

Redefining HPC Observability: Integrating Monitoring, Modeling, and Meaning

Sarah M. Neuwirth Johannes Gutenberg University Mainz, Germany <u>neuwirth@uni-mainz.de</u>

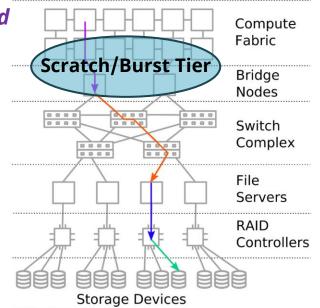
MODA 2025 Workshop, ISC High Performance, June 2025

Motivation

Heterogeneous and Complex HPC Infrastructures

- HPC infrastructure too complex, humans are overwhelmed
- Complexity and scope increase the *urgency*
 - <u>New computational paradigms</u> (AI/ML apps vs. BSP-style HPC)
 - <u>New architectural directions</u> (e.g., IPU, RISC-V, data flow)
 - <u>Heterogeneity overall</u>: node architectures, within the system, storage and parallel file system during application design (e.g., ML within HPC applications)
 - <u>New operations paradigms</u> (e.g., cloud, container)
 - Simplistic approaches to increasing compute demand result in <u>unacceptable power costs</u>
- Difficult for humans to optimally adapt applications to systems and to detect and diagnose vulnerabilities

Carns, P., 2023. *HPC Storage: Adapting to Change*. Keynote at REX-IO'23 Workshop. Ciorba, F., 2023. *Revolutionizing HPC Operations and Research*. Keynote at HPCMASPA'23 Workshop.



B. Settlemyer, G. Amvrosiadis, P. Carns and R. Ross, 2021. *It's Time to Talk About HPC Storage: Perspectives on the Past and Future*, in Computing in Science & Engineering, vol. 23, no. 6, pp. 63-68.



Motivation *Why Redefine Observability?*





HPC is no longer just HPC

- From simulation → simulation + AI + data analytics
- Complex workflows across diverse compute, storage and network layers
- Traditional metrics (e.g., FLOPS, bandwidth) do not capture this complexity



Monitoring ≠ Observability

- Monitoring tells what happened
- Observability helps us understand **why** it happened and what to do next
- We need more than logs and counters we need explainability and context



New Questions need new tools

- Why is my workflow slow now, but not yesterday?
- Is the problem in the application, the file system, or the kernel?
- Can I trust this model's I/O behavior across systems?

<mark>෬</mark>…ဉ ; ; ෭෯; ෫

Actionable Insight Requires Integration

- No single tool sees the full stack
- We need to connect: instrumentation \rightarrow modeling \rightarrow decision-making







Monitoring: The Foundation

Vision: Holistic Monitoring for Intelligent HPC Operations

What We Need:

- Continuous monitoring, archiving, and analysis of operational + performance data
- Unified visibility into applications, system software, and hardware layers

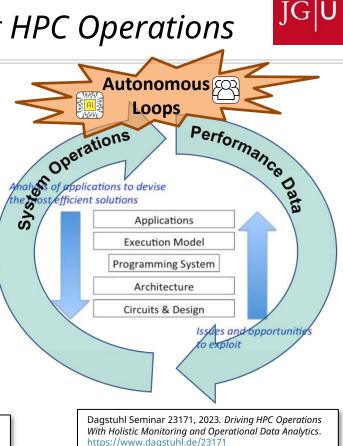
Why It Matters:

- Enables automated feedback loops using AI/ML
- Supports dynamic workload & architecture analysis
- Powers adaptive, actionable responses

Goal: Efficient, explainable HPC operations driven by *autonomous analyze-feedback-response loops*

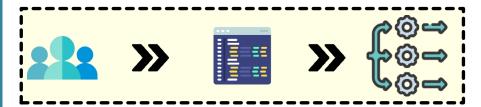
Gentile, A., 2021. *Enabling Application and System Data Fusion*. Keynote at MODA'21 Workshop.

Ciorba, F., 2023. *Revolutionizing HPC Operations and Research*. Keynote at HPCMASPA'23 Workshop.



