

Optimizing Energy with Performance in Mind

Ruofan Wu
December 2nd, 2025



ML.ENERGY

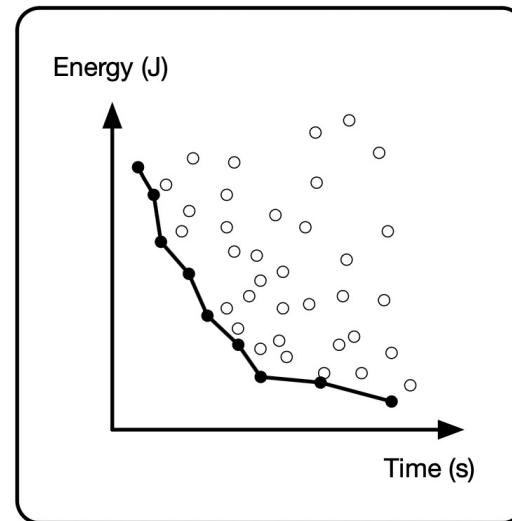


Energy Optimization for AI

Principles

- The **time–energy trade-off frontier** is a key object for reasoning.

- Same computation
- Different ways



Overview of Existing work

Serving

DynamoLLM (HPCA '25)

The ML.ENERGY Benchmark
(NeurIPS '25 D&B)

Training

Zeus (NSDI '23)

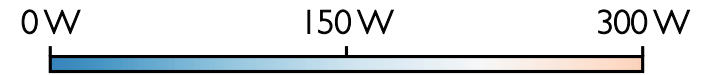
Perseus (SOSP '24)

DynamoLLM: Designing LLM Inference Clusters for Performance and Energy Efficiency

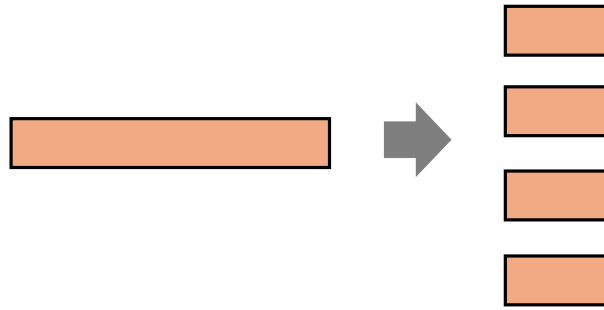
Jovan Stojkovic, Chaojie Zhang, Íñigo Goiri, Josep Torrellas, Esha Choukse

*How to minimize energy
consumption given
latency deadlines?*

Time vs. Energy Trade-off



Model parallelism



Time: 0.8s \rightarrow 0.25s
Energy: 200J \rightarrow 250J

GPU frequency



Time: 0.25s \rightarrow 0.27s
Energy: 75J \rightarrow 60J

LLM Inference

Prefill

What is DynamoLLM ?

Compute-bound

Model parallelism
GPU frequency

Time sensitive

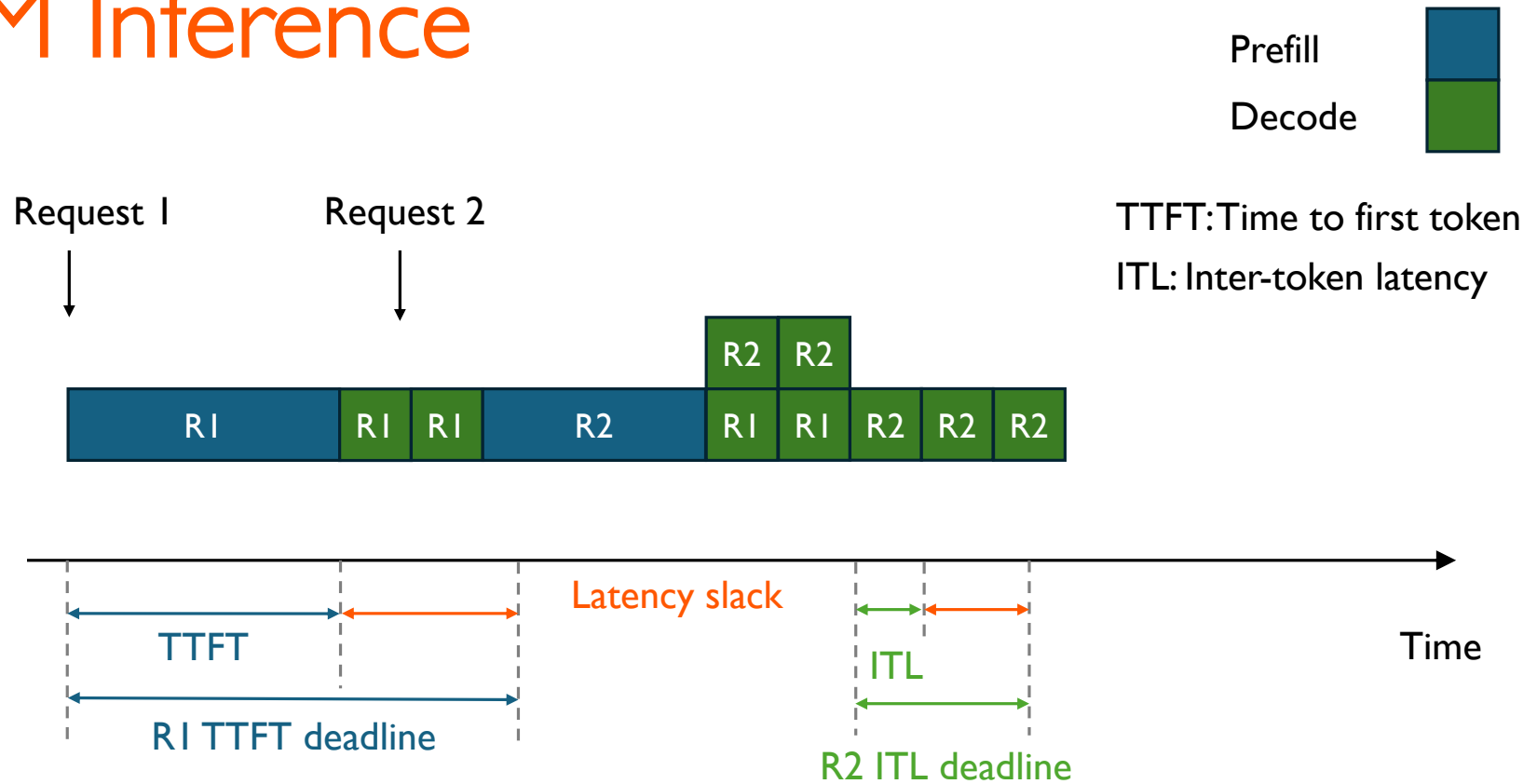
Decode

It's a dynamic energy-management system

Memory-bound

Time insensitive

LLM Inference



Heterogeneous Request Behavior

What is DynamoLLM?

It's a dynamic energy-management system designed for large-scale LLM inference clusters. It observes that different inference requests have vastly different compute and energy characteristics.

