

deep learning with fewer resources: *collocation-aware GPU orchestration*

Pinar Tözün

Associate Professor, IT University of Copenhagen

pito@itu.dk, pinartozun.com, [@pinartozun](https://twitter.com/pinartozun)

UNIL Feb 10, 2026
HEIG-VD Feb 12, 2026

Innovationsfonden

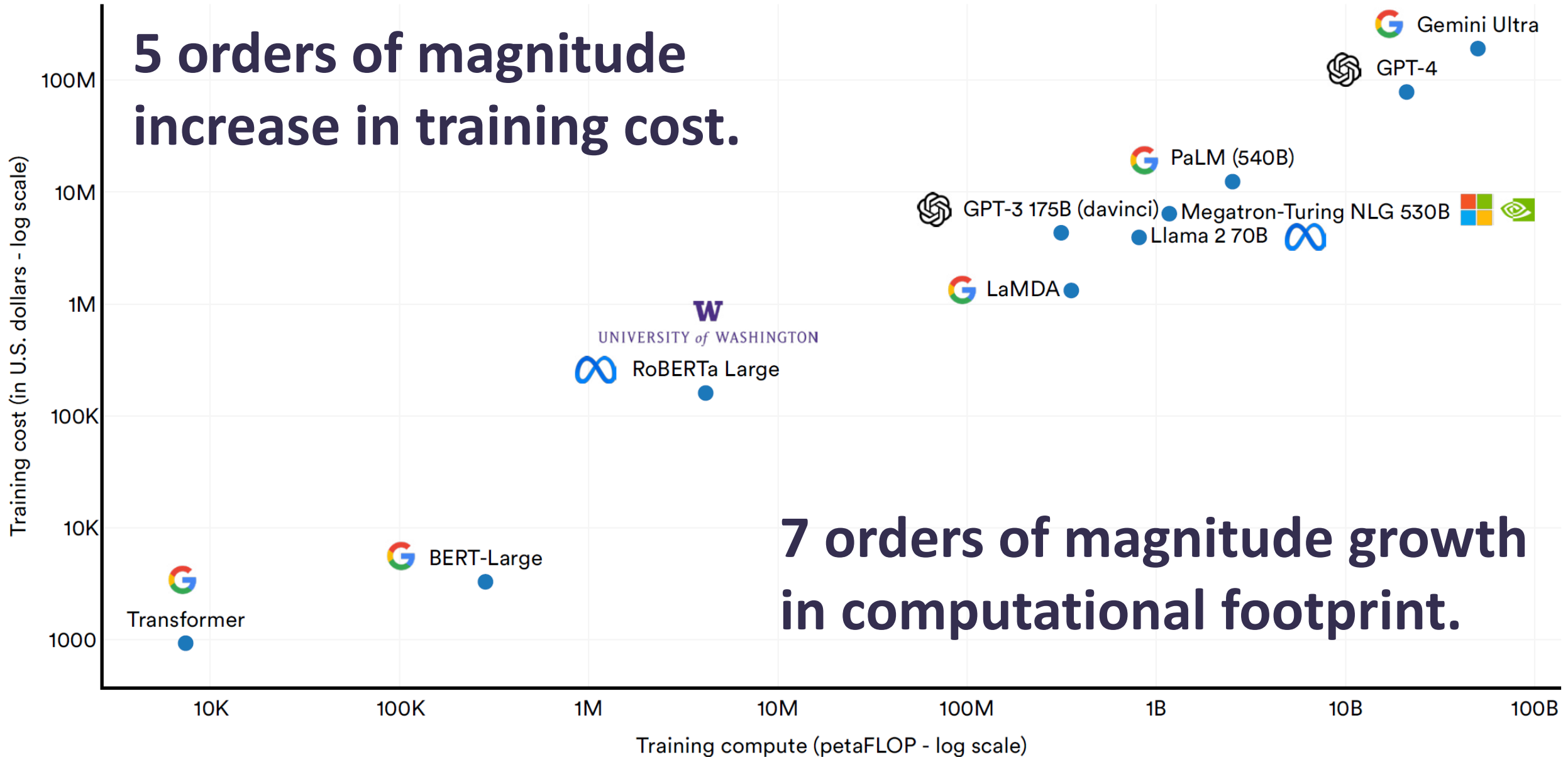


novo nordisk
foundation

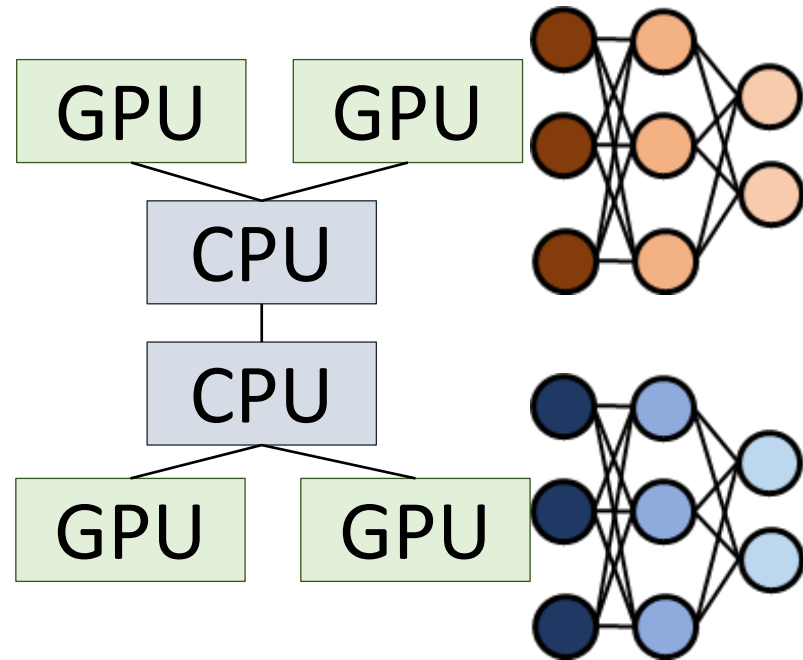
language model training (2017 – today)

**5 orders of magnitude
increase in training cost.**

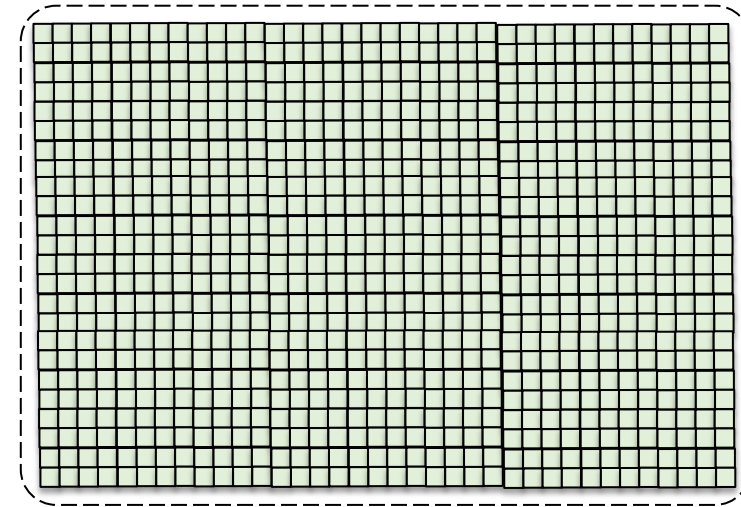
**7 orders of magnitude growth
in computational footprint.**



