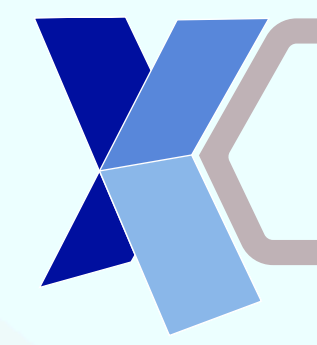




UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA



ICSC

Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

CAPIO

Cross application Programmable
I/O

CAPIO: Cross Application Programmable IO

Scientific HPC in the pre-Exascale era - ITADATA 2024

Marco Edoardo Santimaria

Iacopo Colonnelli

Massimo Torquati

Marco Aldinucci

marcoedoardo.santimaria@unito.it

iacopo.colonnelli@unito.it

massimo.torquati@unipi.it

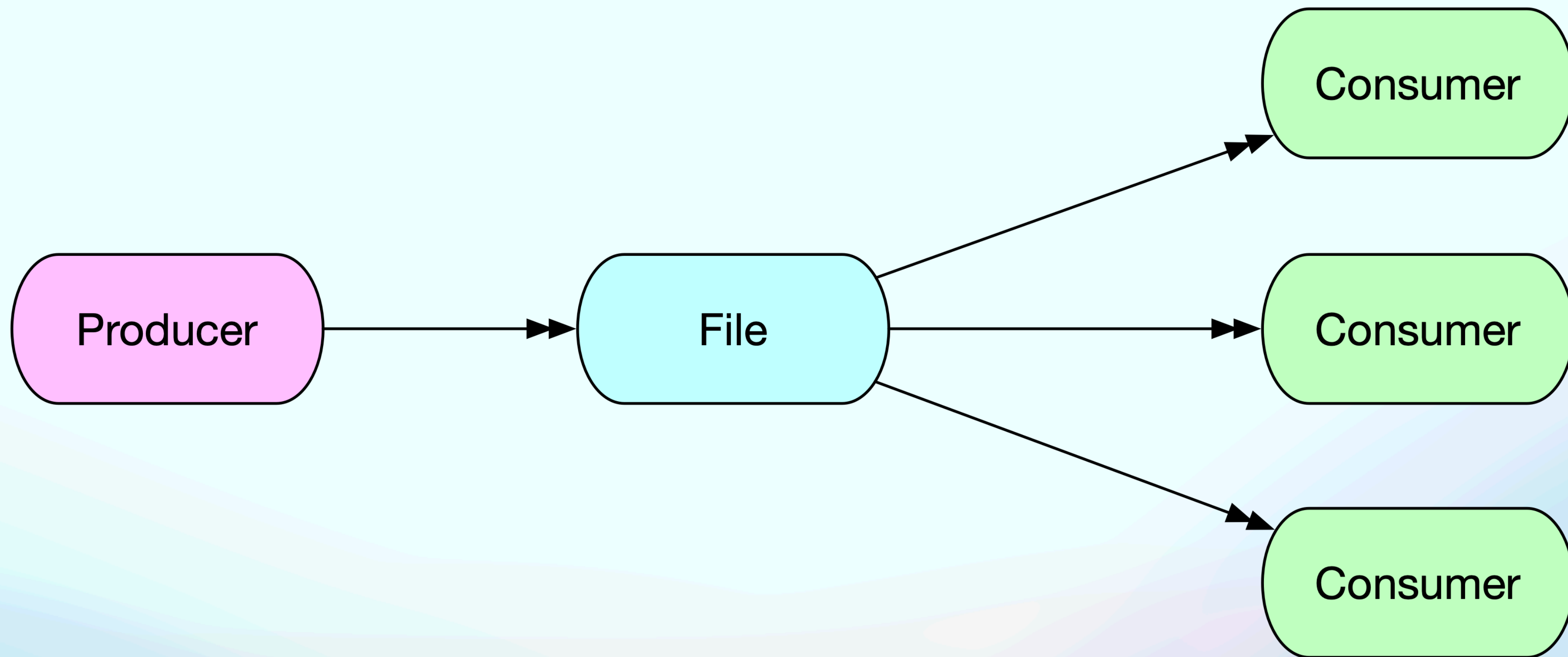
marco.aldinucci@unito.it



Motivations

- The I/O performance gap is still a significant bottleneck in large-scale HPC workflows.
- Parallel file systems struggle to scale linearly with the aggregate disk bandwidth in real scenarios.
- Computation and I/O overlapping are crucial in minimizing I/O overhead.
- Thousands of existing file-based workflows w/o streaming capabilities.

File based workflows



What is CAPIO

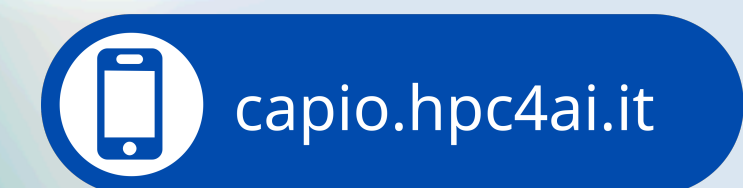


UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA

- User-space data transport middleware enabling I/O coordination in scientific Workflows
- No need to modify the workflow's step
- POSIX I/O System Calls (read/write/seek/stat/...) are transparently intercepted using dynamic linker features (i.e., LD_PRELOAD)



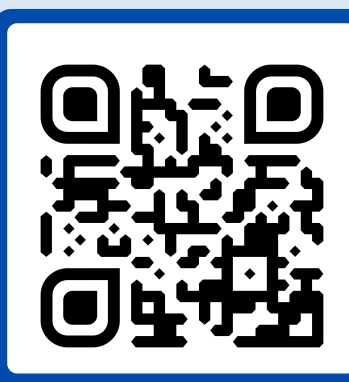
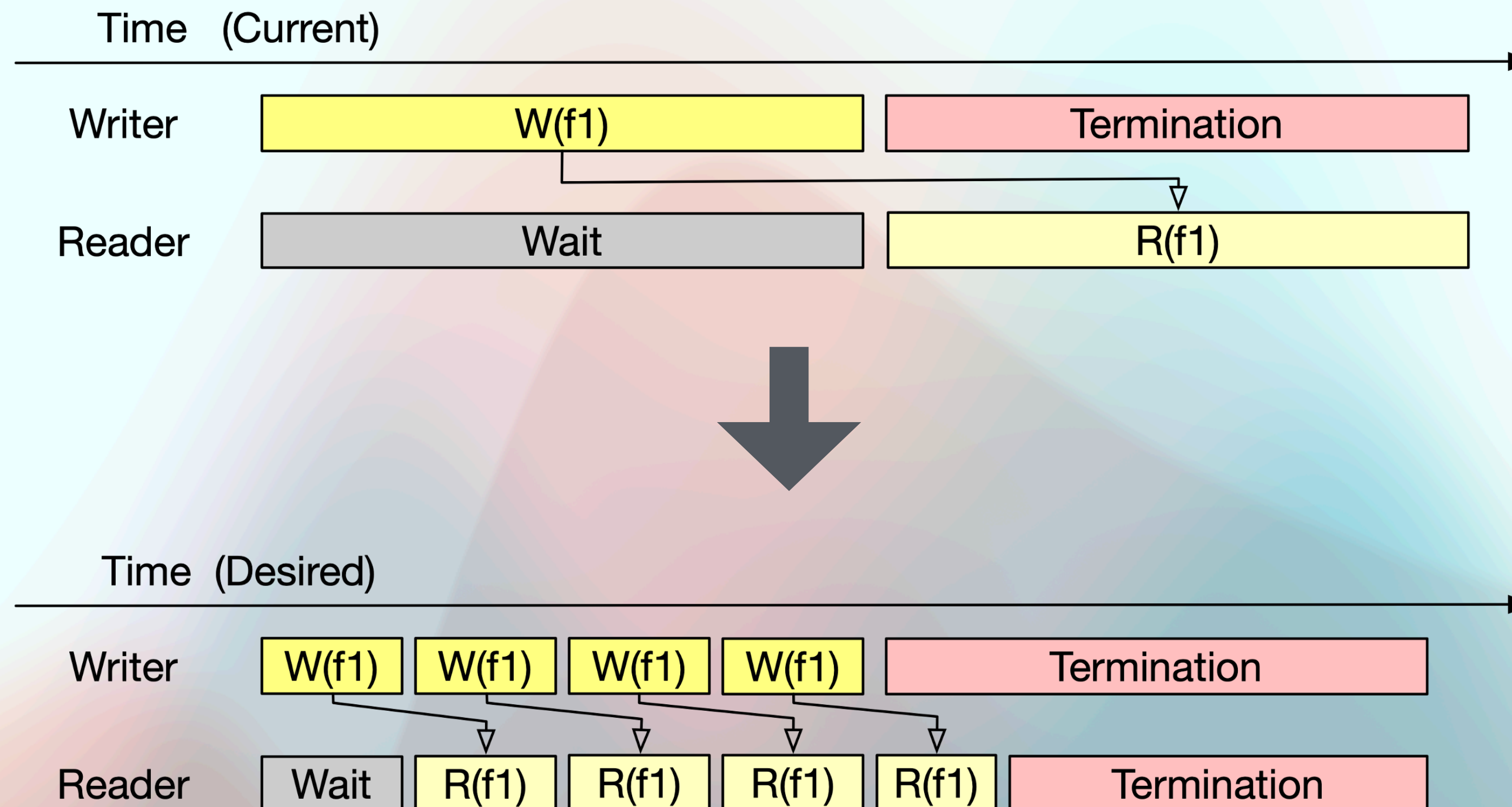
What is CAPIO



UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA



capiro.hpc4ai.it

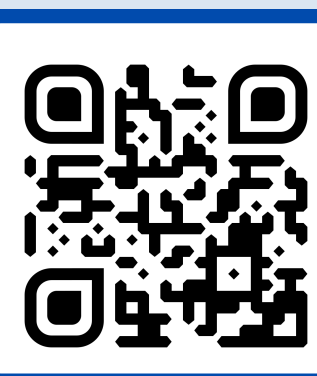
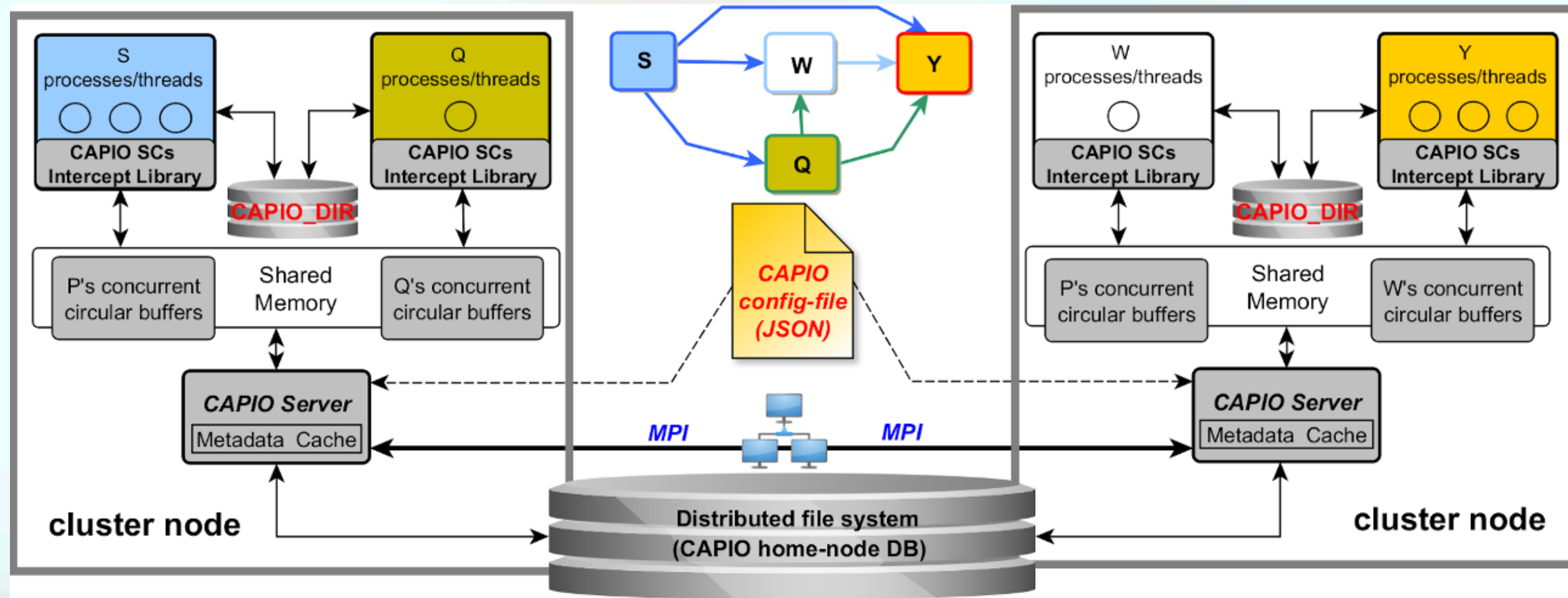
CAPIO Deployment



UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA



capio.hpc4ai.it

What is CAPIO-CL



UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA

- An innovative, fully declarative, I/O coordination language
- Annotates data dependencies within file-based workflows with synchronization semantics
- Allows CAPIO to coordinate streaming injection