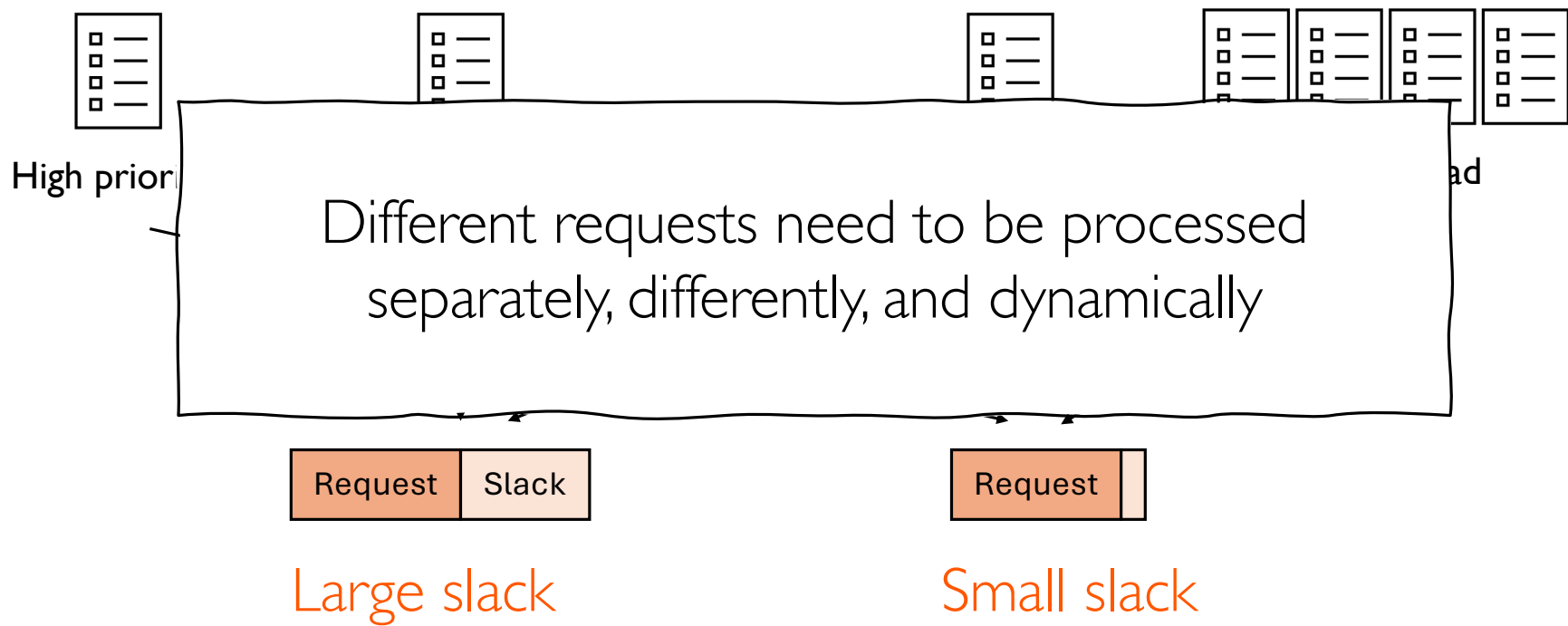
Memory-bound

I have 5 apples. I gave 2 to my friend and then bought 4 more. After that, I ate 3. How many apples do I have now?

4

Compute-bound

Heterogeneous Request Behavior



Hierarchical Control

	Level	Decision	Time scale
Request length Request load Time–energy frontier	Cluster	Number of instances	Minutes
	Pool	Model parallelism	Minute
	Instance	GPU frequency	Seconds

Up to 53% energy reduction
while meeting latency deadlines

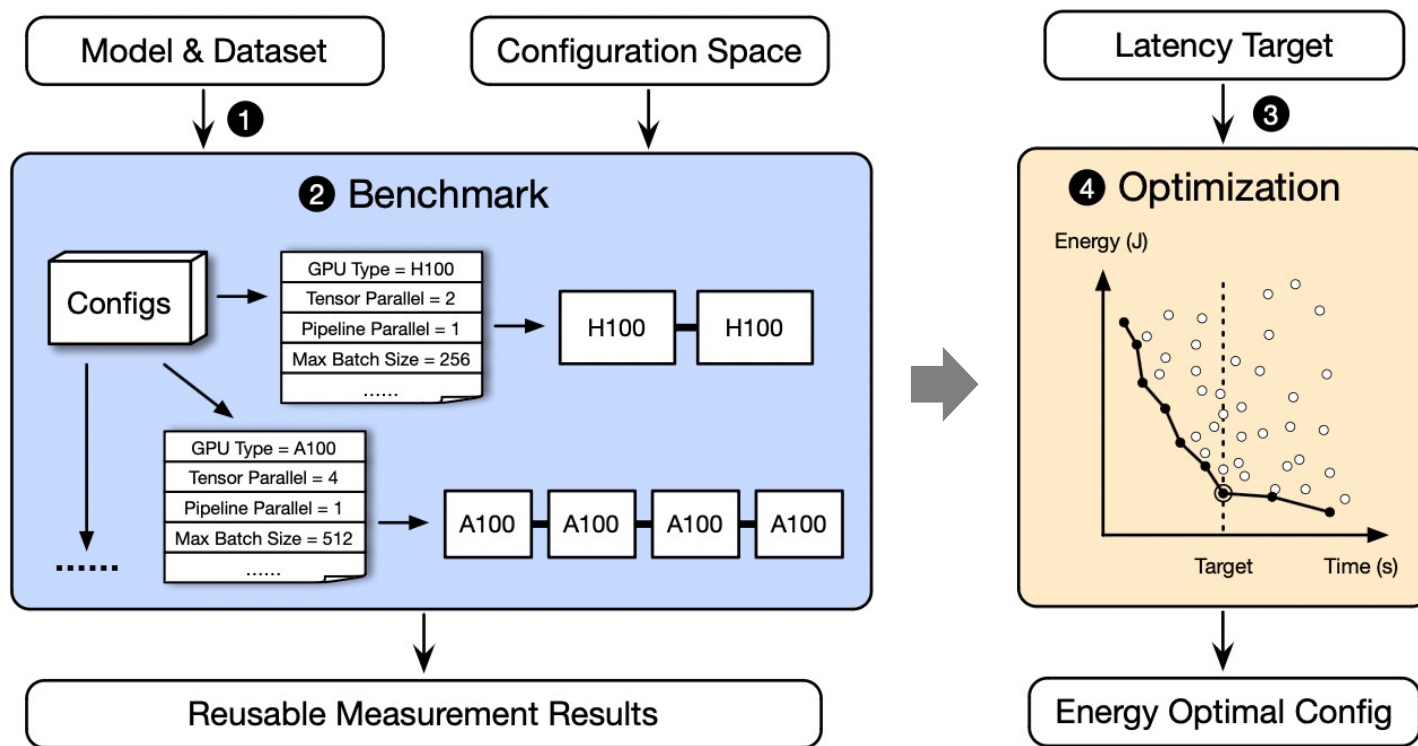
*The ML.ENERGY Benchmark:
Toward Automated Inference
Energy Measurement and
Optimization*

<https://ml.energy/leaderboard/>

*Jae-Won Chung, Jeff J. Ma, Ruofan Wu, Jiachen Liu,
Oh Jun Kweon, Yuxuan Xia, Zhiyu Wu, Mosharaf Chowdhury*

*“What are the
energy implications of
the choices we make?”*

Automated Optimization Recommendation



The ML.ENERGY Leaderboard

The ML.ENERGY Leaderboard

How much time and energy do generative AI models consume?

Version 3.0 / Last updated: November 30, 2025



About

Task

LLM Problem Solving

Text Conversation

Code Completion

MLLM Image Chat

Video Chat

Diffusion Text to Image

Text to Video

Conversational AI Chatbot [About](#)

Median ITL deadline: 200 ms

Per token energy budget: 21.86 J

GPU Models

☒ B200

☒ H100

Reset

Latency target & energy budget

Time vs. Energy Trade-off

Energy-optimal points for each model

18 models satisfy the given constraints (click row for model details).

