# GPEmu: A GPU Emulator for Faster and Cheaper Prototyping and Evaluation of Deep Learning System Research

Meng Wang, Gus Waldspurger, Naufal Ananda, Yuyang Huang,
Kemas Wiharja, John Bent, Swaminathan Sundararaman,
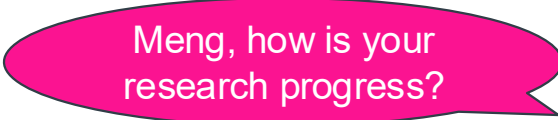Vijay Chidambaram, Haryadi S. Gunawi
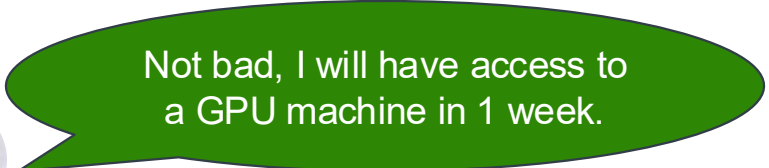
# DL System Research is Hot

❑ Deep Learning (DL) is widely used

❑ Training is time-consuming and inference is latency-sensitive
   ▪ It's important to optimize deep learning system performance

❑ Lots of deep learning system research
   ▪ Data loading
   ▪ Preprocessing
   ▪ Job Scheduling
   ▪ Network communication
   ▪ …

# But GPUs Are Expensive
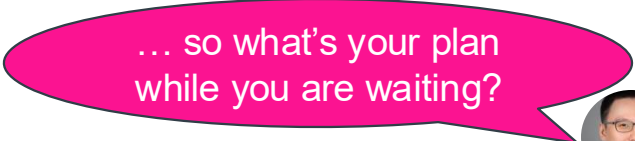
❑ GPUs have limited availability in research cloud
  - GPU machines on *Chameleon cloud* require reservations <u>weeks in advance.</u>

❑ GPUs on commercial clouds are expensive
  - *Cloudbank* provides up to $5,000 per 6-months for system-focused research
  - Can reserve 4 AWS p3.16xlarge instances *<u>only for 51 hours</u>*

Meng, how is your research progress?

Not bad, I will have access to a GPU machine in 1 week.

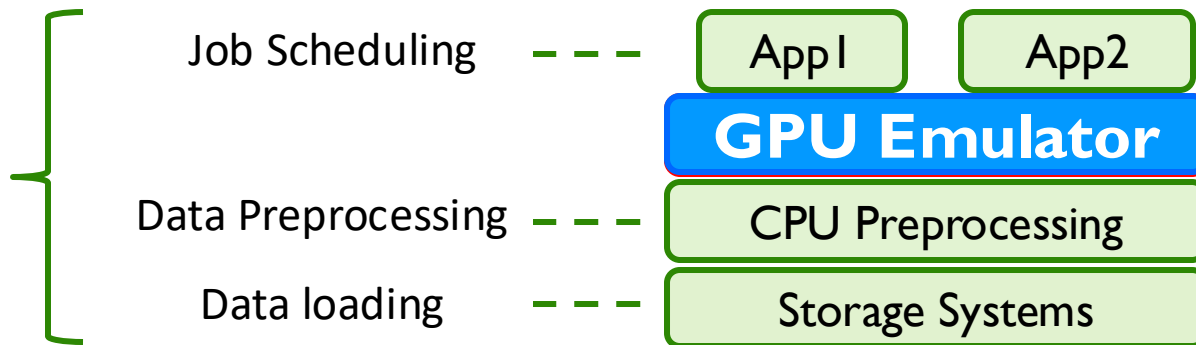… so what's your plan while you are waiting?