commodity hardware for deep learning



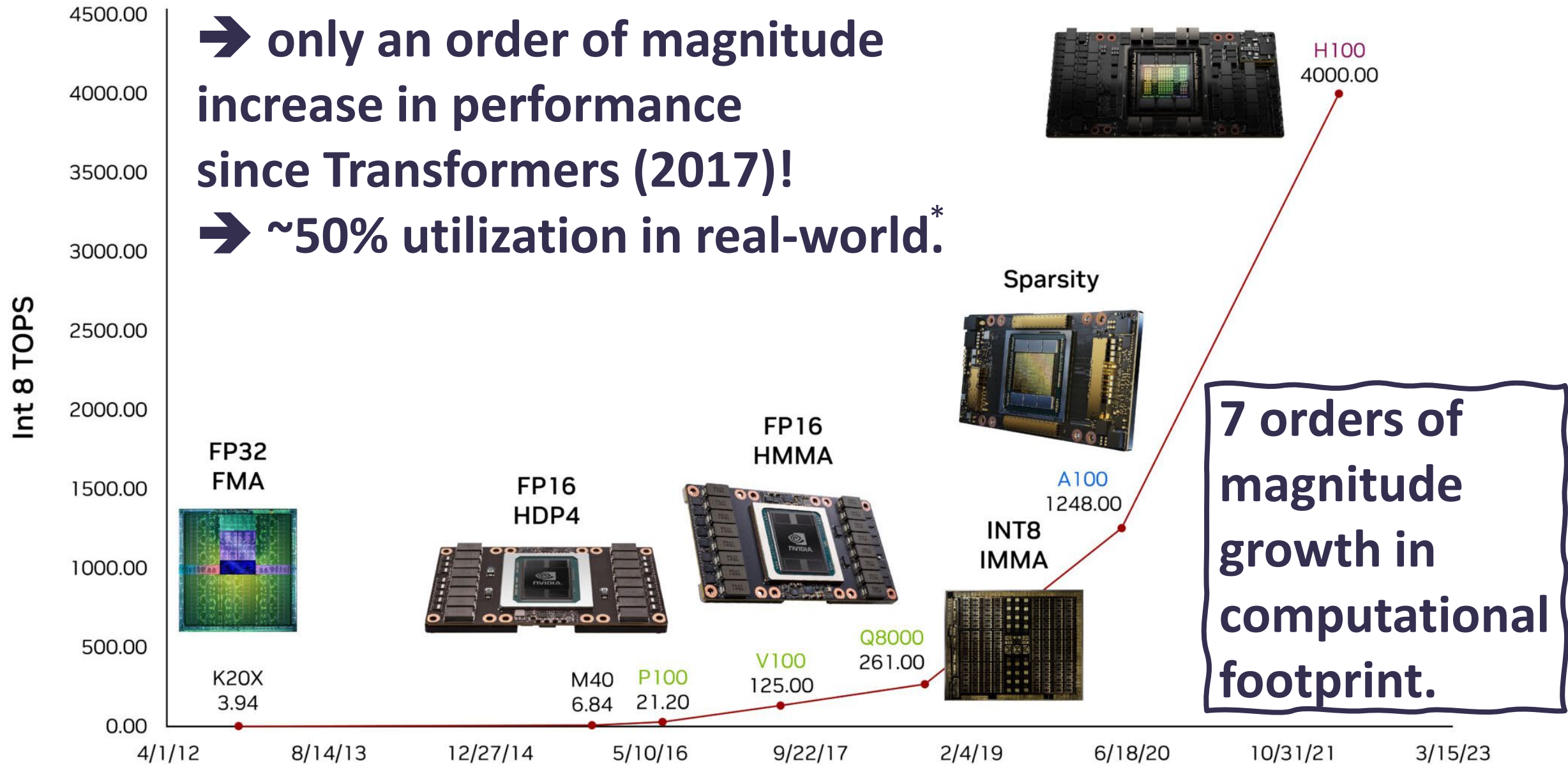
GPU



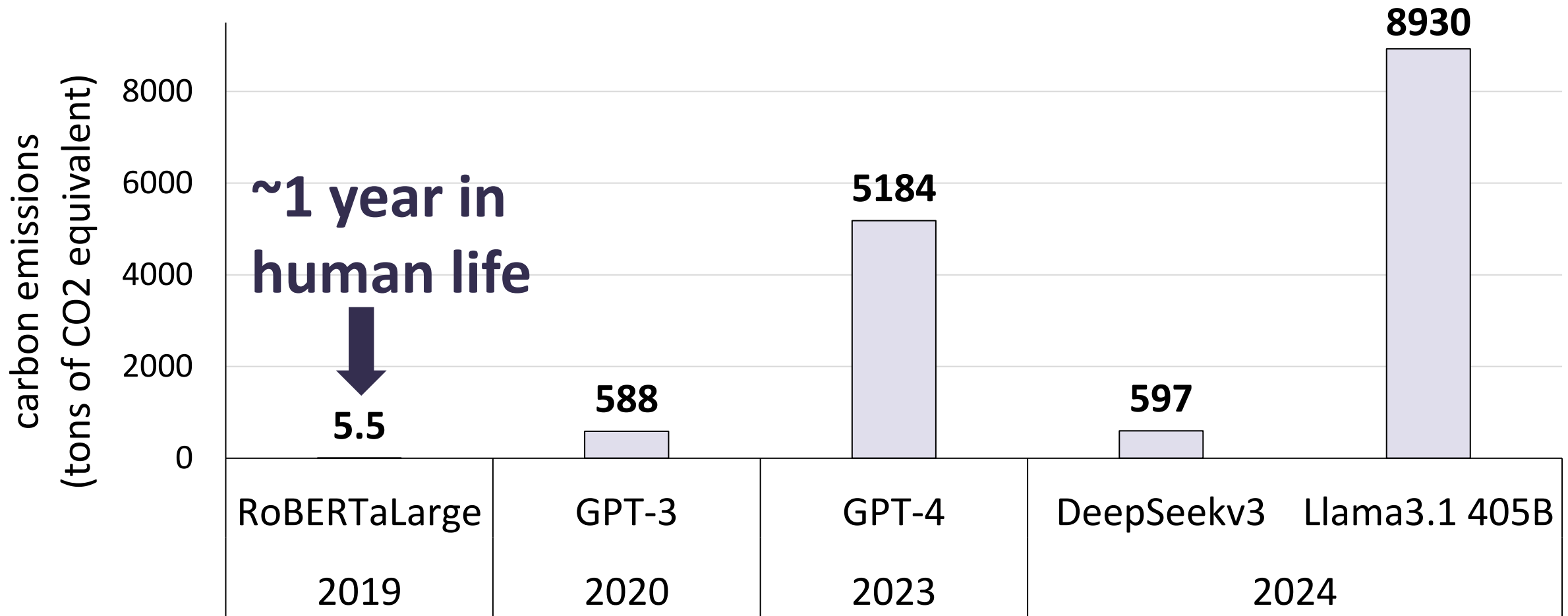
graphics processing unit

- many (simple) cores
- good for throughput-oriented & embarrassingly parallel tasks
 - good for deep learning
 - e.g., large matrix operations

NVIDIA GPUs (2012 – 2023)



carbon footprint of language model training



**can we do better while using fewer resources?
model accuracy cannot be the only metric to aim for!**

workload collocation for model training

- analysis of GPU resource sharing primitives

[An Analysis of Collocation on GPUs for Deep Learning Training](#)

Ties Robroek, Ehsan Yousefzadeh-Asl-Miandoab, Pinar Tözün.
EuroMLSys 2024

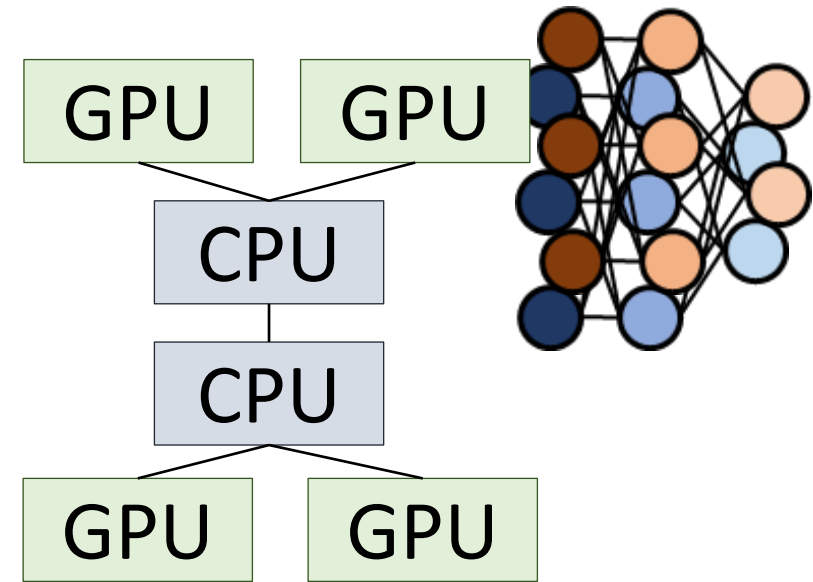
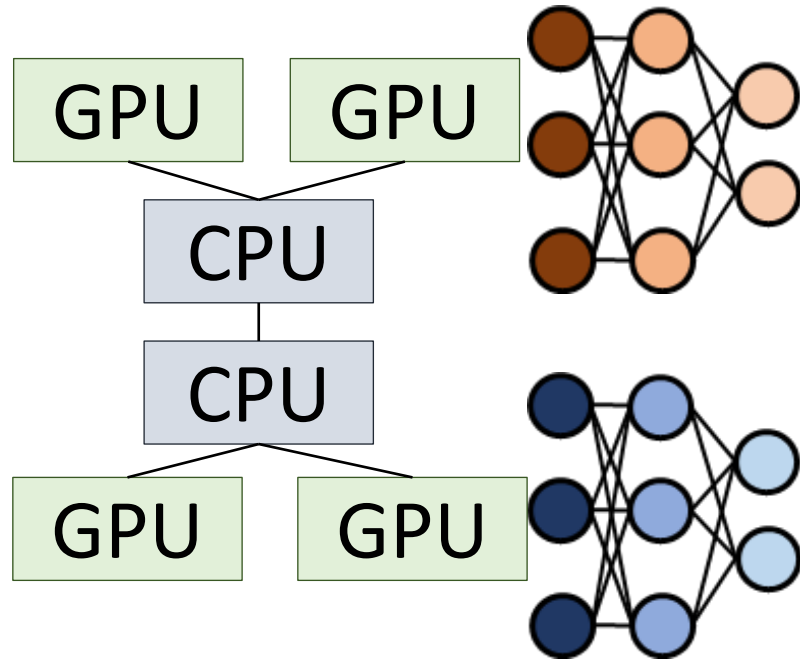
- collocation-aware resource management

[CARMA: Collocation-Aware Resource Manager](#)

Ehsan Yousefzadeh-Asl-Miandoab, Reza Karimzadeh, Bulat Ibragimov, Florina M Ciorba, Pinar Tözün.