Compare	Model	Precision
<input type="checkbox"/>	Llama 3.1 8B Instruct	bfloat16
<input type="checkbox"/>	Qwen 3 8B	bfloat16
<input type="checkbox"/>	Qwen 3 14B	bfloat16
<input type="checkbox"/>	Qwen 3 30B A3B Instruct	bfloat16
<input type="checkbox"/>	Gemma 3 12B	bfloat16
<input type="checkbox"/>	Qwen 3 32B	bfloat16
<input type="checkbox"/>	NVIDIA Nemotron Nano 12B V2	bfloat16
<input type="checkbox"/>	NVIDIA Nemotron Nano 9B V2	bfloat16
<input type="checkbox"/>	Gemma 3 27B	bfloat16
<input type="checkbox"/>	Llama 3.1 70B Instruct	bfloat16
<input type="checkbox"/>	Llama 3.3 70B Instruct	bfloat16

Model Configurations and the Pareto Frontier



<https://ml.energy/leaderboard/>

Zeus: Understanding and Optimizing GPU Energy Consumption of DNN Training

Jie You, Jae-Won Chung*, and Mosharaf Chowdhury*

*“How does energy
interact with time?”*

Understanding GPU Energy Consumption

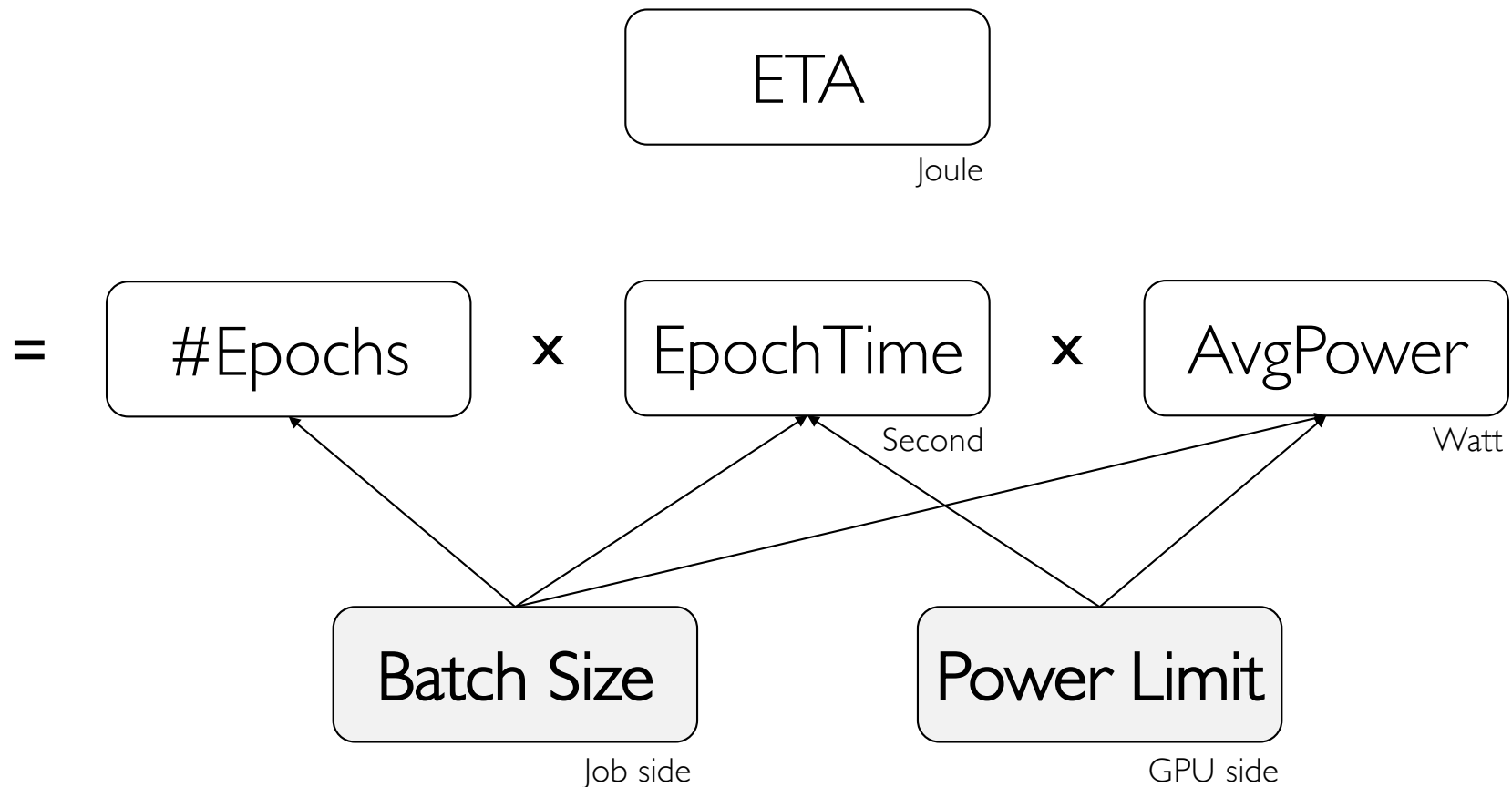
Energy to Accuracy (ETA)

