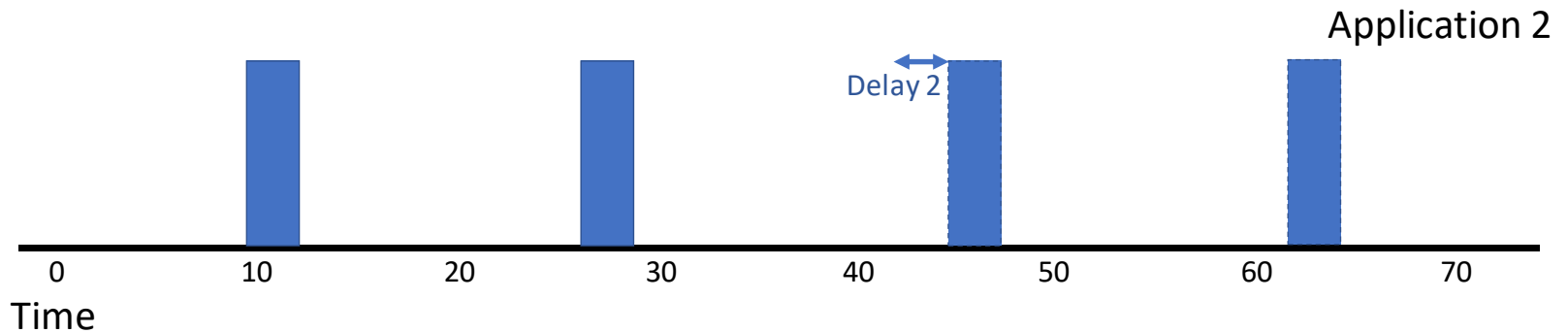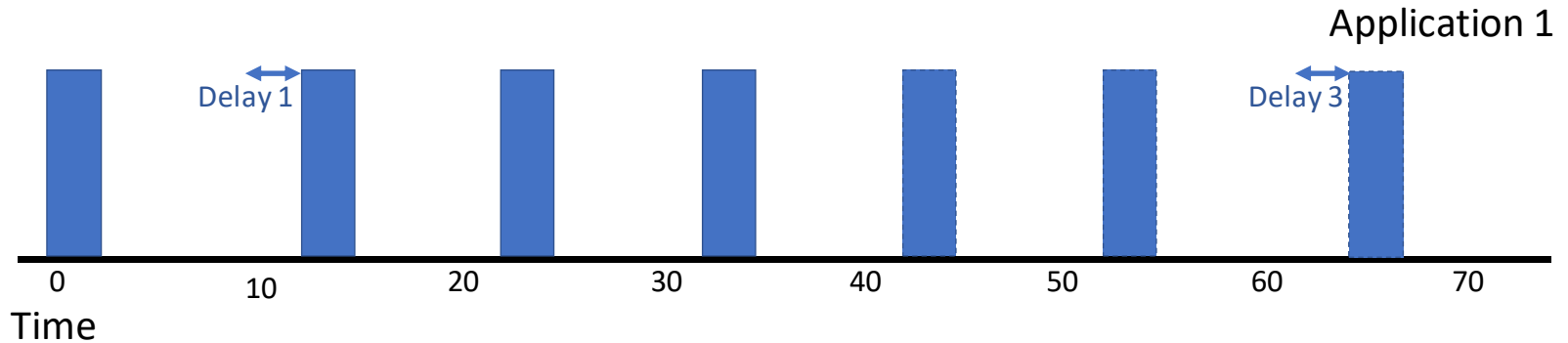