# ZEUS
# Understanding and Optimizing
# GPU Energy Consumption of DNN Training

Jae-Won Chung

October 7th, 2022
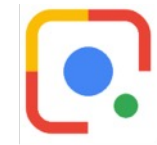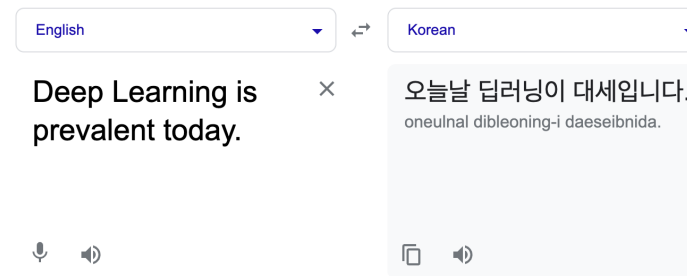
*Work done in collaboration with Jie You and Mosharaf Chowdhury*
*To appear at NSDI '23*

SymbioticLab

UNIVERSITY OF MICHIGAN

# Deep Learning is Prevalent Today

Image processing

Speech recognition

Machine translation

Intelligent assistants

Autonomous driving

Search

Video analytics

# DNN Energy Consumption is Skyrocketing

**DNN**
- Re-training is commonplace (e.g. every hour)[3]

**GPU**
- Dominant power consumer in servers (~70%)[1]
- Training GPT-3 == 120 years of electricity for a household[2]

**Energy**
- Performance optimizations oblivious of energy impact
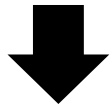
1. Dodge et al. (arXiv '22)   2. U.S. EIA and Google (arXiv '21)   3. Facebook (HPCA '18) and Alibaba (NSDI '22)

# Existing Efforts are not Practical Enough

DNN

$\downarrow$

- New energy-efficient DNN architectures

  SqueezeNext (CVPRW '18), ChamNet (CVPR '19), SkyNet (MLSys '20)

GPU

$\downarrow$

- New energy-efficient HW architectures

  TPU (ISCA '17), EDEN (MICRO '19), LNPU (ISSCC '19)

Energy

- Offline profiling and power model fitting
- Confined to GPU power configuration knobs

  MPC (HPCA '17), ODPP (CCGRID '20), GPOEO (TPDS '22)