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

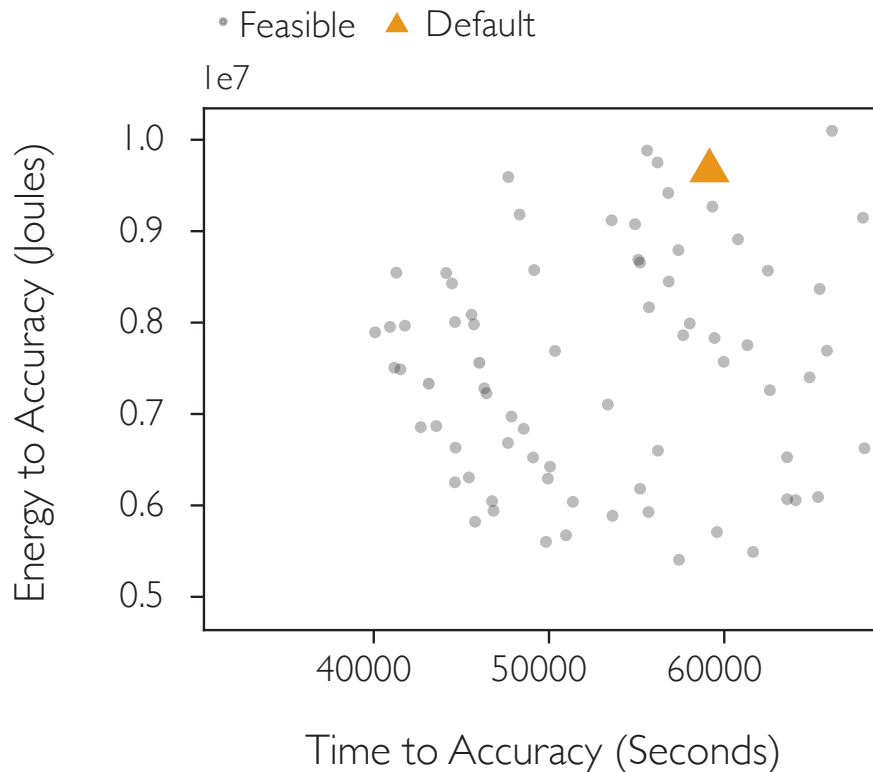
Understanding GPU Energy Consumption

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

Understanding GPU Energy Consumption



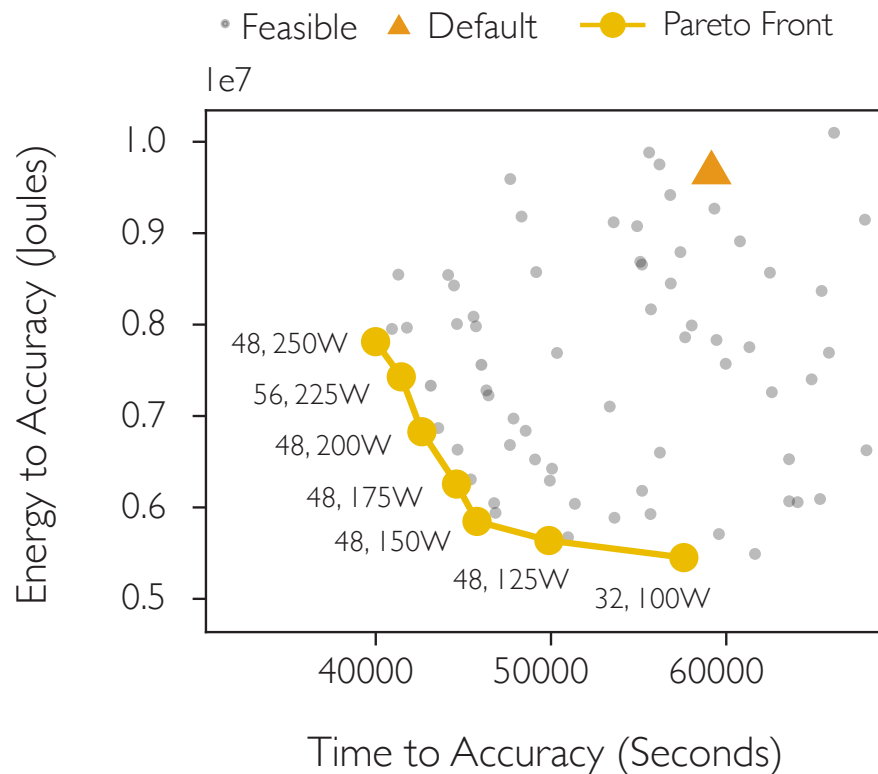
Opportunity for Energy Savings



Training time and total energy affected by
batch size and GPU power limit

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100.
Similar trends found across four GPU generations.

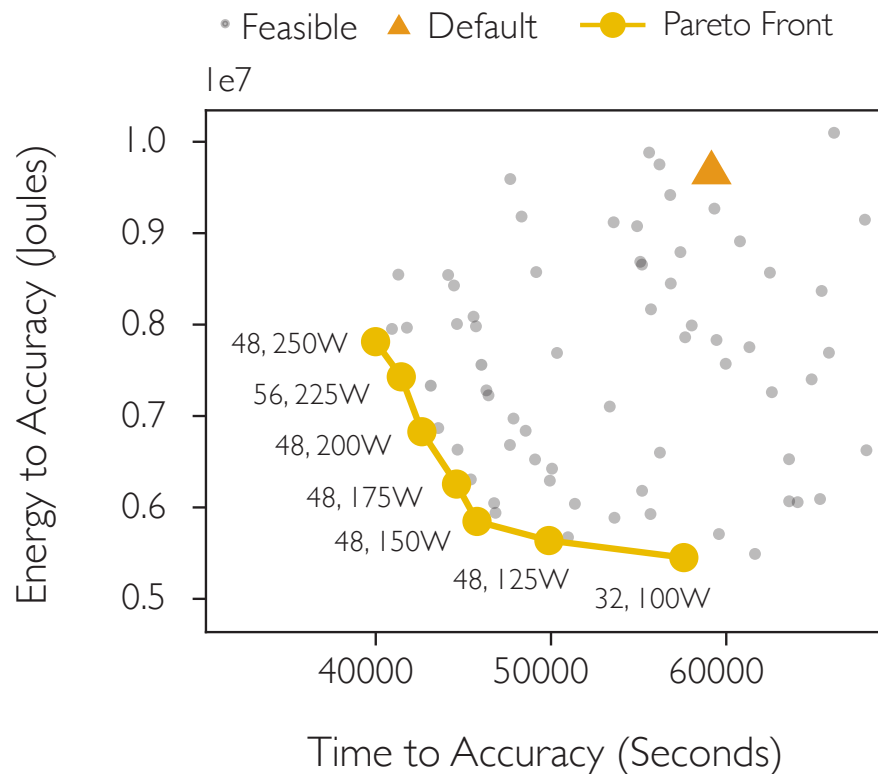
Time vs. Energy Trade-off



Training time and total energy affected by
batch size and GPU power limit

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100.
Similar trends found across four GPU generations.

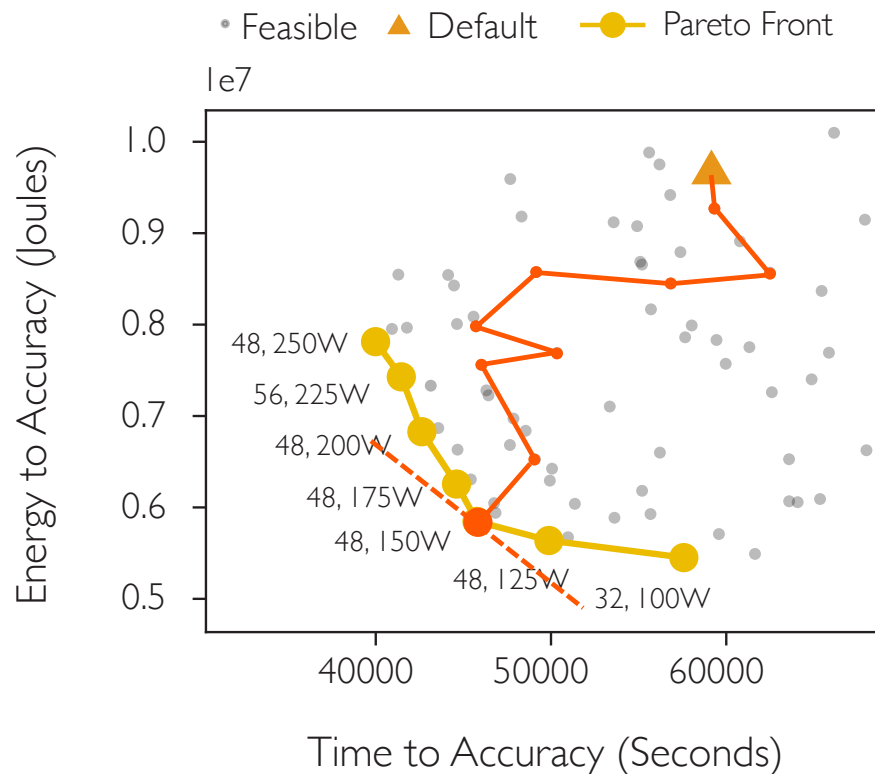
Time vs. Energy Trade-off



Which yellow point is the best?

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100.
Similar trends found across four GPU generations.

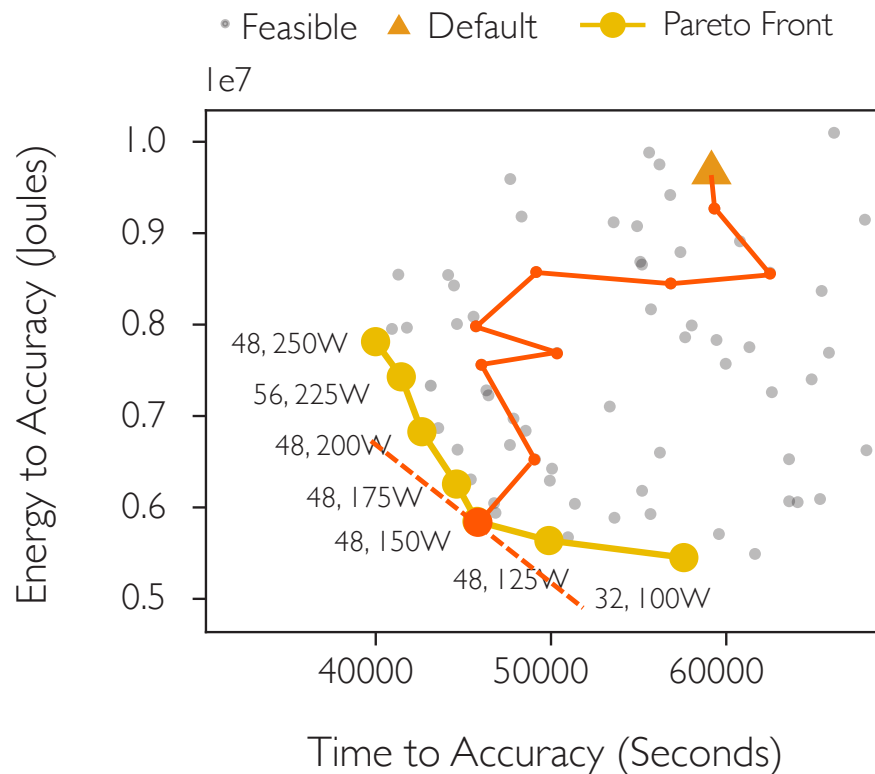
Multi-Armed Bandit Algorithm



- Objective
= Linear combination of time & energy
- Arm = Batch size
- Horizon = Recurring training
- Thompson sampling

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100.
Similar trends found across four GPU generations.

