







## Heuristic-Based Merging of HPC Traces

Extending Hardware Counter Coverage for Performance Prediction



Júlia Orteu Aubach

Fabio Banchelli, Marc Clascà and Marta Garcia-Gasulla

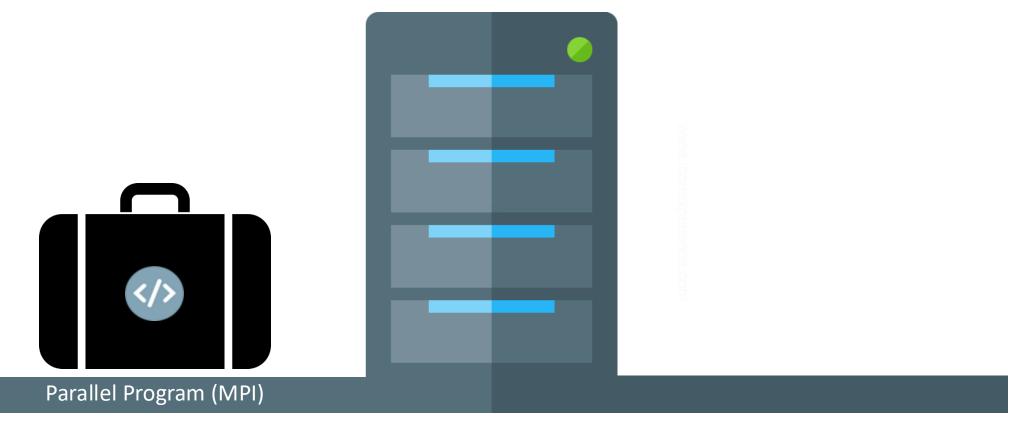
# Contents

- 1. Context
- 2. Objectives
- 3. Methodology
- 4. Experimentation
- 5. Conclusions



High Performance Computing (HPC)

























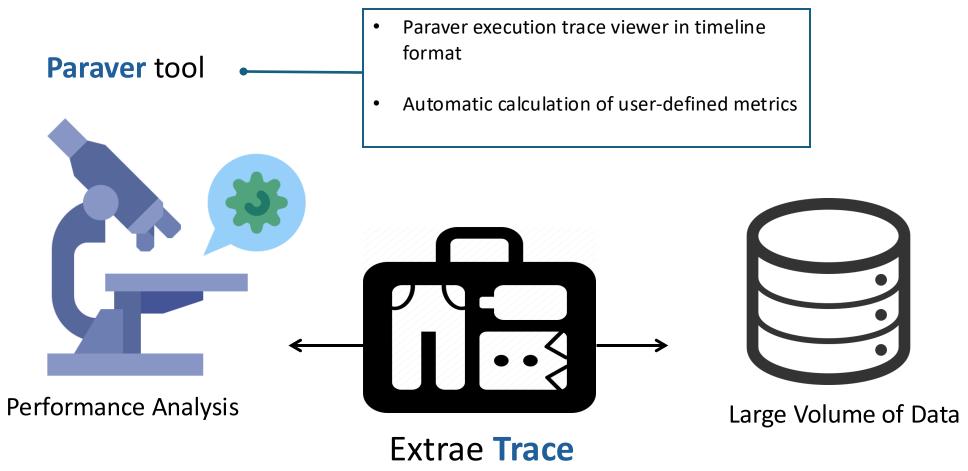
- Instrumentation without code changes or recompilation
- MPI Applications
- C, C++ and Fortran

Parallel Program (MPI)





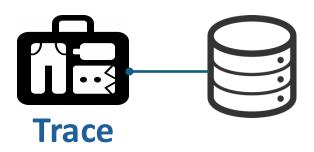












In previous work the objective was to develop and train a machine learning (ML) model using simple parallel programs (benchmarks and kernels) and hardware counters to predict the performance of unseen applications on an HPC machine.

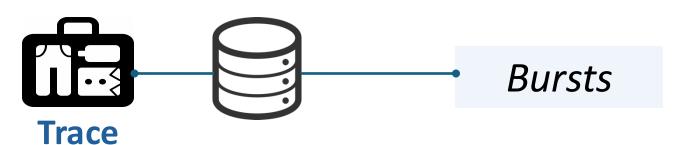
Orteu, J., Clascà, M., Labarta, J., Jennings, E., Andersson, S., Garcia-Gasulla, M. (2025). A Framework and Methodology for Performance Prediction of HPC Workloads. Parallel Processing and Applied Mathematics. PPAM 2024. Lecture Notes in Computer Science, vol 15581. Springer, Cham. <a href="https://doi.org/10.1007/978-3-031-85703-4\_8">https://doi.org/10.1007/978-3-031-85703-4\_8</a>