 capio.hpc4ai.it

Firing rules

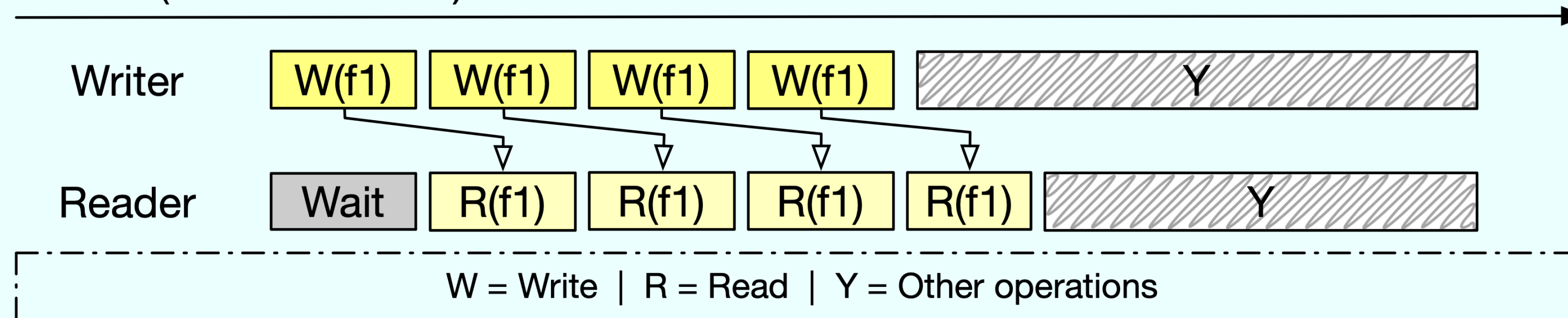


UNIVERSITÀ
DI TORINO

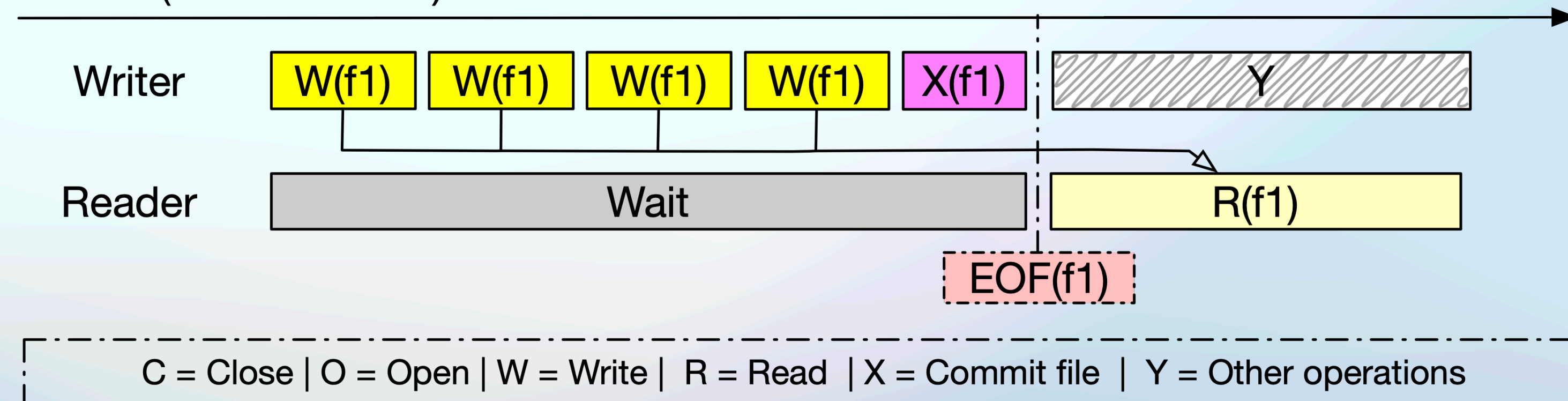


UNIVERSITÀ
DI PISA

Time (FnU semantics)



Time (FoC semantic)



capio.hpc4ai.it



Commit rules

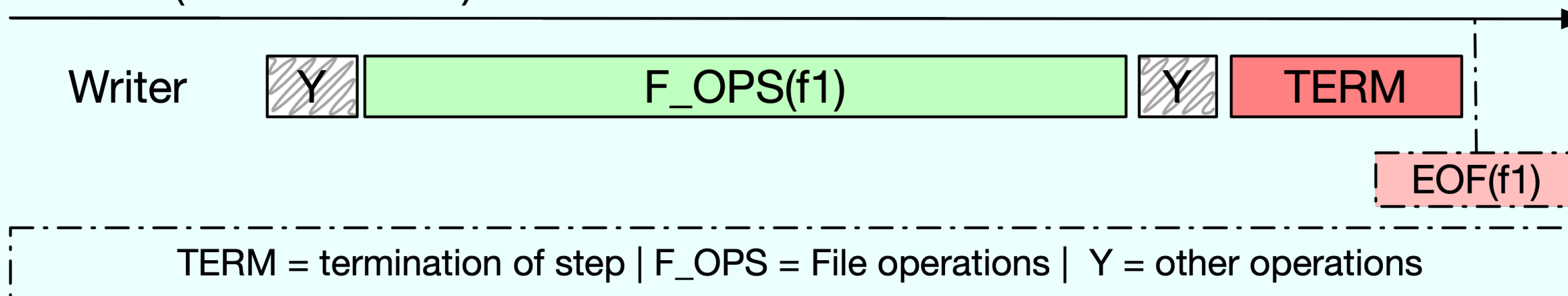


UNIVERSITÀ
DI TORINO

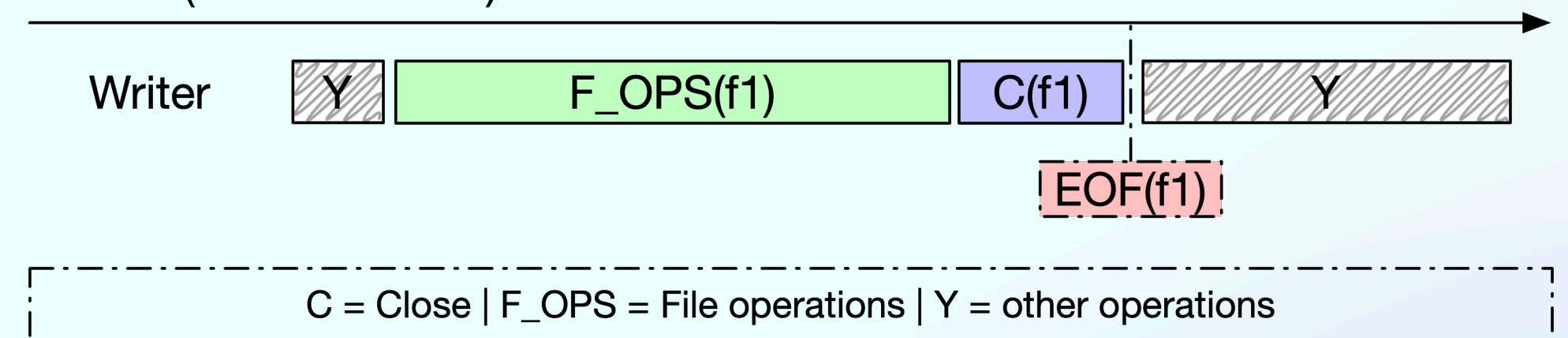


UNIVERSITÀ
DI PISA

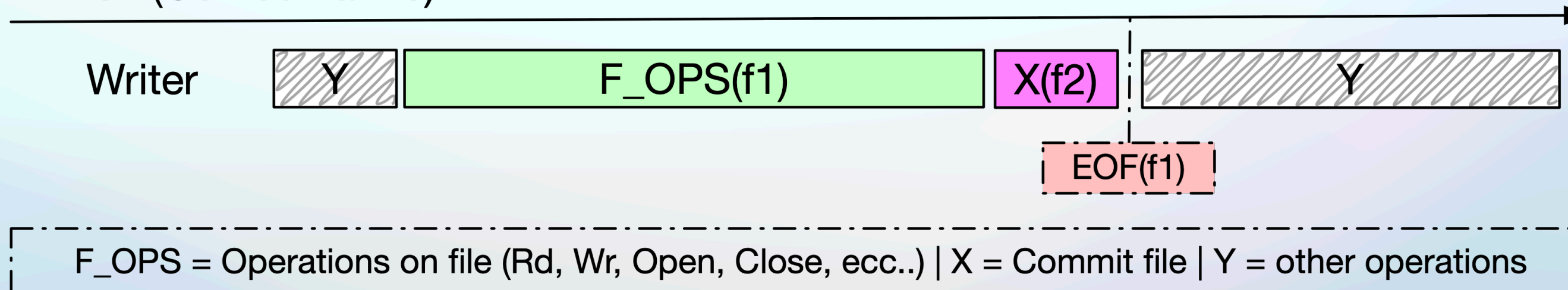
Time (CoT semantic)



Time (CoC semantic)

