



# Modal

## Recycling the World Computer: Fault-Tolerant LLM Training on Idle GPU Capacity

Benjamin Cowen, Jason Mancuso\*, Shariq Mobin

# Agenda

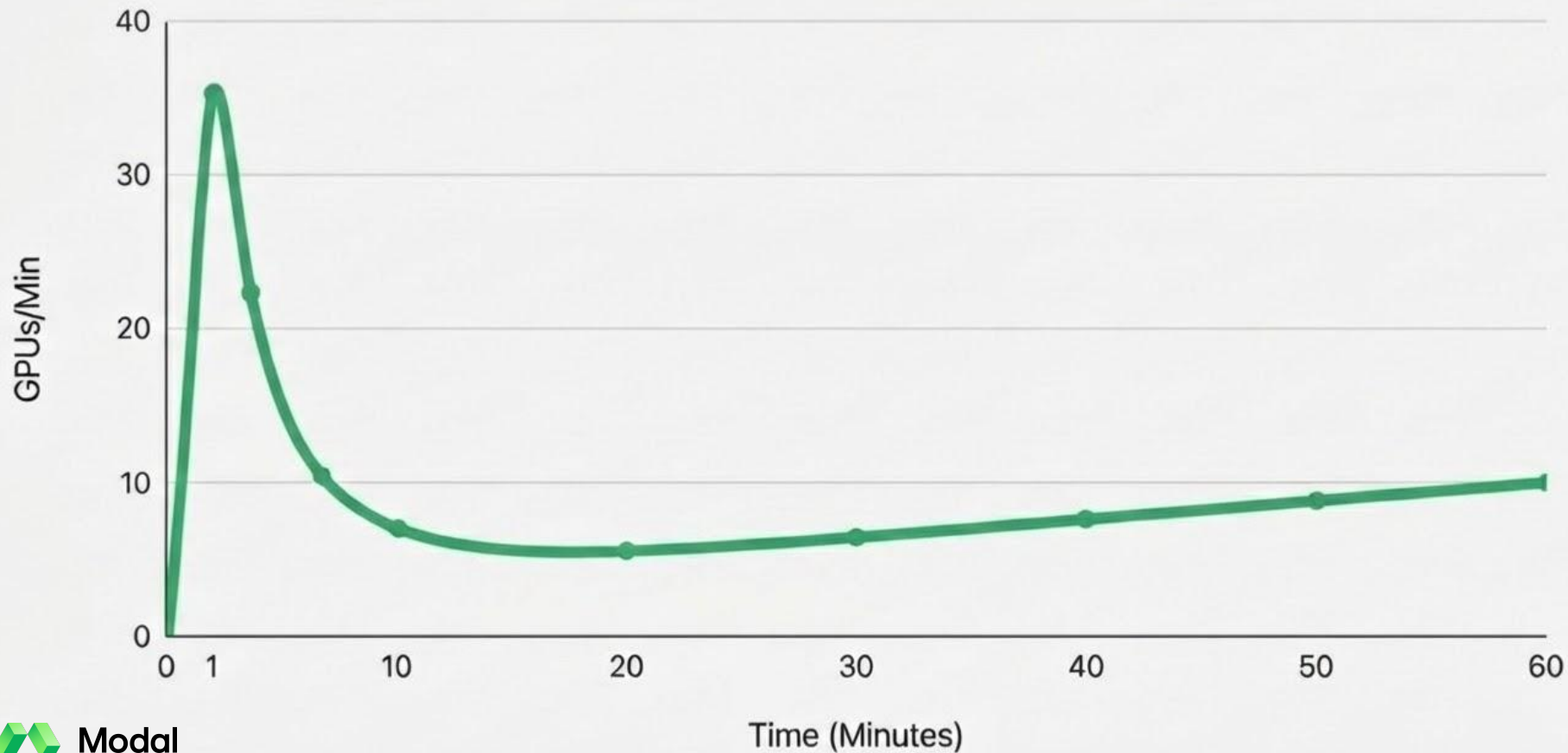
1. What is Modal?
2. A crash course on Modal's core tech
  - a. Linear programming for optimal fleet management
  - b. Why the fleet will always have idle GPUs
3. How to recycle idle GPUs in a multi-cloud, multi-region fleet
4. Which tasks are most suitable for recycled compute?
  - a. Massively parallel + fault-tolerant + reentrant functions with no/low inter-replica communication
5. Pushing the limits of recycled compute toward large-scale LLM training

