# Zico

# Efficient GPU Memory Sharing for Concurrent DNN Training

**Gangmuk Lim**[*], Jeongseob Ahn[§], Wencong Xiao[†], Youngjin Kwon[‡], Myeongjae Jeon[*]

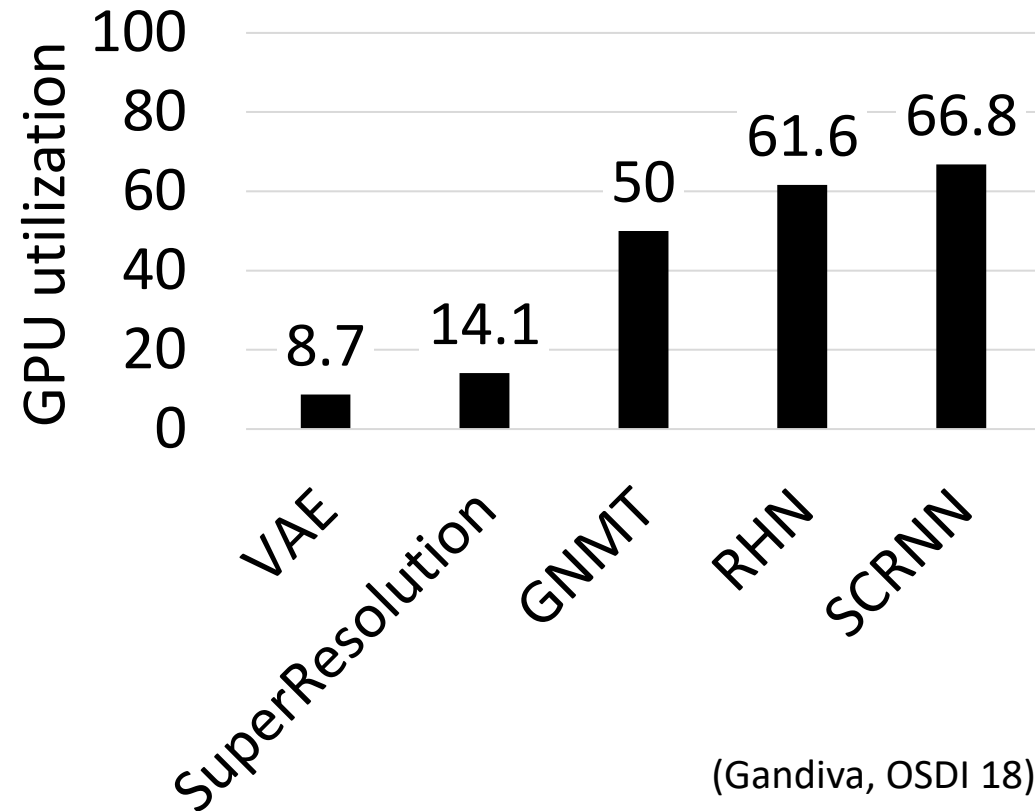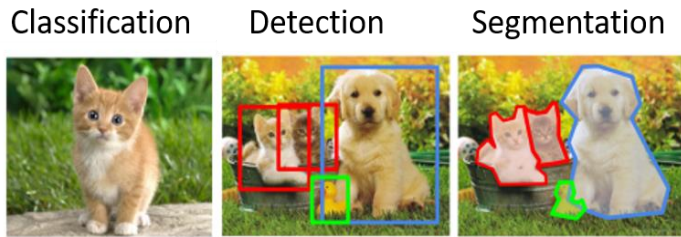[*] UNIST, [§] Ajou University, [†] Alibaba Group, [‡] KAIST

# GPU Utilization in DNN Training

*DNN training jobs require GPU*

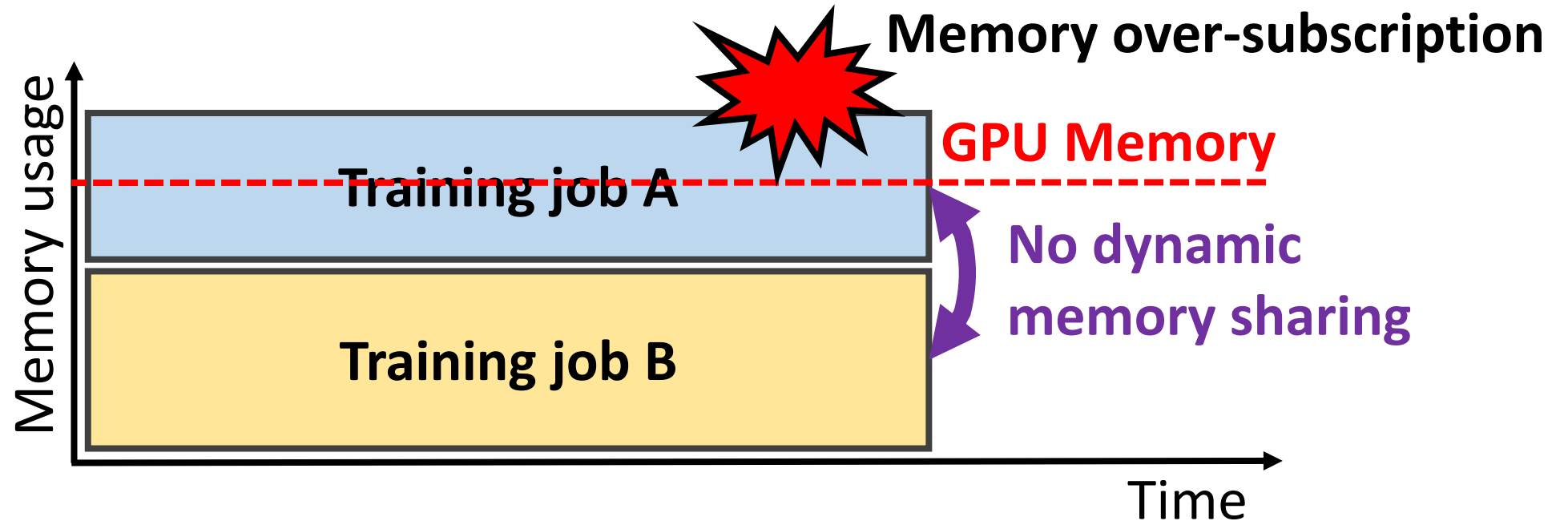**GPU core is often under-utilized**

*[Gandiva OSDI 18, Philly ATC 19, Salus MLSys 20]*
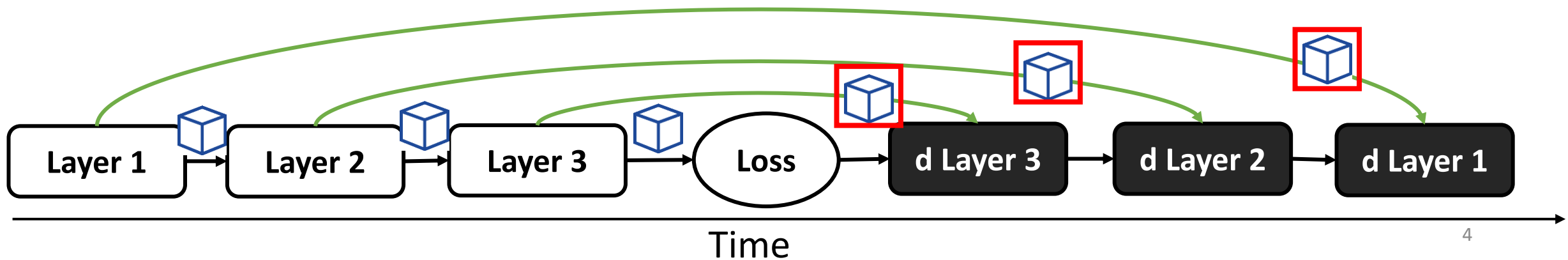


(Gandiva, OSDI 18)
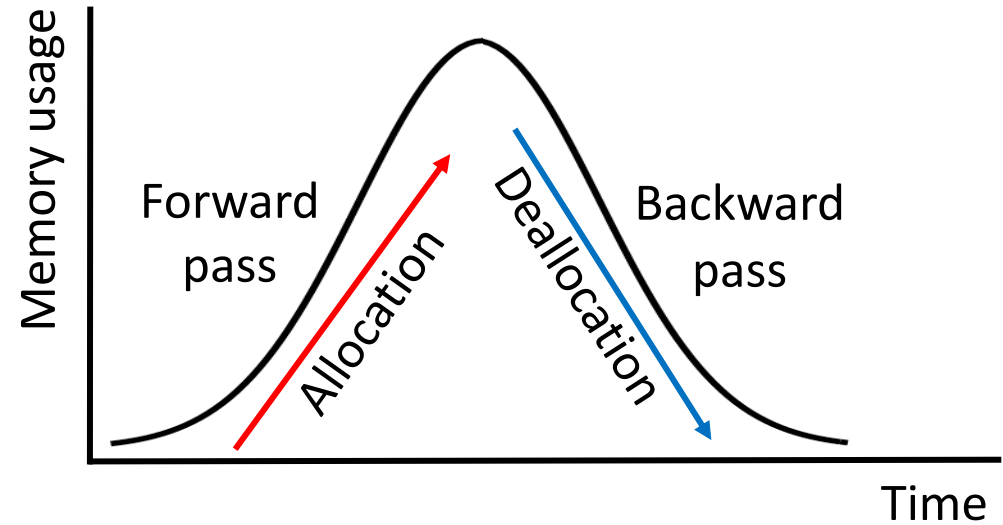
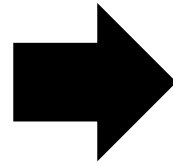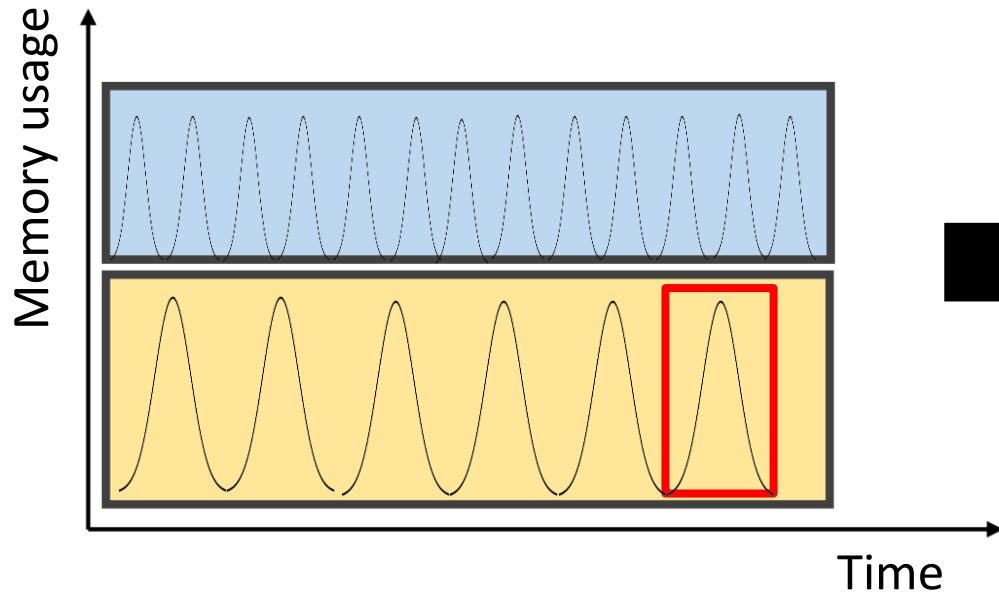# Existing GPU Sharing Solution