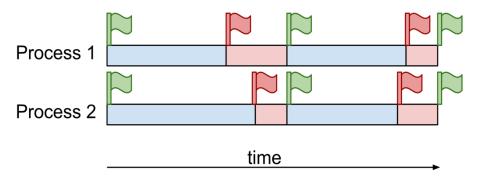


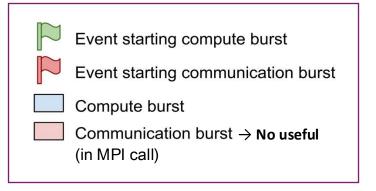


### What type of data?



We define a burst as the time interval between two successive events in a process



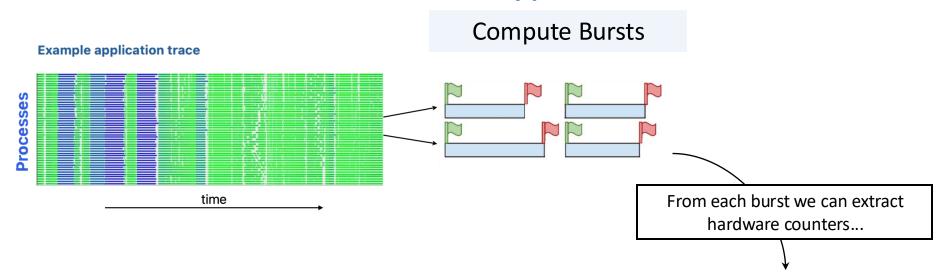








# What type of data?





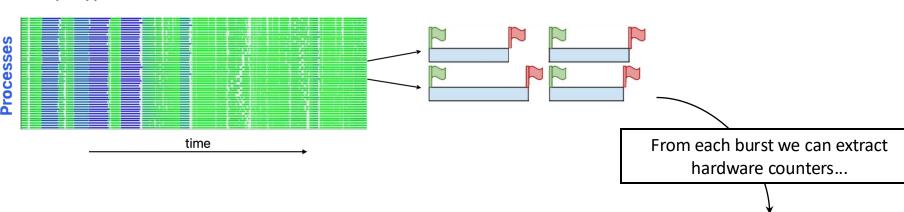




### What type of data?







Performance Application
Programming Interface
(PAPI) provides access to hardware
counters

### PAPI TOT INS

Total Instructions

### PAPI TOT CYC

Total Cycles

### PAPI\_L1\_DCM

L1 Data Cache Misses

### PAPI VEC DP

Vector Double Precision

### PAPI BR INS

Branch Instructions

### PAPI FP OPS

Floating Point Operations







### **Data Format**

We extract a tabular data set for each execution trace.

Row = Bursts

Cols = Trace features

- Timeline information
- MPI contextual information
- Hardware Performance counters

TaskId	Begin_Time	Duration	INS	End_Time	Line	d_IPC	x_PAPI_L1_DCM	x_PAPI_L2_DCM	•••
1	7254812180	71362947	100059606	7.326175e+09	630.0	5.128201	223733.0	185399.0	
1	7326181549	5026	2158	7.326187e+09	1548.0	0.222245	133.0	89.0	
1	7326189542	13978	11780	7.326204e+09	1552.0	0.548392	338.0	178.0	
1	7340619255	6277	3612	7.340626e+09	1668.0	0.481793	167.0	75.0	•••
1	7340635605	1568384	16459717	7.342204e+09	1680.0	5.168772	1098.0	25.0	
		•••		•••	•••				







### What type of data?



- Hundreds of available counters, but only **4-8 can be measured simultaneously** due to architectural constraints
  - This creates an "n-counter ceiling" for our ML models

Performance Application
Programming Interface
(PAPI) provides access to hardware
counters

PAPI\_TOT\_INS

Total Instructions

PAPI\_TOT\_CYC

Total Cycles

PAPI\_L1\_DCM
L1 Data Cache

Misses

PAPI VEC DP

Vector Double Precision PAPI\_BR\_INS

Branch Instructions

PAPI FP OPS

Floating Point Operations







### Existing solutions and their limitations

# Fixed Counter Selection

X Limits feature diversity

# Hardware Multiplexing

- X Introduces timing noise
- X Intermittent sampling









# Objectives

### Research Objectives

- 1. Overcome the n-counter ceiling in HPC performance modeling.
- 2. Maintain burst-level fidelity required for detailed analysis.
- 3. Create synthetic traces compatible with existing tools (Paraver/Extrae).
- 4. Validate methodology across diverse HPC applications.