demo.py

```
1 import os
2 import subprocess
3 import time
4
5 import modal
6 ⚡
7 app = modal.App("h100-scale-test")
8
9
10 @app.function(gpu=["H100", "H200", "B200", "A100", "A100-80GB",
11 def run_on_gpu(i: int):
12     cloud, region, gpu_type = get_container_info()
13     print(
14         f"Running input {i} in {region} on {cloud} "
15         f"with device {gpu_type}..."
16     )
17     time.sleep(60 * 5)
18
19
20 @app.local_entrypoint()
21 def main():
22     _results = list(run_on_gpu.map(range(10_000)))
23
24
25 def get_container_info():
26     cloud, region = os.getenv("MODAL_CLOUD_PROVIDER"), os.getenv('
27
28     # Determine GPU type
29     try:
30         result = subprocess.run(
```

neurips-talk is  v0.1.0 via  v3.13.7 (neurips-talk) on  (us-east-1)  
on  jason.mancuso@modal.com  
> |

## GPU Acquisition Rate Over Time (per Function)



What's our secret?

# What's our secret?

- Near instant cold-starts

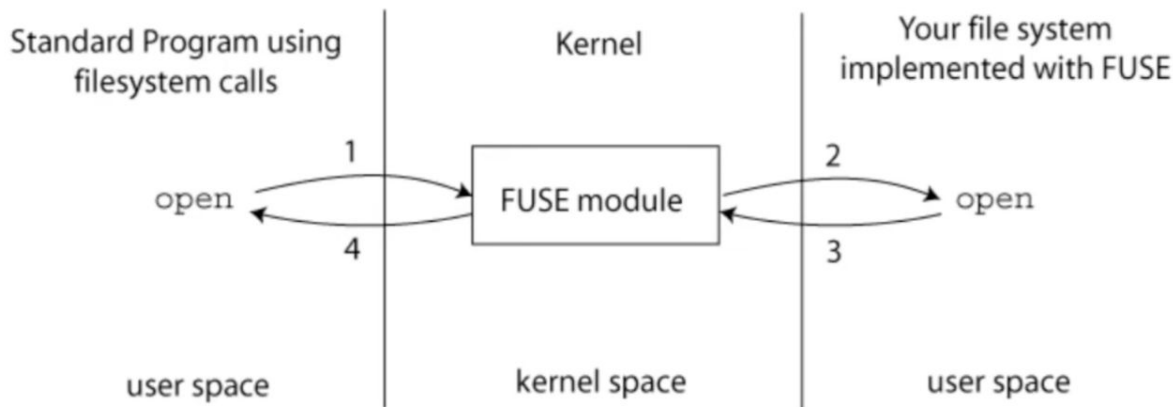
# What's our secret?

- Near instant cold-starts
  - Custom container runtime, built to start from lazy-loaded, content-addressed images



# Lazy loading approach

## Cold-starting our fat Python container



Speaker

Jonathan Belotti

Presentation Title

Fast, Lazy Container  
Loading

Venue Host

**TRAIL**

Fast, lazy container loading on Modal, Jonathon Belotti (via [YouTube](#))

# What's our secret?

- Near instant cold-starts
  - Custom container runtime, built to start from lazy-loaded, content-addressed images

# What's our secret?

- Near instant cold-starts
  - Custom container runtime, built to start from lazy-loaded, content-addressed images
- Large, elastic, global compute pool

# What's our secret?

- Near instant cold-starts
  - Custom container runtime, built to start from lazy-loaded, content-addressed images
- Large, elastic, global compute pool
  - Multi-cloud worker fleet that scales dynamically in response to user workload demand

# What's our secret?

- Near instant cold-starts
  - Custom container runtime, built to start from lazy-loaded, content-addressed images
- Large, elastic, global compute pool
  - Multi-cloud worker fleet that scales dynamically in response to user workload demand

Multi-cloud Substrate



Neoclouds

# What's our secret?

- Near instant cold-starts
  - Custom container runtime, built to start from lazy-loaded, content-addressed images
- Large, elastic, global compute pool
  - Multi-cloud worker fleet that scales dynamically in response to user workload demand

Managing our fleet is not only a logistics problem, but also an arbitrage opportunity.