Monitoring: The Foundation *Different Perspectives*





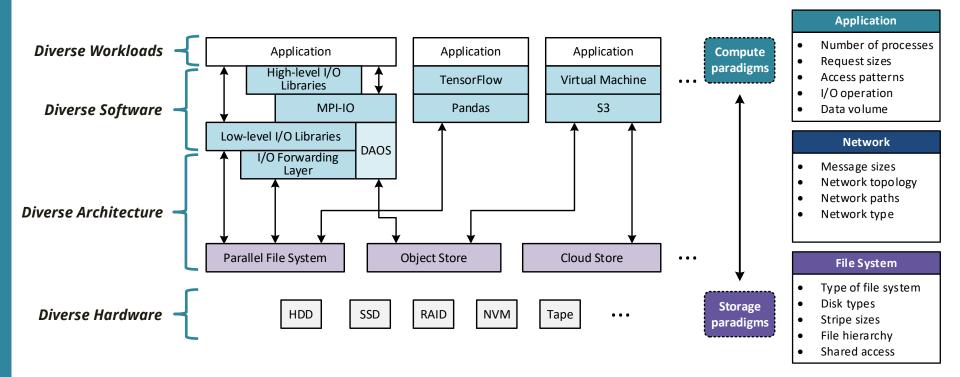
User Interests and Concerns:

- Ease of use
- Application performance
- Portability
- Reproducible science
- Accessibility of the system
- Data persistence
- Core hour usage

System Interests and Concerns:

- Installation, configuration, and operation of the production-ready system, e.g.:
 - Software requirements
 - System configuration
 - High availability service
- System monitoring
- System security
- Benchmarking and anomaly detection

Monitoring: The Foundation Parallel I/O System and Performance Factors



G

Monitoring: The Foundation

Example: Darshan I/O Characterization Tool



Blue Waters, Mira, and Theta popular Darshan log sources used for research:

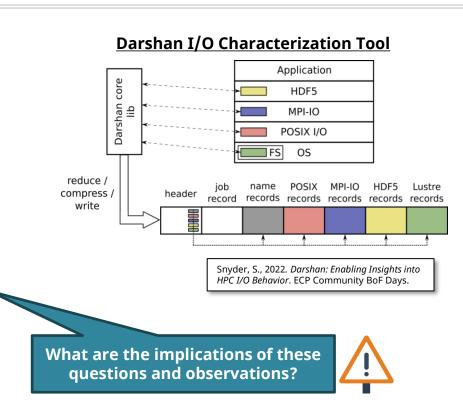
- <u>https://bluewaters.ncsa.illinois.edu/data-sets</u>
- <u>https://reports.alcf.anl.gov/data/</u>
- <u>ftp://ftp.mcs.anl.gov/pub/darshan/data</u>

Open questions:

- How relevant are the logs to current systems?
- How do we know the integrity of the logs?

Community comments:

- "Darshan is one of the first tools to be deactivated in the event of I/O problems."
- "Darshan cannot grasp the complexity of state-ofthe-art parallel storage systems."



Monitoring: The Foundation *Fragile Pipelines – Tool-Driven, Not Insight-Driven*



Category	Examples	Strengths	Limitations
Application-level	Darshan, Recorder, Score-P	Fine-grained function tracing	No visibility into system-wide interactions
System-level	LDMS, DCDB, TACCStats	Aggregated I/O performance metrics	Cannot correlate application performance with system metrics
End-to-end	Ganglia, Nagios, Apollo	Holistic view of system utilization	Lacks deep profiling at kernel and network levels

- Siloed Tool Views: Each tool sees a layer. None explain the whole system.
 => For example, in case of I/O: App-level profilers (e.g., Darshan) vs. system tools (LDMS, DCDB)
- Workflow Scripts Instead of Workflows: Custom scripts per experiment = unscalable, unrepeatable. => Benchmarking tools (e.g., iperf, sockperf) often require client/server logic incompatible with SLURM
- Discarded Insights: Performance data is ephemeral; models are not reused.
 => No structure for reuse → repeated effort, lost opportunities

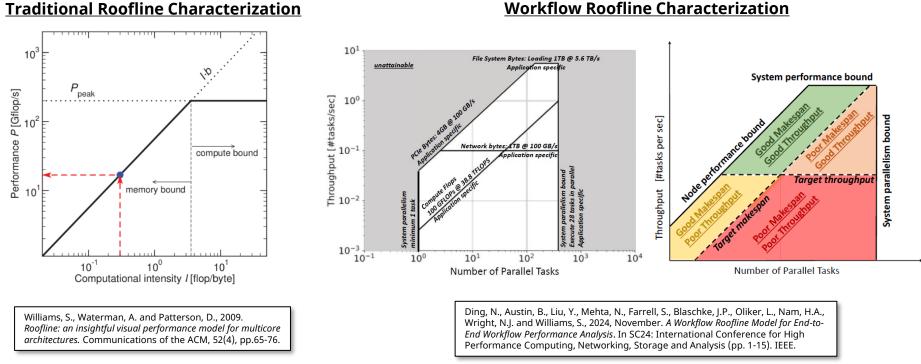






Modeling: From Data to Understanding Roofline Characterization





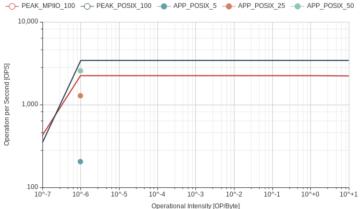
Traditional Roofline Characterization

Modeling: From Data to Understanding *I/O Roofline Characterization: Initial Concept*

