



# Understanding and Optimizing GPU Energy Consumption of DNN Training

Jae-Won Chung

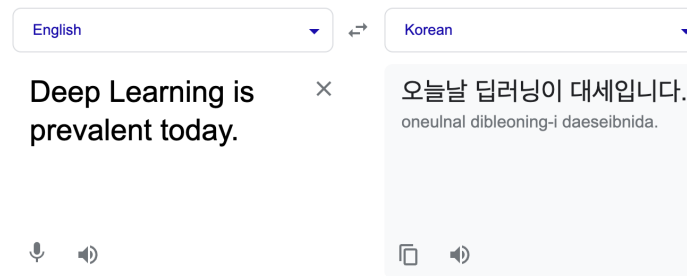
April 17<sup>th</sup>, 2023

*Work done in collaboration with Jie You and Mosharaf Chowdhury*



# Deep Learning is Prevalent Today

Image processing  
Speech recognition  
Machine translation  
Intelligent assistants  
Autonomous driving  
Video analytics  
Image/text generation

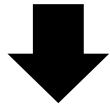


stability.ai

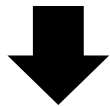


# DNN Energy Consumption is Skyrocketing

DNN



GPU



Energy

- Re-training is commonplace (e.g. every hour)<sup>2</sup>
- Training GPT-3 == 120 years of electricity for a household<sup>1</sup>
- Performance optimizations oblivious of energy impact

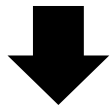
1. U.S. EIA and Google (arXiv '21) 2. Facebook (HPCA '18) and Alibaba (NSDI '22)

# Existing Efforts are not Practical Enough

DNN



GPU



Energy

- New energy-efficient DNN architectures  
SqueezeNext (CVPRW '18), ChamNet (CVPR '19), SkyNet (MLSys '20)
- New energy-efficient HW architectures  
TPU (ISCA '17), EDEN (MICRO '19), LNPU (ISSCC '19)
- Offline profiling and power model fitting
- Confined to GPU power configuration knobs  
MPC (HPCA '17), ODPP (CCGRID '20), GPOEO (TPDS '22)

# Understanding GPU Energy Consumption

## *Energy to Accuracy (ETA)*

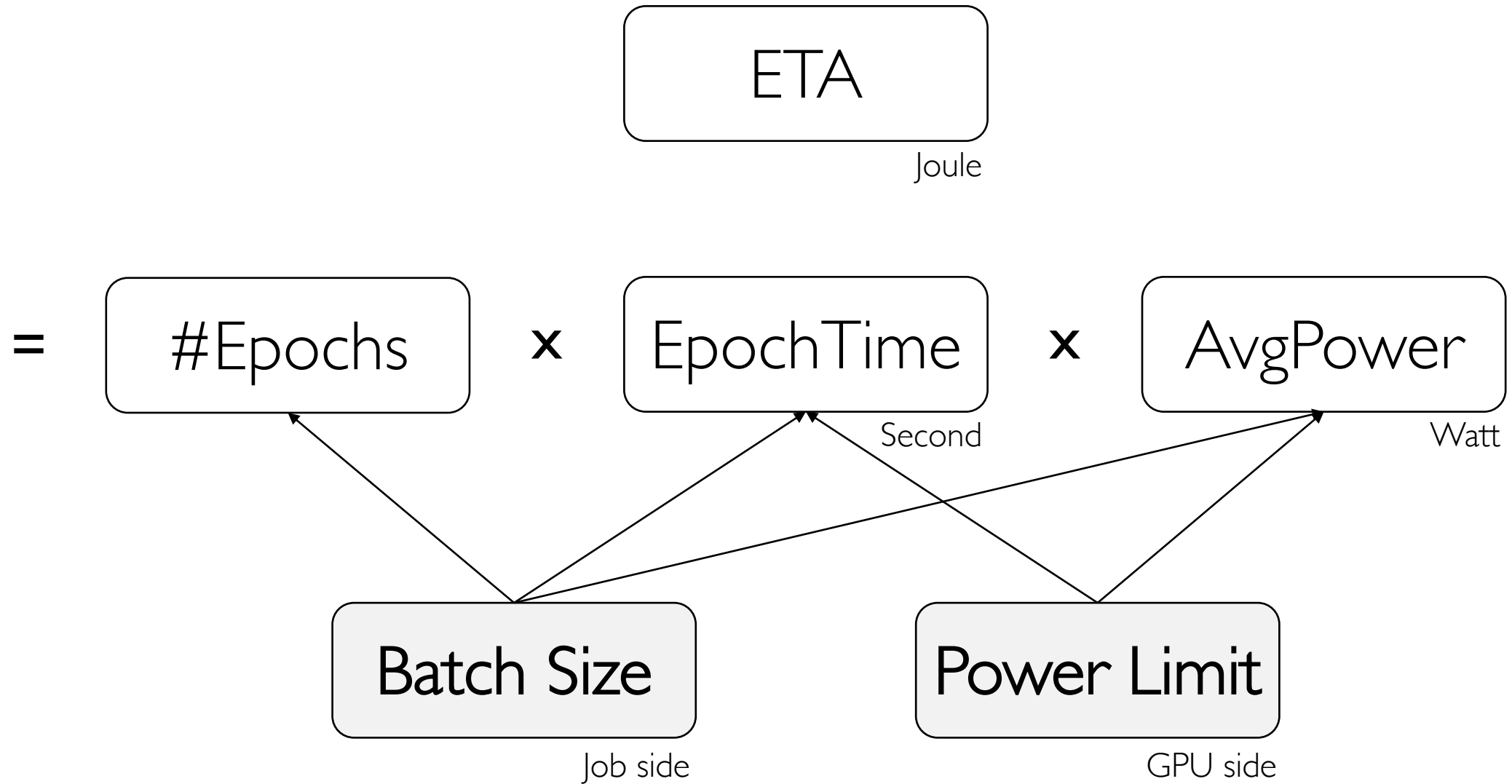
- Energy needed to reach the user-specified **target accuracy**
- Energy-counterpart of *Time to Accuracy (TTA)*

