# Transparent Hot-Swapping of Containerized AI/ML Workloads

Radostin Stoyanov - PhD Student, Scientific Computing Group
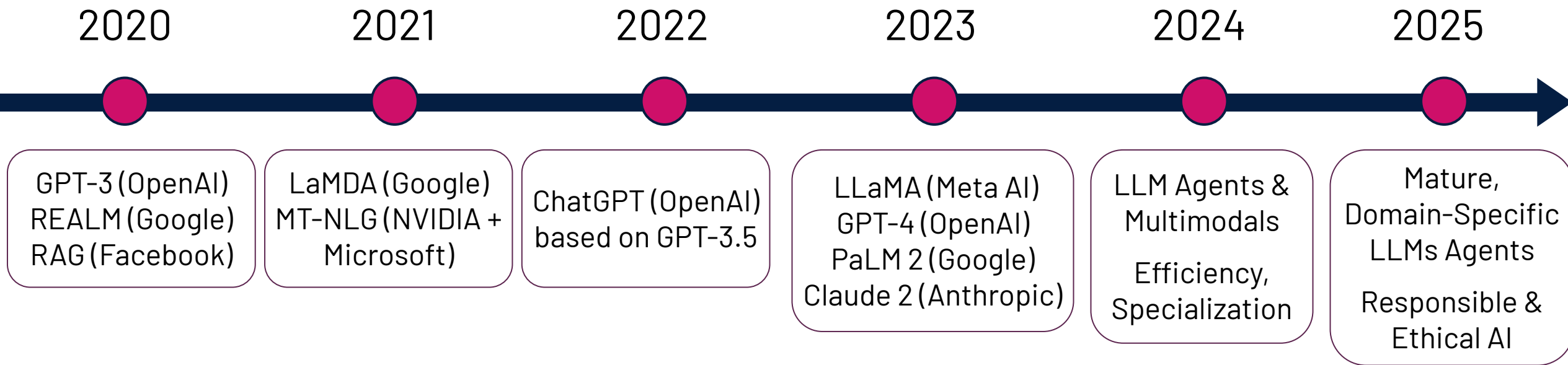
Collaboration with Viktória Spišaková, Marcin Copik, and Adrian Reber

Supervisors: Prof. Rodrigo Bruno, Prof. Wes Armour

# AI/ML Workloads in Production

A vast number of applications are leveraging LLMs

- Demand for **inference** far exceeds that of training
- Large number of inferences necessary to **amortize high training costs**

| 2020 | 2021 | 2022 | 2023 | 2024 | 2025 |
|------|------|------|------|------|------|
| GPT-3 (OpenAI) REALM (Google) RAG (Facebook) | LaMDA (Google) MT-NLG (NVIDIA + Microsoft) | ChatGPT (OpenAI) based on GPT-3.5 | LLaMA (Meta AI) GPT-4 (OpenAI) PaLM 2 (Google) Claude 2 (Anthropic) | LLM Agents & Multimodals  Efficiency, Specialization | Mature, Domain-Specific LLMs Agents  Responsible & Ethical AI |

