Reference Implementation of Smart Scheduler: A CI-Aware, AI-Driven Scheduling Framework for DNN-HPC Workloads

SWATHI VALLABHAJOYULA, SANDEEP SATISH BUDHYA, and RAJIV RAMNATH, The Ohio State University, USA

Many modern scientific workloads in HPC centers rely heavily on AI-driven tasks, particularly deep neural network (DNN) training workloads. Efficiently managing and scheduling these workloads via SLURM interfaces requires users to comprehensively understand available resources, allocation policies, and suitable execution configurations aligned with their models' estimated resource requirements and constraints. Typically, scheduling jobs involves using default configurations, adjusting them as needed, or requesting maximum available limits to ensure uninterrupted execution. However, this approach can lead to job interruptions due to underprovisioning, prolonged wait times, inefficient resource utilization, and increased costs from overprovisioning. These issues ultimately degrade cluster performance, emphasizing the need for a more efficient solution like an AI-enabled Scheduler framework that can profile the DNN workloads and estimate and provision resources dynamically. The existing resource estimation models are trained independently to predict various aspects of batch processing and scheduling, which do not work cohesively to orchestrate a job execution. In our work, we propose to introduce a framework that investigates the feasibility of implementing an iScheduler framework, which transforms the traditional SLURM resource provisioning workflow into an AI-enabled scheduler that plugs different estimators where needed to orchestrate workflow by generating a cyberinfrastructure-aware execution plan, schedules and monitors jobs till completion. We demonstrate the feasibility of our framework by orchestrating a user-specific DNN training workload.

CCS Concepts: • Software and its engineering → Designing software; • Computing methodologies → Cost-sensitive learning.

Additional Key Words and Phrases: AI4CI, AI4OPT, ML, estimation scalability, model, execution time estimation, workflow orchestration

ACM Reference Format:

1 INTRODUCTION

High-performance computing (HPC) environments have become indispensable for handling various scientific workloads, where resources are shared and allocated using batch scheduling based on user requests. Users typically rely on their experience or recommended defaults from Cyberinfrastructure (CI) to estimate resource needs. They aim to secure resources that enable job completion with minimal wait times, interruptions, and, optionally, adjustments to allocation costs.

With the increasing popularity of Deep Neural Network (DNN) models and a huge shift towards Machine Learning (ML) and DNN-driven scientific research is observed and these workflows span from training novel models from scratch to fine-tuning state-of-the-art architectures. This shift has prompted HPC centers to deploy more GPU-powered clusters

Authors' address: Swathi Vallabhajoyula, vallabhajosyula.2@buckeyelink.osu.edu; Sandeep Satish Budhya, budhya.1@buckeyemail.osu.edu; Rajiv Ramnath, ramnath.6@osu.edu, The Ohio State University, Columbus, Ohio, USA, 43201.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM 1

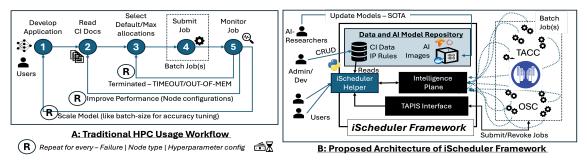


Fig. 1. A,C: The user interaction workflows without and with the iScheduler Framework; B: The Proposed iScheduler Framework with Implementation Software Stack Details

to accommodate the influx of AI workloads. 1. HPC centers have seen a splurge in resource demands as more researchers are initiating AI-driven approaches in their workflows, which are heavily computing-dependent and expensive.

Researchers must comprehend their workloads and fine-tune their allocation requests. Researchers can ensure that everyone receives a fair share of resource allocations by effectively utilizing these resources. A wealth of CI documents, articles, and blogs offer best practices for resource allocation for DNN workloads, considering factors such as batch and epoch scaling, choice of optimizers, and memory requirements. Furthermore, research initiatives have been undertaken to study and understand AI workloads, improve resource allocation, and develop new hardware architectures/chips for compute-intensive AI training, inferencing, and generative AI modeling. Datasets like The MIT Supercloud Dataset ² have been created to facilitate research on workload profiling and estimating resource requirements, particularly for deep neural network (DNN) workloads.

To enhance the efficiency of DNN workloads against the available resources, significant advancements have been made in the development of AI-based resource estimation models and Reinforcement Learning (RL) based schedulers. These models estimate DNN memory requirements and training times, profile user workloads, monitor DNN jobs, and scale frameworks accordingly. They represent a significant step forward in the field, offering researchers new tools to optimize resource utilization. (See section 2 for related literature).