# Understanding GPU Energy Consumption

$$\text{ETA} \text{ (Joule)} = \text{TTA} \text{ (Second)} \times \text{AvgPower} \text{ (Watt)}$$

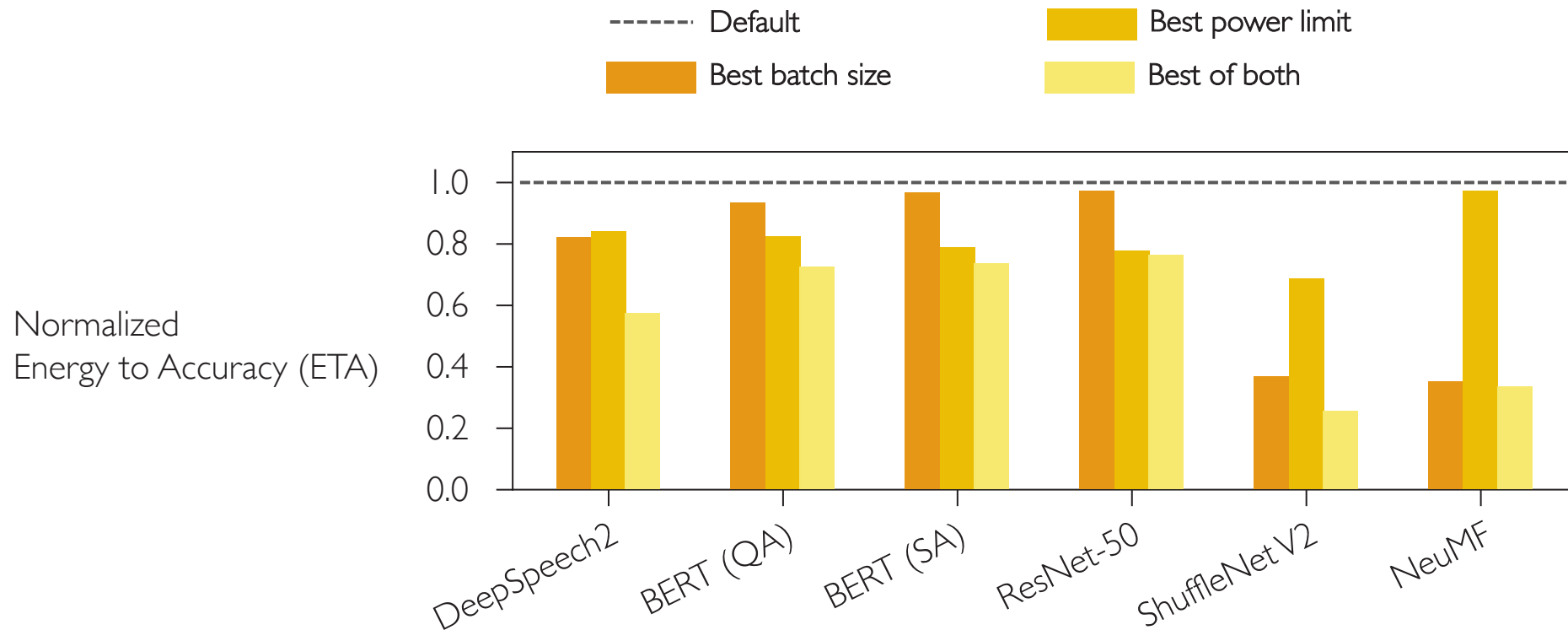
The diagram illustrates the relationship between Energy Time At Hand (ETA), Time To Answer (TTA), and Average Power (AvgPower). It shows the equation:  $\text{ETA} = \text{TTA} \times \text{AvgPower}$ . The units are specified as Joule for ETA, Second for TTA, and Watt for AvgPower.