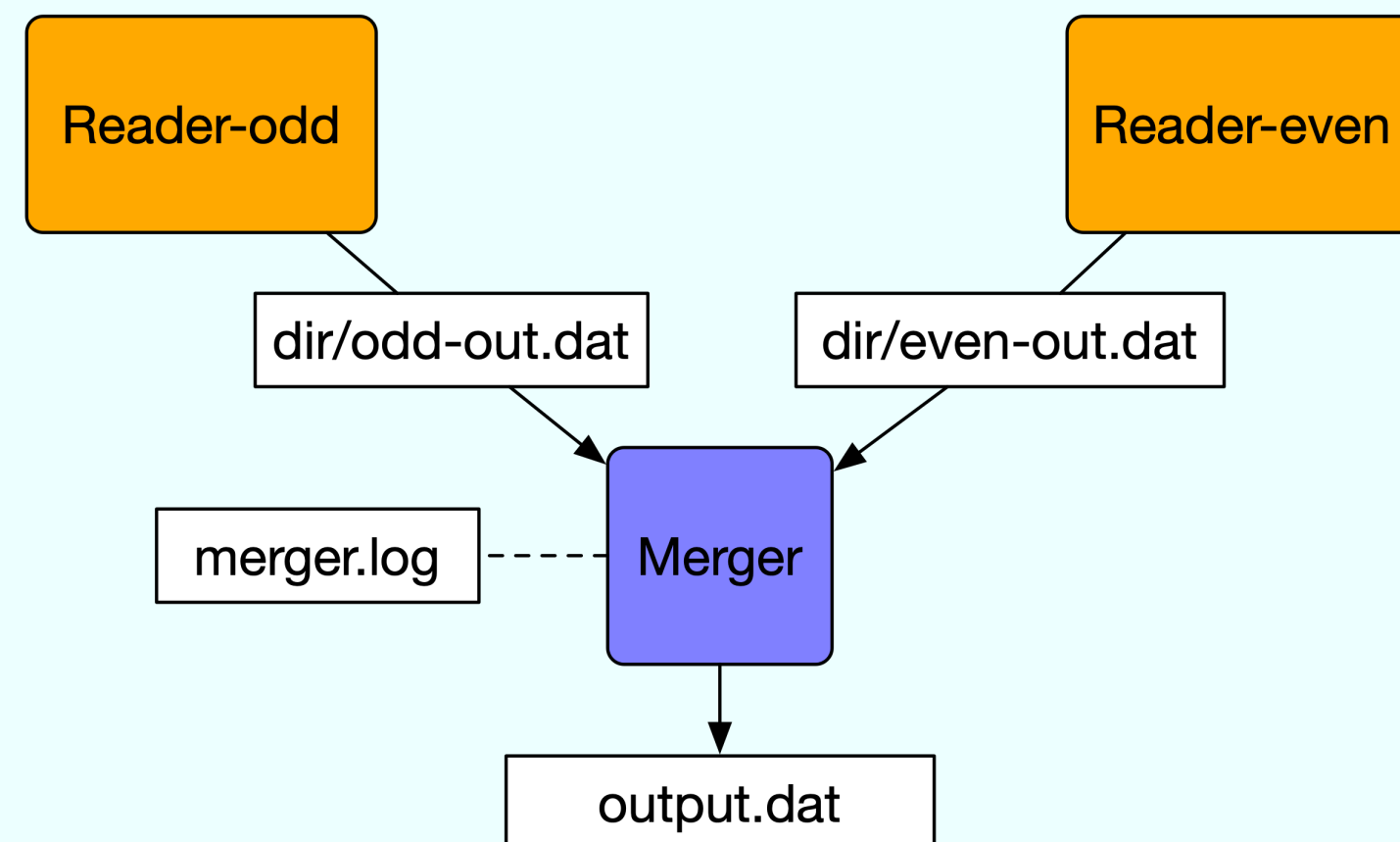


Time (CoF semantic)



capio.hpc4ai.it

Usage example



```
export CAPIO_DIR=$CAPIODIR

srun -N 3 --overlap capio_server -c config.json &

srun -n 1 --export=ALL,LD_PRELOAD=libcapio_posix.so ./Reader-odd &
PID1=$!

srun -n 1 --export=ALL,LD_PRELOAD=libcapio_posix.so ./Reader-even &
PID2=$!

srun -n 1 --export=ALL,LD_PRELOAD=libcapio_posix.so ./Merger

wait $PID1 && wait $PID2
```

```
{
  "name" : "my_workflow",
  "permanent" : [ "output.dat" ],
  "exclude" : [ "merger.log" ],
  "IO_Graph": [
    {
      "name": "reader-even",
      "input_stream": [ "group-even" ],
      "output-stream": [ "even-out.dat" ],
      "streaming" : [{
        "name": [ "even-out.dat" ],
        "committed": "on_close",
        "mode": "no_update"
      }]
    }, {
      "name": "reader-odd",
      "input_stream": [ "group-odd" ],
      "output-stream": [ "odd-out.dat" ],
      "streaming" : [{
        "name": [ "odd-out.dat" ],
        "committed": "on_file:even-out.dat",
        "mode": "update"
      }]
    }, {
      "name": "merger",
      "input_stream": [ "odd-out.dat", "even-out.dat" ],
      "output_stream" : [ "output.dat" ]
    }
  ]
}
```

CAPIO vs ADIOS2



UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA

```
...
#include <adios2.h>

void __write__(){
    std::ofstream out("test.txt");
    out << "Hello world with ADIOS2";
    out.close();
}

void __read__(){
    std::ifstream in("test.txt");
    std::cout << in.rdbuf();
}

void main(){
    std::thread reader_thread(__read__);
    std::thread writer_thread(__write__);
    reader_thread.join();
    writer_thread.join();
}
```

```
...
#include <adios2.h>

void __write__(){
    std::string msg = "Hello world with ADIOS2";
    adios2::fstream ostream(filenameString, adios2::fstream::out);
    ostream.write<std::string>("data", msg);
    ostream.close();
}

void __read__(){
    adios2::fstream istream("test.bp", adios2::fstream::in);
    for (adios2::fstep istep; adios2::getstep(istream, istep);)
        std::string msg = istream.read<std::string>("data").front();
    istream.close()
}

void main(){
    std::thread reader_thread(__read__);
    std::thread writer_thread(__write__);
    reader_thread.join();
    writer_thread.join();
}
```

```
mpirun capio_server -c <config.json> &
CAPIO_DIR=... LD_PRELOAD=libcario_posix.so ./demo_program
```