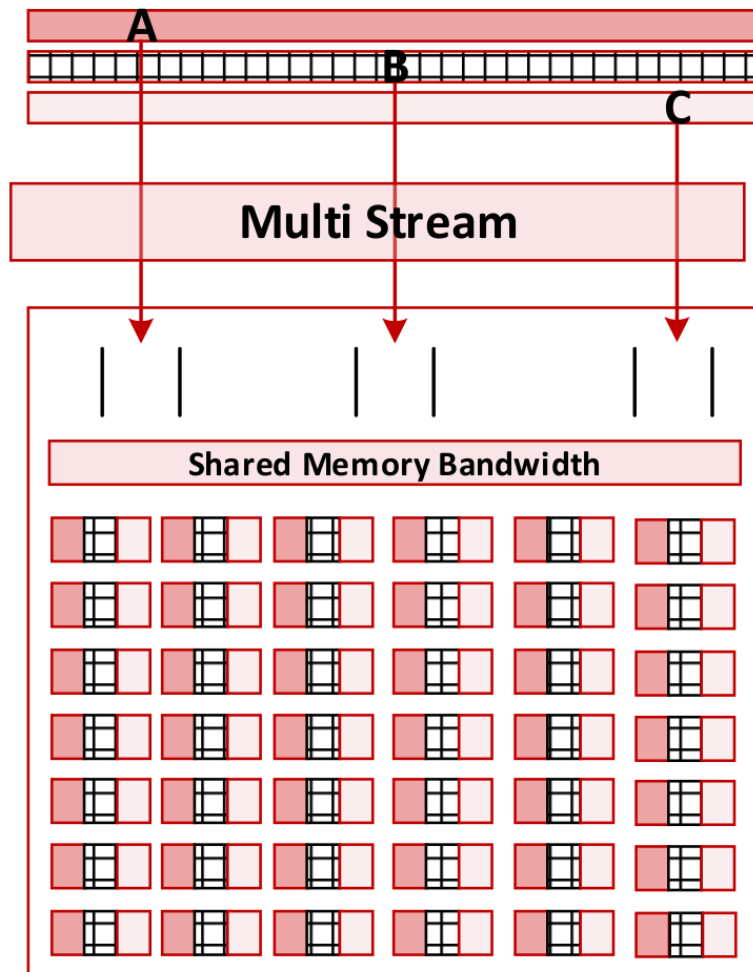


collocated training

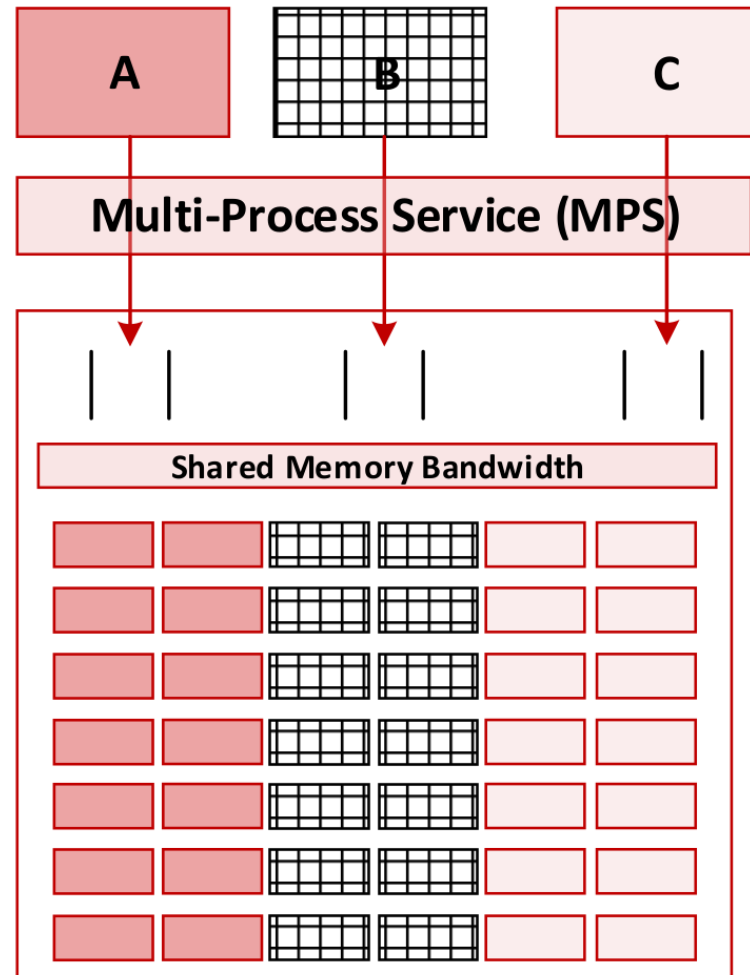


- training more models with fewer hardware resources
- if done well → better hardware utilization & reduces costs
- if not done well → interference & performance degradation