# LLM Inference Serving

**Workloads Characteristics**

- Receive large number of queries with strict SLOs
- Unpredictable output lengths
- Run on expensive GPUs with high energy consumption[1]
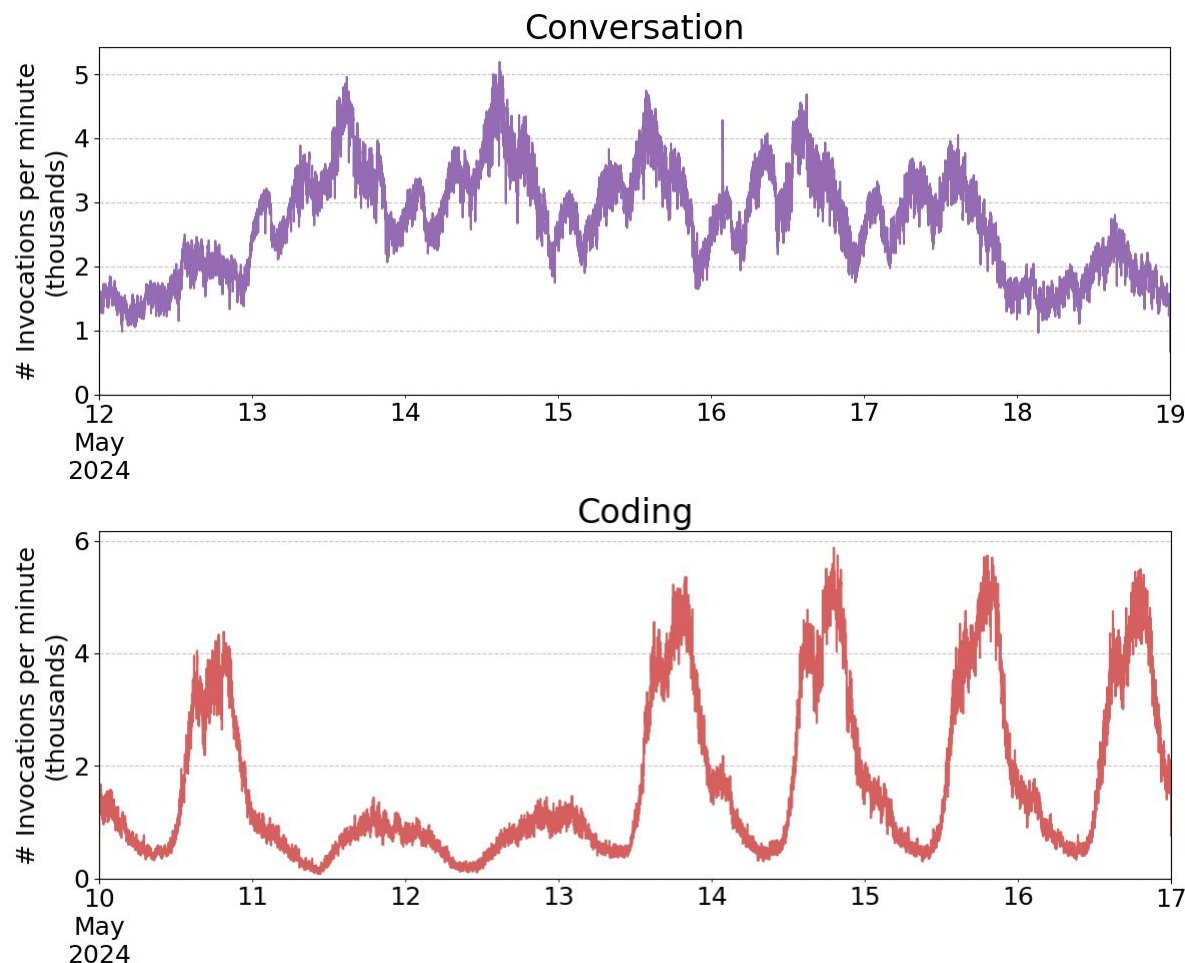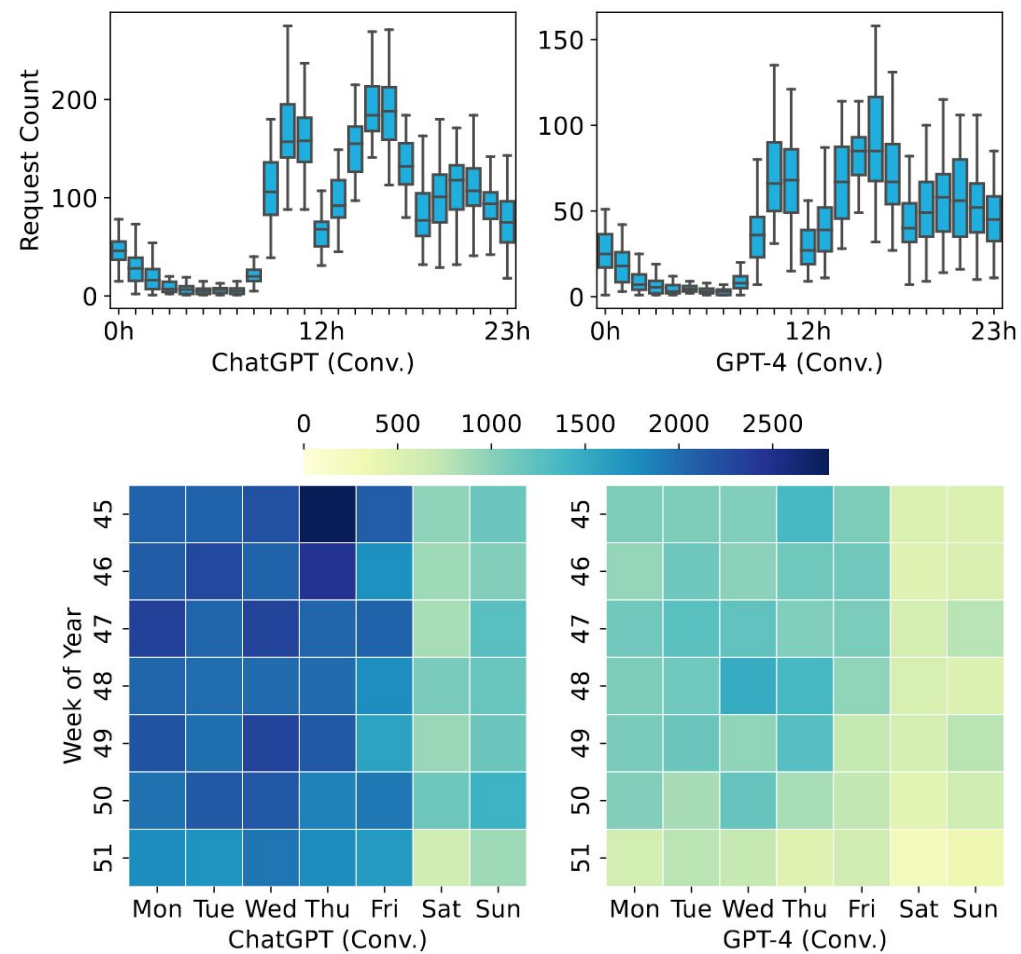- GPU resources in multi-tenant clusters (MLaaS) are often underutilized

