# Find Inefficiencies and Rapid Model Profiling with CentML DeepView

**ASPLOS '23 Workshop, March 25th**

Yubo Gao <ybgao@centml.ai>
CentML Inc, University of Toronto

# Agenda

1.  Introduction to training optimizations    **[1:50pm - 2:30pm]**

    a.  **Why** do we care?
    b.  **What** are the common optimizations?

2.  Performance debugging with DeepView    **[2:30pm - 3pm]**

    a.  Visually identify performance bottlenecks
    b.  Value of performance prediction in optimization workflow

# Why optimize?

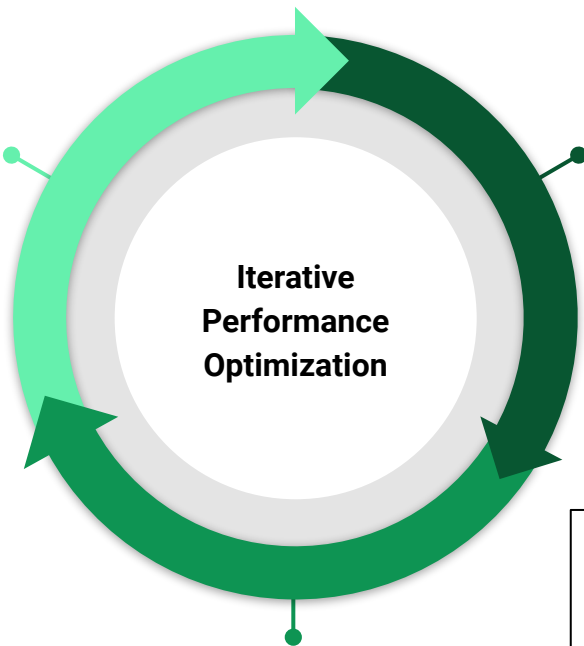Increasing resources required used to train large models.

**Resource underutilization** is a significant problem:

- Observed average GPU utilization below 30% at a large AI research compute cluster.
- Significant resource and energy waste.
- Low utilization leads to lower throughput (increased job completion time).

# Typical training optimization workflow

How much improvement?

Energy improvements?

**Quantify Improvements**

**Iterative Performance Optimization**

**Identify bottlenecks**

Compute bound or memory bound?

Where (which layer, or lines of code) did the time go?

Understand energy and environmental impacts

**Apply training optimizations**

How much improvement to expect?

Pick the most suitable hardware.

# Environment Setup

**DeepView**

Please follow "Environment Setup" at:

https://centml.github.io/asplos23-tutorial/deepview.html

# Interactive Demo

Exploring system optimizations for DL training

# Existing DL Profilers



nvprof
Nsight compute
Nsight Systems
dlprof

Intel vTune

Torch.profile

PyTorch Lightning
profiler

TensorBoard
profiler

1. Incorrect granularity
2. Lack of interactivity
3. Lack of predictive capabilities

# Incorrect granularity



**Forward Pass:** Difficult to attribute runtime to lines of code.

**Backward Pass:** Only layer names are visible, impossible to trace back to code.

**Incorrect granularity makes optimization difficult.**

# Incorrect granularity



NVIDIA DLProf [1]



Tensorflow Profiler [2]

[1]: https://docs.nvidia.com/deeplearning/frameworks/dlprof-user-guide/index.html
[2]: https://www.tensorflow.org/guide/profiler

# Lack of interactivity



NVIDIA Nsight Systems



PyTorch Profiler

# Deep**View**

## Interactive Profiler [3]

**Identifies** performance bottlenecks

**Enables** rapid iterative profiling

**Quantifies** energy consumption and environmental impacts of training jobs.

## Runtime Predictor [4]

**Predicts** a deep neural network's training iteration execution time on a different GPU.

**Recommends** the most cost/time effective hardware option for your workload

[3]: Skyline: Interactive In-Editor Computational Performance Profiling for Deep Neural Network Training, Geoffrey Yu, et. al.
[4]: Habitat: A Runtime-Based Computational Performance Predictor for Deep Neural Network Training, Geoffrey Yu, et. al.