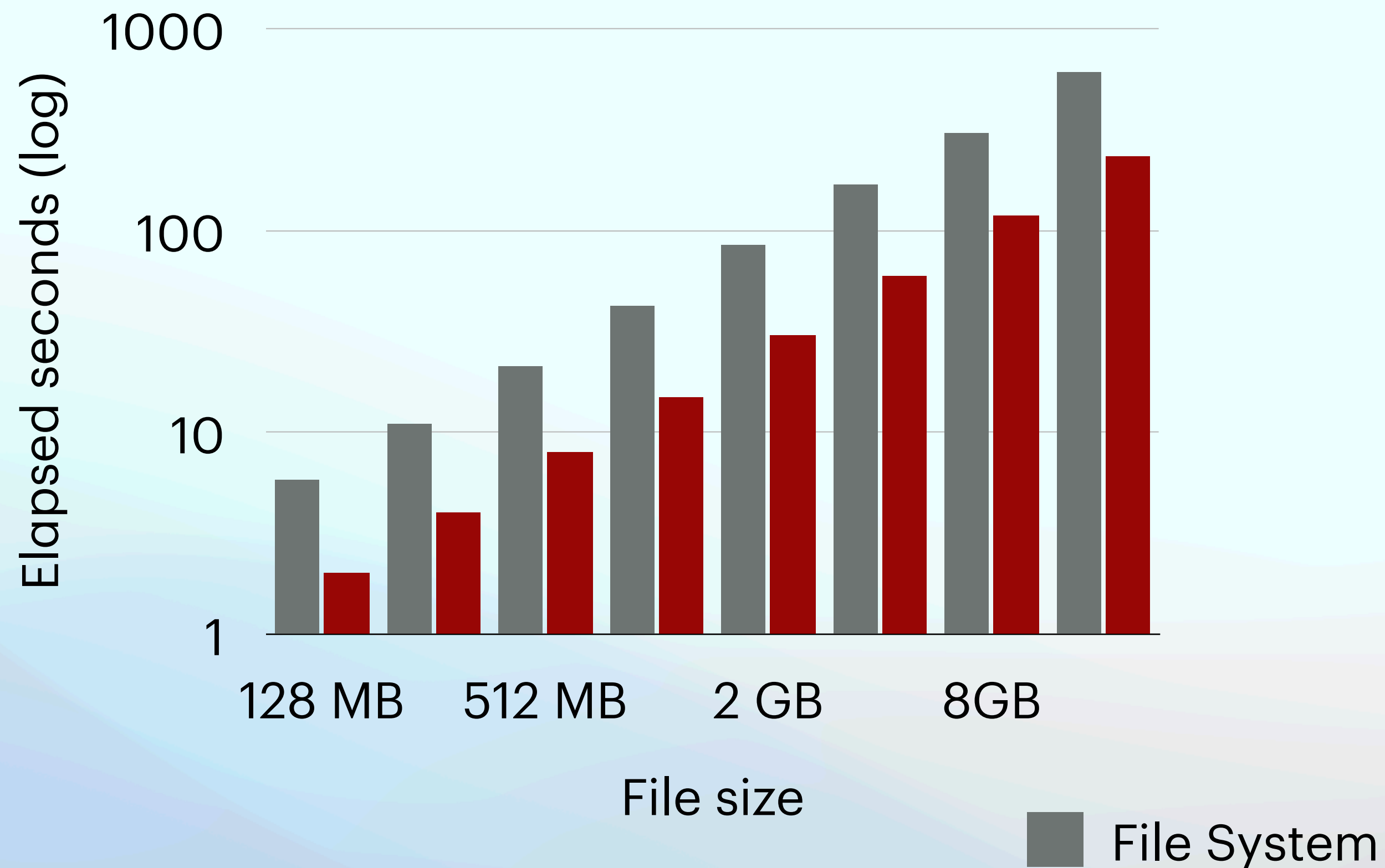
capio.hpc4ai.it

Performance

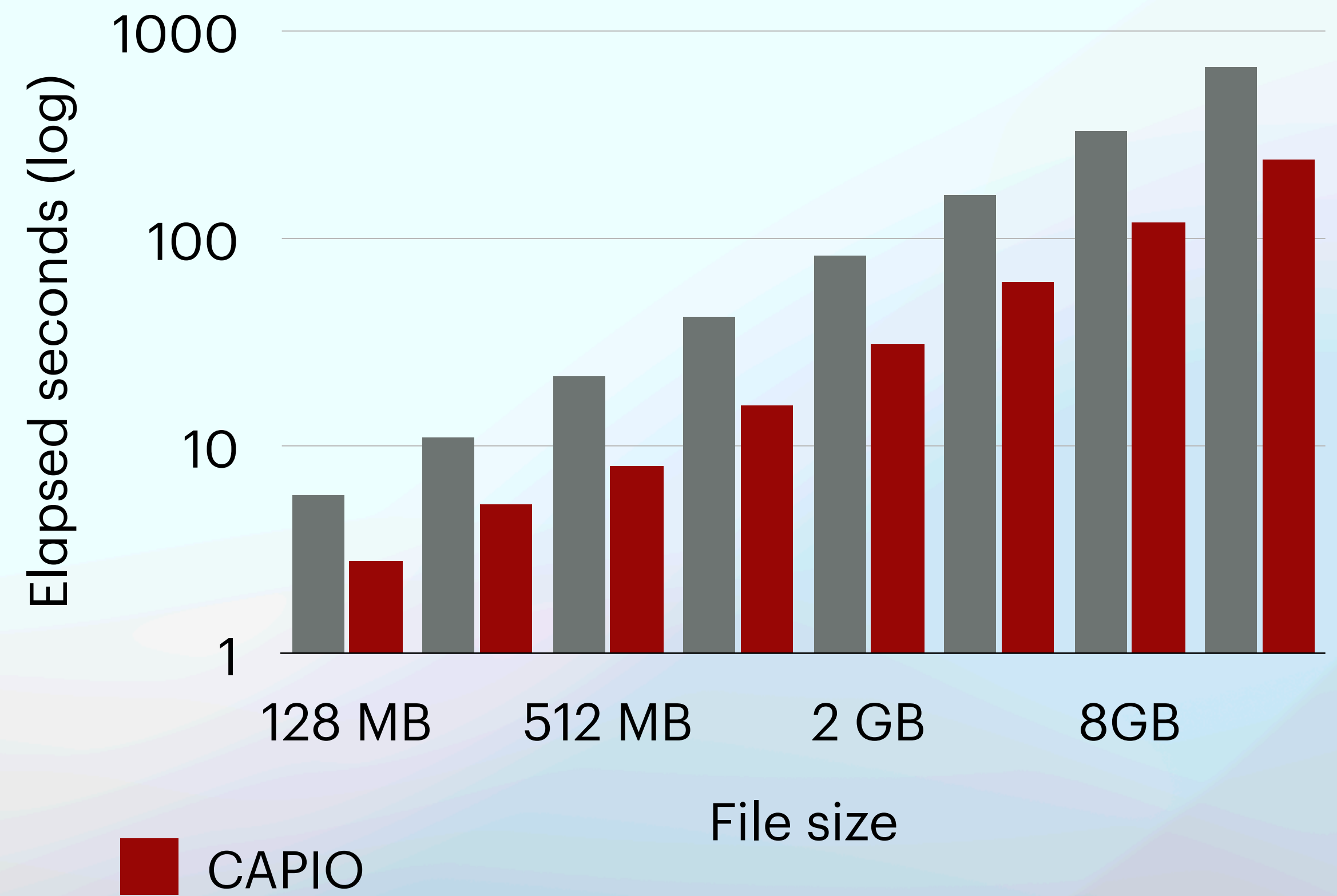
Basic IO patterns

(5 iterations per test)