Multi-Armed Bandit Algorithm

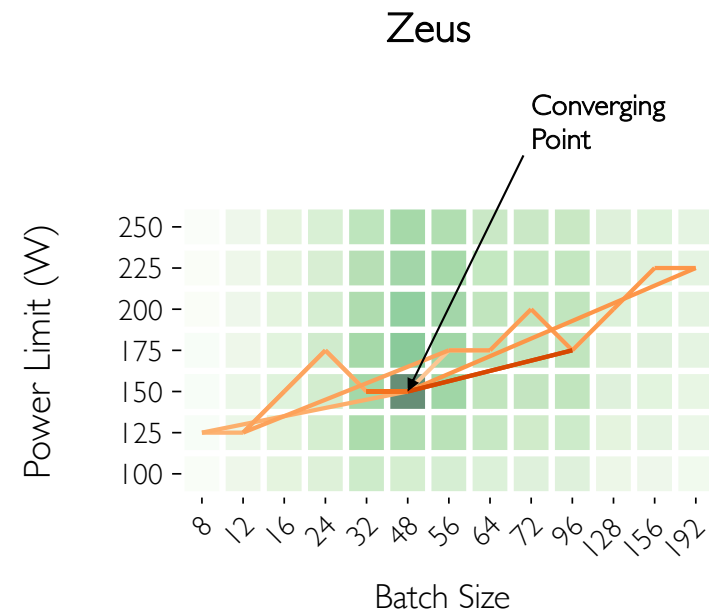
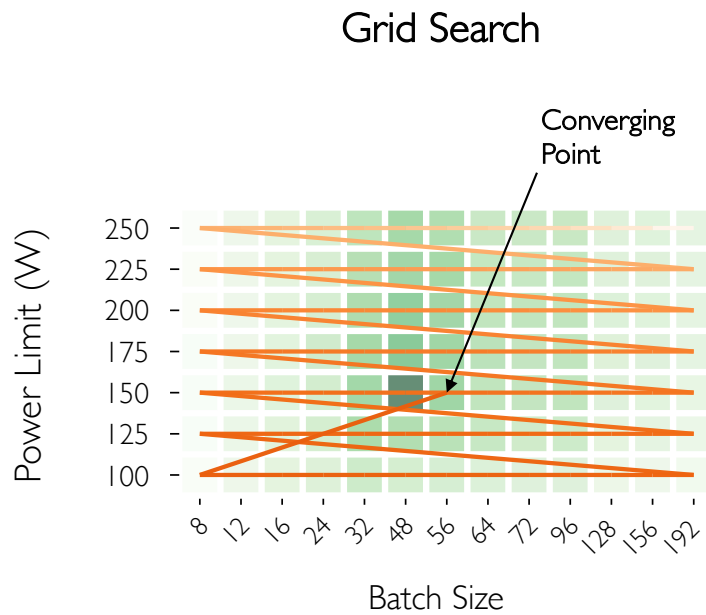


15% to 76% energy reduction

across diverse models
and multiple GPU generations

Results from training DeepSpeech2 on LibriSpeech on an NVIDIA V100.
Similar trends found across four GPU generations.

Zeus in Action



Search Path Training Cost (darker means better)

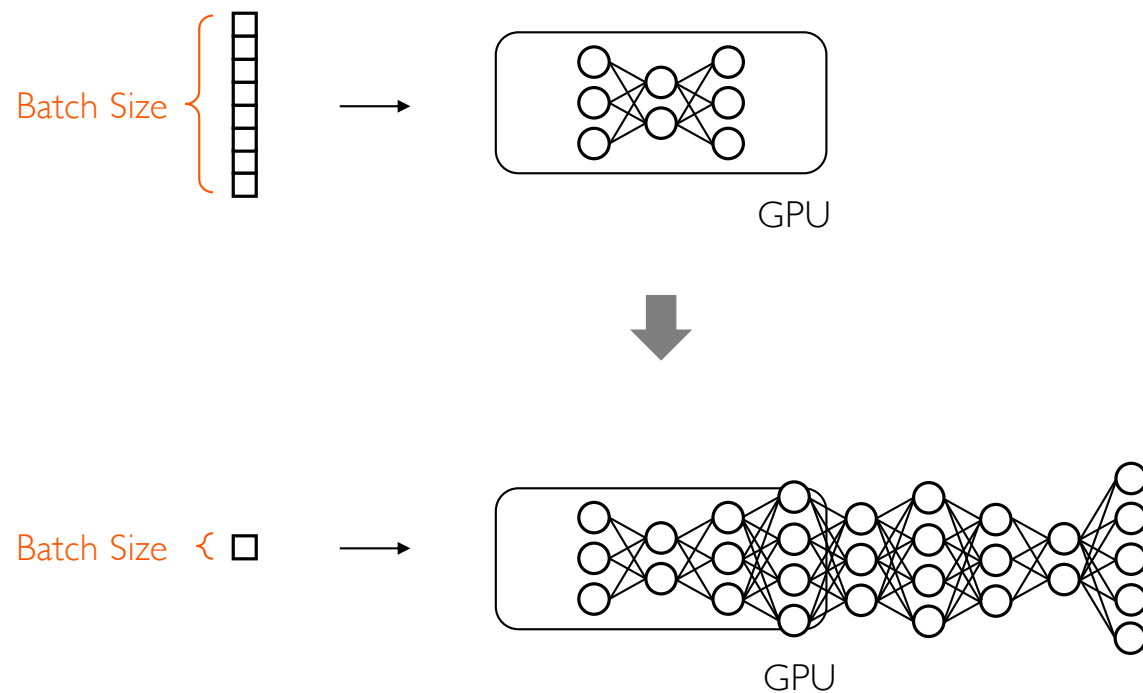
DeepSpeech2 trained on LibriSpeech on an NVIDIA V100 GPU.

Reducing Energy Bloat in Large Model Training

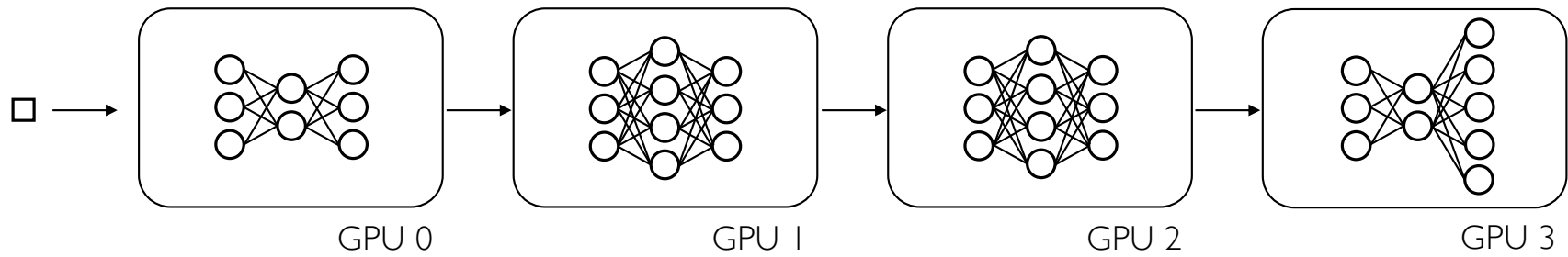
*Jae-Won Chung, Yile Gu, Insu Jang,
Luoxi Meng, Nikhil Bansal, Mosharaf Chowdhury*