# Methodology

### **Benchmark Applications**

Application	Domain	Description
SOD2D	Computational Fluid Dynamics	Spectral high-Order code for solving partial differential equations, primarily used for fluid dynamics and wave propagation problems
SeisSol	Seismic Simula- tion	High-performance seismic wave simulation software for earthquake modeling and ground motion prediction
Stream	Memory Bench- mark	Memory bandwidth benchmark measuring performance of four vector operations: Copy, Scale, Add, and Triad
Alya Solver	Computational Mechanics	Mini-app from HPC mechanics application featuring very fine-grained parallelism for computational mechanics problems
Lulesh	Hydrodynamics	Benchmark tool for shock wave simulations in fluid dynamics, fo- cusing on efficient energy calcula- tions for nuclear fission explosions







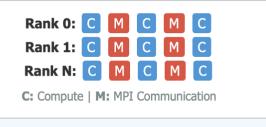
### **Observed Communication Patterns Across Applications**

### 1. Processes with Identical Communication Behavior

### **Characteristics:**

- Consistent structure between different executions
- Each rank may have different burst counts, but this remains constant across runs
- Enables direct, position-based burst matching

**Examples: Stream, Alya, Lulesh** 





**Note:** Different burst counts per rank, but consistent across executions





### **Observed Communication Patterns Across Applications**

### 2. Processes with Structural Variations

### **Characteristics:**

- Different MPI communication patterns between executions
- Variable number of bursts per run (moderate variations)
- Some recurring structures can be identified for pattern-based matching
- Less irregular than highly complex applications, but not fully deterministic

**Example: SOD2D** 



C: Compute | B: Barrier | A: Allreduce | F: Comm\_Free Note: Different MPI calls and burst counts between runs







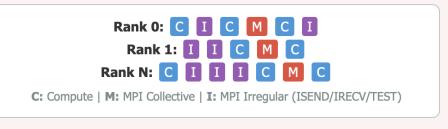
### **Observed Communication Patterns Across Applications**

### **3. Highly Irregular Traces**

### **Characteristics:**

- Complex and heterogeneous communication behavior
- Significant structural irregularities
- Requires sophisticated matching with structural constraints

**Example: SeisSol** 



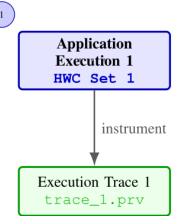
These observed patterns guided the design of our **two-stage matching algorithm**, which progressively increases complexity based on structural similarity detection.

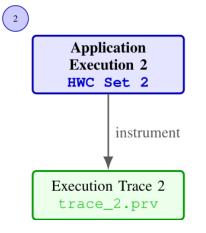


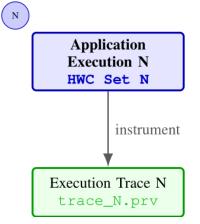




The methodology scales to accommodate multiple executions with diverse counter configurations, limited only by the available hardware counter combinations

