

"Advancing Fault Tolerance, Resource Efficiency, and Cost Optimization for AI Workloads in Cloud Computing"

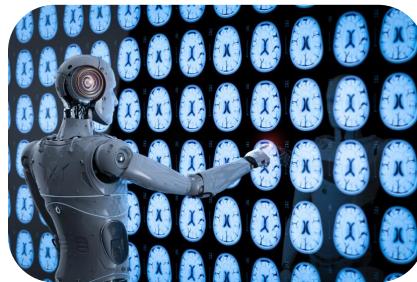
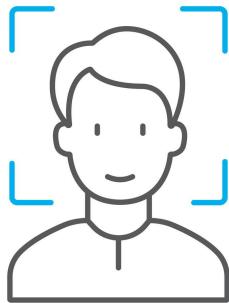
Maziyar Nazari
Jan 2025
Ph.D. Thesis Defense



University of Colorado
Boulder

Introduction

Deep Learning Applications are Everywhere!



Introduction

2 Parties in Machine Learning Application Development Lifecycle

Developer

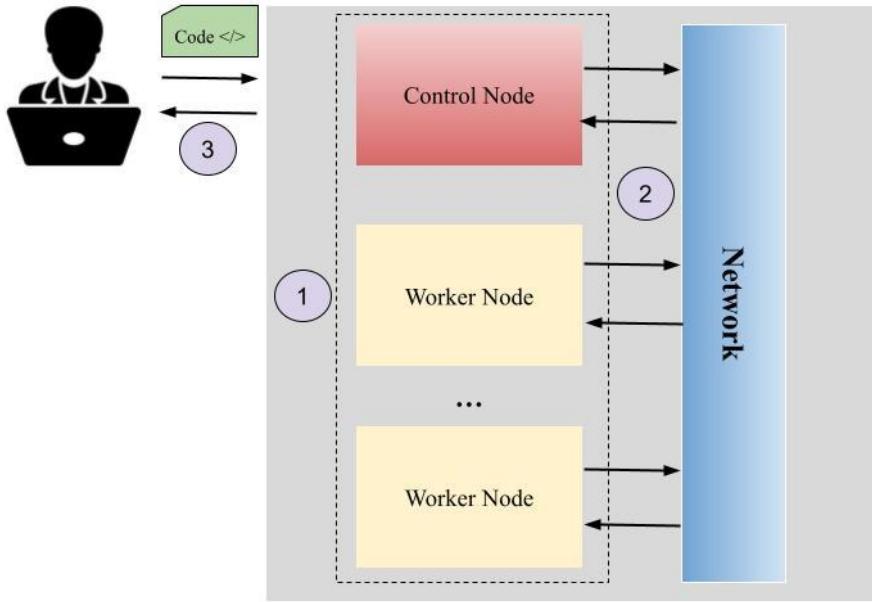
- **Data Collection**
- **Data Preprocessing**
- **Model Development and Deployment**
- **Pick the right cloud configuration**
- **Monitoring**
- **Hyperparameter Tuning**

Cloud Provider

- **Resource Management**
- **Service Level Agreement**
- **Cost Effective**
- **User-friendly Interface**
- **Rich API**



Overview



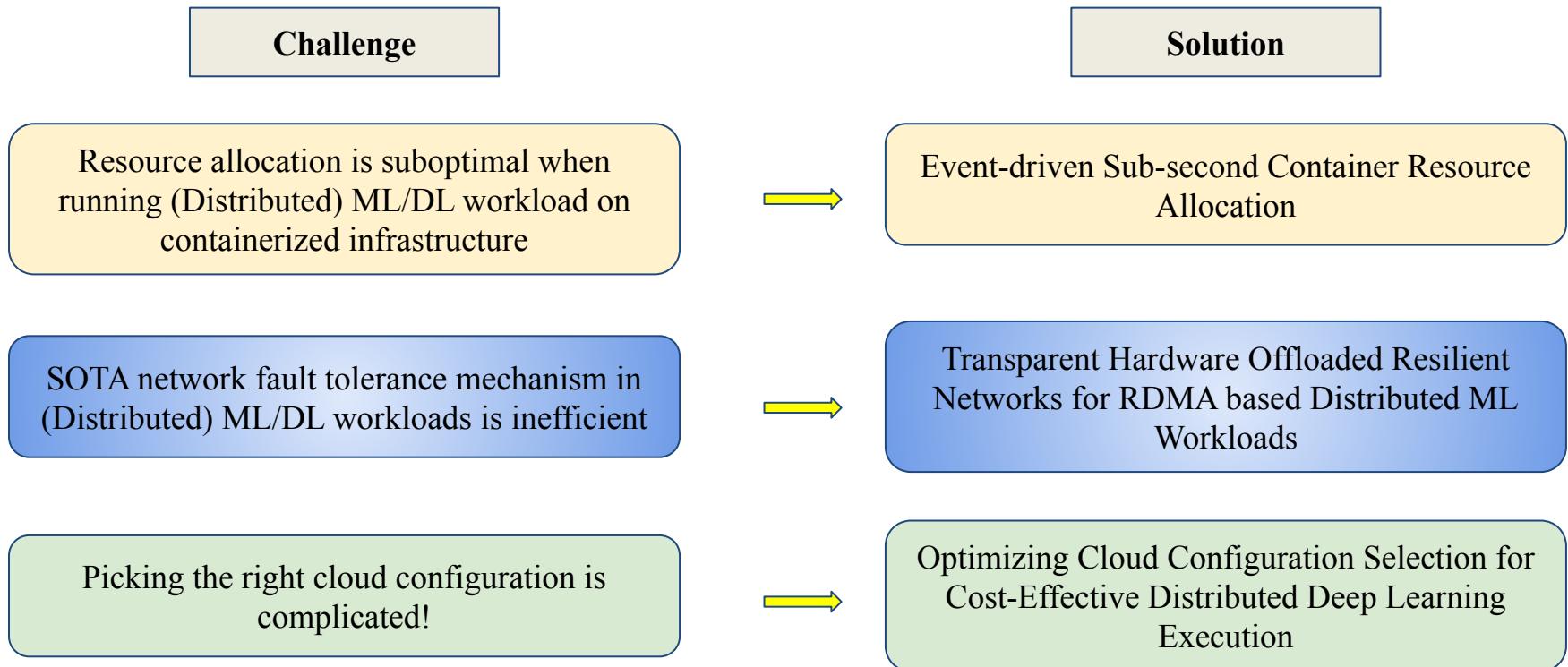
Resource allocation is suboptimal when running (Distributed) ML/DL workload on containerized infrastructure

SOTA network fault tolerance mechanism in (Distributed) ML/DL workloads is inefficient

Picking the right cloud configuration is complicated!



Introduction



Escra: Event-driven Sub-second Container Resource Allocation

[Published in ICDCS 2022]

Escra: Event-driven, Sub-second Container Resource Allocation
Greg Casack, Maziyar Nazari, Sepideh Goodzary, Erika Hunhoff,
Preet Oberai, Eric Keller, Eric Ruzer, Richard Han
University of Colorado Boulder, Boulder, CO USA

Abstract—This paper pushes the limits of automated resource allocation in container environments. Recent work on container systems has mostly limits by averaging scaling behaviors based on past resource usage. However, these systems are slow and run on coarse-grained time scales. In contrast, Escra, a container system which predicts and adapts its behavior to a container's workload, achieves near-optimal performance. Escra's resource allocation for a single container is event-based with fine-grained intervals. Extra intervals allow Escra to manage a large number of containers simultaneously, and it can reduce interference between containers. Escra manages two performance metrics: CPU waste and resource utilization. In microservices environments, our evaluation shows fine-grained and event-based resource allocation can reduce execution latency by up to 9.8% while simultaneously reducing CPU waste by 99.8%. In serverless environments, Escra can increase吞吐量 by up to 10x and reduce CPU waste by over 10x and 100x, respectively. In various environments, Escra reduces CPU reservations by over 2.1x and reduces CPU reservations by over 1.5x while maintaining similar end-to-end performance.

INTRODUCTION
Containerized infrastructure is quickly becoming a preferred method of deploying applications. The light-weight nature of containers coupled with their orchestration tools enable a new way to design, manage, and operate systems that integrate with development workflows. These deployments have container resource limits that are used to limit interference between containers and underlying resource usage.

Setting container resource limits is a trade-off between application performance and efficiency. When resource limits are set low to prioritize efficient resource usage, applications will experience an increased number of CPU threads and out-of-memory (OOM) errors. This results in degraded application performance when resource limits are set high to prioritize application performance. Developers pay the cost when cloud providers charge them based on resources used [1]–[3]. Cloud providers charge them in serverless computing [4]–[7]. Due to this trade-off, setting accurate limits is important. In practice, it is also difficult [1], [8]–[11]. Using profiling to sources are not available which makes application performance hard to predict. Developers pay the cost when cloud providers charge them based on resources used [1]–[3]. Cloud providers charge them in serverless computing [4]–[7].

The aggregate CPU utilization at Twitter is <20% but the maximum reaches up to 60%. Moreover, the utilization is only slightly better at 40–50% when reservations will greatly exceed the usage [10].

characterize application resource requirements will only result in accurate estimation if there is a representative workload. As workloads are often dynamic, the resources needed will change over long timescales (day, pattern, gradual changes in application usage, etc.) and short timescales (changes of coupled systems, etc.). Since estimating an accurate estimate of resource requirements is so complex, developers and operators often resort to static provisioning resources. This results in underutilization or degradation which are often observed by datacenter operators [10], [12]–[14].

[9] by leveraging machine learning to predict future needs and then automatically scaling container resource limits based on these predictions. These works eliminate the developer burden of setting resource limits (e.g., several minutes) to set resource grants more frequently (e.g., several seconds).

Coarse-grained intervals are required because the system has enough information to make resource reservations. This is a poor fit for some workloads where short-lived containers exist as in serverless systems [15]–[18]. Coarse-grained intervals can increase the dynamics of applications can change throughout an interval. Thus, these works still contend with the performance and efficiency trade-off.

In this paper, we argue the performance and efficiency trade-off can be avoided by using a *fine-grained, event-based* resource allocation scheme. To this end, we introduce Escra: a fine-grained, event-based resource allocation system capable of managing containers' resources across multiple nodes. We find resource allocation can easily adapt to sub-second performance within and across hosts, allowing container operators to effectively scale and assign resources. Without performance contention, this scheme uses numerous benefits instead of a constant being killed when it reaches an OOM event, an even-based system can catch the trend and scale the resource allocation to avoid performance degradation over coarse-grained time intervals. A *fine-grained* system can always aim to right-fit allocations to current resource demands, and can quickly respond to instances of CPU throttling.

Escra's design couples a logically centralized controller that schedules resource allocations to containers across servers. Each server is instrumented with kernel hooks and an agent process to apply resource decisions and report container usage to the controller. A *Distributed Container abstraction* enforces per-application resource isolation and allows

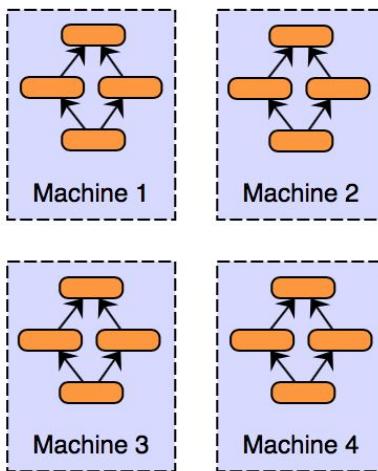


University of Colorado **Boulder**

Distributed Deep Learning

Background

Data Parallelism

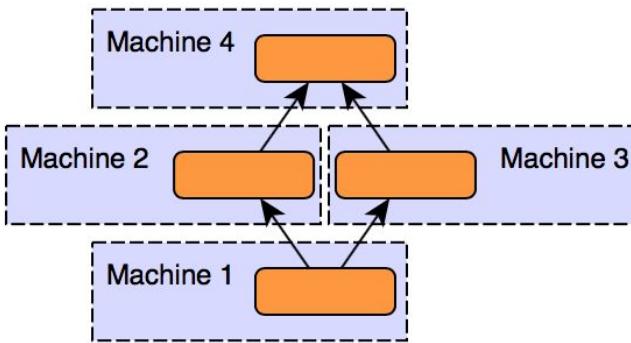


E.g.

 PyTorch

 TensorFlow

Model Parallelism



Background

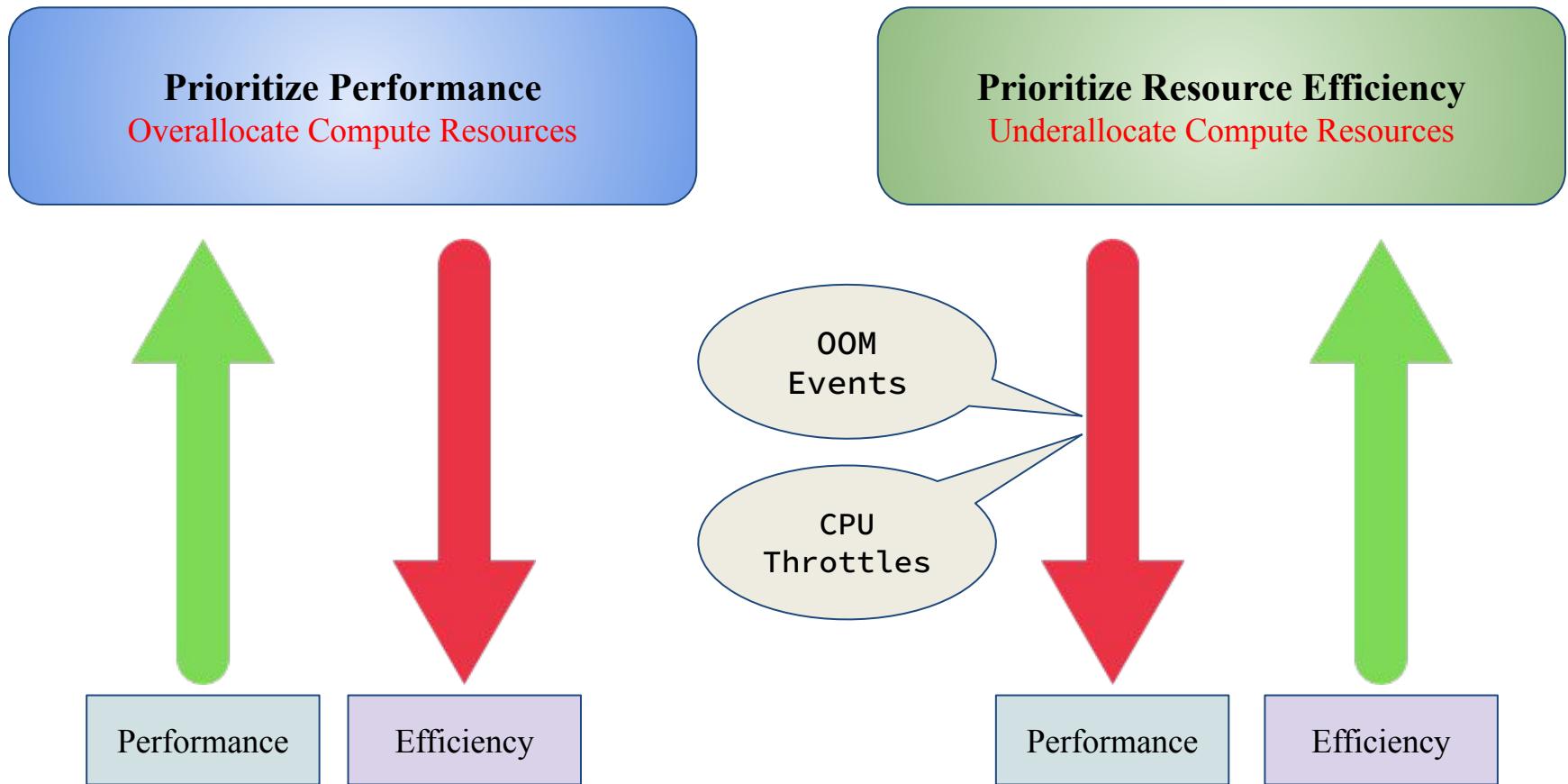


Containerized Infrastructure

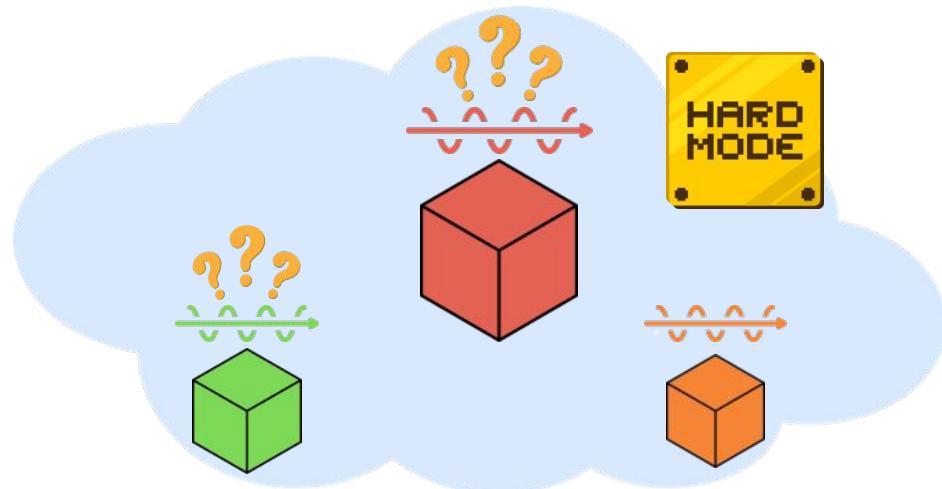
- Light-weight
- Rich Orchestration Systems
- Development Workflow Integration
- Strong Resource Isolation

Introduction: Problem Statement

Developers must make trade-offs!