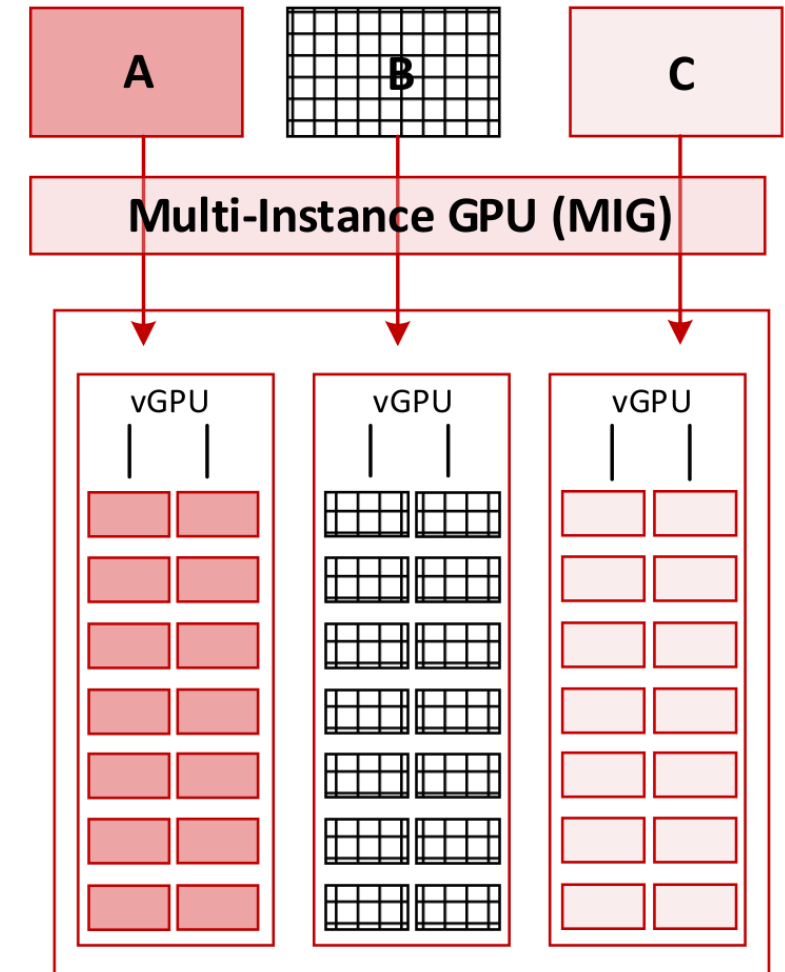
sharing resources on (NVIDIA) GPUs



- most straightforward
- time-multiplexing
- ✗ limited parallelism

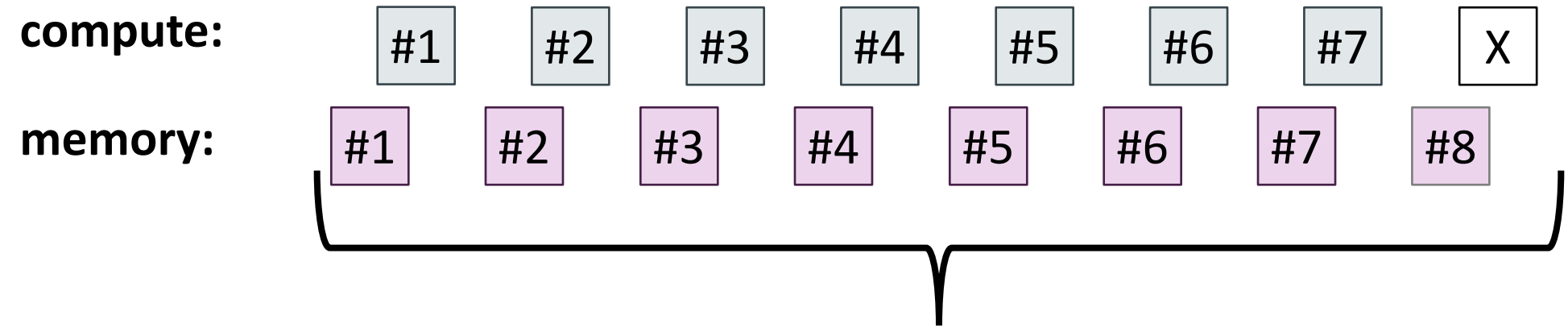






- finer-grained sharing
- ✗ higher chances of interference



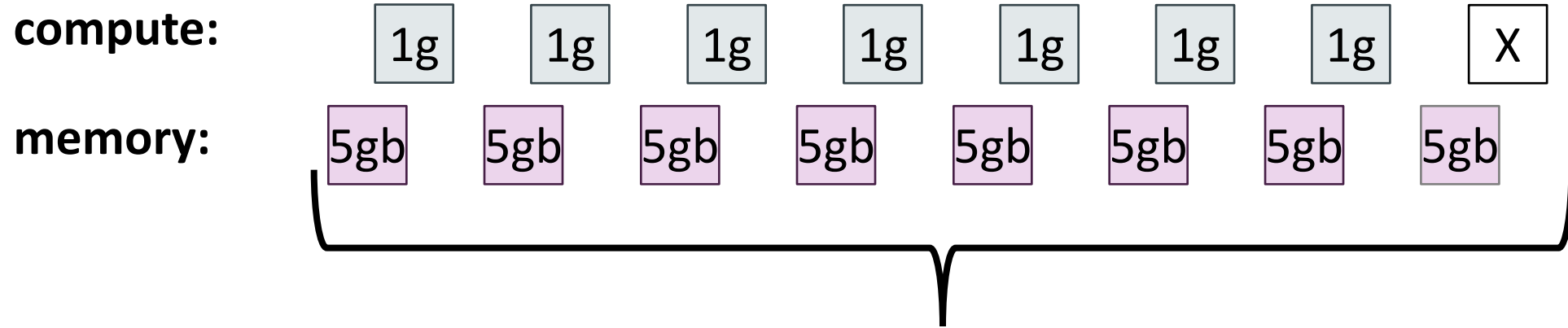
- hardware-support for resource split
- ✗ rigid partitioning

multi-instance GPU



-  1 compute unit
-  1 memory unit
-  unused available (memory/compute) unit
-  unavailable compute unit

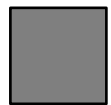
multi-instance GPU on A100 (40GB)



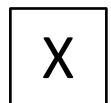
1 compute unit = 1g = 14 SMs



1 memory unit = 5GB



unused available (memory/compute) unit

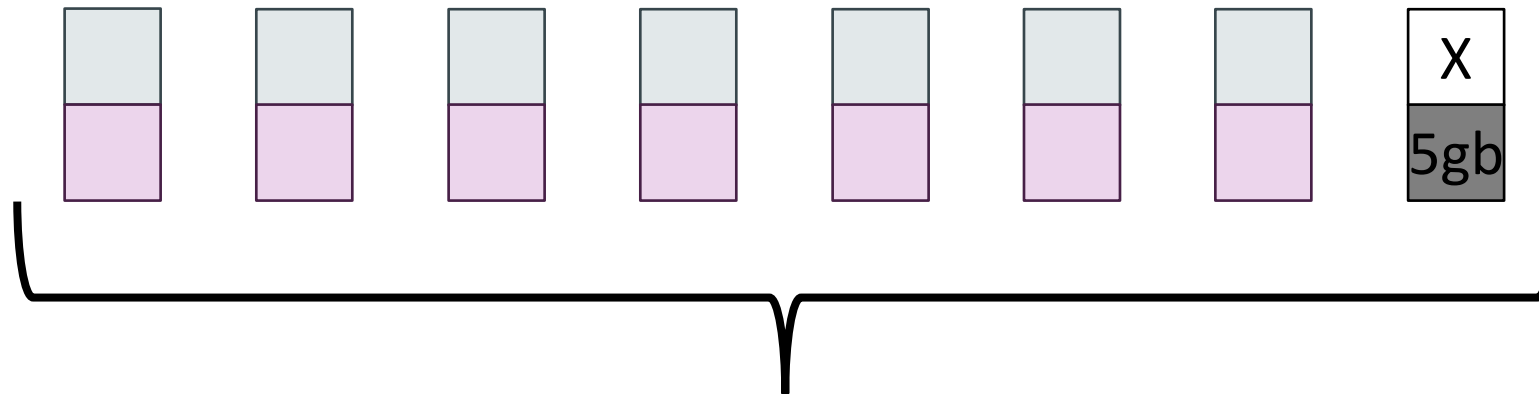


unavailable compute unit = 10 SMs (streaming multiprocessor)

multi-instance GPU on A100 (40GB)

compute:

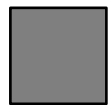
memory:



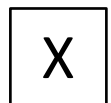
1 compute unit = 1g = 14 SMs



1 memory unit = 5GB



unused available (memory/compute) unit



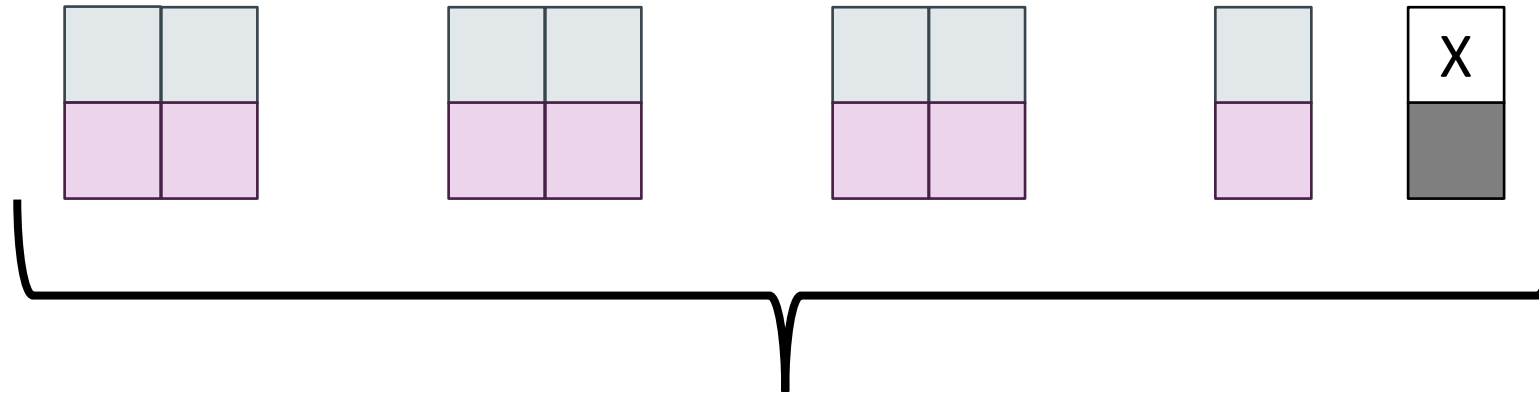
unavailable compute unit = 10 SMs (streaming multiprocessor)

GPU

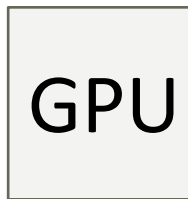
multi-instance GPU on A100 (40GB)

compute:

memory:



1 compute unit = 1g = 14 SMs



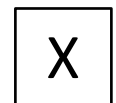
GPU



1 memory unit = 5GB



unused available (memory/compute) unit

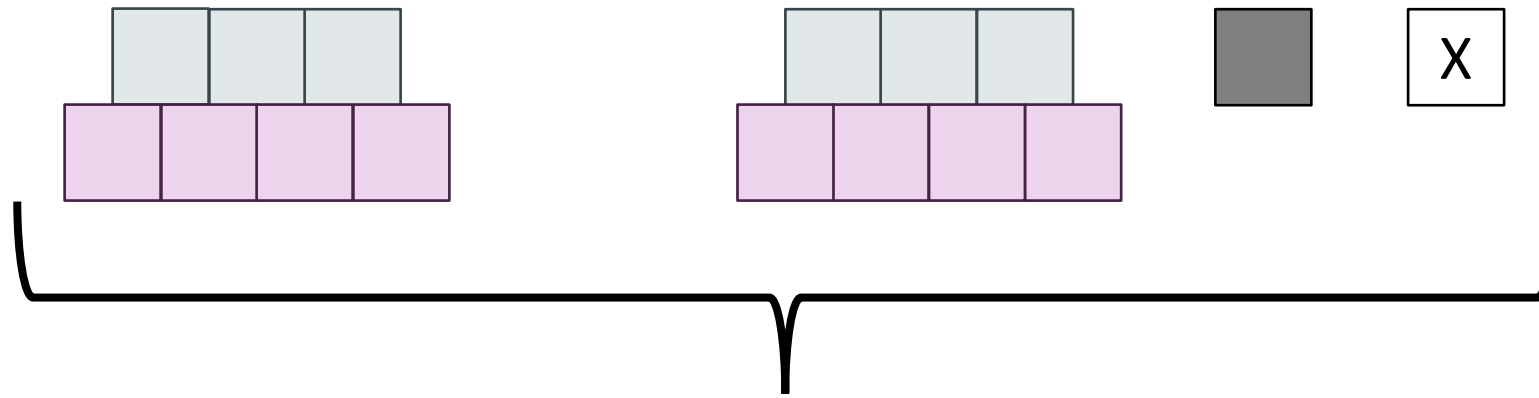


unavailable compute unit = 10 SMs (streaming multiprocessor)

multi-instance GPU on A100 (40GB)

compute:

memory:



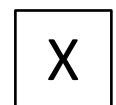
1 compute unit = 1g = 14 SMs



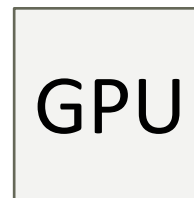
1 memory unit = 5GB



unused available (memory/compute) unit



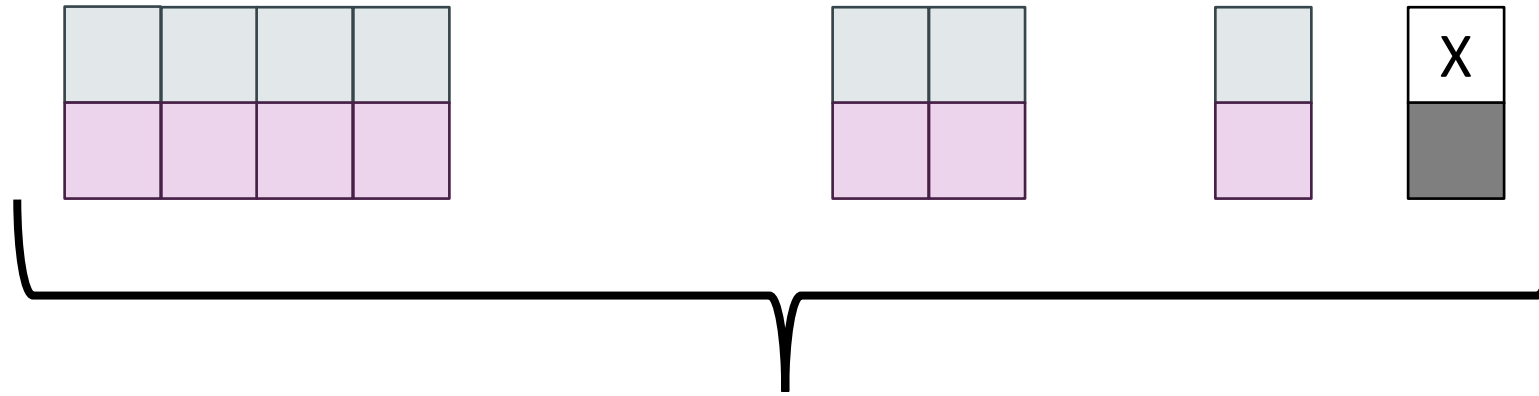
unavailable compute unit = 10 SMs (streaming multiprocessor)



multi-instance GPU on A100 (40GB)

compute:

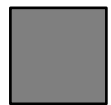
memory:



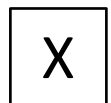
1 compute unit = 1g = 14 SMs



1 memory unit = 5GB



unused available (memory/compute) unit



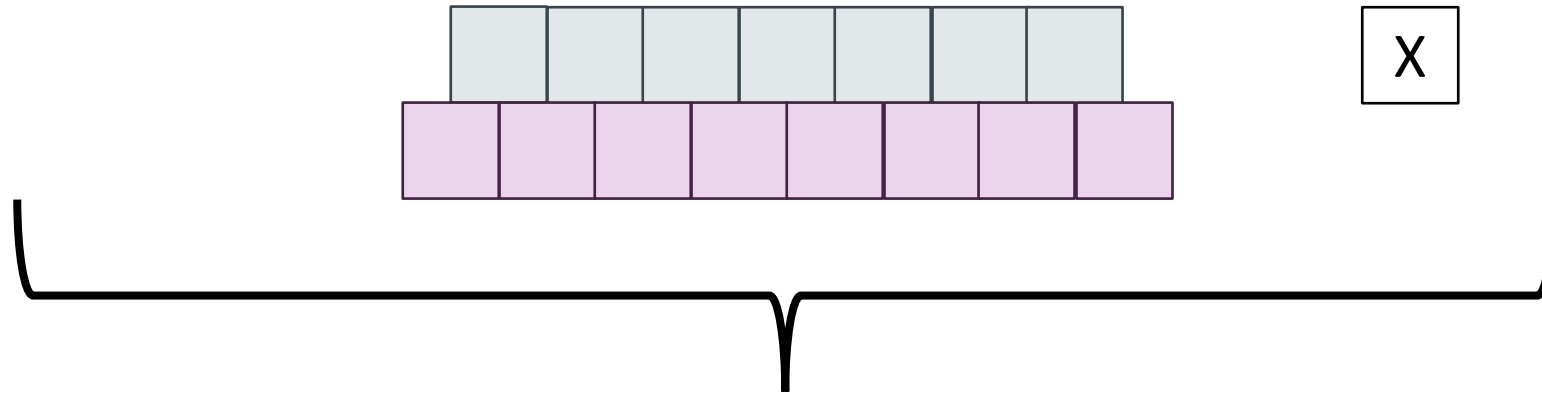
unavailable compute unit = 10 SMs (streaming multiprocessor)

GPU

multi-instance GPU on A100 (40GB)

compute:

memory:



1 compute unit = 1g = 14 SMs



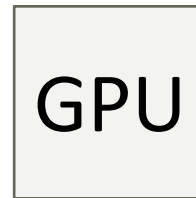
1 memory unit = 5GB



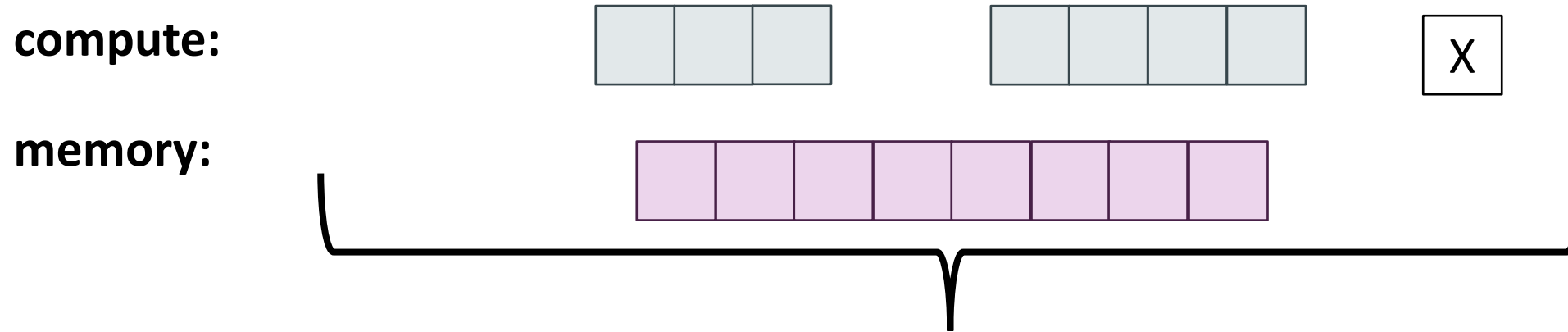
unused available (memory/compute) unit







unavailable compute unit = 10 SMs (streaming multiprocessor)



multi-instance GPU on A100 (40GB)



-  1 compute unit = 1g = 14 SMs
-  1 memory unit = 5GB
-  unused available (memory/compute) unit
-  unavailable compute unit = 10 SMs (streaming multiprocessor)

GPU

performance impact of collocation

NVIDIA DGX Station A100

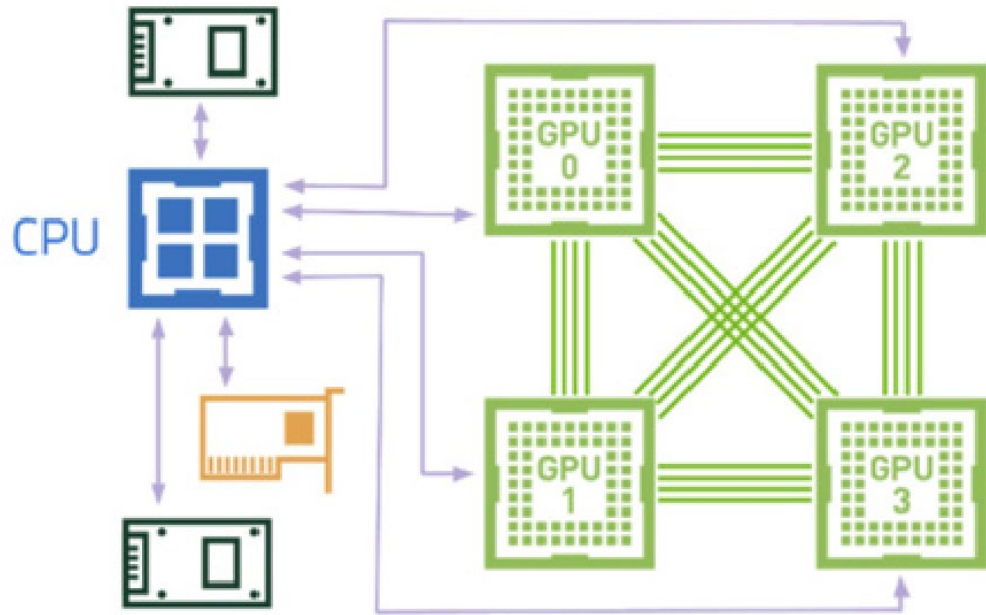


figure [source](#)