# Understanding GPU Energy Consumption



# Opportunity for Energy Savings

## Sweep of feasible batch sizes and power limits



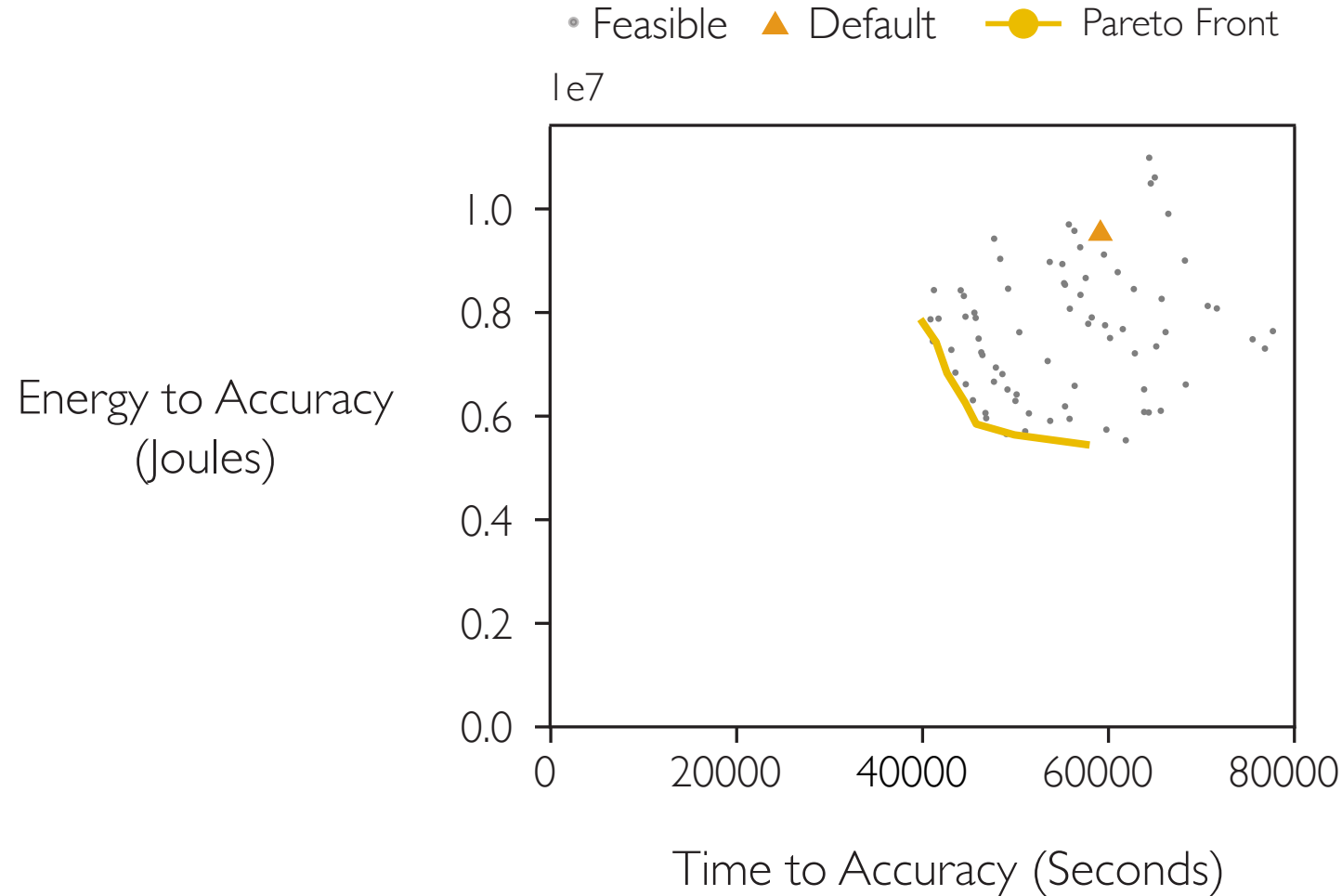
**24 ~ 75%  
energy  
reduction**

Measured on an NVIDIA V100 GPU.

Training terminates when the DNN reaches its original target accuracy.