*“Any way to reduce
energy consumption
without slowdown?”*

Explosive Growth in Model Size



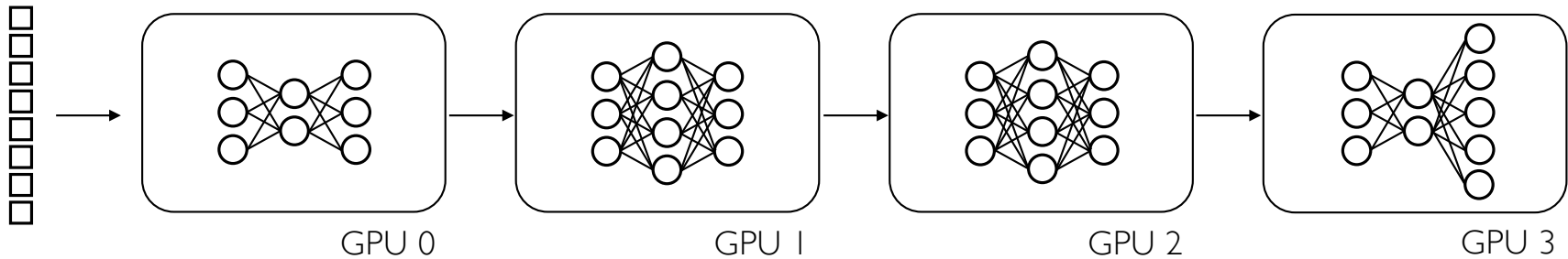
Model Parallelism



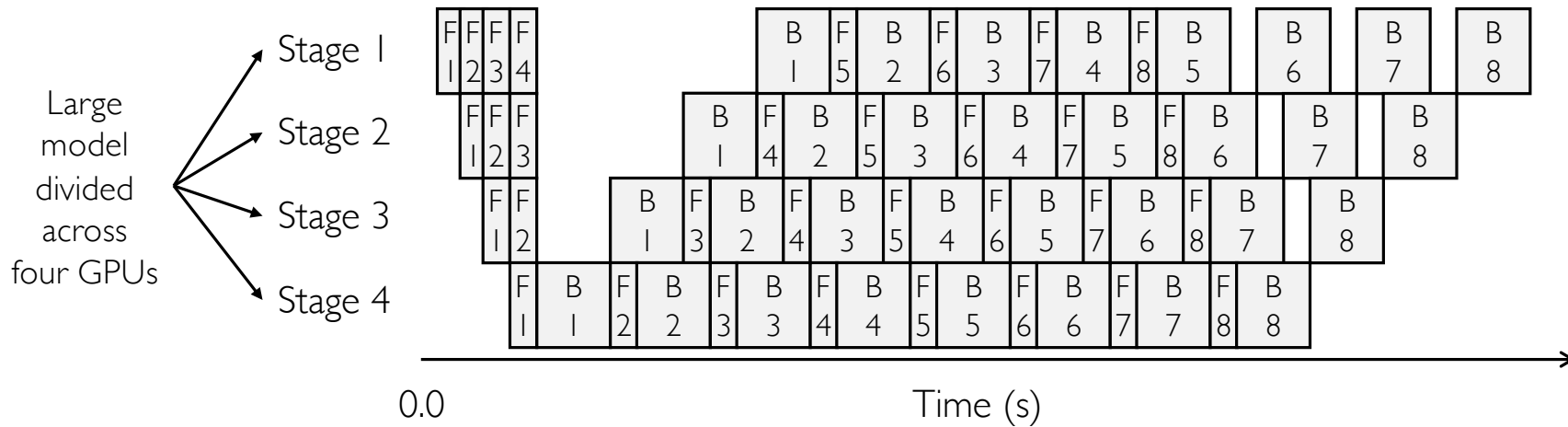
Model Parallelism

Pipeline Parallel training

8 microbatches

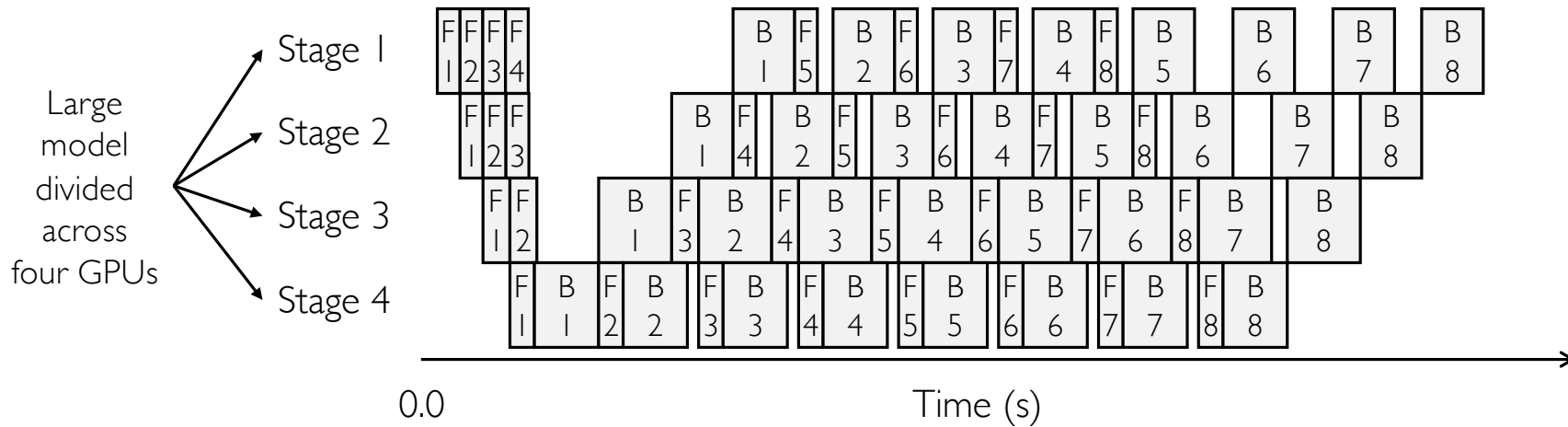


Pipeline Parallel Training



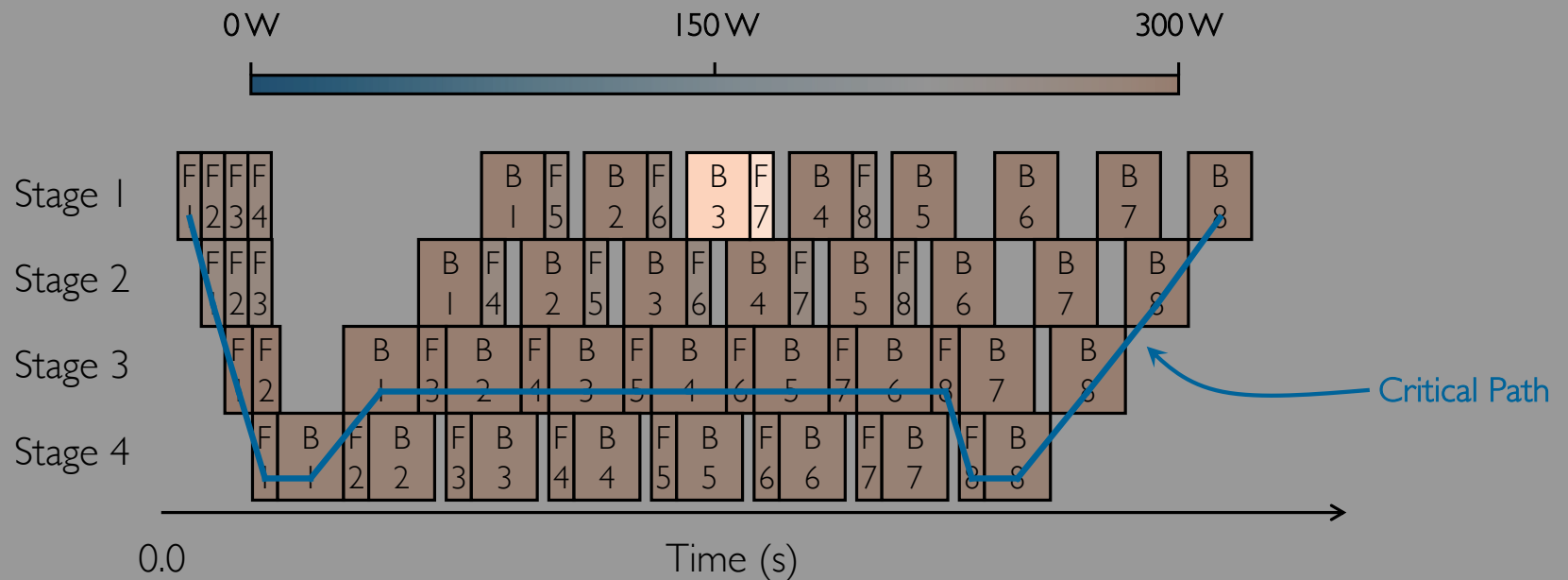
One training iteration with 4 pipeline stages and 8 microbatches (IF1B schedule).

