# MLTF.VU

## **M**achine **L**earning **T**raining **F**acility at **V**anderbilt **U**niversity

Jethro Gaglione, Andrew Melo, Satkar Dhakal, Lindsey Fox,
Matt Heller, Prasanna Koirala, Samuel Ninchritz,
Brandon Soubasis, Alan Tackett

# Motivation and Facility Overview

- Machine learning/AI has been widely adapted for an array of uses in HEP and beyond.

- However, there are considerations about <u>training</u> these models -

    - *<u>Reproducibility</u> - Can someone later consistently produce the same weights?*

    - *<u>Scalability</u> - How does one scale the training from card->node->cluster scale?*

    - *<u>Efficiency</u> - What are the barriers to using high-speed transports and interconnects?*

- MLTF is a prototype facility that enables training w/these considerations in mind -

    - *<u>Hardware</u> - At this stage, capabilities were chosen over large scale*

    - *<u>Software</u> - Easy-to-use and documented hooks/wrappers for clients to enable functionality*

    - *<u>Infrastructure</u> - Services run at the site to connect the client software to hardware resources*

# Motivation and Facility Overview



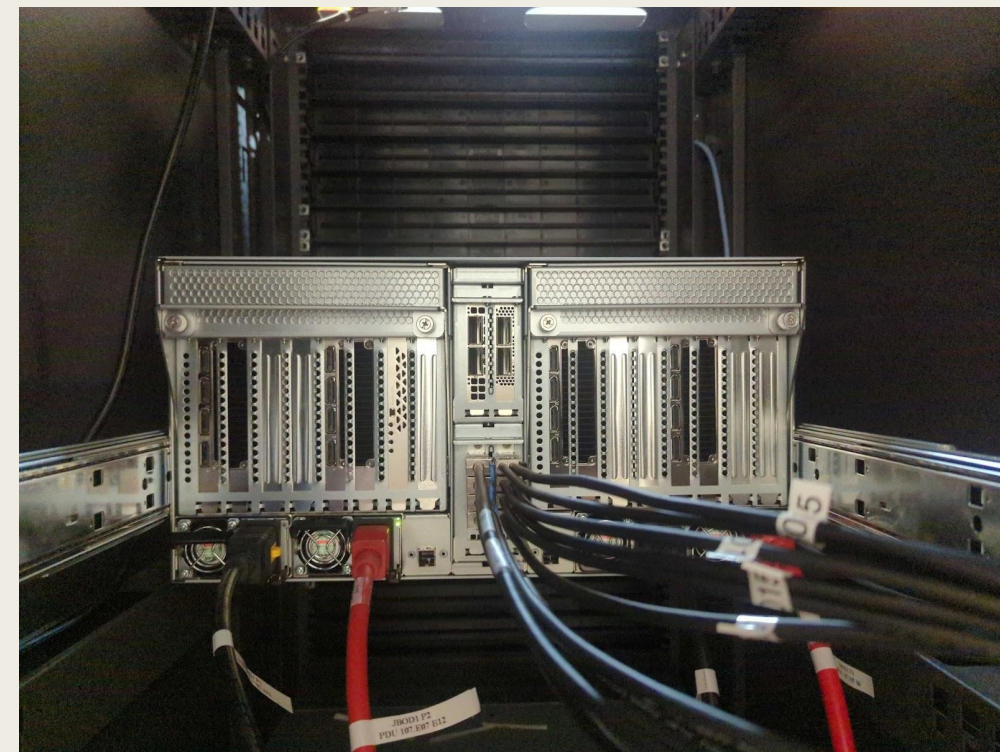https://www.evergreen-line.com/container/images/CNTR_Overive.jpg

- MLTF seeks to provide a higher-level abstraction to the user
  - *Users submit training tasks, not batch jobs*
  - *Tasks describe environment, code, inputs, HW reqs*
  - *Facility takes this description, executes jobs and returns outputs*
- *Roughly* analogous to containerization
  - *Container engine needs to know little about host environment or even execution engine and vice versa*
  - *Similarly, user doesn't have to opt-in*
    - docker run --net=host
- This abstraction benefits both the user and the facility
  - *Facility can use higher-level knowledge to optimize/configure as-needed (e.g. pre-staging input data, configuring NCCL, scheduling decisions)*
  - *User doesn't need to know all implementation details*
- Helps reproducibility - Outputs tied back to task that produced them

