

Frequency Techniques for I/O

Ahmad Tarraf¹, Alexis Bandet², Francieli Boito², Guillaume Pallez², and Felix Wolf¹

¹Technical University of Darmstadt

²University of Bordeaux, INRIA



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming

ADMIRE

malleable data solutions for HPC

ADAPTIVE MULTI-TIER INTELLIGENT
DATA MANAGER FOR EXASCALE

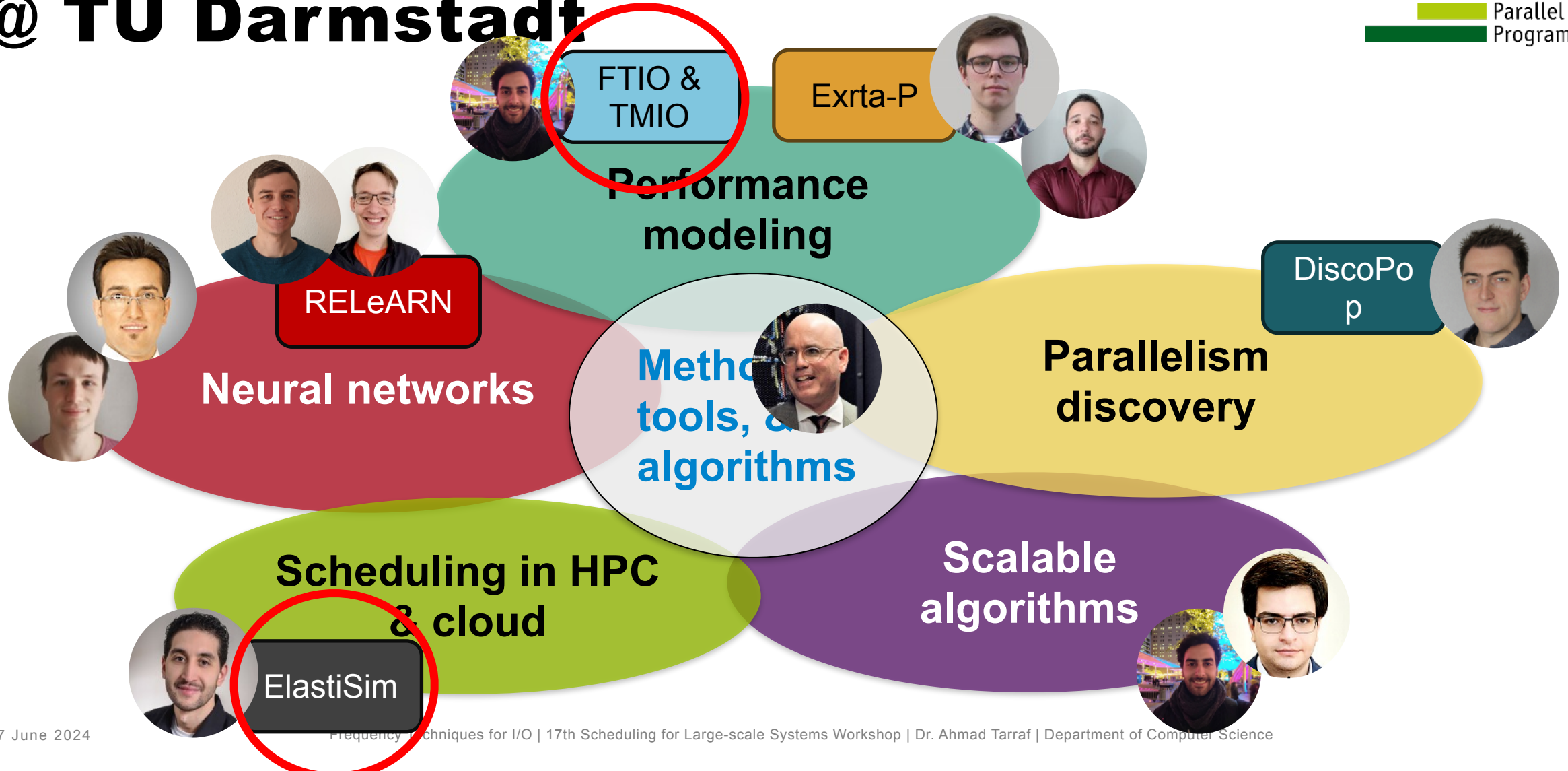
DEEP-SEA



EuroHPC
Joint Undertaking

17th Scheduling for
large-scale systems
workshop

Parallel Programming @ TU Darmstadt



Survey on Malleability in HPC



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming



Collaboration between 4 EuroHPC projects over 2 years

Preprint available on IEEE

IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, VOL. 35, NO. 6, JUN 2024

1

Malleability in Modern HPC Systems: Current Experiences, Challenges, and Future Opportunities

Ahmad Tarraf¹, Martin Schreiber^{2,3}, Alberto Cascajo⁴, Jean-Baptiste Besnard⁵, Marc-André Vef⁶, Dominik Huber³, Sonja Happ⁷, André Brinkmann⁶, David E. Singh⁴, Hans-Christian Hoppe⁸, Alberto Miranda⁹, Antonio J. Peña⁹, Rui Machado¹⁰, Marta Garcia-Gasulla⁹, Martin Schulz³, Paul Carpenter⁹, Simon Pickartz⁷, Tiberiu Rotaru¹⁰, Sergio Iserte⁹, Victor Lopez⁹, Jorge Ejarque⁹, Heena Sirwani¹⁰, Jesus Carretero⁴, Felix Wolf¹

Abstract—With the increase of complex scientific simulations driven by workflows and heterogeneous workload profiles, managing system resources effectively is essential for improving increase or decrease during its execution and resources added or released accordingly. Malleable jobs can also increase en-