Introduction: Problem Statement



Setting Container Limits is Hard!

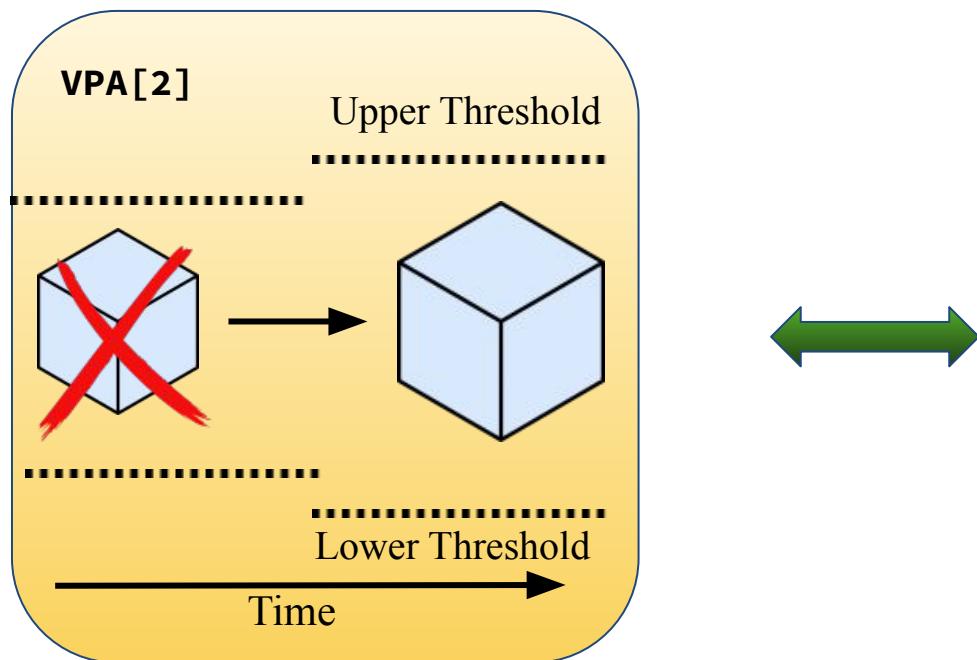
- Requires representative workload
- Resource needs vary over time
- Tools aggregate usage over coarse-grained timescales

→ Overprovisioned Containers



Related Works

Vertical Pod Autoscaler!



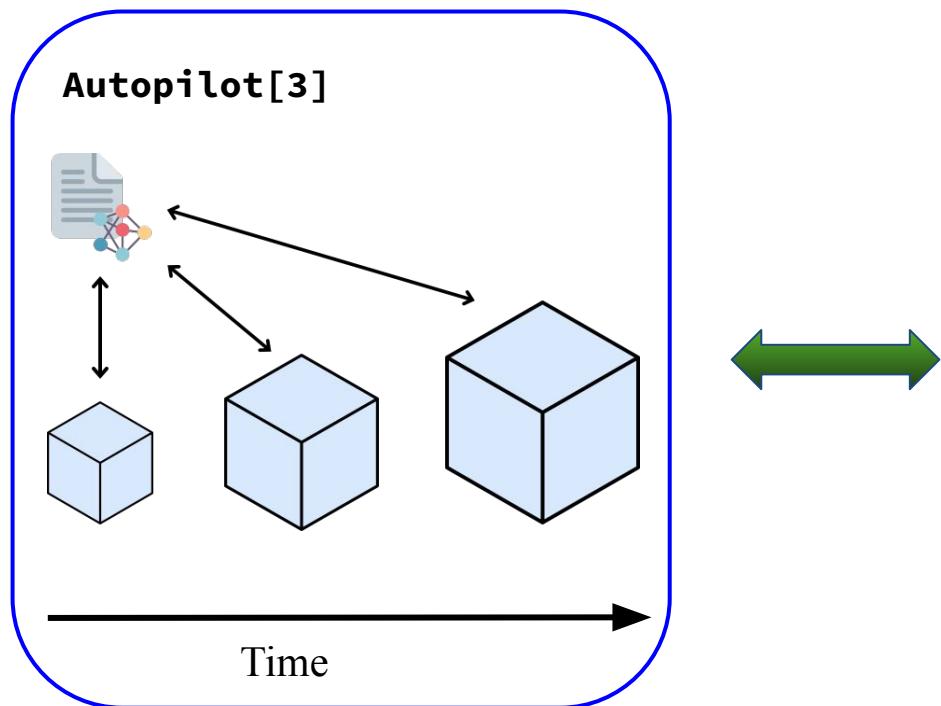
Shortcomings:

- Containers restart to scale
- Scale once per minute
- High slack



Related Works

Google Autopilot: RL-like Algorithm



Shortcomings:

- Cannot respond to quick changes in the workload
- Must allocate to maximum prediction in the next 5-min interval
- Unable to correct inaccurate predictions
- Not suitable for short-lived containers



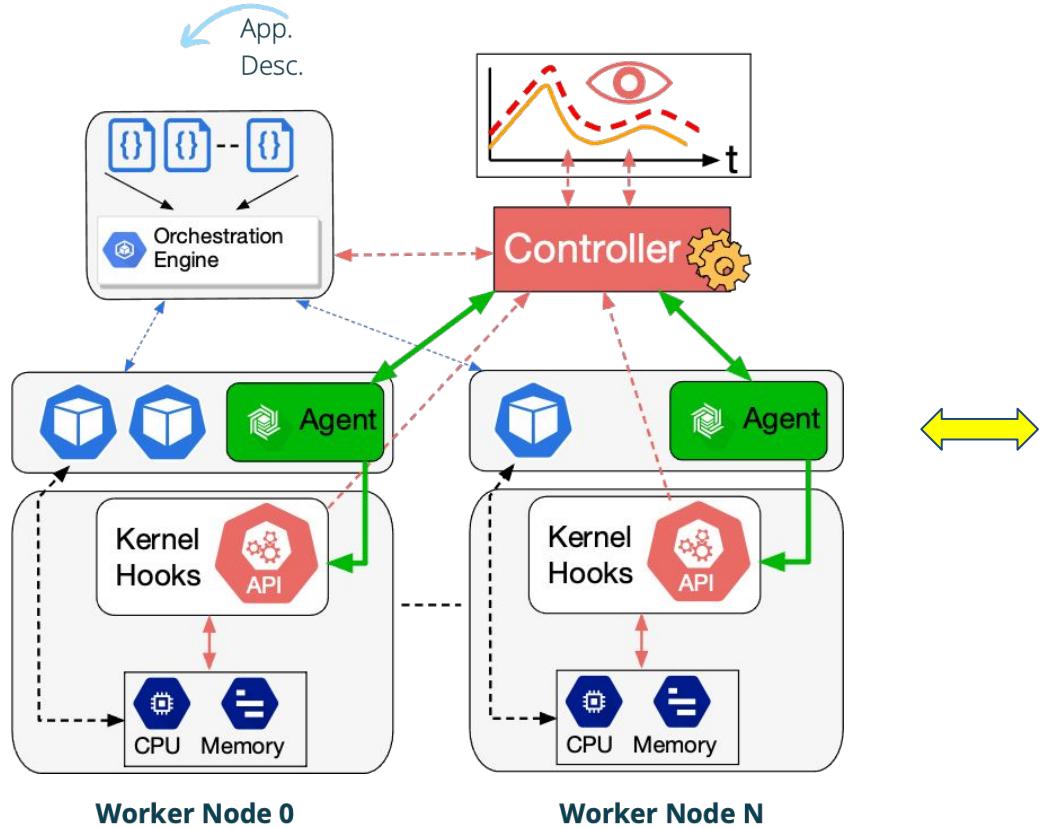
Event-driven Sub-second Container Resource Allocation

Escra:

- Distributed Resource Allocation
- Sub-second interval resource scaling within and across hosts
- OOM prevention and scaling
- Immediate CPU throttle response
- No performance penalty



System Design



Dynamic Resource Allocation Process

Dynamic CPU Allocation:

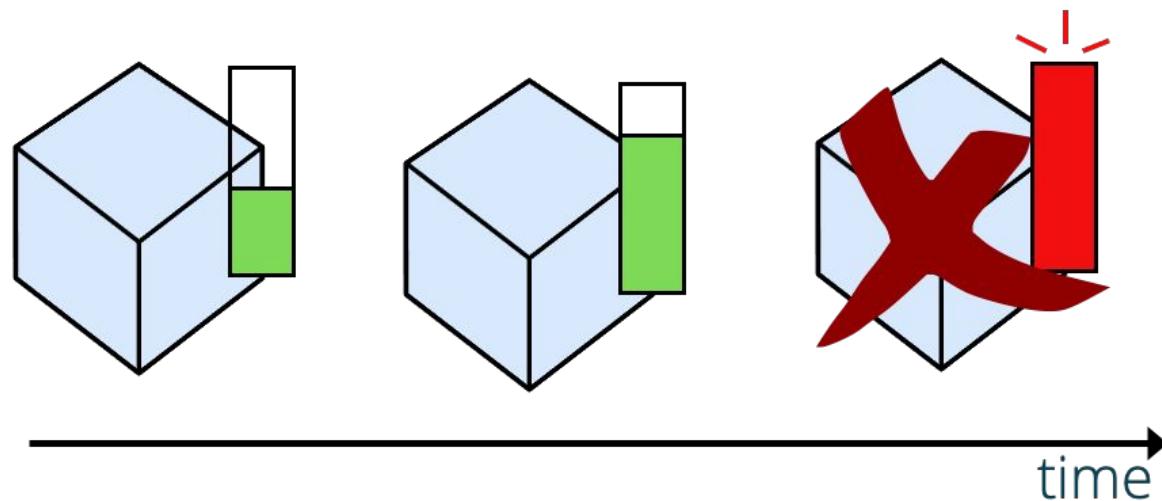
- Scale up and down the CPU quota
- Based on sliding window algorithm over n previous periods
 - Fine-grained, per period interval telemetry data
 - *Past quotas*
 - *If the container is throttled*
 - Refer to the paper for the formula details
- Respects the total cpu limit in our distributed setting



Dynamic Resource Allocation Process

Event-based Memory Allocation:

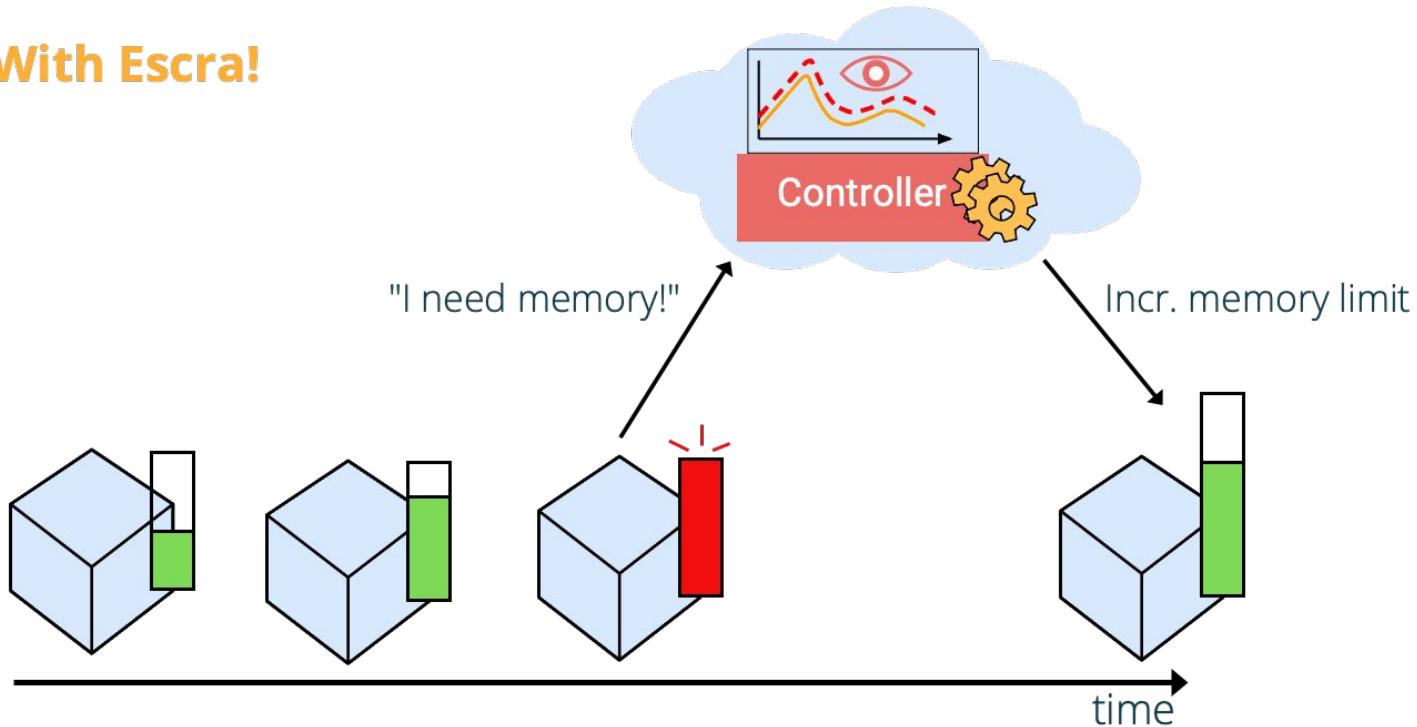
Typical Scenario



Dynamic Resource Allocation Process

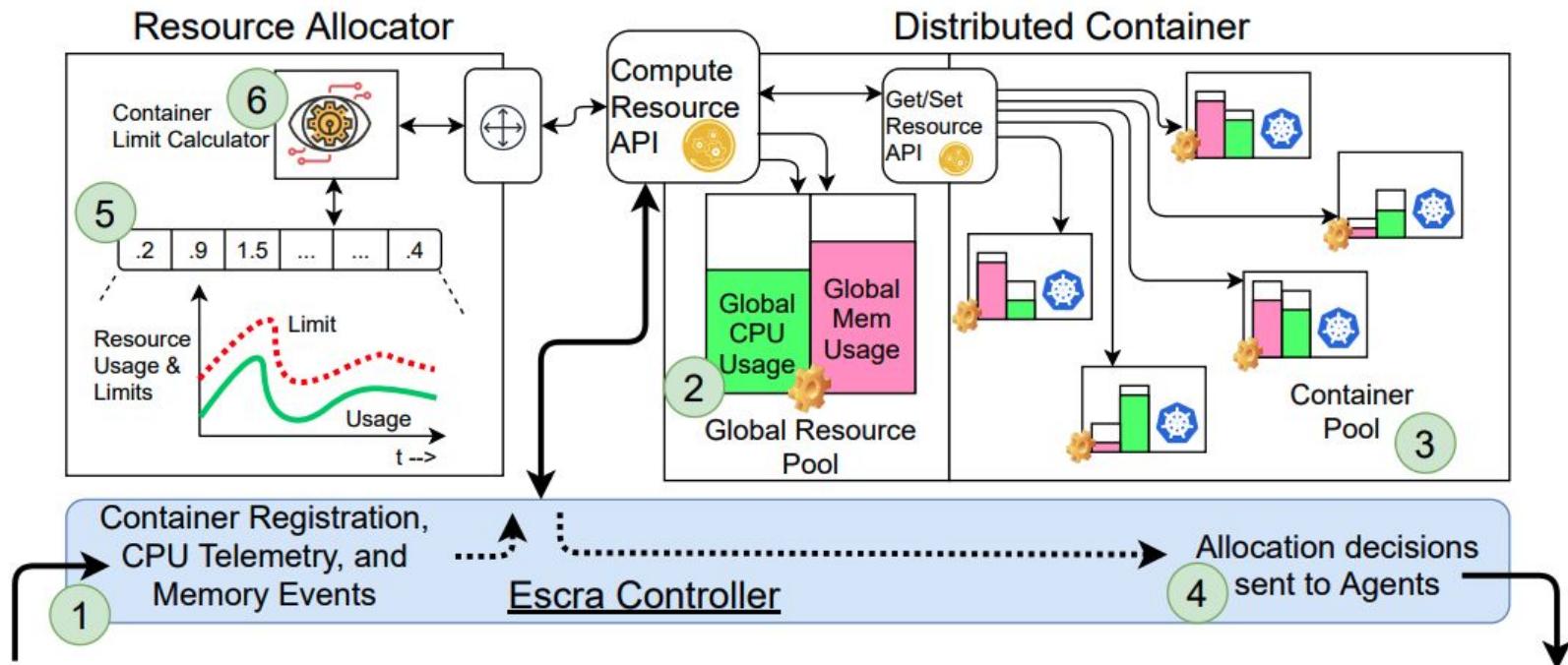
Event-based Memory Allocation:

With Escra!