# Hardware

- The proof-of-concept facility will have 16xA6000s (w/NVLink), RoCE networking, and 16 large NVME disks:
  - RoCE allows RDMA access over an ethernet network for low-latency, high-performance transmission.
  - NVMe allows for direct GPU access to data with high throughput and fast response times.
- The components are housed in PCIe enclosures, which connect to compute nodes via PCIe expander cards
- We would like to acquire PCIe switches, which will allow us to scale out the PCIe fabric to an arbitrary number of enclosures and compute hardware.
- Delays in receiving hardware means a lot yet to explore in terms of functionality 🙁
- There is a possibility for "opportunistic" access to other hardware @ Vanderbilt, including O(100) GPU, DGX hardware, Grace Hopper GH200 machines
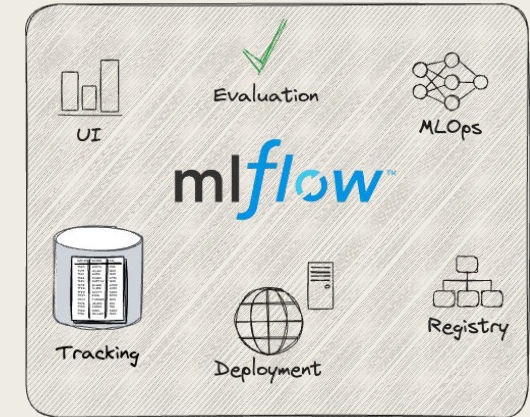


PCIe enclosure. Thick black cables in the middle route to compute hardware.

# Training Framework Considerations

- To begin this work, we need to select a training framework.
  - *Important: this is not a decision about ML library (e.g. Tensorflow/Keras/etc..)*

- With the ongoing gold rush, there's more than a few frameworks available
  - Nearly all of them are "semi-open source"
    - Lock-in / price become a risk if you want to use "interesting" features
  - Many are additionally exclusively Kubernetes-based, meaning the GPUs are provisioned via Kubernetes.
    - This is traditionally not how we deploy resources in HPC centers and would imply segregating (expensive) hardware in many places.
    - K8S scheduling not as robust as Condor/SLURM – the assumption is that if you have more jobs than resources, you just rent more cloud resources. Certainly not usually the case for academia!
- After evaluating several options, MLFlow was chosen as the preferred training framework

# MLflow

■ **MLflow is a solution for the entire lifecycle of ML:**
   – A job submission interface
   – Tracking server which can store:
     ■ Inputs - Python code, environment, input files/datasets
     ■ Outputs - Trained weights
     ■ Diagnostics - CPU/GPU utilization, loss metrics, time per epoch, etc..
   – Inference server - effectively a REST wrapper around a model

■ **Ticks a lot of necessary boxes**
   – *Open-source, with public Helm charts for deployment*
   – *Right mental model: An MLflow "project" is a training "task"*
   – *Minimal changes to existing code required to integrate*
     ■ File describing Project metadata, two lines of python in user code
   – *Robust REST API*
   – *Defined & Documented extension API*

■ **Not the first in academia to come to this same conclusion. NCSA also uses MLflow as their preferred training framework(<u>NCSA docs</u>)**
   – *SLURM & Token authentication plugins developed by NCSA*



<u>https://mlflow.org/docs/latest/index.html</u>

# MLflow Projects Example

- Minimal MLflow Project is simple
  - Directory containing four files: "*MLproject*" entry file, "*my_env.yaml*" environment req. specs, "slurm_config.json" for job req., and your training script.
  - Can be local or in VCS (e.g. Git)
- Multi-step workflows are supported natively - e.g. feature extraction followed by training
- Have adapted an in-production CMS ML production workflow to use MLflow
- Future work:
  - *Support rucio:// URIs for input files*
  - *Support parallelism via Dask/Luigi*

```
name: Project Test

python_env: python_env.yaml
entry_points:
  main:
    command: "python train_sklearn.py"
```

```
# Python version required to run the project.
python: "3.10.4"
# Dependencies required to build packages.
build_dependencies:
  - pip
  - setuptools
  - wheel==0.37.1
# Dependencies required to run the project.
dependencies:
  - mlflow==2.8.1
  - cloudpickle==2.2.1
  - numpy==1.26.0
  - packaging==23.2
  - psutil==5.9.6
  - pyyaml==6.0.1
  - scikit-learn==1.3.1
  - scipy==1.11.3
```

```
{
"account": "mltf_acc",
"mem": "4G",
"partition": "pascal",
"time": "8:00:00",
"gres": "gpu:1",
"mail-user": "jethro.t.gaglione@vanderbilt.edu",
"mail-type": "FAIL",
"ntasks": "1",
"job-name": "MG_process_HiggsW",
"output": "G_NJL_tautau_2J.out",
"sbatch-script-file" : "batchFile",
"modules": ["GCCcore/.11.3.0", "Python/3.10.4"]
}
```