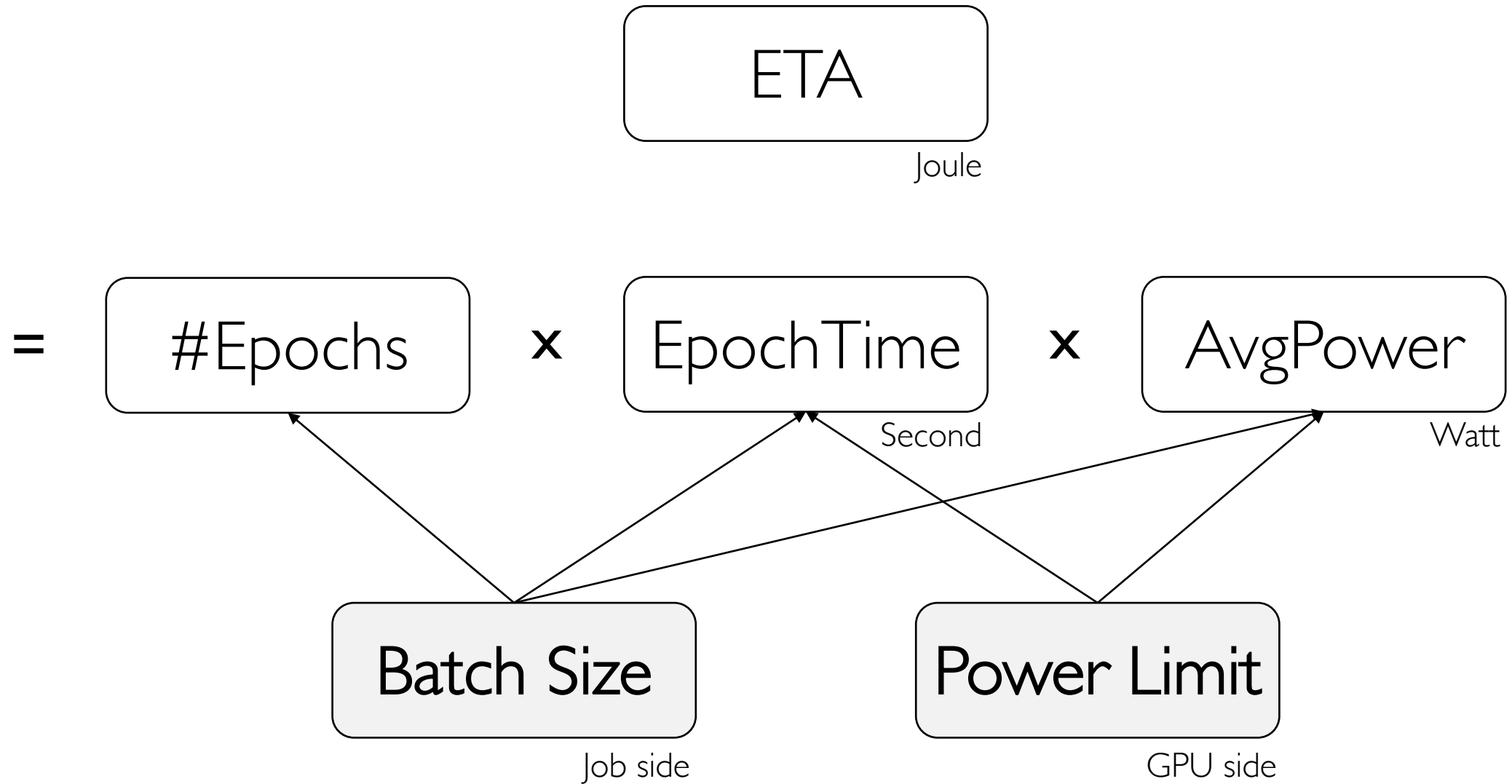# Understanding GPU Energy Consumption

## *Energy to Accuracy* (ETA)

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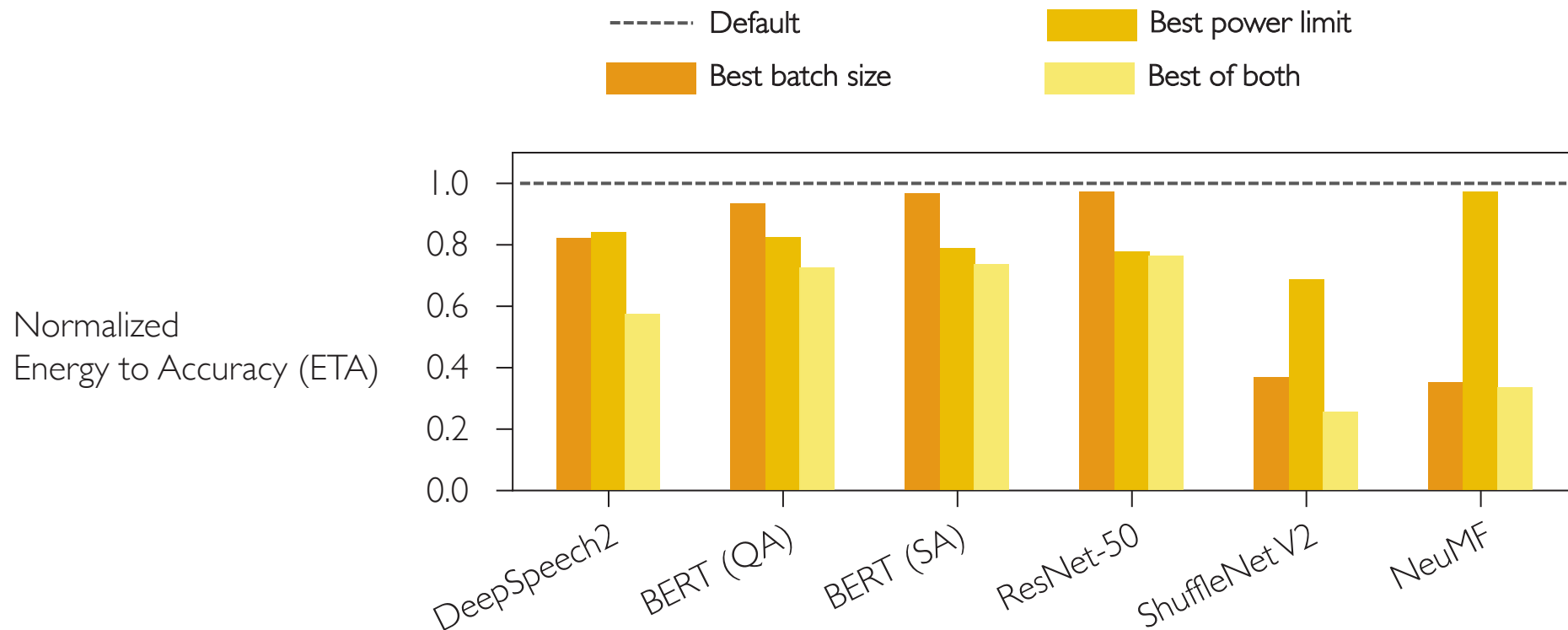
# Understanding GPU Energy Consumption

ETA

Joule

$$= \text{TTA} \times \text{AvgPower}$$

Second

Watt

# Understanding GPU Energy Consumption

ETA

Joule

= #Epochs × EpochTime × AvgPower

Second

Watt

Batch Size

Job side

Power Limit

GPU side

# Opportunity for Energy Savings

## Sweep of feasible batch sizes and power limits



Normalized Energy to Accuracy (ETA)

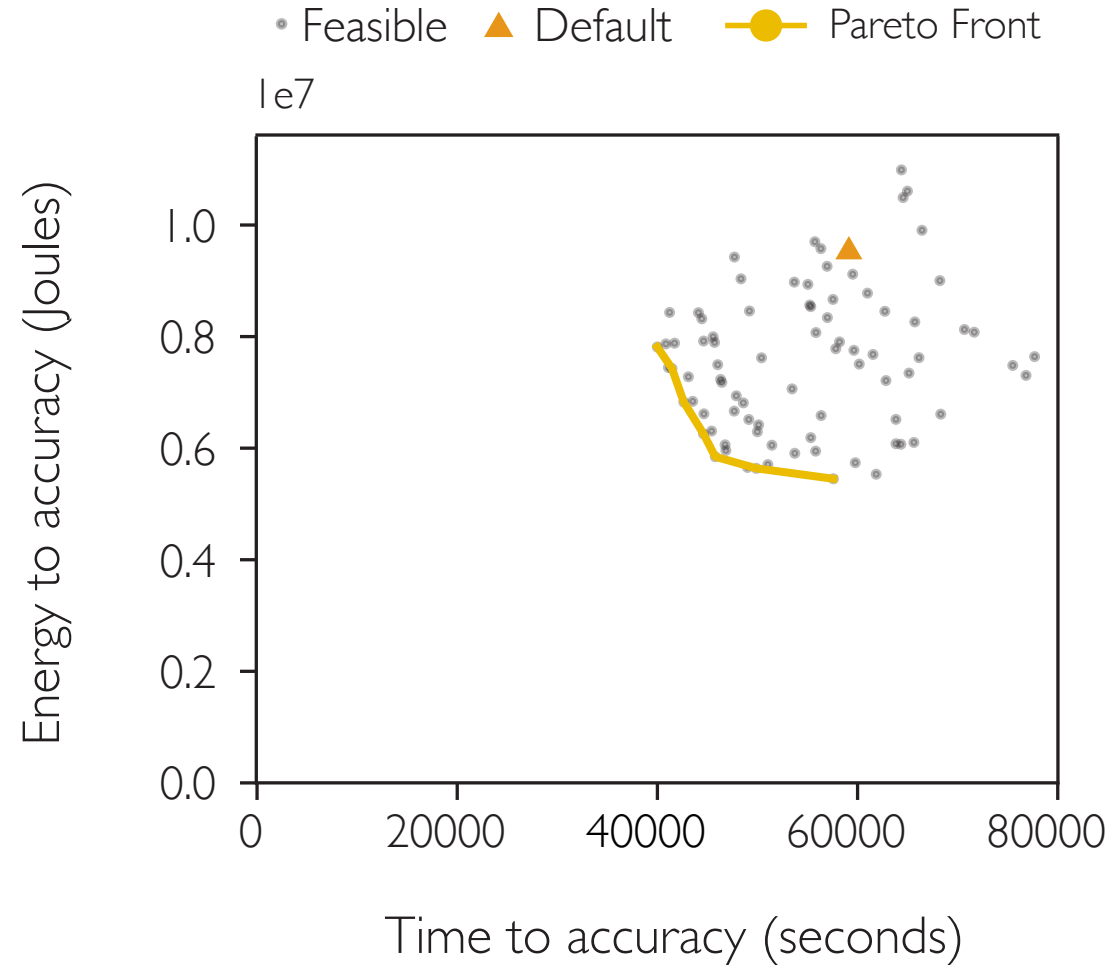Legend: ------- Default · Best batch size · Best power limit · Best of both