Escra

Escra Workflow and Components:



Escra Evaluation

Benchmarks Metrics

Absolute Slack

Container Limit -
Container Usage

Application Throughput

Successful requests/second

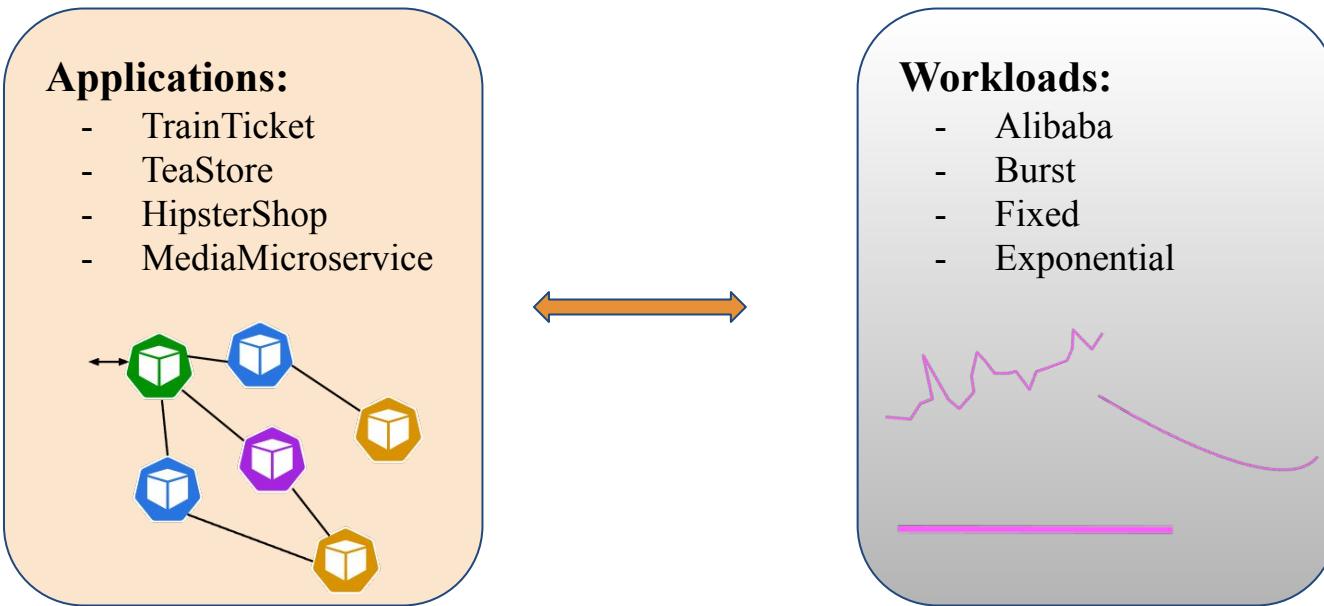
Application 99.9%ile Latency

99.9%ile end-to-end latency



University of Colorado **Boulder**

Applications and Workloads

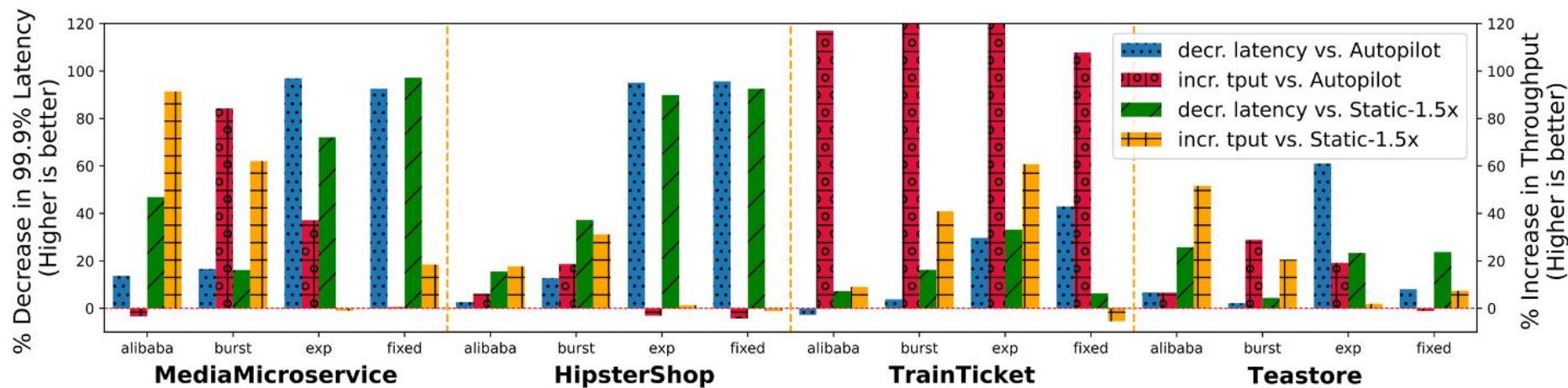


Note: Escra Integrated with Kubernetes to run these apps



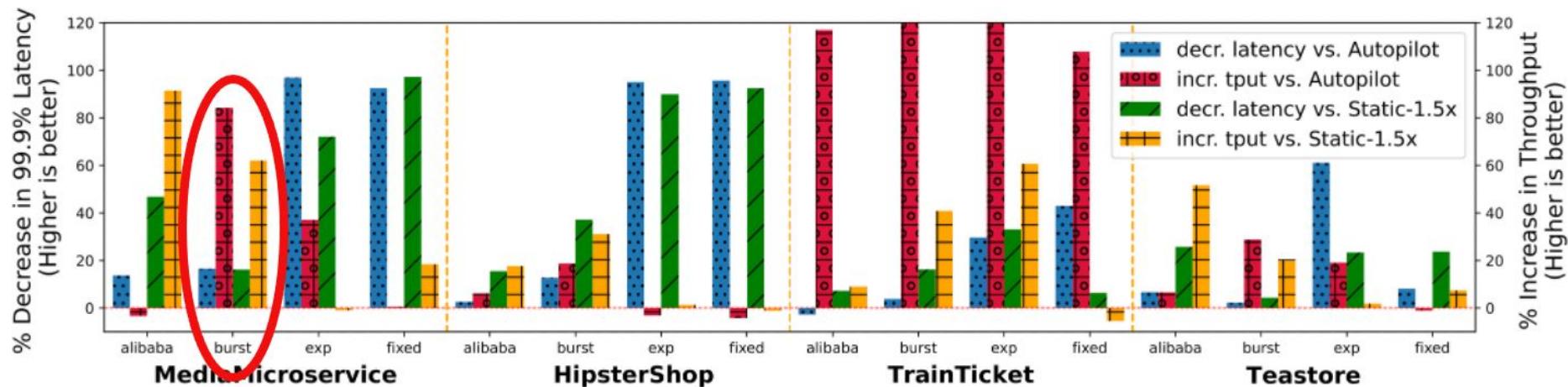
Application Performance Evaluation

Application latency/throughput comparison against Autopilot and static resource allocation



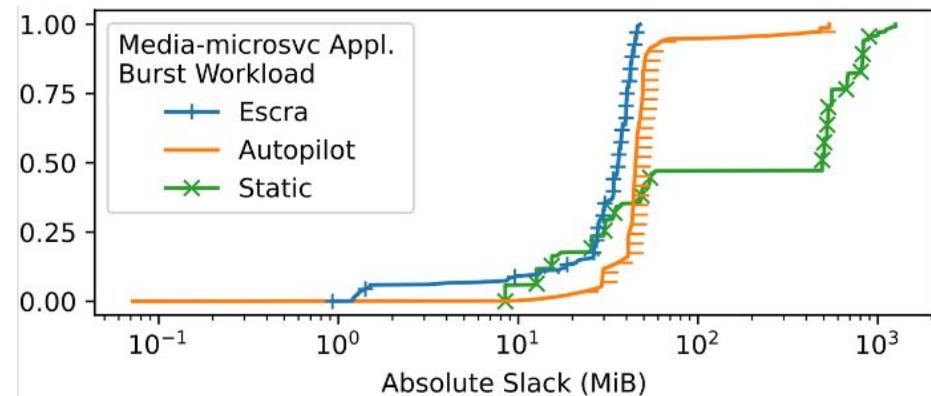
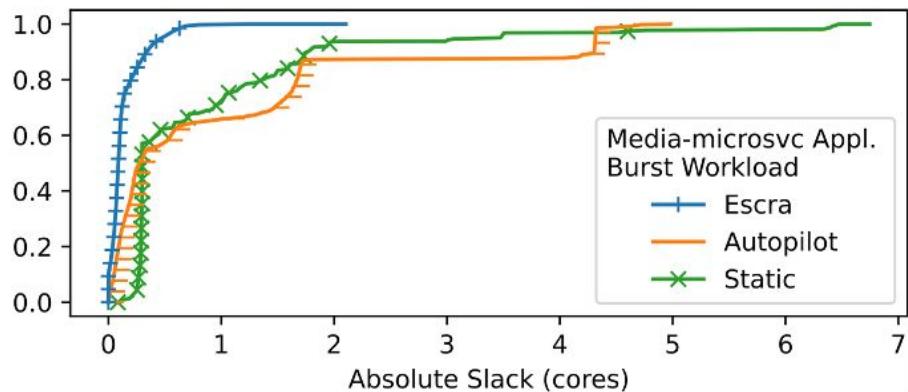
Escra Wins!

Microservice Benchmark: Burst Workload



Escra Reduces Waste

Microservice Benchmark: Burst Workload Escra vs Autopilot



Application performance with Escra vs Autopilot (SOTA):

Throughput increase: 16.6%

Latency Decrease: 84.3%



Escra Evaluation

Serverless Benchmark Metrics

Aggregate Limits
 Σ container limits -
 Σ container usages

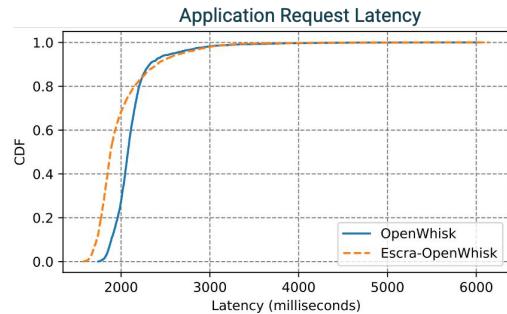
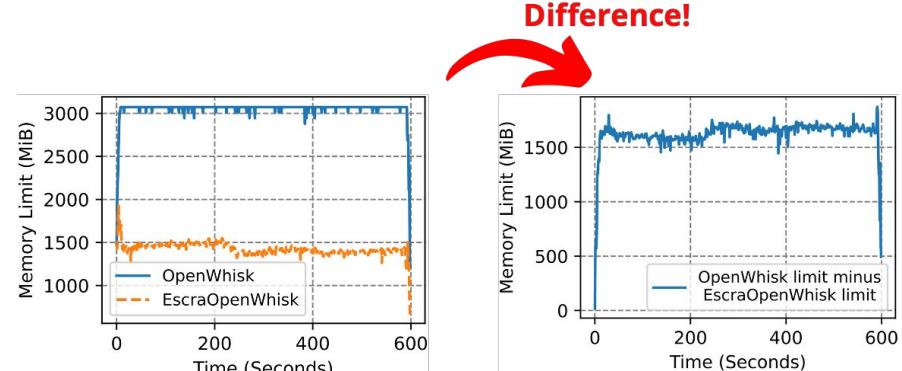
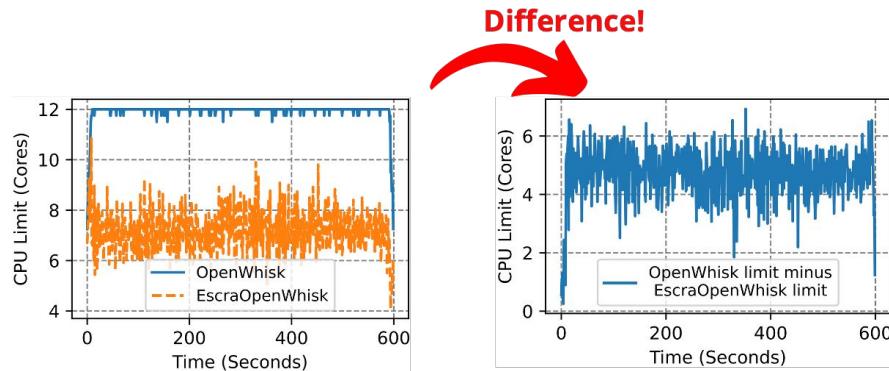
Application Latency
End-to-end latency per
request or job



Escra Integration with OpenWhisk

Serverless Benchmark

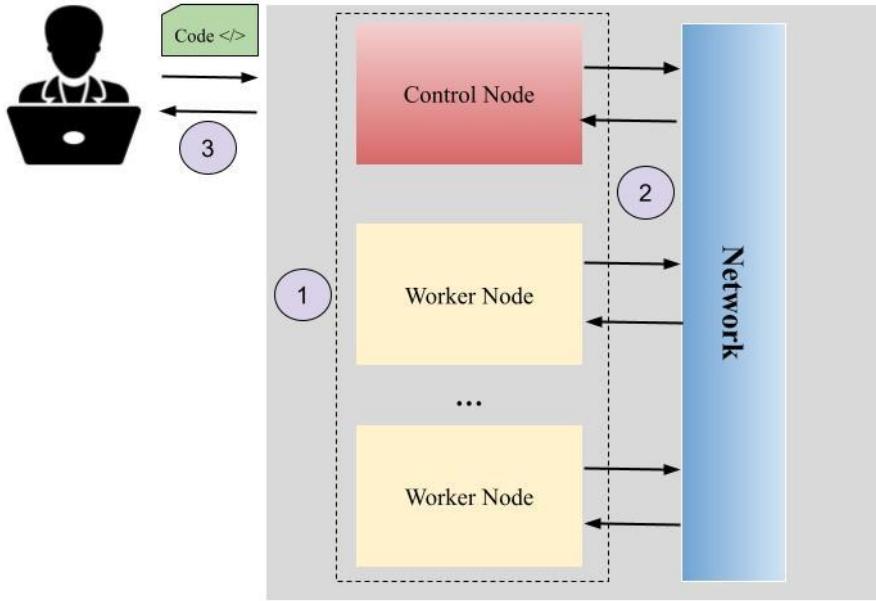
- ImageProcessing Application



- 80% of the times Escra + OW has strong performance gain
- Avg CPU Saving: 5 Cores
- Avf Memory Saving: 1550 MiB



Recap

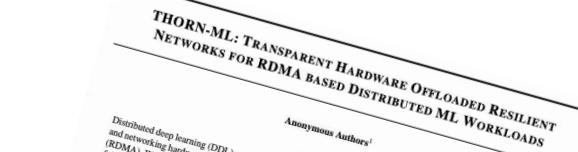


- 1 Escra, Event-driven Sub-second Container Resource Allocation
- 2 THORN-ML: Transparent Hardware Offloaded Resilient Networks for RDMA based Distributed ML Workloads
- 3 Optimizing Cloud Configuration Selection for Cost-Effective Distributed Deep Learning Execution



THORN-ML: Transparent Hardware Offloaded Resilient Networks for RDMA based Distributed ML Workloads

[Under review @ MLSys 2025]



Anonymous Authors¹

ABSTRACT

Distributed deep learning (DDL) requires a great investment in infrastructure, including accelerated compute nodes and networking hardware capable of supporting high performance networking, e.g., Remote Direct Memory Access (RDMA). When a host running a DDL application becomes unreachable, the cost can be high at application-level; failure recovery is slow and disruptive. When the host is unreachable due to a core data center network fail, we argue that the cost is avoidable. This paper introduces THORN-ML, a hardware-offloaded resilient network architecture that is completely transparent to DDL applications and works with commodity hardware. We evaluate THORN-ML on a cluster of 5 nodes with Nvidia A100 GPUs and Mellanox ConnectX-5 NICs, with several applications leveraging model parallelism and/or data parallelism, and find that THORN-ML reduces disruption from minutes (impacting the whole cluster) to milliseconds (impacting packets that can be re-transmitted).

In the event of a failure, the DDL training job can restart from the last checkpoint, avoiding the need to start from scratch. However, this mechanism comes with checkpoints overhead, as the workers must revert to the checkpoints in our measurements on a 2 node GPU cluster training the GPT-2 model and ResNet101. A single epoch takes approximately 4 to 80 seconds, and the time to initialize all nodes runs from 4 to around 10 seconds (details of Section 2). Combined, this is the time the resources required for this job cannot perform any useful computation. Our measurements assume taking a checkpoint every epoch, which has I/O storage and compute overhead, even though these are likely to be performed at longer intervals, and (2) these measurements were gathered on a small cluster, and are likely to be higher on bigger clusters.

(1) These measurements assume taking a checkpoint every epoch, which has I/O storage and compute overhead, even though these are likely to be performed at longer intervals, and (2) these measurements were gathered on a small cluster, and are likely to be higher on bigger clusters.