- Traditional Roofline Model...
 - is based on looking at the <u>relationship between work and traffic</u>
 - provides intuitive approach through simple bound and bottleneck analysis
- <u>I/O Roofline Characterization</u> is based on IOPS and the I/O bandwidth => <u>I/O interface specific</u>, i.e., POSIX, MPIIO, etc.
 - <u>IOPS</u>: number of reads and writes that a storage system can perform per second
 - <u>Bandwidth</u>: total amount of data read or written per second
- **X-axis:** *I/O Operational Intensity* = $\frac{\text{Total I/O Operations}}{\text{Read Bytes+Write Bytes}}$
- Y-axis: P = min(Peak IOPS, Peak I/O Bandwidth × I/O Intensity) where P is the attainable perf., P_{peak} is the peak perf., b is the peak bandwidth, and I is the arithmetic intensity

Zhu, Z., Bartelheimer, N. and Neuwirth, S., 2023. *An Empirical Roofline Model for Extreme-Scale I/O Workload Analysis*. IPDPSW'23.

> Zhu, Z., and Neuwirth, S., 2023. *Characterization* of Large-Scale HPC Workloads With Non-Naïve I/O Roofline Modeling and Scoring. ICPADS'23.



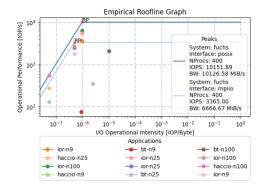


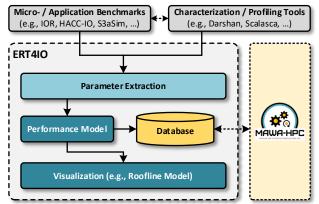
Modeling: From Data to Understanding *I/O Roofline Characterization: Workflow Implementation*



- <u>ERT4IO (Empirical Roofline Tool for I/O)</u> provides automated
 I/O Roofline characterization
 - Parameter extraction from Darshan logs
 - Generates Roofline visualization
- Forwards results to MAWA-HPC framework
 - Enables further data analysis
 - Preserves and shares knowledge via performance history database with the HPC community
- Applicable to various use cases
 - Systems with different configurations and hardware can be compared and evaluated
 - Intuitive estimation of an application's I/O performance
 - Identifying performance bottlenecks and anomalies

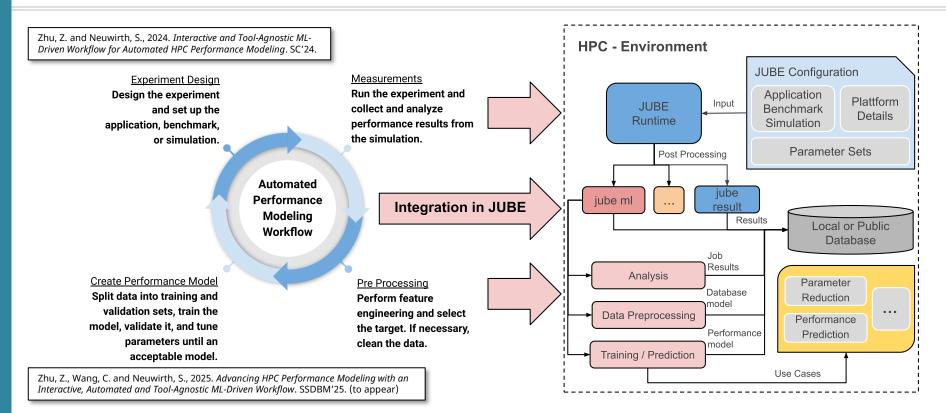






Modeling: From Data to Understanding *Automated ML-Driven Performance Modeling Workflow*