24 ~ 75% energy reduction

Measured on an NVIDIA V100 GPU.
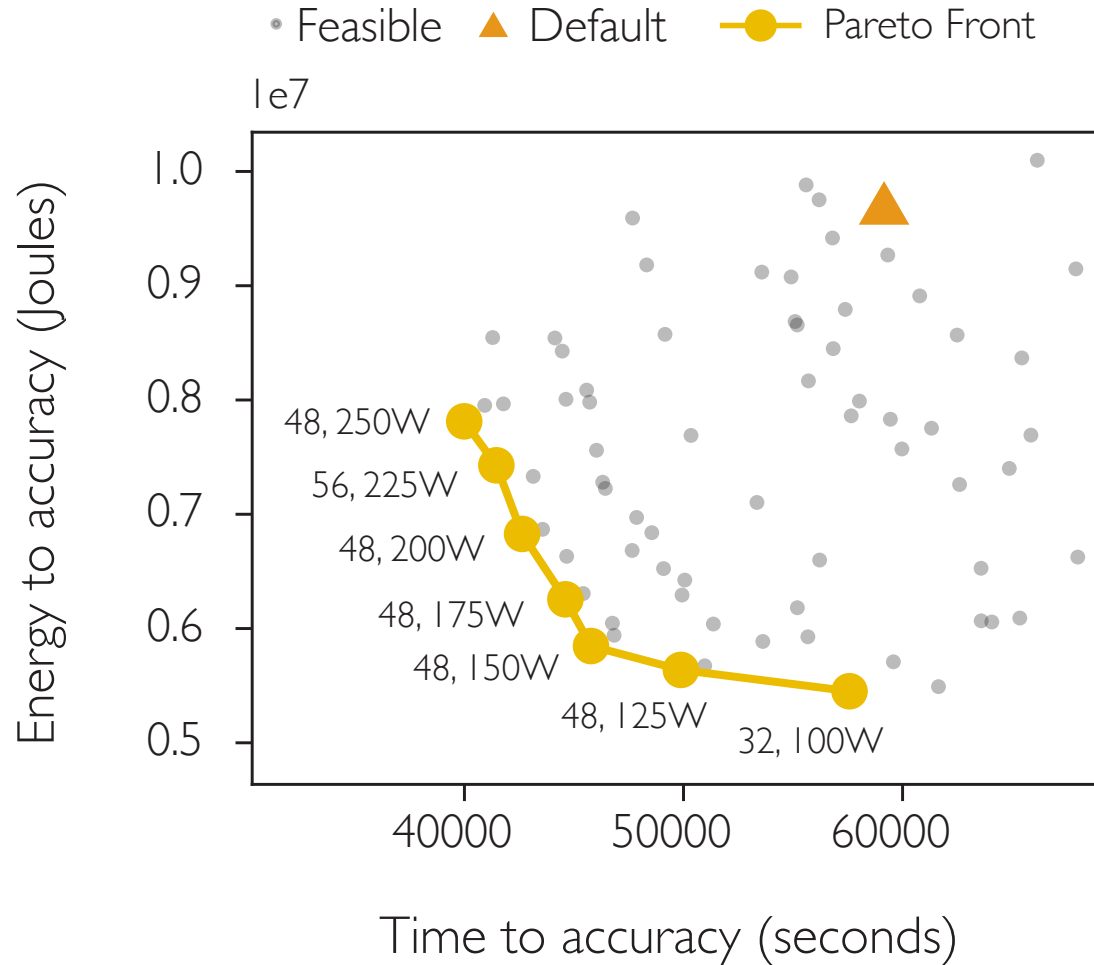Training terminates when the DNN reaches its original target accuracy.

# Relationship Between Time and Energy



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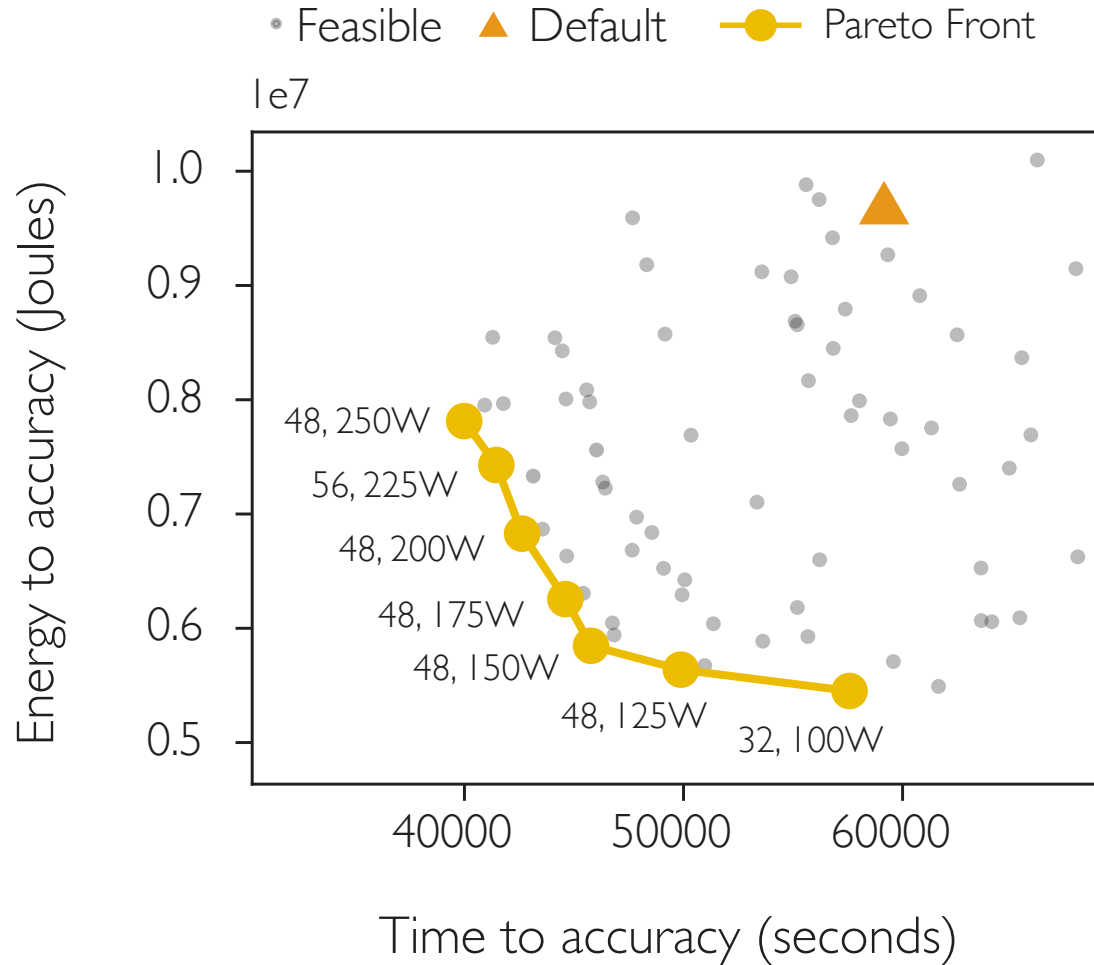Similar trends found over 6 DL workloads and 4 GPU generations.