The large (and increasing) number of components increases the likelihood of failure of individual components. Given how failures are handled, this simultaneously increases the impact of a failure. Consider the case where a host fails DDL frameworks, such as TensorFlow (Tensorflow), Ray (Ray), and PyTorch (Pytorch), handle host-level failures by using distributed checkpoints. Training usually involves several iterations or epochs, where forward-backward propagation occur in a neural network. A checkpoint of the state of the application (which is distributed across the cluster) is then recorded after each epoch. Specified intervals based on application configuration,

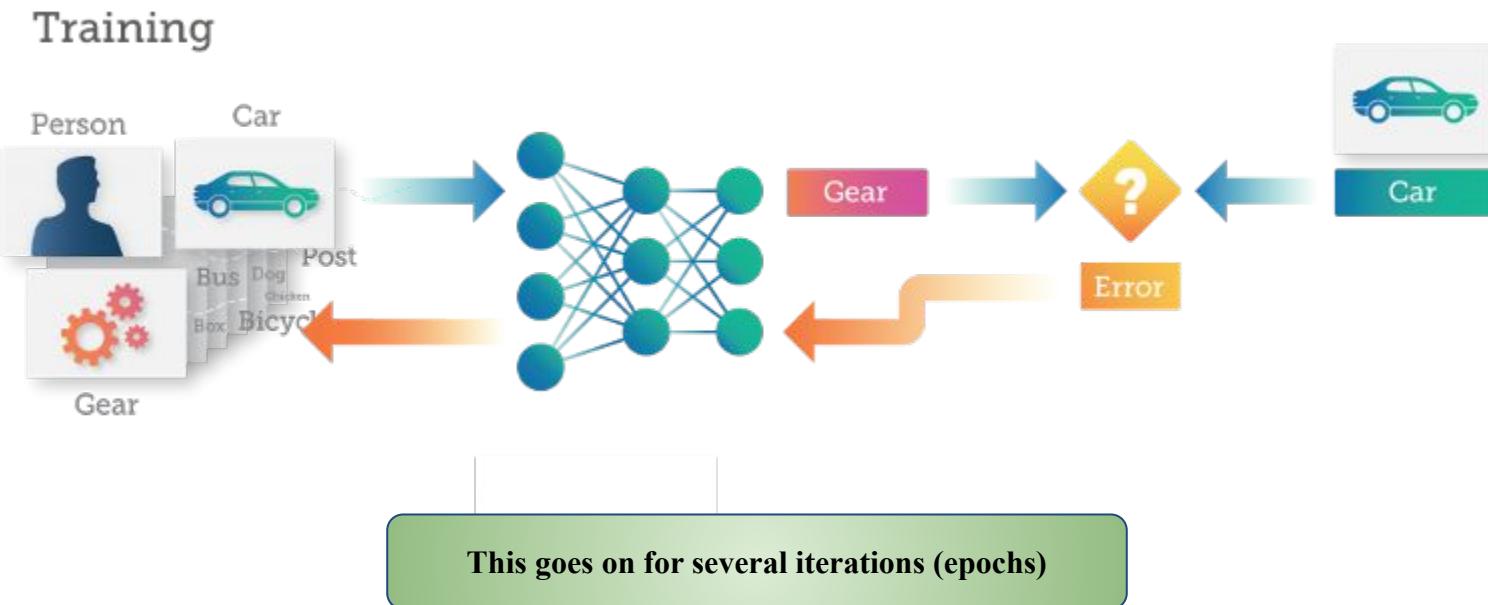
¹Anonymous Institute, Anonymous City, Anonymous Region, <anon@emaildomain.com>, Correspondence to: Anonymous Author, Preliminary work. Under review by the Machine Learning and Systems (MLSys) Conference. Do not distribute.

On the network side, core data center networks are designed around redundancy and routing protocols to detect and route around failures with little disruption (Zengberg et al., 2009; Kouvelis et al., 2022; Zengberg, 2018). However, redundancy and protocols do not extend to the edge of the network, leaving the connectivity between the hosts and the ToR switches vulnerable to failure (accounting for nearly 10% of over-disruption based on measurements from Mellanox Dubey et al., 2024). The components at this connection point include the NIC on each host, the ToR switches themselves, and the cables that connect the NICs to the ToR switches. Current practice detects the hosts as unreachable, whether due to the host failure or failure at the network edge, and as such, treat these failures equally.



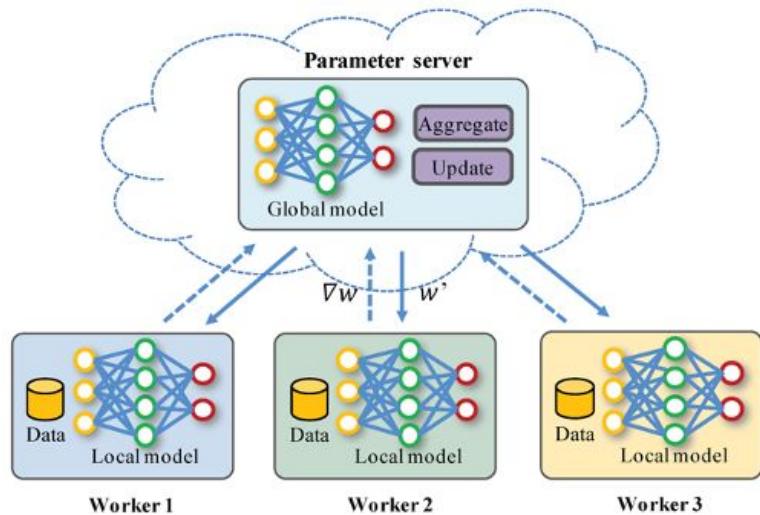
University of Colorado **Boulder**

Background - Deep Learning



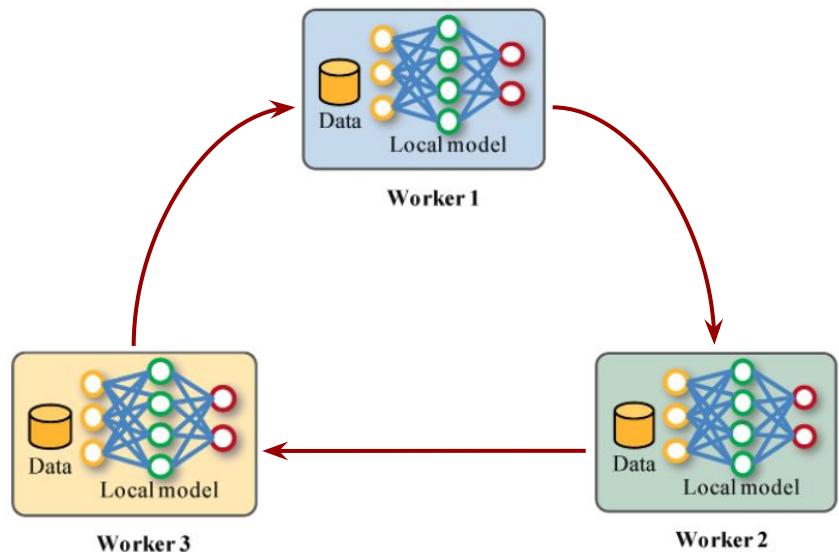
Background - Distributed Deep Learning

Data Parallelism



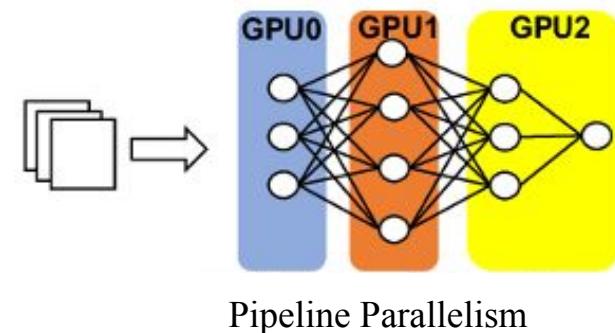
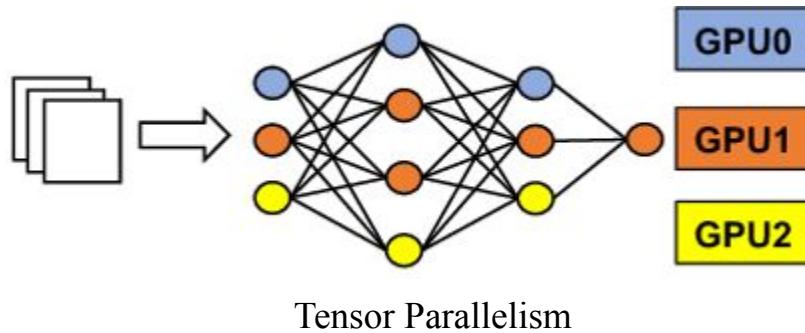
Parameter Server Strategy

Ring (All-reduce)

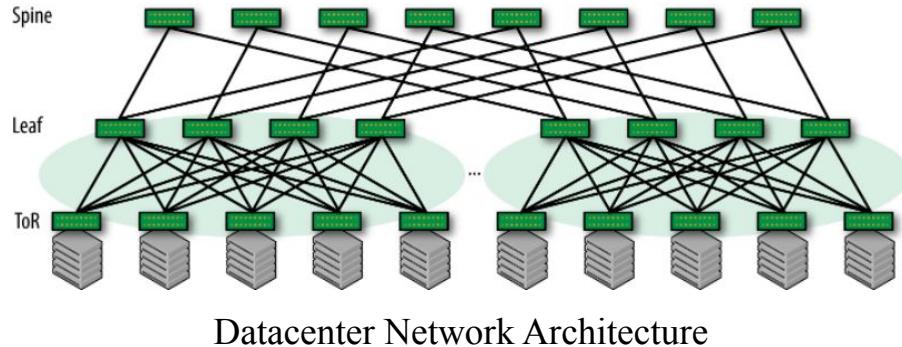


Background - Distributed Deep Learning

Model Parallelism

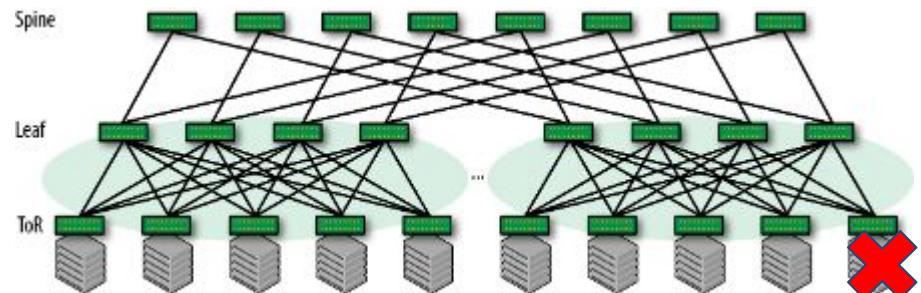


Background - Impact of Failure

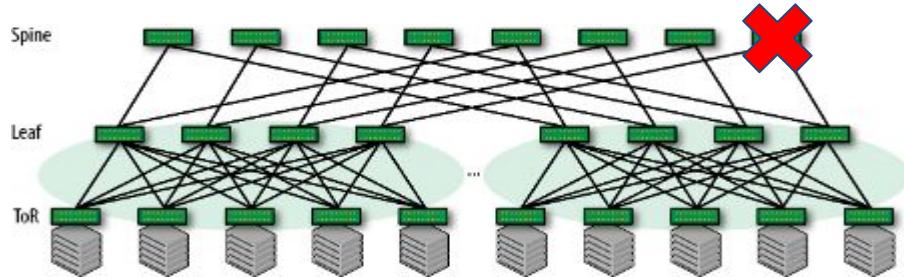


Popular Deep Learning Frameworks:

- Take checkpoints every X epochs (Hyperparameter)
- Restart the training from the very recent epoch

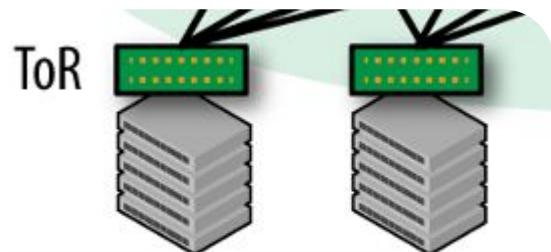


Background - Impact of Failure



Failure at network core is handled with minimal disruption!

- Redundancy
- Advance routing protocols



How about network edge?

- ToR switch
- NIC
- Cable in between



Background - Impact of Failure

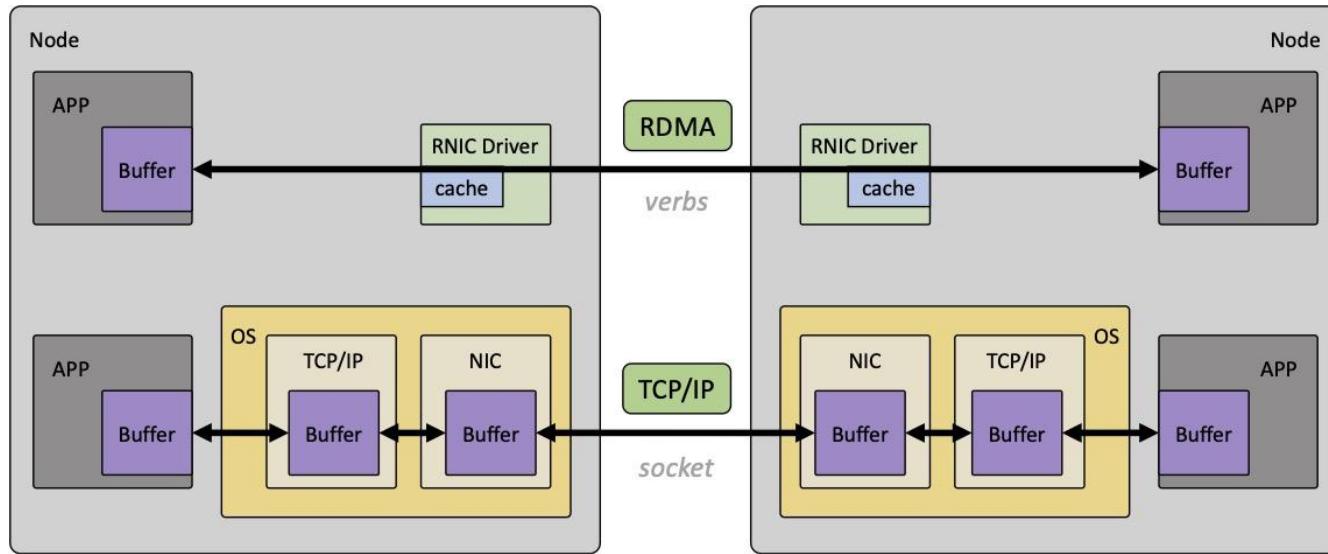
Arch	Checkpointing (s)	Model Init (s)	Epoch time (s)
Data parallel	0.77	22.61	4.47
Pipeline parallel	0.53	21.28	8.10
Pipeline & Tensor parallel	0.34	20.91	9.90

Checkpointing and initialization latency/overhead in GPT-2 training on CodeParrot dataset using Megatron-LM

- Average of 1-2 minute disruption (entire cluster) to deal with server failure (at scale of 5 GPUs – higher for larger scale)
- When failures happen, disruptions are costly (e.g., For Meta, a cluster is 32,000 GPUs)
- Network failures account for 8.4% (switches and cables), and 1.7% (NIC) of overall disruptions [10]
- Authors in MegaScale [11]: the init time for Megatron-LM for 2048 GPUs is 1047s => Very Disruptive



Background - RDMA



<https://wiki.huangxt.cn/papers/network/RDMA-Related>

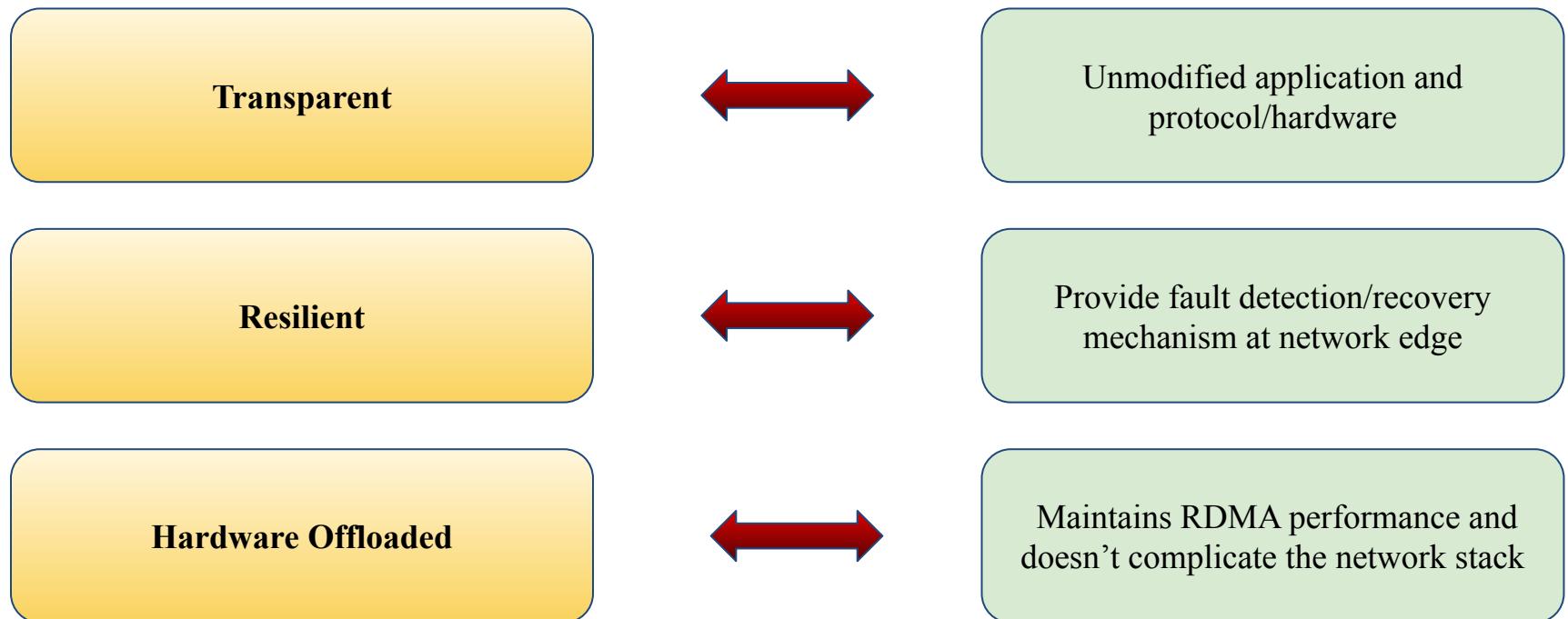
How to use RDMA

- Get RDMA Devices
- Open a Device
- Register Memory Region
- Setup queue pairs
- Talk!



In this project ...

