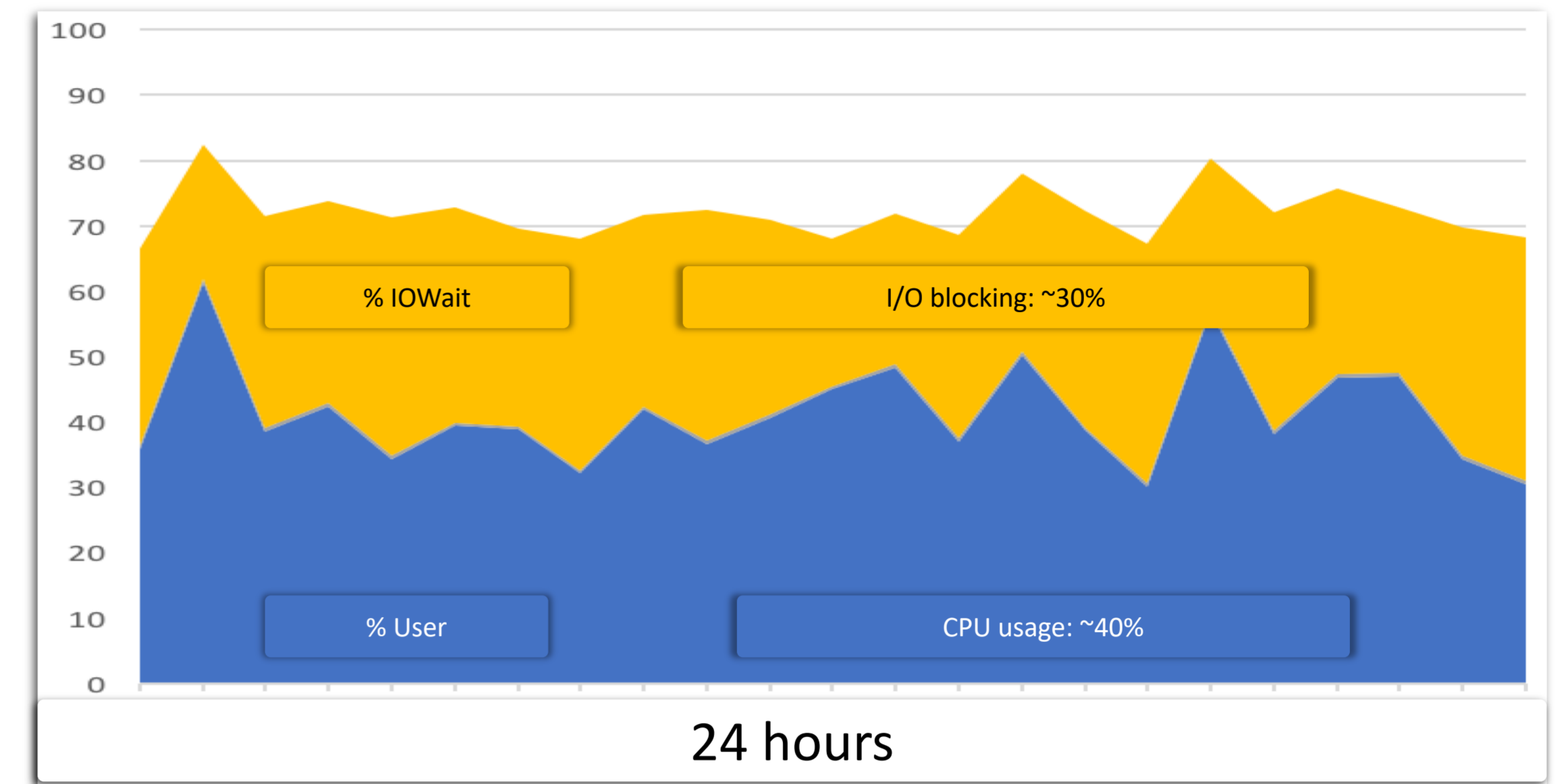# Application of Deep Learning on Integrating Prediction, Provenance, and Optimization

**Malachi Schram[1], Nathan Tallent[1], Ryan Friese[1], and Alok Singh[2]**
[1]Pacific Northwest National Laboratory
[2]University of California, San Diego