# 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.

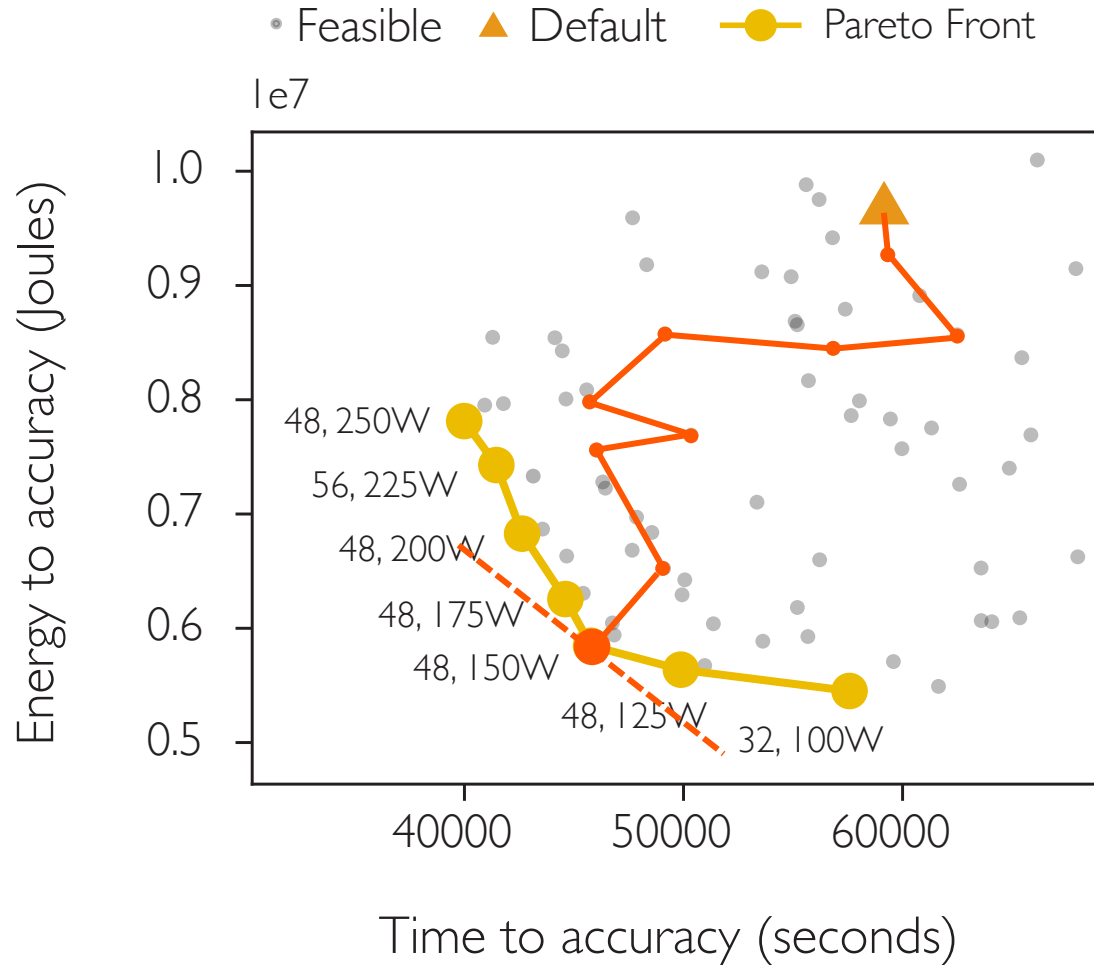# Relationship Between Time and Energy



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.

# Challenge #1: Average Power

ETA

| #Epochs | × | EpochTime | × | AvgPower |

## GPU is a black box

- Confidential hardware architecture
- Unknown internal voltage/frequency control algorithm

## Power modelling lacks practicality

- Requires offline profiling
- Does not generalize to other DNNs and GPUs

# Challenge #2: Time to Accuracy (TTA)

**Difficult to predict number of epochs**

- Batch size affects model accuracy
- We would be solving HPO if we can predict TTA

ETA

| #Epochs | × | EpochTime | × | AvgPower |

**DNN training is stochastic**

- Parameter initialization and batch order are random
- TTA varies even when we train with the same config