Multi-cloud Substrate



Neoclouds

**Parameters:**

$A :=$  total GPUs requested by users

$B :=$  count of GPUs in the buffer

$I = \{I_1, I_2, \dots, I_n\} :=$  possible instance types

$\mathbf{C} = [C_1, \dots, C_n] :=$  cost of instance types

$\mathbf{L} = [L_1, \dots, L_n] :=$  scaling limits for each type

**Output:**

$\mathbf{X} = [x_1, \dots, x_n] :=$  instances of each type to launch

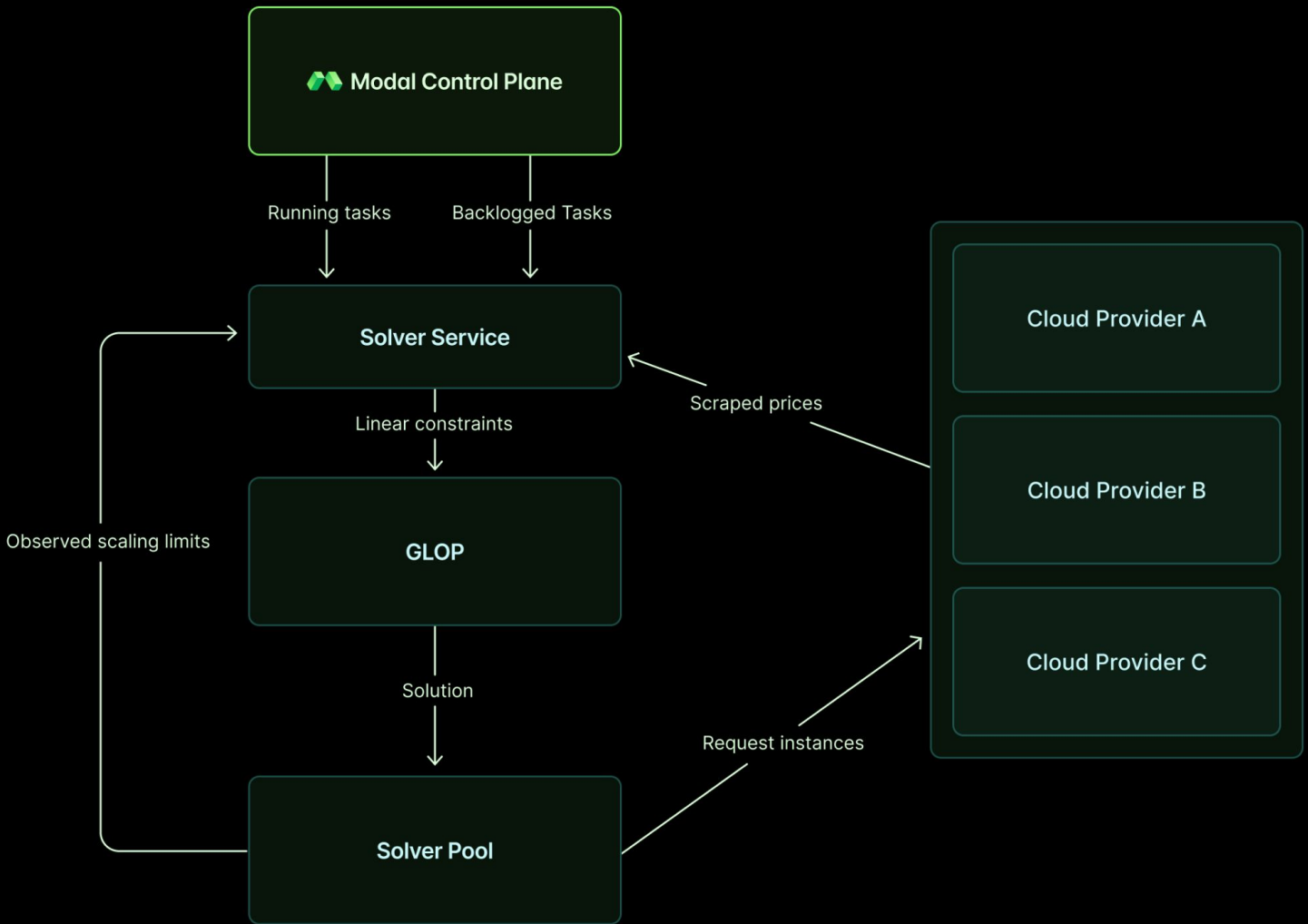
**Objective:**

$$\text{minimize } \mathbf{C}^T \mathbf{X} = \sum_{i=1}^n x_i C_i$$

**Constraints:**

$$8 \sum_{i=1}^n x_i \geq A + B$$

$$x_i \leq L_i \quad (i = 1, \dots, n)$$



# Sources of Idle Compute in the Fleet

# Sources of Idle Compute in the Fleet

- Buffer slots on active workers

# Sources of Idle Compute in the Fleet

- Buffer slots on active workers
  - These are intentional. We explicitly overprovision so that we can accommodate bursts in our user's workloads.

# Sources of Idle Compute in the Fleet

- Buffer slots on active workers
  - These are intentional. We explicitly overprovision so that we can accommodate bursts in our user's workloads.
- Idle slots on draining workers

# Sources of Idle Compute in the Fleet

- Buffer slots on active workers
  - These are intentional. We explicitly overprovision so that we can accommodate bursts in our user's workloads.
- Idle slots on draining workers
  - When the solver decides to prune a worker, we can't move it over until user tasks have fresh containers elsewhere in the fleet.
  - Unused GPUs on draining workers can't accept new tasks, so they remain idle for the duration.

## Recycled Compute

Running background tasks on these “idle” GPUs, whether in the buffer or in the drain.

For background tasks in the buffer, we need to be able to SIGKILL in ~10s or 100s of milliseconds.

## How much Recycled Compute is available?

<b>GPU Types (Inclusive)</b>	<b>Min VRAM</b>	<b>Estimated* Idle Compute</b>	<b>% Draining</b>
B200,H200	141GB+	100k GPU-hrs/month	10%+
H100, A100-80 + above	80GB+	300k GPU-hrs/month	10%+
A100-40, L40S + above	40GB+	500k GPU-hrs/month	20%+
A10G, L4 + above	24GB+	1.2M GPU-hrs/month	30%+