**Pacific Northwest**
**NATIONAL LABORATORY**

**Motivation**: Current computing grid scheduling tools do not protect computing resources from "bad" behaviors. There is limited mechanism to handle input/output (I/O), memory, and networking contention. In recent grid production campaigns, CPU usage had up to 60% wasted with 30% being blocked by I/O (figure at right). The pilot job mechanism is ideal for dynamically determining the status of the worker node and forecast contention before any new jobs are submitted.
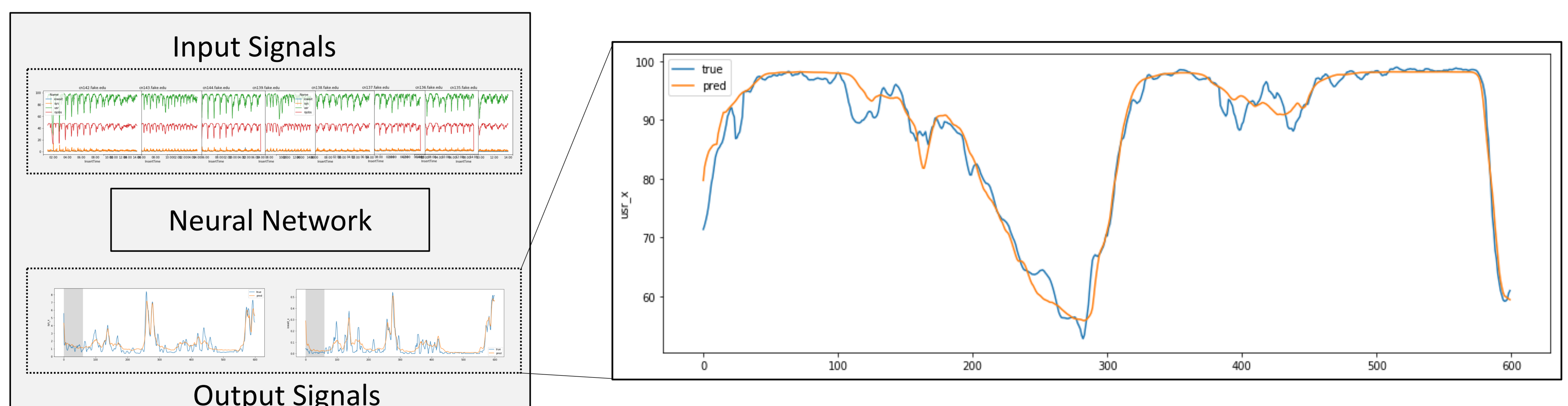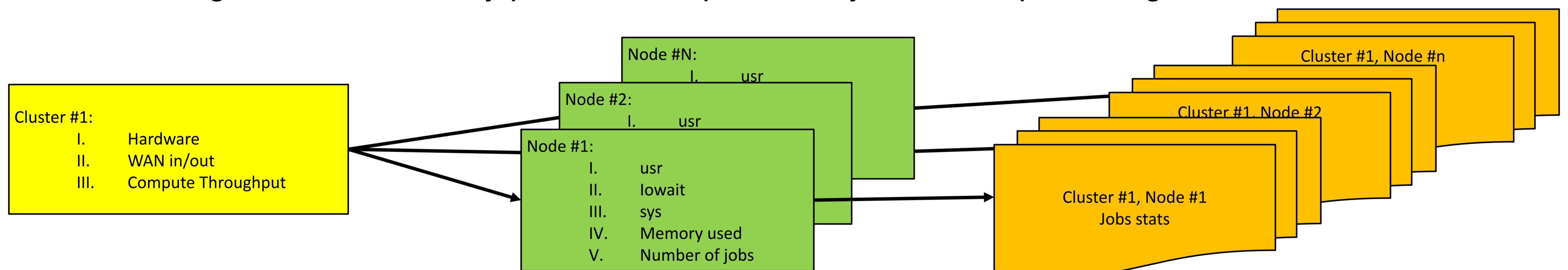


**Goal**: We investigate using machine learning (ML) within the grid pilot job mechanism to improve the efficient use of computing resources. In this research, we investigate two methods to detect job anomalies and/or contention:

1) Preemptive job scheduling using binomial classification long short-term memory (LSTM) networks
2) Forecasting intra-node computing loads from the active jobs and additional job(s)

**Approach 1**: We use LSTM to predict failures based on two years of labeled data, which correspond to 10 million jobs and 910 million samples. The model is composed of two LSTM layers with a sigmoid activation layer. The model has been trained with two labels, "Done" and "Failed", resulting in a 88% precision. Using this model to preemptively predict failures can potentially introduce a 14% speedup.

**Approach 2**: We capture higher-fidelity time series information at the job level to build a multilevel machine learning model, which only presents the preliminary results of predicting the node-level loads.





**Conclusions and future work**: We investigated two methods to detect job anomalies and/or contention.
- The initial results from Approach 1 suggests a 14% speedup by preemptively killing "bad" jobs.
- Preliminary results for Approach 2 suggests the ability to accurately predict future load.
- As part of future work under Approach 2, we intend to investigate:
  - Multi-job scenarios, error analysis and anomaly threshold, forecasting window size, and multilevel machine learning model.