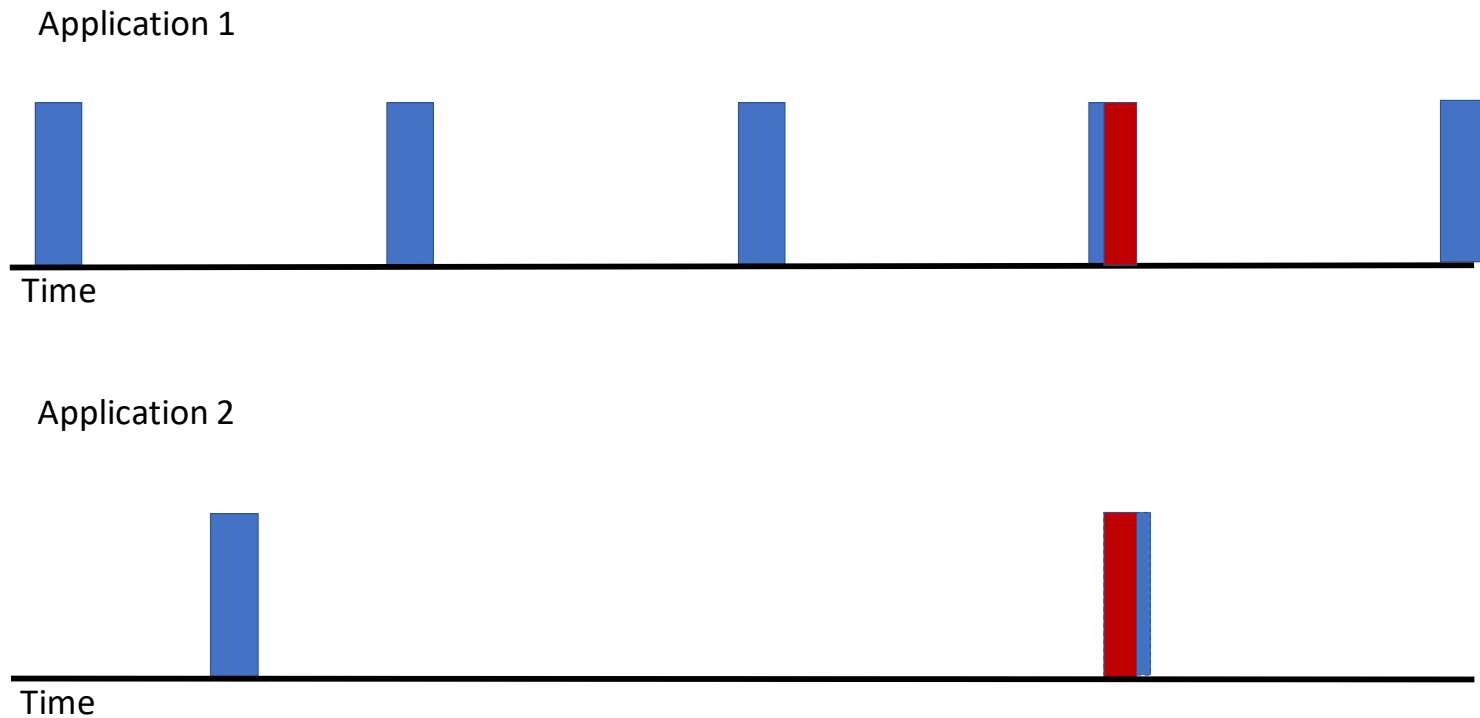- ‣ New control logic coordinates both runtimes

