

Target-based Resource Allocation for Deep Learning Applications in a Multi-tenancy System

Wenjia Zheng¹, Yun Song¹, Zihao Guo¹, Yongchen Cui¹, Suwen Gu¹, Ying Mao¹,
Long Cheng²

Fordham University¹ - University College Dublin²



Outline

- Introduction & Background
- Motivation
- *TRADL* Modules
- System Evaluation
- Conclusion

Deep Learning is Reshaping Our World



Self-Driving Cars



Face Recognition

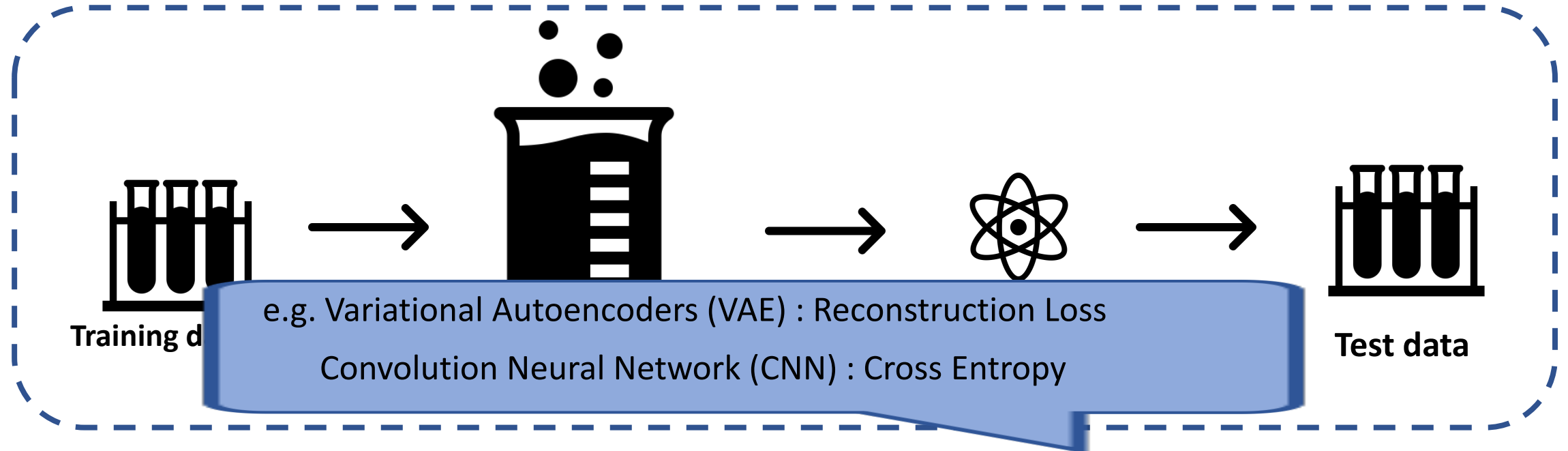


Voice Controlled Assistance



Automatic Language Translation

Training process of Deep Learning



Evaluation: Loss function (predefined)



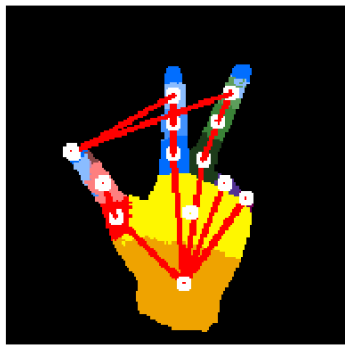
Target: Global minimum

Error-tolerated DL application

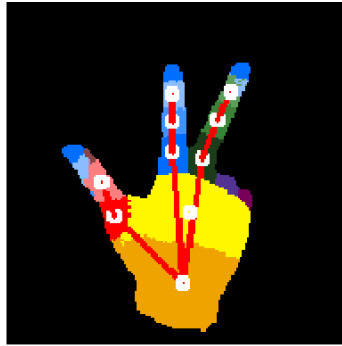
Hand pose estimation



(a)



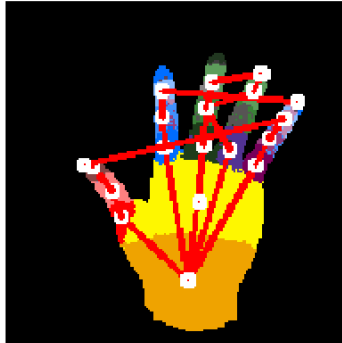
(b)



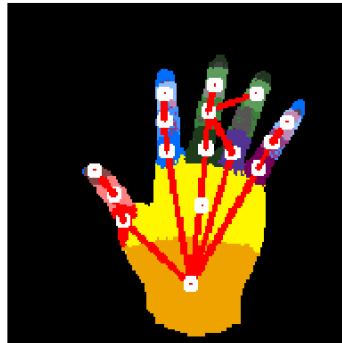
(c)



(d)



(e)



(f)