However, accessing these models remains challenging for end-users who request resources and schedule jobs due to the absence of reproducible code or ready-to-use components. To our knowledge, no framework is available to seamlessly integrate these models for automating end-to-end job scheduling (from development to execution). Our work aims to bridge this gap by introducing an iScheduler reference framework. This framework orchestrates job execution by leveraging various AI models designed to allocate resources and monitor jobs via helper functions and API calls (to interact with HPC centers).

The iScheduler Framework (as shown in Figure 1 B) offers the following features to orchestrate the job execution³:

- Provides an end-user helper tool, the iScheduler Helper Python module, to interact with cloud-hosted estimators for fetching resource requirements and orchestrating job submissions.
- Utilizes a database (DB) to store system information and policies, aiding predictions and decision-making before job submission using SLURM commands.4

 $^{^1}$ Clusters like 'cardinal' at OSC are deployed, and chipsets like NVIDIA Blackwell are being released

²https://registry.opendata.aws/dcc/

³The entire iScheduler framework (from Figure 1, B) can be configured on one server (or node), and the users need to deploy only the iScheduler Helper python module on their work environment ⁴We used the CI documentation and SLURM commands to create the DB

- Offers intermediate predictions (memory, walltime, queueing times), execution status, and cost feasibility against
 available systems, allowing manual usage or job delegation (via Intelligence Plane).
- Includes an Intelligence Plane (IP) Service for job submission, monitoring, and tracking on behalf of the user (using TAPIS APIs)
- Supports plug-and-play upgrades of AI estimators as container images, running as TAPIS apps.

We delineate framework components in Section 3 and illustrate the workflow for training vision models in Section 4. Leveraging TAPIS, we execute estimators and run the target training workload as applications. While certain components have broad applicability, we focus on prototyping a workflow for training DNN workloads for three main reasons: a) the increasing adoption of neural networks in workloads, b) the deterministic nature of DNN workloads, facilitating the development of AI-based estimators with wider applicability and reducing the need for individual workload profiling, and c) the simplicity of checkpointing and re-training models, enabling dynamic scaling and reallocation.

The code for all components and an example can be found on GitHub ⁵

2 BACKGROUND

Conventional User Job Execution Lifecycle: A typical user interacts with HPC environments (depicted in Figure 1, A) by developing their applications, consulting CI documentation for suitable resources, submitting batch allocation requests, and monitoring job progress. If a job fails, users diagnose issues and analyze runtime requirements before resubmitting. This iterative process, known as online analytics, aims to address submission errors, enhance performance, and experiment with parameters for scalability and accuracy improvement. However, it is both costly and time-consuming.

2.1 Existing Research

Various AI-driven models have emerged for enhanced scheduling, each addressing different aspects. For instance, DNNMem[2] focuses on estimating memory requirements, while TPUGraphs[3] predict training times per epoch. Frameworks like HARP[7] conduct offline profiling of DNN training loops to understand resource needs. Conversely, models like Scavenger[4] perform online profiling, adjusting allocations based on consumption. Some models predict job start times, while others, like Mirage[1], use reinforcement learning for scheduling GPU jobs. Despite advancements, these models exhibit high estimation errors and require assistance in resource allocation. For example, our walltime estimator for DNNs has a 20% error rate. While models like Scavenger offer better accuracy (4-20% errors), they can only profile jobs while executing.

In our previous research, HARP[5–7], demonstrated offline profiling's viability for user-centric workflows and developed specialized AI-based walltime estimators. We compared default job configurations with our estimators' predictions, assessing feasibility and cost-effectiveness. HARP involved: 1) pre-profiling user workloads against predefined configurations to generate training data; 2) training off-the-shelf estimators, dynamically selecting optimal models based on validation accuracies; and 3) manual configuration of CI policies to estimate predictions' cost and feasibility. However, users needed to run all three steps once to create and store estimators locally, requiring local HARP instances and familiarity with its modules.

⁵ https://github.com/manikyaswathi/iSchedulerFramework

2.2 Opportunities for HPC Job Execution Workflow Enhancement

Current Bottlenecks: 1: Manual fine-tuning is time-consuming, expensive, and resource-intensive, as depicted in Figure 1A.; 2: Independent use of current models still needs more accurate resource requirement estimation, hindering end-to-end job automation.; 3: Utilizing available queues rather than waiting for ideal configurations demands diligent job monitoring and frequent rescheduling. For example, development queues offer similar computing capabilities to longer allocation queues but with limited allocation time (2 hours), requiring users to fragment their execution and submit multiple jobs ⁶.