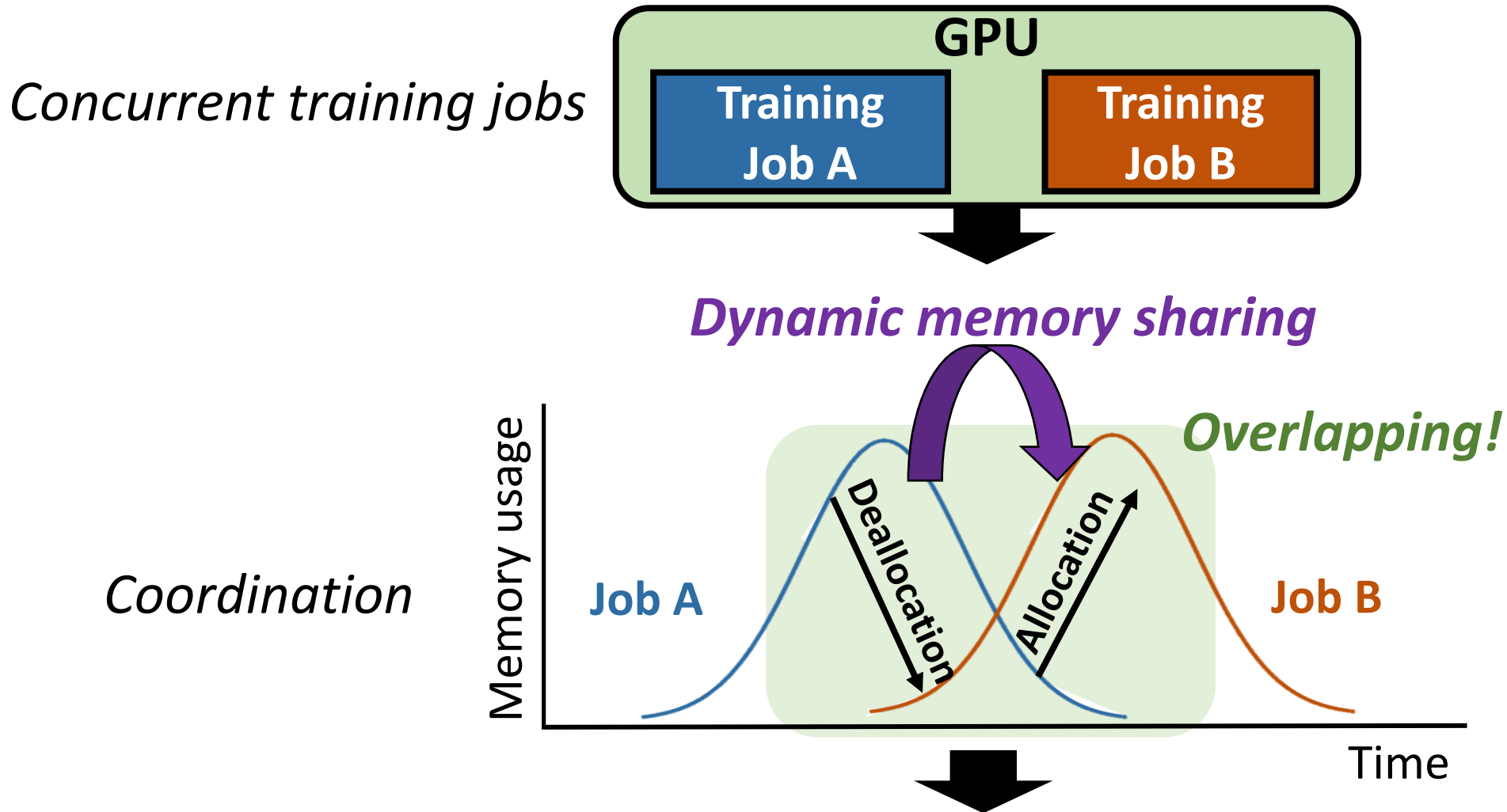## *Working set in concurrent training easily exceeds GPU memory*



e.g. NVIDIA MPS, NVIDIA MIG, Salus

# Cyclic Memory Usage Pattern

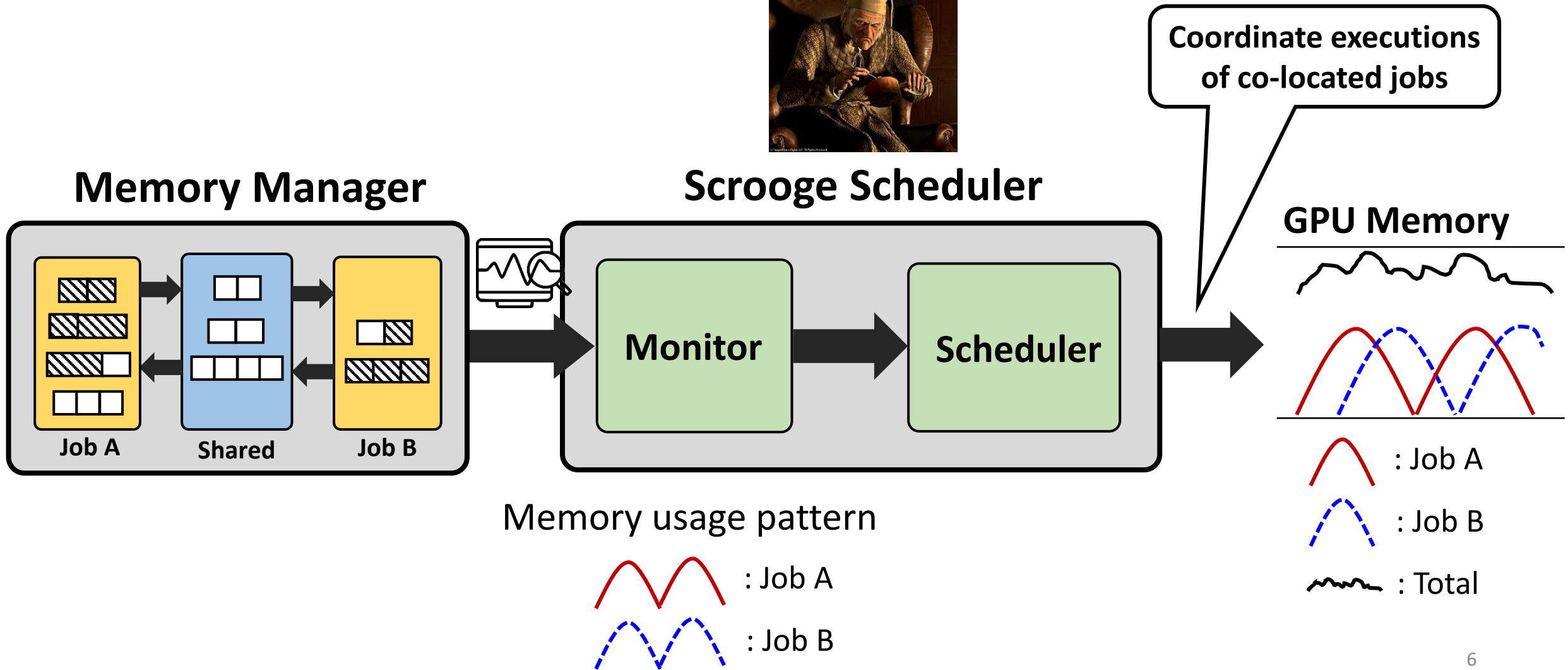*DNN training job shows cyclic memory usage pattern*

# Memory Sharing Opportunity

**GPU**

*Concurrent training jobs*

**Training Job A**  **Training Job B**

**Dynamic memory sharing**

**Overlapping!**

*Coordination*

Memory usage

Job A

Deallocation

Allocation

Job B

Time

***Efficiently reducing the system-wide memory footprint***

# Zico Overview



**Memory Manager**

Job A   Shared   Job B

**Scrooge Scheduler**

Monitor → Scheduler

Memory usage pattern

—— : Job A

- - - : Job B

Coordinate executions
of co-located jobs

**GPU Memory**

—— : Job A

- - - : Job B

~~ : Total

# Contributions

**Safe and efficient memory management**
> Handling asynchrony between CPU and GPU
>
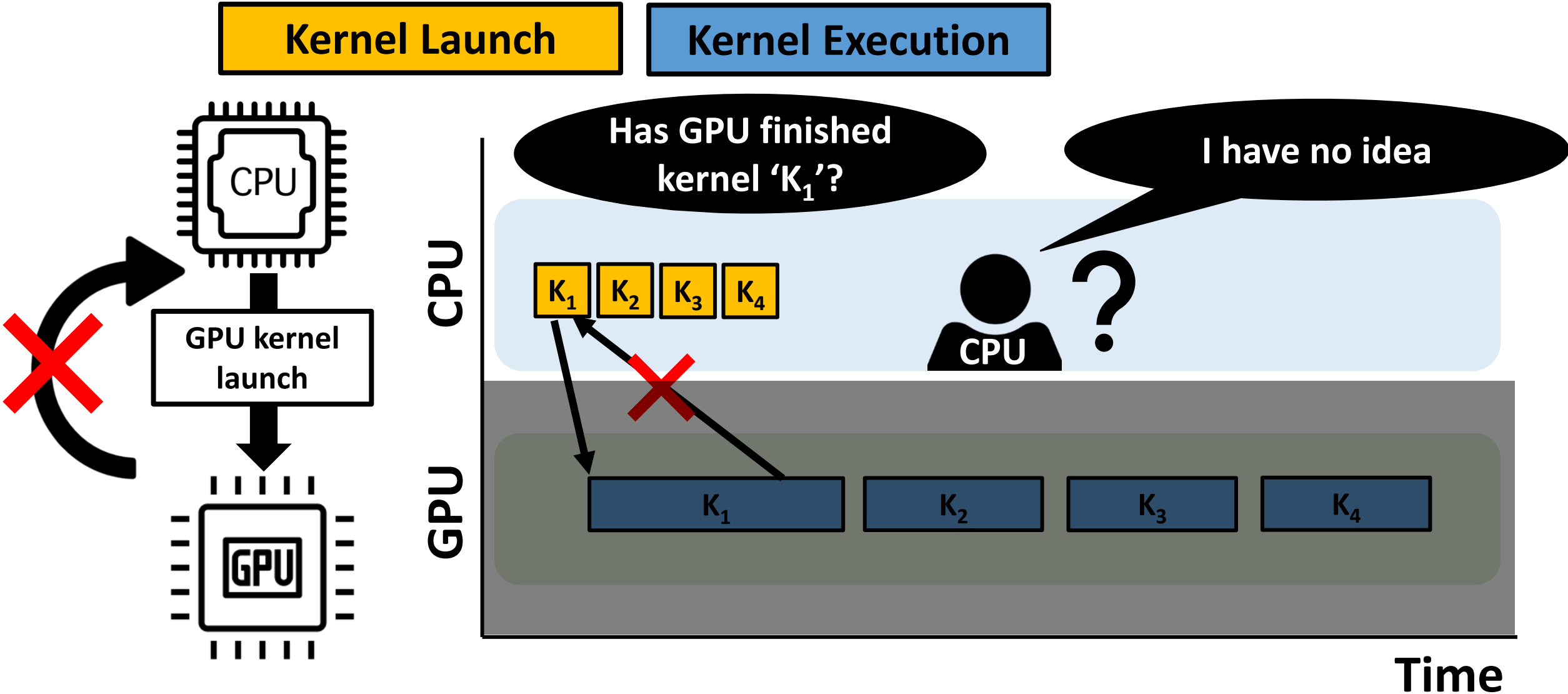> Preventing early allocation