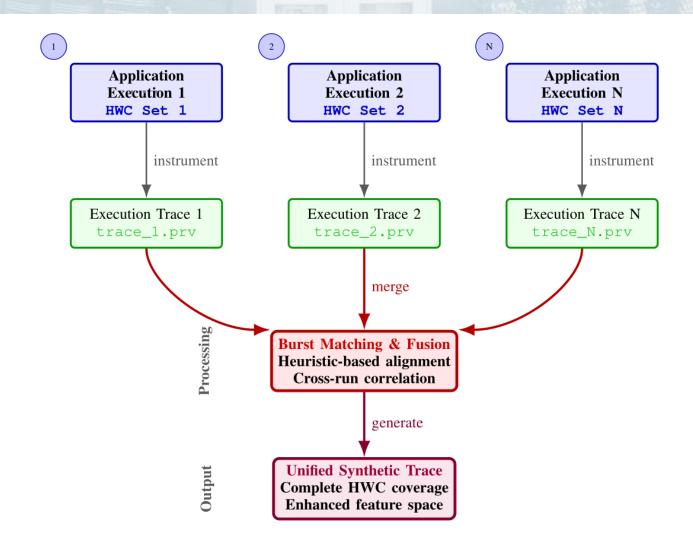






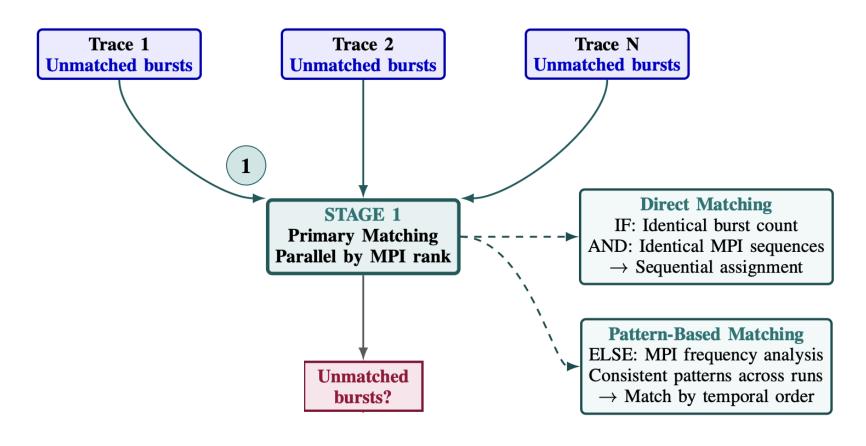








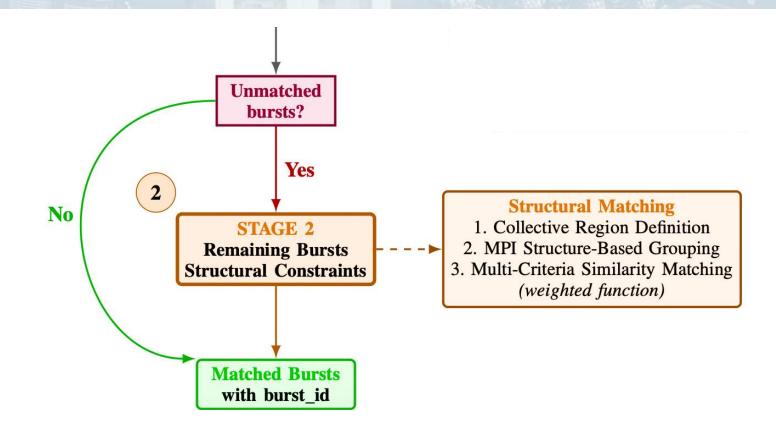


















### **Structural Matching**

- 1. Collective Region Definition
- 2. MPI Structure-Based Grouping
- 3. Multi-Criteria Similarity Matching (weighted function)