The proposed framework (Figure 1B) seamlessly integrates plug-and-play resource estimators by leveraging CI awareness through a centralized database for end-to-end job execution. Key components of our architecture are (a) CI database: Stores all CI policies and aids in Intelligence Plane decisions. (b) Model repository: Hosts AI estimator models. (c) TAPIS Interface: Facilitates job scheduling across diverse CIs and clusters via APIs. (d) Intelligence Plane (IP): Manages job scheduling, monitors execution progress, and adjusts schedules as needed. Additionally, our iScheduler Helper module offers various functions for users to invoke suitable AI models for estimations based on their configurations. The subsequent sections discuss further details on these components and an example workflow.

3 ISCHEDULER: COMPONENTS AND IMPLEMENTATION DETAILS

This section elaborates on the iScheduler components, including their objectives, implementation details, and rationale behind the selected software stack. The current prototype is designed to facilitate DNN vision model training, and the software stack choice aims to complement other ICICLE⁷ modules, such as cloud-to-edge orchestration.

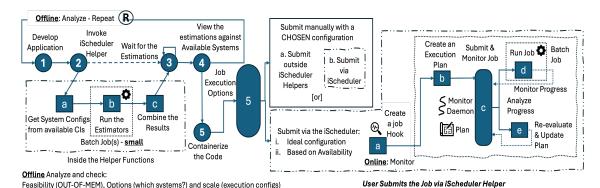


Fig. 2. The user interaction with the iScheduler

4 ISCHEDULER WORKFLOW WITH A DNN TRAINING JOB EXECUTION SCENARIO

Figure 2 illustrates the user interaction flow with the iScheduler Helper module to ascertain their job's resource requirements against available systems. It shows the orchestration of job execution with the Intelligence Plane.

Use Case: A user trains a ResNet50 vision model with 20k images and 50 epochs to improve accuracy by testing various batch sizes. They utilize the iScheduler helper to manage resource requirements efficiently. Workloads are

Manuscript submitted to ACM

 $^{^6}$ Refer to the cluster batch limitations outlined in the respective CI documentation.

⁷https://icicle.osu.edu/

Table 1. Components of iScheduler Framework with software requirements, implementation details, and functionalities

iScheduler Helper tool	Database (S)
C/arr. Darth are seen dulls	S/w: MongoDB
S/w: Python module	-Chosen for its lightweight and NoSQL nature
Easy to use with TensorFlow and PyTorch.	-Ideal for prototype development, flexibility in schema
Actors:	changes
-Developers: Create functions	Actors:
to add workflows.	-Developers/Admins manage the database,
-Users: Utilize functions to invoke i	adding data automatically via a Python script
nference models.	-Users interact with the database using the iScheduler
Functionality:	Functionality:
It offers an interface for end users to access the central	The DB stores CI details to estimate resource needs against
server functionalities via API calls.	node configurations, validating allocation feasibility
Users configure their access tokens (TAPIS/DB).	and computing costs against batch limitations.
Intelligence Plane (S)	AI Models (S)
S/W: Kafka, Spring, Flask	
-Kafka is used to capture the profiling and progress	S/w: Images (Docker/Apptainer)
status of the executing jobs for online analytics.	-Facilitates on-the-fly inferencing deployment
-Flask: The consumers of the monitoring data that	-Seamless integration with TAPIS for support
invoke IP rules and make appropriate decisions.	-Supports plug-and-play for developing new models
Actors:	Actors:
Developers - add more IP features	-Developers/Admins create iScheduler helper functions
Users interact with IP using the iSchdulerHelper.	to add new workflow estimation models
Functionality:	-Users invoke these models via the helper functions
-Submits a TAPIS job on behalf of the user by creating	Functionality:
a job hook to monitor job progress.	-Replaces human estimations and incorporates
-Maintains independent daemons for tracking job	existing work (online and offline
progress and making ad-hoc decisions for rescheduling.	analytic models) that could be invoked based on need
-The current prototype supports rescheduling jobs	in the framework workflow at appropriate times
based on progress and remaining allocation time	
The Framework user TAPIS to invoke AI models and automate the Job execution. Tapis framework serves as a cloud-hosted	

API that enables the configuration of various systems on multiple CIs, application registration, and job scheduling.

(S) - A server is needed to deploy the framework as a centralized service

single-node allocations managed through TAPIS, with checkpointing and callbacks ensuring job continuity and progress monitoring.

Workflow: The workflow between the user and the Smart Scheduler is as follows:

- After configuring the DNN architecture, the user invokes the appropriate iScheduler Helper function, providing model details and desired batch sizes (32, 64).
- The Helper Function retrieves node configurations (CPU and GPU note configs for Pitzer Cluster) from the CI database and runs estimators to provide resource predictions.
- The user reviews the estimations and selects an execution plan, either manually submitting the job or delegating it to the Intelligence Plane. Result: A Pitzer CPU node can execute ResNet50 for all batch sizes with maximum execution times of 9 and 9 hours and corresponding costs of \$0.756 and \$0.756 for batch sizes 32 and 64, respectively. Additionally, Pitzer GPU (with a maximum of 32 GB memory) can handle only batch size 32 with an estimated time of 2 hours and costs of \$0.348.