# MLflow Gateway

- Remote submission of MLflow projects opens interesting use-cases
  - *No need to get local accounts @ target cluster*
  - *CI/CD applications*
- MLflow has built-in support for remote submission via DataBricks
  - *Extension API allows this to be pluggable*
  - *Support for SLURM submission exists due to MLflow-SLURM*
- We are developing MLflow Gateway, a Client/Server interface to job submission
  - *Client side is an MLflow backend plugin*
  - *Server side is a REST service which interacts with local cluster*
- Since the gateway is an "intermediary" between the client and the facility, this provides a hook to add customization to the jobs before they are submitted to the batch system
  - *Difficult to do in direct SLURM submission*
  - *A task requesting a Rucio dataset as an input can trigger the Rucio rule, and wait for the replication to complete before submitting jobs to batch system*

# MIflow Gateway



- Most components are existing, available, open-source projects.

- Authentication is token-based, using familiar SSO pages.

- All user code is run within a container to ensure isolation.

  - Either user-provided container or user python env within a container

# Authentication Flow

# MLflow Deployement

- Central components (MLFlow tracking, databases, artifact storage) are deployed in k8s, while GPU remains in SLURM cluster
- Currently we have SLURM submission functional, which requires access to Vanderbilt's cluster
- Goal is to have the whole stack (MLflow deployment + Gateway) packaged as Helm chart so it can be deployed elsewhere w/a k8s cluster and SLURM

# Training

# Submitting Jobs and Viewing Results

- After setting the tracking server,

```
export MLFLOW_TRACKING_URI=https://mlflow-test.mltf.k8s.accre.vanderbilt.edu
```

activating the required token,

```
export $(mlflow-token)
```

and running:

```
mlflow run --backend slurm --backend-config <path to slurm_config.json> <path to MLflow Project dir.>
```

- The resulting model weights/parameters, tracked model metrics, system metrics, and a range of other selected artifacts can be viewed/accessed in the MLflow UI.

- Code to activate/customize tracking and logging functionalities is inserted in the training script. Popular ML packages (Scikit-learn, Keras, Pytorch Lightning) have autolog support, used by simple adding the following to the Python training script:

```python
import mlflow
mlflow.autolog()
```

# MIflow User Interface



- Detailed description of training conditions, e.g. date/time, duration, source code filename, user…
- Option to save datasets
- Automatic or custom logging of parameters and metrics of interest.

# MIflow User Interface

- Complete set of information on the trained model.

- Environment information along with "requirements.txt" file.

# MlflowUser Interface



- Easy access and training execution of logged models for efficient and accurate reproducibility.

# Preliminary Testing and Performance

We've tested our software infrastructure for flexibility across various training scenarios with :

- Data-distributed training (PyTorch DDP)

- Model-distributed training.

- Hyperparameter Optimization (Optuna)

- In conjunction with established full-suite training frameworks
    - *We've wrapped MLTF functionality around the CMS BTV POG's b-hive modular object-tagging framework.*
    - *Aside from the benefit of being able to train at MLTF, this also brings the enhancement of the MLflow suite to log system metrics, track models and parameters, and stage the trained model for inference/reproduction.*

# Documentation

- Our goal is to simplify the user experience – a system that is hard to use will often remain unused
- We have begun to develop user-facing documentation for the facility - https://docs.mltf.vu/
  - We currently cover PyTorch-based training.
  - Plan to eventually explicitly document all popular frameworks.
- Still a lot of work to be done on this front

# Additional Plans and Future Work

(in no particular order)

- We plan to test out interesting features like GPUDirect Storage (NVME<->GPU communication) and large-scale distributed training (10s of GPUs spread across multiple nodes).

- Development of MLflow Gateway to facilitate MLTF usage for external users.

- Add Luigi/Dask support for parallelization

  – *e.g. the user can request 2 Dask jobs, and MLflow will handle requesting them from SLURM properly*

- Deploy and demonstrate this deployment at another site
- Package and release Gateway client/server
- Develop a class/tutorial to train new users on how to use this facility