Modeling: From Data to Understanding JUBE-ML Prototype Implementation

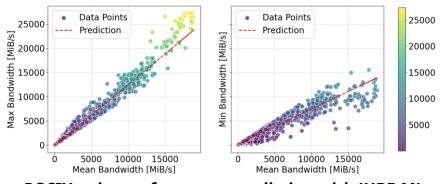
- Prototype extends the JUBE framework with support for automatic ML-based performance modeling and supports variety of ML algorithms
- JUBE-ML is enhanced by the sqlite-ml extension
- JUBE-ML stage provides *four key functionalities*:
 - Data analysis: insights into performance results (e.g., min, max, mean) and lambda functions
 - Data preprocessing: cleans and filters data for ML, creates new tables with selected features & targets
 - *ML model training:* applies ML models (regression or classification) to build a performance model
 - Prediction: validates the model & enables different analysis scenarios, such as identifying irrelevant parameters and predicting system performance



Modeling: From Data to Understanding JUBE-ML Case Study: I/O Bandwidth Prediction

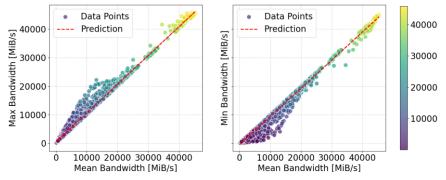


- <u>Use Case</u>: I/O bandwidth modeling and prediction to demonstrate the JUBE-ML workflow
 - IOR benchmark to simulate various I/O workloads and generate performance data
 - Linear regression: 25% for training, 75% for validation
- Despite the small training sample, the model predicts both minimum & maximum bandwidths well



POSIX write performance prediction with JUBE-ML.

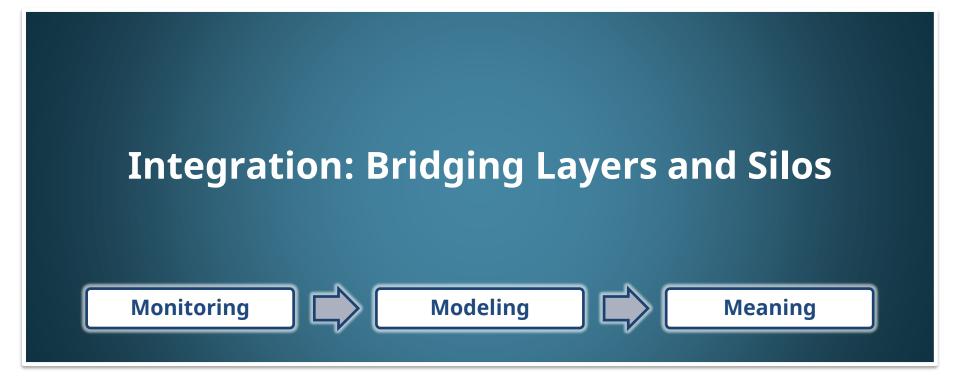
Block size	64MB, 256MB, 512MB	
Transfer size	1MB, 2MB, 4MB, 8MB	
Lustre striping count	0, 2, 4, 8	
Lustre striping size	1MB, 2MB, 4MB, 8MB	
Nodes	1, 2, 4, 8	
Tasks per node	1, 2, 4, 8	



POSIX read performance prediction with JUBE-ML.







Integration: Bridging Layers and Silos *Goal: Holistic & Automated Monitoring and Analysis Cycle*



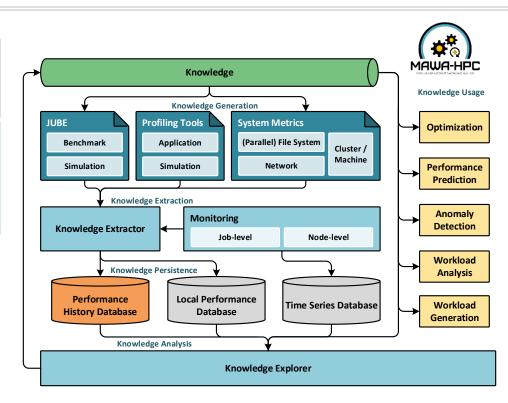
Idea: Develop and implement standardized and tool-independent approach for HPC workload and application analysis

Goal: Establish a *performance history database* to categorize systems, workload behaviors, and characteristic patterns for different science domains

Zhu, Z., Bartelheimer, N. and Neuwirth, S., 2023. *MAWA-HPC: Modular and Automated Workload Analysis for HPC Systems*. ISC'23.

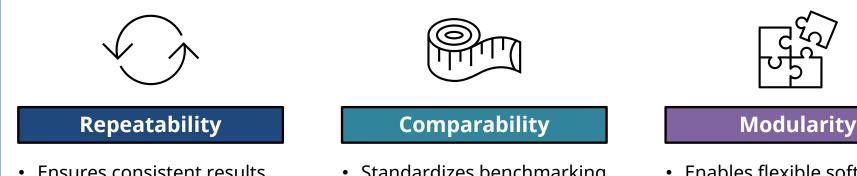
Bartelheimer, N., Zhu, Z., and Neuwirth, S., 2023. Toward a Modular Workflow for Network Performance Characterization. IPDPSW'23.