# GPU is NOT Always Needed

**GPU is NOT NEEDED when...**

DL **System** Stack

Working on **the layers above the GPU** to increase GPU utilization

Job Scheduling — — — | App1 | App2 |

**GPU Emulator**

Data Preprocessing — — — | CPU Preprocessing |
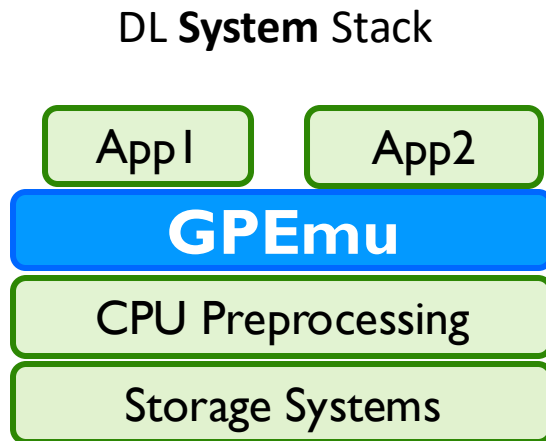
Data loading — — — | Storage Systems |

Only **GPU Performance** is needed

# GPEmu Overview

**A GPU Emulator for Faster** and **Cheaper** prototyping and evaluation of DL system research.

DL **System** Stack

**+** Without GPUs

**+** 30+ models and 6 GPUs

**+** Easy to use

**+** Support both Single node & Distributed setups

App1    App2

**GPEmu**

CPU Preprocessing

Storage Systems

- 💡 **Time Emulation**
- 💡 **Memory Emulation**
- 💡 **Distributed Support**
- 💡 **Sharing Support**

- **Effectively reveal** DL system **bottlenecks** during emulations
- **Quickly show** the **benefits** of new system optimizations across a comprehensive spectrum of research areas

THE UNIVERSITY OF
CHICAGO

# GPEmu Features

**Time Emulation**

**Sleep-based approach**

(a) CDF of batch compute time
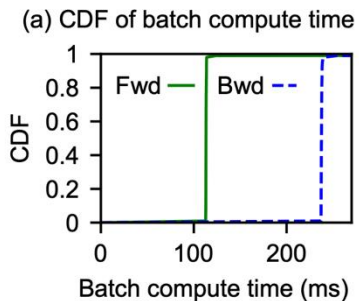


Profile its GPU time
$T_{compute}$

Consistent with the same setup

model.compute() → sleep($T_{compute}$)