Fundamental Computation Imbalance



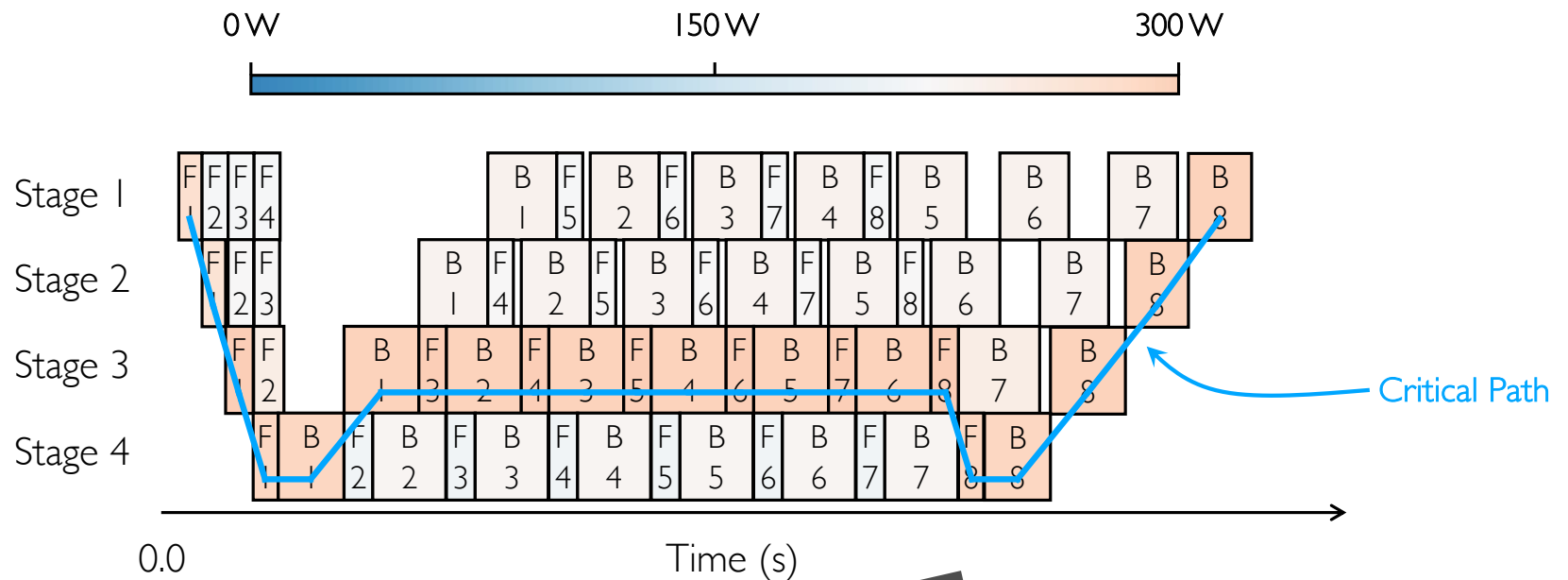
One training iteration with 4 pipeline stages and 8 microbatches (IF1B schedule).
Drawn to scale for GPT-3 1.3B on NVIDIA A100 GPUs.

Where Do the Joules Go?



One training iteration with 4 pipeline stages and 8 microbatches (IFIB schedule).
Drawn to scale for GPT-3 1.3B on NVIDIA A100 GPUs.

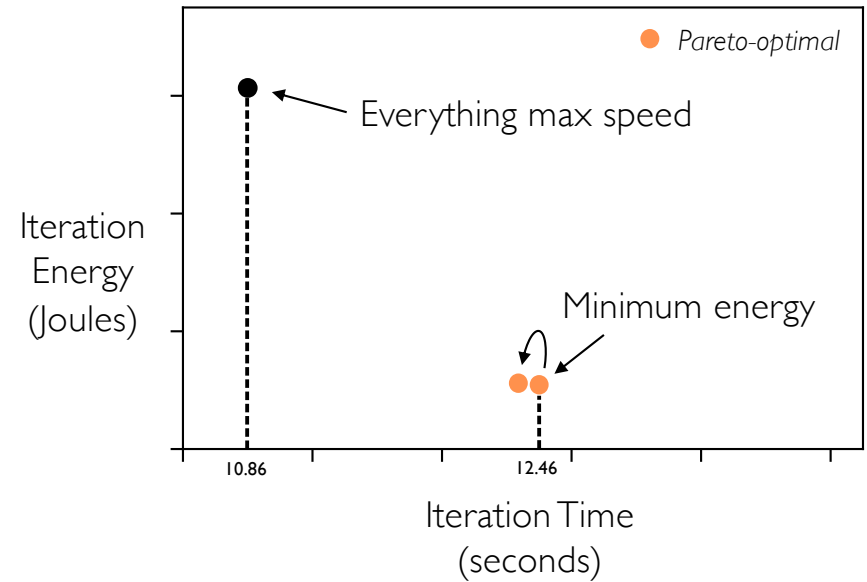
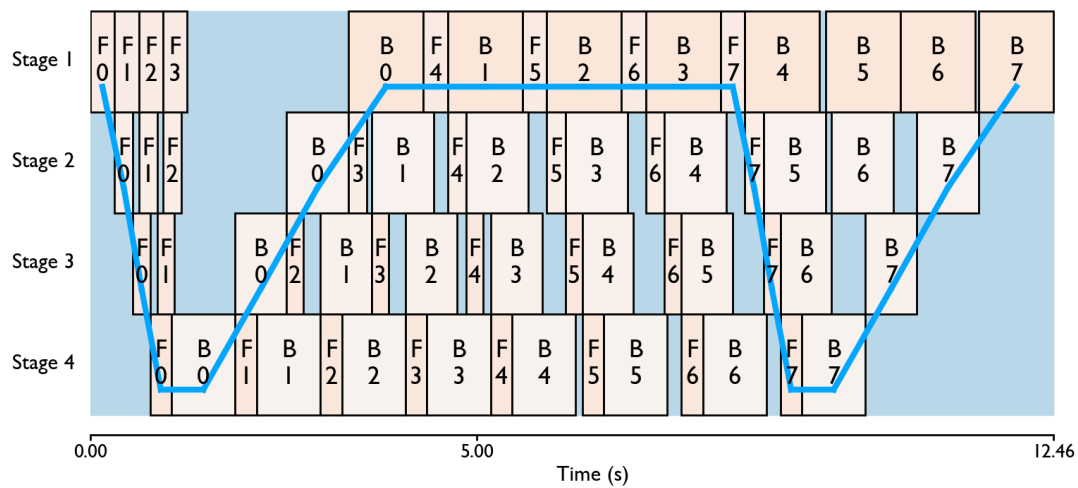
Cutting 30% Energy Bloat