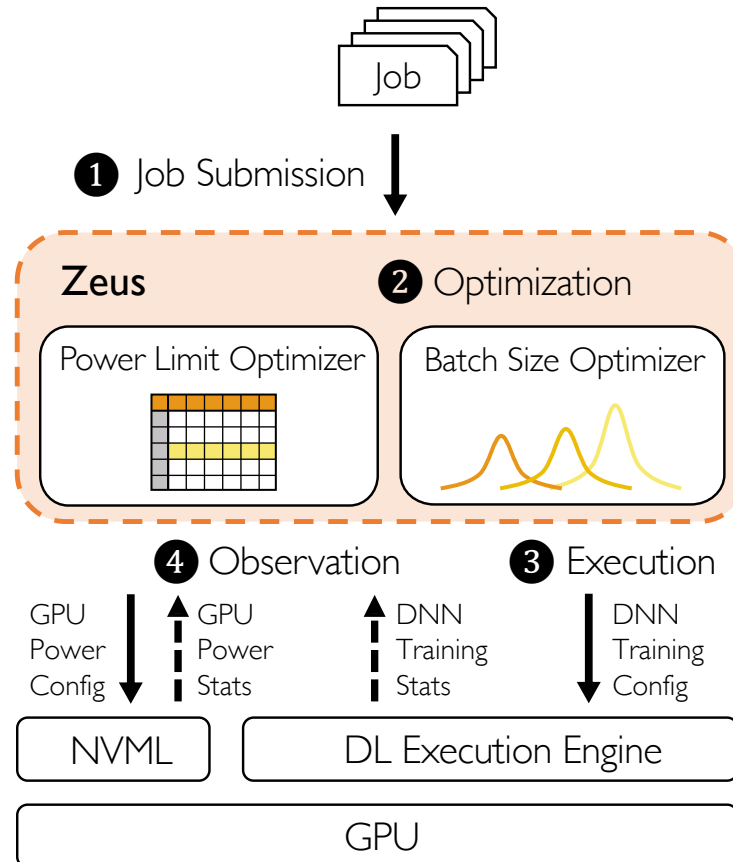*An Energy Optimization Framework for DNN Training*

**Optimizes the cost**

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

**without any**

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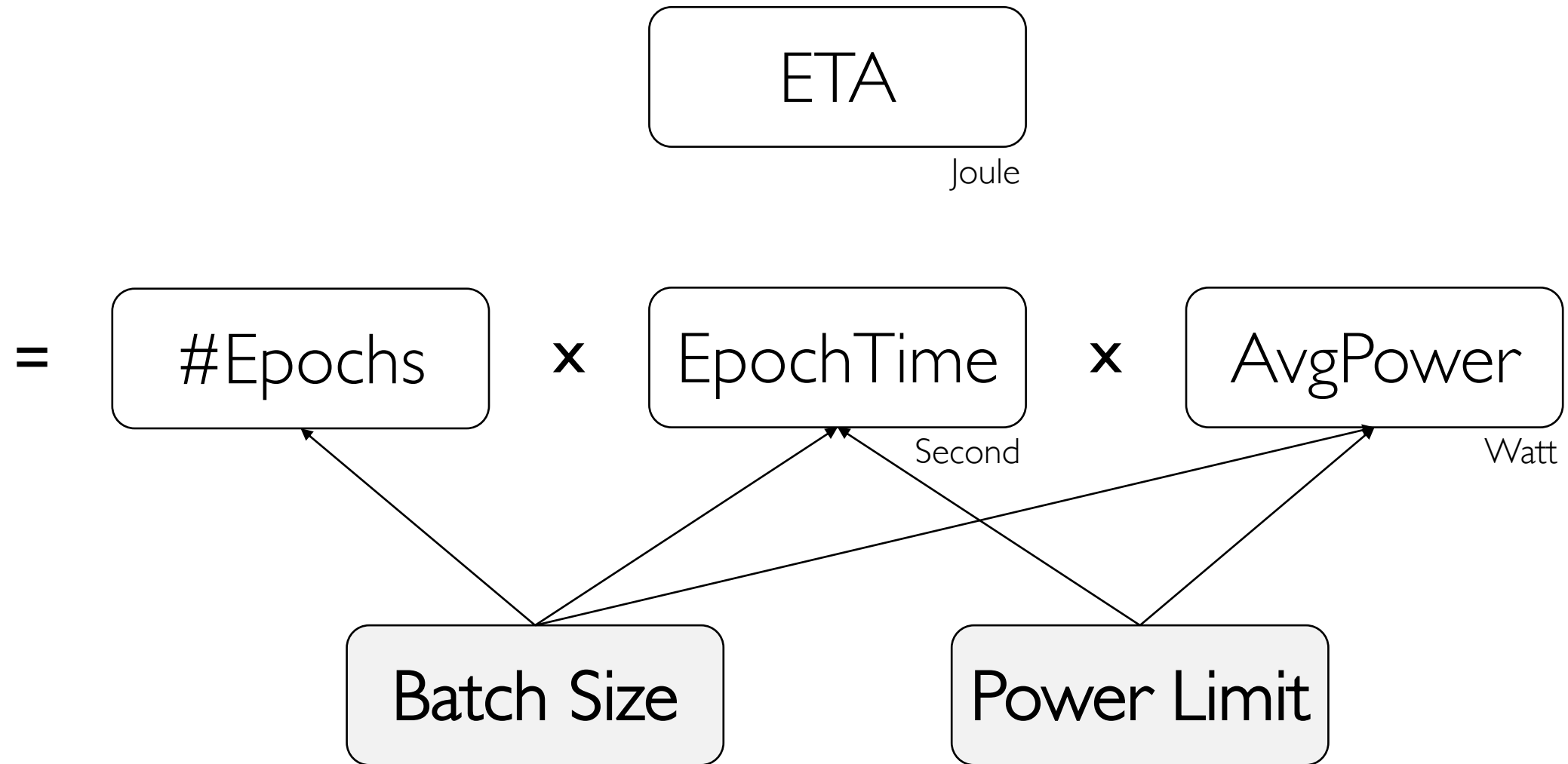- accuracy degradation

