



Architecting an AI-Centric Workload Management System: Evolving from Expert-Driven Complexity to an AI-Empowered Approach



PI: Tadashi Maeno

CoPIs: Sara Mason (CFN), Alexei Klimentov (CDF), Paul Nilsson

ECA eligibility: No

Proposal term from: Oct 2026 to: Sep 2028

Annual funding: FY27 \$250k FY28 \$250k



Proposal Description: Main Goals

- Design and study a next-generation Workload Management System (WMS) that transforms an "expert-only" black box to adaptive, learning-based control. The system will integrate AI agents, ML models, and LLMs as first-class components, building on BNL-led foundational efforts such as PanDA, AskPanDA, and predictive ML pipelines.
- Replace static heuristics and manual coordination with AI agents that learn from telemetry, historical trends, and system feedback, and interact with the system through MCP.
- Enable advanced automation for diagnostics, resource prediction, anomaly detection, proactive error mitigation, and real-time coordination of complex, multi-step workflows.
- Lower knowledge barriers through LLM/RAG-based natural-language interfaces that allow operators to query system state and receive guided troubleshooting support.
- Leverage LLMs to automate the database lifecycle, including schema evolution and dynamic query optimization based on historical access patterns.

Proposal Description: Timeline and Success Metrics

Timeline

Year 1: Exploratory Studies & Architectural Design

- Identify high-impact AI/ML/LLM integration points within large-scale WMS.
- Define reference architecture and MCP interaction model.
- Perform feasibility studies on workload diagnostics and resource prediction using historical datasets.

Year 2: Proof-of-Concept Instantiation and Evaluation

- Construct research-grade prototypes instantiating the AI-centric architecture and design guidelines.
- Evaluate end-to-end scenarios using representative workloads from the CFN compute facility.

Success Metrics

- **Operational Efficiency:** Reduction in manual intervention and troubleshooting labor hours, assessed through feedback from operations teams (e.g., ATLAS distributed production and analysis operation team) based on experience in the real production environment.
- **Usability and Accessibility:** Successful use of natural-language interfaces, allowing operators to diagnose issues without direct log inspection.
- **Strategic Alignment:** Delivery of a validated reference architecture aligned with EIC, FCC, and HL-LHC computing requirements.

Intellectual Merit Summary

- Formalizes AI-driven control principles and platform-agnostic design patterns for autonomous large-scale workload orchestration.
- Defines safety boundaries and hybrid architectures for coexisting learning-based and deterministic rule-based schedulers.
- Examines frameworks to synthesize operational telemetry into robust, generalizable system models.
- Establishes methods to ground natural-language interaction in machine-interpretable control, ensuring transparency and traceability in system behavior.

Return on Investment or Potential Future Funding

- Reduces manual troubleshooting and expert intervention through AI-assisted diagnostics and control, lowering operational risk and recurring labor costs.
- Improves utilization of heterogeneous, geographically distributed compute and storage resources via adaptive coordination.
- Captures institutional knowledge within AI agents, ensuring long-term continuity for multi-decade projects like the EIC, FCC, and HL-LHC.
- Aligns with DOE/ASCR AI-focused initiatives (e.g., AmSC) and BNL strategic goals, providing a foundation for follow-on funding within these emerging programs.
- Ensures interoperability with DOE–industry AI partnerships by standardizing on MCP and AI agents, positioning this work for future IRI funding focused on autonomous facility management.

The Broader Impact on the Laboratory

- Accelerates the onboarding of the scientific computing workforce through LLM-assisted diagnostics, lowering the entry barrier for early-career personnel.
- Converges multi-directorate collaboration from NPP, CFN, and CDF, and unites stakeholders from EIC, FCC, sPHENIX, and HL-LHC to build a universal solution.
- Positions BNL as a leader in AI for Science, aligning with DOE's priority for AI-driven infrastructure and autonomous laboratories.
- Directly improves the efficiency of the CFN compute facility through smarter, data-driven resource allocation.
- Reusable AI agent structures, ML pipelines, and LLM-driven tools to enhance BNL-wide R&D.

Names of Suggested BNL Reviewers

- Meifeng Lin (CDF, CSI Char)
- Yihui (Ray) Ren (CDF, AI Codesign Group Lead)
- Michael Begel (NPP, Omega Group Lead)
- Michel Enrique Hernandez Villanueva (NPP, Belle II, HSF)

If ECA Eligible – How this LDRD Benefits Your Application

Tasnuva Chowdhury is an ECA-eligible participant whose expertise in ML and MLOps is critical to the success of this project. Her experience in developing, deploying, and maintaining ML pipelines in production environments directly supports the project's goals of building AI-centric WMSs. Her role includes taking core technical responsibility within the project, providing meaningful leadership experience and supporting the development of early-career researchers in advanced scientific computing.

LDRD Funding Table from BOM

Working with Silvan Minogue based on preliminary labor list

Name	Total Person Months	Role
Tadashi Maeno	2.40	Overall project lead.
Sara Mason	1.20	CFN contributions coordination.
Alexei Klimentov	1.20	SCDF contributions coordination.
Paul Nilsson	2.40	LLM research lead.
Torre Wenaus	1.20	EIC liaison.
Johannes Elmsheuser	1.20	HL-LHC liaison.
Ofer Rind	1.20	Computing infrastructure.
Kolja Kauder	1.20	sPHENIX liaison.
Tasnuva Chowdhury	2.40	ML and MLOps research.
New hire (Postdoc)	24.00	AI Agent-focused research.