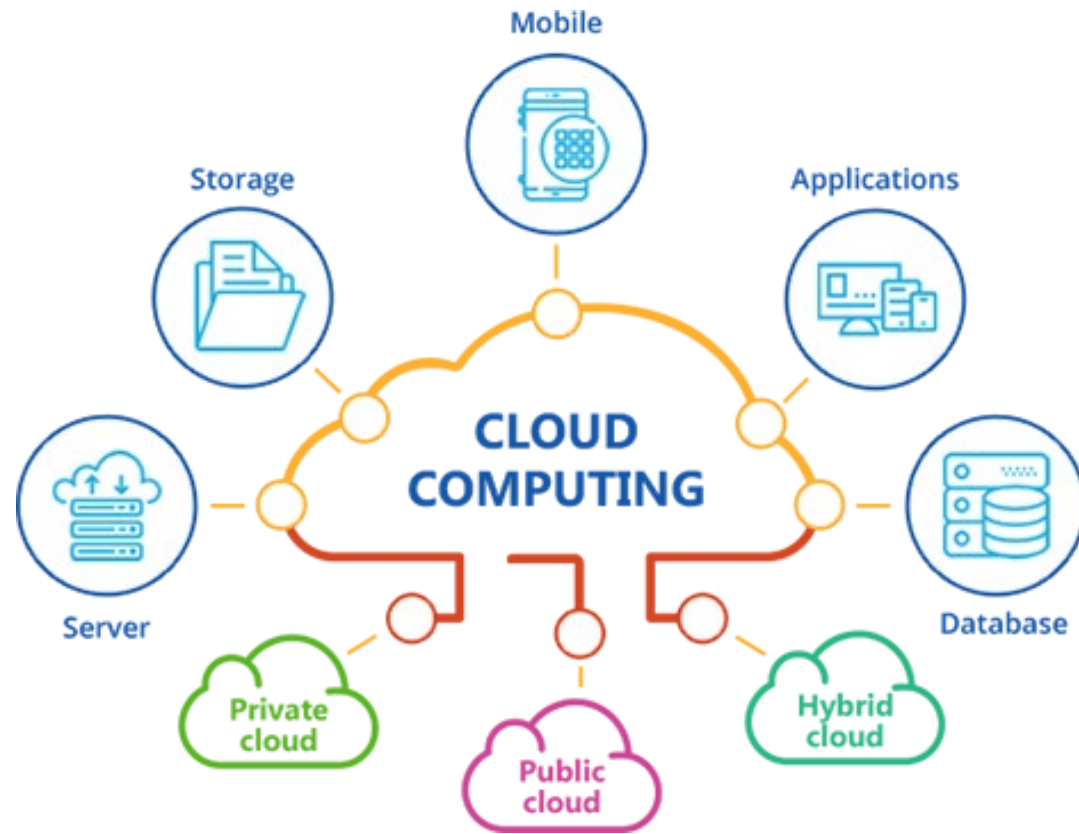
It need to achieve 100% accuracy



assumption:

users can tolerate a certain level of error rate

Limited Resources on the Cloud



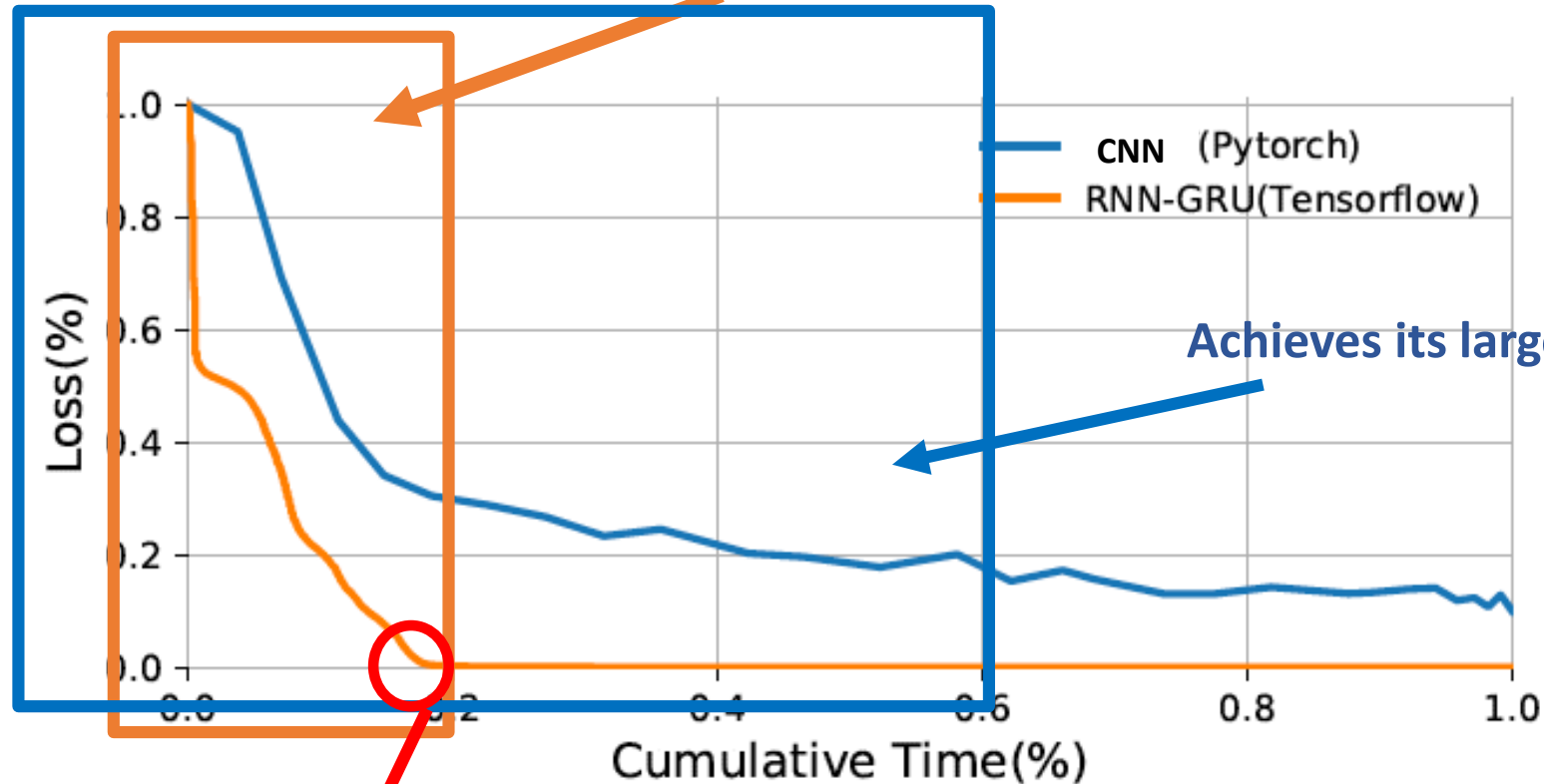


Outline

- Introduction & Background
- Motivation
- *TRADL* Modules
- System Evaluation
- Conclusion