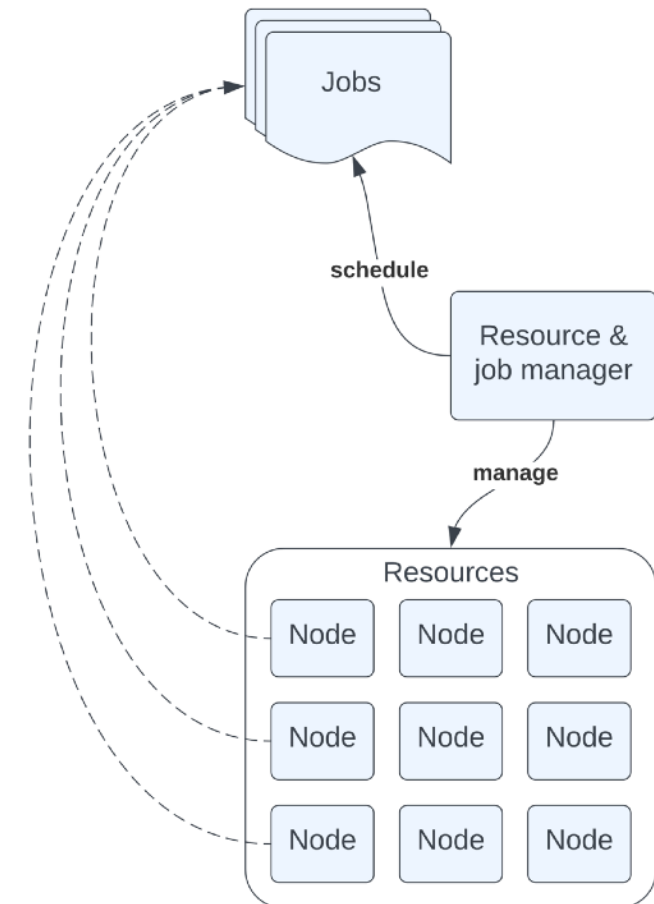
TECHNISCHE
UNIVERSITÄT
DARMSTADT

 Parallel
 Programming

ElastiSim: A Batch-System Simulator for Malleable and Evolving Workloads

Resource & job management

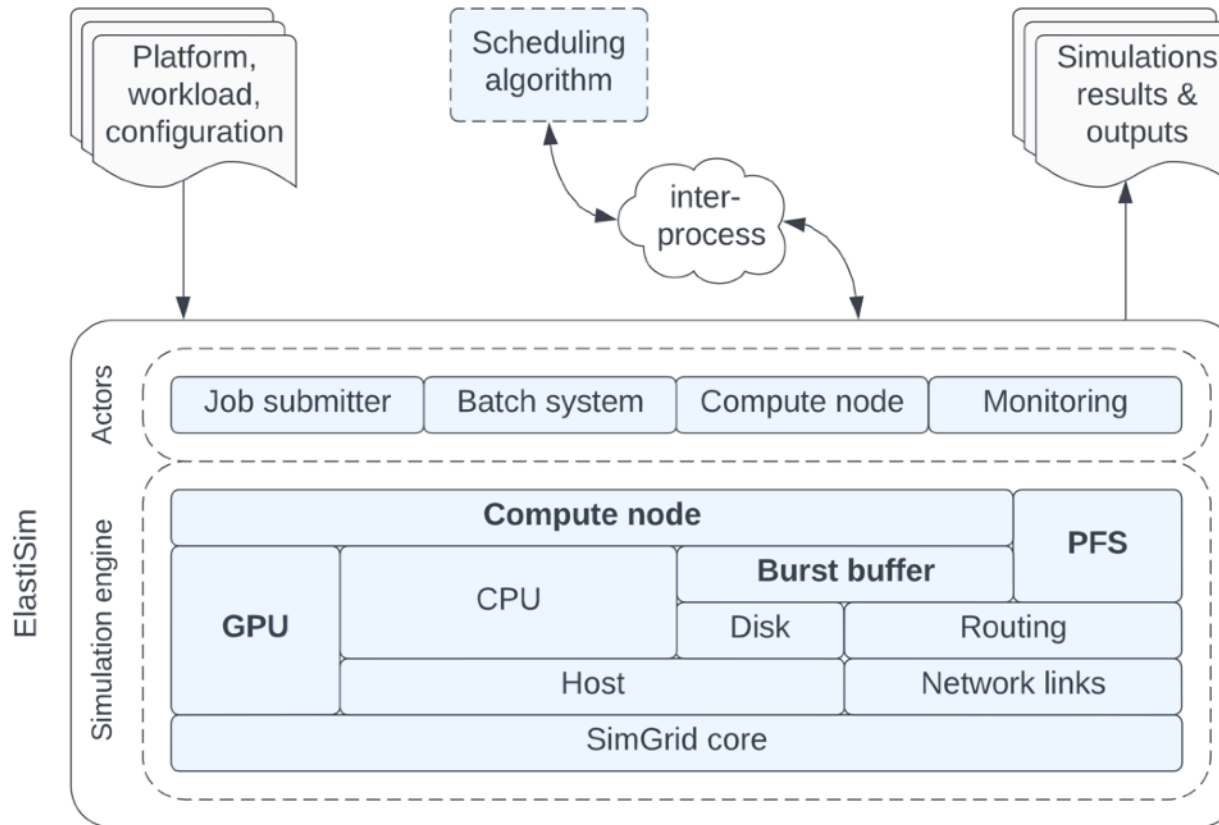
- Resource and job management systems (also often called batch systems) schedule jobs and provide resources in large-scale computing environments
- Depending on the objective, batch systems aim to maximize system efficiency and decrease job completion times
- Scheduling algorithms are key components to improve system performance



What is ElastiSim?

- ElastiSim is a simulator that simulates
 - jobs and applications,
 - the batch system supporting rigid, moldable, malleable, and evolving workloads,
 - the scheduling algorithm (as part of the batch system),
 - the platform (powered by SimGrid).
- Typical use case: evaluating algorithms for the combined scheduling of rigid, moldable, malleable, and evolving jobs

ElastiSim architecture



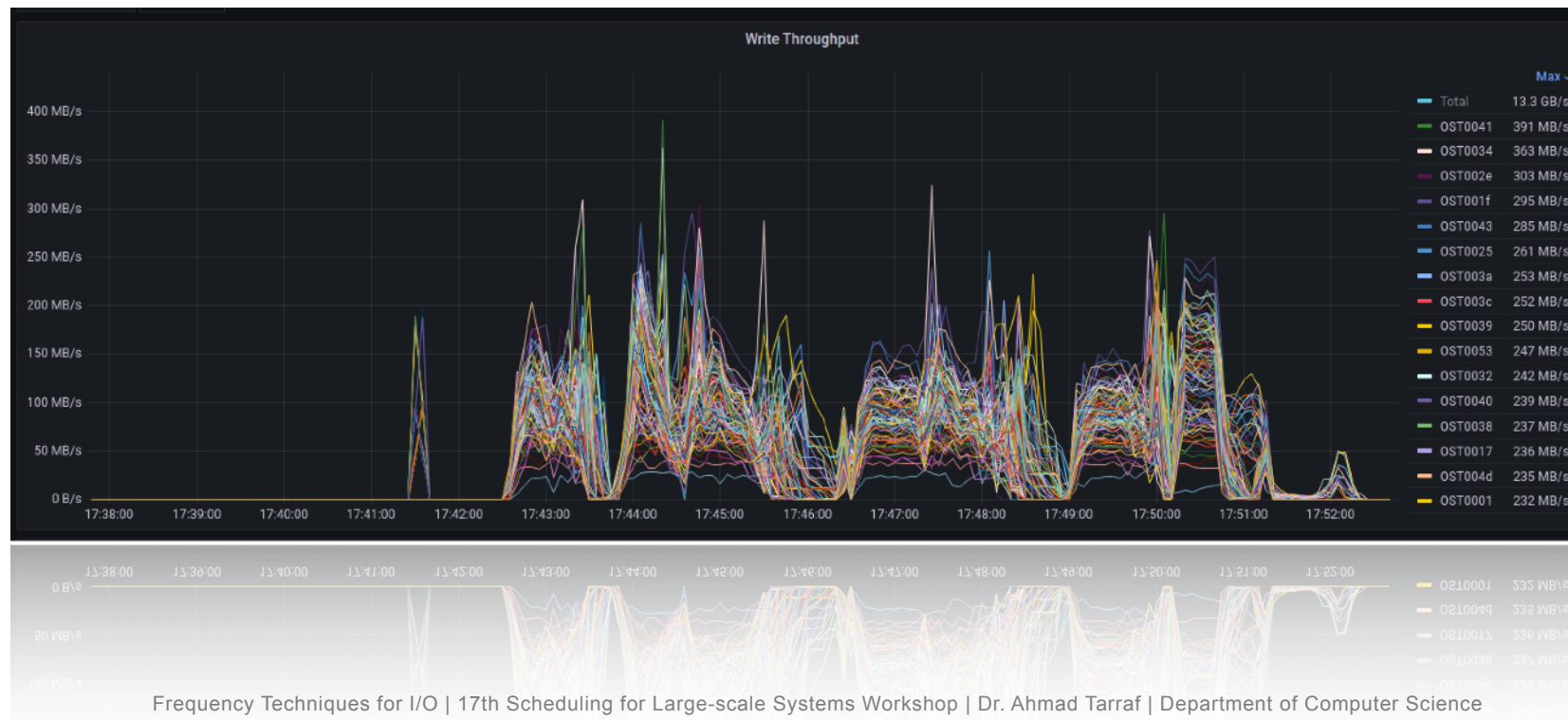
Summary

- Website: <https://elastisim.github.io>
(includes Slack invitation)
- GitHub: <https://github.com/elastisim>
- Contact me (Taylan Özden):
taylan.oezden@tu-darmstadt.de



ElastiSim

Frequency Techniques for I/O



Motivation

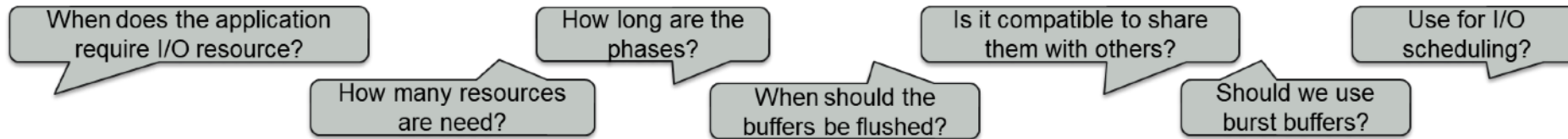
- HPC applications usually alternate between compute and I/O phases (e.g., Checkpointing)
- Compute resources are allocated **exclusive**
- I/O bandwidth is a **shared resource**; which often suffers from:
 - **Variability**: I/O performance depends on what others are doing
 - **Contention**: causes lower overall I/O performance
 - **Lower utilization**: compute resources are often “wasted” while waiting for I/O



SC22 after Jack Dongarra's presentation in the Dallas ballroom

Motivation (cont')

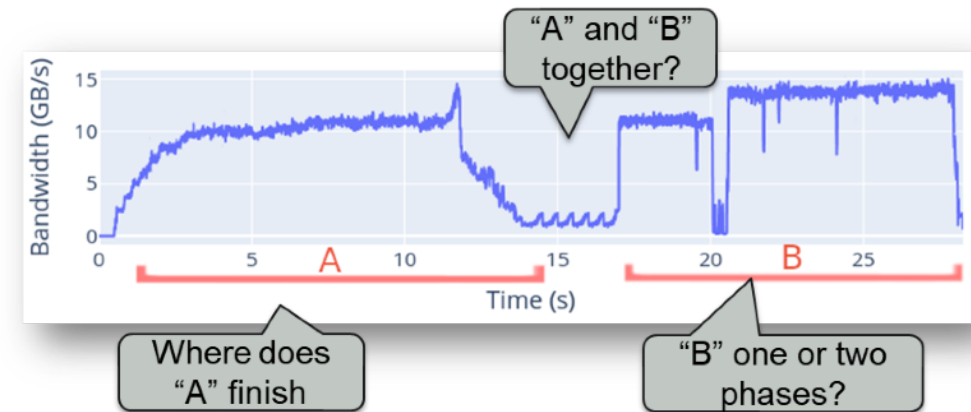
- Several **solutions** exist: I/O scheduling, I/O-aware batch scheduling, burst buffers/caches,
- **But they often required knowledge about the application's I/O behavior**
- Especially the **temporal I/O behavior** can be useful if proved online, to answer questions like:



- **Difficult to get this information online with a low overhead!**
- **Especially, if we think in terms of phases!**

- I/O phases composed of many I/O requests
- Not all I/O is interesting
- Borders of I/O phases? Threshold?