 $S = 0.6 \times D_{temporal} + 0.2 \times D_{size} + 0.2 \times D_{partner}$ 

### **Temporal Position**

Relative position within collective regions

Weight: 60%

### **Communication Size**

Message sizes for MPI operations

Weight: 20%

### **Communication Partner**

Rank identifiers for point-topoint

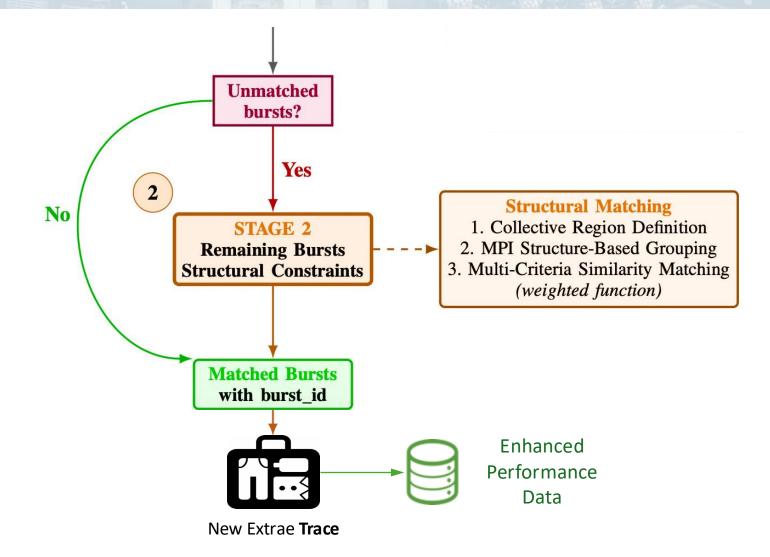
Weight: 20%

**Matching Threshold:** S < 0.3 for genuine similarity















# Experimentation

### MareNostrum5 Configuration

### **Execution Parameters**

Processes: 100-112 MPI processes

•Repetitions: 2-10 runs per

configuration

•Noise Reduction: System isolation

### **Hardware & System**

•Platform: Intel Sapphire Rapids CPUs

•Frequency: Fixed at 2 GHz

Allocation: Exclusive node access





### Three Different Hardware Counter Configurations

### **INS\_MIX**

**Focus:** Instruction mix & cache hierarchy

PAPI\_TOT\_INS, PAPI\_TOT\_CYC,
PAPI\_LD\_INS, PAPI\_SR\_INS, PAPI\_BR\_INS,
PAPI\_L3\_TCM, PAPI\_L1\_DCM, PAPI\_L2\_DCM

### **OPS\_SET**

**Focus:** Floating-point operations & vectorization

PAPI\_TOT\_INS, PAPI\_VEC\_INS, PAPI\_FP\_INS, PAPI\_FP\_OPS, PAPI\_DP\_OPS, PAPI\_SP\_OPS, PAPI\_VEC\_DP

### OPS\_CYC

**Focus:** Computational performance

PAPI\_TOT\_INS, PAPI\_TOT\_CYC,
PAPI\_VEC\_DP, PAPI\_VEC\_SP, PAPI\_DP\_OPS







### Applications with Identical Communication Behaviour

Rank 0: C M C M C

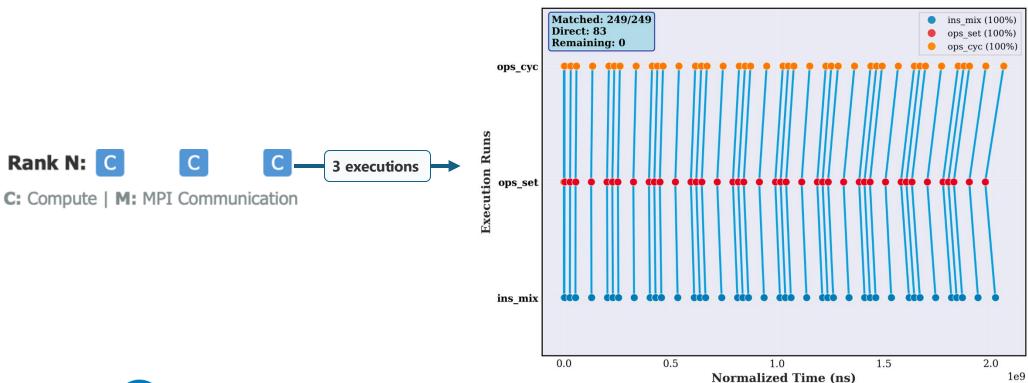
Rank 1: C M C M C

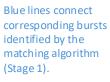
Rank N: C M C M C

C: Compute | M: MPI Communication



### Applications with Identical Communication Behaviour









### **Applications with Identical Communication Behaviour**

**Stream** 

100%

9,296 bursts | 112 MPI ranks

**Direct Matching** 

Memory bandwidth benchmark with identical burst patterns

Alya

100%

16,428 bursts | 48 MPI ranks

**Direct Matching** 

Computational mechanics with deterministic structure

Lulesh

100%

43,077 bursts | 27 MPI ranks

**Direct Matching** 

Hydrodynamics with consistent MPI patterns







Processes with Estructural Variations Between Executions

Execution 1 - Rank 0: C B C A C

Execution 2 - Rank 0: C B C F C A

C: Compute | B: Barrier | A: Allreduce | F: Comm\_Free Note: Different MPI calls and burst counts between runs