Motivation

Achieves its largest reduction at 20% of the time

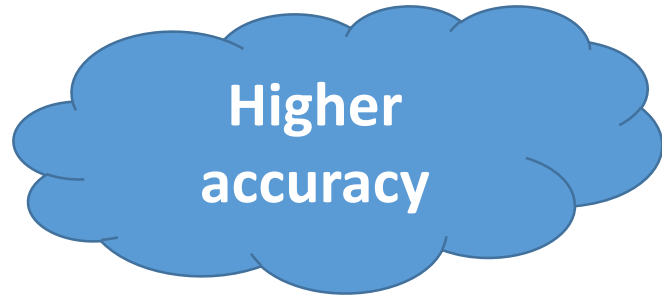


Achieves its largest reduction at 60% of the time



Allocate less resource & shift resource to other learning tasks

Trade-off relationship



client

Assumption:
Users can tolerate a certain
level of error rate

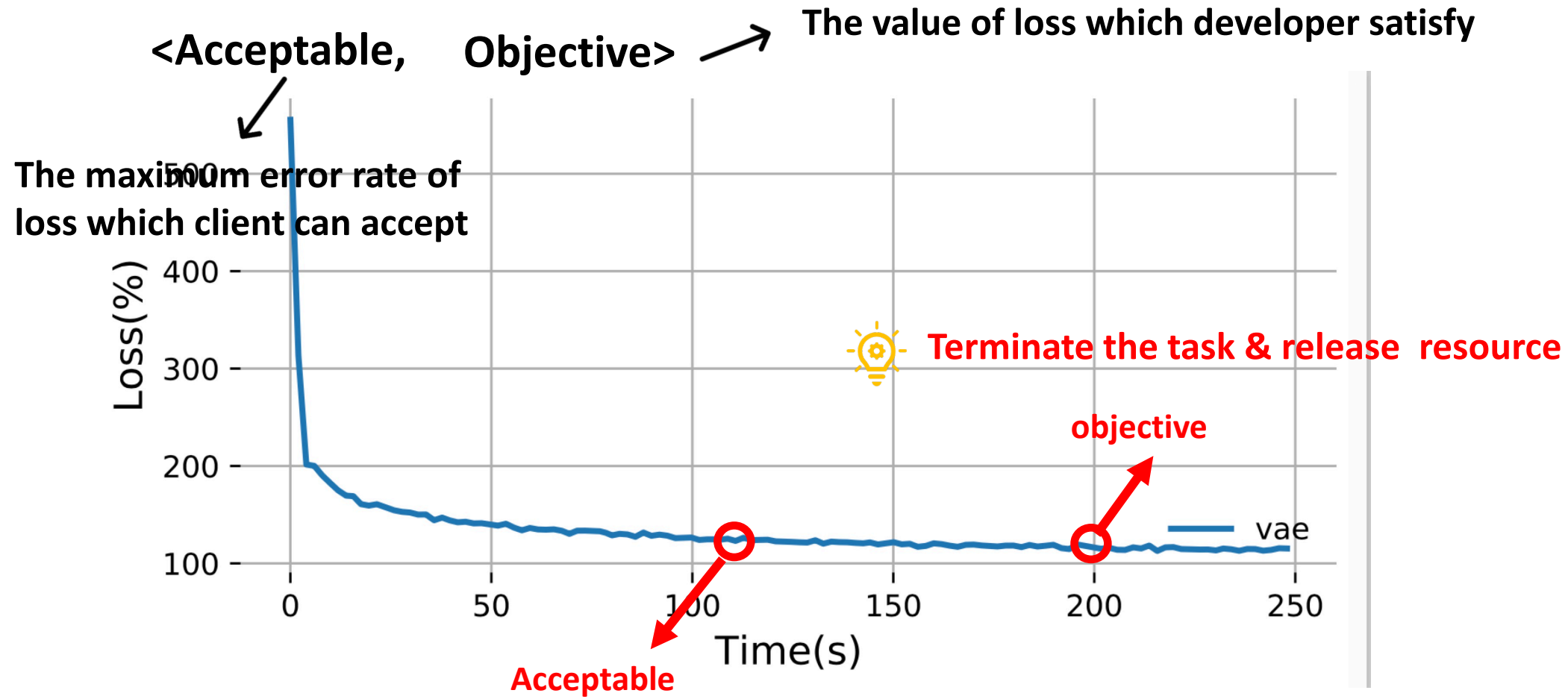


developer



Two-tier target

Two-tier target



Allocate less resource & shift resource to other learning tasks

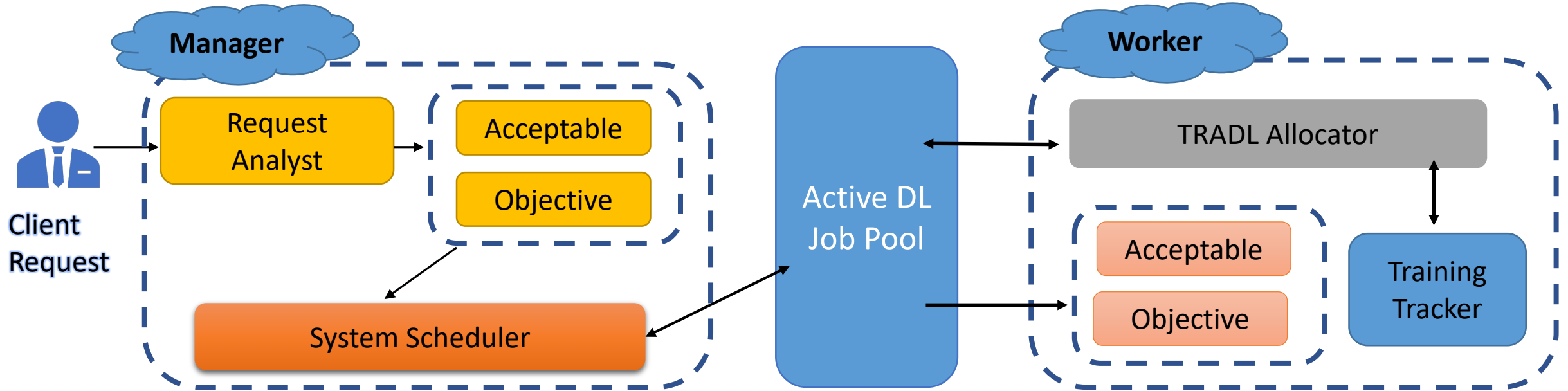


Outline

- Introduction & Background
- Motivation
- *TRADL* Modules
- System Evaluation
- Conclusion

TRADL Modules

「Interact with clients and assign tasks to the workers」
「Execute the tasks and generate the results」