**Key Performance Metrics:** Latency & Throughput

**Optimization Strategies**

- Efficient KV cache management with distributed inference
- Iteration-level scheduling with continuous batching
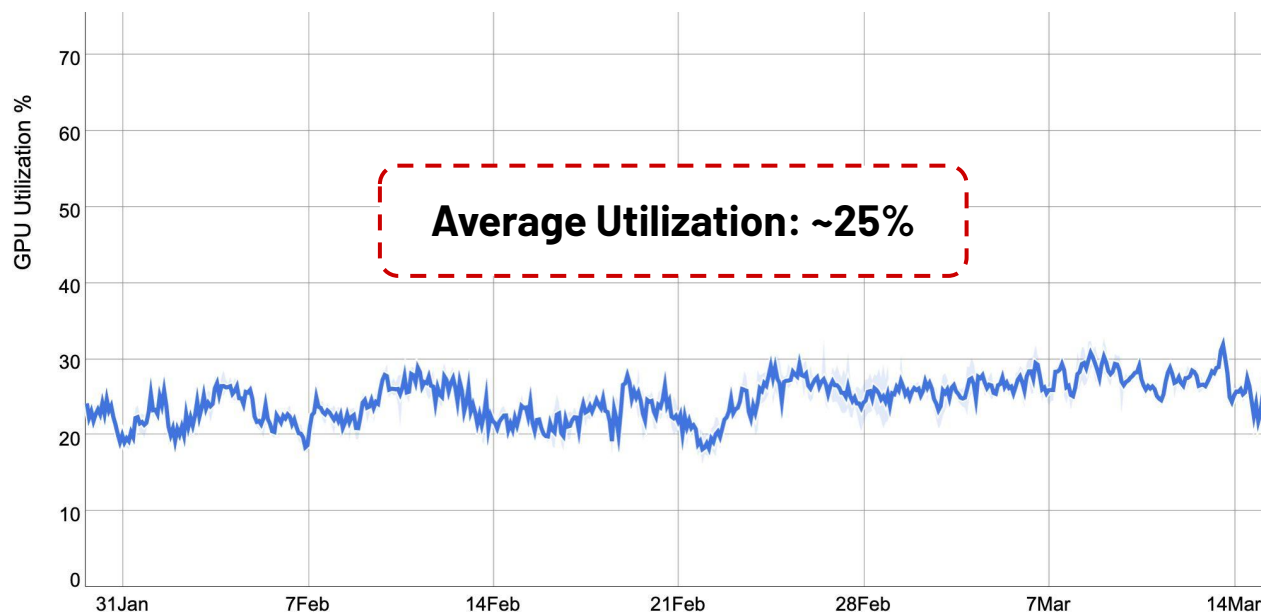- Smaller (specialized) models are faster & less expensive to run