**Parallel application**

Application rank 1

MPI I/O (1) wrapped calls

(2) MPI calls

(3) MPI non-I/O wrapped calls

Clarisse Control Point Logic

FlexMPI Control Point Logic

MPI library

(4) Clarisse application communication

(5) FlexMPI application communication

Application rank 0

MPI I/O (1) wrapped calls

(2) MPI calls

(3) MPI non-I/O wrapped calls

Clarisse Control Point Logic

FlexMPI Control Point Logic

Clarisse Application Controller

MPI library

FlexMPI Application Controller

(6) Clarisse monitoring and coordination

(7) FlexMPI monitoring and coordination

Application 1

Period 1

Time

Total I/O Conflict

Partial I/O conflict

Application 2

Period 2

Time

▸ Solutions:
  ▸ **I/O scheduling**: blocks one I/O operation using publish-subscribe support

Application 1

Delay 1

Delay 3

0    10    20    30    40    50    60    70
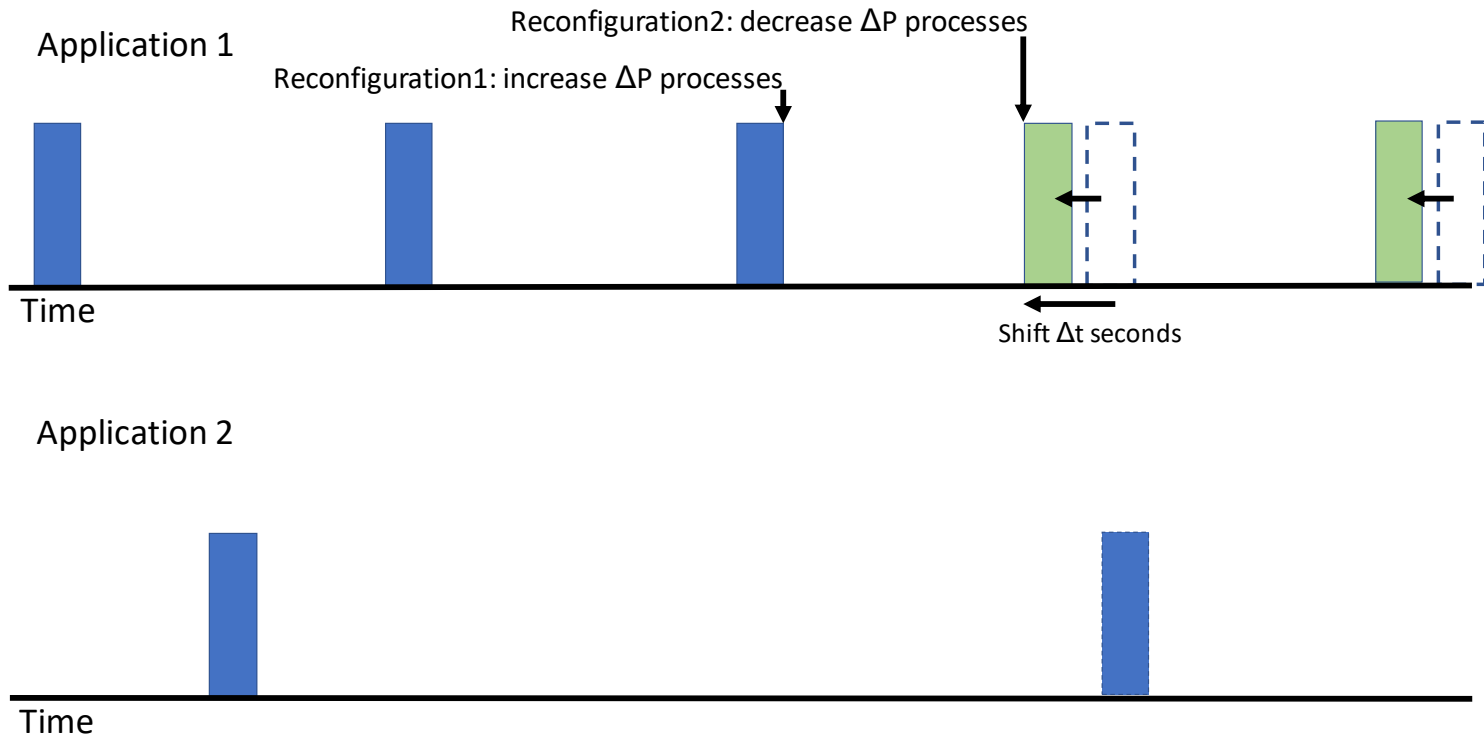
Time