**Memory-aware scheduling**
> Minimizing time delay while maximizing throughput
>
> Widely applicable (identical jobs, non-identical jobs)
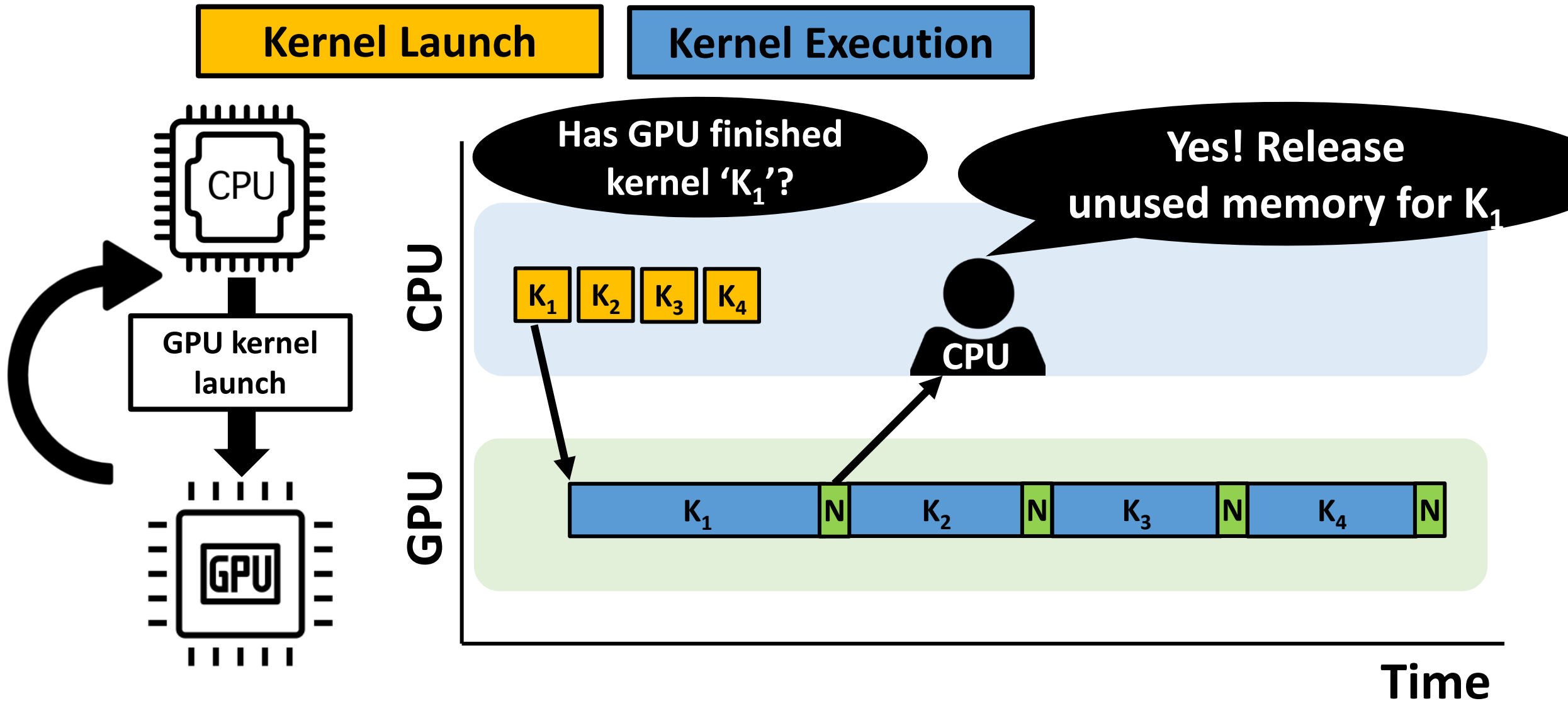
# Asynchrony between CPU and GPU
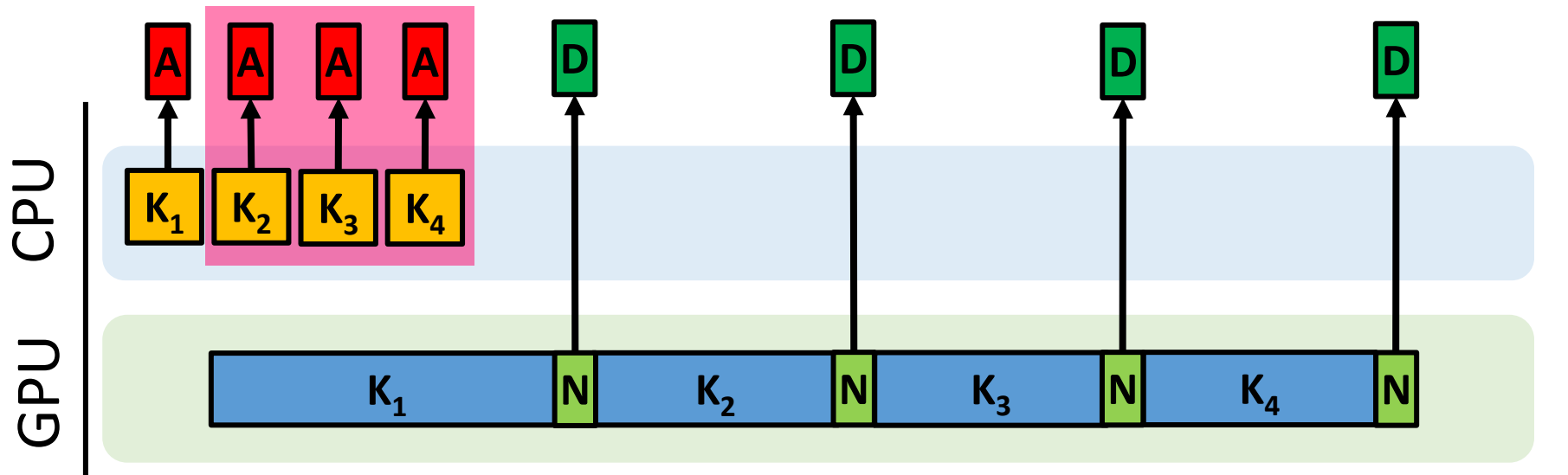
# Early Memory Allocation

CPU kernel launch speed > GPU execution speed

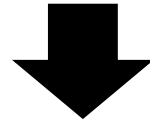*Early memory allocations for kernels which have not started its execution yet*

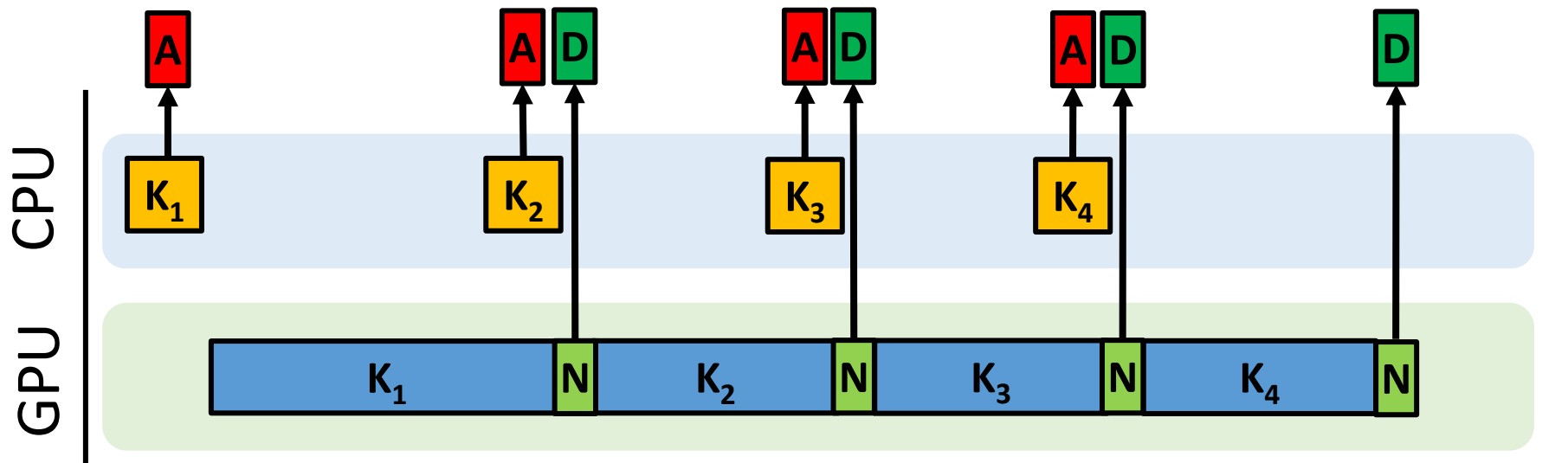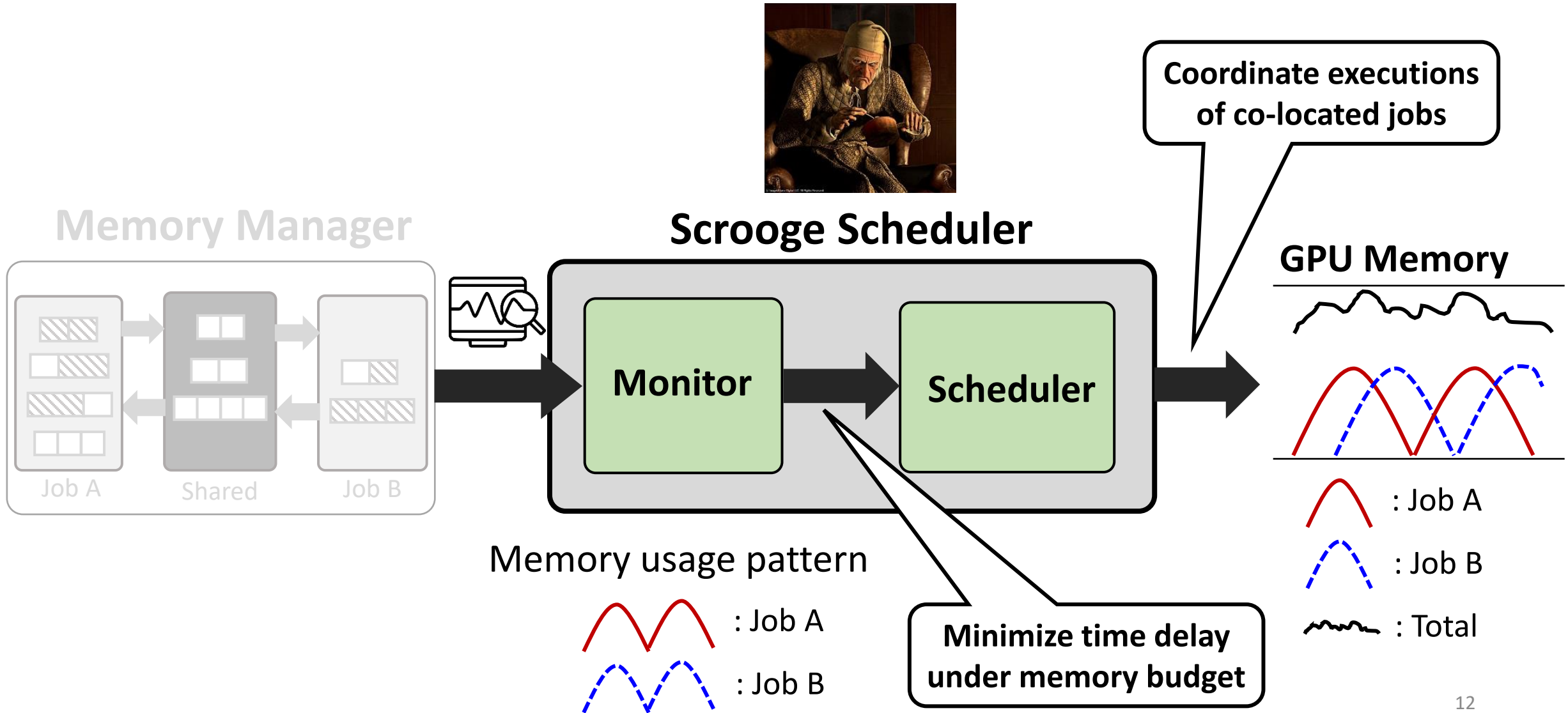***Increasing memory consumption unnecessarily***

# Controlling Inflight Kernel

Controlling the number of inflight kernel

Preventing early allocation

# Roadmap



**Memory Manager**

**Scrooge Scheduler**

Monitor

Scheduler

**GPU Memory**

Coordinate executions
of co-located jobs

Job A      Shared      Job B

Memory usage pattern

— : Job A

--- : Job B

Minimize time delay
under memory budget

— : Job A

--- : Job B

— : Total

# Naïve Scheduling



*Wavelet: Efficient DNN Training with Tick-Tock Scheduling, MLSys 21*
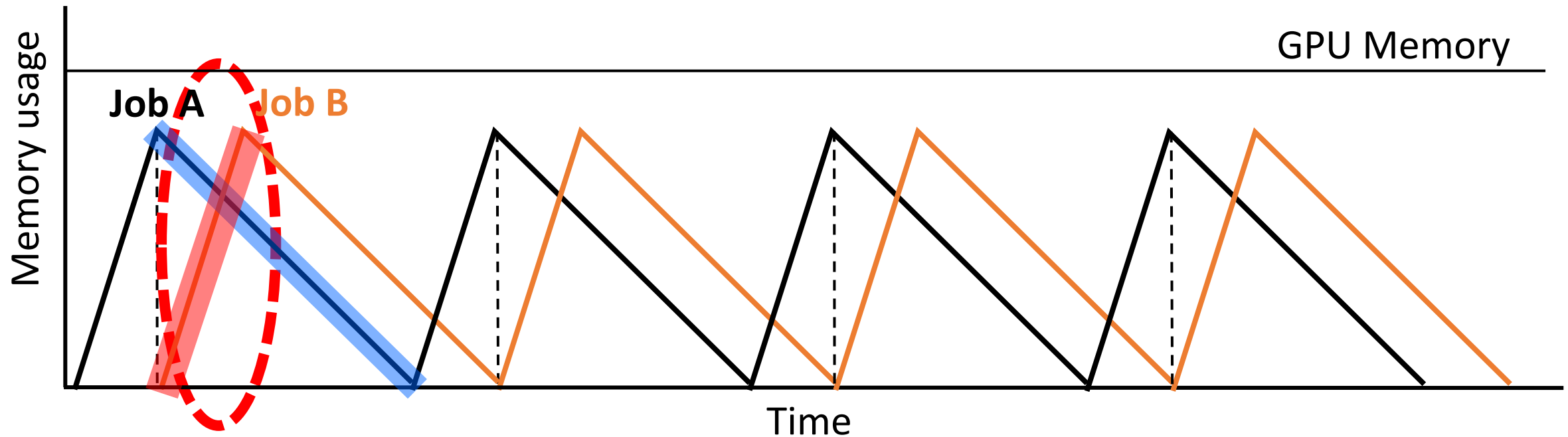
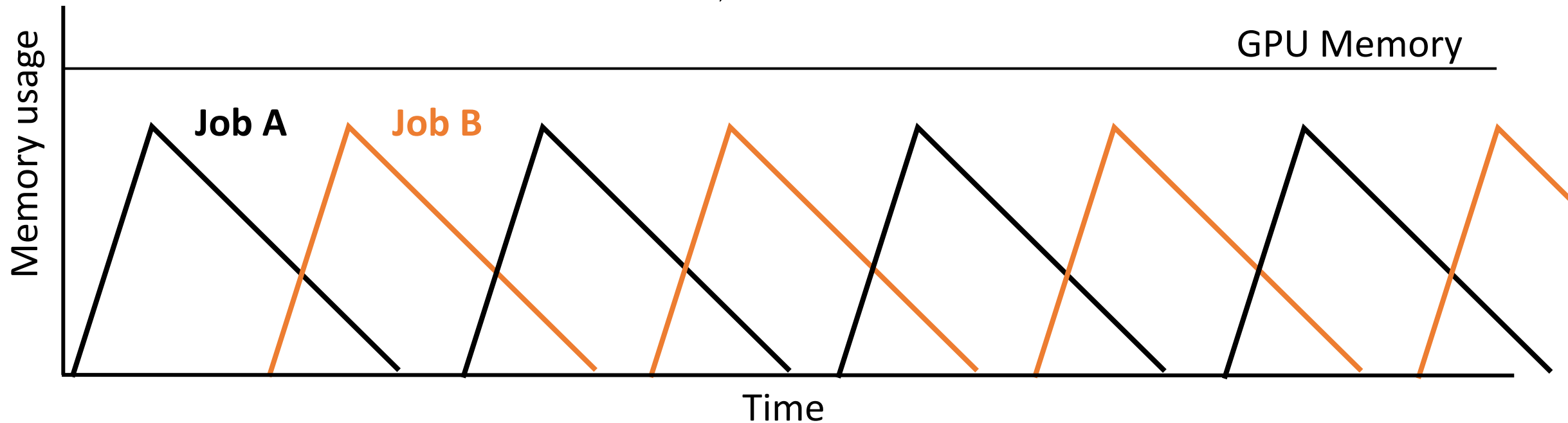# Limitation of Naïve Scheduling