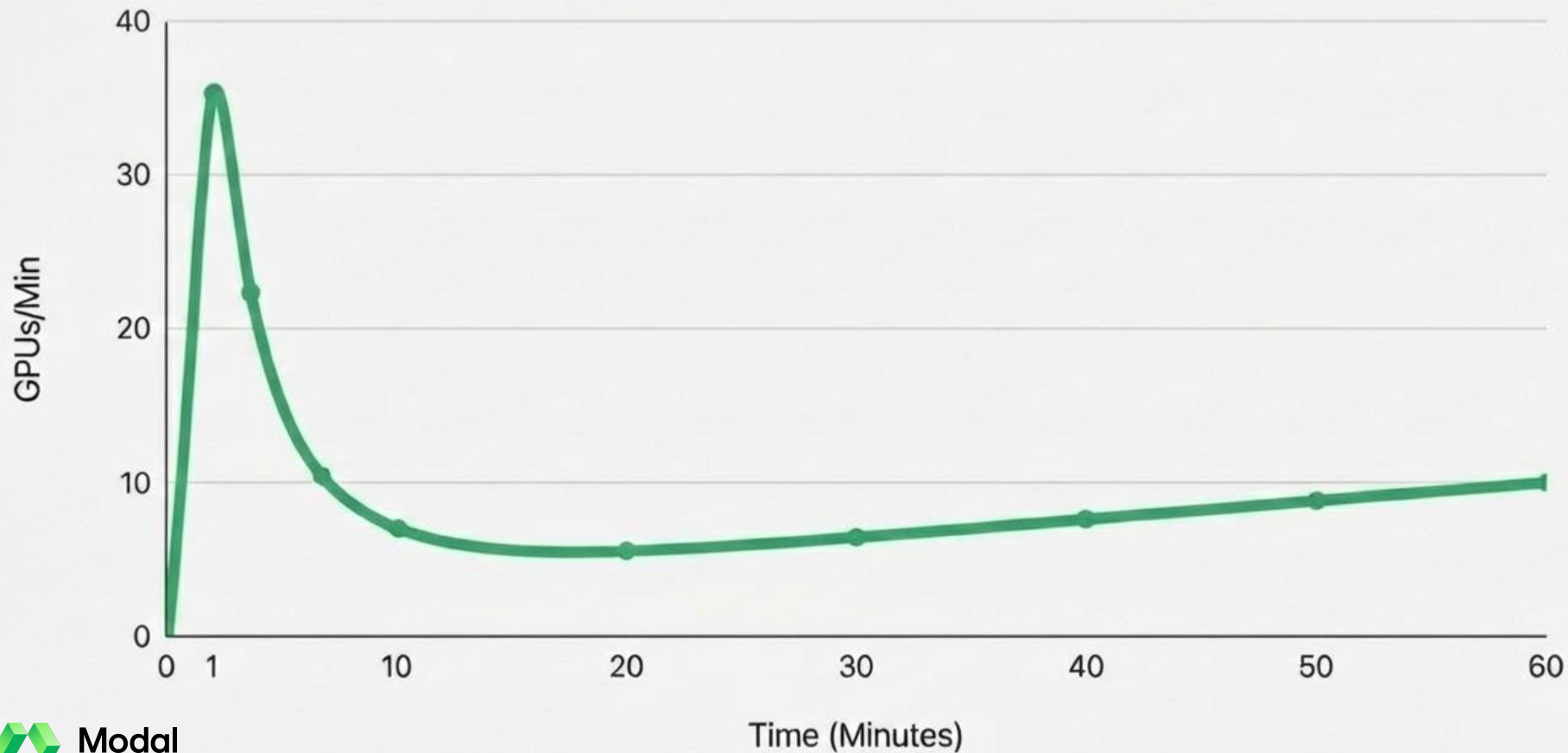
# How much Recycled Compute is available?

GPU Types (Inclusive)	Min VRAM	Estimated* Idle Compute	% Draining
B200,H200	141GB+	100k GPU-hrs/month	10%+
H100, A100-80 + above	80GB+	300k GPU-hrs/month	10%+
A100-40, L40S + above	40GB+	500k GPU-hrs/month	20%+
A10G, L4 + above	24GB+	1.2M GPU-hrs/month	30%+