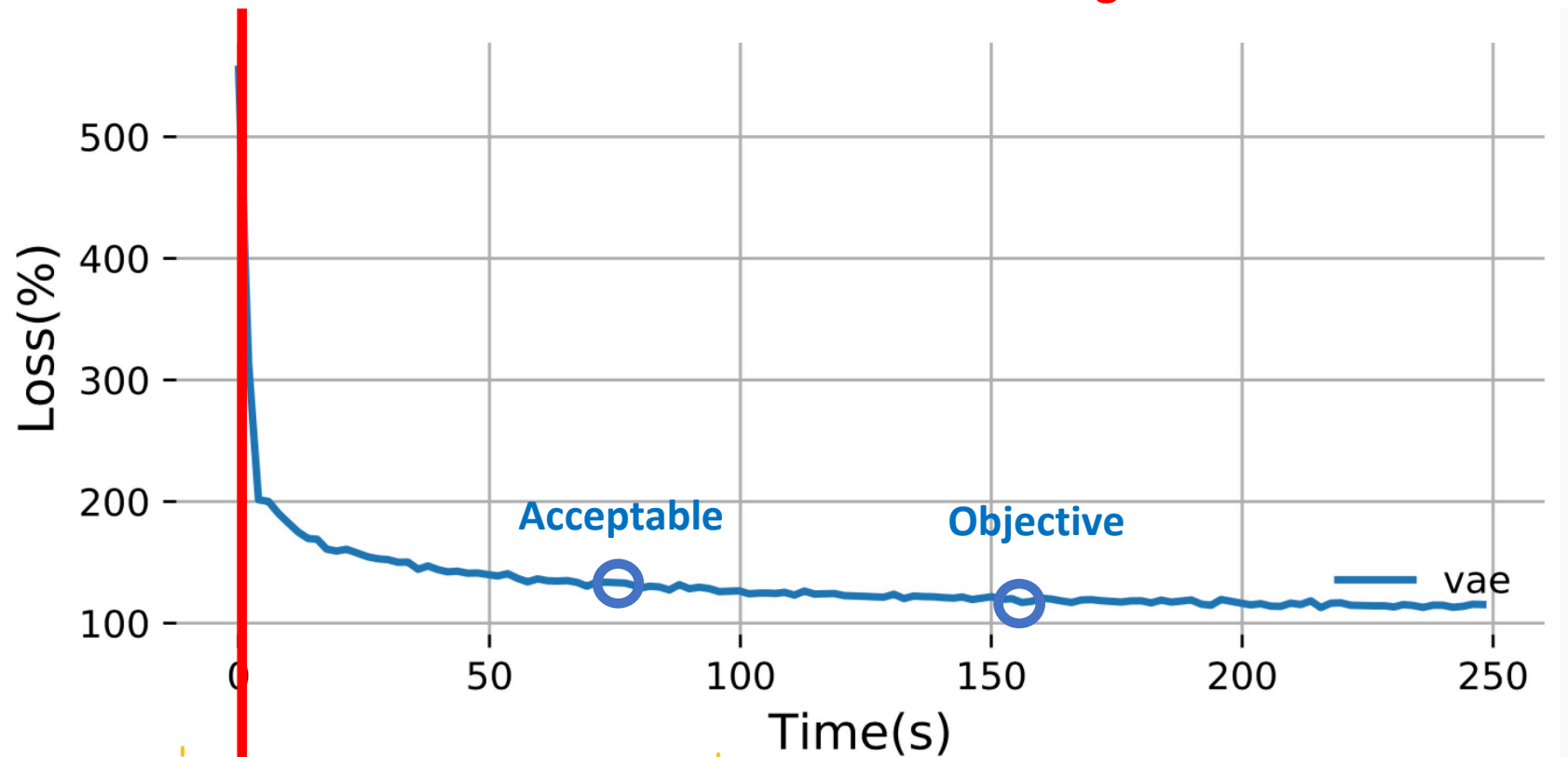
「Control the resource allocation globally & label the training job with Acceptable and Objective」