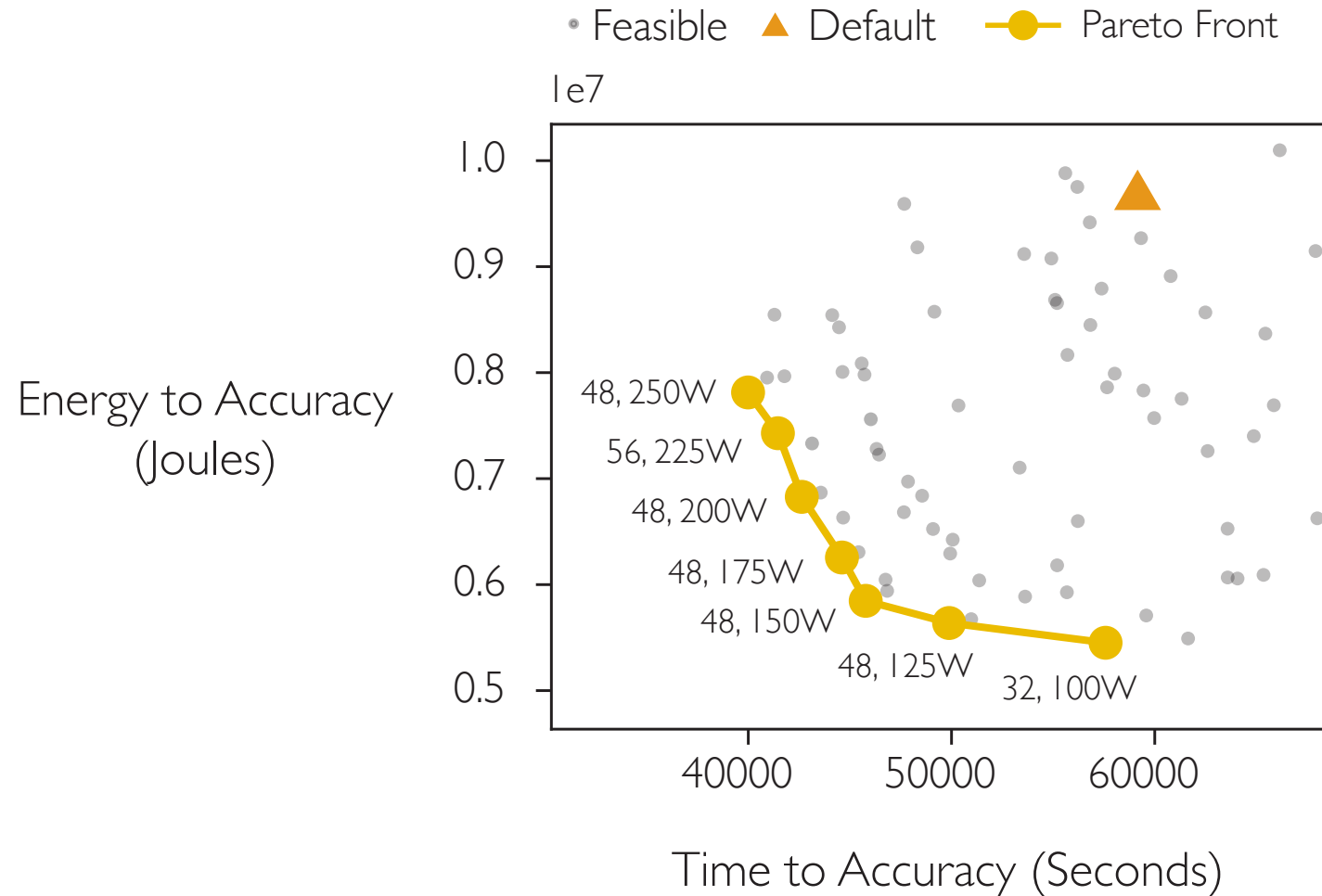


# Relationship Between Time and Energy



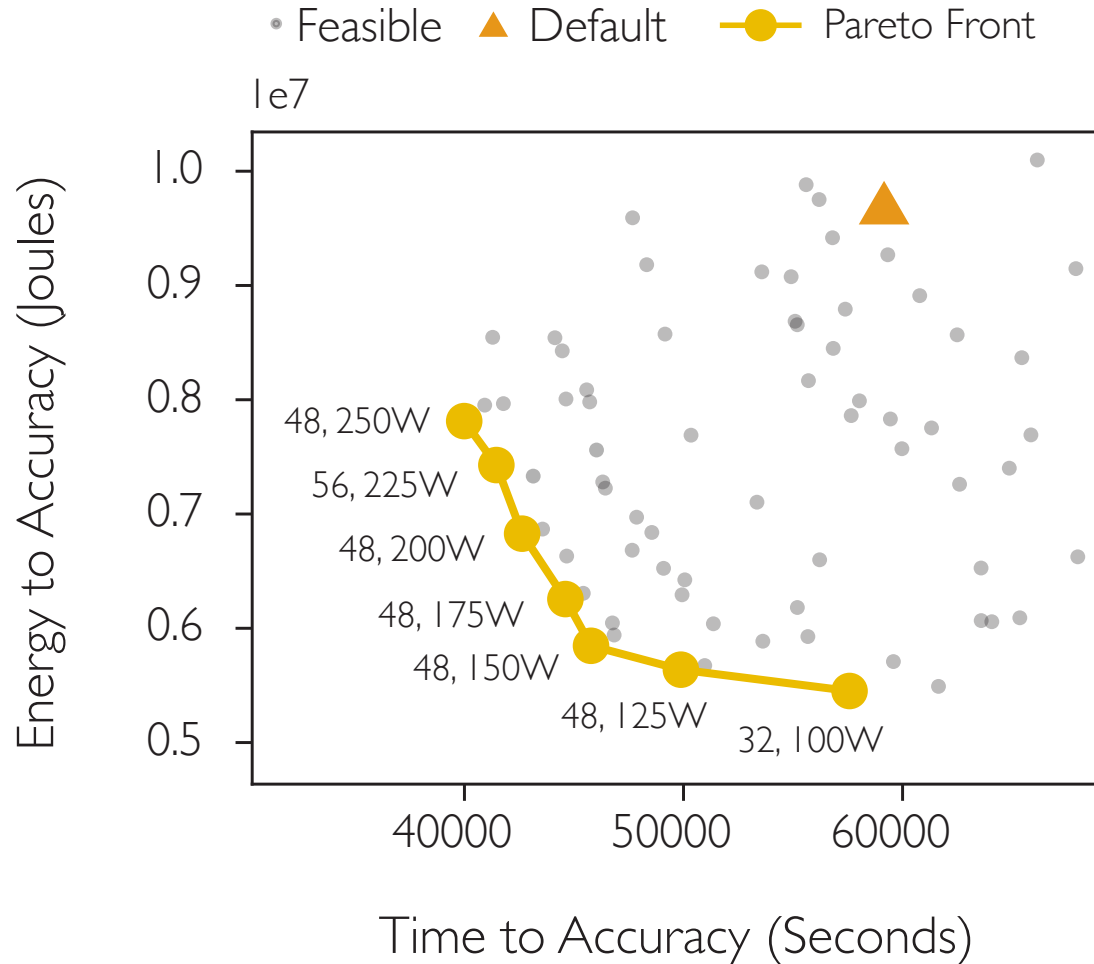
Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100 GPU.  
Similar trends found across 6 DL workloads and 4 GPU generations.

# Relationship Between Time and Energy



Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100 GPU.  
Similar trends found across 6 DL workloads and 4 GPU generations.