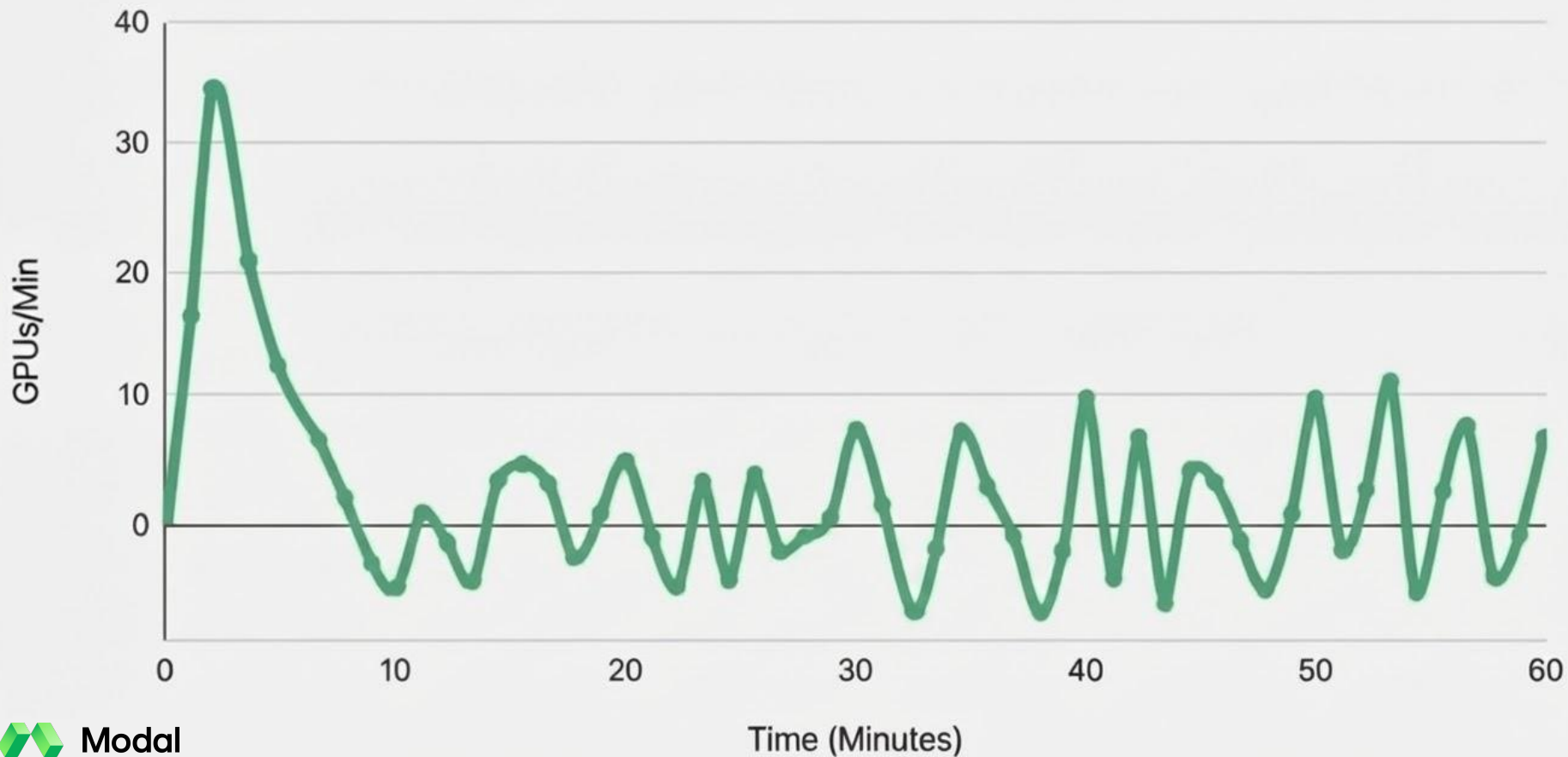
Training Costs	Pre-Training	Context Extension	Post-Training	Total
in H800 GPU Hours	2664K	119K	5K	2788K
in USD	\$5.328M	\$0.238M	\$0.01M	\$5.576M

Table 1 | Training costs of DeepSeek-V3, assuming the rental price of H800 is \$2 per GPU hour.

## GPU Acquisition Rate Over Time (per Function)



## GPU Acquisition Rate Over Time (per **Recycled** Function)






















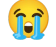
# The Shape of Recycled Compute

- Massively parallel
- Skewed hardware heterogeneity
- High availability in aggregate, but low for any one device
  - High preemption likelihood
- High variance interconnect
  - Inter-container network could be low-latency, high-bandwidth if within single AZ
    - But can't rely on Infiniband/NVLink/RDMA :(
  - Could be low-latency if within single geographic region
  - Could also be in completely different regions, communicating over unreliable WAN

# The Shape of *Ideal Tasks* for Recycled Compute

- Massively parallel
- Hardware-agnostic
- Reentrant
  - Easy to checkpoint + resume on fresh container
- Parallel-wise fault-tolerant
  - Can still make progress if many individual replicas are preempted
- Little to no inter-replica communication

# Ideal Tasks for Recycled Compute

	Hardware agnostic?	Reentrant?	Fault-tolerant parallelism?	Low comm?
Classical protein folding				
Molecular physics simulation				
Offline batch inference				
Synthetic data (agentic rollouts)				
Distributed LLM training				

Low-Communication Distributed Training, aka Federated Learning



# DiLoCo: Distributed Low-Communication Training of Language Models

Arthur Douillard<sup>1</sup>, Qixuan Feng<sup>1</sup>, Andrei A. Rusu<sup>1</sup>, Rachita Chhaparia<sup>1</sup>, Yani Donchev<sup>1</sup>, Adhiguna Kuncoro<sup>1</sup>,  
Marc'Aurelio Ranzato<sup>1</sup>, Arthur Szlam<sup>1</sup> and Jiajun Shen<sup>1</sup>