# **Interactive Demo**

Iterative Profiling with CentML DeepView

# GPU Runtime Predictor (DeepView.Predict)



Pick the best GPU for your training job, whether it is:

- Deciding which new GPU to purchase for your local workstation

- Which cloud GPU instance to pick

- Efficiently schedule jobs in a heterogeneous GPU cluster

Tedious to benchmark model on all the available GPUs.

# GPU Runtime Predictor (DeepView.Predict)



**Wave Scaling** - If the same kernels execute on the source and target GPUs, then scale based on hardware parameters.

**Pre-Trained MLPs** - If not, learn runtimes of different operators with a pretrained model.

# GPU Runtime Predictor (DeepView.Predict)

**Why predict?** Why not ....

➔   Measure the performance directly?

➔   Apply heuristics?

➔   Use standard benchmarks?

➔   Always use the best available GPU?

# Aren't more powerful GPUs better?

Not all the time! Not alwa~~ys~~ [Small **inputs** underutilize GPU.] ~~th~~e A40.

# Aren't more powerful GPUs better?



Small **models** under-utilize GPU.

# Which GPUs are supported?

| Generation \ Use Case | Desktop/Consumer | Workstation/Server |
|---|---|---|
| Pascal | GTX1080Ti | Quadro P4000<br>P4<br>P100 |
| Volta | | V100 |
| Turing | RTX2070<br>RTX2080Ti | Quadro RTX4000 |
| Tesla | | T4 |
| Ampere | RTX3090 | A100<br>A40<br>A4000 |
| Hopper | | H100 (coming soon) |

- Beta

- Deprecated / Not available

# How accurate is DeepView.Predict?



Prediction errors are generally no more than 10%.

# Cloud Deployment

Deployment

## Deployment Target

Estimation for `4 million` total iterations

**Providers**

Filter by provider

All

Filter by GPU

All

Filter Max Number of GPUs:

1 2 4 **all**



● google
● azure
● aws

**Deployment Plan**

# n1-standard-1

**Estimated Cost: $20**

**Estimated Training Time: 3.4 Hours**

| GPU | Num. GPU | VRAM |
|---|---|---|
| **p100** | **4** | **16 GB** |

20

# CentML DeepView

CentML DeepView provides an integrated experience which allows ML practioners to:

- Visually identify model bottlenecks
- Perform rapid iterative profiling
- Understand energy consumption and environmental impacts of training jobs
- Predict deployment time and cost to cloud hardware



## Getting Started

Follow the instructions depending on your setup

- Local workstation GPU
- Remote GPU with SSH access
- Clusters, containers, and other setups where SSH is not possible

### DeepView.Profile

Deepview.Profile is our free and open source tool provides easier way to identify bottlenecks and perform rapid iterative profiling for Deep Learning. To find out more, visit DeepView.Profile

### DeepView.Predict

DeepView.Predict is a tool that predicts a deep neural network's training iteration execution time on a given GPU. It currently supports PyTorch. To find out more, visit DeepView.Predict

**Get started with DeepView at** `docs.centml.ai`**!**

# What we do: system-level optimizations



**CentML expertise**

Profiling Tool:
DeepView →

System-Level
Optimizations →

Hardware-Specific
Optimizations →

Models

| ResNet | GPT3 | BERT | SSD |

Algorithms

SGD without momentum    SGD with momentum

Frameworks

TensorFlow    PYTORCH    mxnet

Libraries and ML Compilers

cuDNN    TVM    XLA

Hardware

GPU    A100    TPU V3

aws
Microsoft Azure
Google Cloud

# Additional Optimizations: Horizontal Fusion

Horizontally fused training array for efficient training

- Best for training small models + hyper-parameter tuning

# Thank you!

To learn more about DeepView, visit
`docs.centml.ai`

Also check us out at **centml.ai**,

or contact us at **ybgao@centml.ai**