We propose a solution that is



Design

High-level design decisions in THORN to achieve aforementioned goals

Single virtual device that the application can bind to and use RDMA API w/o modification

Tunneling the traffic across hosts to hide the underlying network redundancy

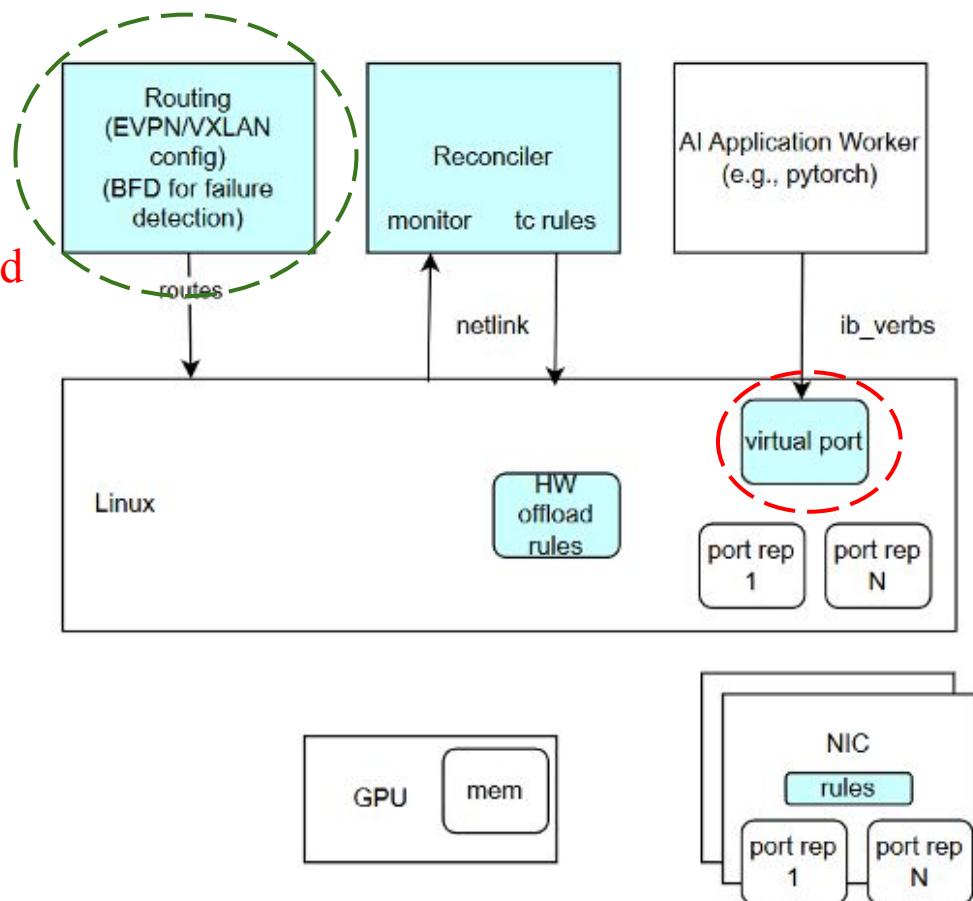
Use of standard routing protocols to detect failures and select alternate paths



Architecture Overview

Unmodified Application:

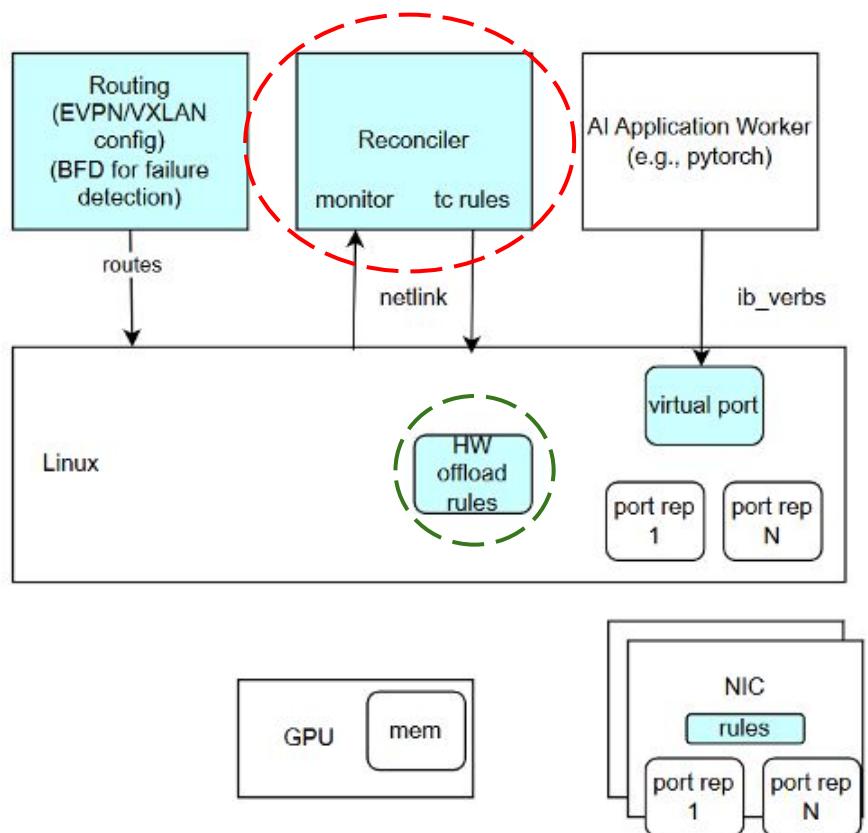
- Single virtual device (SR-IOV) exposed to the application to bind
- Routing software for failure detection and re-routing



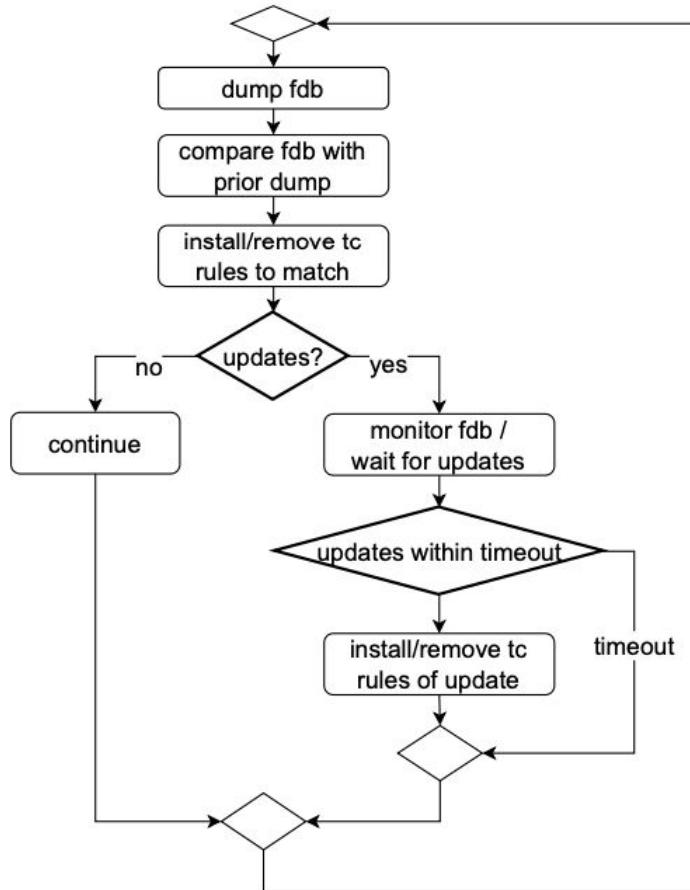
Architecture Overview

Unmodified Hardware

- Reconciler that converts network routes to local rules that direct traffic between virtual port and port representors
- Hardware offload to NIC (programmable)



Architecture Overview: Reconciler



Reconciliation loop for offloading the routing table updates

- Uses *Netlink* to talk to Linux
- Ensures correctness



Evaluation

Testbed Summary:

Cluster Composition: 5 hosts and 2 Mellanox SN2700 switches (48-port, 100Gbps).

Host Configuration:

- Intel Xeon Silver 64-core CPU (2.10GHz).
- Nvidia A100 GPU (16 GB memory).
- Dual-port, 100Gbps Mellanox ConnectX-5 network card (RDMA support).

Network Topology:

- ConnectX-5 card ports connected to separate switches.
- Switches interlinked via an active-backup link pair.

Storage System:

- 1 TB storage node per switch.
- High-throughput, low-latency log collection and dataset distribution.

Macrobenchmark:

ResNet101:

- Dataset: CIFAR10 (60000 32x32 Images)
- Output: 10 Classes
- Framework/Architecture: Pytorch DDP

GPT-2 (Language Model):

- Dataset: CodeParrot (Python Code Generation)
- Framework/Architecture: Megatron-LM
 - Pipeline Parallel
 - Pipeline & Tensor Parallel
 - Data Parallel

NCCL (Nvidia Collective Communication Lib):

- Operation Bandwidth Measurement

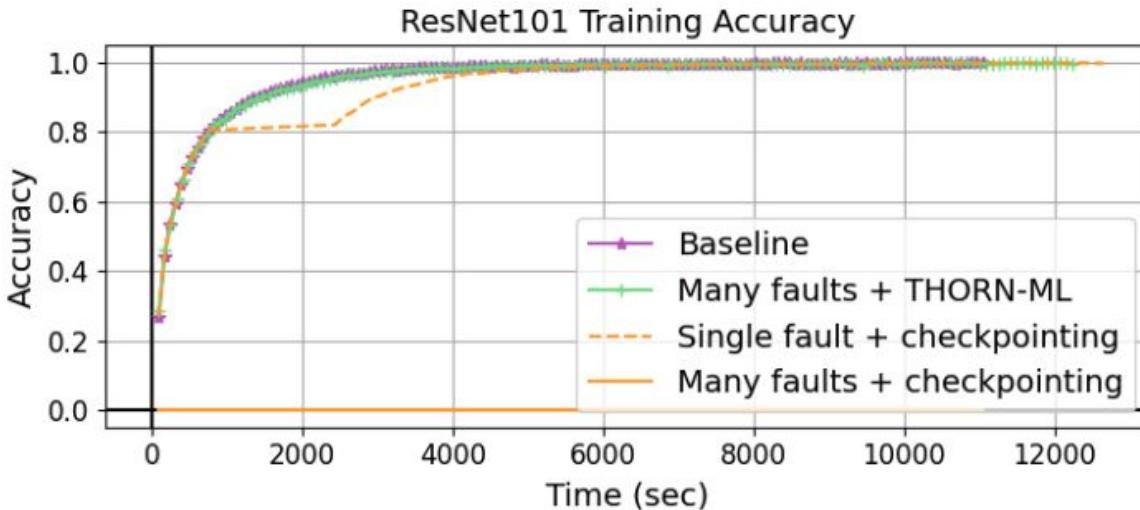
Microbenchmark:

- RDMA performance test (*ib_send_lat*)



Evaluation: Macro-benchmark

THORN is able to keep the overall training job even in extreme fault scenarios w/o noticeable difference with baseline condition



ResNet101 training on CIFAR-10 dataset using Pytorch DDP API on 4 nodes.

Representation of time to accuracy

Note on Single fault with checkpointing:

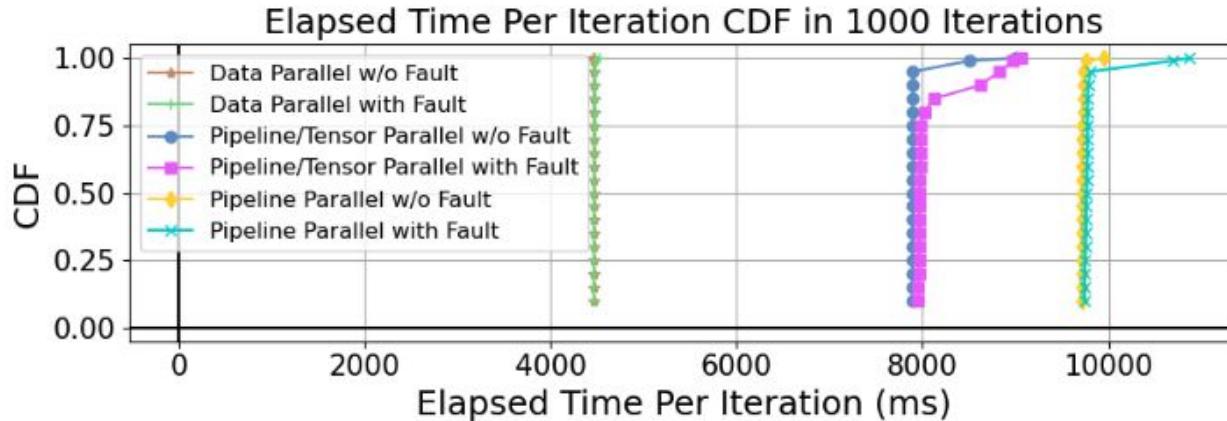
- Checkpoints happen every 10 epochs
- Accuracy remains the same amount for roughly 13 mins



University of Colorado **Boulder**

Evaluation: Macro-benchmark

THORN handles network faults with the overhead of ms scale => order of magnitude less than checkpointing



Distributed GPT-2 training on CodeParrot dataset performance. Each data point is averaged over 10 iterations.

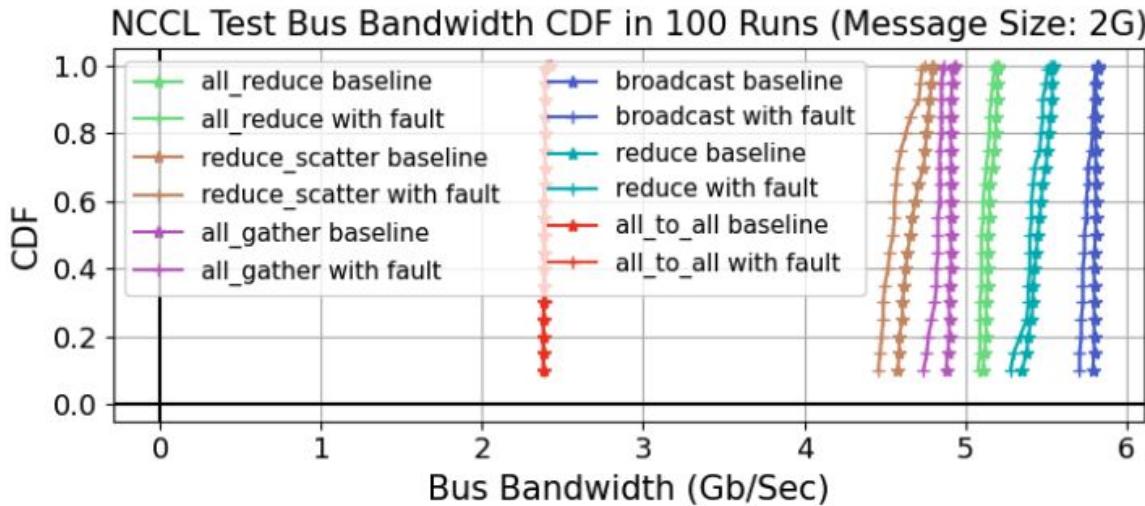
Latency Impact Summary:

- **Data Parallel Architecture:**
 - Added latency at the 90th percentile: **5.3ms** (1.1% overhead).
 - Added latency at the 95th percentile: **7.3ms** (1.6% overhead).
- **Pipeline Parallel Architecture:**
 - Added latency at the 90th percentile: **54.9ms** (5.6% overhead).
 - Added latency at the 95th percentile: **76.7ms** (7.8% overhead).
- **Combined Pipeline and Tensor Parallel Architecture:**
 - Added latency at the 90th percentile: **729.4ms** (9.2% overhead).
 - Added latency at the 95th percentile: **928.4ms** (11.7% overhead)



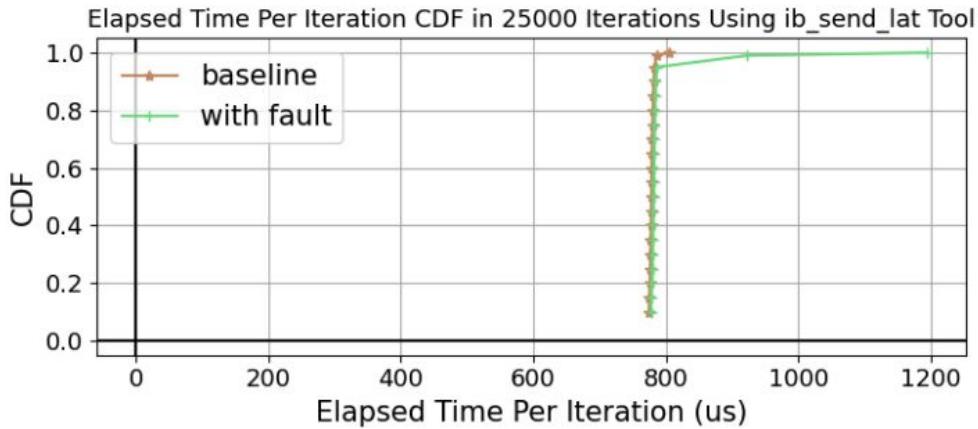
Evaluation: Macro-benchmark

THORN negligible difference (less than 2%) in bus bandwidth of nccl collective operations



- NCCL test collective operations performance report with a message size of 2GB for a baseline without THORN-ML and without faults compared to with THORN-ML and with faults.
 - Similar tests were run with 1GB, 256MB, and 1MB, and each showed that
 - provide nearly the same throughput
- **80th Percentile:** Maximum difference of **0.129Gbps** between fault and non-fault scenarios.