+ GPU Compute Time

+ Host-to-GPU Transfer Time

+ GPU-driven Preprocessing Time

+ 30+ models

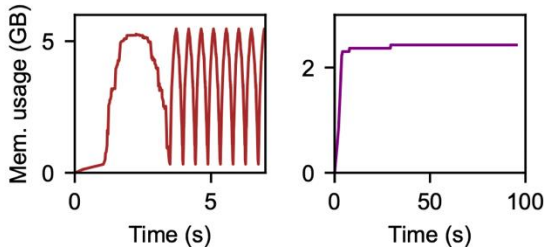+ 6 GPUs

# GPEmu Features

**💡 Memory Emulation**

## 1. Emulate GPU Memory Usage

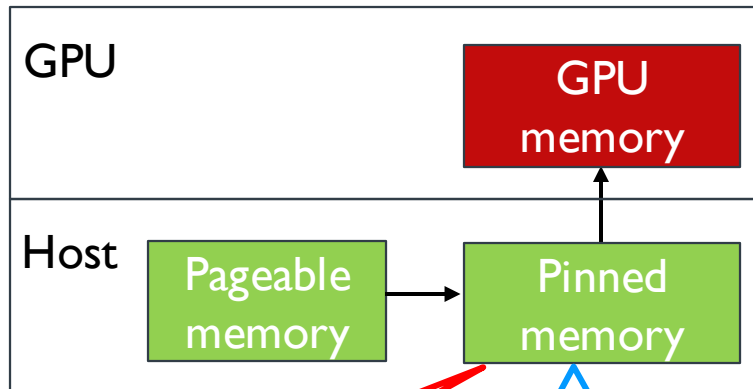**+** By GPU Compute

**+** By GPU-driven preprocessing

(a) RN50 GPU Mem. Usage    (b) DALI Prep. GPU Mem

Emulate whether a DL workload can fit into a GPU's memory

## 2. Emulate Pinned Memory

**GPU** | GPU memory

**Host** | Pageable memory → Pinned memory

Was managed by **CUDA**

Now managed by **our own manager**

# GPEmu Features



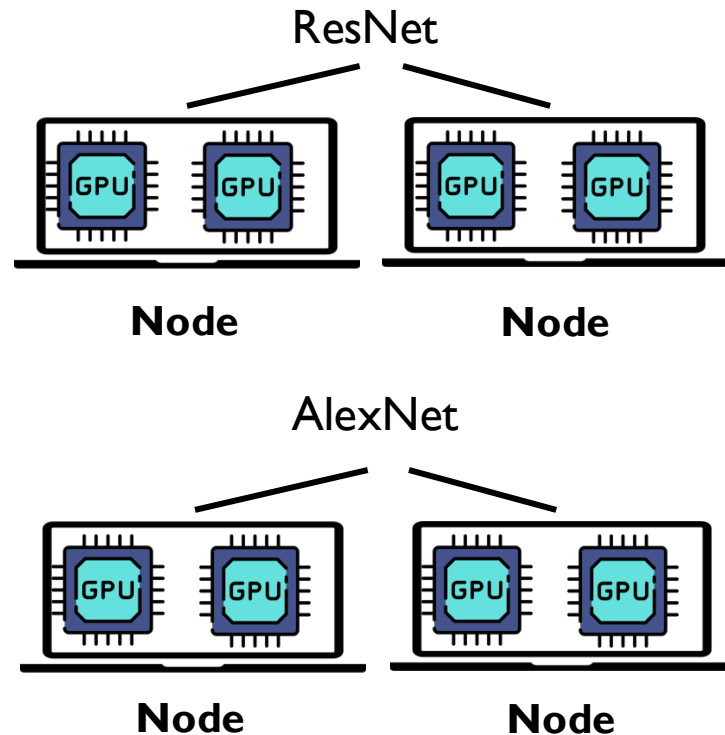**Cluster**

ResNet

Node          Node

AlexNet

Node          Node

> 💡 **Distributed Support**

**+ Multi-GPU** Training

**+ Multi-Node** Training

- Support Distributed data-parallel training

**+ Multi-Job** Scheduling

- emuGPU device plugin
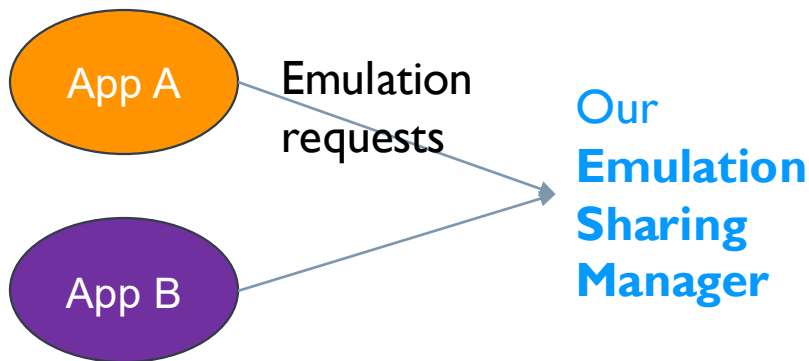- Evaluate scheduling jobs wrapped in GPEmu containers