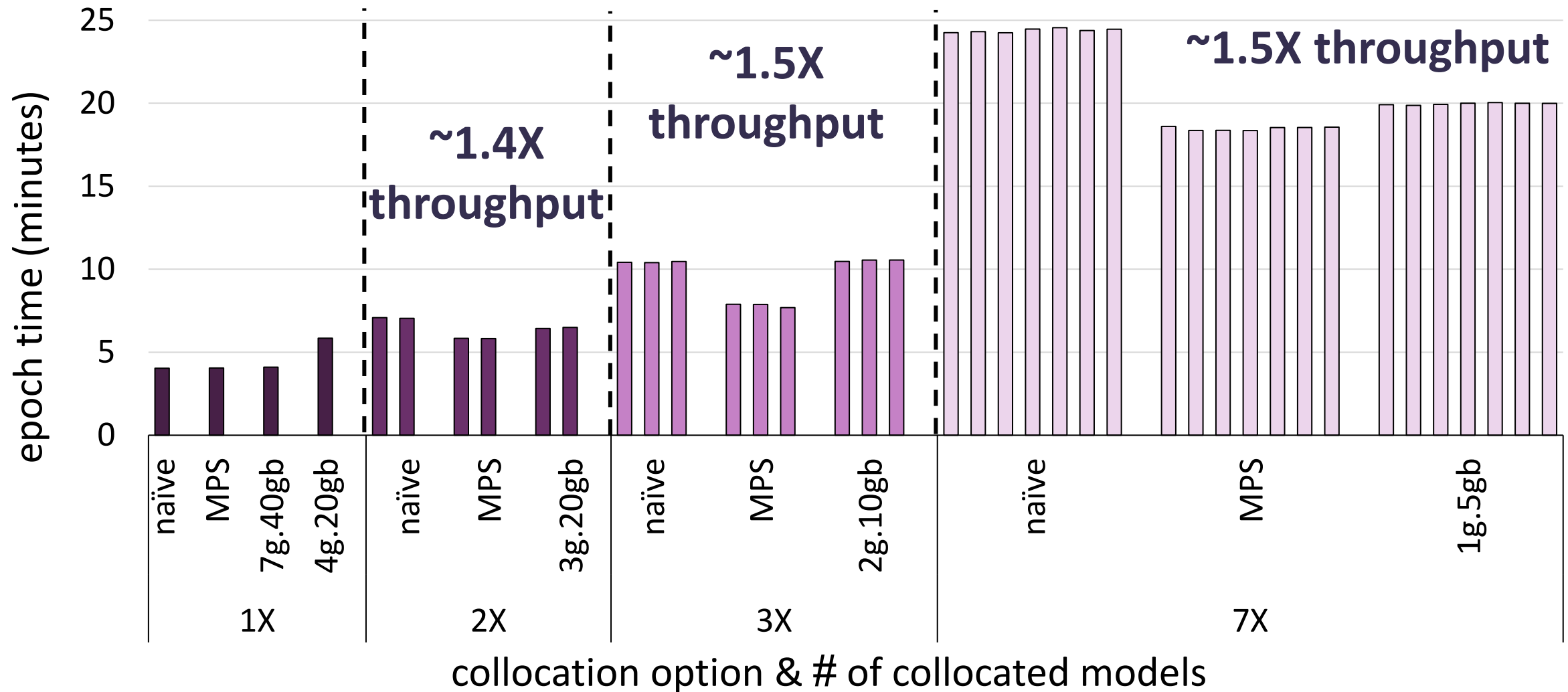
 Display GPU
  NVMe
  PCIe
  NVLink

CPU = AMD 7742 – 512 GB RAM
 64 physical cores
 GPU = NVIDIA A100 – 40 GB RAM

workloads	model	batch size	dataset
small	ResNet26 EfficientNet	128	CIFAR-10
medium	ResNet50 EfficientNet	128	downsampled ImageNet*
large	ResNet152 CaiT	32 128	ImageNet (2012)
xlarge	DLRM	1	Criteo Terabyte

- image models: CNN & transformers recommender model
- on single GPU with PyTorch v2.0
- results reported from 2nd epoch of training

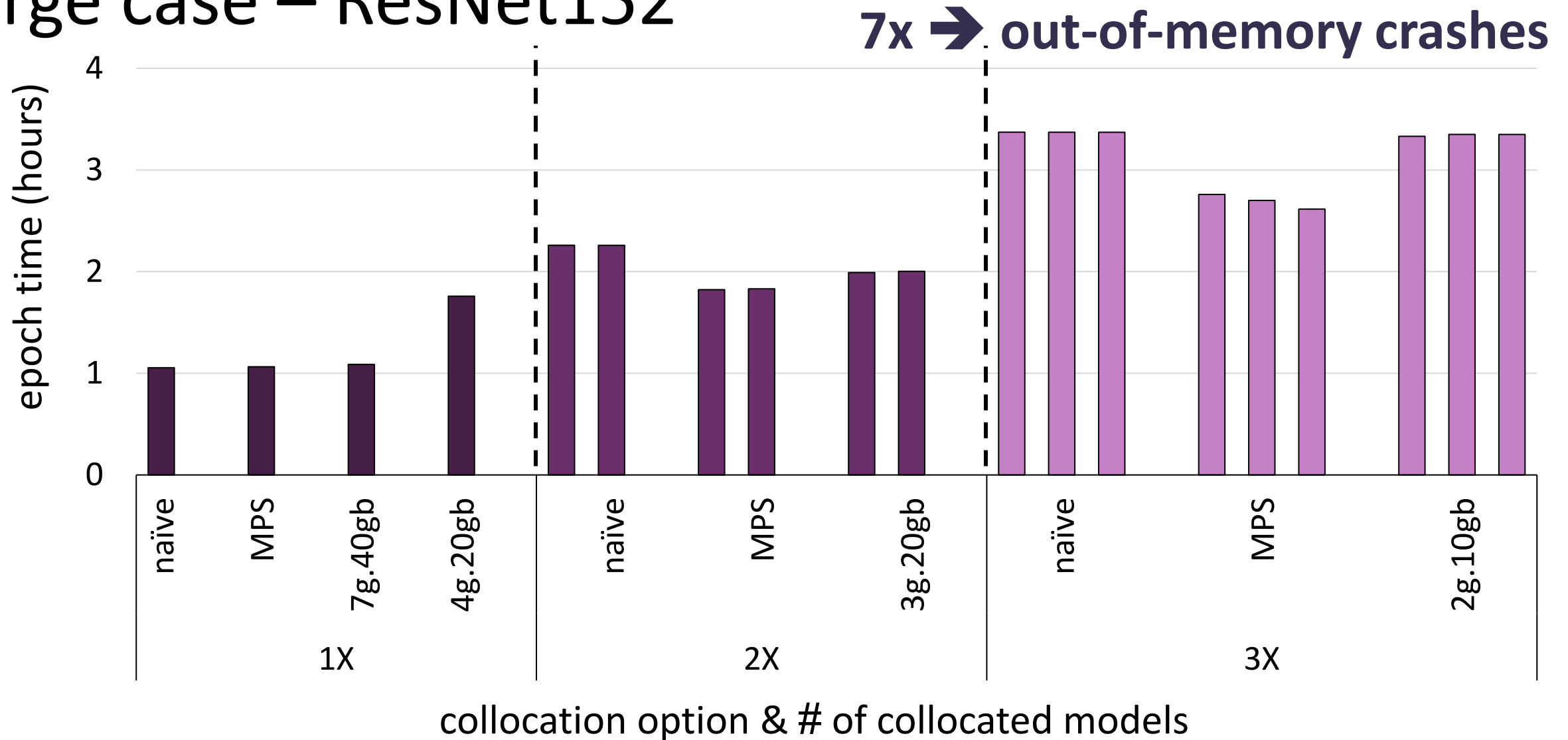
medium case – ResNet50



still some throughput benefits

but diminishing returns for increased collocation

large case – ResNet152



**no more throughput benefits – 80% utilization when training alone
better to collocate with smaller or less compute heavy tasks**

mixed workloads: compute- & memory-heavy

	DLRM time per training block	ResNet152 time per epoch	sm activity	memory footprint
DLRM alone			5%	29.14 GB
ResNet152 alone			82%	8.47 GB

mixed workloads: compute- & memory-heavy

	DLRM time per training block	ResNet152 time per epoch	sm activity	memory footprint
DLRM alone			5%	29.14 GB
ResNet152 alone			82%	8.47 GB
MPS			81%	37.62 GB

mixed workloads: compute- & memory-heavy

	DLRM time per training block	ResNet152 time per epoch	sm activity	memory footprint
DLRM alone	5.36 h	-	5%	29.14 GB
ResNet152 alone	-	1.05 h	82%	8.47 GB
MPS			81%	37.62 GB

mixed workloads: compute- & memory-heavy

	DLRM time per training block	ResNet152 time per epoch	sm activity	memory footprint
DLRM alone	5.36 h	-	5%	29.14 GB
ResNet152 alone	-	1.05 h	82%	8.47 GB
MPS	5.57 h (+5%)	1.10 h (+4%)	81%	37.62 GB