→ Application and system



FTIO: Frequency Techniques for I/O

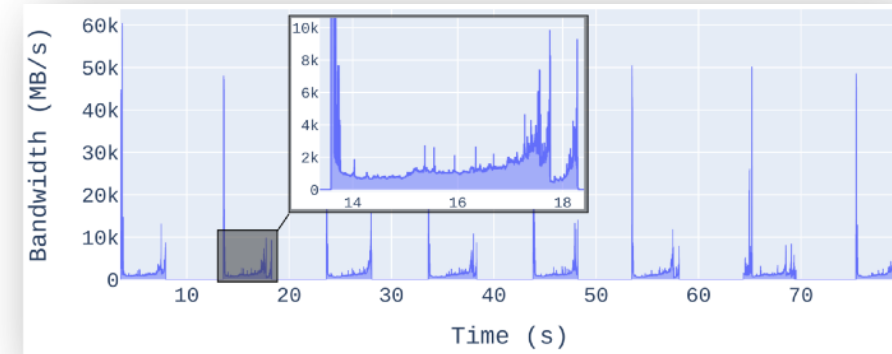


→ Periodic I/O is often encountered in HPC!

Information about applications' periodicity, **even if not perfectly precise**, leads to good contention-avoidance techniques [1, 2, 3]

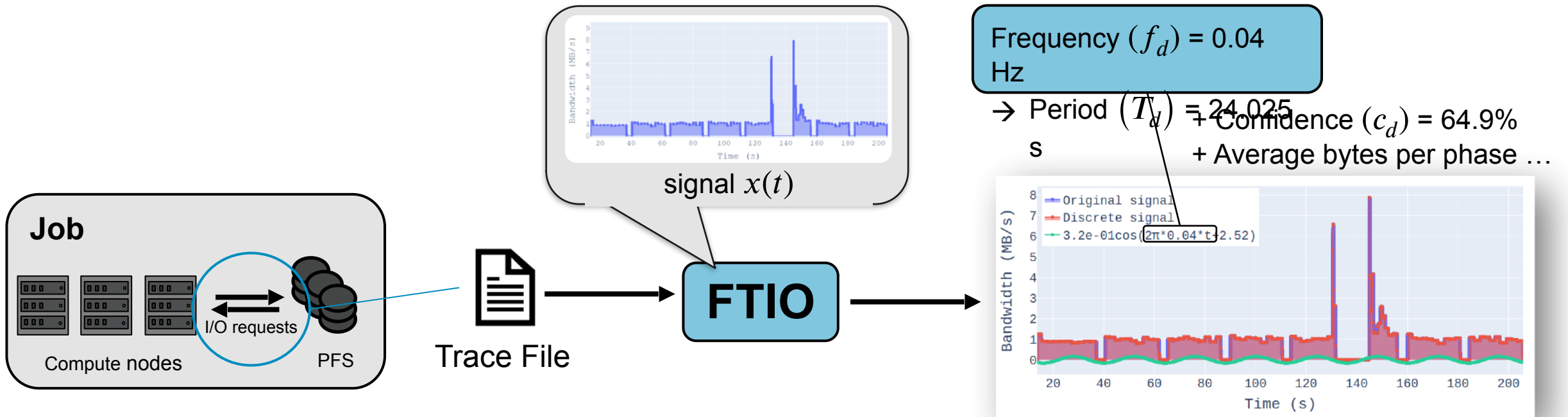
→ Frequency Techniques for I/O:

- Examine the I/O behavior in the **frequency domain**
- Describes the temporal behavior of the **I/O phases** through a single metric, namely the **period (T_d)**
- Online (**prediction**) and offline (**detection**) realizations with low overhead
- Additional metrics quantify the **confidence** in the results and **further characterize** the I/O behavior



Period (T_d) of I/O phases:
The time between the start of consecutive I/O phases

FTIO: In a Nutshell



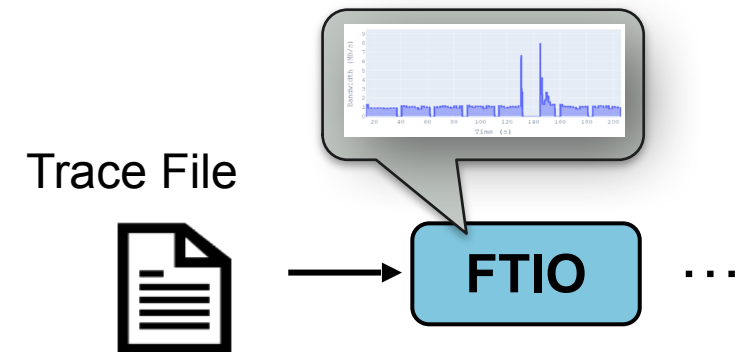
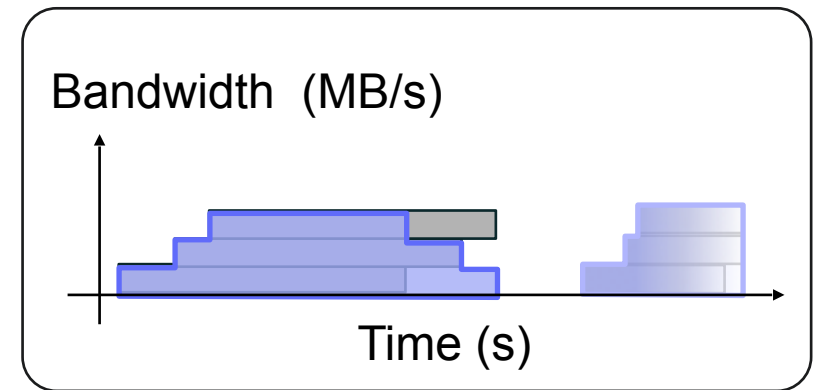
FTIO: Required Input

Trace file containing:

- Bandwidth per rank
 - Time (start and end) when the bandwidth changed
- FTIO calculates internally the **application-level bandwidth** by **overlapping** the rank-level metrics


Any metric can be predicted, not only I/O!
E.g., network consumption, scheduling points, energy, etc.

Application-level bandwidth and (start) time can also be provided directly → Any level is ok!





FTIO: Required Input (cont')

Supported Formats/Tools for **online prediction**:

- **TMIO** (JSONL, MessagePack, ZeroMQ)
- ADMIRE Monitoring Proxy (ISC24) 

Supported Formats/Tools for **offline detection**:

- Darshan
- Recorder (folder)
- TMIO (JSON, JSONL, MessagePack, ZeroMQ) 
- ADMIRE Monitoring Proxy (ISC24)
- Custom file format 

TMIO:

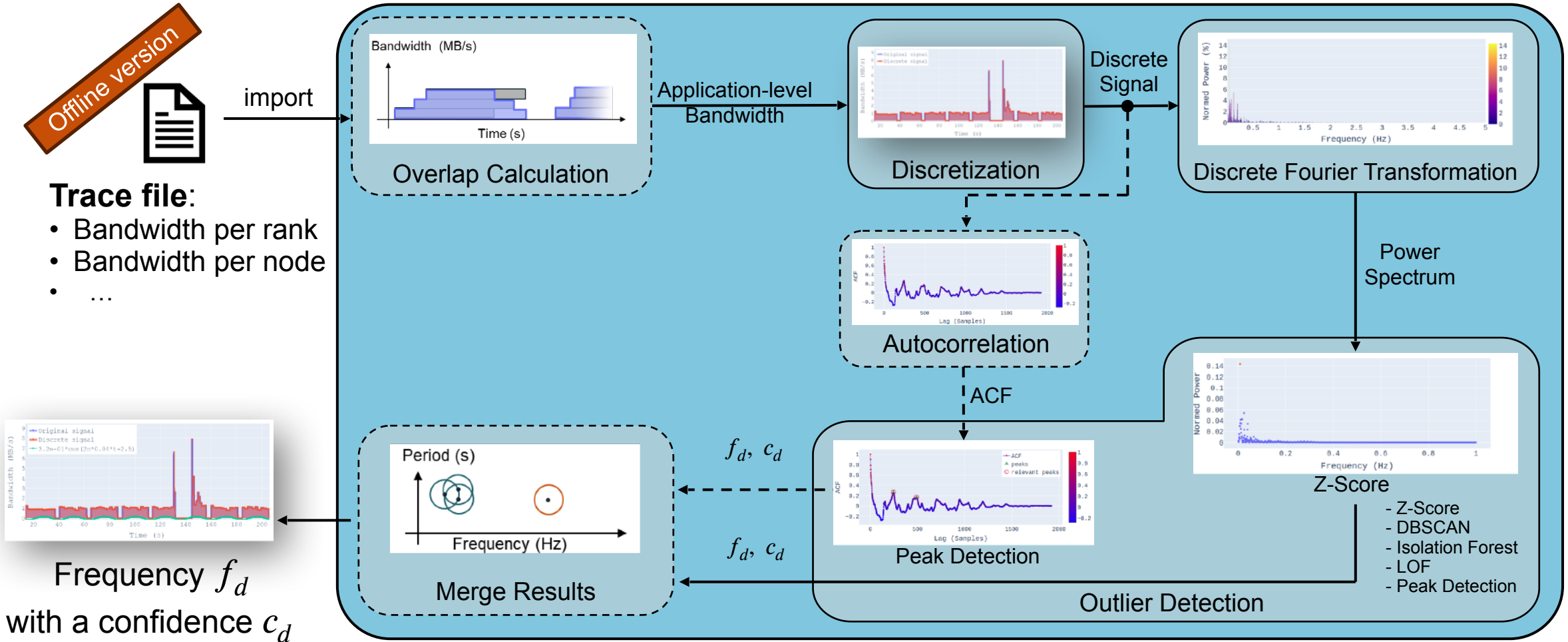
- Tracing **MPI-IO**
- C++ library that uses the PMPI interface
- Flushes I/O data online
- Can be easily attached to existing code
- Will be made publicly available

```
demo.json
12  "bandwidth": {
13    "b_rank_avr": [1.496276, 2.013454, 2.062243, ...],
14    "t_rank_start": [1.950272, 1.964889, 1.975749, ...],
15    "t_rank_e": [1.964871, 1.975739, 1.986342, ...]
16  }
```