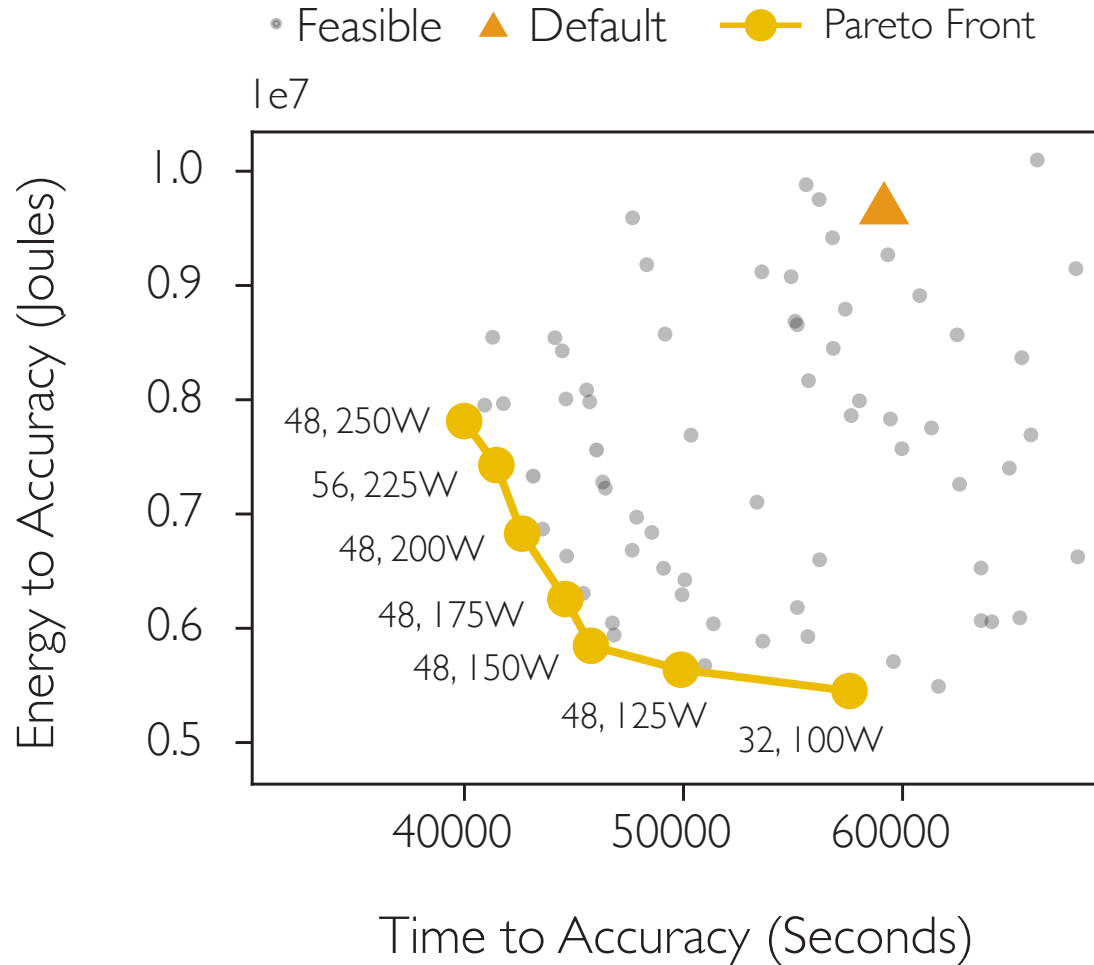
# Relationship Between Time and Energy



1. Time and energy minimized by different knobs
2. Efficient time and energy show a **trade-off**

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100 GPU.  
Similar trends found across 6 DL workloads and 4 GPU generations.

# Relationship Between Time and Energy

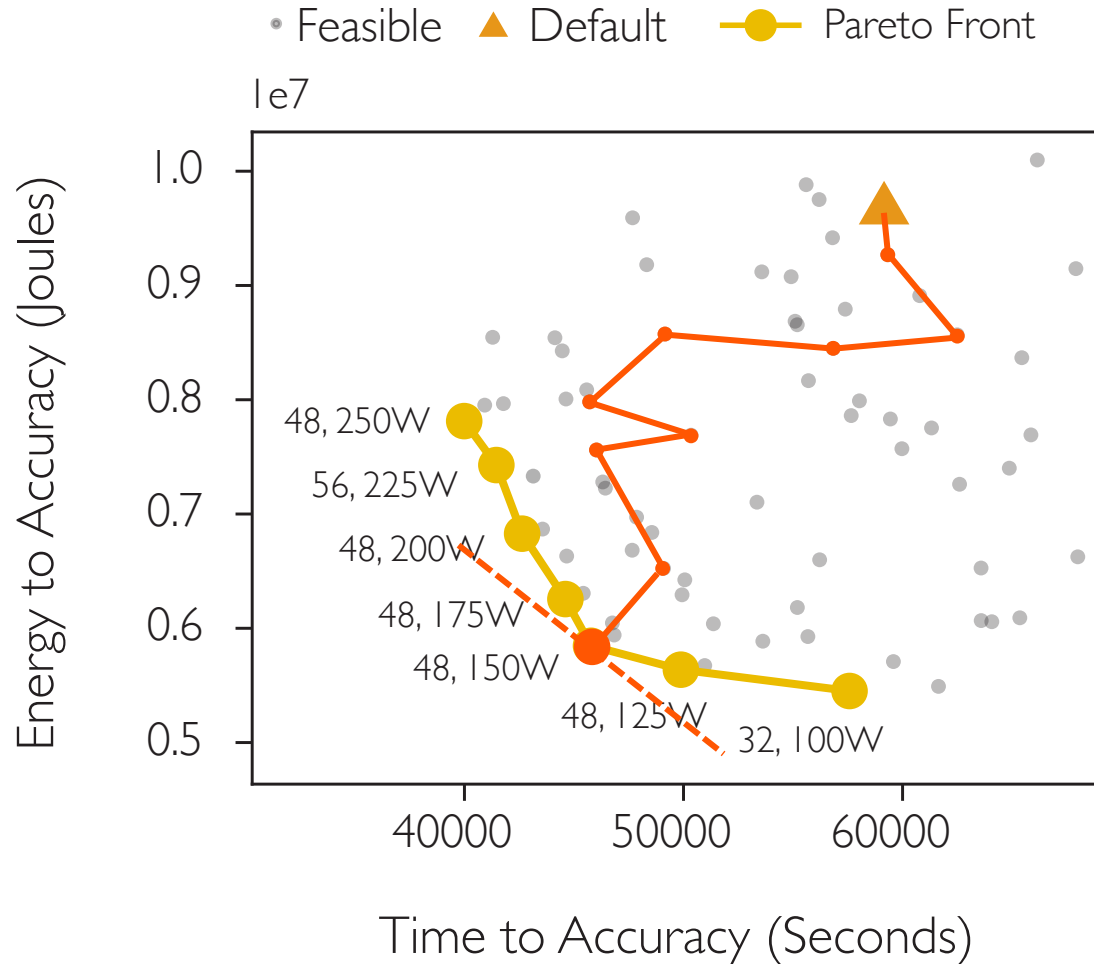


Which yellow point is the best?

$$\text{Cost} = \eta \cdot \text{ETA} + (1 - \eta) \cdot \text{MaxPower} \cdot \text{TTA}$$

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100 GPU.  
Similar trends found across 6 DL workloads and 4 GPU generations.