# Overall Workflow
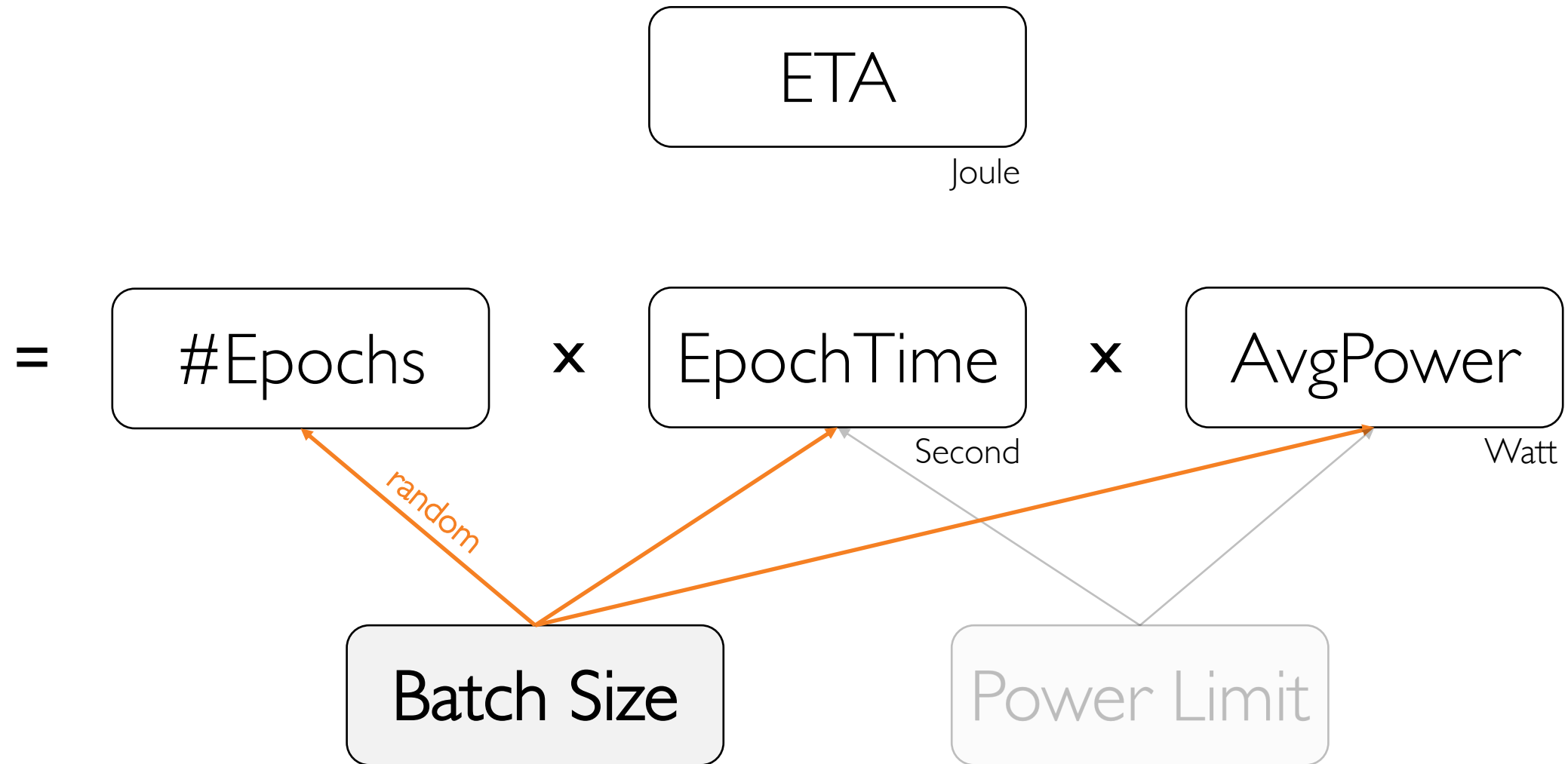
Re-training jobs are opportunity for exploration!



1. Decoupling
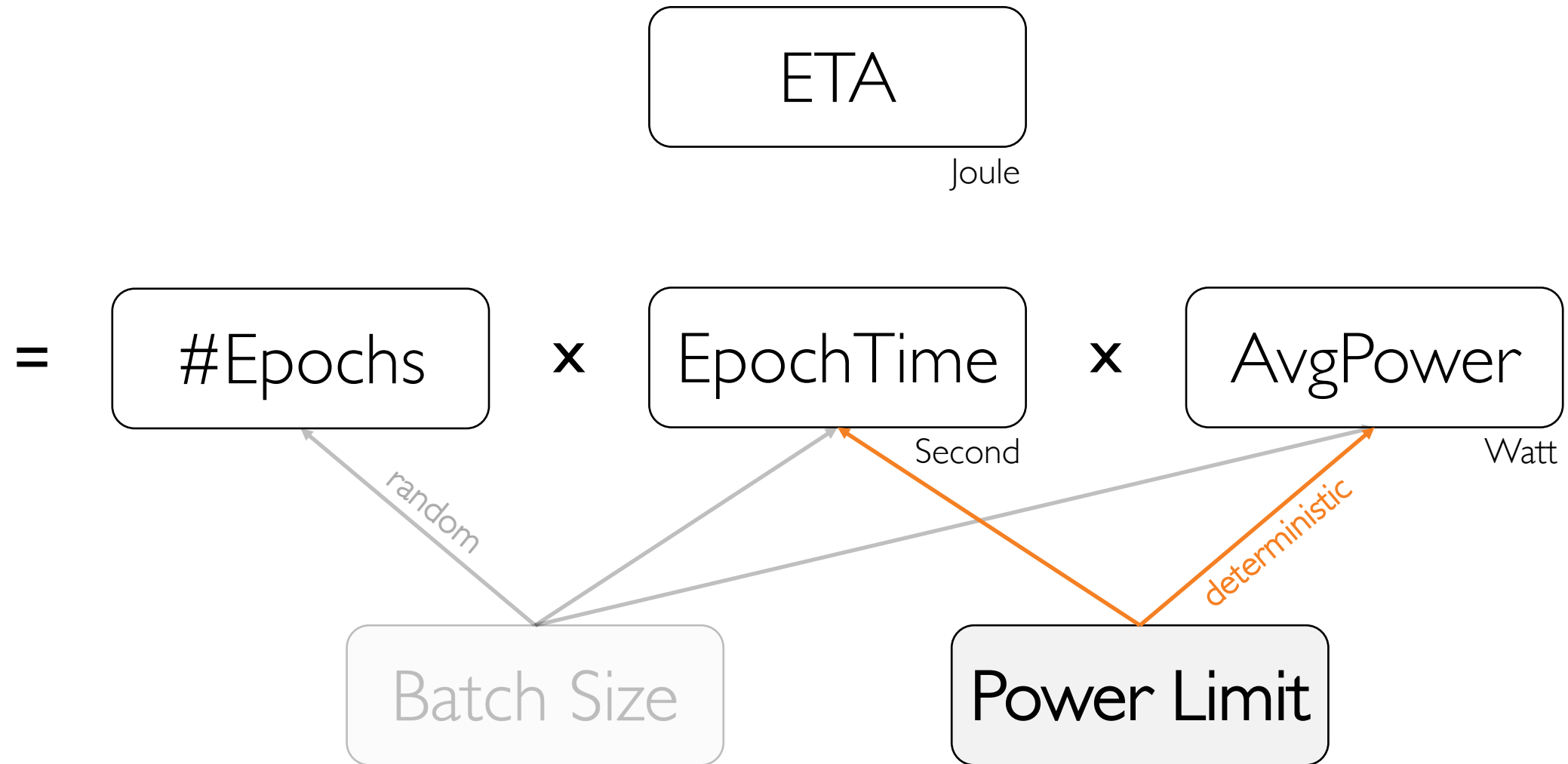