Resource Assignment Algorithm

One job

1. Resources consumed by other processes

→ Using resource $< \alpha * \text{Total Resource}$ ↗ Free compete



Allocate less resource & terminate the task & release resource

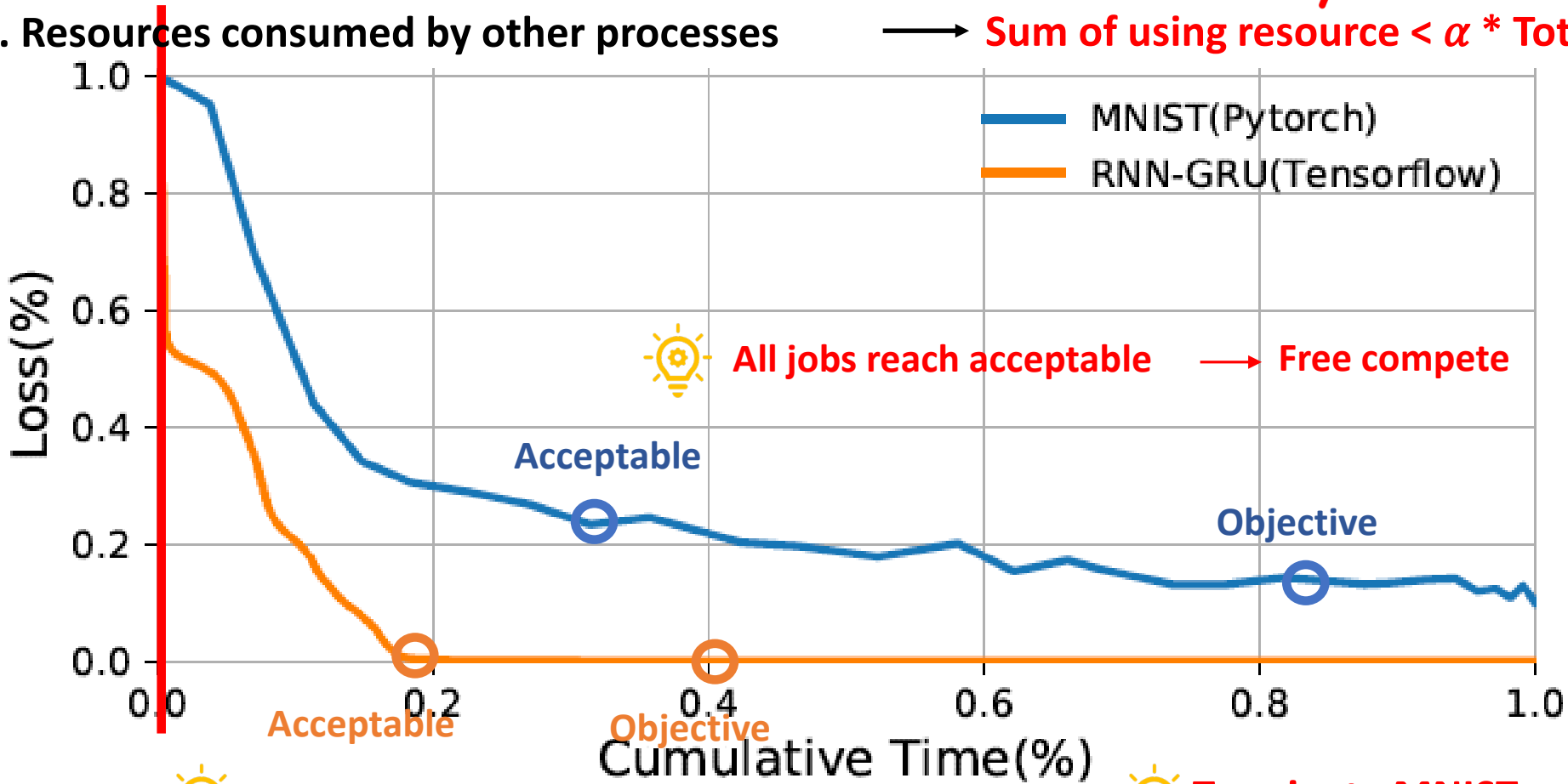


Resource Assignment Algorithm

Multi jobs

1. Resources consumed by other processes

Sum of using resource < α * Total Resource



Allocate less resource to RNN-GRU
Terminate RNN-GRU
Terminate MNIST



Outline

- Introduction & Background
- Motivation
- *TRADL* Modules
- System Evaluation
- Conclusion



System Evaluation

- **Workload**

Model	Loss Function	Two-tier Target
Variational Autoencoders (VAE)	Recon. Loss	120, 110
Convolution Neural Network (CNN)	Cross Entropy	0.3, 0.25
Gated Recurrent Unit (GRU)	Quadratic Loss	0.01, 0.008
Long Short-Term Memory (LSTM-CFC)	Softmax	0.02, 0.015
Long Short-Term Memory (LSTM-CRF)	Squared Loss	10.5, 10.2
Bidirectional-RNN	Softmax	0.6, 0.5
Recurrent Network (RNN)	Softmax	0.55, 0.5
Dynamic RNN	Softmax	0.15, 0.1

System Evaluation

- **Evaluation scenarios:**

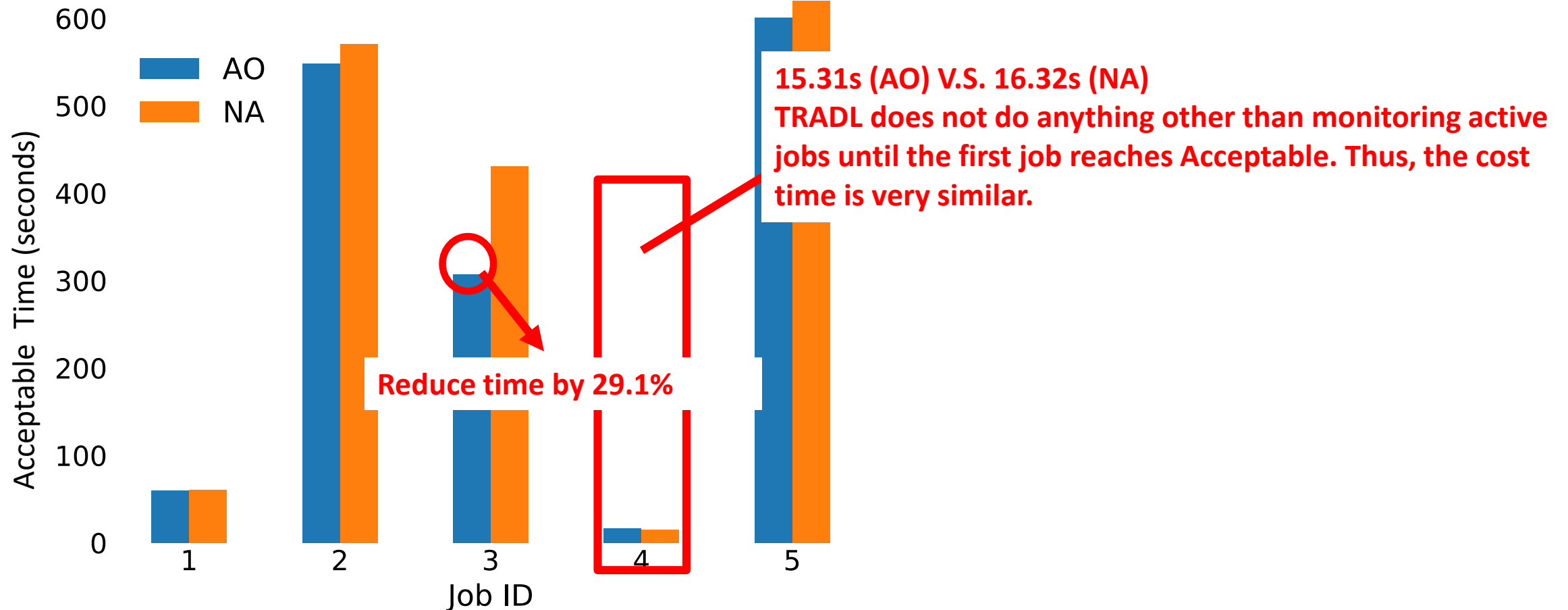
1. **Acceptable Only:** $\langle \text{Acceptable}, _ \rangle$ (denoted as **AO**)
2. **Full two-tier Target:** $\langle \text{Acceptable}, \text{Objective} \rangle$ (denoted as **FT**)
3. **No Algorithm:** $\langle _, _ \rangle$ (denoted as **NA**)

- **Evaluation metrics:**

1. **Acceptable and Objective time:** the time cost that each job achieve the two-tier target.
2. **CPU Usage:** the CPU usage for the deep learning applications