[1] *Double-Exponential Increases in Inference Energy: The Cost of the Race for Accuracy.* Zeyu Yang, Karel Adamek, Wesley Armour.
https://arxiv.org/abs/2412.09731
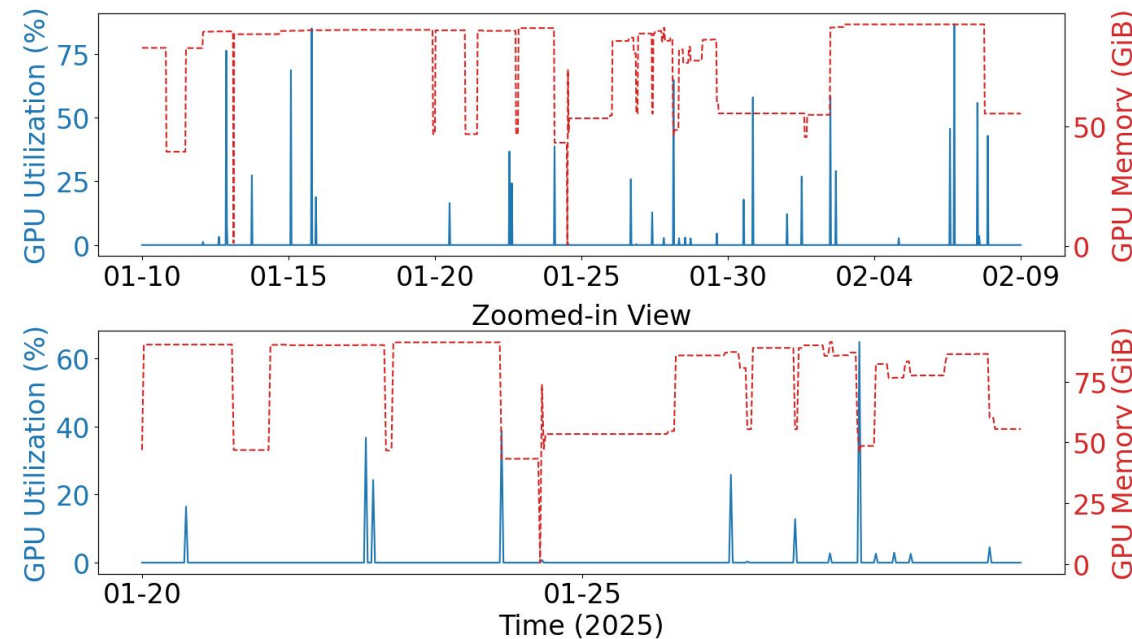
# Real-World LLM Deployments



BurstGPT: A Real-world Workload Dataset to Optimize LLM Serving Systems
https://github.com/HPMLL/BurstGPT

Azure LLM Inference Dataset trace
https://github.com/Azure/AzurePublicDataset

# Problem: Low GPU Utilization

**Inference Serving / Notebooks / Visualization Workloads**



Average Utilization: ~25%

*"Improving GPU Utilization using Kubernetes Engine"*
Maulin Patel, Pradeep Venkatachalam (Google)

**e-INFRA CZ Kubernetes Cluster** [1,2,3]
**GPU Utilization of LLM Inference Workloads**



H100 NVL (94GB) serving LLaMA 3.3, Qwen 2.5-Coder, Command R7B, Phi-4, QWQ, and MXBai-Embed-Large

[1] *Efficient Transparent Checkpointing of AI/ML Workloads in Kubernetes.* Viktória Spišaková, Radostin Stoyanov, Adrian Reber. KubeCon 2025
[2] *Optimizing Resource Utilization for Interactive GPU Workloads with Container Checkpointing.* Viktória Spišaková, Radostin Stoyanov. FOSDEM 2025
[3] *Kubernetes Scheduling with Checkpoint/Restore: Challenges & Open Problems.* Viktória Spišaková et. al., JSSPP 2025

# Transparent Checkpointing

# Transparent Checkpointing

## Checkpoint/Restore in Userspace

- Transparent checkpointing of CPU-GPU state

- Supports both AMD and NVIDIA GPUs

- Integrated with Docker, Podman, Kubernetes

github.com/checkpoint-restore/criu
github.com/nvidia/cuda-checkpoint

*CRIUgpu: Transparent Checkpointing of GPU-Accelerated Workloads*. Radostin Stoyanov, Viktória Spišaková, Jesús Ramos, Steven Gurfinkel, Andrei Vagin, Adrian Reber, Wesley Armour, Rodrigo Bruno. (2025). https://arxiv.org/abs/2502.16631

# Challenges

**cuda-checkpoint**

- Only for GPU state; CPU processes continue running

- Checkpoints a single process, not a process tree (container)