# GPEmu Features

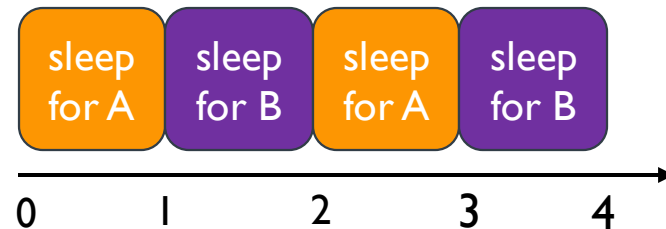**Sharing Support**

**+ Support single-node**

**+ Support K8S**

App A

App B

Emulation requests

Our **Emulation Sharing Manager**

## GPU Time Sharing

Model on GPU

| A | B | A | B |

0    1    2    3    4

| sleep for A | sleep for B | sleep for A | sleep for B |

0    1    2    3    4

# Supported Research

+ **Data stall analysis**

+ **Preprocessing disaggregation**

+ **Data loader optimization**

+ **Distributed training optimization**

+ **GPU scheduling**

+ **GPU sharing**

**Reproduced 9 papers**

DataStall [49], VLDB '21
TF-DS [22], SoCC '23
FastFlow [60], VLDB '23
FFCV [42], CVPR '23
LADL [66], HiPC '19
Synergy [48], OSDI '22
Allox [41], EuroSys '20
Salus [67], MLSys '20
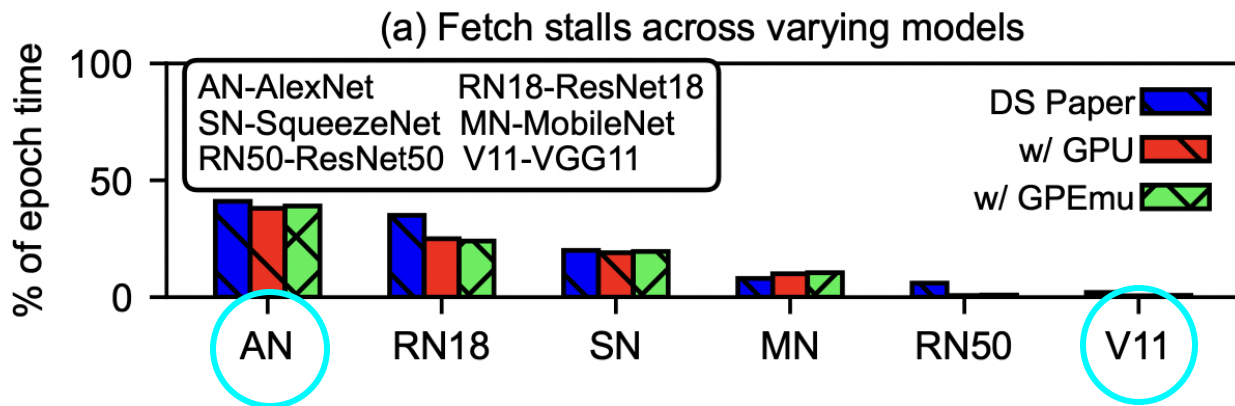Muri [74], SIGCOMM '22

**Prototyped 3 new optimizations**

Small file first caching
Async data loading
File grouping

# Case Study 1: Data Stall Analysis

❑  DataStall @ VLDB 2022

❑  Fetch Stall: time spent waiting for data fetching

**Pattern is consistent**: Different models have different fetch stall percentages
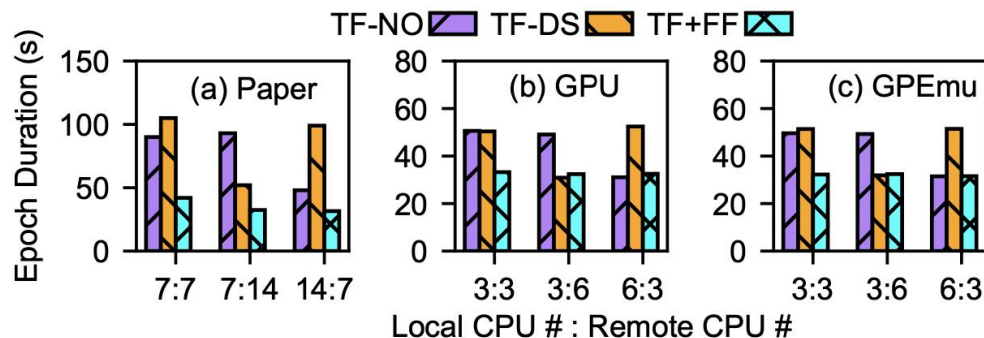


(a) Fetch stalls across varying models

AN-AlexNet        RN18-ResNet18
SN-SqueezeNet   MN-MobileNet
RN50-ResNet50   V11-VGG11