Evaluation: Micro-benchmark



- Using the *ib_send_lat* command with a message size of 2^{23} bytes
- data points with higher delay are due to re-transmission, and the number of packets that need to be re-transmitted indicates the fail-over time
- detection and fail-over time is so small that it only applies in the 99.9th percentile



Related Work

Fast Distributed Deep Learning over RDMA [4]

- Claims that grpc is suboptimal
- Not fault tolerant

Ultima [5] mitigates the straggler/fault impact in DDL in software level

- Provide fault tolerance at some level
- Not leveraging RDMA connections and hardware offload

Popular frameworks like Tensorflow handle faults in software level using checkpointing mechanism

- High overhead

Maestro [6] Proposes multipath RDMA for DDL workload at datacenters

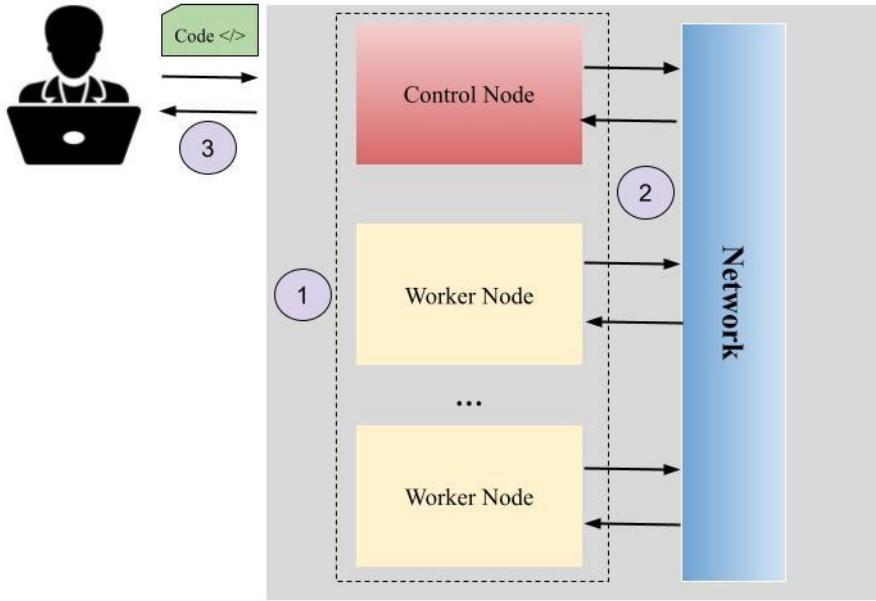
- Implemented in user-space/software
 - Modified Application
 - Low performance

MP-RDMA[12] handles congestion in RDMA networks (address ECMP problem)

- Uses specialized hardware (modified hw/protocol)
- not applicable on the edge of the network



Recap



- 1 Escra, Event-driven Sub-second Container Resource Allocation
- 2 THORN-ML: Transparent Hardware Offloaded Resilient Networks for RDMA based Distributed ML Workloads
- 3 Optimizing Cloud Configuration Selection for Cost-Effective Distributed Deep Learning Execution



Optimizing Cloud Configuration Selection for Cost-Effective Distributed Deep Learning Execution

[In submission to ICML 2025]



University of Colorado **Boulder**

Everybody uses Deep Learning

- Including non-expert users

Distributed Deep Learning is Popular

- Model and Data is getting larger

Numerous options are available in cloud platform

- e.g. Tens of HW Configs



Problem Statement

How to choose the best configuration in the cloud?

which one is more expensive and which one is faster to run my training task?



Problem Statement

Naive solution leads to suboptimal results!

go for the cheapest instance type (\$/hour)



Doesn't guarantee lower cost

Go for the most expensive instance type

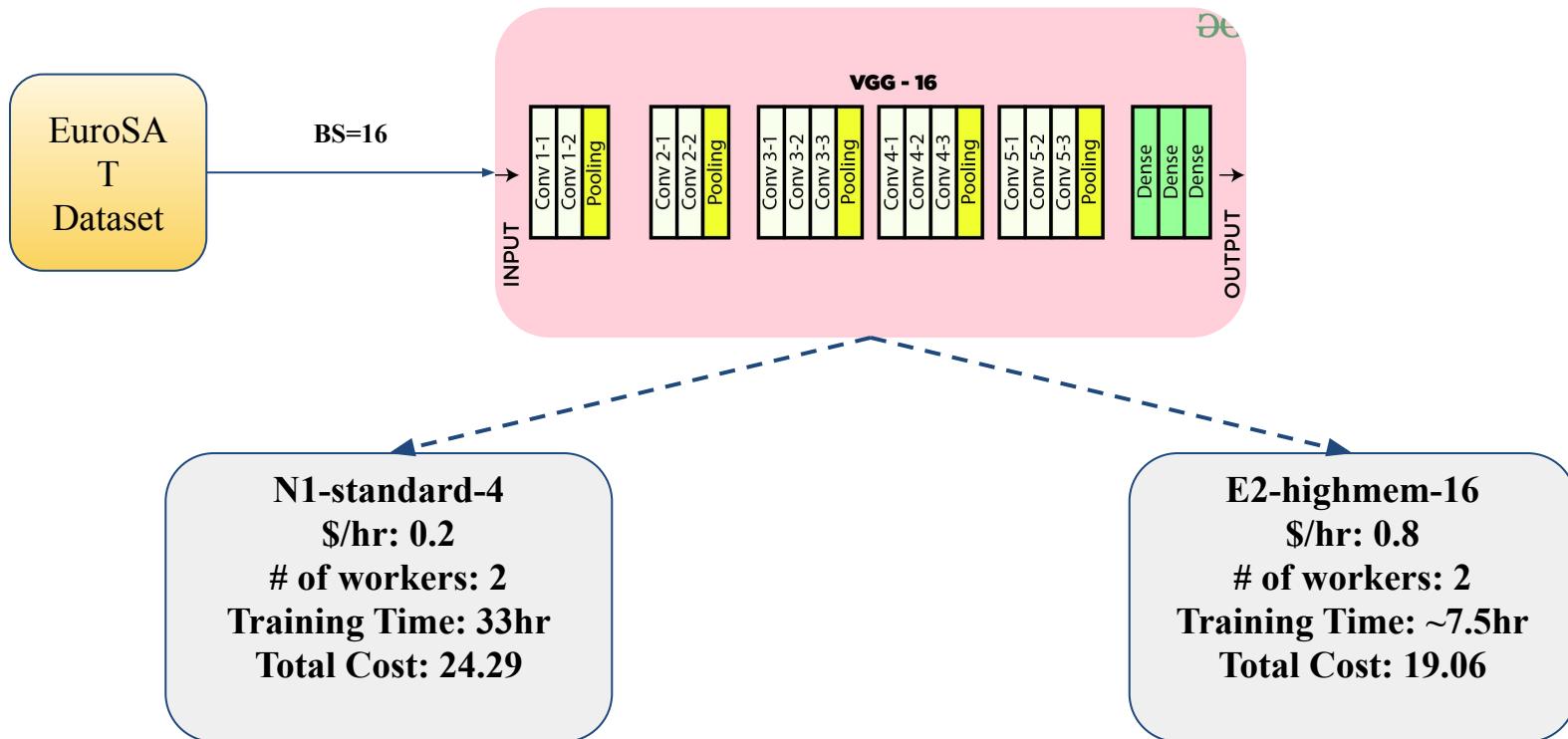


Unnecessary cost to get reasonable performance

There's a sweet spot



Motivating Example



Background Work

CherryPick [8]

- Leverage Bayesian Statistics to predict the best cloud configurations
- Still needs a few sample runs in the cloud
 - *mostly works for recurrent jobs*
- Not tailored for Distributed Deep Learning

Ernest [9]

- Creates resource consumption model of machine learning jobs running on
- Needs retraining for every given VM type
- Not tailored for Distributed/Deep Learning

None of the previous work directly address the problem of cost and training time prediction in distributed deep learning platforms



Our Approach

In many cases, the training task can run partially on a local weak setup (i.e laptop)!

Why don't we leverage a minimal local run as a representative workload
and predict the cost and training time on a powerful cloud?



Our Approach

Deep Learning based solution that

- Input
 - Partial resource usage on local setup
 - Training job features such as parameters/hyperparameters
 - Target desired configuration and setup in the cloud (e.g. 3 workers using e2-highmem-16 instance)
- Predict
 - training time and cost for specific configuration

Compare the performance with traditional machine learning algorithms

Let's hope the user make informed decisions with this results



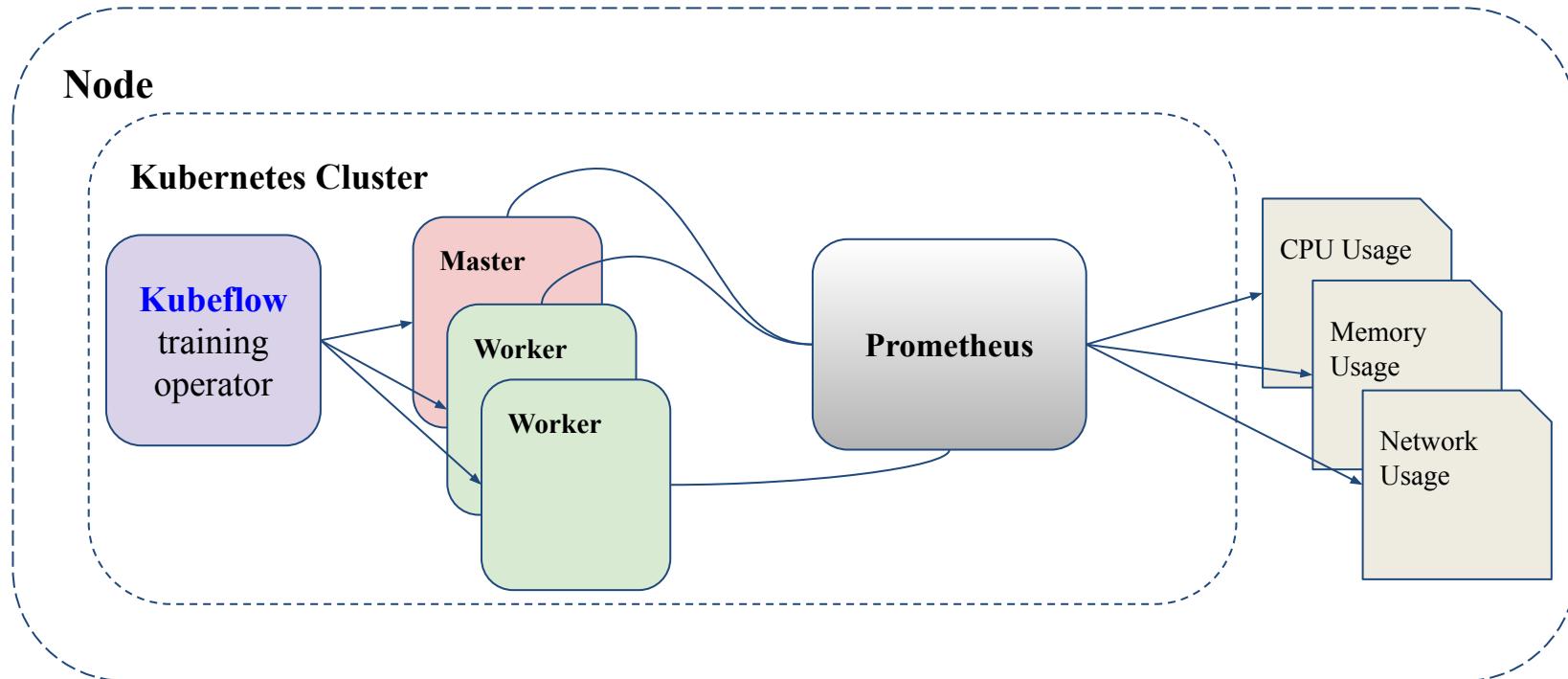
Data Collection

Local

- Represents a partial run on a local laptop (single machine)
- All jobs are run in a 3-worker architecture (1 master, 2 workers)
- Collected partial usage data on 9 different deep learning tasks
 - Resource usage time series
 - *cpu usage*
 - *memory usage*
 - *inbound network usage*
 - *and outbound network usage*



Partial Training Job Data Collection



Data Collection

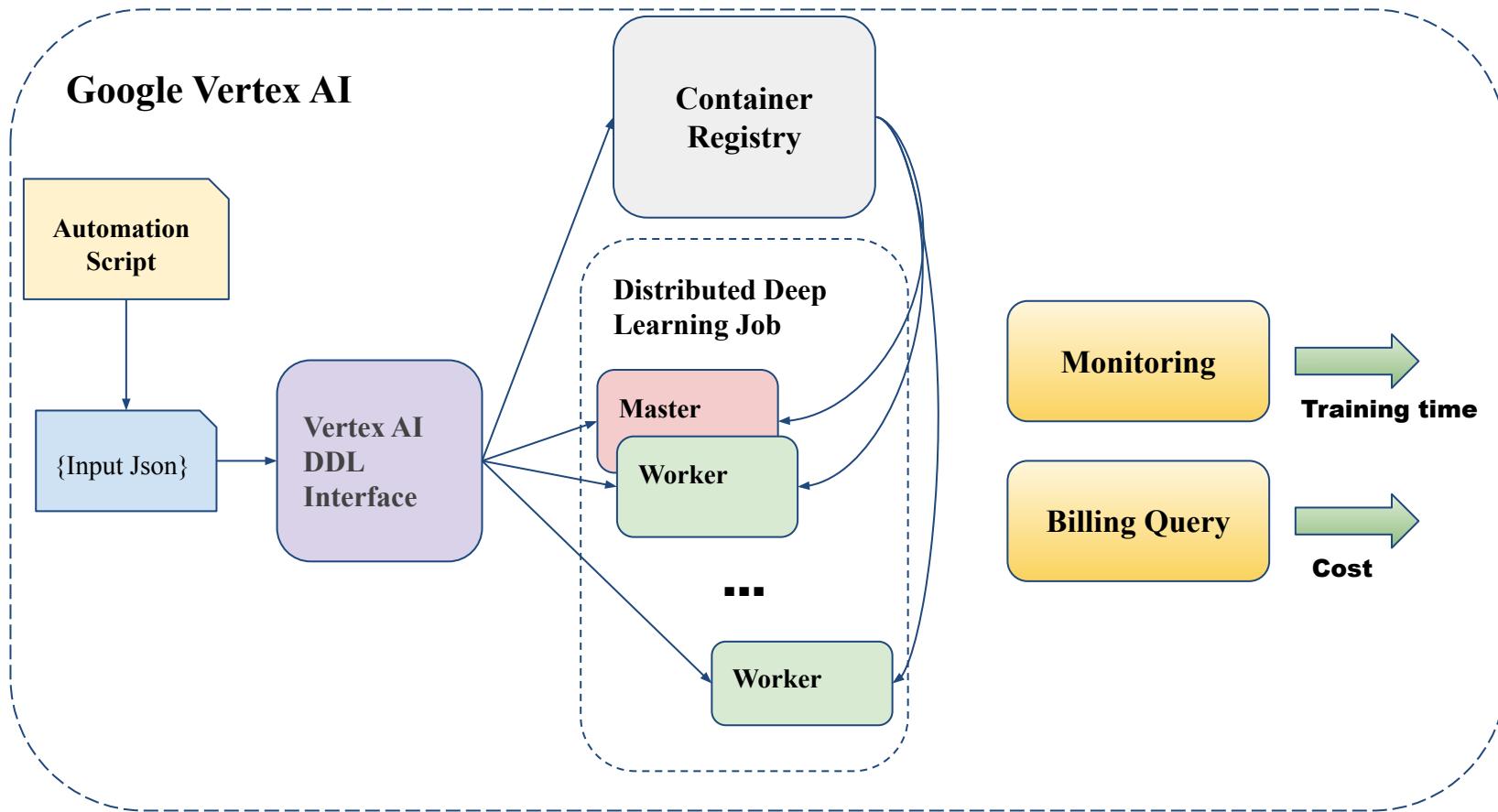
Google Vertex AI

- GCP service offering to run distributed deep learning tasks easily using custom docker containers
- Run each deep learning job with several different configurations
 - Change Model hyperparameters (e.g. *batch size*)
 - Change Resource configuration (e.g. *instance type & number of workers*)
- Record cost and training time for every training job