One to One



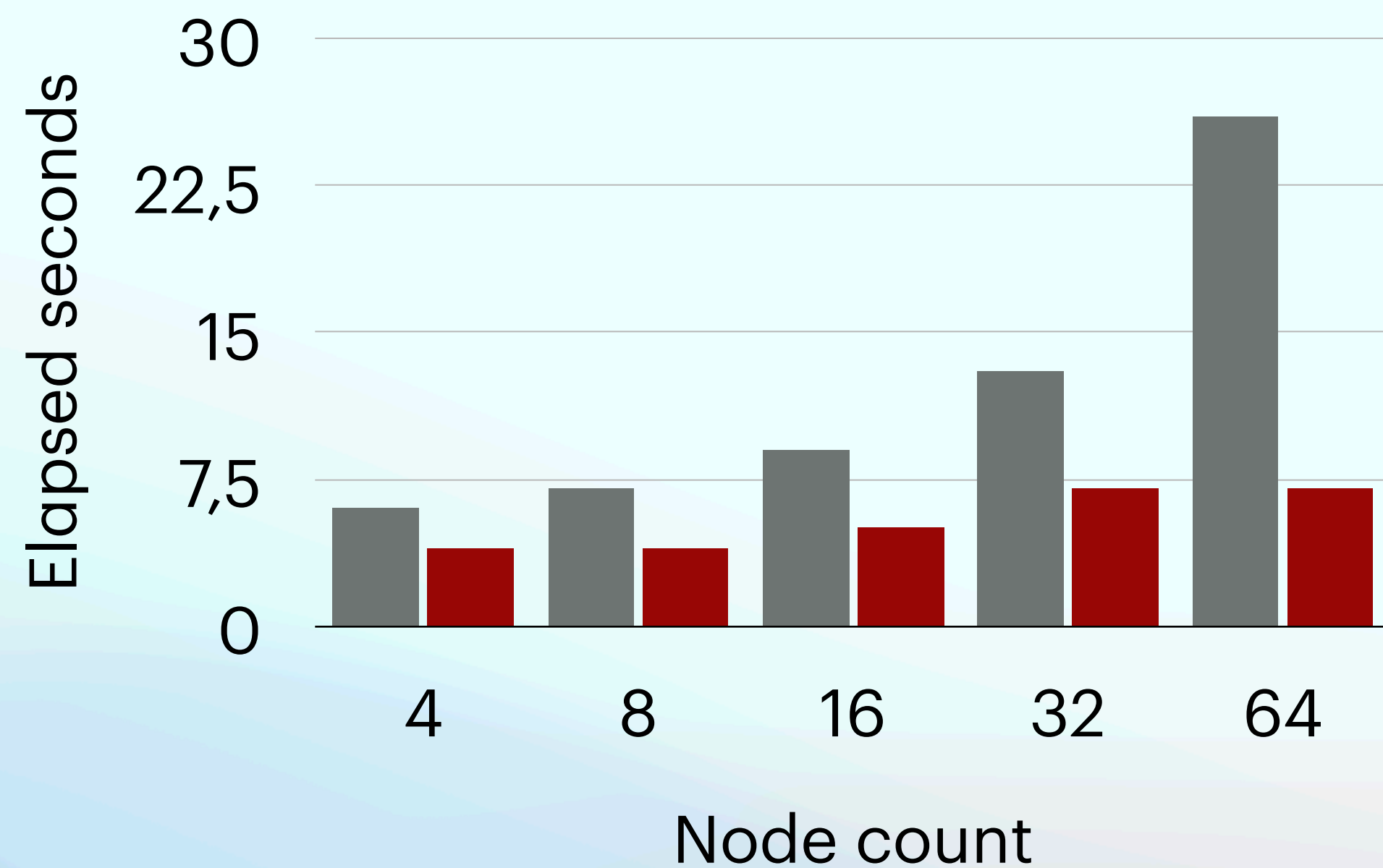
Broadcast (1 prod - 5 cons)



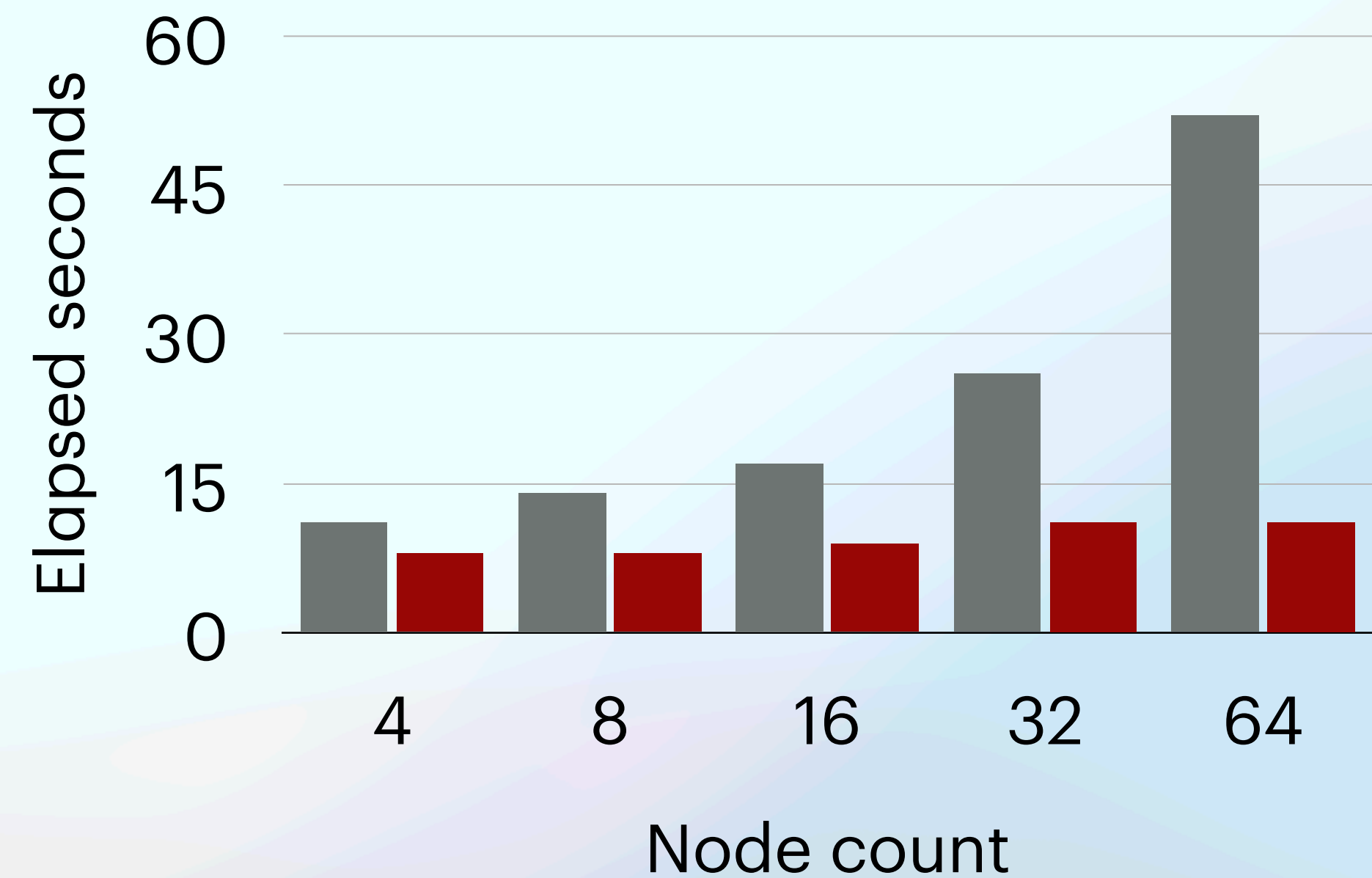
Scaling

(Broadcast, 5 iterations per test, 1 producer - n consumers)