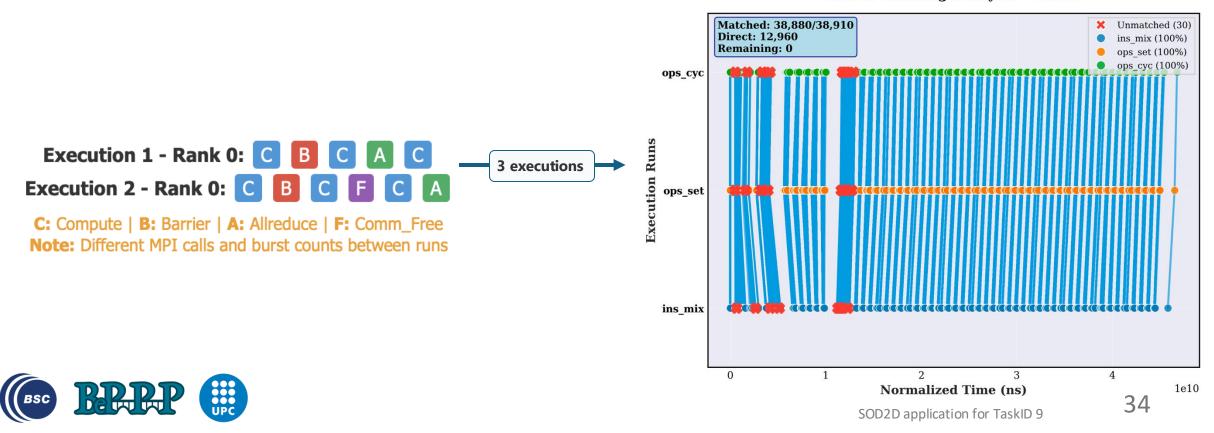




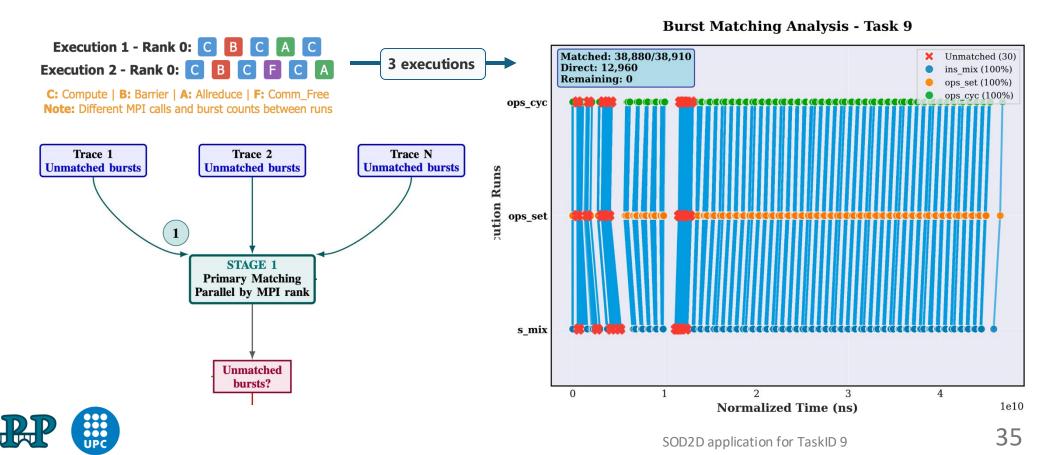


### **Processes with Estructural Variations Between Executions**

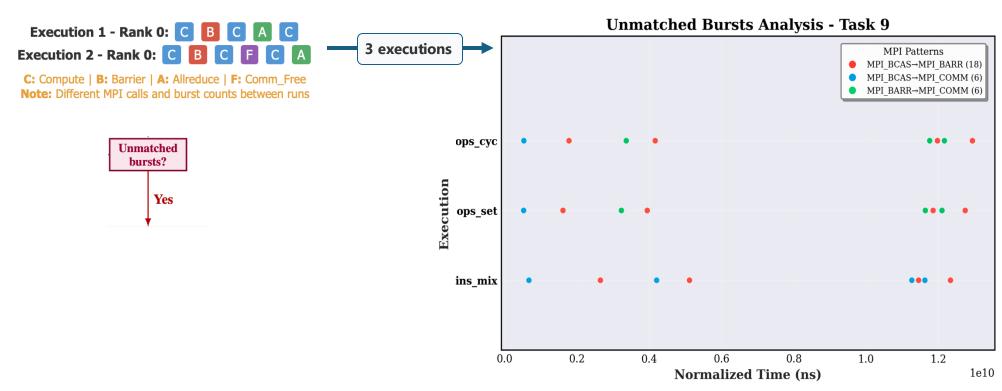
### **Burst Matching Analysis - Task 9**



### Processes with Estructural Variations Between Executions



#### Processes with Estructural Variations Between Executions

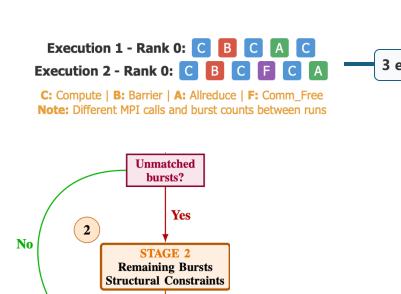




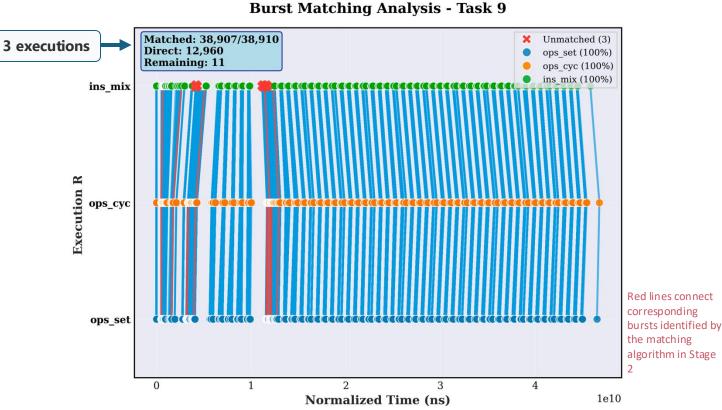




#### Processes with Estructural Variations Between Executions



Matched Bursts with burst\_id







#### Processes with Estructural Variations Between Executions

SOD2D

99.99%

~1.5M bursts total

**Pattern-Based** 

CFD with minor structural variations (0.02% difference)