Application 2

Delay 2

0    10    20    30    40    50    60    70

Time

- ▸ Avoiding I/O conflicts with Clarisse + FlexMPI
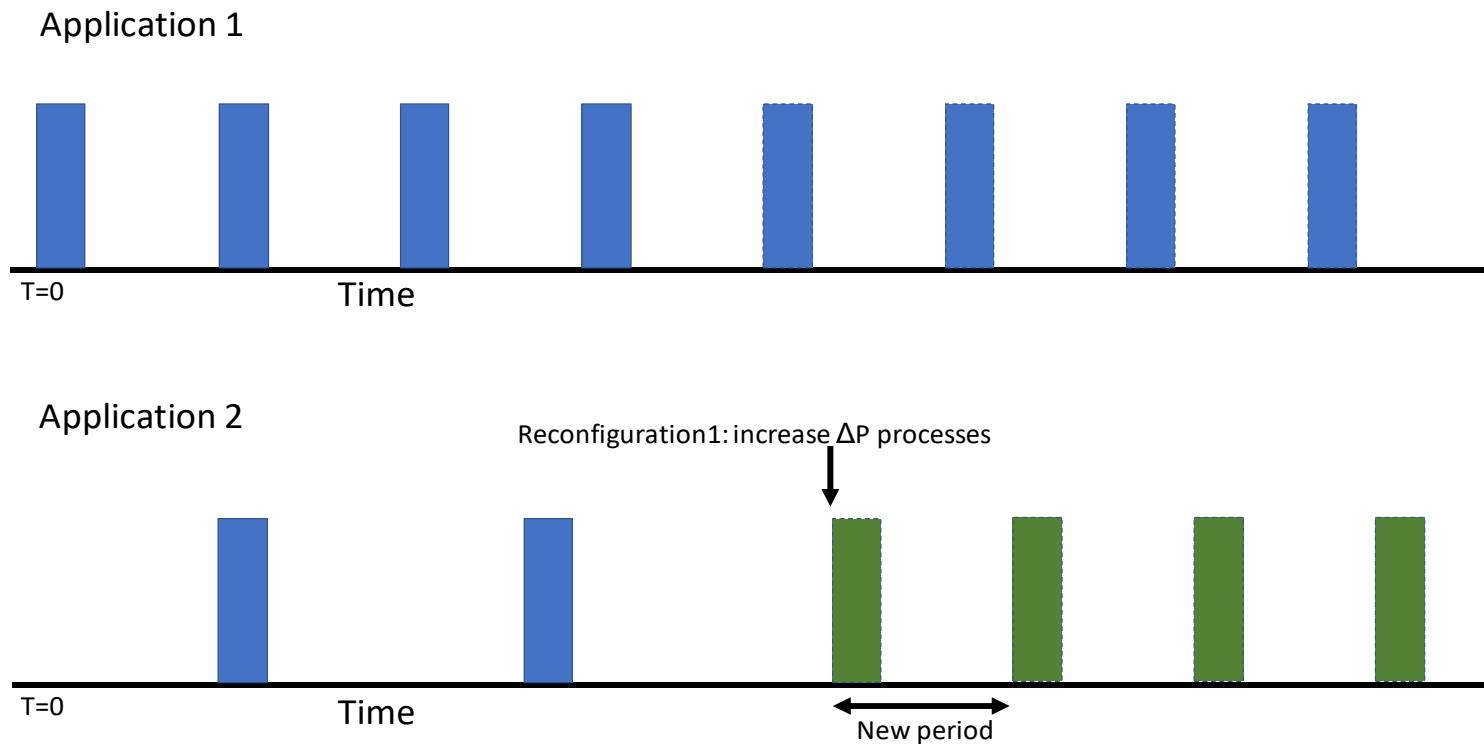  - ▸ Leverage malleability for changing the I/O time stamp
  - ▸ Prediction of the I/O interference

Application 1

Time

Application 2

Time

- **Phase shifting**
  - Leverage malleability for changing the I/O access time (phase)
  - Temporary use of computational resources