Zhu, Z., and Neuwirth, S., 2023. *Characterization of Large-Scale HPC Workloads With Non-Naïve I/O Roofline Modeling and Scoring.* ICPADS'23.



Integration: Bridging Layers and Silos

Reproducible Benchmarking: Integration is Key



- Ensures consistent results from the same setup
- Builds trust in computational outcomes
- Forms the basis for scientific validation

Schifrin, A., 2023. Automated Performance Characterization of HPC Systems. Bachelor thesis, Goethe University Frankfurt.

Standardizes benchmarking across models or methods

- Allows fair evaluation of new approaches
- Helps identify performance • trade-offs

Bartelheimer, N. and Neuwirth, S., 2023, Toward Reproducible Benchmarking of PGAS and MPI Communication Schemes. ICPADS'23.

- Enables flexible software architecture
- Facilitates integration of AI/ML components
- Supports code reuse and maintainability

Zhu, Z., Wang, C. and Neuwirth, S., 2025. Advancing HPC Performance Modeling with an Interactive, Automated and Tool-Agnostic ML-Driven Workflow, SSDBM'25, (to appear)

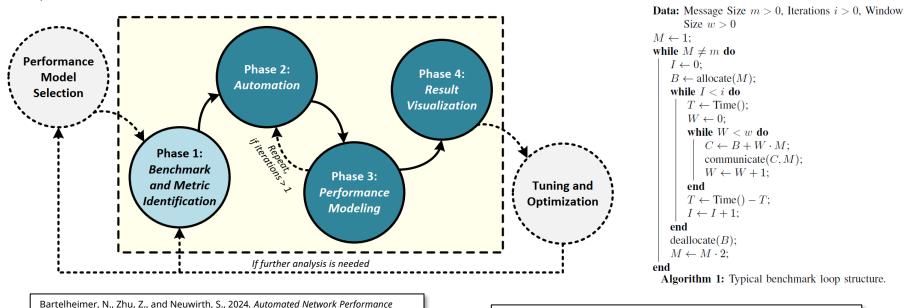




Integration: Bridging Layers and Silos *Workflow Design for Reproducible Benchmarking*



Core component of the benchmarking framework is the JUBE Benchmarking Environment



Characterization for HPC Systems. International Journal of Networking and Computing.

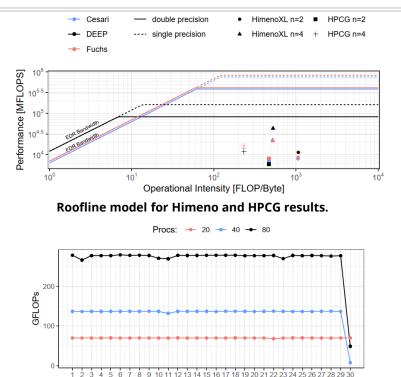
JUBE Documentation: https://apps.fz-juelich.de/jsc/jube/jube2/docu/index.html

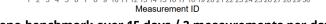
Integration: Bridging Layers and Silos *Reproducible Benchmarking: Example Configuration*



Benchmark-independent, Platform-specific platform.xml	<pre><parameterset name="systemParameter"> <parameter name="nodes" type="int">2</parameter> <parameter name="taskspernode" type="int">1</parameter> <parameter name="threadspertask" type="int">1</parameter> <parameter name="threadspertask" type="int">1</parameter> <parameter mode="python" name="tasks" type="int">\$nodes * \$taskspernode</parameter> <parameter name="timelimit">00:30:00</parameter> 00:30:00 00:30:00 </parameterset></pre>		<pre><parameterset name="executeset"> <parameter name="submit">sbatch</parameter> <parameter name="submit_script">submit_job</parameter> <parameter name="starter">srun</parameter> </parameterset></pre>
Benchmark-specific, Platform-independent likwid-specs.xml	<pre><parameterset name="copy_params"> <parameter name="benchmark_set_cp" tag="copy"></parameter></parameterset></pre>	ameter> <pre></pre>	be_pat_int
Benchmark-specific, Platform-specific platform-likwid- specs.xml	<pre><pre><pre><pre><pre><pre><pre>cparameter name="thread_domain_params"> <pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre>		meter can be changed and represents the parameter part after the <size> uffix must start with the ":" (colon) character</size>

Integration: Bridging Layers and Silos *Automated Performance Characterization*

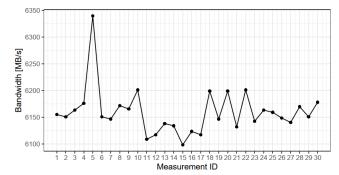




Himeno benchmark over 15 days / 2 measurements per day.