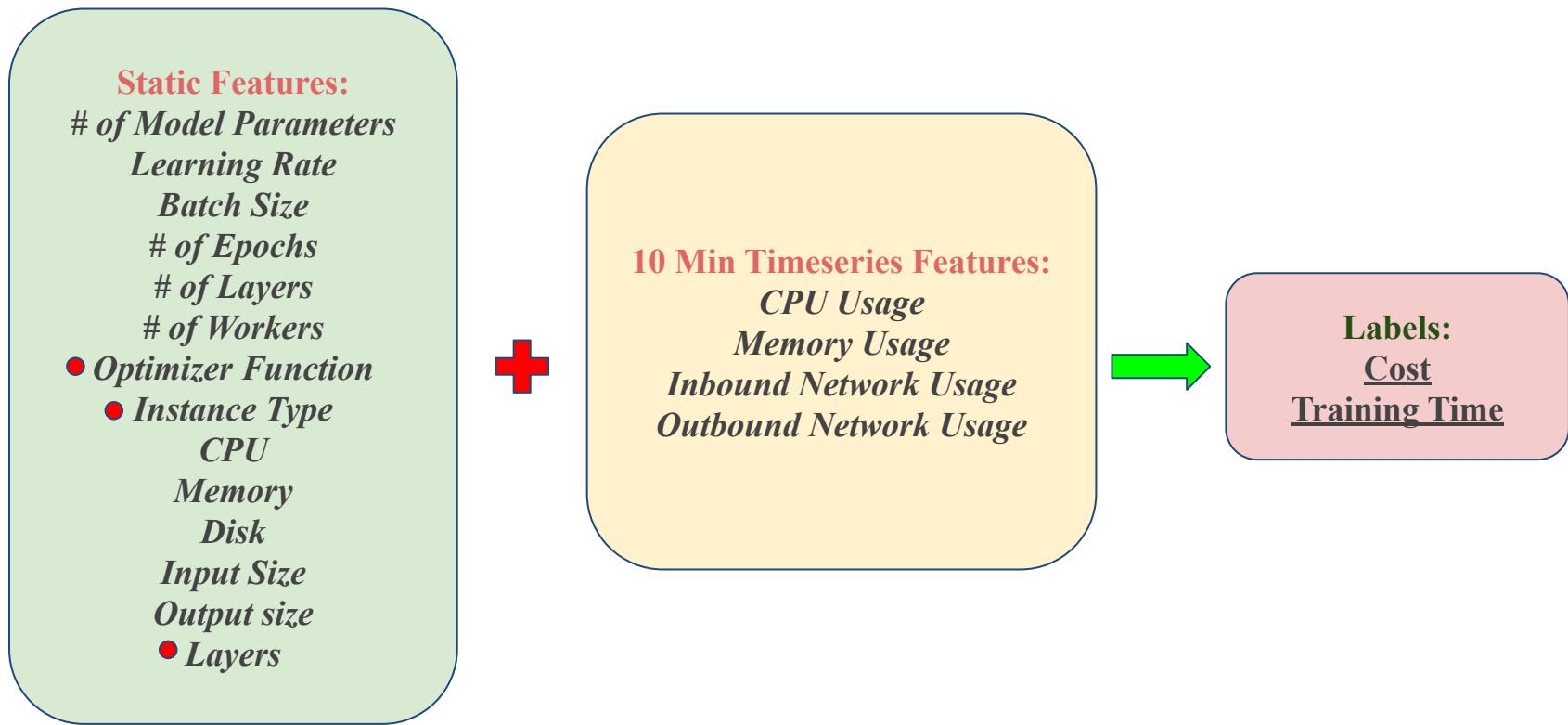


Training Job Data Collection

Dataset Collection:

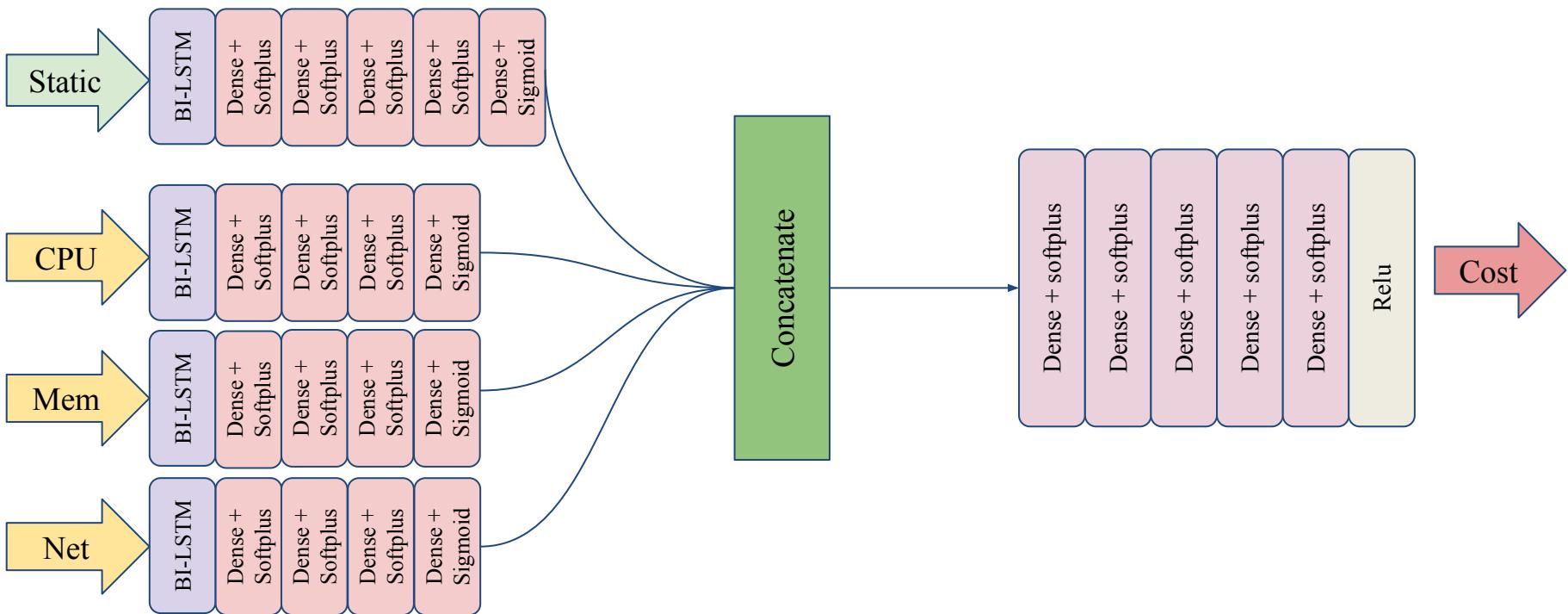


Features and Data Preprocessing



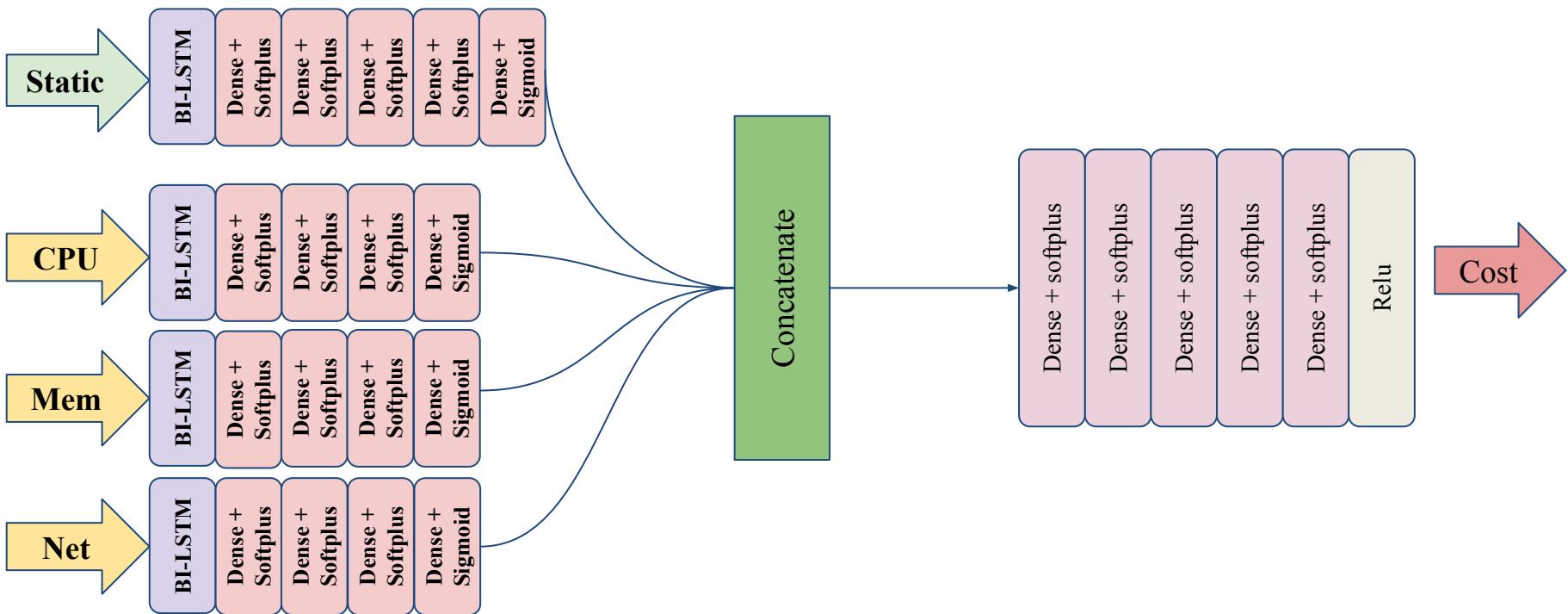
Deep Learning Model Design

Cost Model Development:



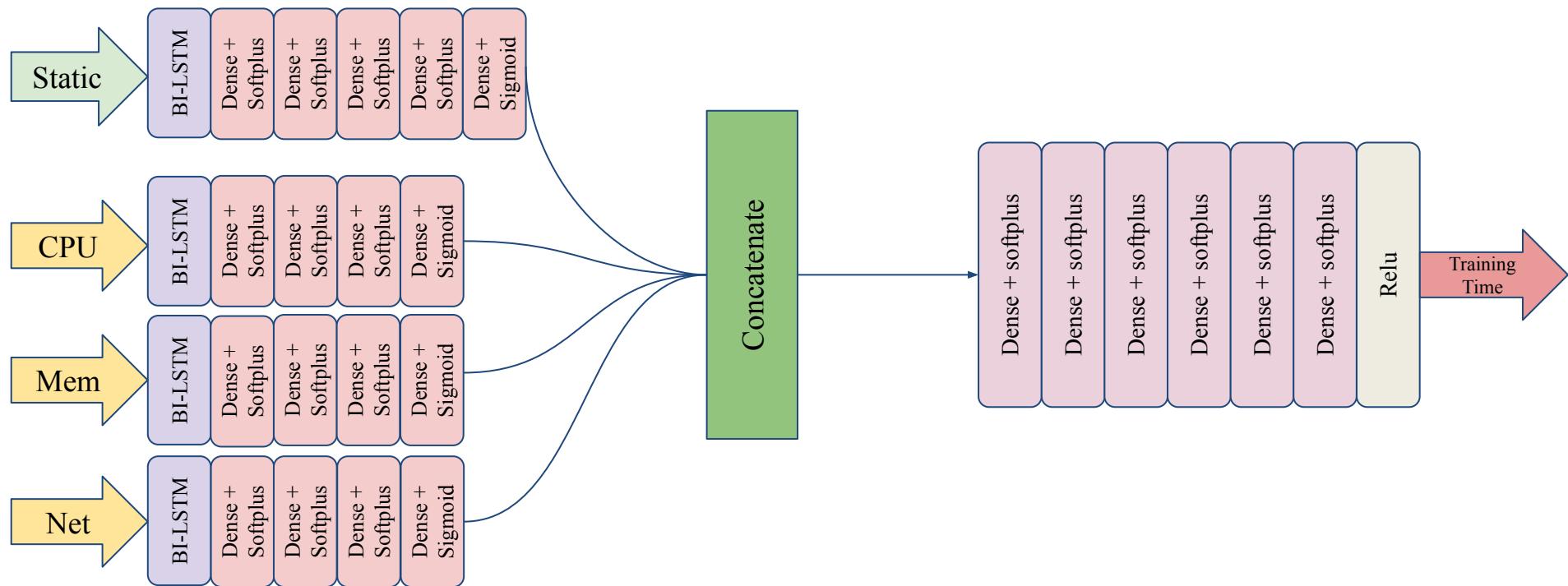
Deep Learning Model Design

Cost Model Development:



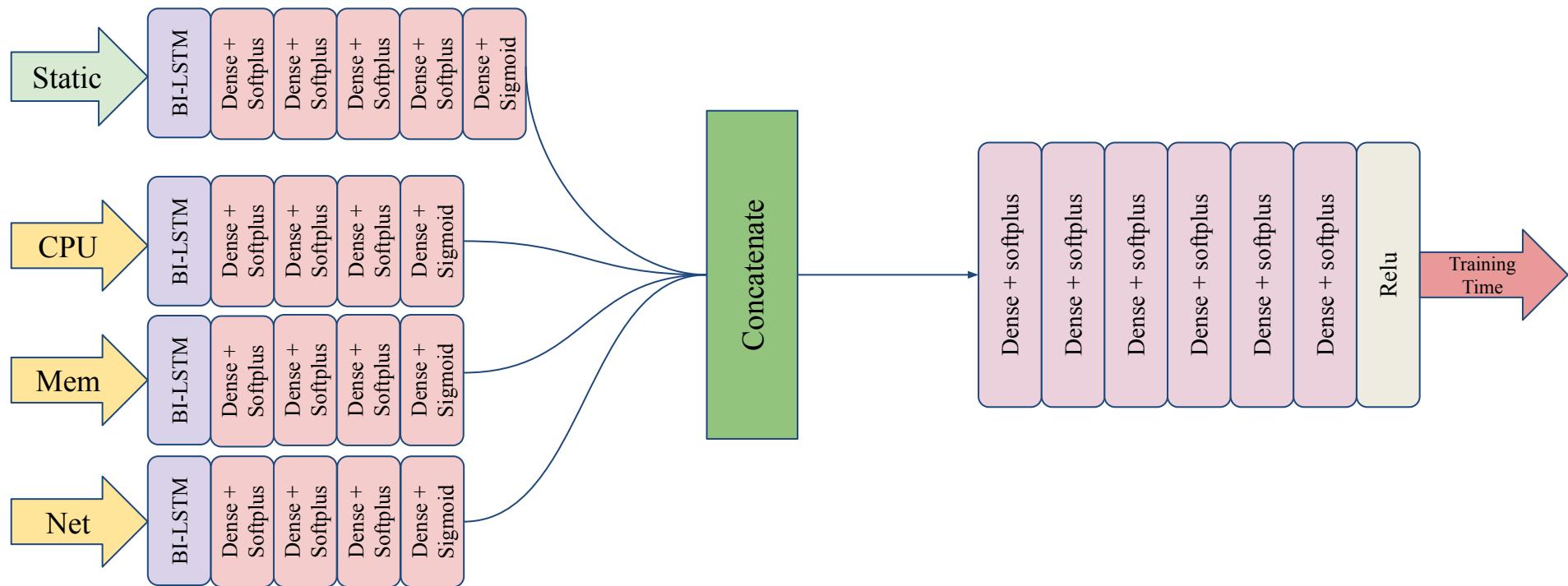
Deep Learning Model Design

Training Time Model Development:



Deep Learning Model Design

Training Time Model Development:



Deep Learning Training Performance (Cost)

- No signs of abnormality in training
- Converge to acceptable RMSE (11% / 13% of the average)
- Achieving R2-score close to ~1 shows a pretty good fit!

Model Hyperparameters:

Loss: Huber (delta=1.0)

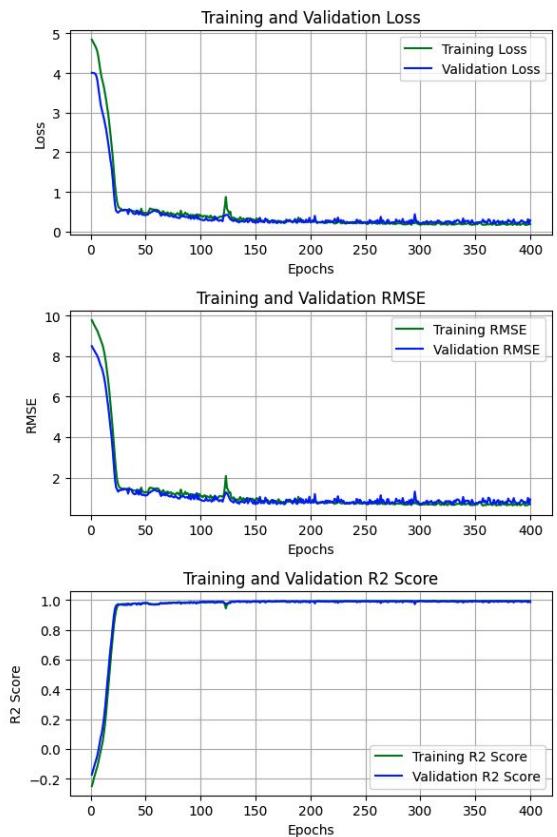
Optimizer: Adam (lr=0.0001, beta=0.5)

Epochs: 400

Batch Size: 32

Train/Validation Split: 80%/20%

	Training	Validation
Loss	0.1175	0.2046
RMSE	0.5602	0.7668
R2 Score	0.9959	0.9904



Deep Learning Training Performance (Training Time)

- No signs of abnormality in training
- Converge to acceptable RMSE (16% / 13% of the average)
- Achieving R2-score close to ~1 shows a pretty good fit!