# Relationship Between Time and Energy



Which yellow point is the best?

$$\text{Cost} = \eta \cdot \text{ETA} + (1 - \eta) \cdot \text{MaxPower} \cdot \text{TTA}$$

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100 GPU.  
Similar trends found across 6 DL workloads and 4 GPU generations.



*An Energy Optimization Framework  
for DNN Training*

## **Optimizes the cost**

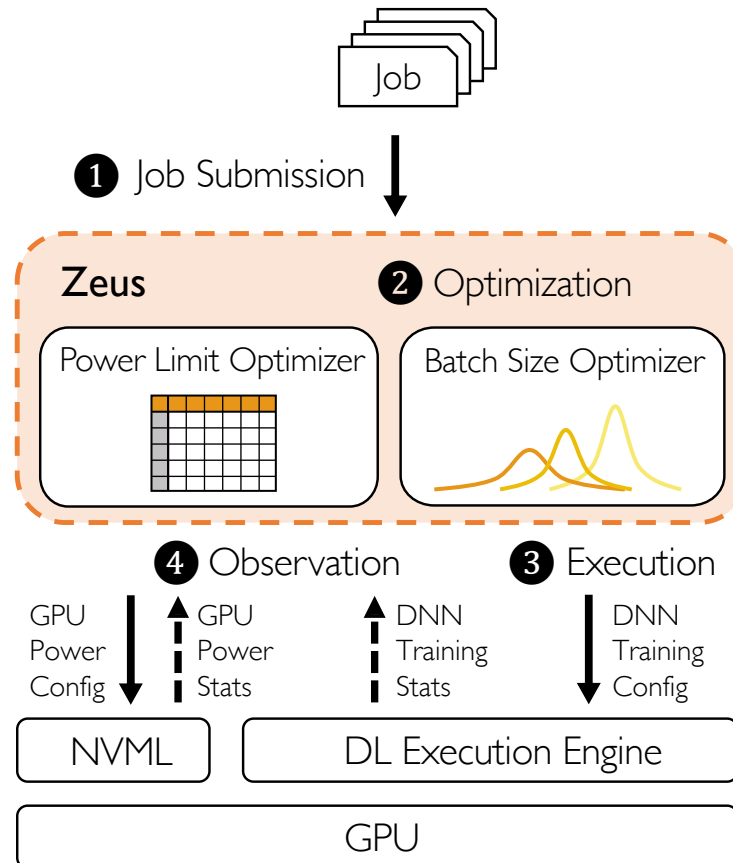
- of an arbitrary DNN model
- on an arbitrary GPU type
- in an efficient manner

## **without any**

- offline profiling,
- hardware modification, or
- accuracy degradation

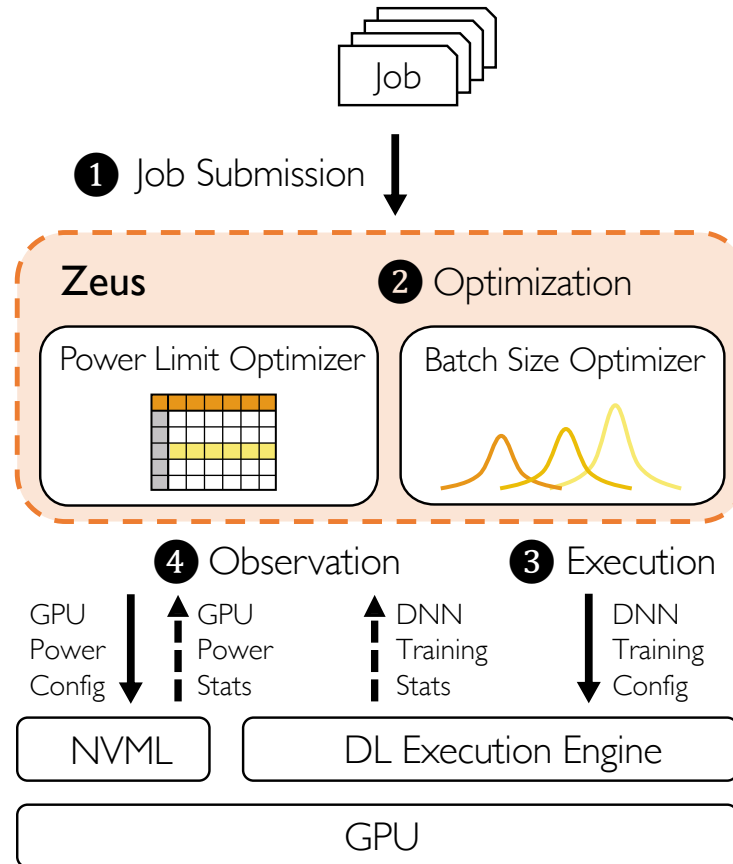
# Overall Workflow

Re-training jobs are opportunity for exploration!