DS Paper
w/ GPU
w/ GPEmu

**GPEmu** values differ from DS paper but **align with GPU** results

# Case Study 2: Preprocessing Disaggregation

❏ FastFlow @ VLDB 2023

❏ FastFlow: dynamically splits preprocessing pipeline between local and remote CPUs

**Pattern is consistent**: FastFlow outperforms others in all scenarios.



**GPEmu** values differ from FastFlow paper but **align with GPU** results

# Supported Research

+ **Data stall analysis**

+ **Preprocessing disaggregation**

+ **Data loader optimization**

+ **Distributed training optimization**

+ **GPU scheduling**

+ **GPU sharing**

### Reproduced 9 papers

DataStall [49], VLDB '21
TF-DS [22], SoCC '23
FastFlow [60], VLDB '23
FFCV [42], CVPR '23
LADL [66], HiPC '19
Synergy [48], OSDI '22
Allox [41], EuroSys '20
Salus [67], MLSys '20
Muri [74], SIGCOMM '22

### Prototyped 3 new optimizations

Small file first caching
Async data loading
File grouping

**More on the Paper!**

### Design Considerations and Analysis of Multi-Level Erasure Coding in Large-Scale Data Centers

Meng Wang
University of Chicago
Chicago, IL, USA
wangm12@uchicago.edu

Jiajun Mao
University of Chicago
Chicago, IL, USA
jiajunm@uchicago.edu

Rajdeep Rana
University of Chicago
Chicago, IL, USA
rajrana22@uchicago.edu

John Bent
Los Alamos National Laboratory
Los Alamos, NM, USA
johnbent@gmail.com

Serkay Olmez
Seagate Research
Longmont, CO, USA
serkay.olmez@seagate.com

Anjus George
Oak Ridge National Laboratory
Oak Ridge, TN, USA
georgea@ornl.gov

Garrett Wilson Ransom
Los Alamos National Laboratory
Los Alamos, NM, USA
gransom@lanl.gov

Jun Li
CUNY Queens College and Graduate
Center
New York, NY, USA
jun.li@qc.cuny.edu

Haryadi S. Gunawi
University of Chicago
Chicago, IL, USA
haryadi@cs.uchicago.edu

**ABSTRACT**
Multi-level erasure coding (MLEC) has seen large deployments in
the field, but there is no in-depth study of design considerations
for MLEC at scale. In this paper, we provide comprehensive design
considerations and analysis of MLEC at scale. We introduce the
design space of MLEC in multiple dimensions, including various
code parameter selections, chunk placement schemes, and various
repair methods. We quantify their performance and durability, and
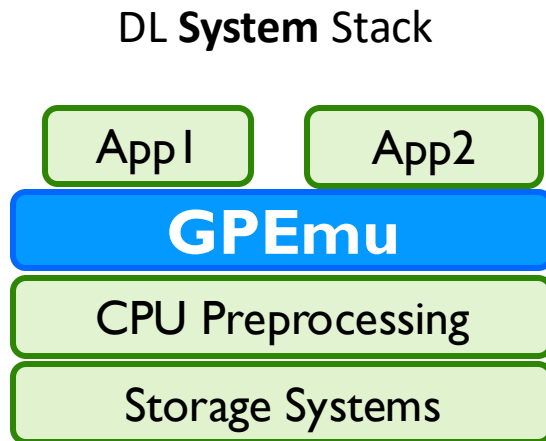
# Conclusion

**A GPU Emulator for Faster** and **Cheaper** prototyping and evaluation of DL system research.

DL **System** Stack

**+** Without GPUs

**+** 30+ models and 6 GPUs

**+** Easy to use

**+** Support both Single node & Distributed setups

App1   App2

**GPEmu**

CPU Preprocessing

Storage Systems

- **Time Emulation**
- **Memory Emulation**
- **Distributed Support**
- **Sharing Support**

**Effectively reveal** DL system **bottlenecks** and **quickly show** the **benefits** of new system optimizations

# Thank you!
# Questions?