▸ **Phase coupling**

  ▸ Leverage malleability for changing the I/O period
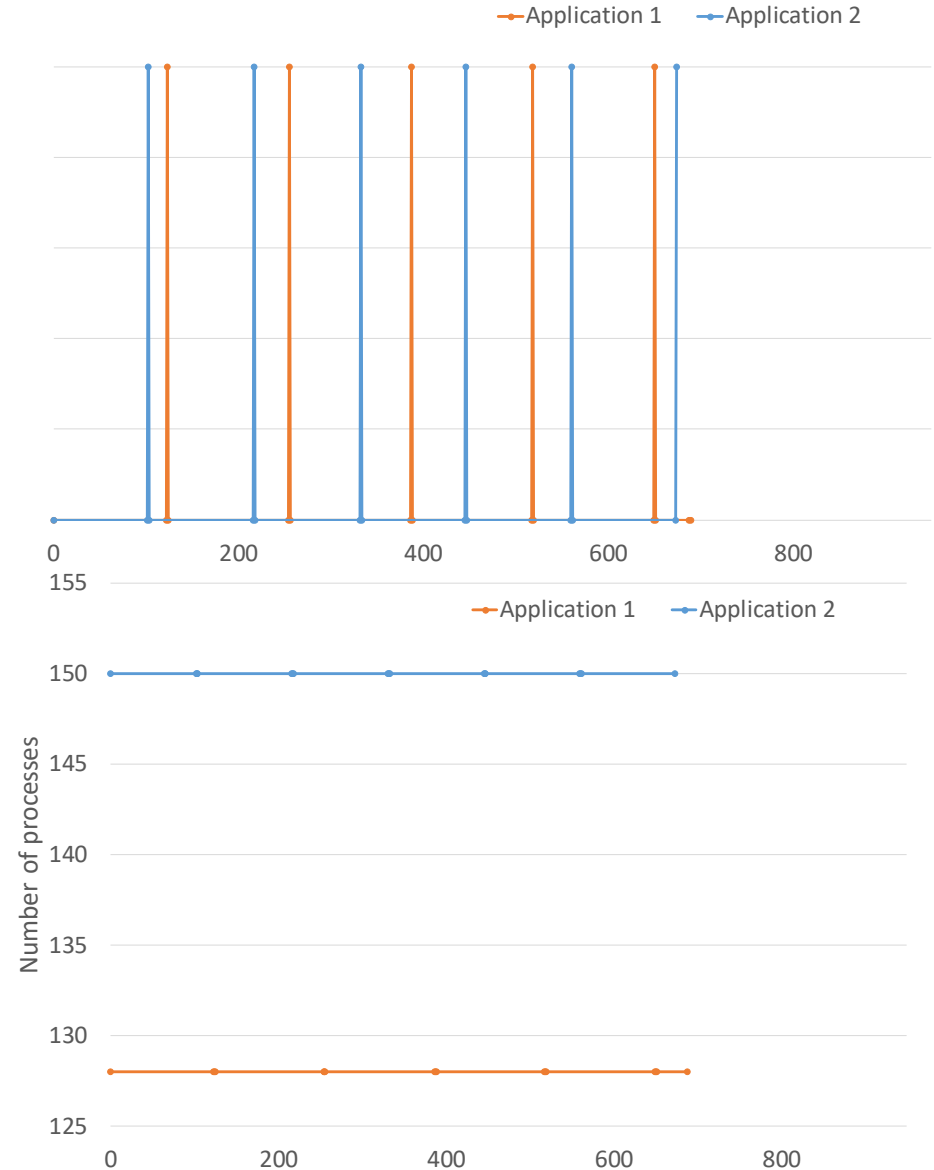
  ▸ Long-term use of computational resources

Application 1

T=0            Time

Application 2

Reconfiguration1: increase ΔP processes

T=0            Time

New period

## **Results**

- Bebop cluster
- Two applications
- 150 and 128 processes

## Results

- Bebop cluster
- Two applications
- 150 and 128 processes

**Results**

- Bebop cluster
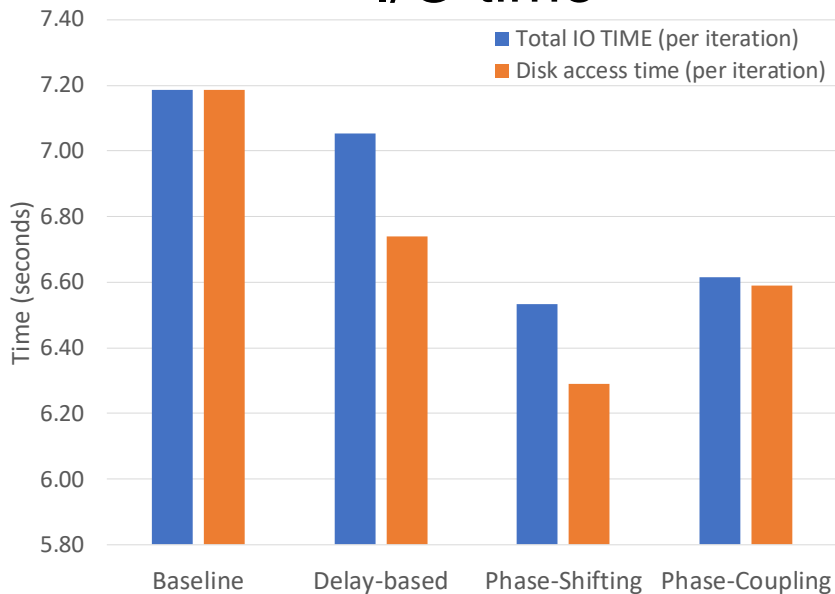- Two applications
- 150 and 128 processes