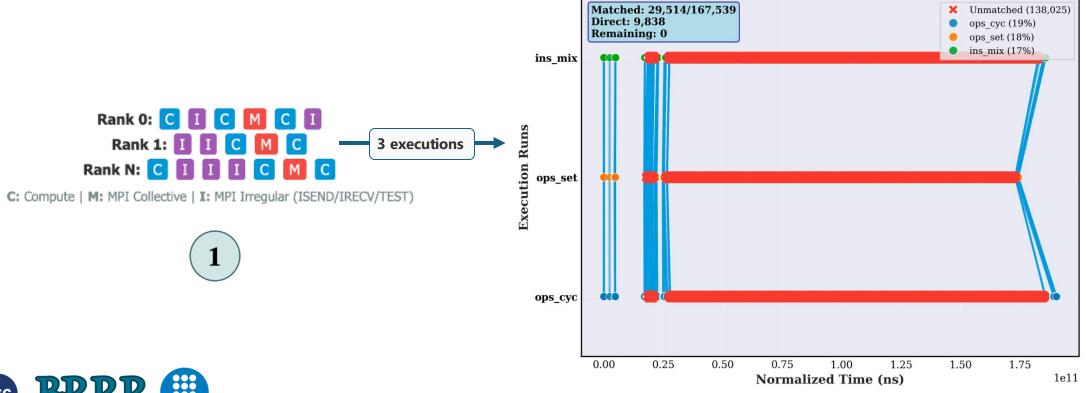






#### Highly Irregular Traces

#### **Burst Matching Analysis - Task 9**



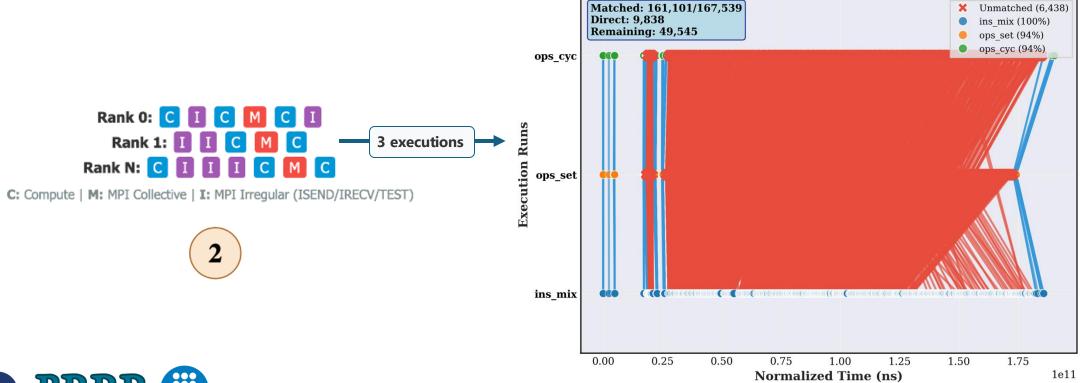






#### Highly Irregular Traces

#### **Burst Matching Analysis - Task 9**









#### Highly Irregular Traces

#### SeisSol

90-99%

>5.5M bursts | 128 MPI patterns

#### **Complex Matching**

Seismic simulation with significant nondeterminism







#### Validation Framework

Apply the same heuristic matching algorithm to N traces with **identical counter sets** to quantify matching precision and establish baseline accuracy.



#### **Controlled Setup**

N executions with identical counter sets per application



#### Algorithm Application

Apply matching heuristics to identify burst correspondences



## **Quality Assessment**

Calculate three validation metrics for matched bursts



## Statistical Analysis

Correlation matrices, distributions, and visual validation









#### Quality Assessment

#### **Pearson Correlation**

Correlation coefficient between base trace and matched traces

#### **Mean Absolute Error**

Absolute deviation in original counter units

$$MAE = \frac{1}{M} \sum_{i=1}^{M} |b_i - \mu_i|,$$

where M is the number of the matched bursts

 $b_i$  base value for burst i

$$\mu_i = \frac{1}{N-1} \sum_{j \neq base} t_{j,i}$$

#### **Relative Difference**

Relative deviation for non-zero base values

< 30% Acceptance Rate

RelDiff = 
$$\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \frac{|b_i - \mu_i|}{b_i}$$
,  $\mathcal{I} = \{i : b_i > 0\}$ .