2. Power Limit Optimizer
3. Batch Size Optimizer
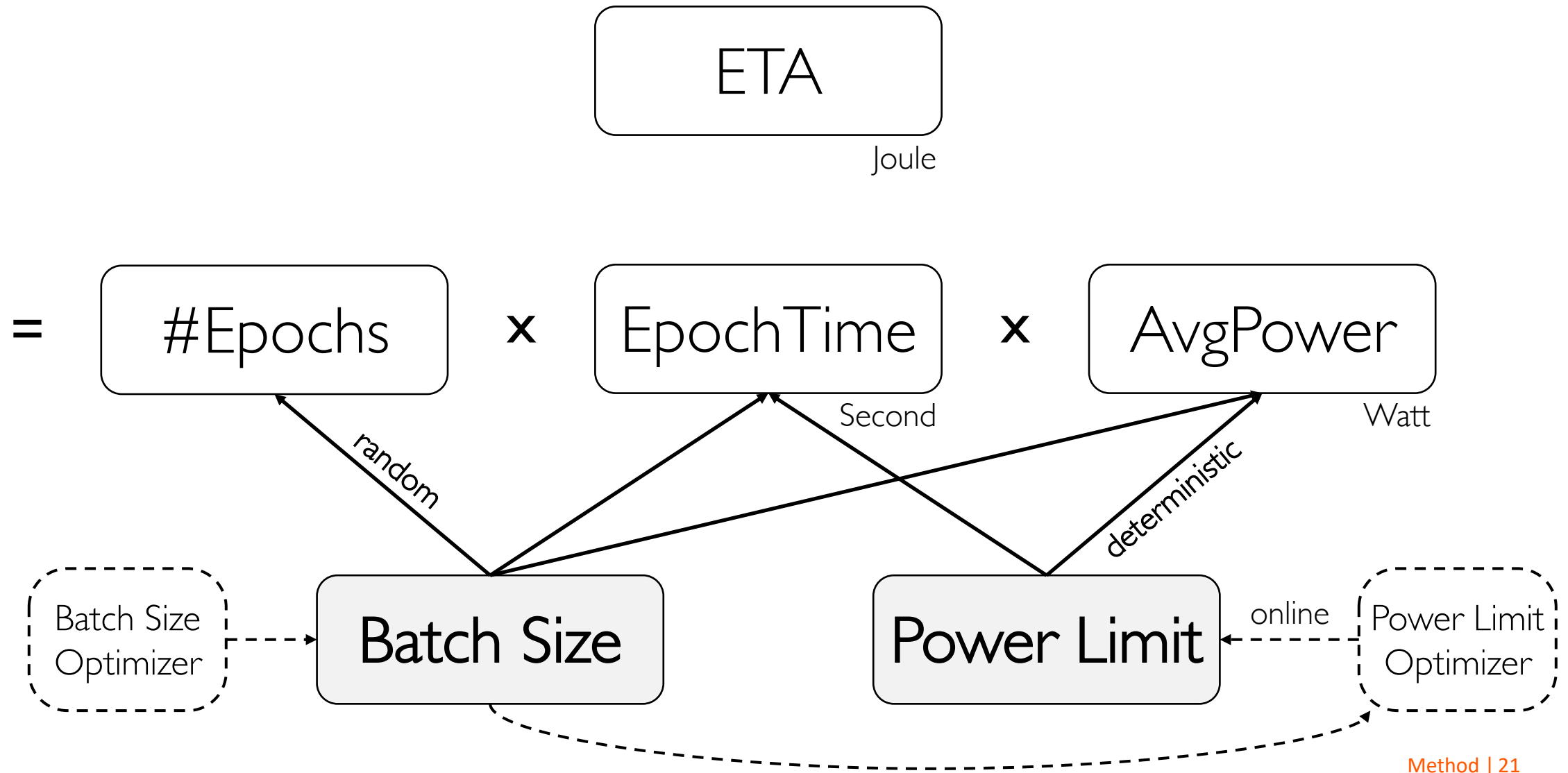
# 1. Decoupling Batch Size and Power Limit