- The user has two options after reviewing the estimation results: a) Manual Submission, where they configure a SLURM script to submit the job themselves, and b) Delegating the job execution and monitoring To the Intelligence Plane. The prototype demonstrates two feasibilities submitting with the best configurations (ideal system configurations for minimal interruptions) and based on system availability (CI queue state: number of idle nodes, number of waiting jobs).
- Containerize the code, publish it to a publicly accessible location, and create a TAPIS Application.
- The Intelligence Plane manages job execution, submitting, and monitoring progress through callbacks. Jobs may run as single tasks or cascading ones if exceeding maximum walltime. It reassesses plans after each sub-job, adjusting or rescheduling as needed, triggering AI-finetuning loops for memory and walltime estimations. For instance, submitting a job (batch size 32) with an ideal configuration (Pitzer GPU) has an estimated job completion of 28 hours (26 hours wait time + 2 hours execution) and \$0.348 cost. Alternatively, executing based on availability yields an estimated job completion of 9 hours (no waiting time as resources are allocated immediately) with \$0.756 cost.

Offloading the job execution-monitoring loop to the IP enables the utilization of various estimation models, including offline profilers, estimated allocations based on cluster queue information, and online profiles with predefined rules for job execution. Opting for availability allows users to execute their jobs partially or entirely on less preferred nodes until a desired allocation becomes available, reducing overall execution time, including wait times. Development queues, with identical node configurations as full queues, are typically idle at night, making them ideal for utilization with the scheduler framework without requiring manual monitoring of short cascading jobs. Employing an intelligent scheduler can enhance the productivity of such compute-intensive workflows and diminish the human-in-the-loop effort.

Our future work will focus on refining the current prototype to support DNN workflow orchestration, collaborating closely with domain experts. This involves validating and designing each component, such as ensuring the CI database mirrors the SLURM DB, requiring collaboration with CI Admins for effective development.

REFERENCES

- Qiyang Ding, Pengfei Zheng, Shreyas Kudari, Shivaram Venkataraman, and Zhao Zhang. 2023. Mirage: Towards Low-interruption Services on Batch GPU Clusters with Reinforcement Learning. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. 1–13.
- [2] Yanjie Gao, Yu Liu, Hongyu Zhang, Zhengxian Li, Yonghao Zhu, Haoxiang Lin, and Mao Yang. 2020. Estimating gpu memory consumption of deep learning models. In Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 1342–1352.
- [3] Mangpo Phothilimthana, Sami Abu-El-Haija, Kaidi Cao, Bahare Fatemi, Michael Burrows, Charith Mendis, and Bryan Perozzi. 2024. TpuGraphs: A Performance Prediction Dataset on Large Tensor Computational Graphs. Advances in Neural Information Processing Systems 36 (2024).
- [4] Sahil Tyagi and Prateek Sharma. 2023. Scavenger: A Cloud Service for Optimizing Cost and Performance of ML Training. In 2023 IEEE/ACM 23rd International Symposium on Cluster, Cloud and Internet Computing (CCGrid). IEEE, 403–413.
- [5] Manikya Swathi Vallabhajosyula and Rajiv Ramnath. 2022. Towards Practical, Generalizable Machine-Learning Training Pipelines to Build Regression Models for Predicting Application Resource Needs on HPC Systems. In Practice and Experience in Advanced Research Computing (Boston, MA, USA) (PEARC '22). Association for Computing Machinery, New York, NY, USA, Article 43, 5 pages. https://doi.org/10.1145/3491418.3535172
- [6] Manikya Swathi Vallabhajosyula and Rajiv Ramnath. 2023. Insights from the HARP Framework: Using an AI-Driven Approach for Efficient Resource Allocation in HPC Scientific Workflows. In Practice and Experience in Advanced Research Computing. 341–344.
- [7] Swathi Vallabhajosyula and Rajiv Ramnath. 2022. Establishing a Generalizable Framework for Generating Cost-Aware Training Data and Building Unique Context-Aware Walltime Prediction Regression Models. In 2022 IEEE Intl Conf on Parallel Distributed Processing with Applications, Big Data Cloud Computing, Sustainable Computing Communications, Social Computing Networking (ISPA/BDCloud/SocialCom/SustainCom). 497–506. https://doi.org/10.1109/ISPA-BDCloud-SocialCom-SustainCom57177.2022.00070