<sup>1</sup>Google DeepMind

arxiv: 2311.08105

# Low-Communication Distributed Training, aka Federated Learning



## Streaming DiLoCo with overlapping communication: Towards a Distributed Free Lunch 🥪

Arthur Douillard<sup>\*,1</sup>, Yanislav Donchev<sup>\*,1</sup>, Keith Rush<sup>2</sup>, Satyen Kale<sup>†,2</sup>, Zachary Charles<sup>2</sup>, Zachary Garrett<sup>2</sup>, Gabriel Teston<sup>3</sup>, Dave Lacey<sup>1</sup>, Ross McIlroy<sup>1</sup>, Jiajun Shen<sup>1</sup>, Alexandre Ramé<sup>1</sup>, Arthur Szlam<sup>1</sup>, Marc'Aurelio Ranzato<sup>1</sup> and Paul Barham<sup>1</sup>

<sup>1</sup>Google DeepMind, <sup>2</sup>Google Research, <sup>3</sup>Google, <sup>\*</sup>Equal core contributions, <sup>†</sup>Currently at Apple.

Low-Communication Distributed Training, aka Federated Learning

Communication-Efficient Language Model Training Scales Reliably  
and Robustly:

## Scaling Laws for DiLoCo

Zachary Charles<sup>1\*</sup>   Gabriel Teston<sup>2</sup>   Lucio Dery<sup>3</sup>   Keith Rush<sup>1</sup>  
Nova Fallen<sup>1</sup>   Zachary Garrett<sup>1</sup>   Arthur Szlam<sup>3</sup>   Arthur Douillard<sup>3</sup>

<sup>1</sup>Google Research   <sup>2</sup>Google Search   <sup>3</sup>Google DeepMind

# Low-Communication Distributed Training, aka Federated Learning

## Communication-Efficient Language Model Training Scales Reliably

**Harder:** DiLoCo's hyperparameters are robust and predictable across model scales.

**Better:** DiLoCo further improves over data-parallel training as model size increases.

**Faster:** DiLoCo uses orders of magnitude less bandwidth than data-parallel training.

**Stronger:** DiLoCo tolerates a significantly larger batch size than data-parallel training.

But what about fault-tolerance?!