▸ **Results**

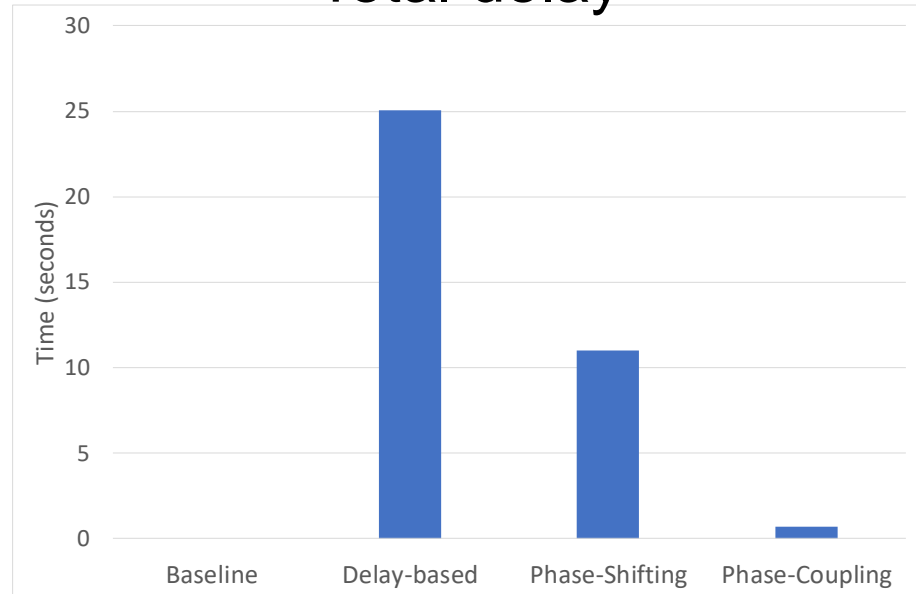  ▸ Two identical applications executed at the same time.
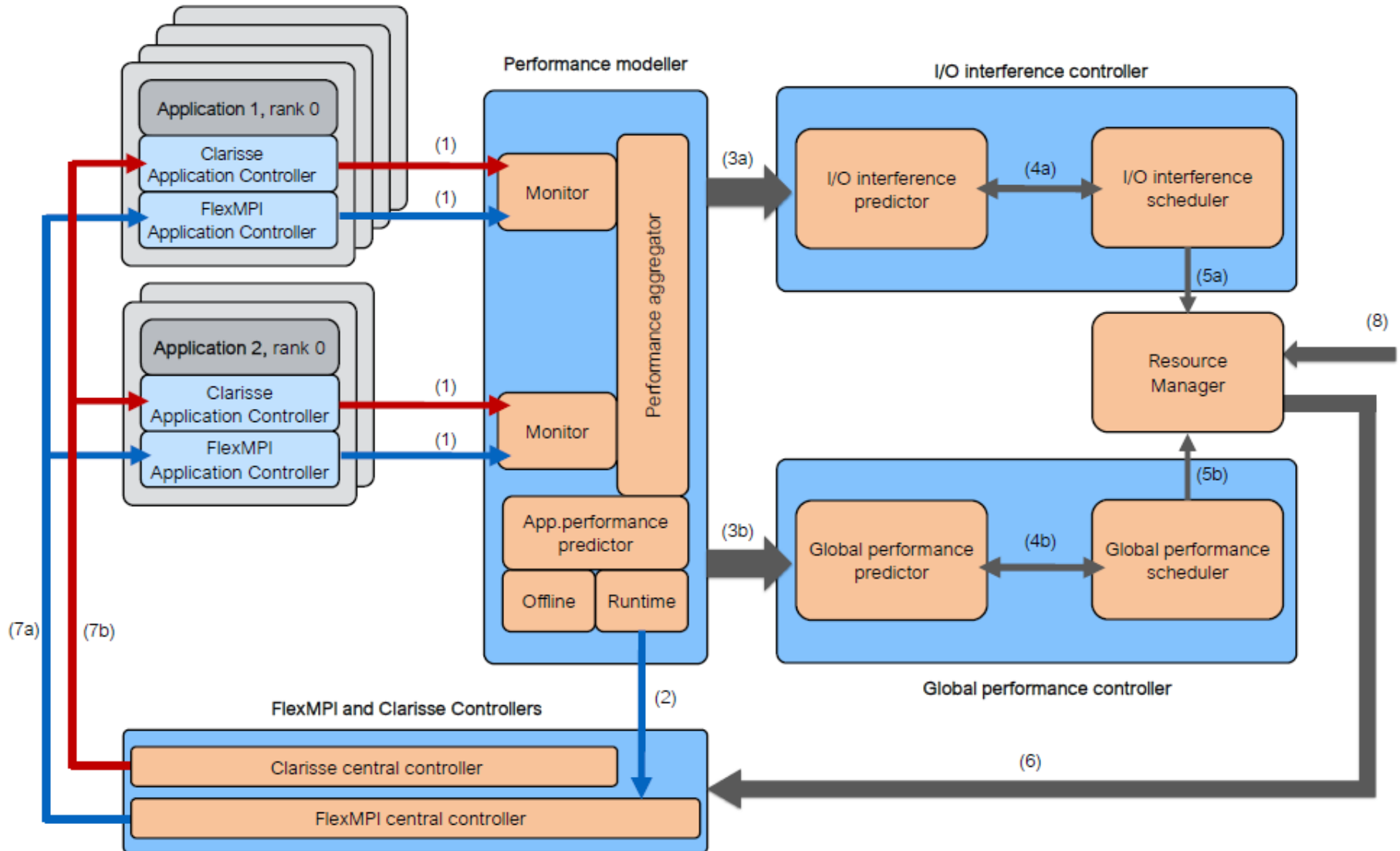
  ▸ 64 processes