ETA
Joule

= #Epochs × EpochTime × AvgPower
         Second        Watt

Batch Size    Power Limit

# 1. Decoupling Batch Size and Power Limit

# 1. Decoupling Batch Size and Power Limit

# 1. 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

$$\text{Cost} = \eta \cdot \text{ETA} + (1 - \eta) \cdot \text{MaxPower} \cdot \text{TTA}$$
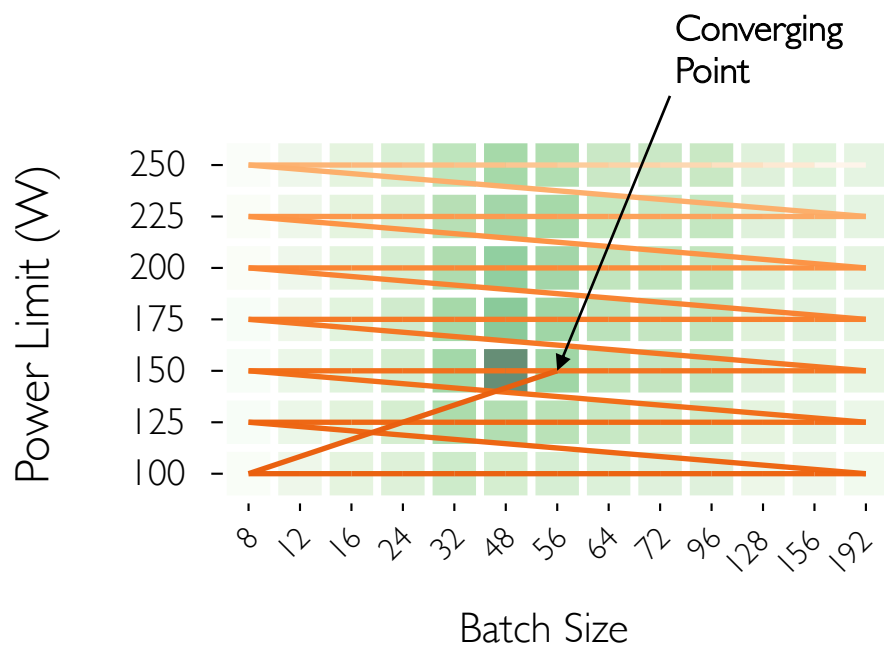
**Multi-Armed Bandit**
1. Models cost as a Gaussian random variable
2. Automatically controls exploration and exploitation
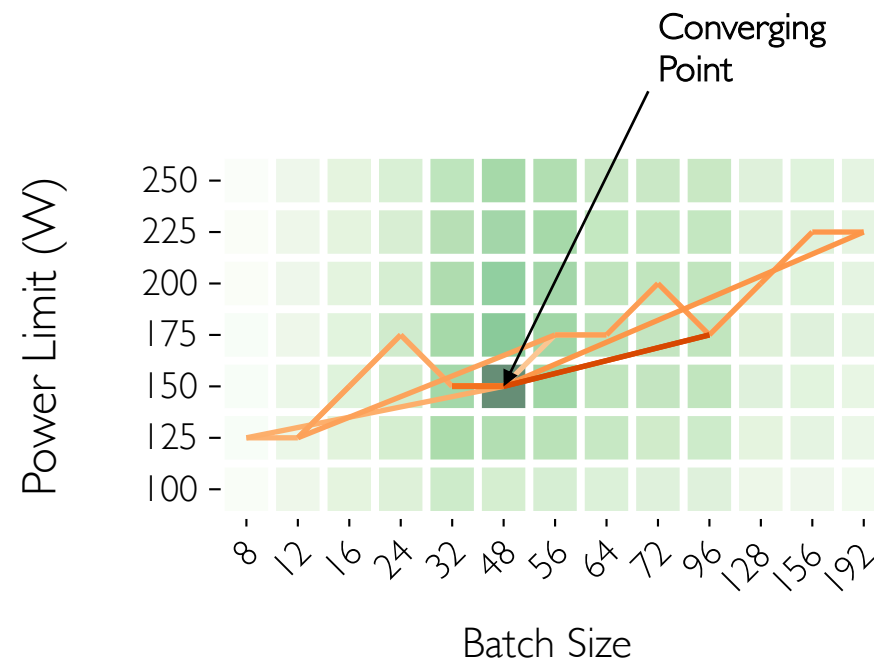
# Workloads and GPU Generations

| 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



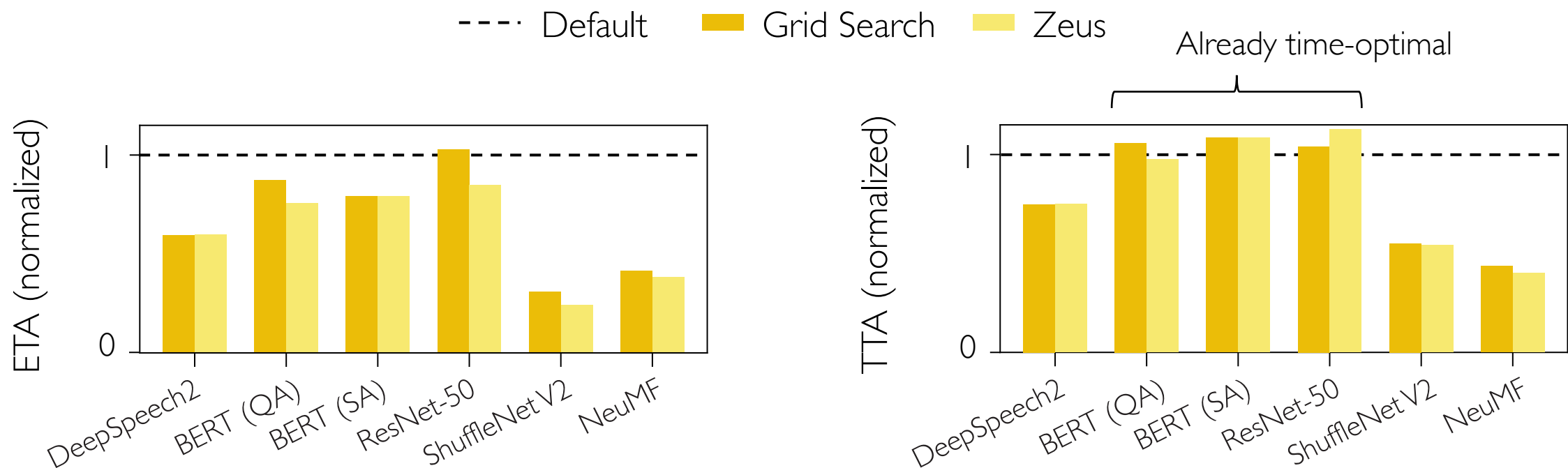DeepSpeech2 trained on LibriSpeech on an NVIDIA V100 GPU.

# Zeus Quickly Converges to the Optimum



ResNet50 trained on ImageNet on an NVIDIA V100 GPU

DeepSpeech2 trained on LibriSpeech on an NVIDIA V100 GPU
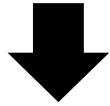
# Zeus Leads to Large Benefits



**15 ~ 76% energy reduction**
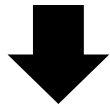
**Up to 60% time reduction**

Results obtained on an NVIDIA V100 GPU

# Conclusion

DNN
- Works on arbitrary DNN models

GPU
- Works without modifying existing hardware

Energy
- Fully online with JIT profiling and MAB
- Jointly optimizes both job- and GPU-side configurations

https://ml.energy/zeus