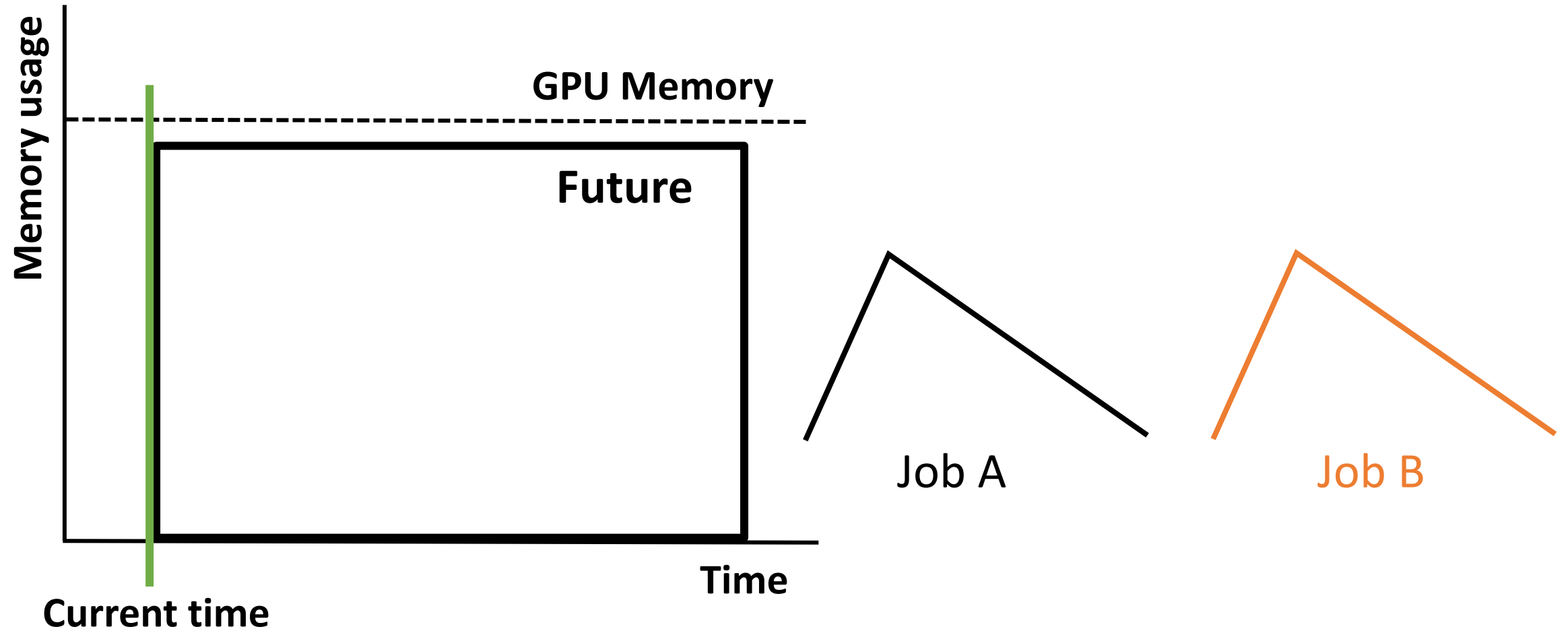☹ **Memory over-subscription** ➡ **Hurting the throughput**

# What We Want to Achieve
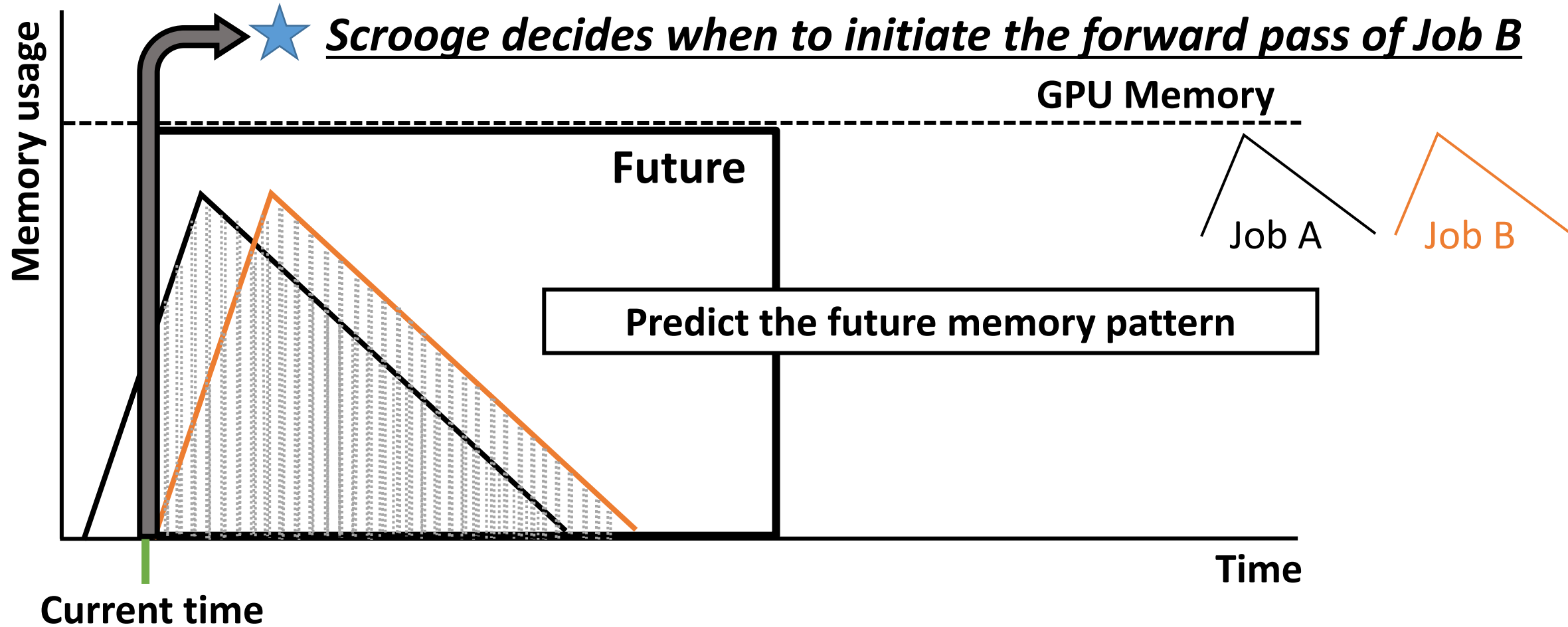
☹ **Memory over-subscription** ➡ **Hurting the throughput**

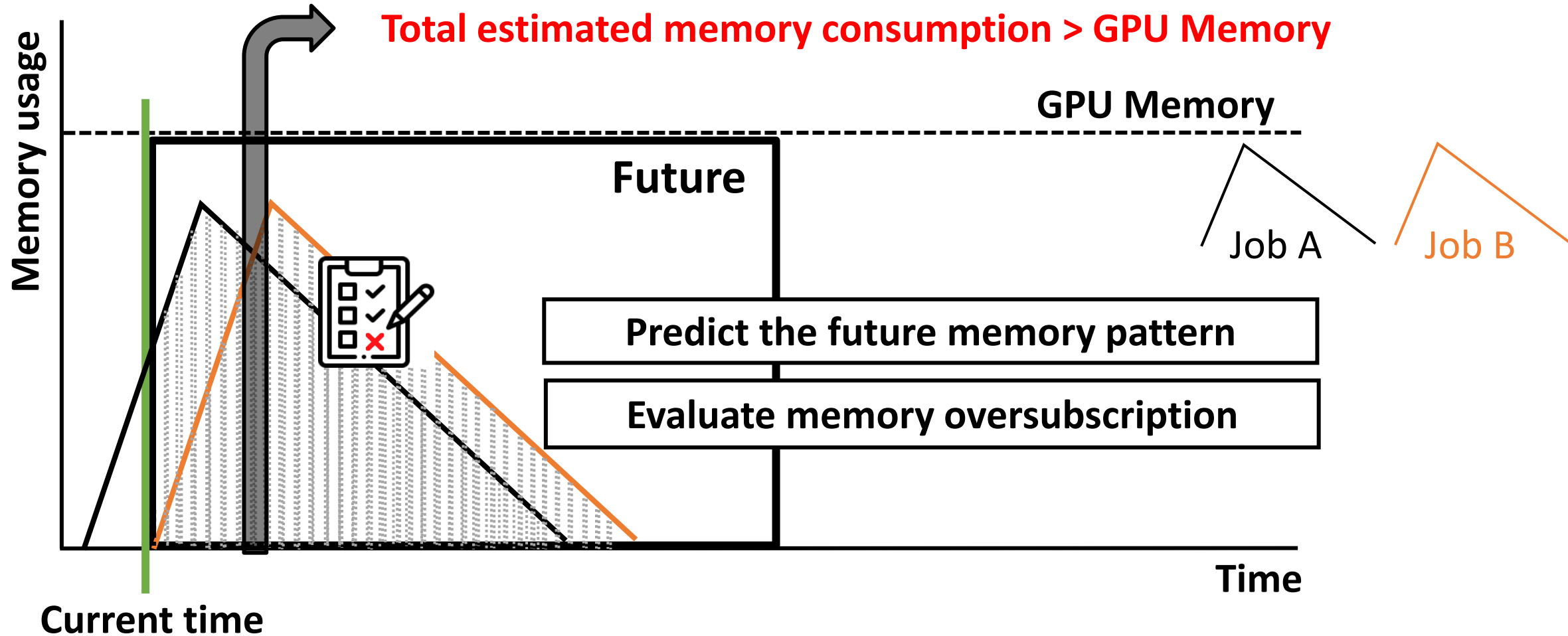☺ **Proper coordination** ➡ **Preventing over-subscription**

# Scrooge Scheduler

# Scrooge Scheduler



★ ***Scrooge decides when to initiate the forward pass of Job B***

GPU Memory

Memory usage

Future

Predict the future memory pattern

Job A

Job B

Current time

Time

# Scrooge Scheduler

# Scrooge Scheduler

# Scrooge Scheduler



No over-subscription will occur
Schedule a new iteration of Job B !