128 MB Files



256 MB Files

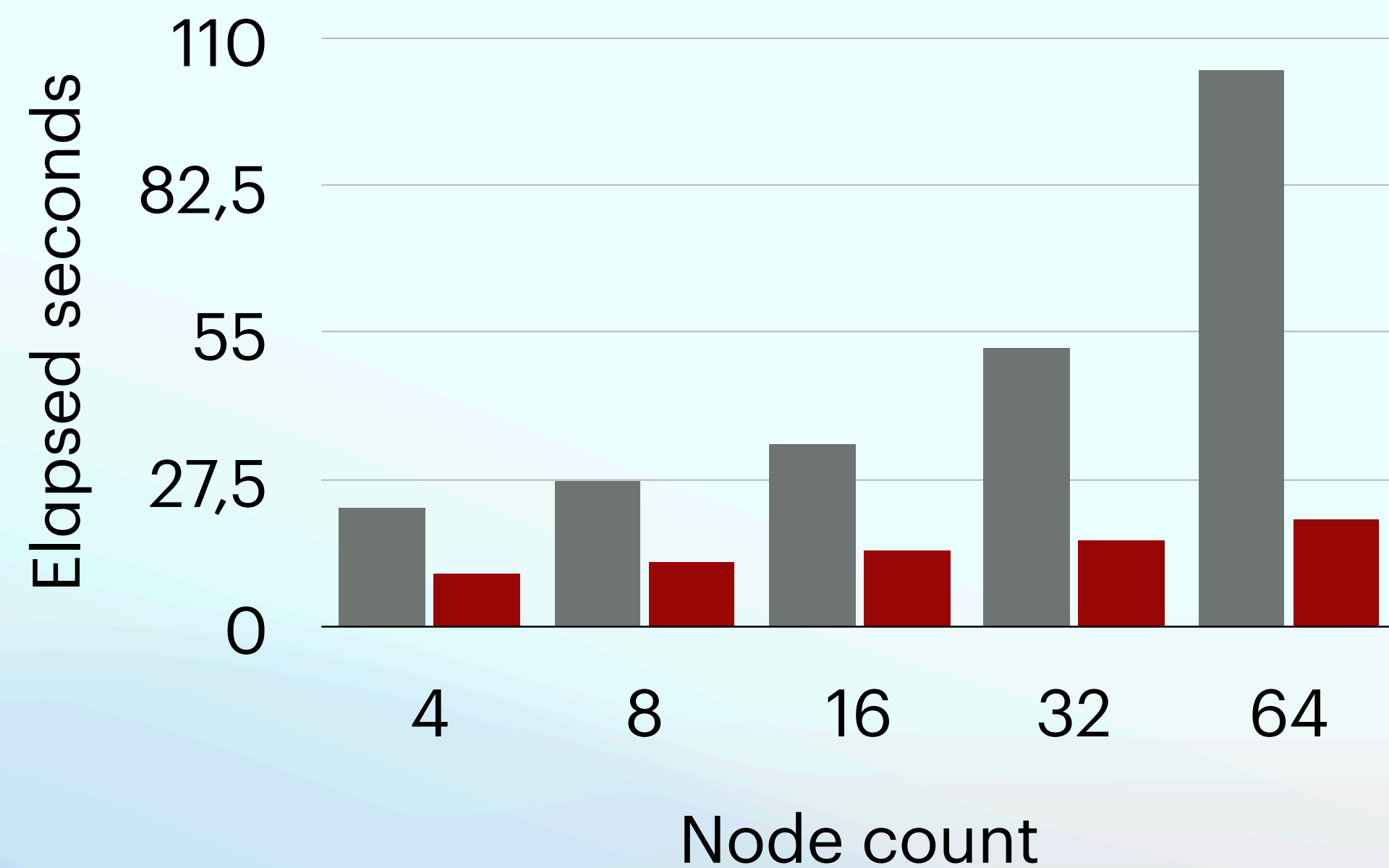


■ File System ■ CAPIO

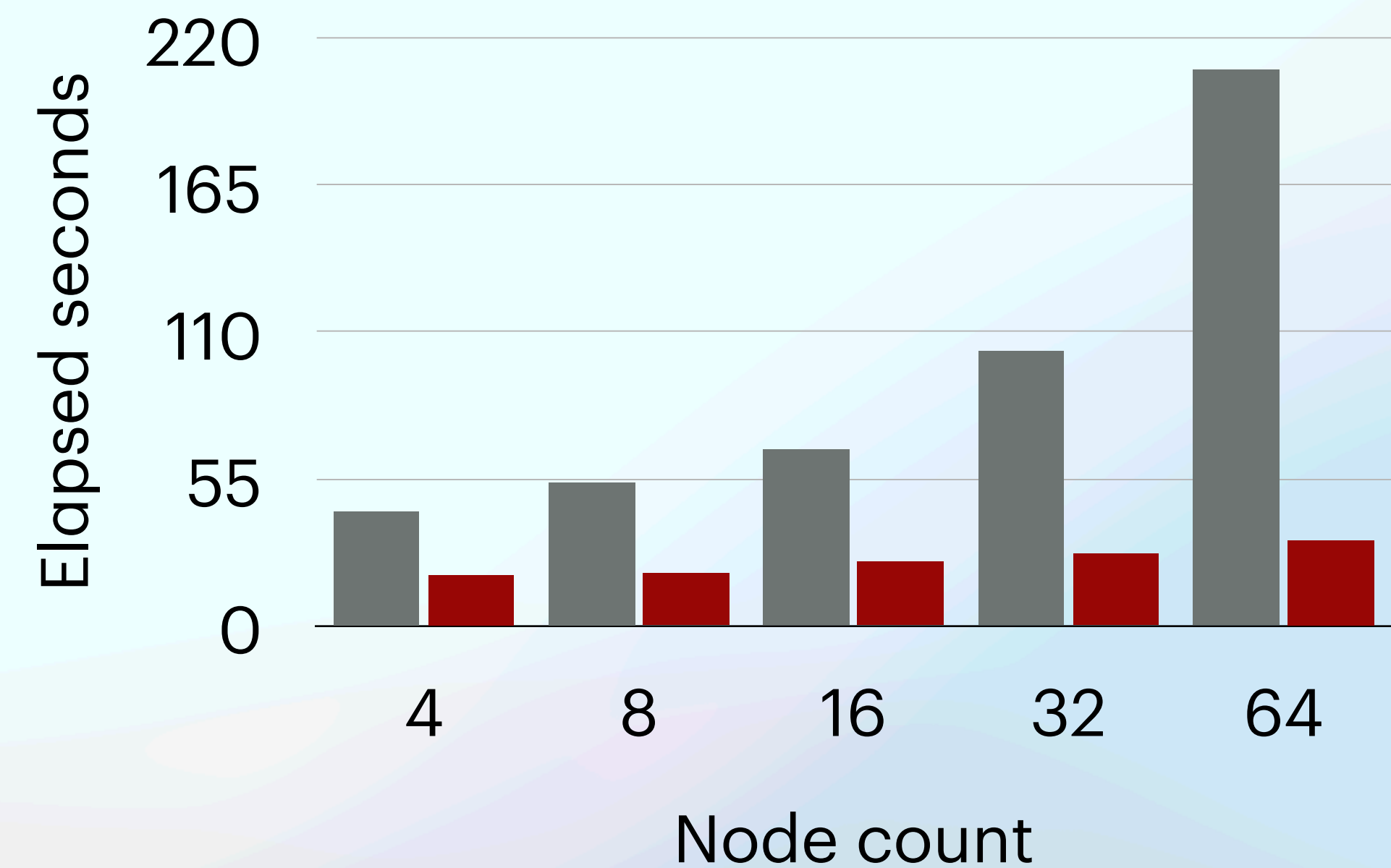
Scaling

(Broadcast, 5 iterations per test, 1 producer - n consumers)