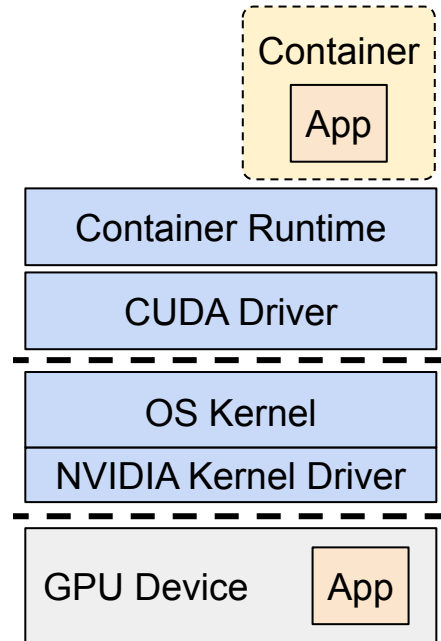- Without CRIU, it does not support migration

**CRIU + GPU Plugins**

- Saving checkpoint data to disk can be very slow

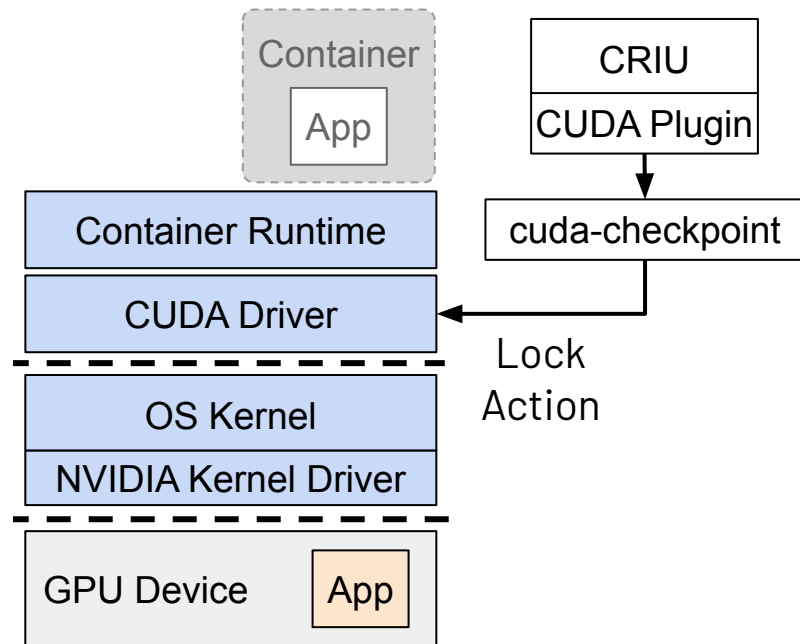- Container engines are not optimized for large container checkpoints[1]