https://github.com/tuda-parallel/FTIO/blob/main/docs/file_formats.md

FTIO: The Core

Parallel Programming



FTIO: User Interface



Parallel
Programming

```
Discretization
Time window : 192.223398 s
Frequency step: 0.005202 Hz
Sampling frequency: 10.0 Hz
Expected samples: 1922
Abstraction error: 0.000788

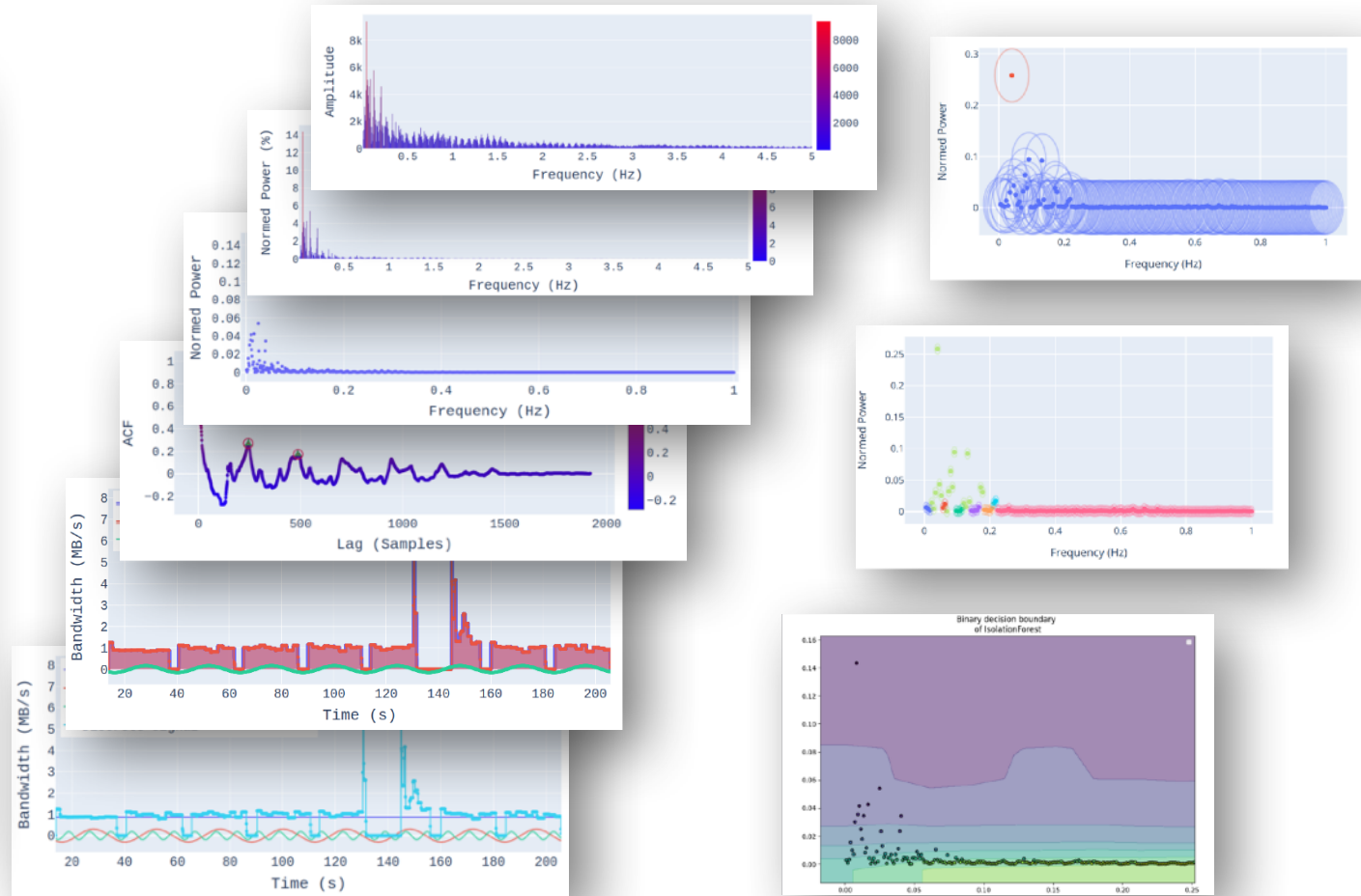
Discretization finished: 0.054 s
Executing: DFT + Z-score

DFT
Ranks: 1536
Start time: 13.76 s
End time: 205.98 s
Ignored bytes: 0.800000

Z-score
Spectrum: Power spectrum
mean: 1.041e-03
std: 1.874e-03
Frequencies with Z-score > 3 -> 11 candidates
+ Z > Z_max*80.0% > 3 -> 1 candidates
Dominant frequency at: 4.162e-02 Hz (t = 24.025 s, k = 6) -> confidence: 64.953%
Precision of 0.04 Hz is 11.16% (Positive only: 11.33%)

DFT + Z-score finished: 0.006 s
Total elapsed time: 0.838 s

Prediction results:
Frequency: 4.162e-02 Hz -> 24.025 s
Confidence: 64.95 %
```



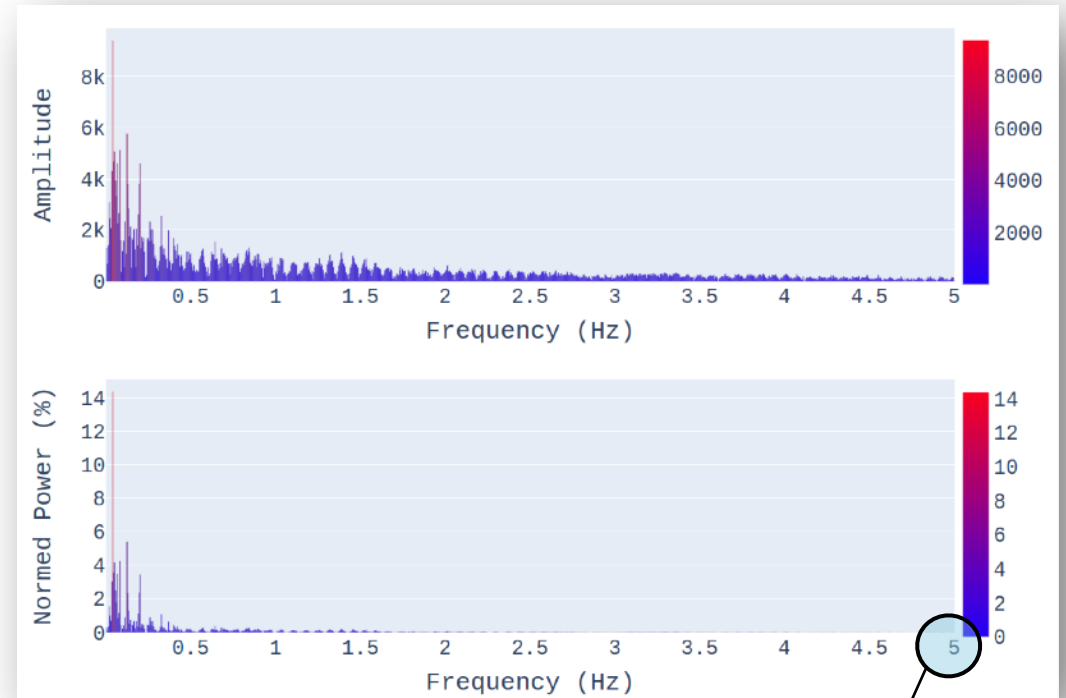
FTIO: Summary

Periodicity detection:

- **DFT + outlier detection** (Z-score, DB-Scan, Isolation forest, peak detection, or LOF)
- **Optionally: Autocorrelation + Peak detection**
- Merge results from both predictions (DB-Scan)

Properties

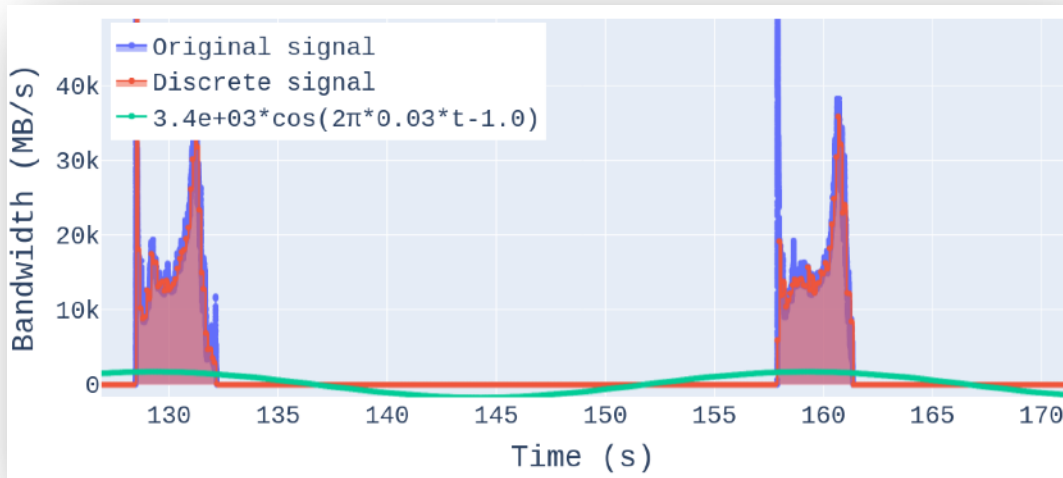
- Filters noise (e.g., using the power spectrum)
- **Several parameters control the accuracy** (sampling frequency f_s , time window Δt , and number of samples N)
- **Online version** offers two methods to adapt to changing behavior
- Optimized Python code that uses true multiprocessing (pools or manual process creation)



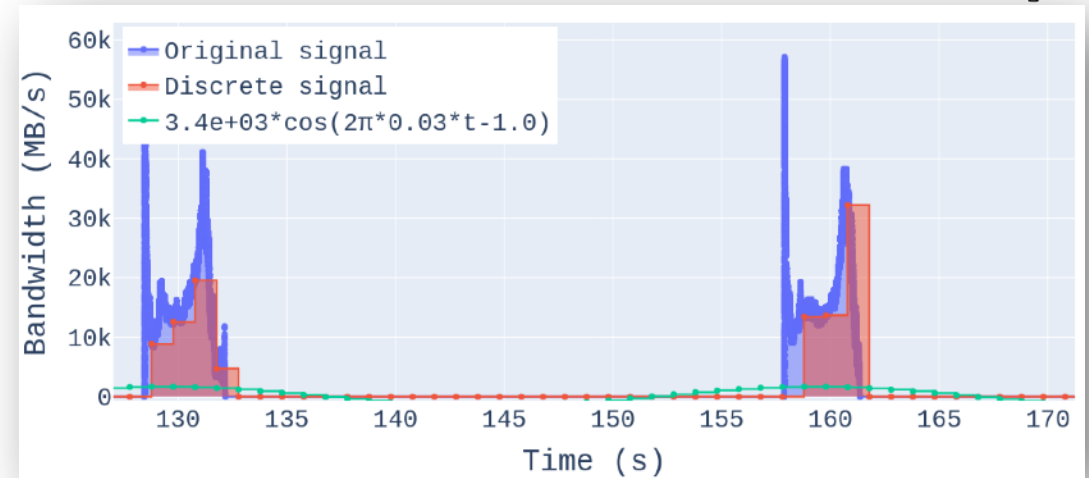
With $f_s = 10$ Hz, the limit is at $\frac{f_s}{2}$

→ DFT is symmetric, only half the spectrum is needed

FTIO: Sampling Frequency



$$f_s = 10 \text{ Hz}$$



$$f_s = 1 \text{ Hz}$$

Sampling frequency (f_s):

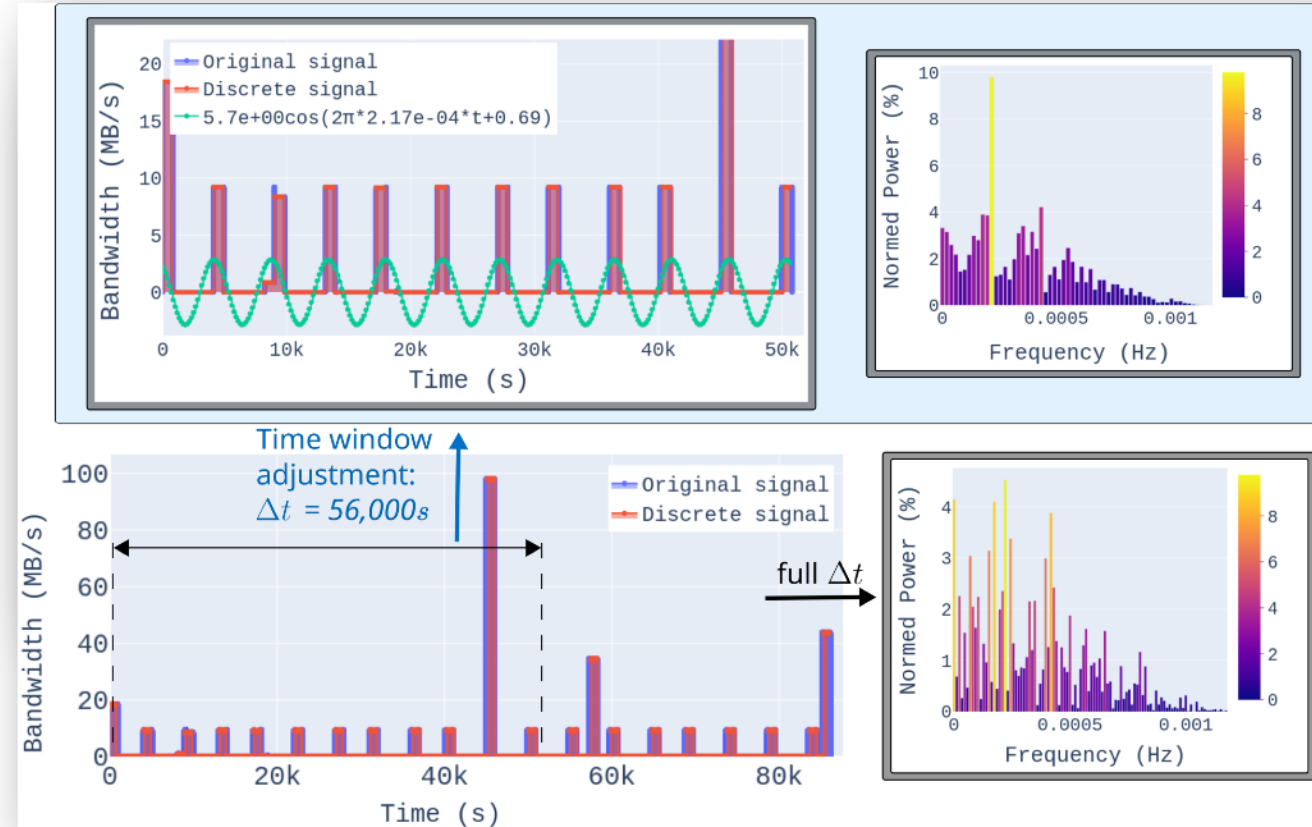
- Used to control the granularity at which the data is captured
- Specifies the range of frequencies of interest (Nyquist:

$$\left[0, \frac{f_s}{2}\right)$$

FTIO: Detection Example



- Nek5000 with 2048 ranks on the Mogon II (downloaded from the [HPCIO analysis website](#))
 - FTIO automatically sets f_s to 0.0006 Hz (bin widths in seconds)
 - FTIO detected I/O phases are **not periodic** in the full-time window due to **irregular** I/O phases:
 - Phases at 0 s and 45,000 s write 13 and 75 GB
 - Phases at 57,000 and 85,000 s write ~30 GB each
 - Other phases write 7 GB
- When Δt is set to 56,000 s, FTIO detects a dominant period of **4642.1 s** with **85.4 %** confidence