Server Server

Heat map of the allocated nodes (overall benchmark runs).



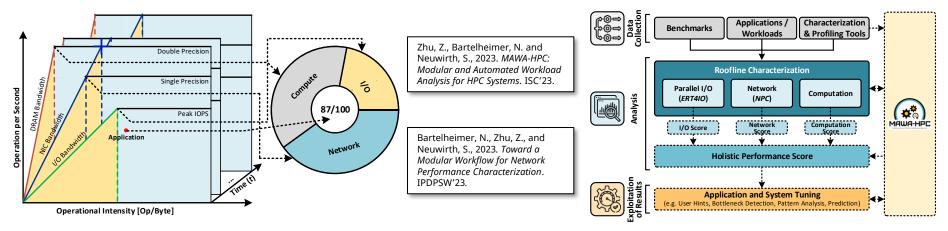
RDMA point-to-point performance over 15 days / 2 measurements per day.

Redefining HPC Observability: Integrating Monitoring, Modeling, and Meaning • ©Sarah M. Neuwirth • Johannes Gutenberg University Mainz

22

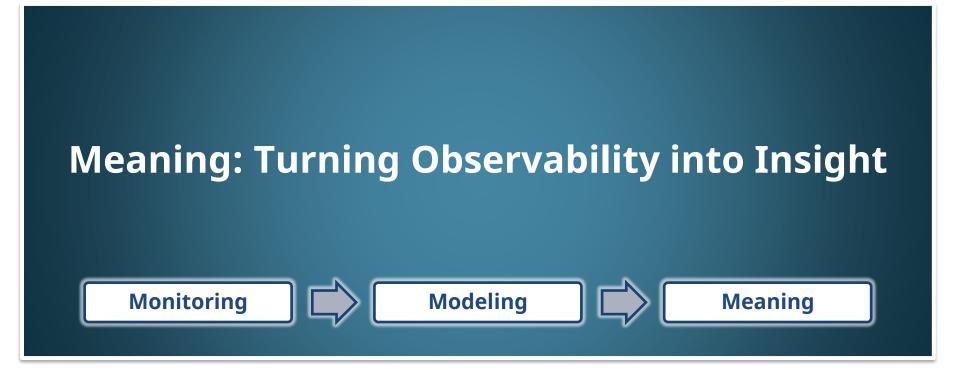
Integration: Bridging Layers and Silos *Multi-dimensional Performance Modeling*

- JG
- <u>Goal</u>: provide a comprehensive view of application and system performance ⇒*emerging workloads*
- Multi-dimensional performance models, for example Roofline model, to account for multiple performance factors (e.g. network, compute power, and parallel I/O)
- Including time as an additional dimension, the Roofline model can provide insight into an application's performance over time, enabling the identification of performance anomalies









Meaning: Turning Observability into Insight *Toward eXplainable I/O*



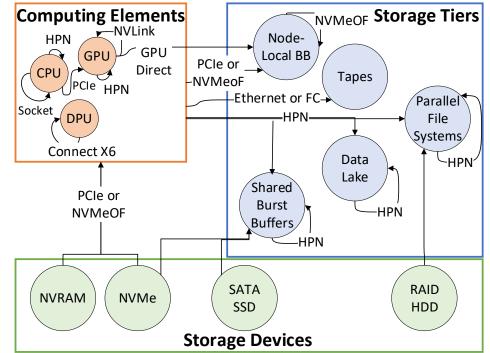


Neuwirth, S. and Devarajan, H., Wang, C., and Lofstead, J.., 2025. XIO: Toward eXplainable I/O for HPC Systems. SSDBM'25. (to appear)