[1] https://github.com/checkpoint-restore/criu/issues/2519

# Transparent Hot-Swapping

# Transparent Hot-Swapping

Container

App

Container Runtime

CUDA Driver

OS Kernel

NVIDIA Kernel Driver

GPU Device    App

**Starting Container A**

- Mounting GPU libraries & device files in container

- Loading model, allocating GPU memory, KV Cache
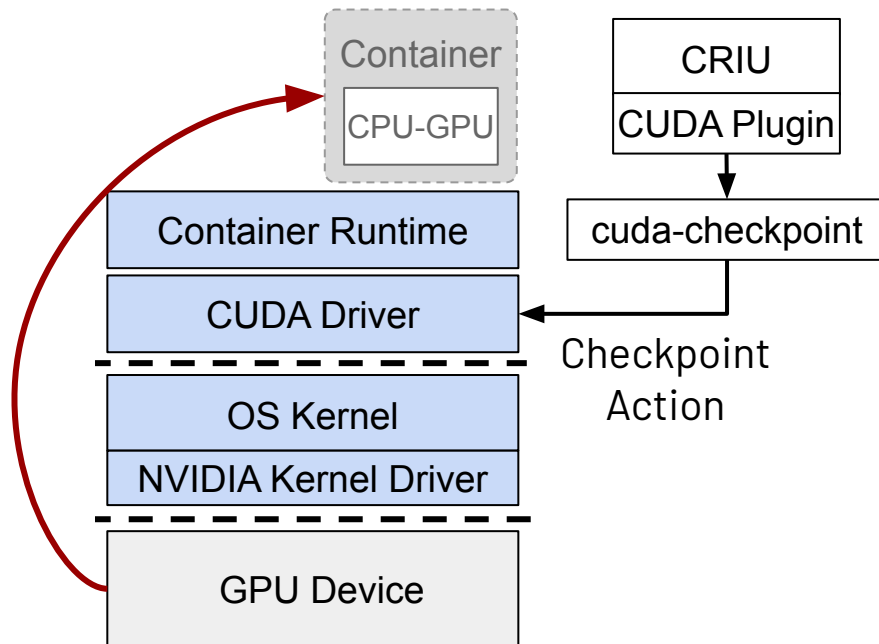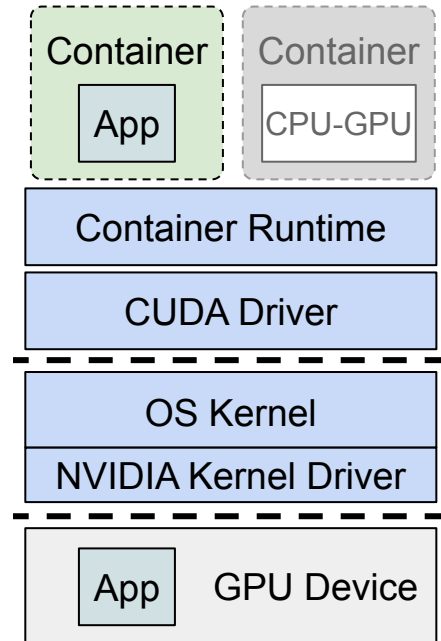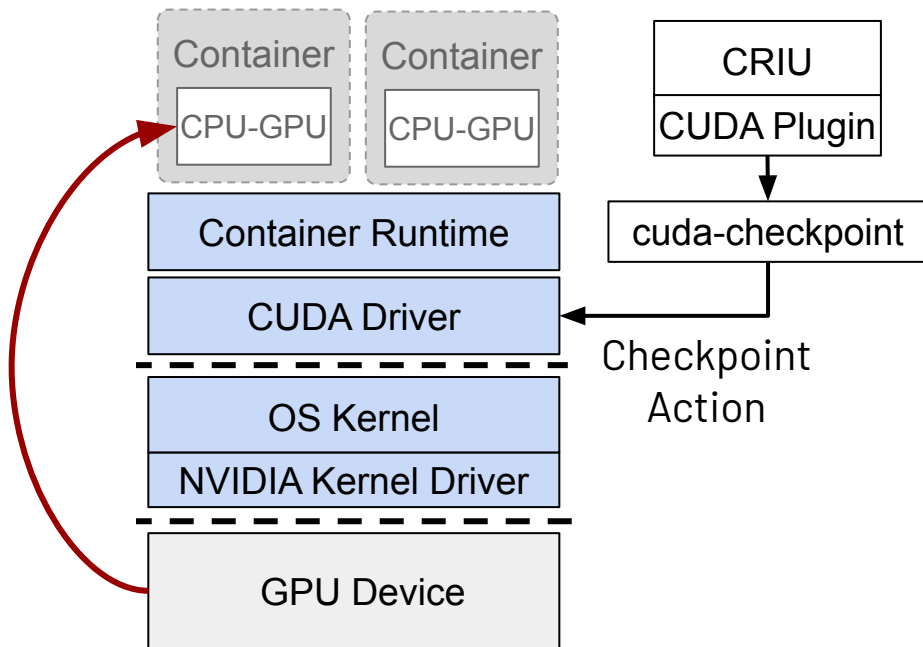
- Ready to accept user requests

# Transparent Hot-Swapping

**Transparent Preemption of GPU Workloads**

- Stop execution of CPU processes:

  `ptrace(SEIZE+INTERRUPT)`

- Lock CUDA driver APIs

## Diagram

Container
App

Container Runtime

CUDA Driver

OS Kernel

NVIDIA Kernel Driver

GPU Device — App

CRIU
CUDA Plugin

cuda-checkpoint

Lock Action

# Transparent Hot-Swapping

Container
CPU-GPU

Container Runtime

CUDA Driver

OS Kernel

NVIDIA Kernel Driver

GPU Device

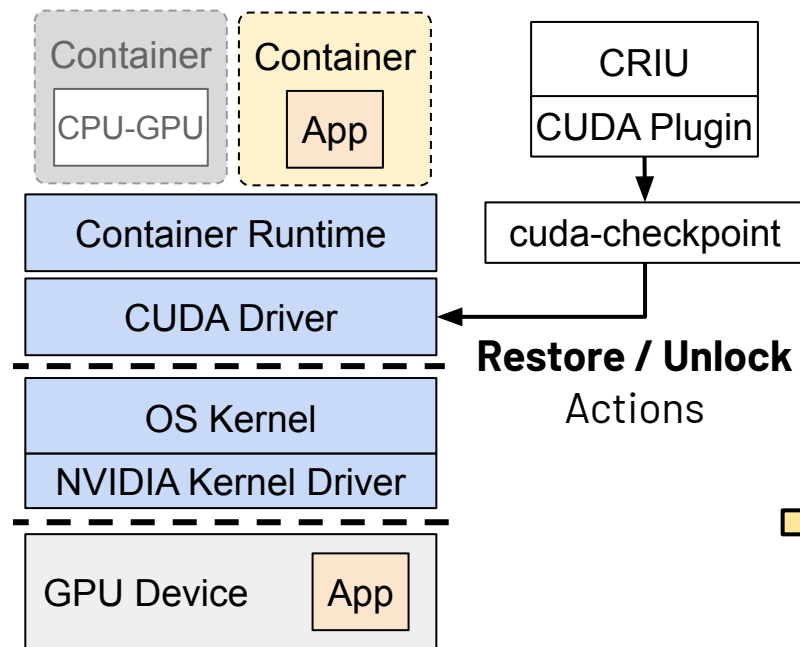CRIU
CUDA Plugin

cuda-checkpoint

Checkpoint
Action

**Transparent Preemption of GPU Workloads**

- Stop execution of CPU processes:

  `ptrace(SEIZE+INTERRUPT)`

- Lock CUDA driver APIs

- Checkpoint GPU state to host memory

- Leave container container A in a "stopped" state

⇨ GPU resources have been released

# Transparent Hot-Swapping



**Starting Container B**

- Mounting GPU libraries & device files in container

- Loading model, allocating GPU memory, KV Cache

⟹ Container B is ready to accept user requests

# Transparent Hot-Swapping

**Transparent Preemption of GPU Workloads**

- Stop execution of CPU processes:

  `ptrace(SEIZE+INTERRUPT)`

- Lock CUDA driver APIs

- Checkpoint GPU state to host memory

- Leave container container B in a "stopped" state

⟹ GPU resources have been released

# Transparent Hot-Swapping

| Container | Container |
|---|---|
| CPU-GPU | App |

| CRIU |
|---|
| CUDA Plugin |

| cuda-checkpoint |
|---|

| Container Runtime |
|---|
| CUDA Driver |

**Restore / Unlock**
Actions

| OS Kernel |
|---|
| NVIDIA Kernel Driver |

| GPU Device | App |
|---|---|

**Resuming Preempted GPU Workloads**

- Resume execution of CPU processes

- Restore GPU state from host memory

- Unlock CUDA driver APIs

⟹ Container A is ready to accept user requests

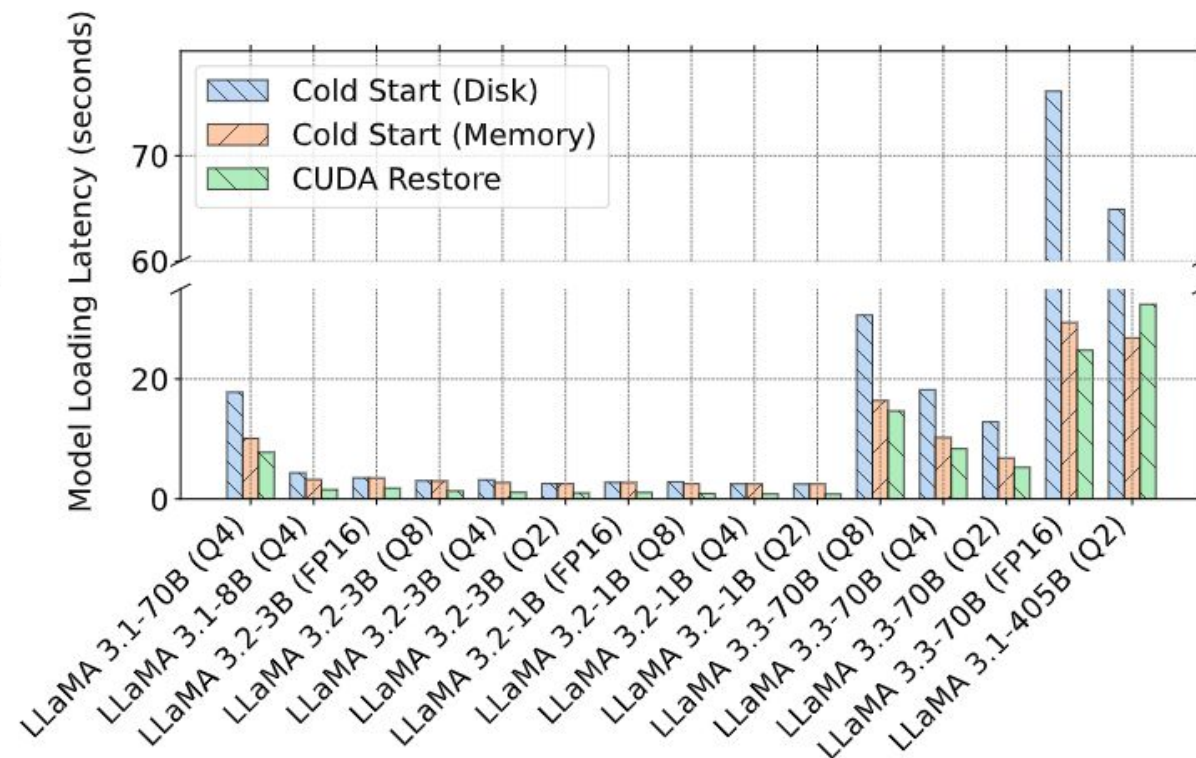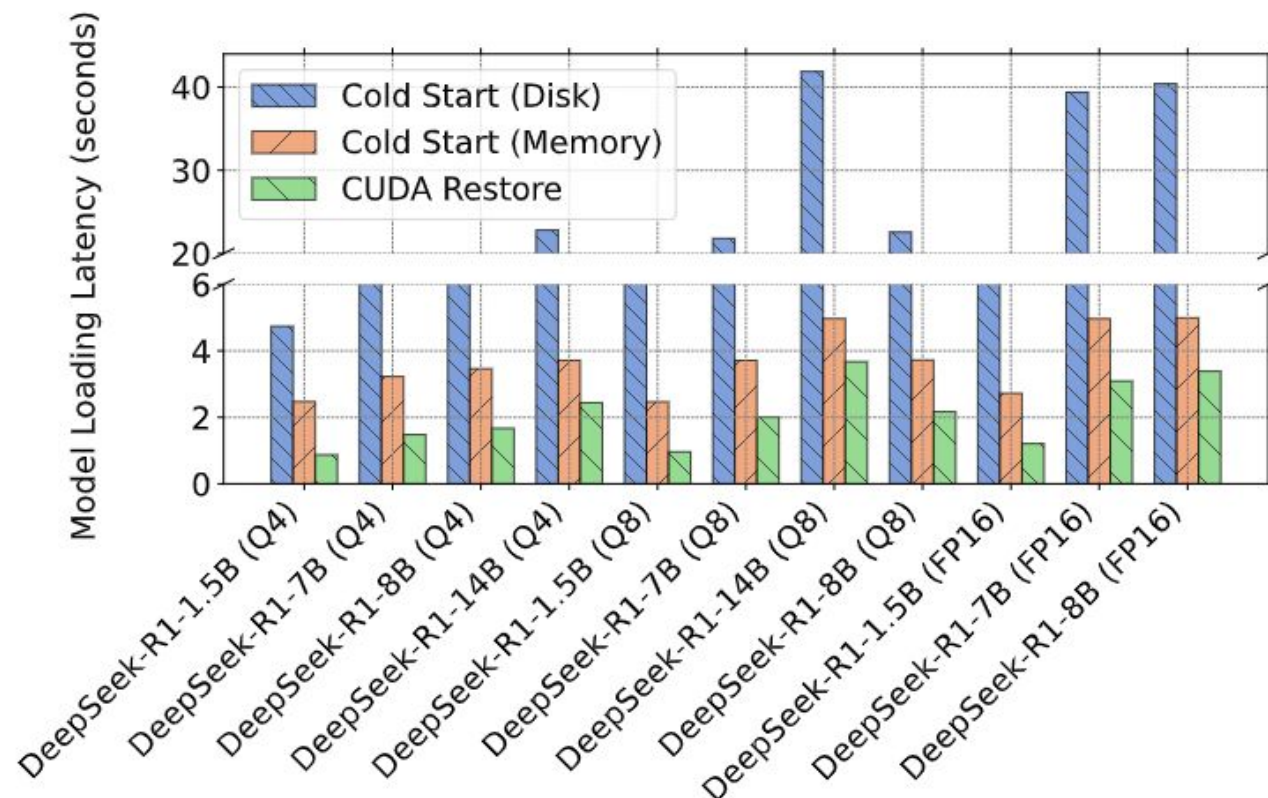# Transparent Hot-Swapping Demo



## Part 2: Transparent Hot-Swapping of vLLM Containers

### Part 1: vLLM Initialization + Checkpoint Creation

# Evaluation



Ollama LLM Inference Workloads on NVIDIA A100 (SXM4 80GB)

# Summary & Questions

- Transparent Hot-Swapping of GPU Workloads

- Optimized Resource Utilization for LLM Inference

- Out-of-the-box Integration with Container Platforms

github.com/checkpoint-restore/criu

github.com/nvidia/cuda-checkpoint