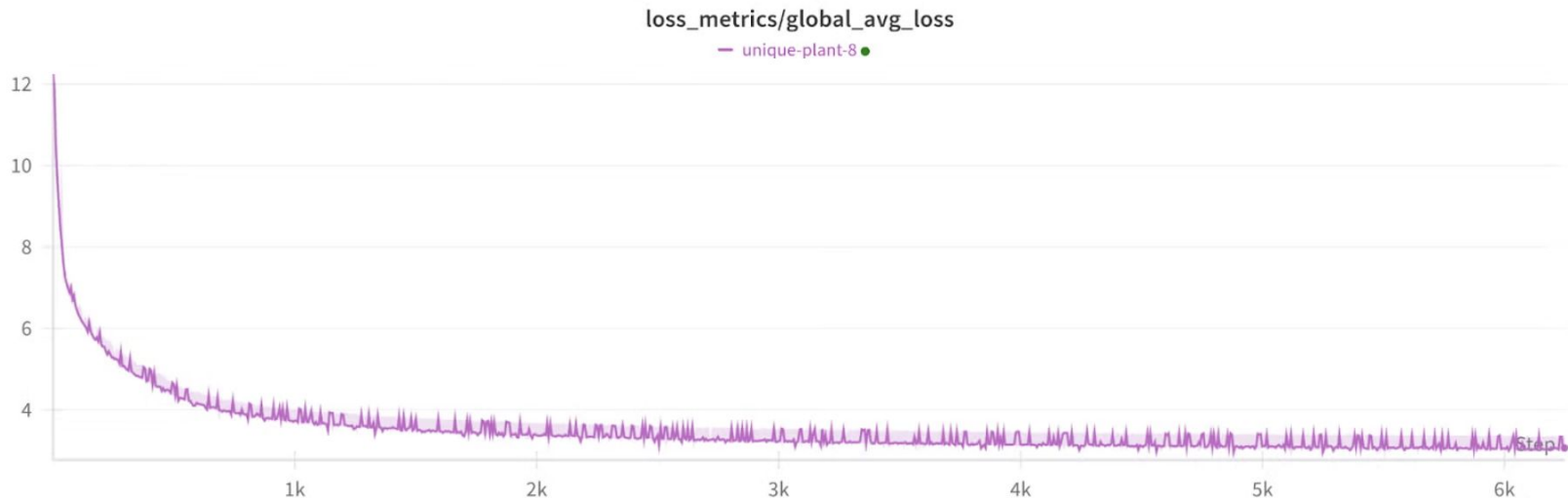
Blog

# Fault Tolerant Llama: training with 2000 synthetic failures every ~15 seconds and no checkpoints on Crusoe L40S

By Tristan Rice, Howard Huang | June 20, 2025

Via [pytorch.org](https://pytorch.org)

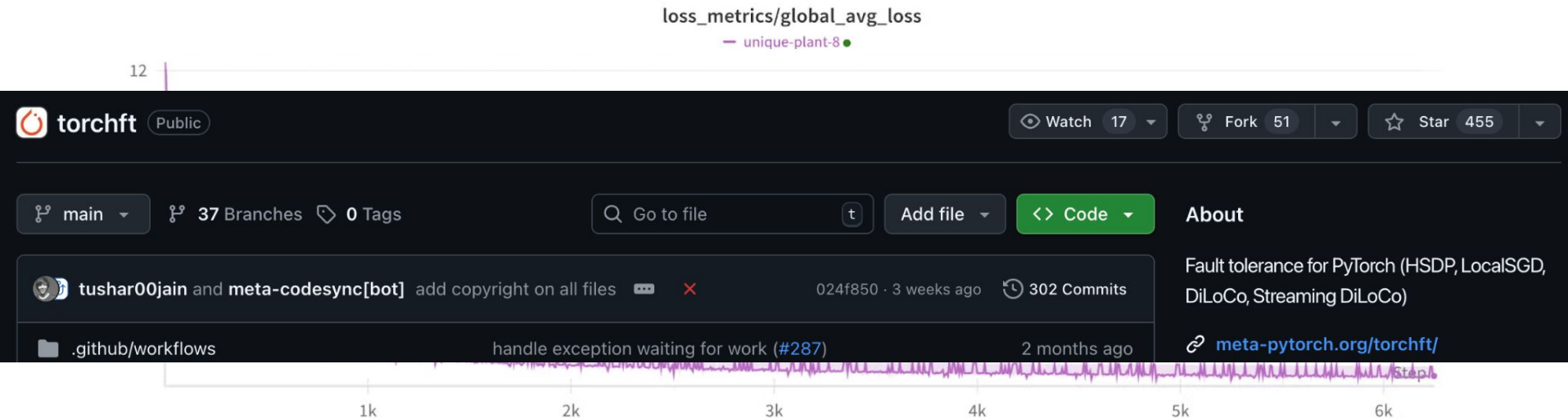
# But what about fault-tolerance?!



Training loss across 1200 failures with no checkpoints.

*NOTE: Each small spike is a non-participating worker recovering which affects the metrics but not the model*

# But what about fault-tolerance?!



Training loss across 1200 failures with no checkpoints.

*NOTE: Each small spike is a non-participating worker recovering which affects the metrics but not the model*

Don't Try This at Home

# Don't Try This at Home

TorchFT ProcessGroup issues

# Don't Try This at Home

TorchFT ProcessGroup issues

Explosion of egress costs

→ Need to customize collective communication ops to be geography-aware

# Don't Try This at Home

TorchFT ProcessGroup issues

- We use Tunnels for inter-container networking
- ProcessGroup requires

Explosion of egress costs

- Need to customize collective communication ops to be geography-aware

So, back to inference and evolutionary algorithms we go!

Thanks!

We're hiring: <https://modal.jobs>, or stop by Kiosk 9 in Expo Hall B!

\$30 free Modal credits: <https://modal.com/signup>



**Modal**