# Overall Workflow

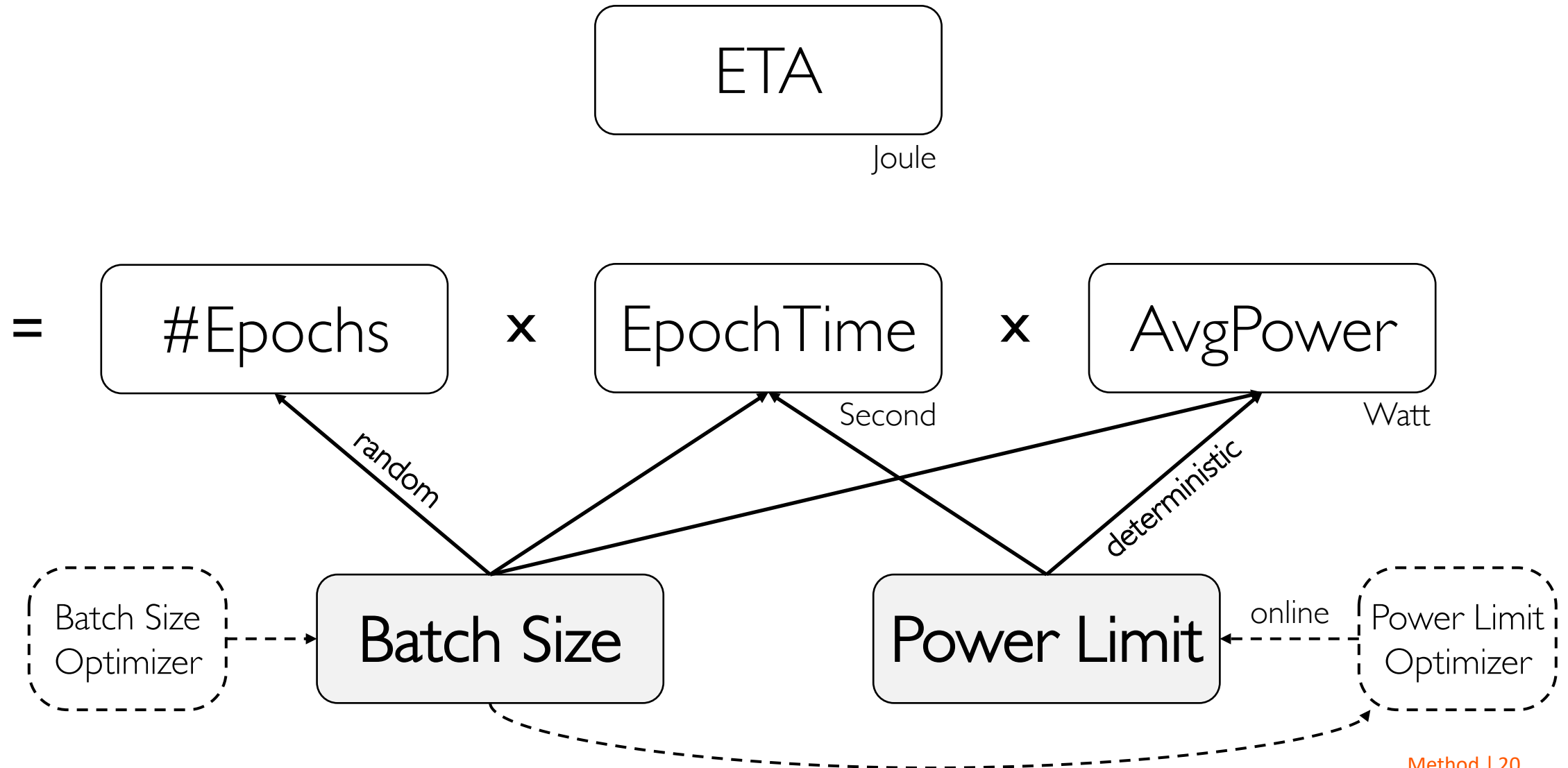
Re-training jobs are opportunity for exploration!



1. Decoupling Variables
2. Power Limit Optimizer
3. Batch Size Optimizer



# I. Decoupling Batch Size and Power Limit



# 2. Power Limit Optimizer

## Just-in-time online profiler

- Profiles the **power** and **throughput** of each power limit
- **Five seconds** per power limit is enough

## Low overhead

- Profile **only once** for each batch size
- Profiling **contributes** to the training process

# 3. Batch Size Optimizer

A good solution must

1. incorporate the **stochasticity** of DNN training, and
2. intelligently **trade-off** exploration and exploitation

$$\text{Cost} = \eta \cdot \text{ETA} + (1 - \eta) \cdot \text{MaxPower} \cdot \text{TTA}$$

**Multi-Armed Bandit**

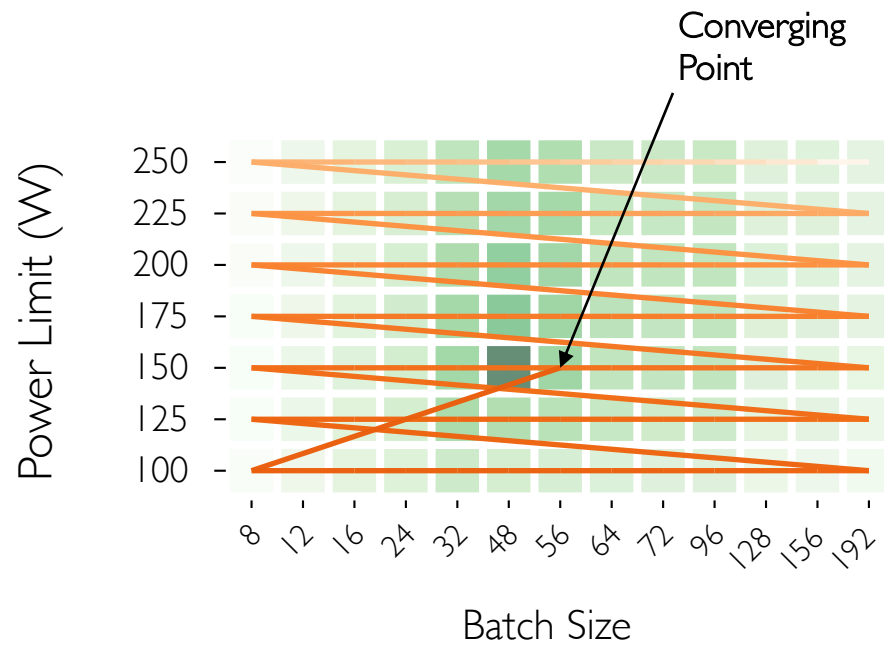
1. Models **cost** as a Gaussian random variable
2. Automatically controls exploration and exploitation