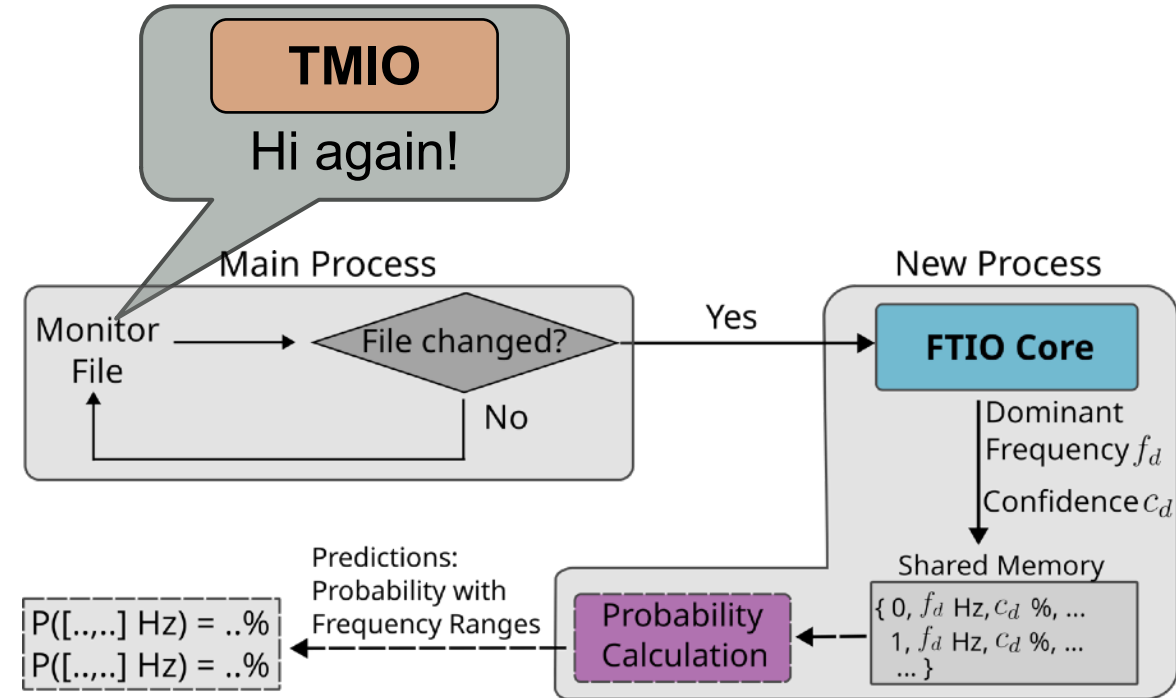


FTIO: Online Version

- Predicts the period **during the execution** of an application
- Monitors a file for changes, whenever a changes is detected, a new prediction **process** is launched

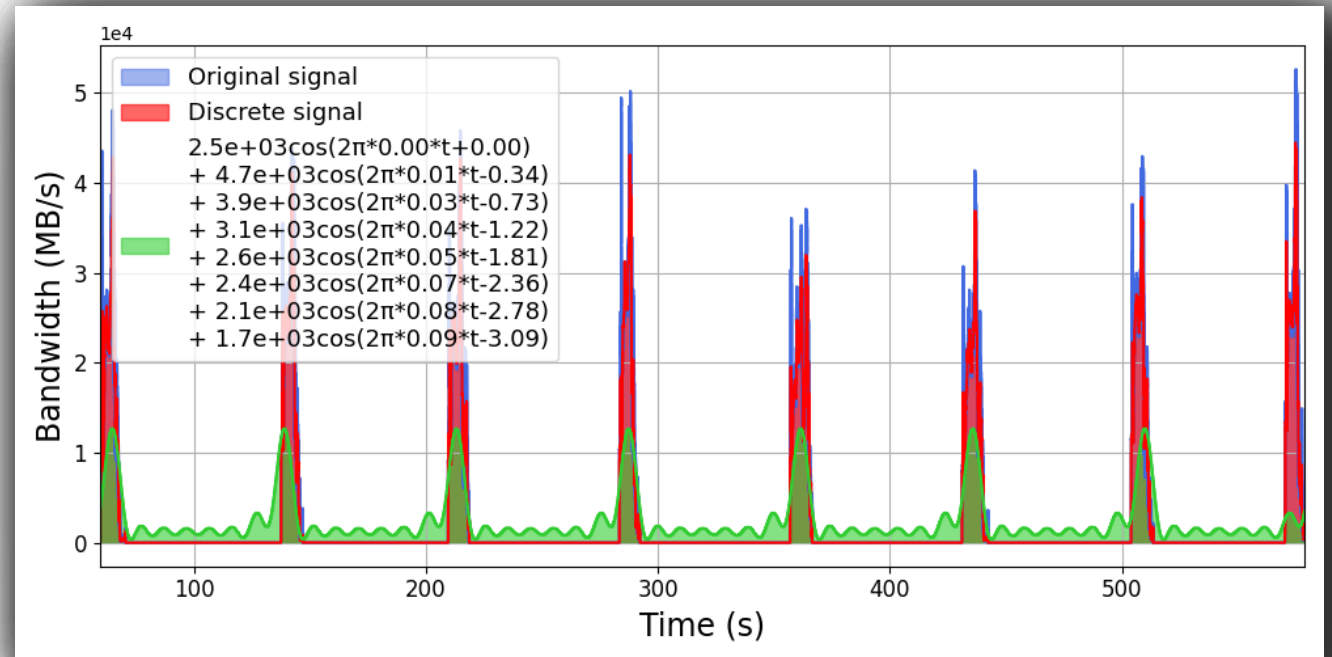
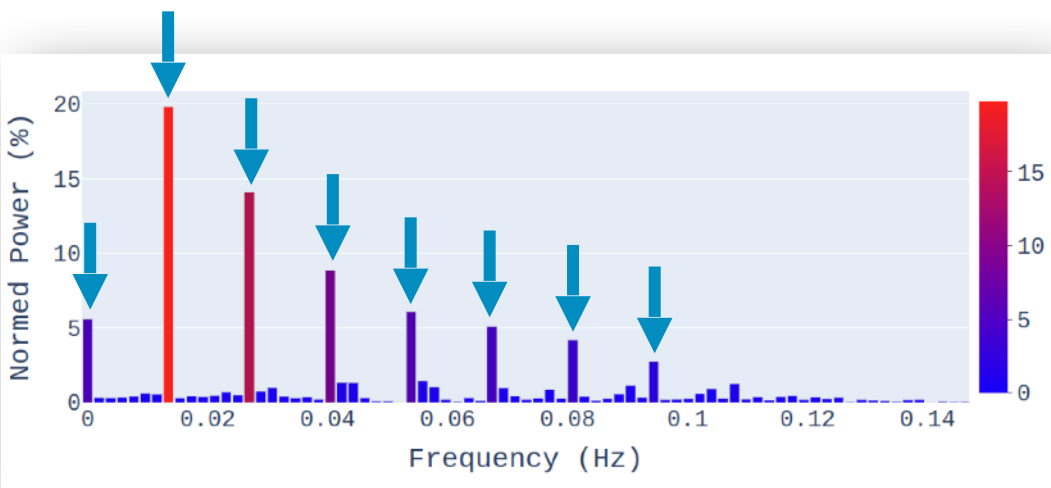
To adapt to **changing I/O behavior** FTIO offers:

1. Adapting time windows (discards the old data at some point); the width of the time window is found by FTIO based on the found period
2. Probability calculations with frequency intervals



NEW: Support for TCP messages using ZeroMQ

FTIO and I/O Bursts: IOR with 7680 Ranks

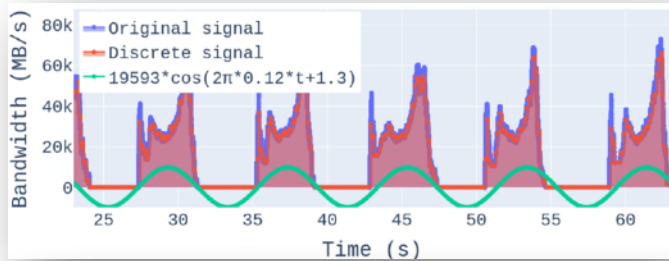


More Examples and Extensions

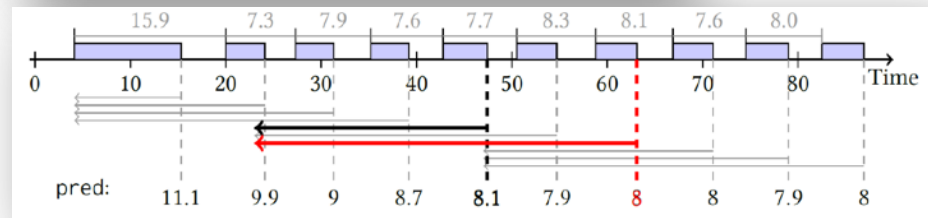
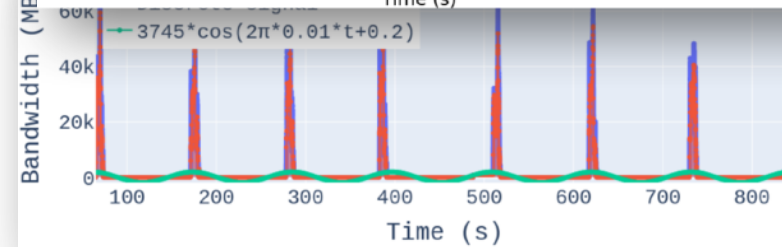
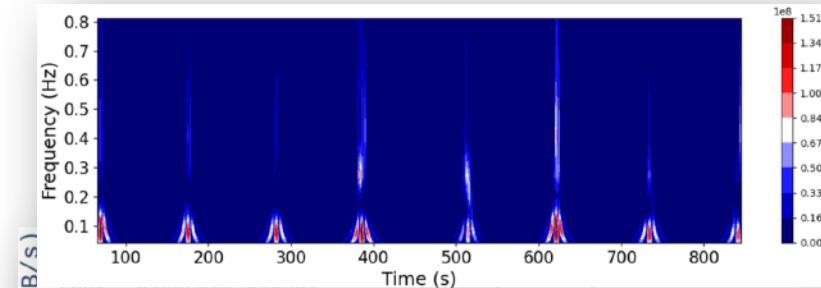
Parallel
Programming