GPU Memory

Future

Memory usage

Current time

Time

Predict the future memory pattern

Evaluate memory oversubscription

Job A    Job B

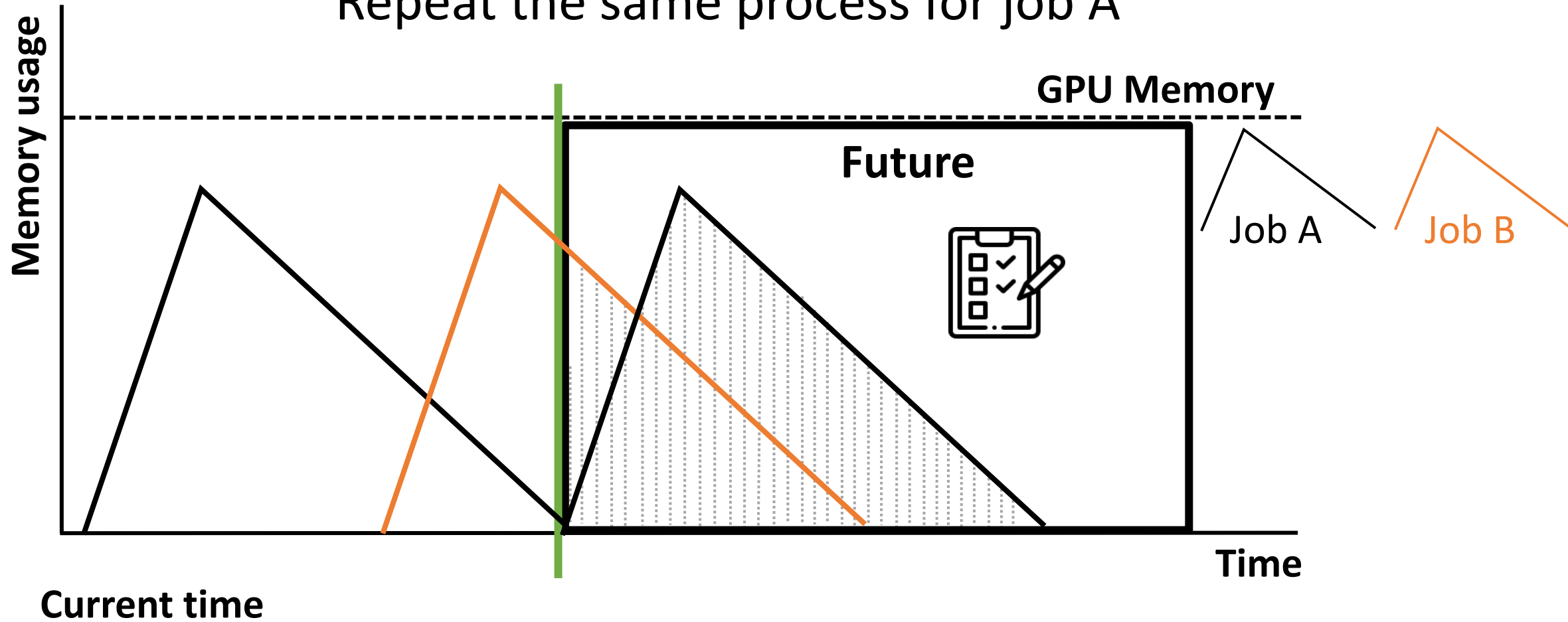# Scrooge Scheduler

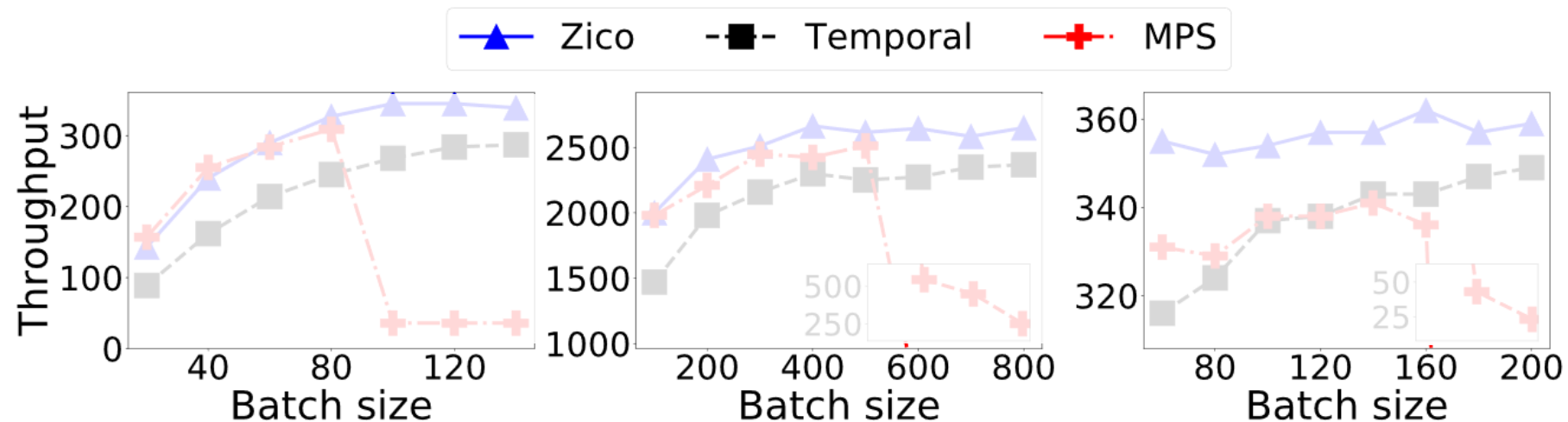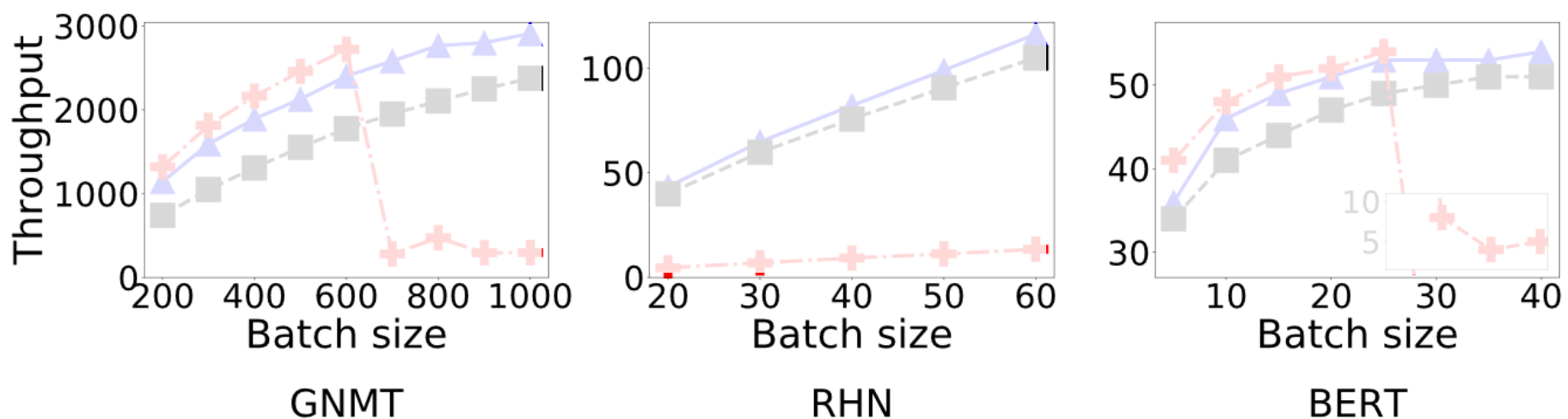## Repeat the same process for job A

# Evaluation

- Machine
  - GPU: Tesla V100 GPU, RTX 2080 Ti GPU
  - CPU: 3.8 GHz Intel Xeon(R)Gold 5222 4 CPU cores
  - RAM: 64 GB

- Benchmark
  - NASNet, ResNet-110, ResNet-50, GNMT, BERT, RHN

- Policies in comparison
  - Temporal: Ideal temporal sharing (no job switching overhead)
  - Spatial: NVIDIA MPS (no dynamic memory sharing)
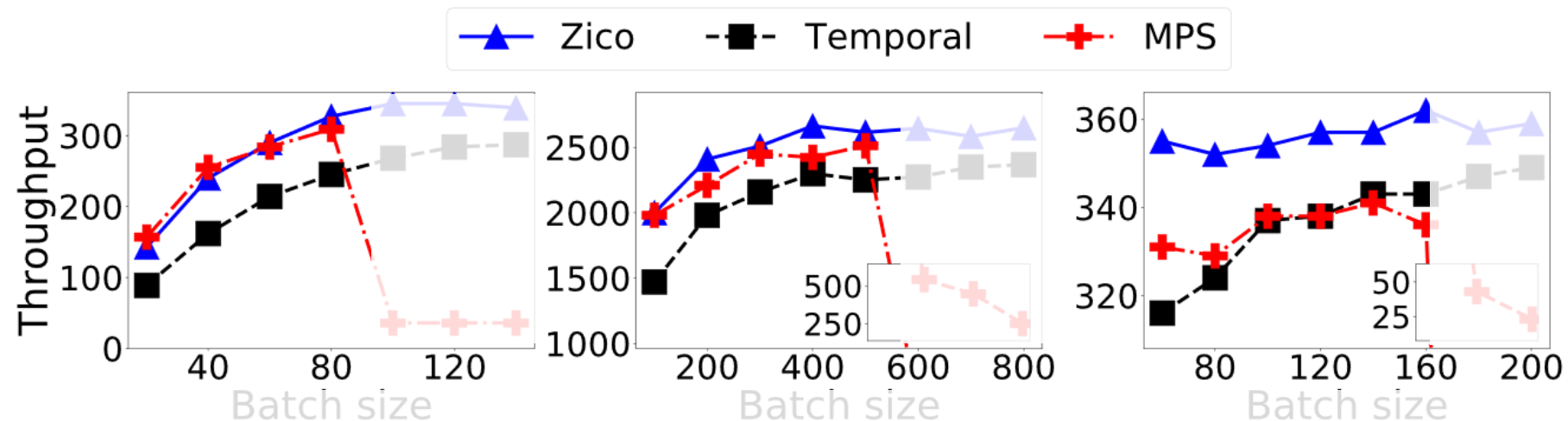
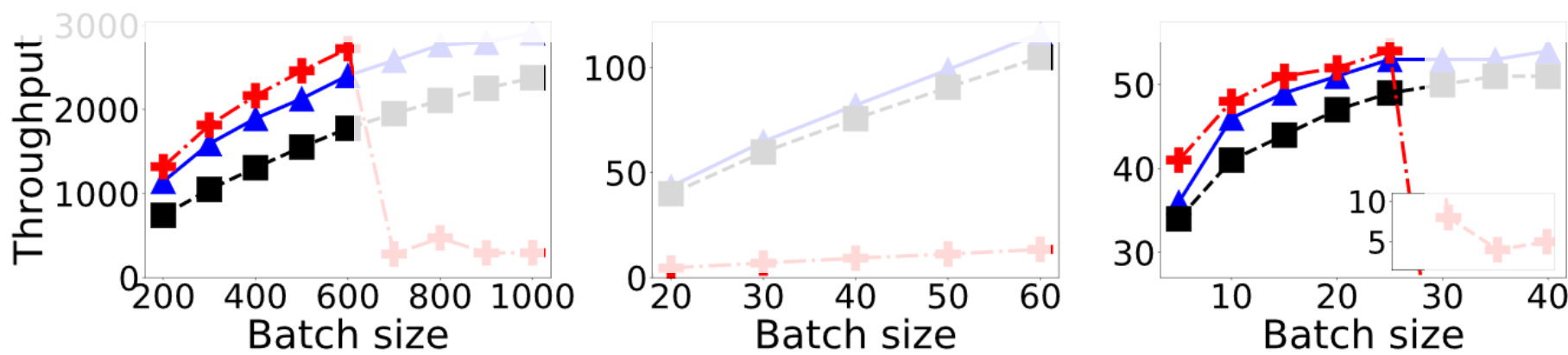- Base framework: TensorFlow v1.13.1