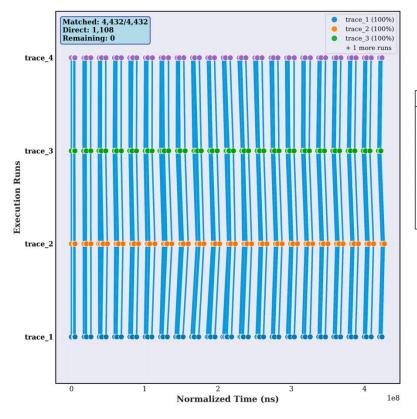






#### Applications with Identical Communication Behaviour

Validation approach using Lulesh with identical counter sets extracted from 4 different executions.



#### LULESH INS\_MIX VALIDATION

Counter	Correlation	MAE	Rel Diff	< 30% Diff
PAPI_TOT_INS	1.000	39.789	0.011	100.0%
PAPI_SR_INS	1.000	6.914	0.014	100.0%
PAPI_LD_INS	1.000	12.193	0.011	100.0%
PAPI_BR_INS	1.000	4.690	0.007	100.0%
PAPI_TOT_CYC	1.000	349.238	0.081	96.7%
PAPI_L1_DCM	1.000	12.651	0.204	76.3%
PAPI_L2_DCM	1.000	25.171	0.390	61.7%
PAPI_L3_TCM	0.973	11.387	0.499	51.6%

For the OPS SET and OPS CYC counter sets, detailed validation tables are not presented since these applications achieved **perfect deterministic matching** with no duplicate columns generated during trace fusion.







## Processes with Estructural Variations Between Executions

#### SOD2D INS\_MIX VALIDATION RESULTS

Counter	Correlation	MAE	Rel Diff	< 30% Diff
PAPI_L2_DCM	1.000	2.729	0.512	50.8%
PAPI_L1_DCM	0.999	7.165	0.214	76.1%
PAPI_L3_TCM	0.999	2.000	0.488	51.3%
PAPI_SR_INS	0.963	7.387	0.010	99.9%
PAPI_LD_INS	0.876	13.693	0.009	99.7%
PAPI_BR_INS	0.875	6.258	0.005	99.7%
PAPI_TOT_INS	0.876	20.070	0.003	99.7%

#### Highly Irregular Traces

SEISSOL INS\_MIX VALIDATION

Counter	Correlation	MAE	Rel Diff	< 30% Diff
PAPI_SR_INS	0.994	4.942	0.015	97.2%
PAPI_LD_INS	0.993	11.373	0.015	97.3%
PAPI_TOT_INS	0.993	41.992	0.019	97.3%
PAPI_L1_DCM	0.991	4.738	0.669	38.5%
PAPI_BR_INS	0.989	4.725	0.008	97.4%
PAPI_L2_DCM	0.980	0.885	0.092	87.2%
PAPI_L3_TCM	0.976	1.129	0.174	75.7%

SEISSOL OPS\_SET VALIDATION RESULTS

Counter	Correlation	MAE	Rel Diff	< 30% Diff
PAPI_VEC_SP	1.000	0.000	0.000	100.0%
PAPI_SP_OPS	1.000	0.000	0.000	100.0%
PAPI_VEC_DP	0.993	0.000	0.000	99.8%
PAPI_VEC_INS	0.993	0.000	0.000	99.8%
PAPI_DP_OPS	0.993	0.000	0.000	99.4%
PAPI_FP_OPS	0.993	0.000	0.000	99.4%
PAPI_FP_INS	0.992	0.000	0.000	99.4%







# Conclusions and Future Work



## Conclusions

#### Research Objectives

1. Overcome the n-counter ceiling in HPC performance modeling.



2. Maintain burst-level fidelity required for detailed analysis.



3. Create synthetic traces compatible with existing tools (Paraver/Extrae).



4. Validate methodology across diverse HPC applications.









### Conclusions

#### Excellent Performance

- Instruction counters: >97% acceptable correspondence
- Floating-point operations: Perfect precision even in non-deterministic scenarios

#### Variable Performance

- Cache hierarchy: 38.5-87.2% acceptance rates
- 1
- Treat cache data as supplementary evidence requiring interpretation





## **Future Work**

- 1. Training performance prediction models on expanded feature spaces and expand our prior work on performance analysis of HPC workloads.
- 2. Production deployment requires evaluation across broader application domains and HPC architectures to confirm generalization.





## Thank you!