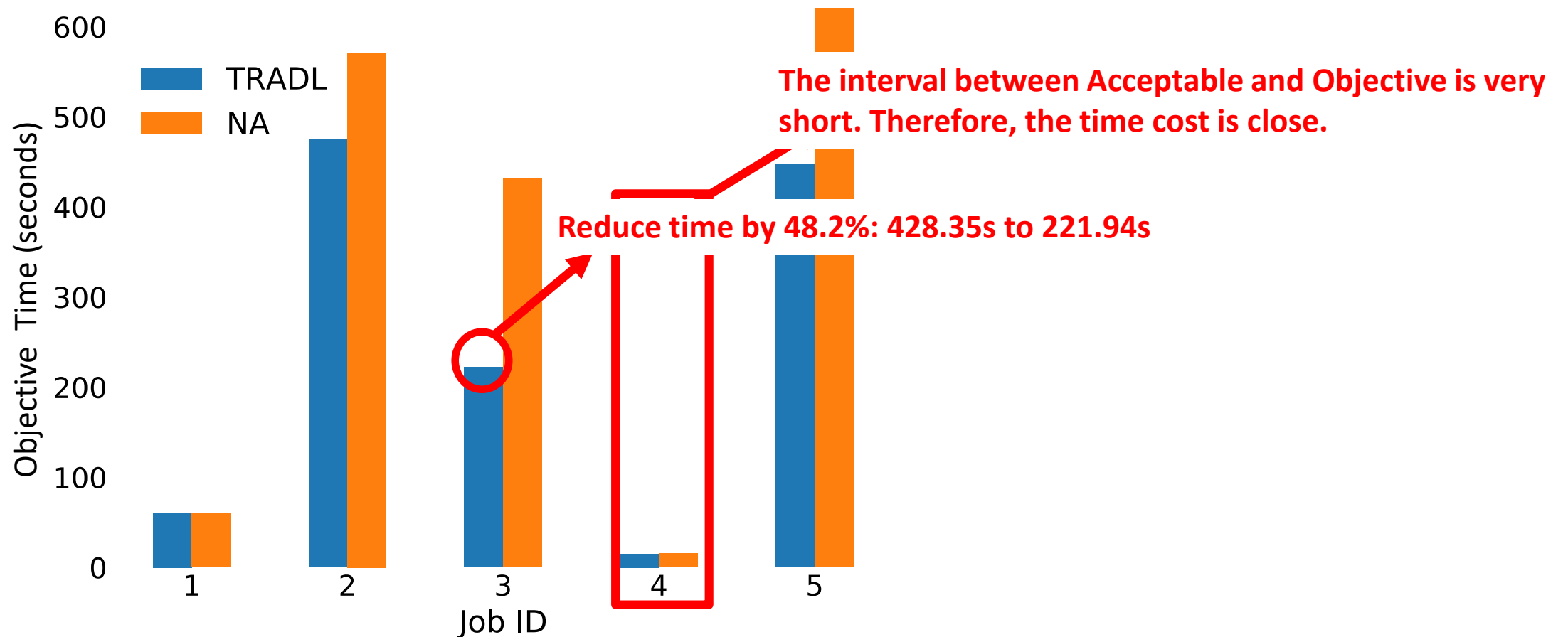
Acceptable and Objective time

- Acceptable Only

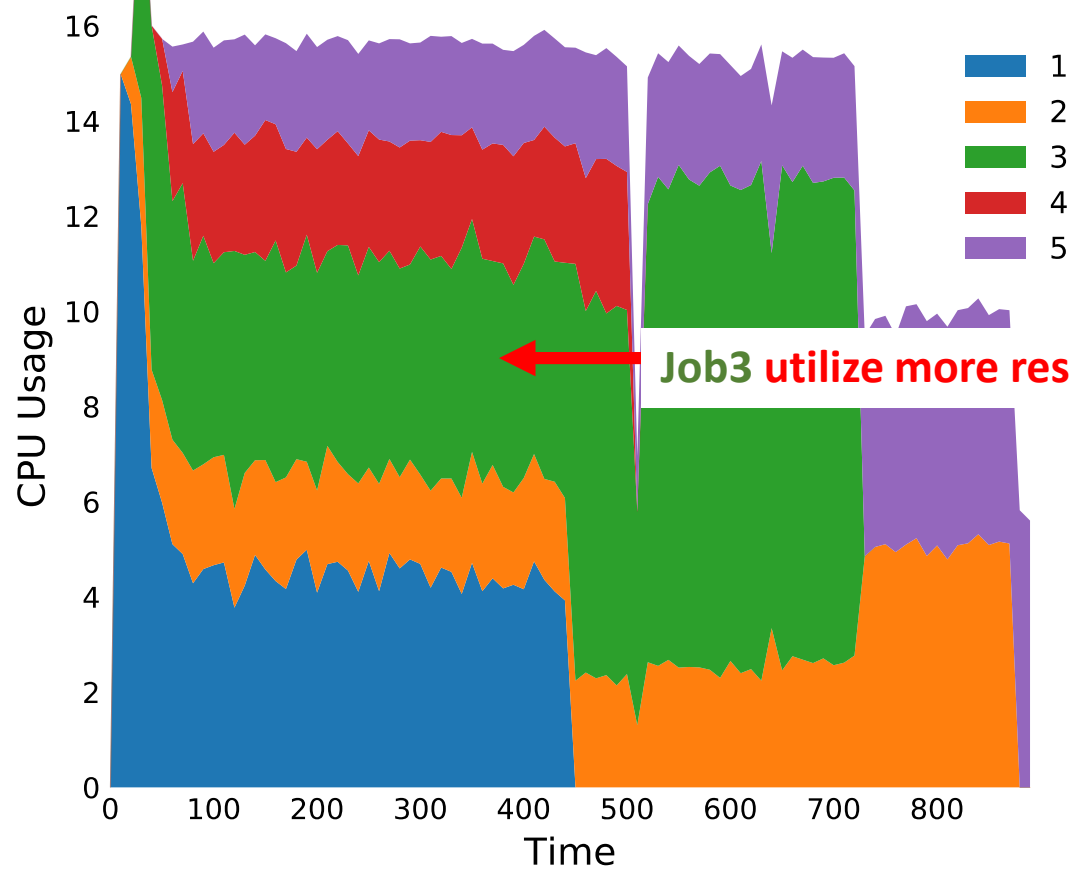


Acceptable and Objective time

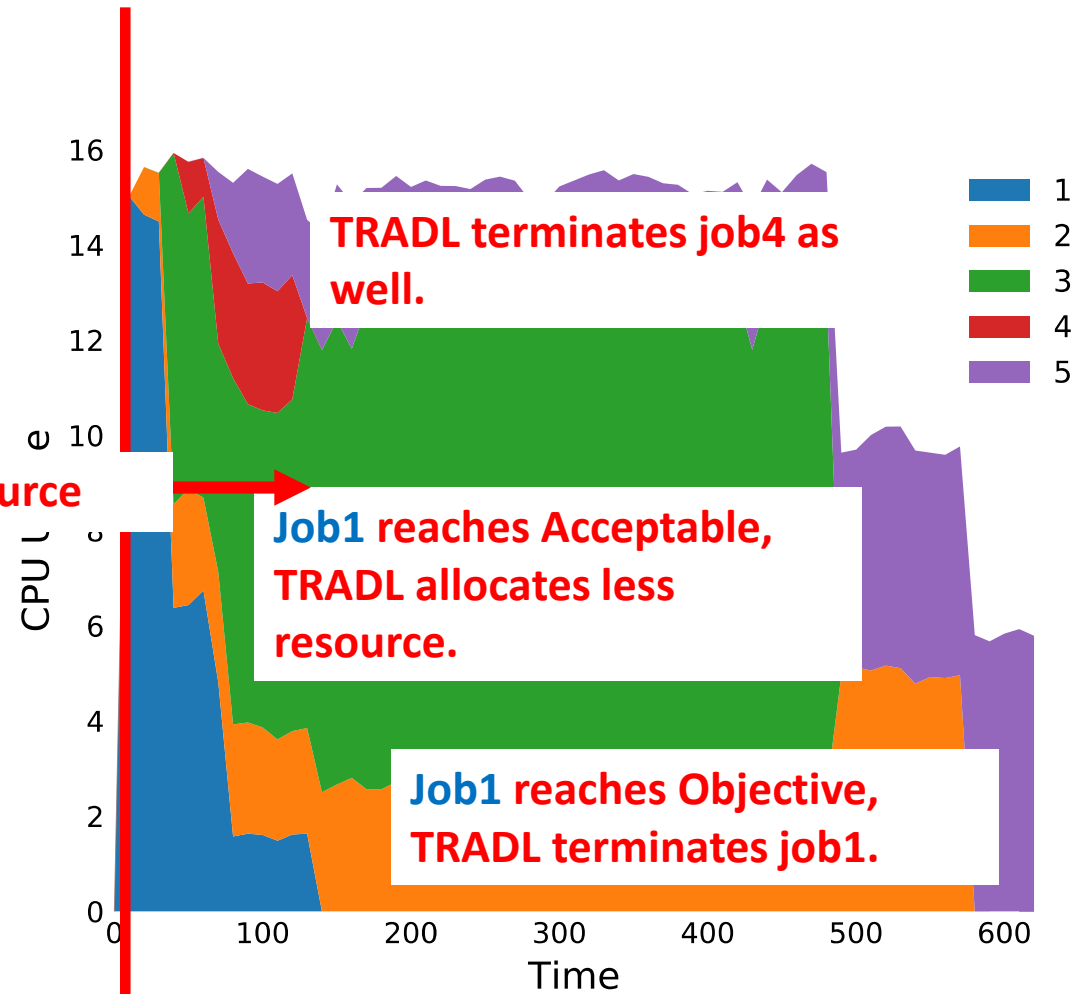
- Full two-tier Target



CPU Usage



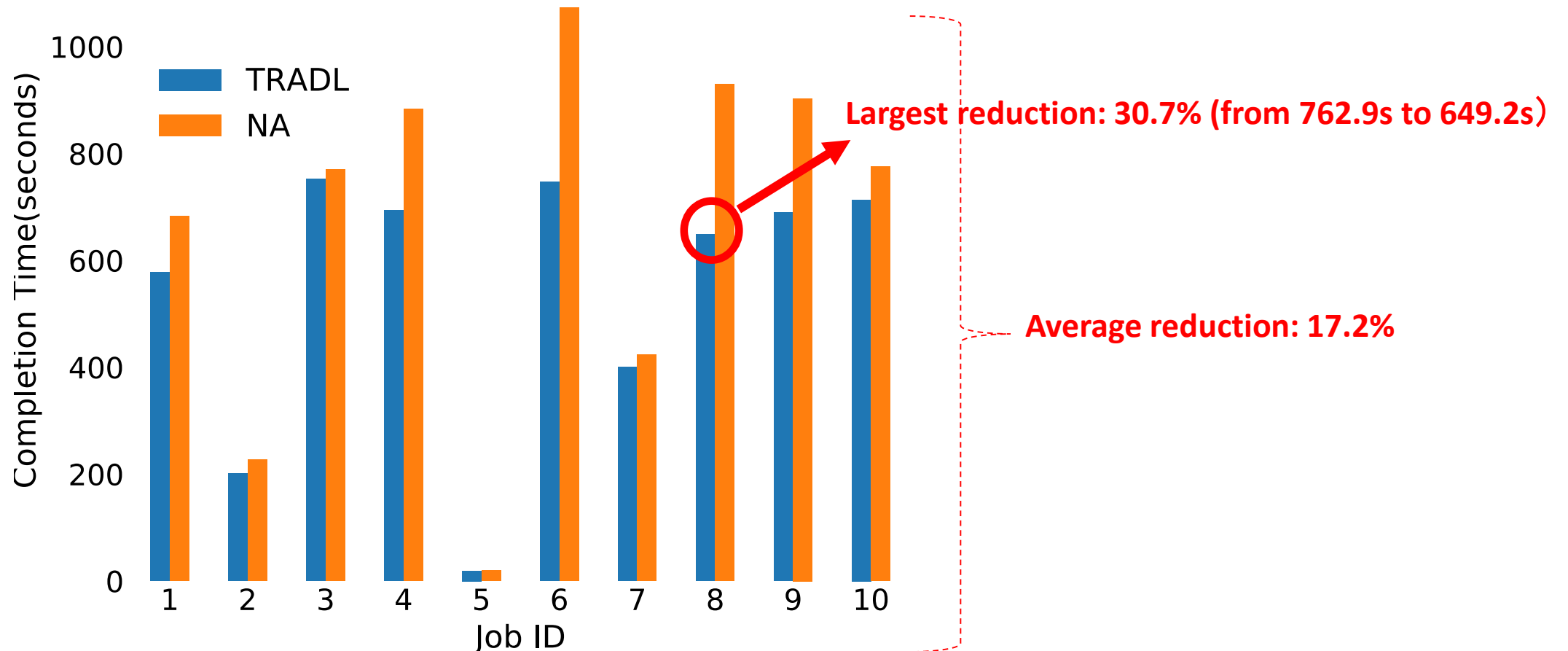
NA



TRADL

Scalability

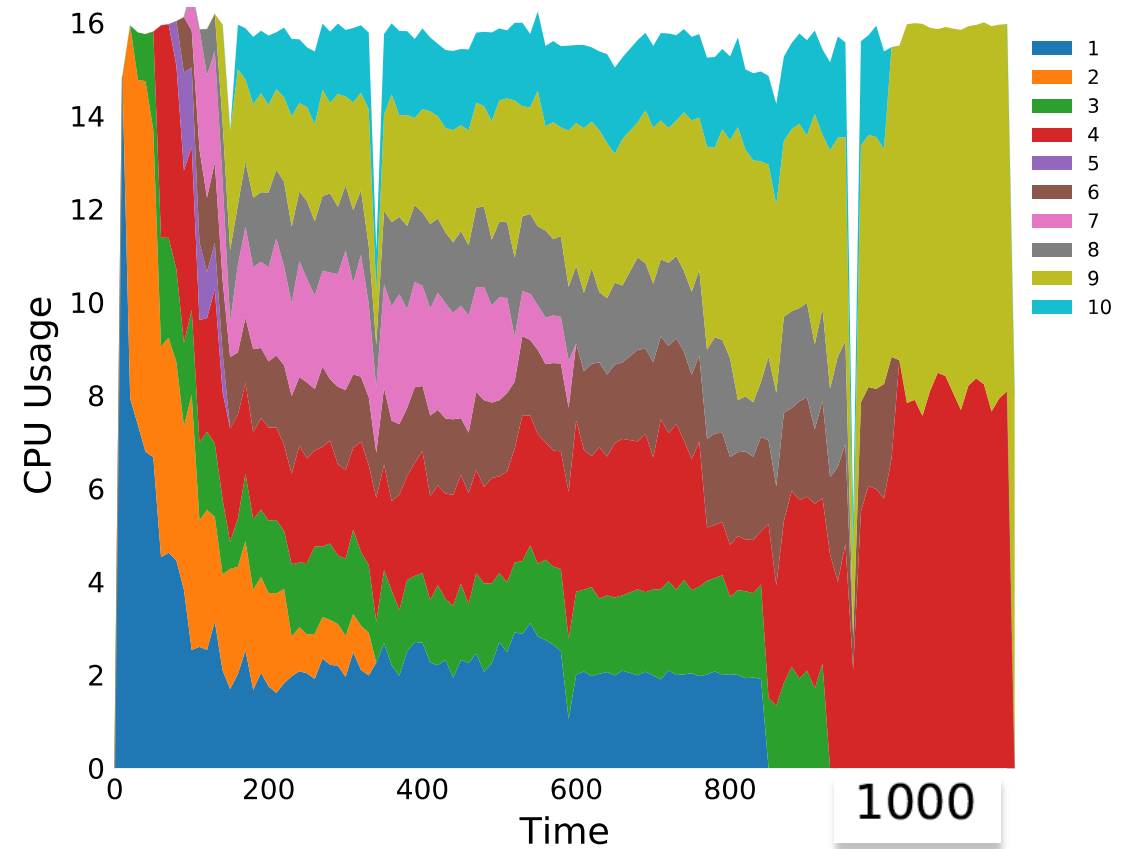
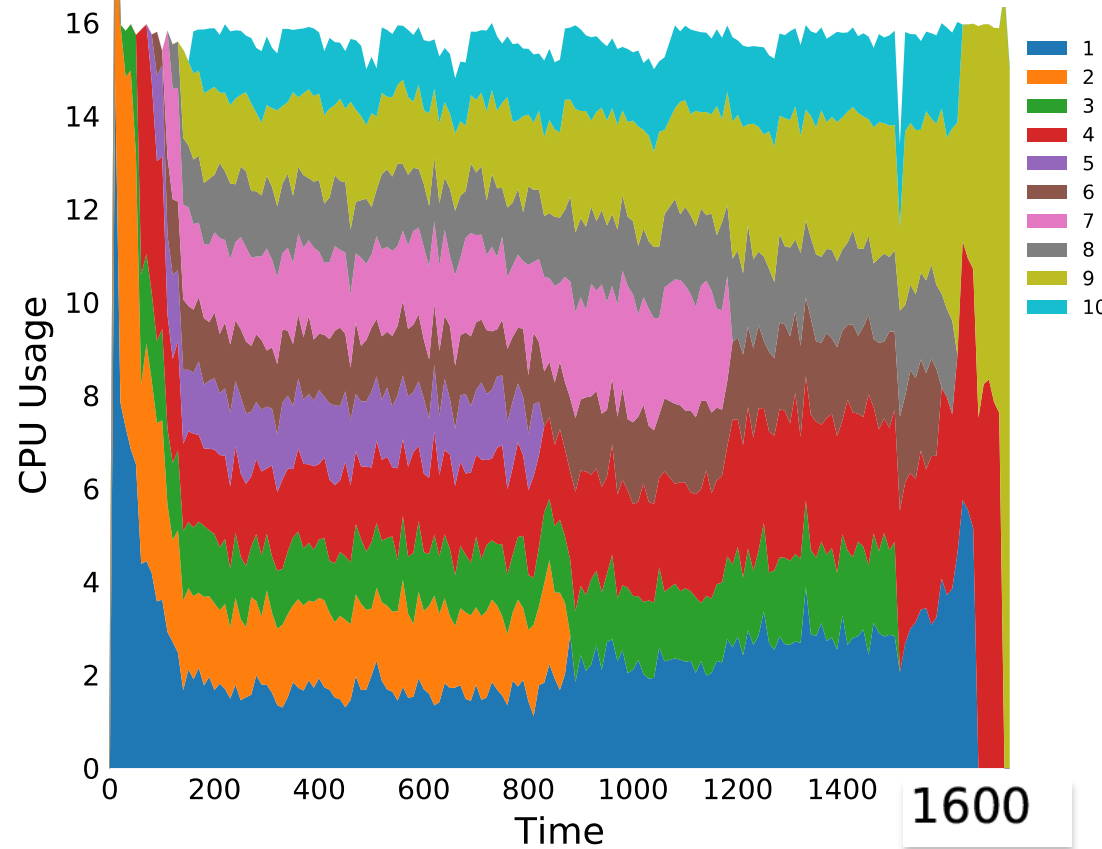
- Acceptable and Objective time



Scalability

- CPU Usage

Job4 and **Job9** utilize more resource than NA.



The whole process is accelerated.

Outline

- Introduction & Background
- Motivation
- *TRADL* Modules
- Evaluation
- **Conclusion**

Conclusion

- Based on the assumption that users can tolerate a certain level of error rate, we develop *TRADL*, which dynamically updates the resources configurations while the jobs are running to accelerate the training process.
- *TRADL* shows a significant improvement on the time cost of achieving the targets, and the largest reduction is 48.2%.



Thank you!
Q&A

if $RU(c_i) > \frac{1}{n}$ **then**
 $c_i.set_limits(\frac{1}{n \times 2})$
else if $RU(c_i) \leq \frac{1}{n}$ **then**
 $c_i.set_limits(\frac{1}{n+1})$