# Throughput: Identical Jobs



**Zico/Temporal: 1.03x ~ 1.6x**

# Throughput: Identical Jobs
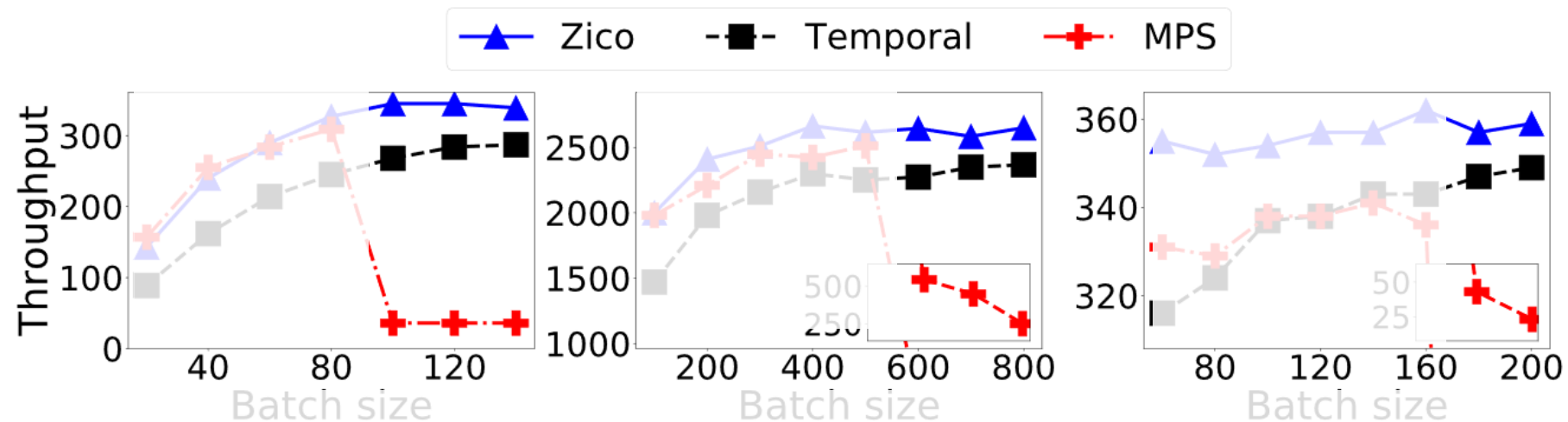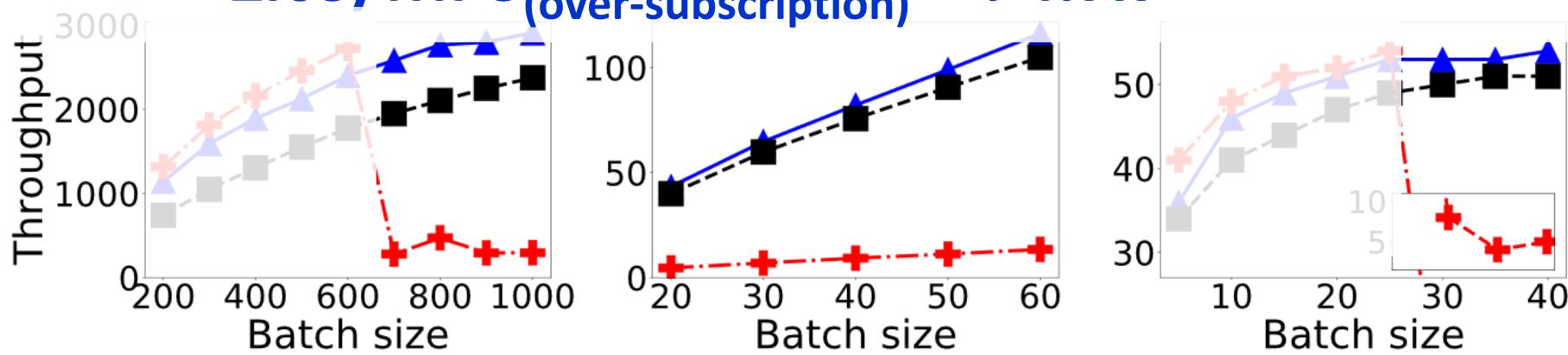


**Zico/MPS(no over-subscription): Similar throughput**

# Throughput: Identical Jobs



**Zico/MPS(no over-subscription): Similar throughput**
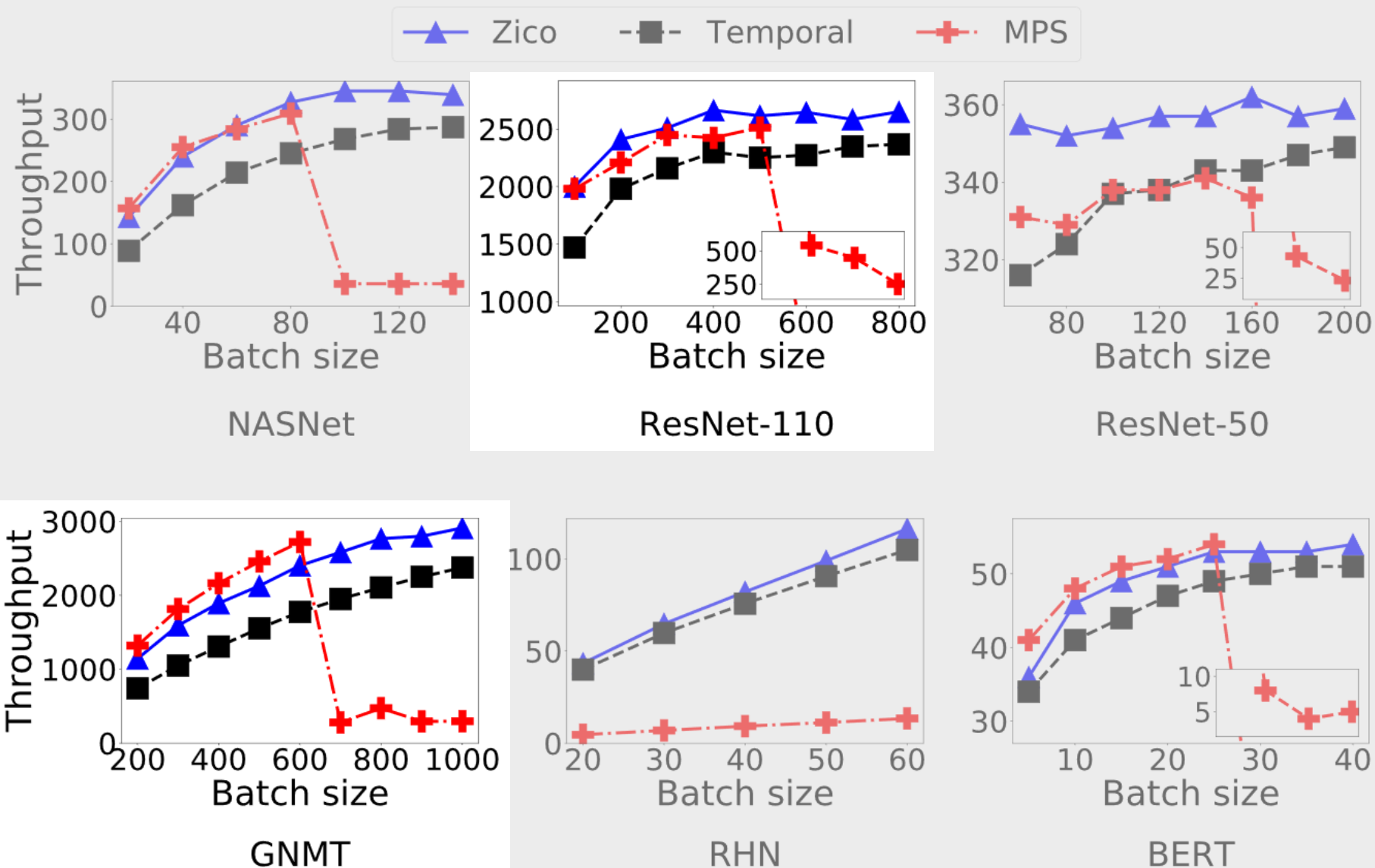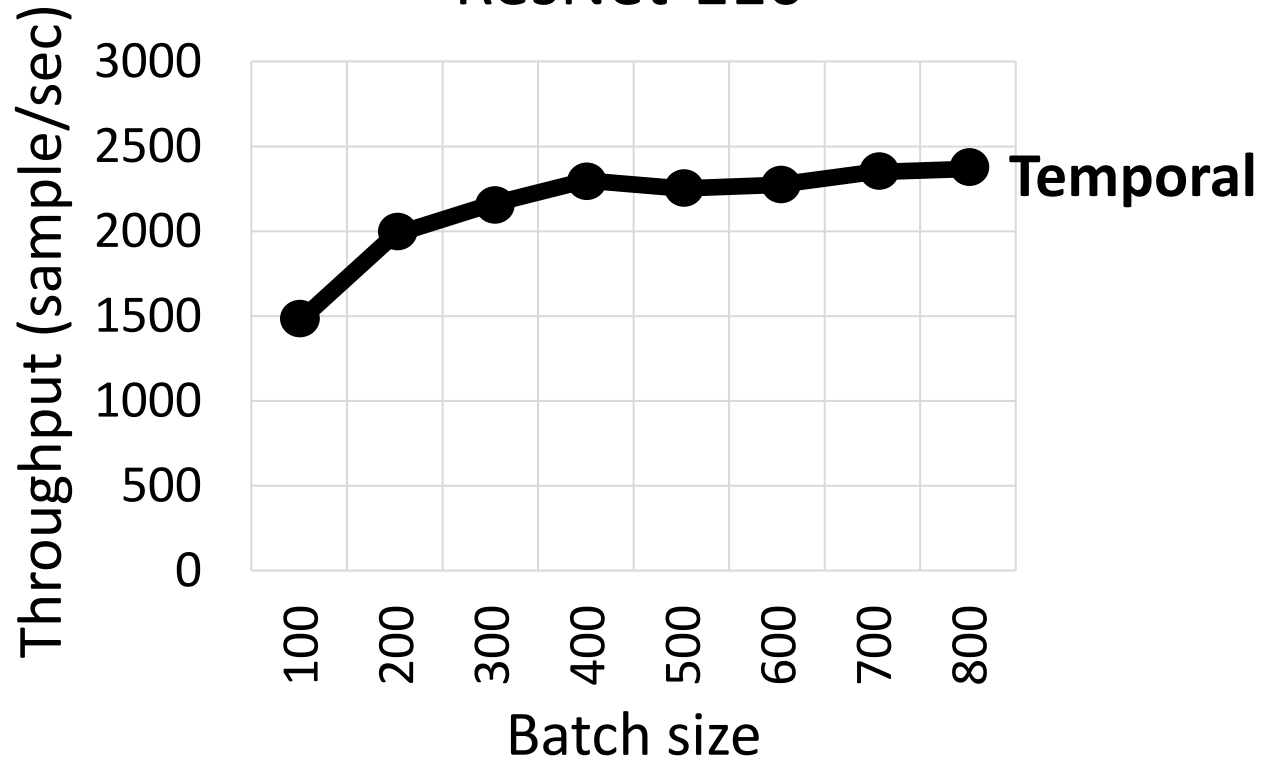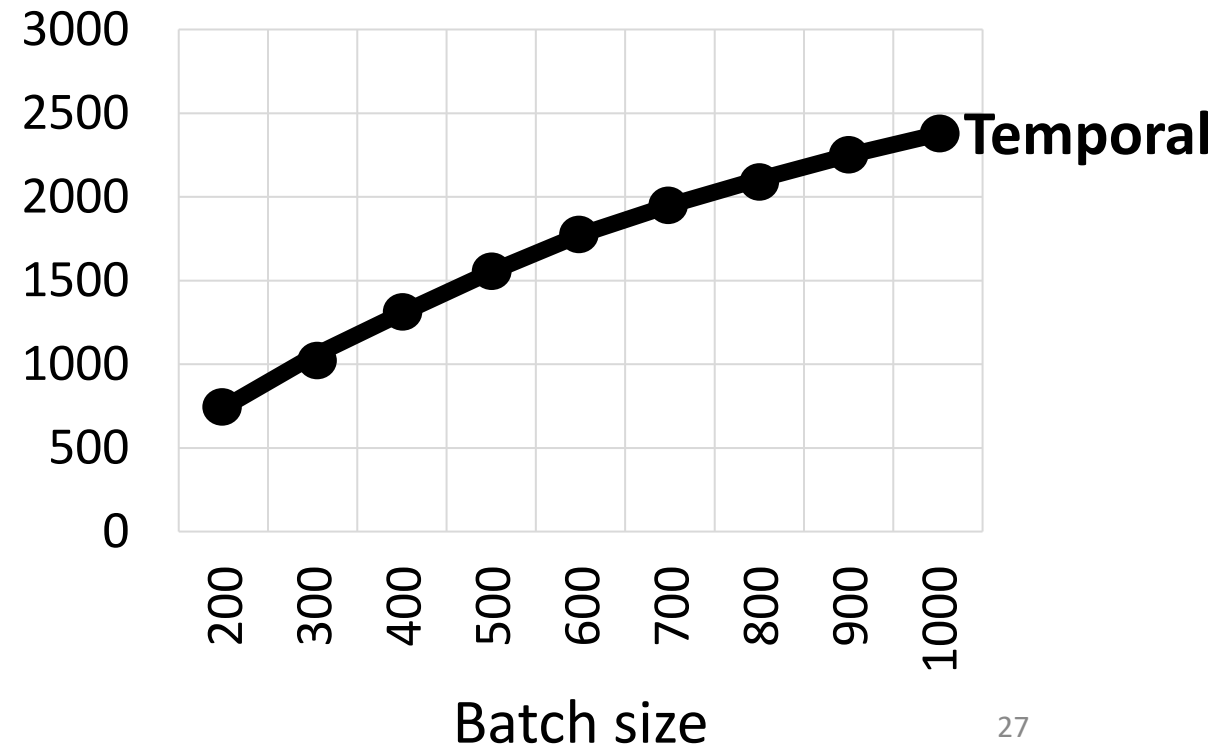**Zico/MPS(over-subscription): 4.7x**