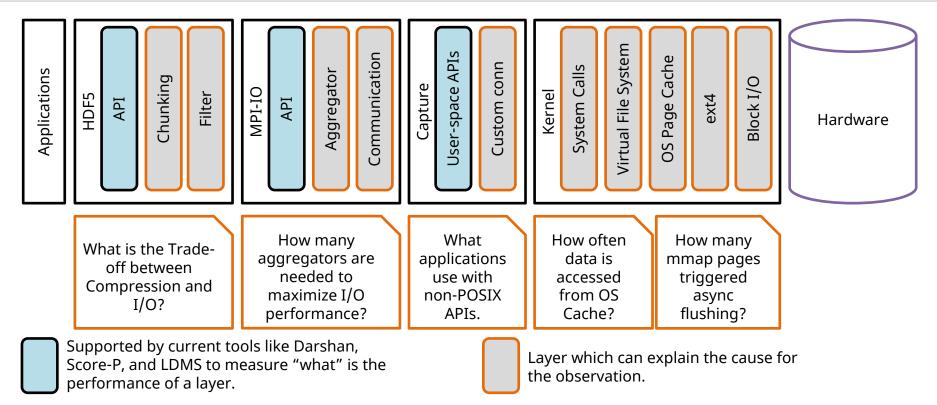
Meaning: Turning Observability into Insight Master Architectural Plan (MAP)



- Derived from current TOP500 and IO500
- Primary objective of MAP is to provide a <u>standardized yet adaptable blueprint</u> to align software and monitoring tools across various HPC configurations
- <u>Graph-based representation</u> of components and interconnects for modeling data flows in HPC systems
- Emphasizes data flow, not just system layout
 => <u>connects observation points across the stack</u>



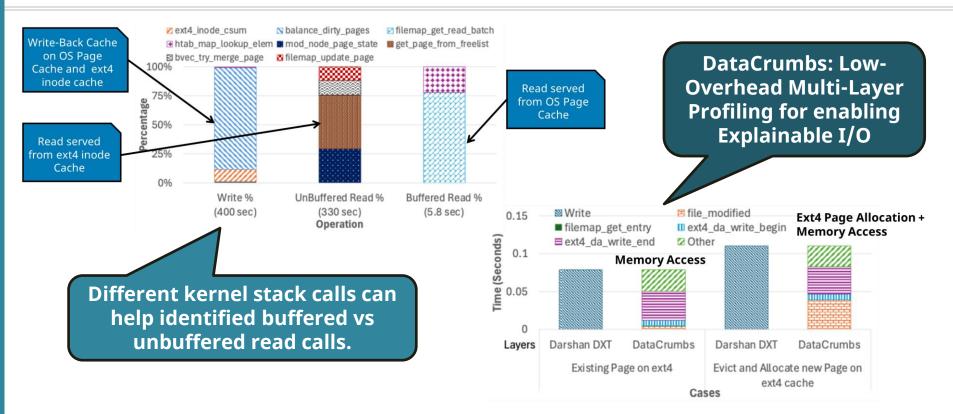
Meaning: Turning Observability into Insight Observability vs. I/O Performance Variability



Redefining HPC Observability: Integrating Monitoring, Modeling, and Meaning • ©Sarah M. Neuwirth • Johannes Gutenberg University Mainz

JG|U

Meaning: Turning Observability into Insight *Toward better Observability with DataCrumbs*



Meaning: Turning Observability into Insight Verifying Performance Meaningfully, Beyond Metrics



Program Trace File - Rank 0 MPI File open(MPI COMM WORLD, ...); int fd = open("./test", O RDWR); MPI Comm comm = MPI COMM WORLD; MPI Info info = MPI INFO NULL: MPI File write at(fh, 0, &data, ...); MPI File fh; pwrite(fd, buf, 4, 0); MPI Status status; MPI_File_sync(fh); MPI File open(comm, "./test", fsync(fd); MPI MODE RDWR, info, &fh); MPI Barrier(MPI COMM WORLD); if (rank == 0) { MPI File close(fh); int data = 7; close(fd); MPI File write at(fh, 0, &data, 1, MPI INT, &status): MPI_File_sync(fh); Trace File - Rank 1 MPI Barrier(comm); if (rank == 1) { int fd = open("./test", O RDWR); int data: MPI File sync(fh); MPI File sync(fh); fsvnc(fd); MPI Barrier(comm); MPI Barrier(MPI COMM WORLD); MPI File read at (fh, 0, &data, 1,

MPI INT, &status); MPI File close(&fh);

MPI File open (MPI COMM WORLD, ...); MPI File read at (fh, 0, &data, ...); pread(fd, buf, 4, 0); MPI File close(fh); close(fd);

Step 1: Generating Execution Trace

	Step 2: Detecting Conflic	ts				
	Conflict	Conflicts				
	<rank 0:="" pwrite(".="" t<br="">Rank 1: pread("./te</rank>					
•	Step 3: Establishing Happens-before Order					
	MPI_File_open +	MPI File open				
	open	open				
	MPI_File_write_at	MPI_File_sync				
Ξ.	↓ /					
	pwrite	fsync				
	· · · · · · · · · · · · · · · · · · ·	↓				
	MPI_File_sync	MPI_Barrier				
		/ · · ·				
	fsync	MPI_File_read_at				
		`				
	MPI_Barrier	pread				
		\				
	MPI_File_close +	MPI_File_close				
	↓	¥				
	close	close				

Step 4: Verifying Consistency Semantics Consistency Properly Synchronized Semantics POSIX \checkmark Commit \checkmark Х Session Х MPI-IO

Explanation:

The execution is properly synchronized under POSIX because there exists a path between the two conflicting operations, suggesting pwrite happens-before pread, which is sufficient for POSIX consistency

Additionally, this path contains a commit (fsync) operation, which satisfies the Commit semantics requirement.