512 MB Files

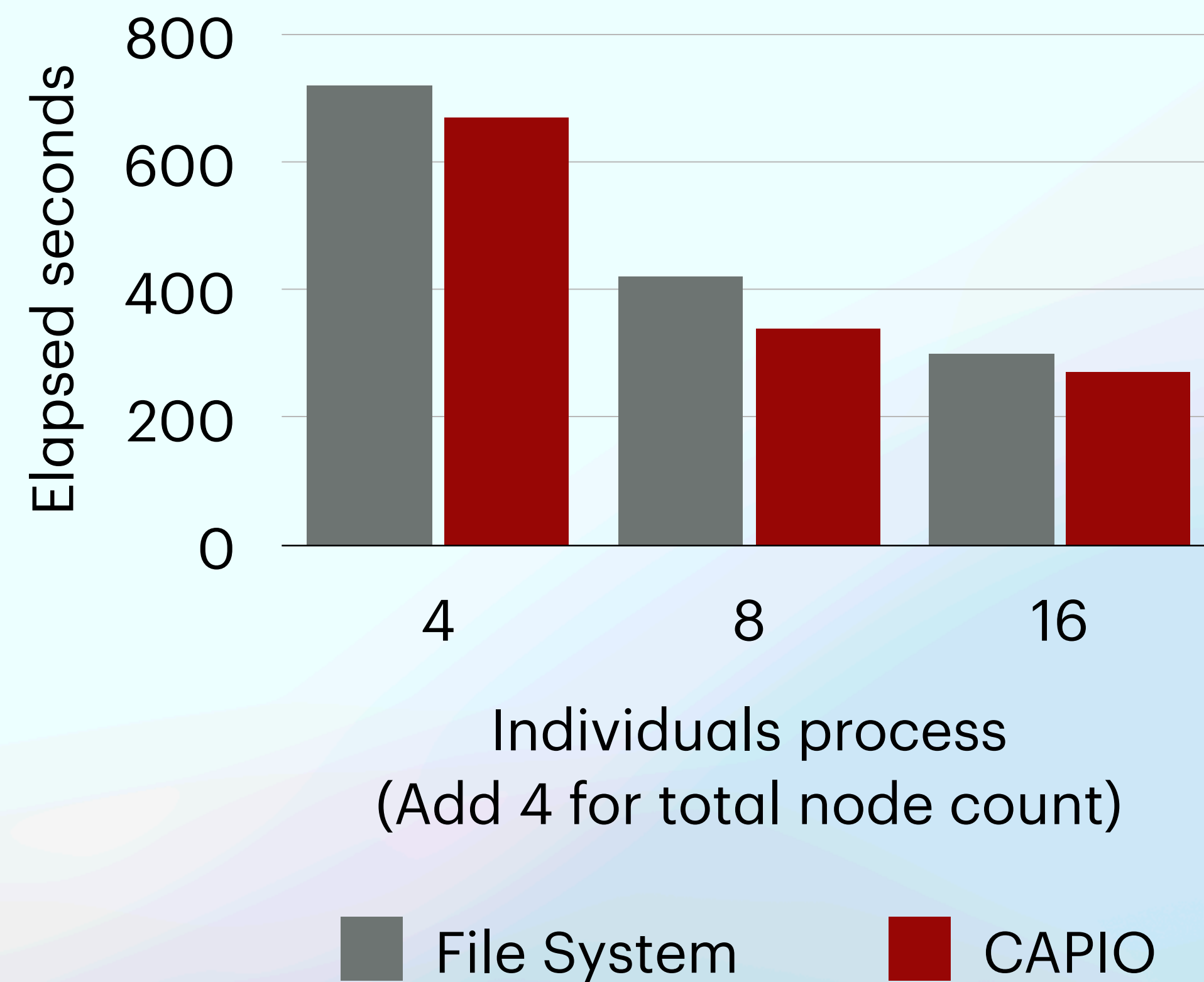
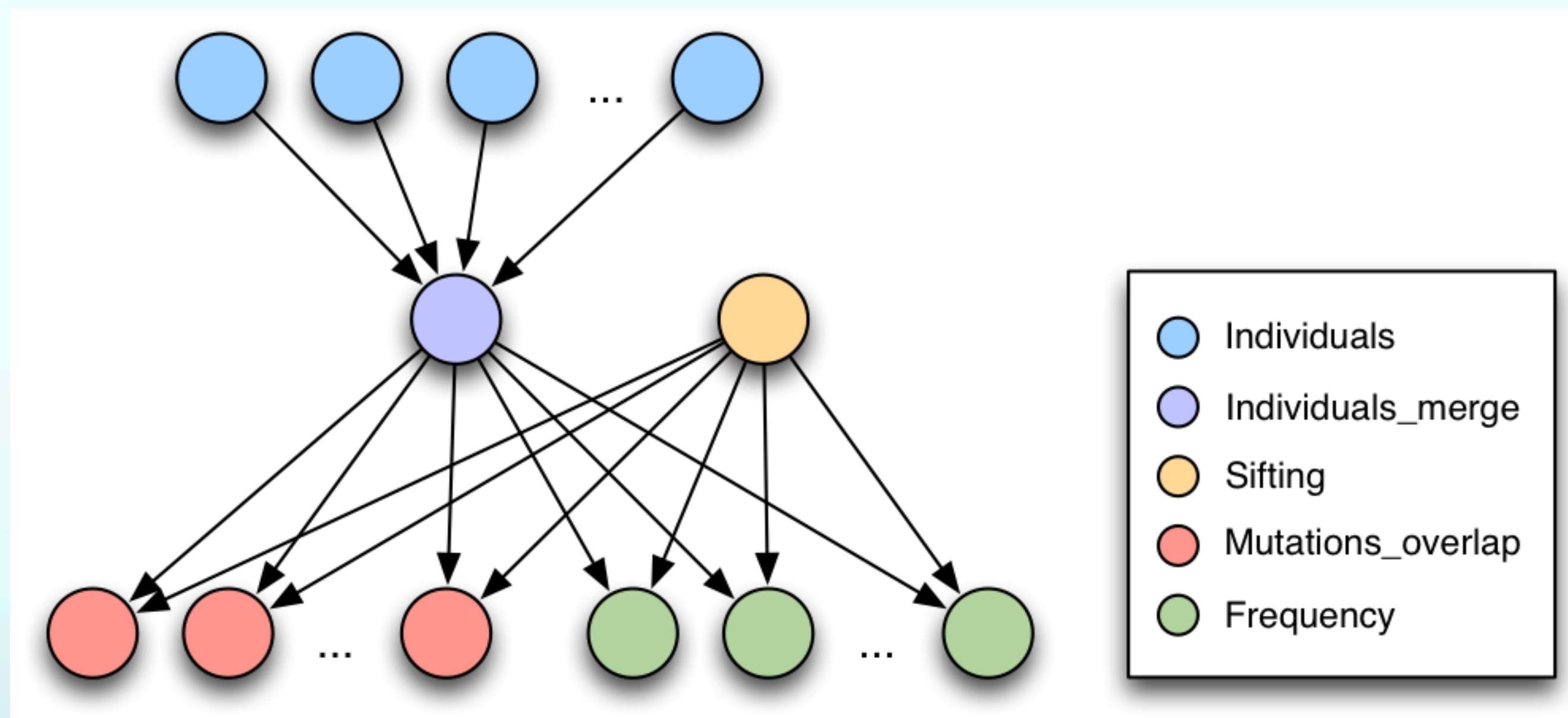


1024 MB Files



■ File System ■ CAPIO

1000 Genome workflow



Future works

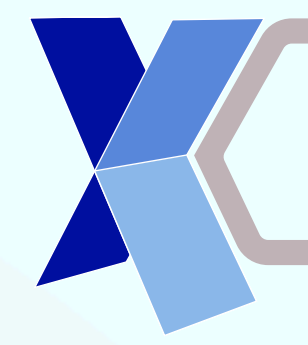
- Remove MPI from CAPIO
- Integration with WFM



UNIVERSITÀ
DI TORINO



UNIVERSITÀ
DI PISA



ICSC

Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

CAPIO
Cross application Programmable
I/O

Questions?

Marco Edoardo Santimaria

Iacopo Colonnelli

Massimo Torquati

Marco Aldinucci

marcoedoardo.santimaria@unito.it

iacopo.colonnelli@unito.it

massimo.torquati@unipi.it

marco.aldinucci@unito.it