# Throughput: Identical Jobs

# Throughput: Identical Models



ResNet-110 — Throughput (sample/sec) vs Batch size. GNMT — Throughput (sample/sec) vs Batch size. Both curves labeled "Temporal".
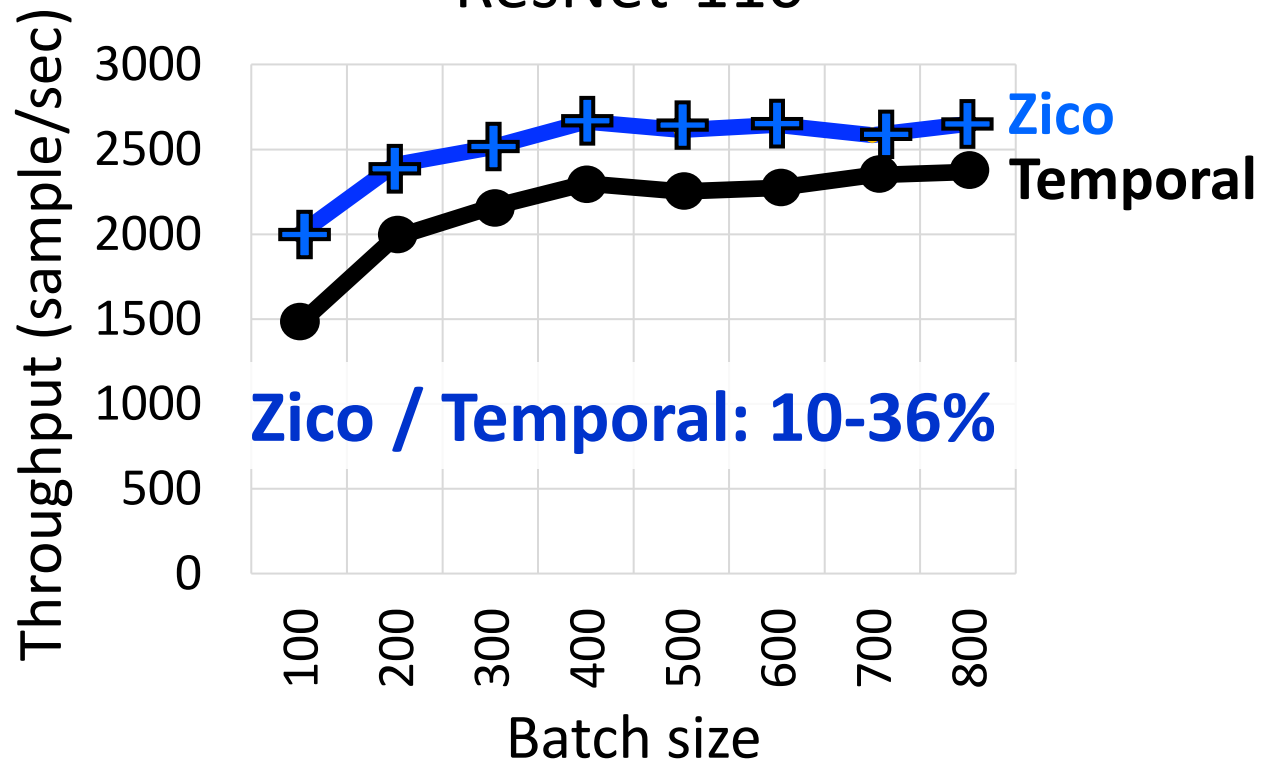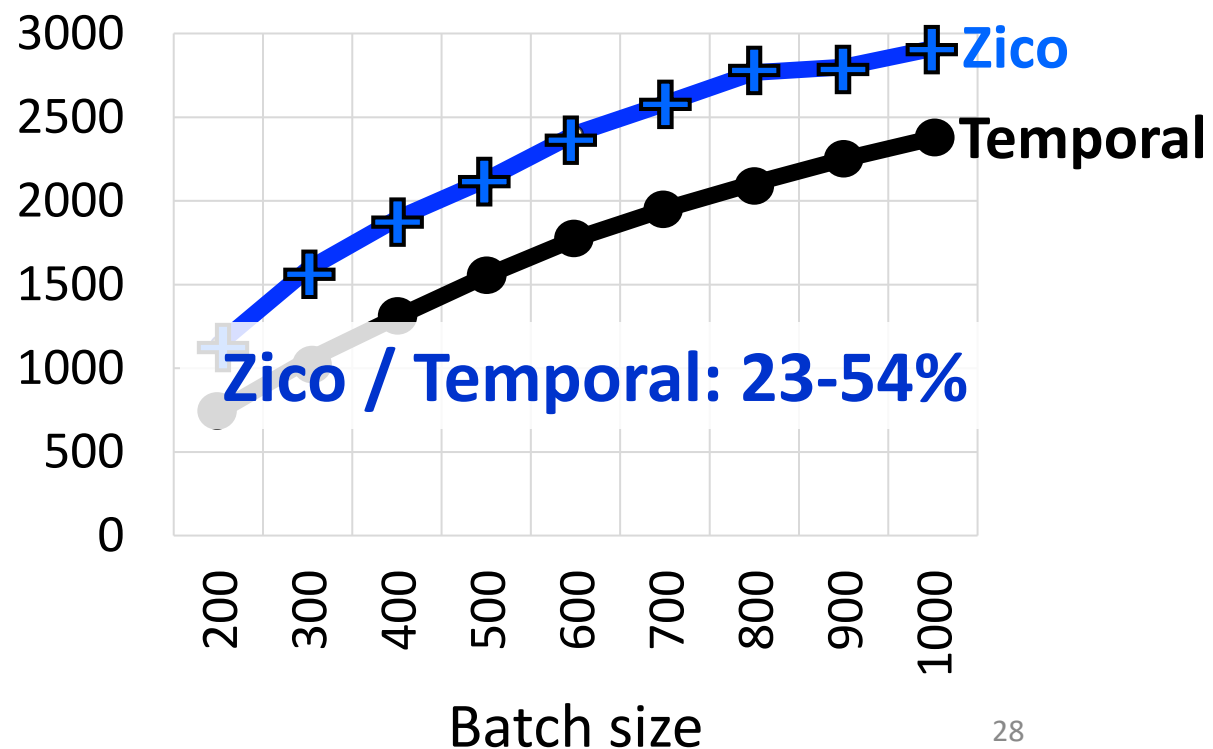
# Throughput: Identical Models

> *Temporal doesn't fully utilize GPU*
>
> *Zico always outperforms Temporal!*



ResNet-110

GNMT
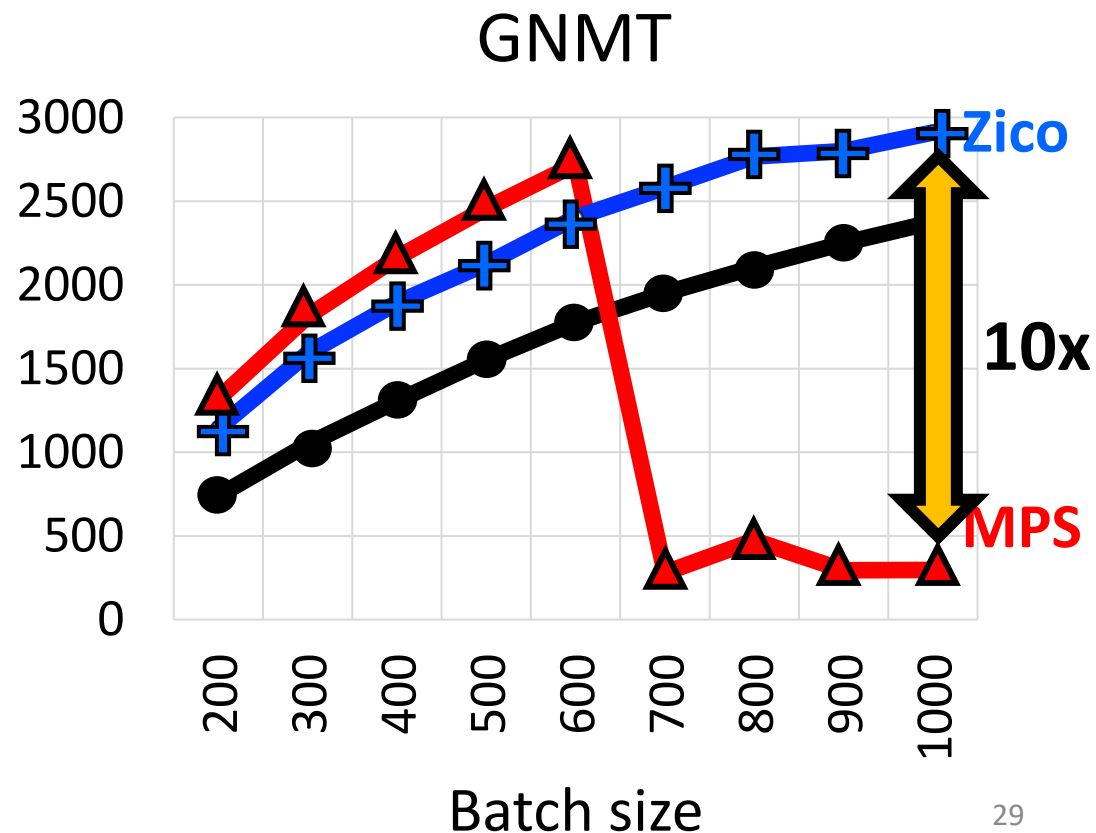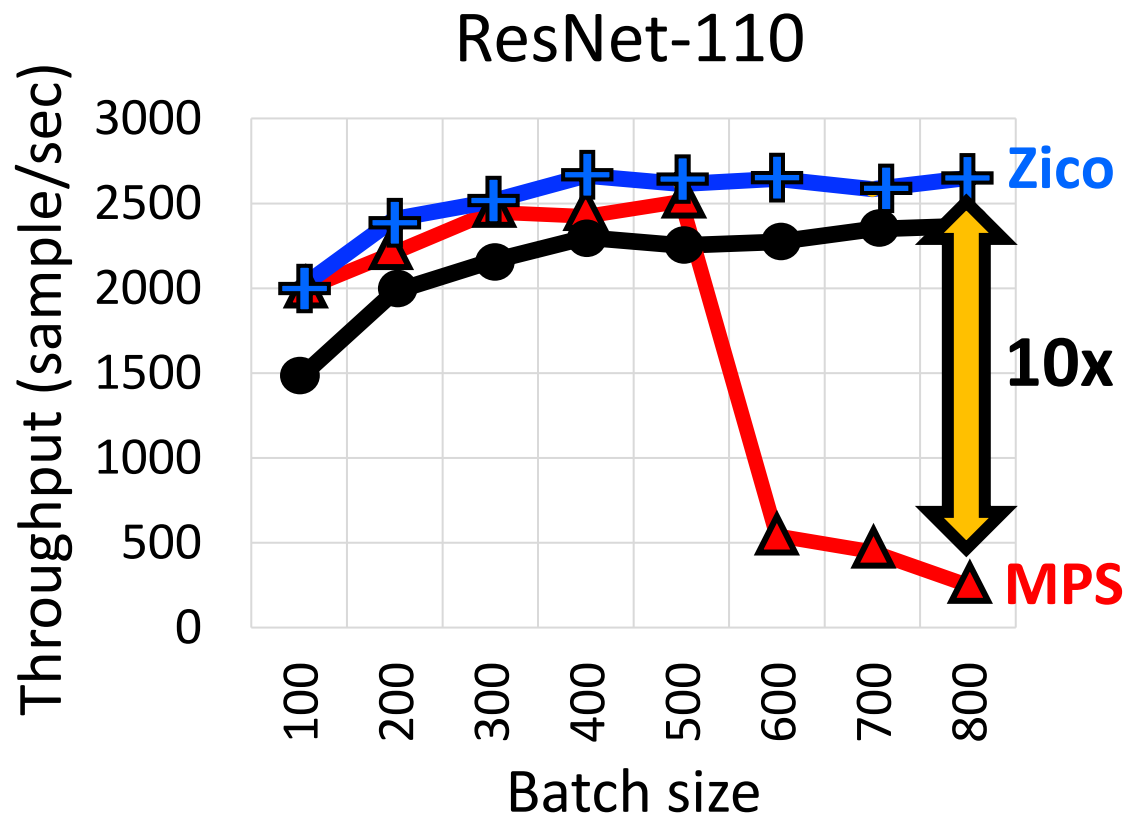
**Zico / Temporal: 10-36%**