## Results

- Two different applications executed at the same time.
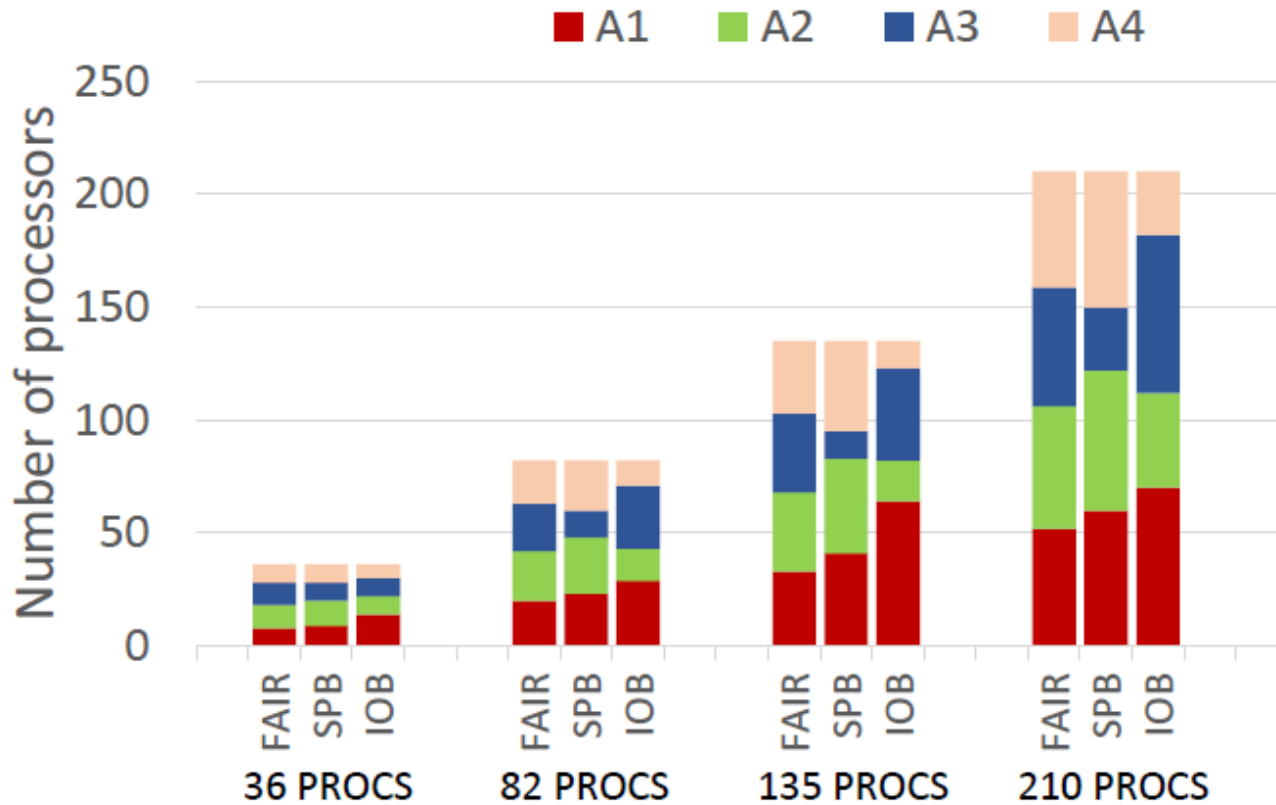- 64 and 50 processes