**collocation can lead to (almost) free lunch
when large models stress different hardware resources**

workload collocation for model training

- analysis of GPU resource sharing primitives

[An Analysis of Collocation on GPUs for Deep Learning Training](#)

Ties Robroek, Ehsan Yousefzadeh-Asl-Miandoab, Pinar Tözün.
EuroMLSys 2024

- collocation-aware resource management

[CARMA: Collocation-Aware Resource Manager](#)

Ehsan Yousefzadeh-Asl-Miandoab, Reza Karimzadeh, Bulat Ibragimov, Florina M Ciorba, Pinar Tözün.

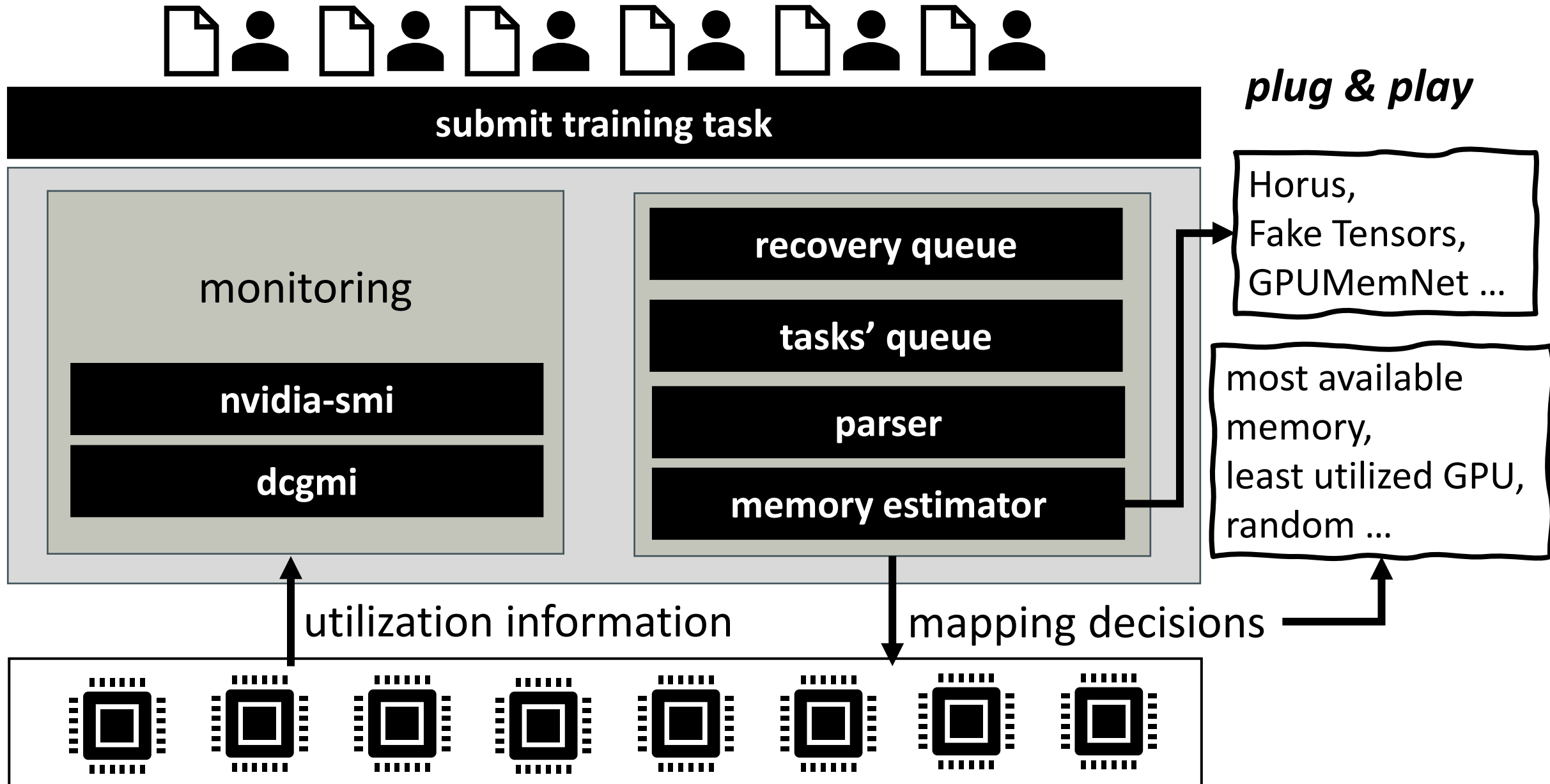


collocation-aware resource management

requirements:

- (1) don't overload the GPU compute
→ degrades performance**
- (2) minimize & recover from
out-of-memory crashes**

CARMA: collocation-aware resource manager



CARMA evaluation

trace mix

- **heavy:** 1-2GPUs, epoch time: 7.5mins to >1hour, memory: 9GB to 30GB
XLNet (base & large), BERT (base & large), GPT2 (large)
- **medium:** 1 GPU, epoch time: 1min to >1hour, memory: 1GB to 30GB
EfficientNet, ResNet (50), MobileNet, VGG ...
- **light:** 1 GPU, epoch time: up to 1min, memory: up to 1GB
ResNet (18, 34), MobileNet (small) ...

2 traces of 60 training tasks:

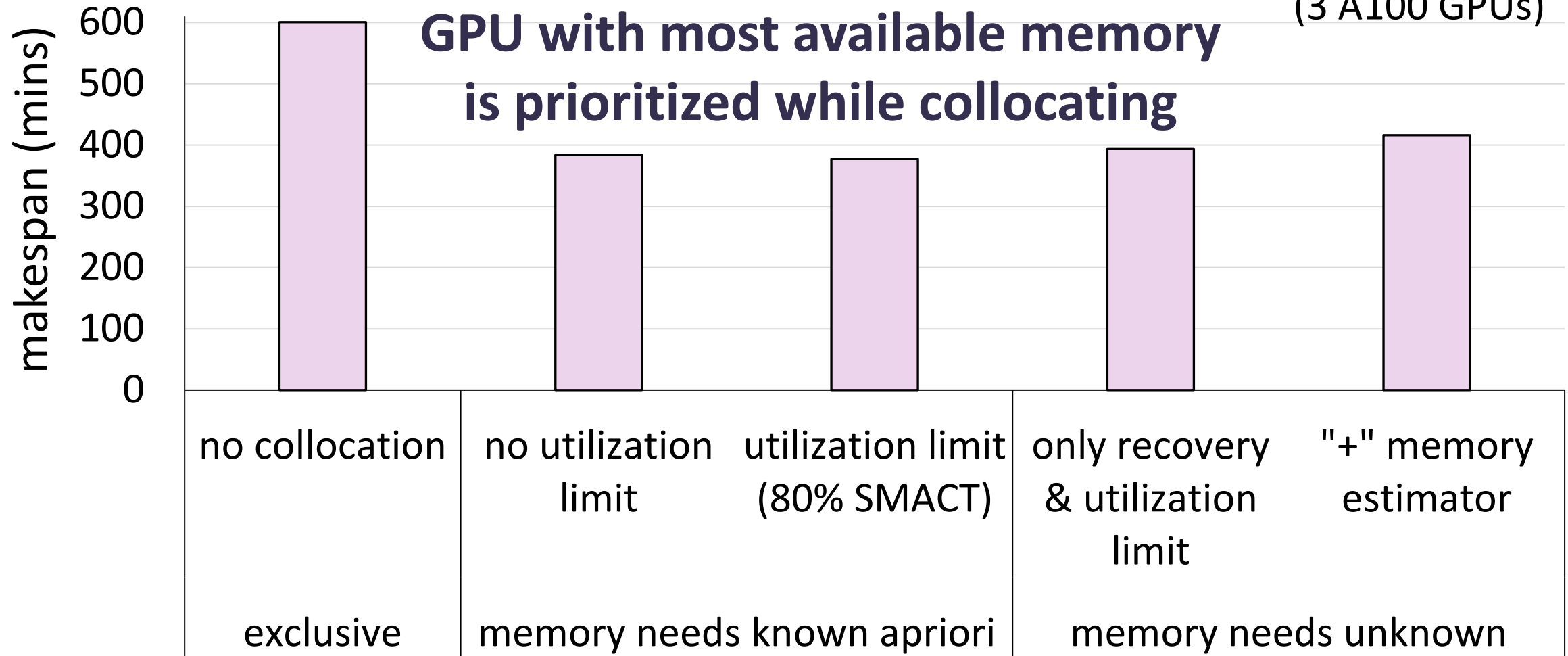
- 30% light, 60% medium/heavy, and 10% heavy 2 GPU models

based on real-world workload trace analysis

- task submission times: [“Analysis of Large-Scale Multi-Tenant GPU Clusters for DNN Training Workloads.”](#) ATC 2019
- task distribution: [“An Empirical Study on Low GPU Utilization of Deep Learning Jobs.”](#) TPDS 2022