IOR with 9216 ranks:
Continuous Wavelet
transformation

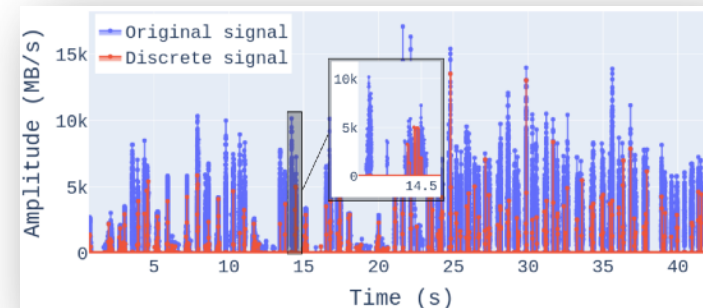
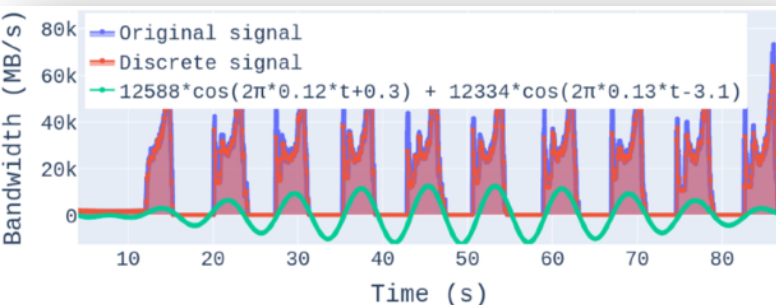
Mini-IO with 144
ranks: Non-periodic
signal ($f_s = 100$
Hz)



HACC-IO with 3072
ranks: Online Time
window adaptation



HACC-IO with 3072
ranks: Merge
Frequencies for
higher accuracy



FTIO Meets I/O-Sets

- **IO-Sets:** Method for I/O scheduling and the **Set-10** heuristic
 - Places applications on classes according to their time between the start of consecutive I/O phases (w_{iter})
 - Priority depends on class
 - F. Boito, G. Pallez, L. Teylo, and N. Vidal, IEEE TPDS 2023, <https://inria.hal.science/hal-03648225>
- **FTIO:** Frequency techniques to characterize temporal I/O behavior
 - Finds the frequency (f_d) of I/O phases ($\frac{1}{period}$, or $\frac{1}{w_{iter}}$)
 - A. Tarraf, A. Bandet, F. Boito, G. Pallez, and F. Wolf, IPDPS 2024 (to appear)
- **Set-10 implementation on BeeGFS servers**
 - The BeeGFS client sends the application's priority together with each I/O request to the servers
 - C. Barthelemy, F. Boito, E. Jeannot, G. Pallez, and L. Teylo, (unpublished for now)

FTIO Meets I/O-Sets (cont')

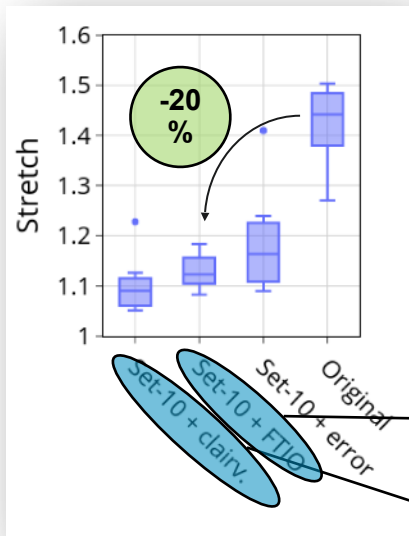
- **FTIO** in prediction mode watches over the trace of each application:
 - **TMIO** monitors an application and continuously appends to a trace
 - **FTIO outputs frequency and confidences** whenever new traces are available
 - A wrapper code, which called FTIO, recovers the output, and calculates the priority according to Set-10 heuristic and writes it to a per-application `/sys/kernel/config/` file
 - Before the first FTIO prediction, use a default value
 - **Whenever FTIO cannot answer (low confidence < 50%), keep the previously given priority**
- Used the **Grid'5000 French infrastructure** (as we needed root access)
- BeeGFS with a single OSS and a single OST, writes to a local hard disk
- BeeGFS client recovers the priority from the file when sending requests

Experimental Methodology

- Applications are generated using the IOR benchmarking:
 - Using fsync option(to have stable performance without caching)
 - Using MPI-IO API with file-per-process write access
 - Modified IOR:
 - To get start and end timestamps of I/O phases (to calculate metrics)
 - Included **TMIO**
 - 16 applications, each with 8 processes, all on the same client node
 - 15 with low-frequency: 10 iterations of sleep (compute) for 360 s then write 320MB (period of ~384s)
 - 1 with high-frequency: 200 iterations of sleep (compute) for 18 s then write 16MB (period of ~19.2s)
- Basically, we recreated the experiment from the **IEEE TPDS paper where Set-10 had excellent results** (while adapting to a different platform)

I/O Scheduling Results

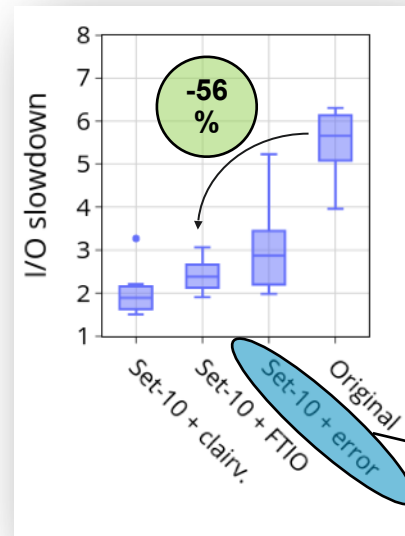
Parallel
Programming



The lower,
the better

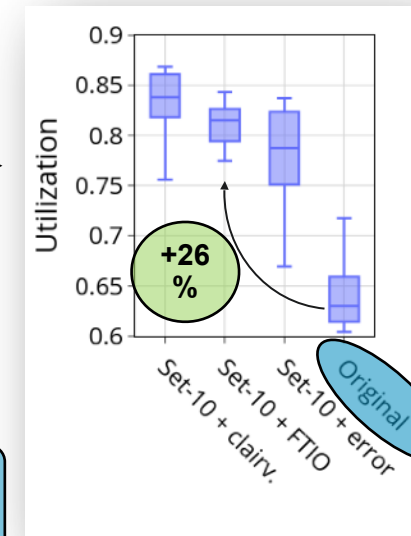
FTIO provides **actual**
periods to Set-10

Ideal periods are
hardcoded



The higher,
the better

Error (50%) is injected
in the result from FTIO



BeeGFS without
any modifications

Stretch:

For each application, how much **it was slowed-down by others compared to running in isolation** itself (min=1)

I/O-Slowdown:

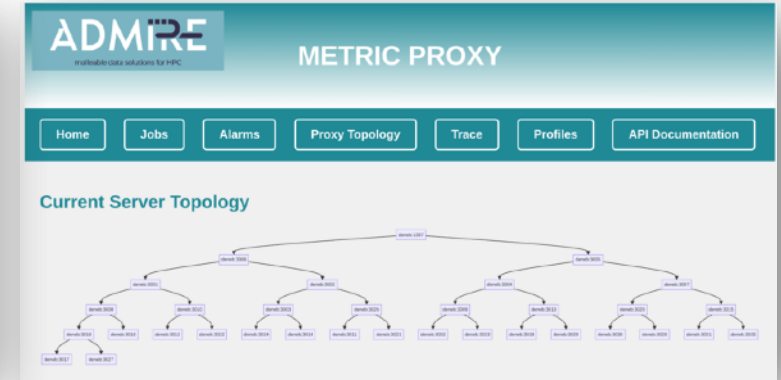
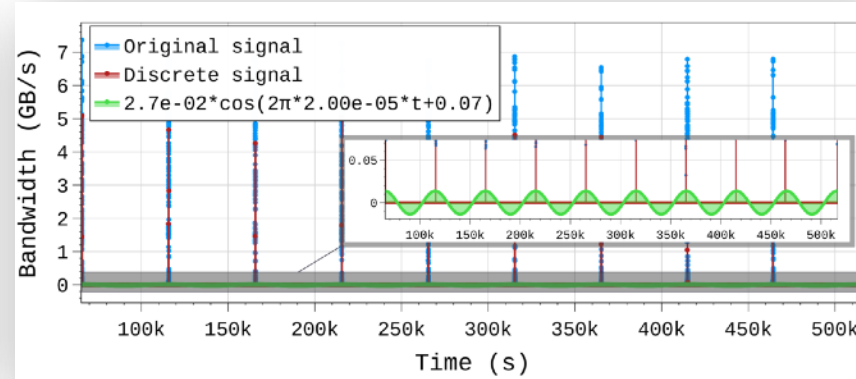
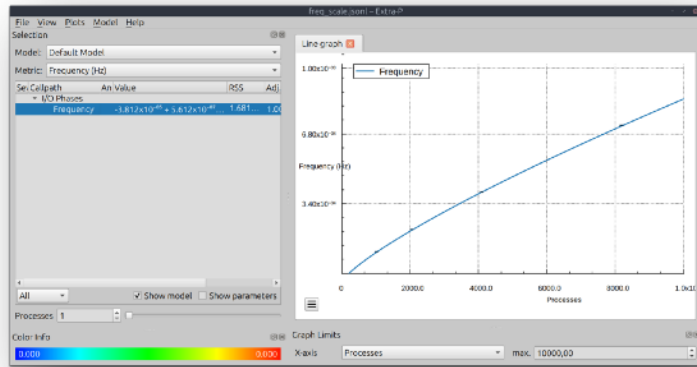
For each application, how much **its I/O was slower-down compared to running in isolation** (min=1)

Utilization:

How much of the **system time was spent on compute**, i.e., not doing or waiting for I/O ($\in [0, 1]$)

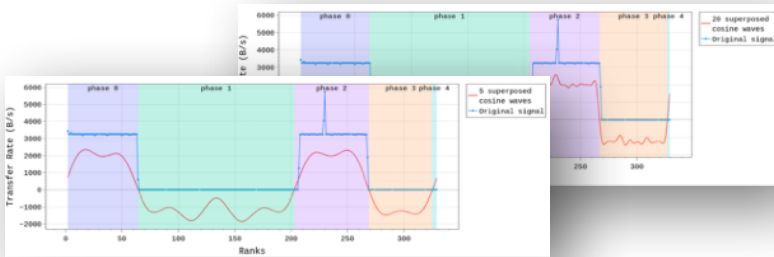
→ FTIO makes Set-10 possible in practice, where the period is not known in advance!

Future Work



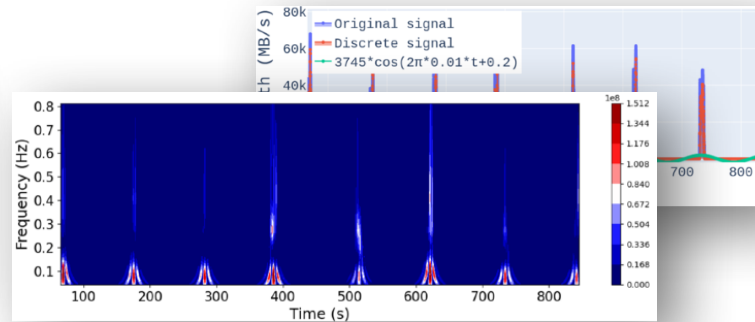
Phase models with Extra-P + FTIO

See: https://admire-eurohpc.eu/wp-content/uploads/2024/02/D5_5_admire.pdf



Classifying Phase with FTIO

GekkoFS and FTIO



Extension: Enhancing predictions

Integration with the monitoring framework Metric Proxy (ISC24)

... and many more!

Contact us for collaborations
ahmad.tarraf@tu-Darmstadt.de

Open-Source Tool



- Available on GitHub: <https://github.com/tuda-parallel/FTIO>
- Python code, can be easily installed: `pip install ftio-hpc`



Users, contributions, and collaborations are welcomed!

```
https://github.com/tuda-parallel/FTIO/blob/main/examples/API/test_api.py
1 import numpy as np
2 from ftio.cli.ftio_core import core
3 from ftio.parse.args import parse_args
4 from ftio.freq_dft import display_prediction
5 from ftio.freq.freq_plot_core import convert_and_plot
6 from ftio.parse.bandwidth import overlap
7
8 # Set up data
9 # 1) overlap for rank level metrics
10 b_rank = [0.0,0.0,1000.0,1000.0,0.0,0.0,1000.0,1000.0,0.0,0.0,1000.0,1000.0,0.0,0.0]
11 t_rank_s = [0.5,0.0,10.5,10.0,20.5,20.0,30.5,30.0,40.5,40.0,50.5,50.0,60.5,60]
12 t_rank_e = [5.0,4.5,15.0,14.5,25.0,24.5,35.0,34.5,45.0,44.5,55.0,54.5,65.0,64.5]
13 b,t = overlap(b_rank,t_rank_s, t_rank_e)
14 # or 2) Directly specify the app level metrics
15 # t = [10.0, 20.1, 30.0, 40.2, 50.3, 60, 70, 80.0,]
16 # b = [10, 0, 10, 0, 10, 0, 10, 0]
17
18 # pass arguments
19 argv = ["-e", "no"] #no plot
20 # set up data
21 data = {
22     "time": np.array(t),
23     "bandwidth": np.array(b),
24     "total_bytes": 0,
25     "ranks": 10
26 }
27 #parse args
28 args = parse_args(argv,"ftio")
29 # perform prediction
30 prediction, dfs = core([data], args)
31 # plot and print info
32 convert_and_plot(data, dfs, args)
33 display_prediction("ftio", prediction)
```


Conclusion

- Presented an approach to characterize and predict the I/O phases of an application with a simple metric: its period, obtained using DFT
- Additional metrics describe the **confidence** in the results and allow for further characterization
- **Online** and **offline** realization
- Demonstrated its usefulness for I/O Scheduling
- Several **parameters** can be changed to enhance the results obtained
- ***Is available on GitHub!***

- Contact: ahmad.tarraf@tu-darmstadt.de

References

1. Anne Benoit, Thomas Herault, Lucas Perotin, Yves Robert, and Frédéric Vivien. 2023. Revisiting I/O bandwidth-sharing strategies for HPC applications. Technical Report RR-9502. INRIA. 56 pages. <https://hal.inria.fr/hal-04038011>
2. Matthieu Dorier, Gabriel Antoniu, Rob Ross, Dries Kimpe, and Shadi Ibrahim. 2014. CALCioM: Mitigating I/O interference in HPC systems through crossapplication coordination. In IPDPS'14. IEEE, 155–164.
3. Emmanuel Jeannot, Guillaume Pallez, and Nicolas Vidal. 2021. Scheduling periodic I/O access with bi-colored chains: models and algorithms. J. of Scheduling 24, 5 (2021), 469–481.

Acknowledgment

- We acknowledge the support of the European Commission and the German Federal Ministry of Education and Research (BMBF) under the EuroHPC Programmes DEEP-SEA (GA no. 955606, BMBF funding no. 16HPC015) and ADMIRE (GA no. 956748, BMBF funding no. 16HPC006K), which receive support from the European Union's Horizon 2020 programme and DE, FR, ES, GR, BE, SE, GB, CH (DEEP-SEA) or DE, FR, ES, IT, PL, SE (ADMIRE). This work was also supported by the French National Research Agency (ANR) in the frame of DASH (ANR-17-CE25-0004), by the Project Région Nouvelle Aquitaine 2018-1R50119 "HPC scalable ecosystem"
- We gratefully acknowledge the computing time provided on the high-performance computer Lichtenberg at the NHR Centers NHR4CES at TU Darmstadt. This is funded by the Federal Ministry of Education and Research, and the state governments participating on the basis of the resolutions of the GWK for national high performance computing at universities (www.nhr-verein.de/unsere-partner). Moreover, we also gratefully acknowledge the computing time on the PlaFRIM experimental testbed supported by Inria, CNRS (LABRI and IMB), Université de Bordeaux, Bordeaux INP and Conseil Régional d'Aquitaine (<https://www.plafrim.fr>)





TECHNISCHE
UNIVERSITÄT
DARMSTADT

Parallel
Programming

**Reproduce all our experiments
from the paper!**



<https://github.com/tuda-parallel/FTIO/tree/main/artifacts/ipdps24>

27 June 2024

Frequency Techniques for I/O | 17th Scheduling for Large-scale Systems Workshop | Dr. Ahmad Tarraf | Department of Computer Science

34

Questions?

Thank you for your attention!