### I/O time



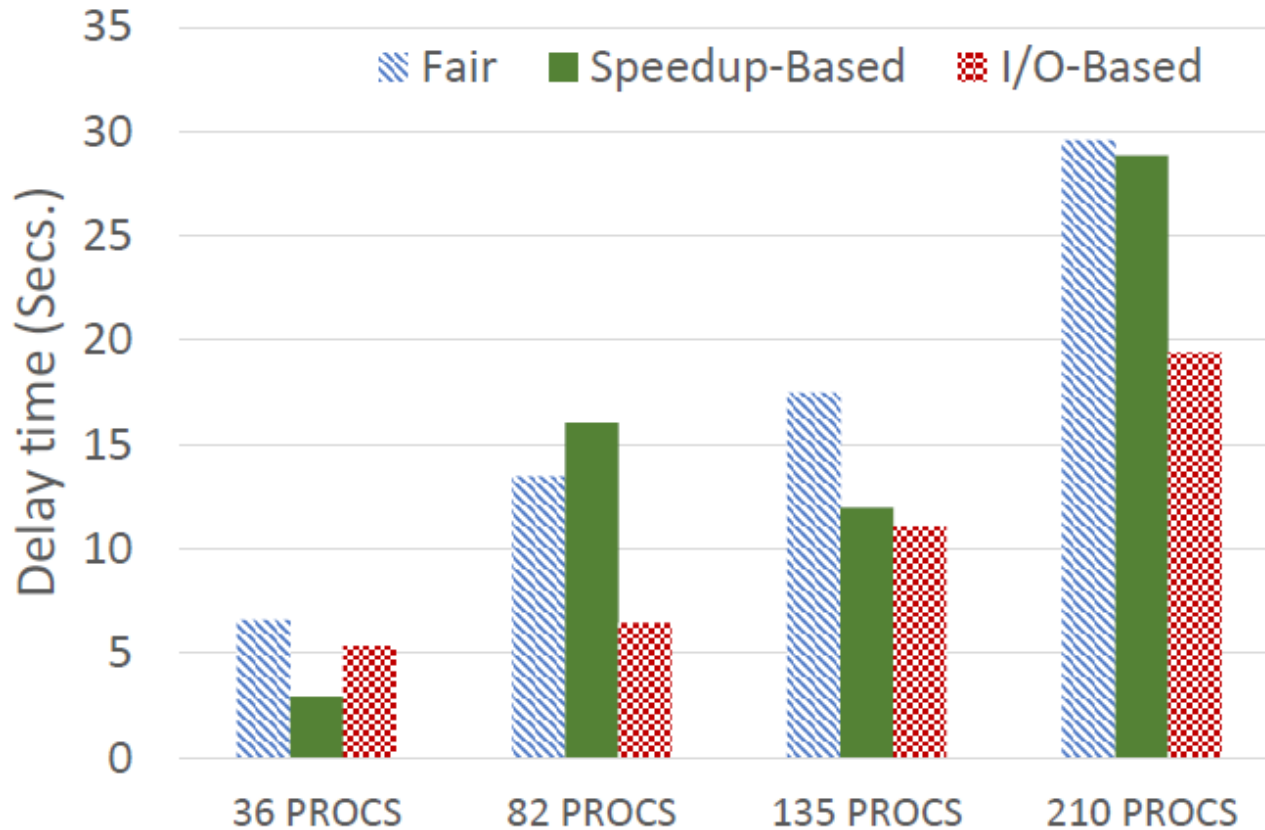### Total delay

▸ We only consider running applications

▸ Assign the available processors to the running applications

▸ Two baseline schedulers:
  ▸ Fair
  ▸ Speedup-based
  ▸ I/O-aware

- ▸ Extended Clarisse and FlexMPI coordination

- ▸ Automated learning techniques for I/O scheduling

- ▸ Integration with a system-wide monitoring tool

- ▸ Application modelling