CARMA on a training workload trace

on NVIDIA DGX Station
(3 A100 GPUs)



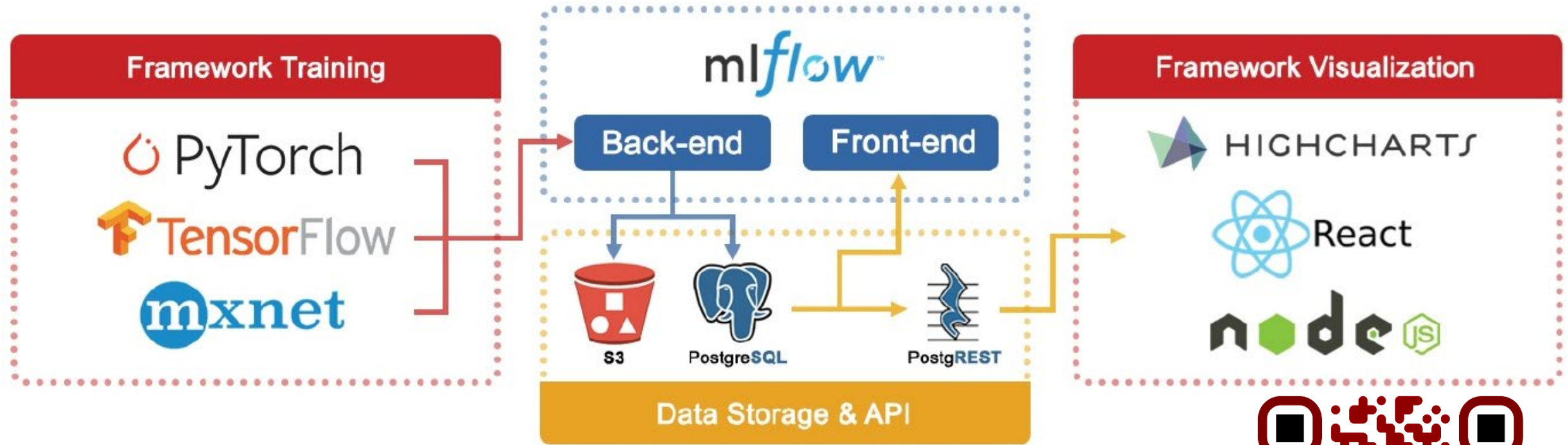
30-37% reduction in end-to-end trace completion time leads to ~15% reduction in total energy need

workload collocation for model training

- not all training needs all the resources of a single GPU
- collocation on GPUs benefits when the aggregate compute & memory needs of the collocated training runs fit in the GPU
- a collocation-aware resource manager help reduce time & energy required for a training workload with various models

need to build collocation-aware resource managers for deep learning targeting both small & large scales!

how to monitor hardware? – radT



- easy, extensible, and scalable tracking of hardware metrics (GPU utilization, storage access, carbon footprint ...)
- frontend for data exploration

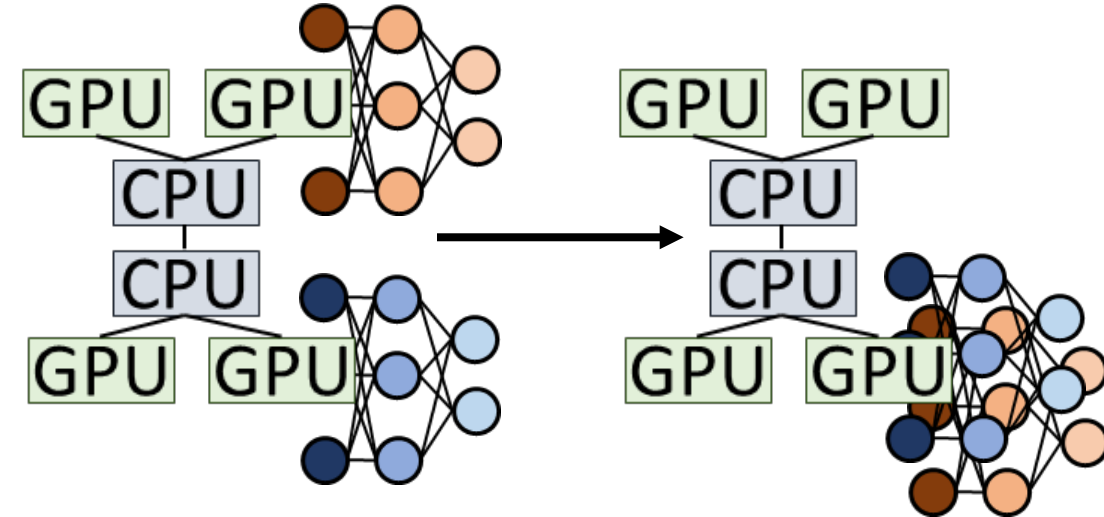
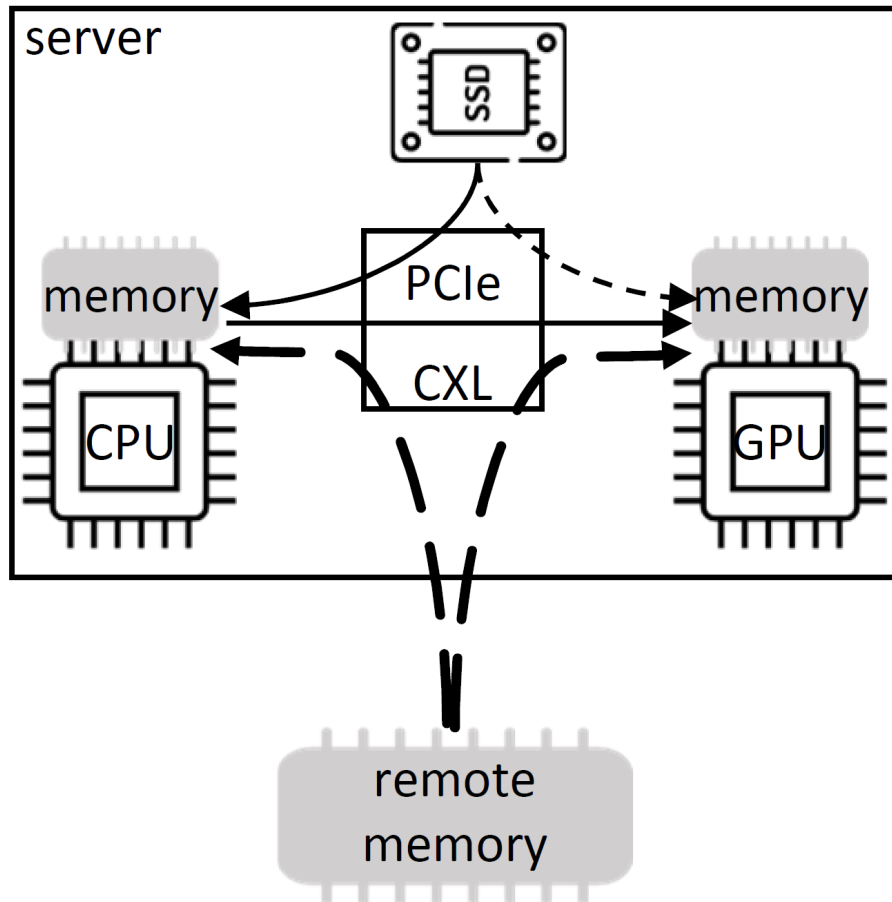


used by our group & data scientists @ITU for systematic benchmarking of deep learning training

RAD - resource-aware data systems**postdocs**Ties
RobroekEhsan
Yousefzadeh-
Asl-Miandoab**phd students**Robert
BayerJens Birk
Andersen**collaborators**Pamela Delgado
HEIG-VDTilmann Rabl
HPIAna Klimovic
ETHJulian Priest
ITU

- can we utilize modern hardware well?
- can we do more with less?

impact of storage hierarchy on data-intensive systems



workload collocation on GPUs



data management & processing on tiny hardware @ the edge

workload collocation for model training

thank you!

- not all training needs all the resources of a single GPU
- collocation on GPUs benefits when the aggregate compute & memory needs of the collocated training runs fit in the GPU
- a collocation-aware resource manager help reduce time & energy required for a training workload with various models

**need to build collocation-aware resource managers
for deep learning targeting both small & large scales!**