# Workloads and GPU Generations

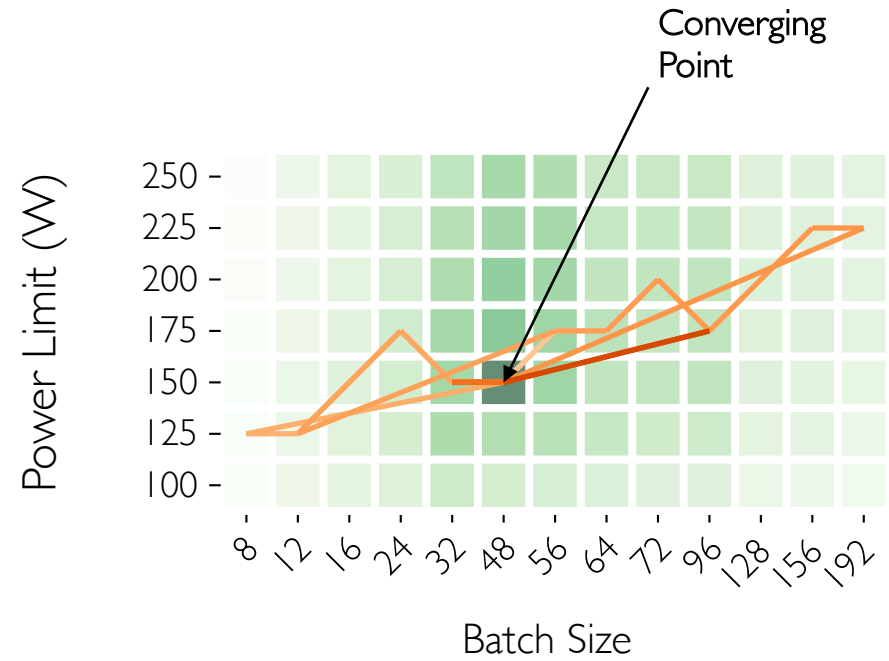
Task	Dataset	DNN	GPU	Arch
Speech Recognition	LibriSpeech	DeepSpeech2	NVIDIA A40	Ampere
Question Answering	SQuAD	BERT	NVIDIA V100	Volta
Sentiment Analysis	Sentiment140	BERT	NVIDIA RTX6000	Turing
Image Classification	ImageNet	ResNet-50	NVIDIA P100	Pascal
Image Classification	CIFAR-100	ShuffleNet-v2		
Recommendation	MovieLens-1M	NeuMF		

# Zeus in Action

## Grid Search



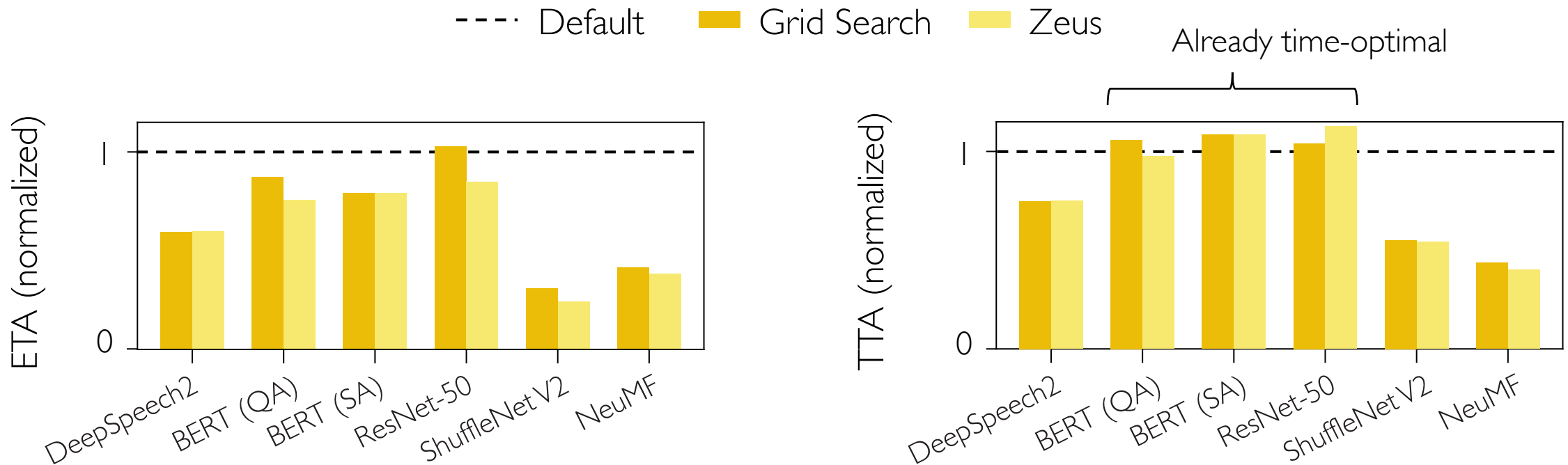
## Zeus



Search Path Training Cost (darker means better)

DeepSpeech2 trained on LibriSpeech on an NVIDIA V100 GPU.

# Zeus Leads to Large Benefits



**15 ~ 76% energy reduction**

**Up to 60% time reduction**

Results obtained on an NVIDIA V100 GPU

# Demo: Stable Diffusion

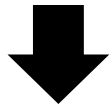
<https://youtu.be/MzIF5XNRSJY>

# Conclusion

DNN



GPU



Energy

- Works on arbitrary DNN models
- Works without modifying existing hardware
- Fully online with JIT profiling and MAB
- Jointly optimizes both job- and GPU-side configurations





<https://ml.energy/zeus>