**Zico / Temporal: 23-54%**
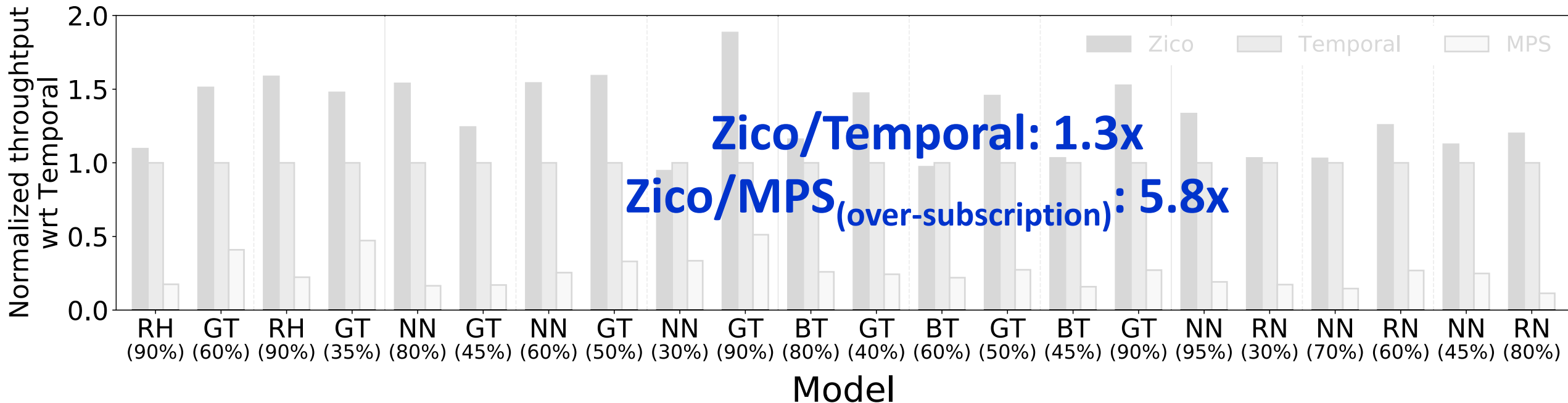
# Throughput: Identical Models

MPS suffers from memory over-subscription

Zico successfully co-locates two jobs w/o over-subscription
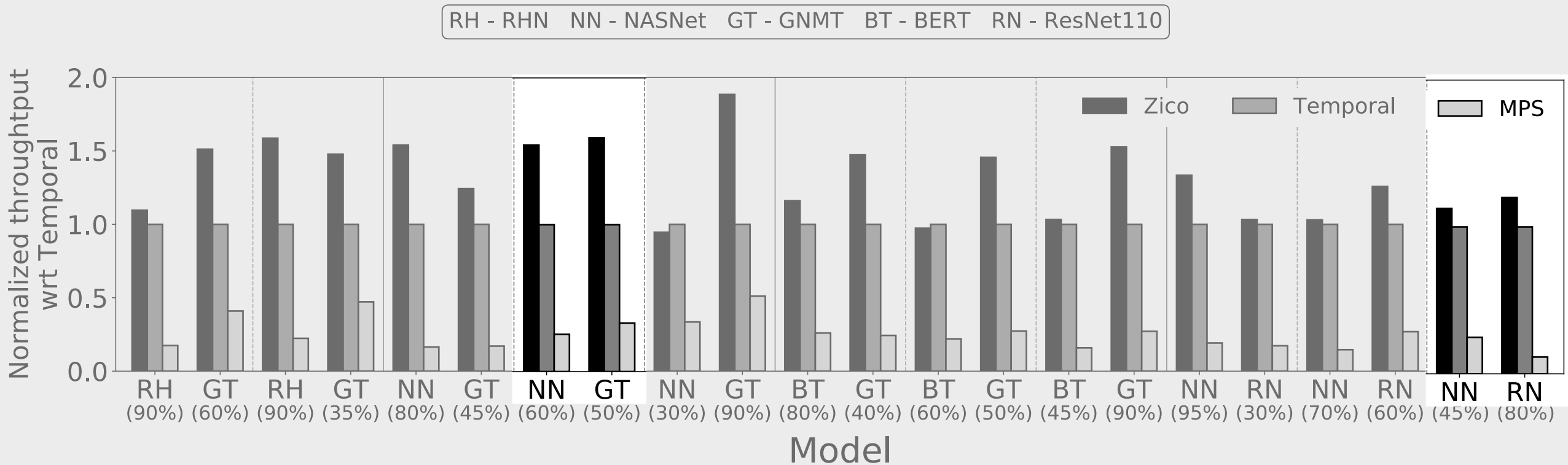


ResNet-110

GNMT

# Throughput: Non-identical Jobs

Note: In non-identical jobs experiment, MPS is set to always over-subscribe the memory.

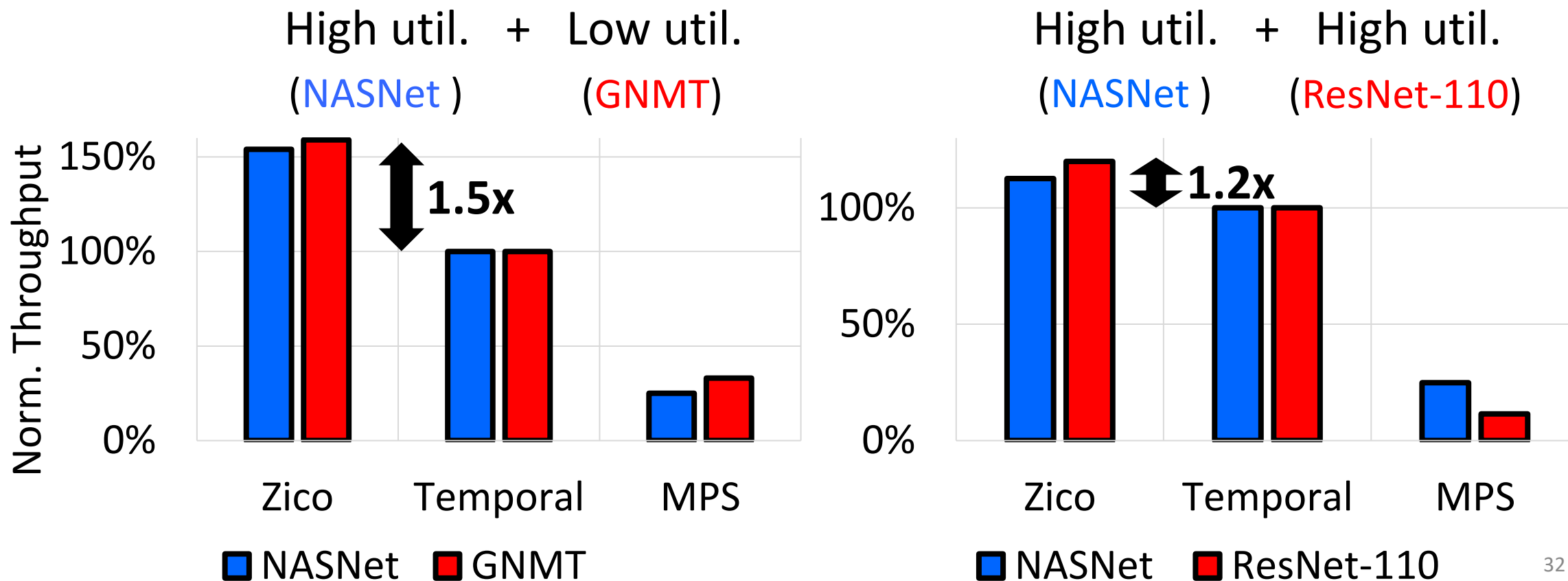RH - RHN   NN - NASNet   GT - GNMT   BT - BERT   RN - ResNet110



**Zico/Temporal: 1.3x**
**Zico/MPS(over-subscription): 5.8x**

# Throughput: Non-identical Jobs



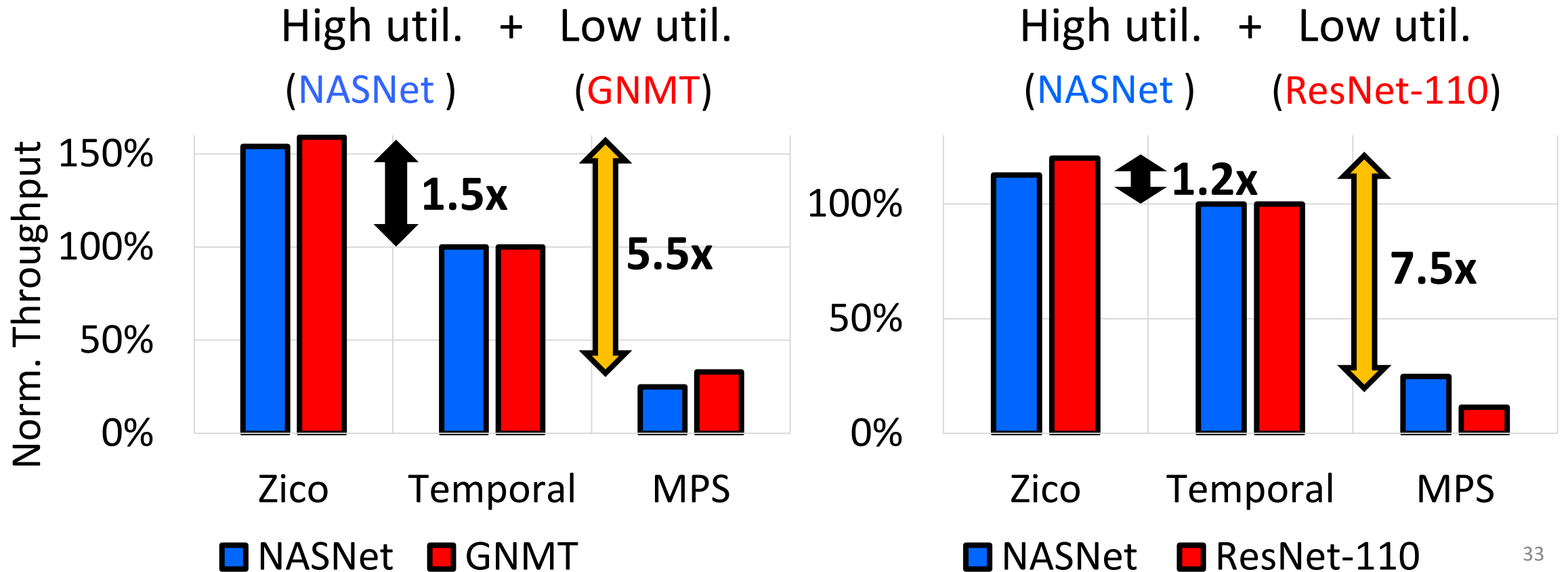RH - RHN   NN - NASNet   GT - GNMT   BT - BERT   RN - ResNet110

# Throughput: Non-identical Models

*Zico successfully co-locates non-identical models!*
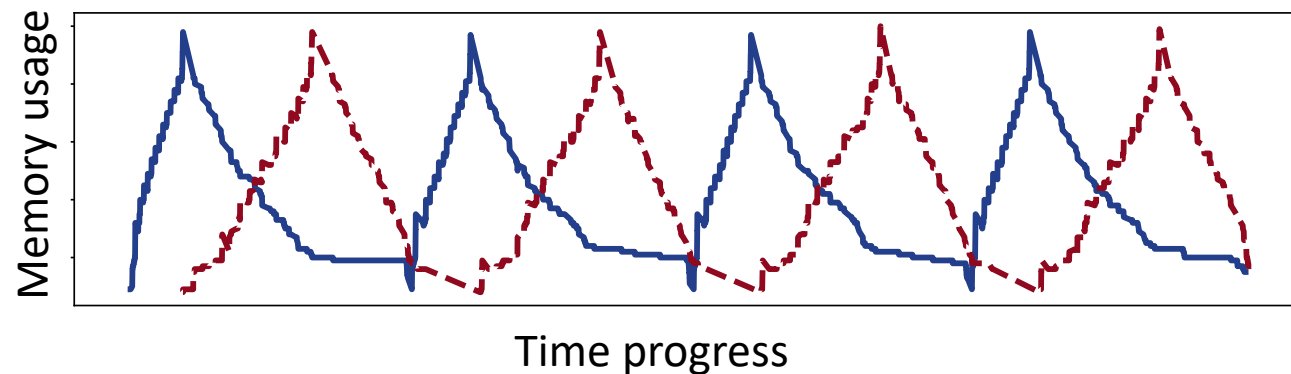
*More improvement when low utilization model is co-located*

# Throughput: Non-identical Jobs



High util.  +  Low util.
(NASNet )      (GNMT)

High util.  +  Low util.
(NASNet )      (ResNet-110)

# Scheduling Example

## Identical model co-location

BERT + BERT
(budge: 32GB)

# Scheduling Example

## Identical model co-location
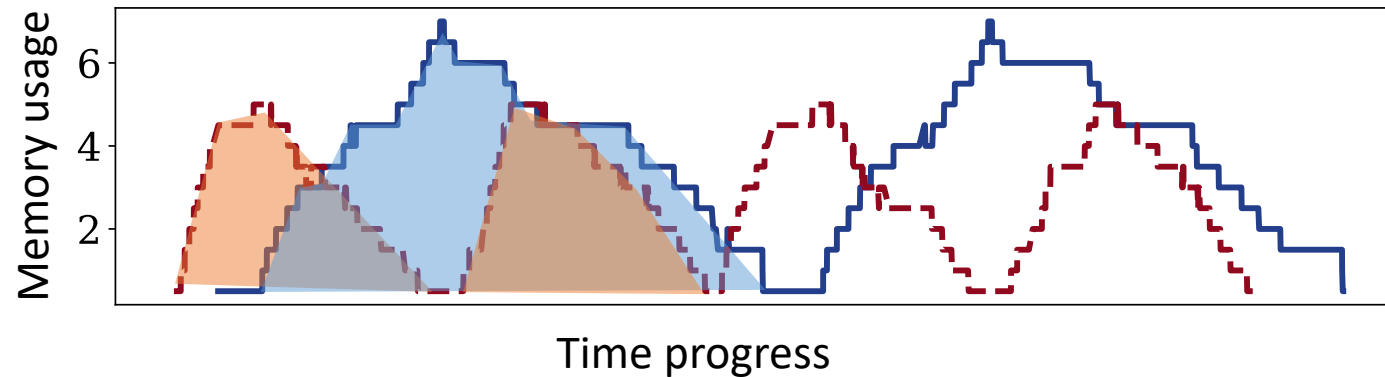
BERT + BERT
(budge: 32GB)



## Non-identical model co-location

NASNet + ResNet-110
(budge: 11GB)

# Summary

- Zico is the first introducing memory-aware scheduler

- Zico proposes widely applicable GPU sharing techniques for training