However, there is no close-to-open pair and sync-barrier-sync construct in this path, thus it is not properly synchronized under Session consistency and MPI-IO consistency.

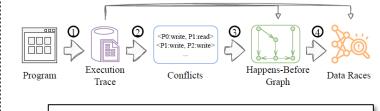
Goals:

Ensure correctness on relaxed consistency systems

- Diagnose portability issues across backends
- Support future file systems beyond POSIX
- Enable developers to detect and fix semantic violations

What is VerifyIO? An open-source tool for trace-based verification of I/O consistency semantics in HPC applications.

Why does it matter? Emerging file systems and libraries relax consistency models (e.g., MPI-IO, Session), which can silently break application correctness.

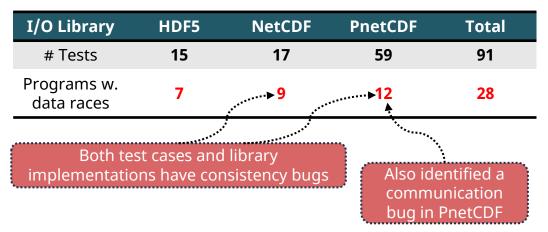


Wang, C., Zhu, Z., Mohror, K., Neuwirth, S., and Snir, M., 2025. VerifyIO: Verifying Adherence to Parallel I/O Consistency Semantics. IPDPS'25.

Meaning: Turning Observability into Insight Verifying Adherence to Consistency Semantics – Results

JG

- 91 built-in test programs from three productionlevel I/O libraries: HDF5, NetCDF, PnetCDF
- 4 consistency models: POSIX, Commit, MPI-IO, Session
- Do I/O libraries comply with the MPI standard?





Redefining HPC Observability: Integrating Monitoring, Modeling, and Meaning • ©Sarah M. Neuwirth • Johannes Gutenberg University Mainz

MPI-IO

Meaning: Turning Observability into Insight FlexBench – Sandbox for What-If Analysis



FlexBench is...

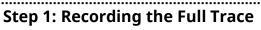
- A benchmark generator that reconstructs and manipulates I/O patterns from execution traces
- Built on top of Recorder+, leveraging Context-Free Grammars (CFGs) to compress and describe I/O behavior

Why do we need FlexBench?

- Traditional tracing tools like Darshan or Recorder capture detailed I/O but lack replay and what-if capabilities
- FlexBench enables users to replay, analyze, and tune the I/O behavior of applications, even without modifying code

Core Goals:

- 1. Use CFGs to precisely describe application I/O patterns
- 2. Reproduce original I/O performance from filtered traces
- 3. Expose optimization opportunities (e.g., parameter tuning)





Step 2: Filtering for a Clean Trace



Step 3: Deriving the Benchmark



Step 4: Performance Optimization







Vision and Closing

Vision and Closing Toward Self-Aware HPC Systems

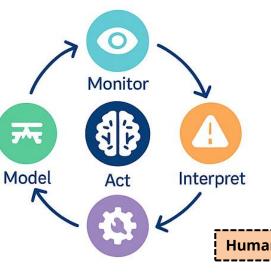


What is a Self-Aware HPC System?

- A system that can monitor itself, understand performance context, and adapt intelligently
- Observability + models + feedback = self-optimization

A Community Vision:

• Requires shared efforts: open tools, standardized metadata, cross-layer collaboration



Core Capabilities:

Why it matters:

- Enables resilience: early detection of failure patterns
- Drives efficiency: energyaware scheduling, performance tuning
- Builds trust: traceable, explainable decisions

Humans on top of the loop

"We do not need smarter tools – we need smarter systems."

Thank you for your Attention!





Dr. Sarah M. Neuwirth

Professor of Computer Science Johannes Gutenberg University Mainz Email: neuwirth@uni-mainz.de Website: https://www.hpca-group.de/ NHR South-West HPC Center: https://nhrsw.de/