Model Hyperparameters:

Loss: Poisson

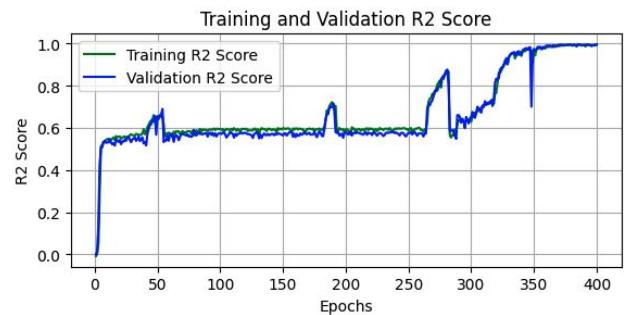
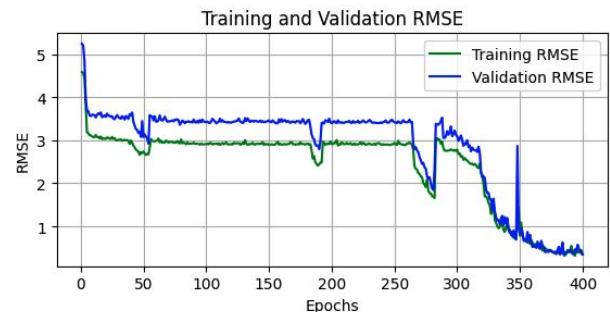
Optimizer: Adam (lr=0.0008, beta=0.5)

Epochs: 400

Batch Size: 32

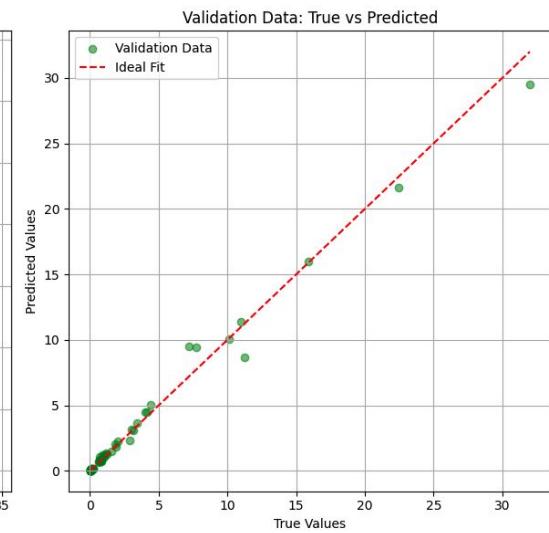
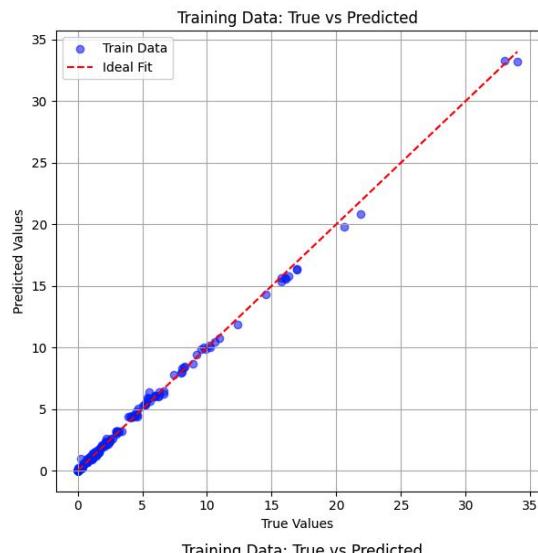
Train/Validation Split: 80%/20%

	Training	Validation
Loss	-2.2585	-2.6542
RMSE	0.4095	0.3273
R2 Score	0.9920	0.9961

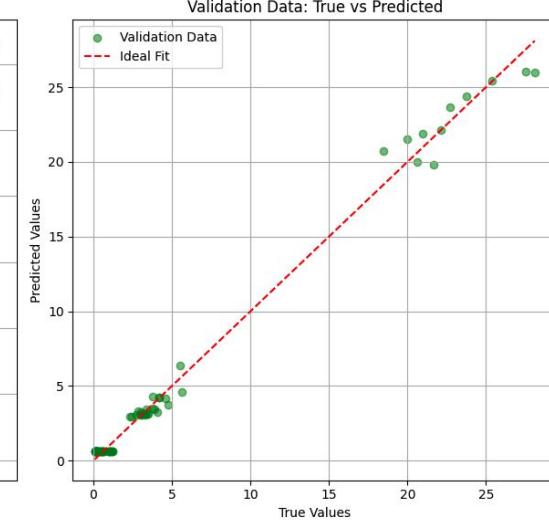
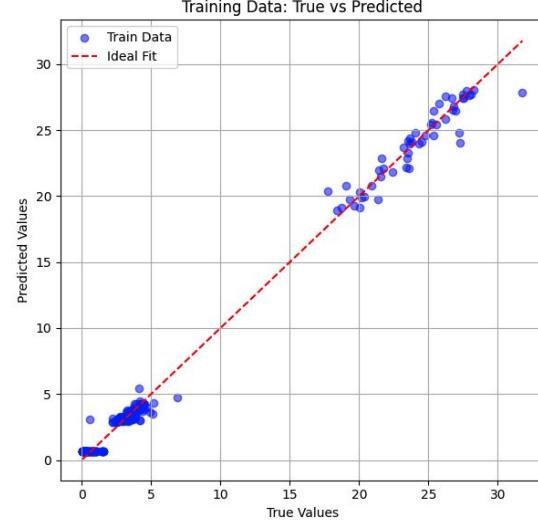


Confusion Plots vs Perfect Fit

Training Time
confusion plot

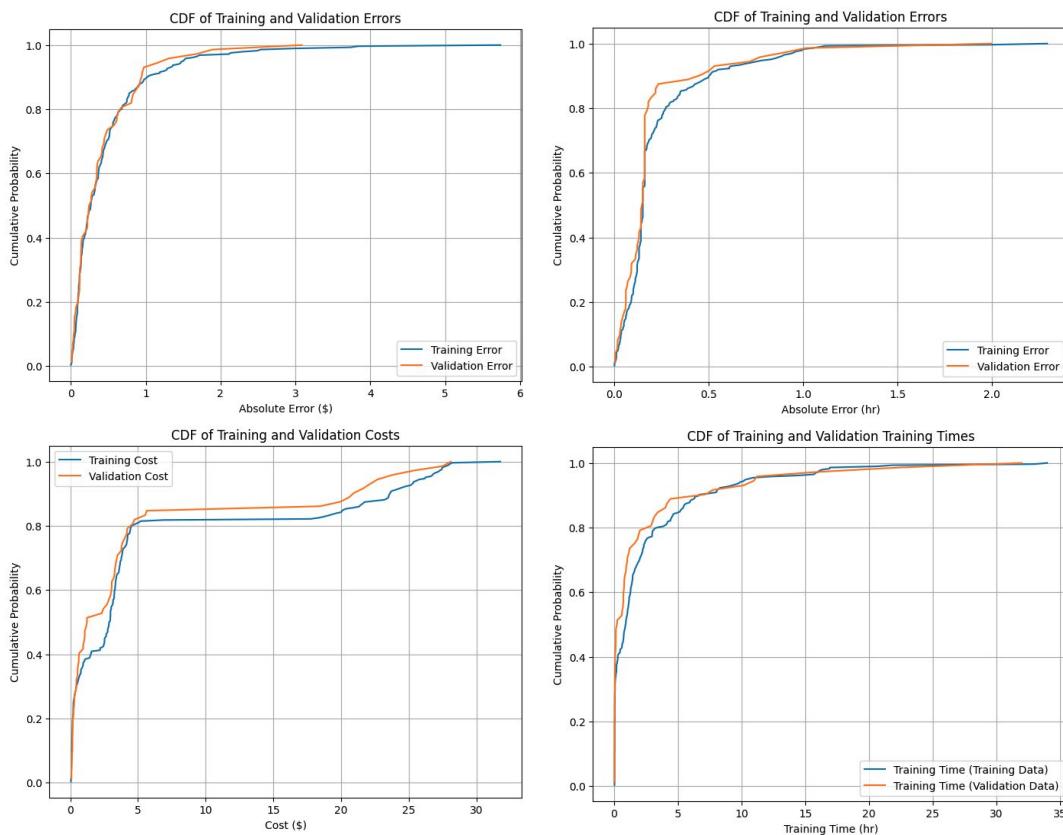


Cost confusion plot



Error Distribution

- Cost Value 90 %ile: 23.6\$ and 20.9\$,
- Prediction Error 90 %ile: 0.78\$ and 0.93\$
- Training time 90 %ile: 6.65 and 6.89 hr
- Prediction Error 90 %ile: 1.5 hr



Deep Learning vs Traditional ML Algorithms

ML Algorithms:

- Linear Regression/Poisson Regression for Cost/Training Time prediction 
- Not showing promising results compared to DL
- Random Forest Regression
- XGBoost

Feature Engineering:

- Random Forest Top 5 features
 - **Cost Model:**
 - Partial resource usage
 - An increase in resource usage during the later stages of the time series correlates with a higher cost.
 - **Training Time Model:**
 - Increase in the # of workers/machines correlates with decrease in training time
 - More CPU correlates with less training times
 - n1-standard-4 instance has a wider range in predictions



Machine Learning Models Comparison (refer to thesis)

Random Forest Regressor and XGBoost show comparable performance to deep learning

Cost Prediction

	Training	Validation
<i>RMSE</i>	1.43	1.34
<i>R2 Score</i>	0.97	0.97

Training Time Prediction

	Training	Validation
<i>RMSE</i>	0.28	0.31
<i>R2 Score</i>	0.96	0.96

Linear/Poisson Regression

Random Forest Regressor

XGBoost

Training

Validation

	Training	Validation
<i>RMSE</i>	0.33	0.70
<i>R2 Score</i>	0.99	0.99

Training

Validation

	Training	Validation
<i>RMSE</i>	0.28	0.34
<i>R2 Score</i>	0.996	0.995

Training

Validation

	Training	Validation
<i>RMSE</i>	0.19	0.70
<i>R2 Score</i>	0.999	0.99

Training

Validation

	Training	Validation
<i>RMSE</i>	0.08	0.26
<i>R2 Score</i>	0.999	0.997



Deep Learning vs Traditional ML Algorithms

Deep Learning

- Dataset Size and Complexity
- Automated Feature Engineering
- Transfer Learning
- Sequential Data
- Computationally intensive
- Generalize well!

ML Algorithms:

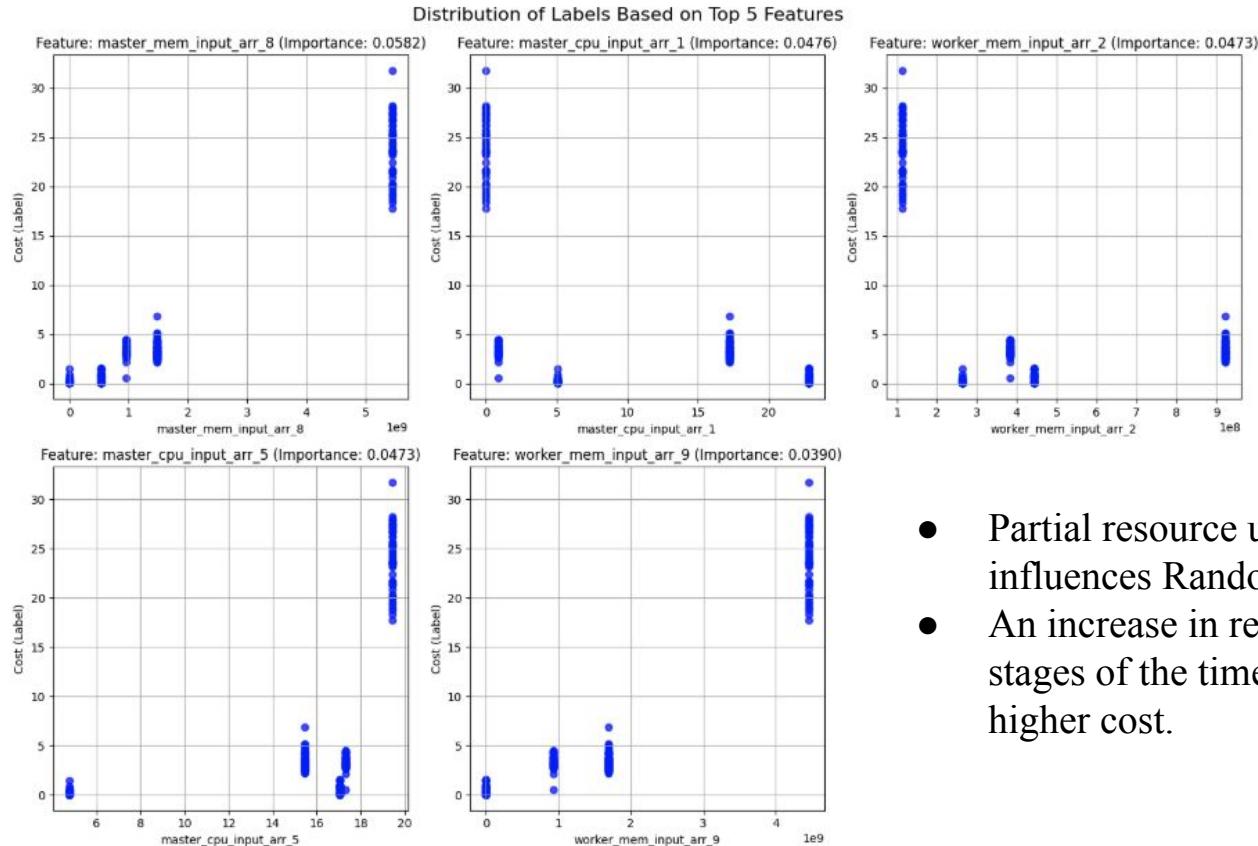
- Less computationally intensive
- Cannot naturally handle sequential data
- Need manual feature engineering
- Lose to DL when it comes to large scale tasks and complex datasets
- ALgorithms like XGBoost usually appear to overfit the training data

Hybrid Methods:

- Model Stacking
- Ensemble Methods



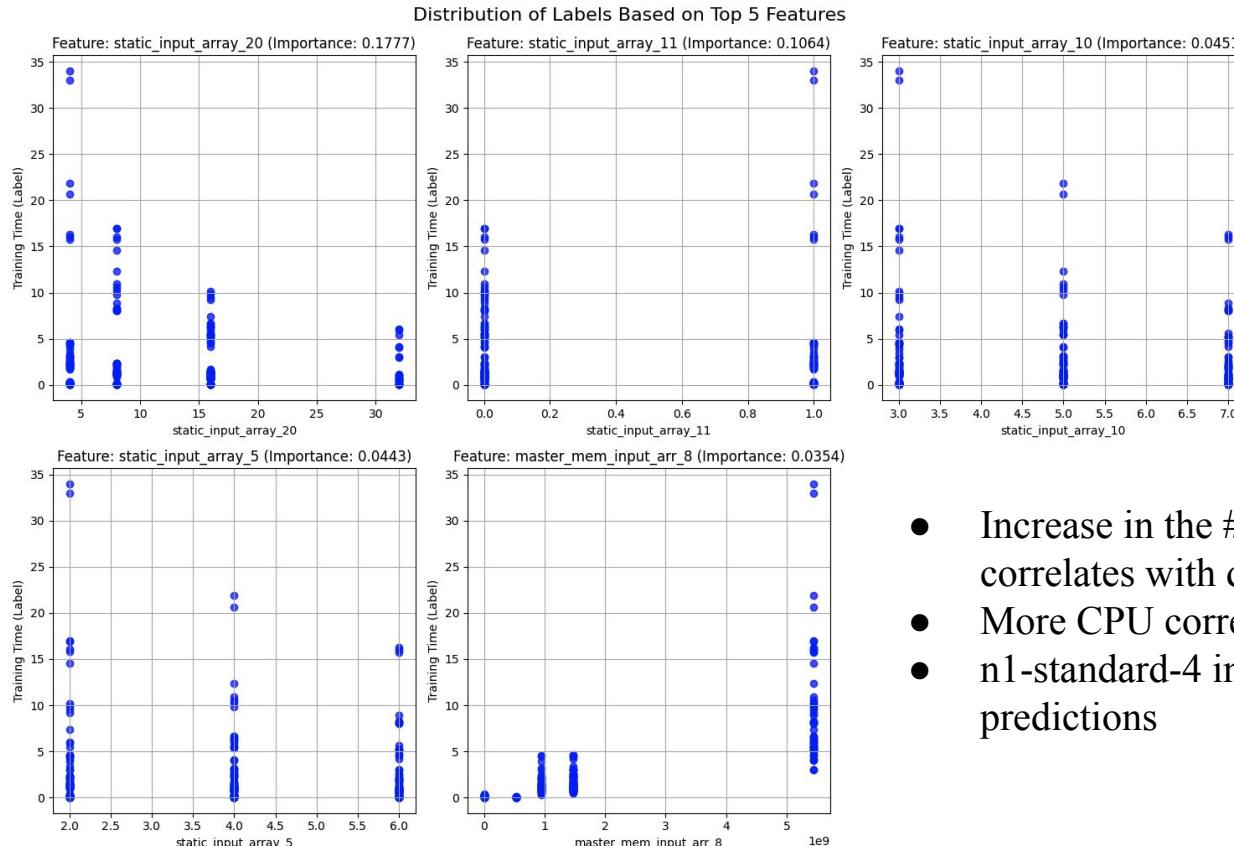
RF model feature importance analysis



- Partial resource usage significantly influences Random Forest predictions.
- An increase in resource usage during the later stages of the time series correlates with a higher cost.



RF model feature importance analysis



Static_input_arr_11: n1-standard-4 instance type
Static_input_arr_10: # of Machines
Static_input_arr_20: Available CPU
Static_input_arr_5: # of workers

- Increase in the # of workers/machines correlates with decrease in training time
- More CPU correlates with less training times
- n1-standard-4 instance has a wider range in predictions



Conclusion

In this thesis:

- We studied the advantages of fine-grained resource allocation in containerized environments
- And proposed a fast resilient network architecture for state-of-the-art distributed deep learning workload
- Finally, we investigated distributed deep learning cost and training time predictions in the cloud to help DDL developers make informed decisions on choosing their cloud configurations
- Minimizing waste:
 - In overprovisioning through efficient auto-scaling
 - of entire clusters for extended periods due to isolated network edge failures
 - Resources in instance selection, preventing overuse of large instances when smaller, cost-effective options can suffice



Thank You!



University of Colorado **Boulder**

References

- [1] [Google Vertex AI](#)
- [2] [VPA](#)
- [3] [Autopilot](#)
- [4] [Fast DDL over RDMA](#)
- [5] [Ultima](#)
- [6] [Maestro](#)
- [7] [tc-flower](#)
- [8] [CherryPick](#)
- [9] [Ernest](#)
- [10] [Llama3](#)
- [11] [MegaScale](#)
- [12] [MPRDMA](#)