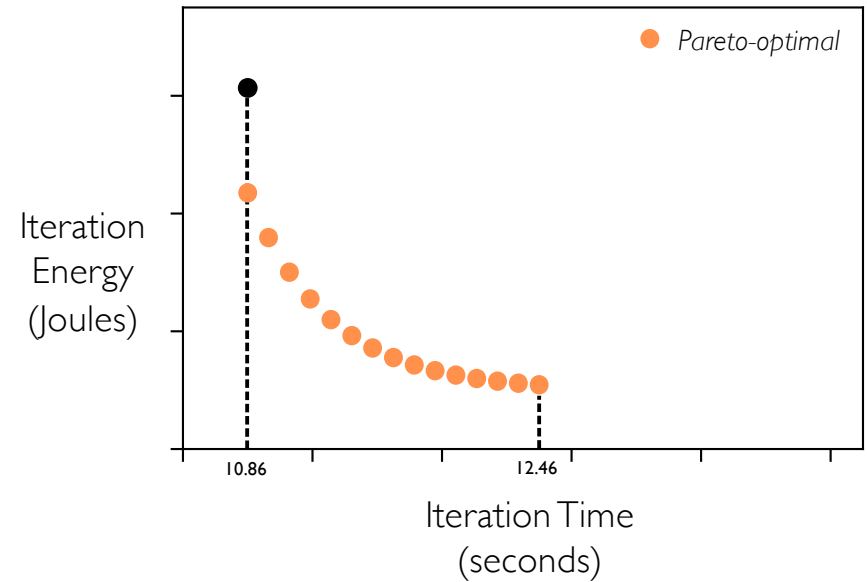
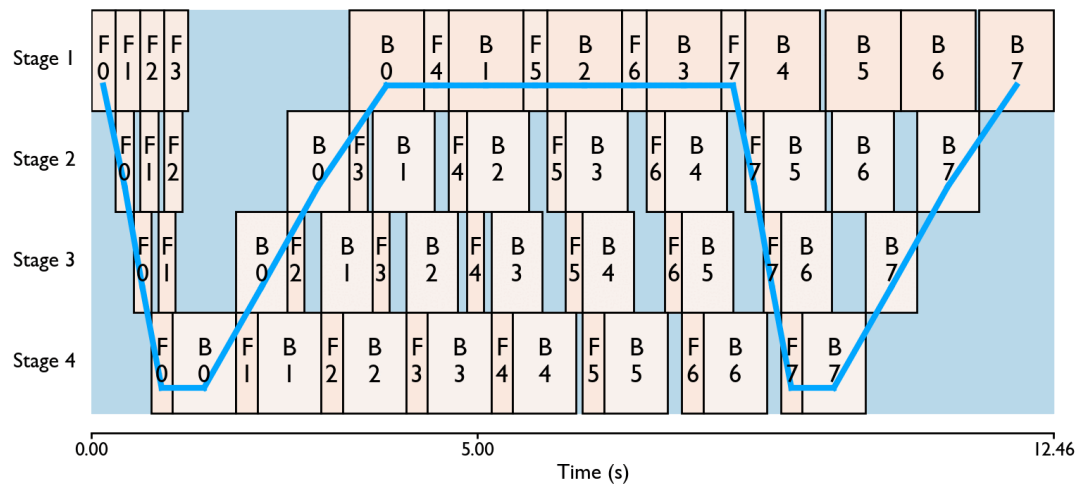
One training iteration of GPT-1. B and F denote backpropagation and forward passes, respectively. The chart shows the execution of the training iteration across four pipeline stages and eight microbatches on NVL A100 GPUs, drawn to scale.

NP-Hard

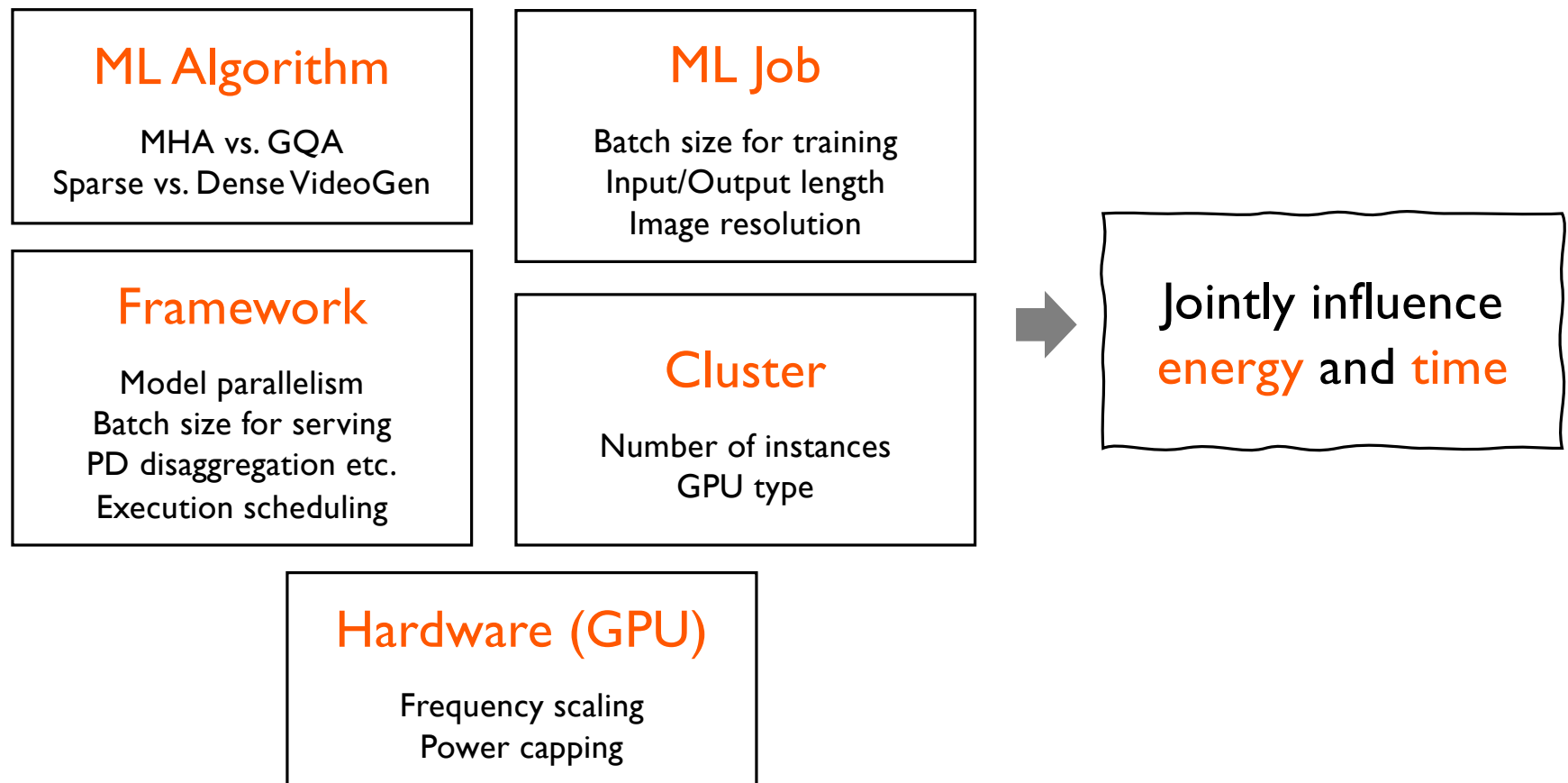
Time vs. Energy Trade-off



Time vs. Energy Trade-off



Summary: Decisions Across the AI Stack



Summary

- Optimize energy along the time–energy trade-off frontier.
- Leverage available latency slack for energy savings.

