

# An Analysis of Collocation on GPUs for Deep Learning Training

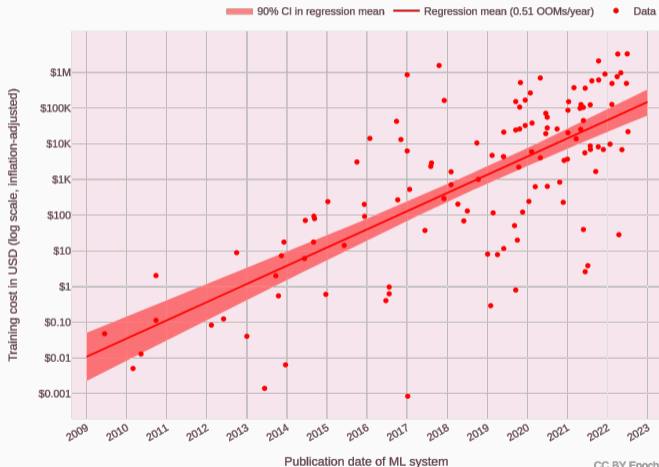
Ties Robroek, Ehsan Yousefzadeh-Asl-Miandoab, Pınar Tözün  
(IT University of Copenhagen)

---

# Need to utilize resources efficiently

- ▶ Training resource consumption is increasingly relevant
- ▶ GPUs are only utilized 52% on average for 100,000 jobs\*

Estimated training compute cost in USD: using price-performance trend



\* Jeon et al. "Analysis of Large-Scale Multi-Tenant GPU Clusters for DNN Training Workloads." USENIX ATC 2019.

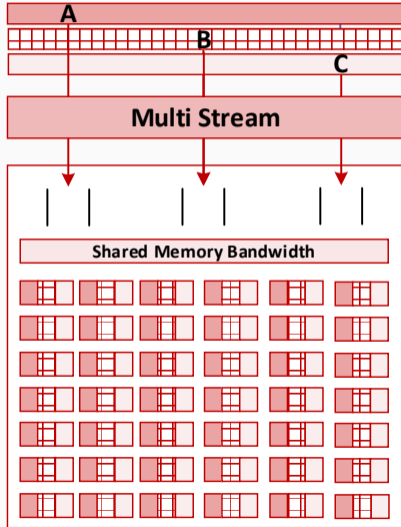
# Content

1. Collocation
2. Homogeneous Collocation
3. Mixed Collocation
4. Conclusion

# Collocation

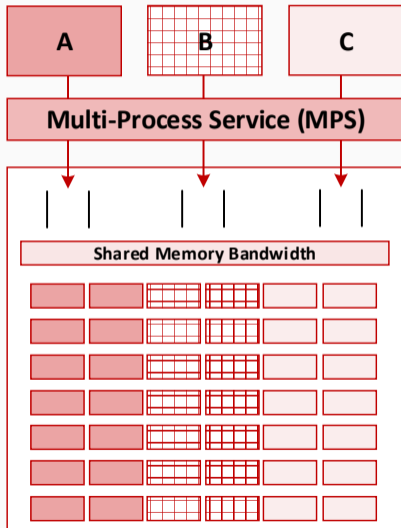
---

# Collocation on GPUs: Multi Stream (Naïve)



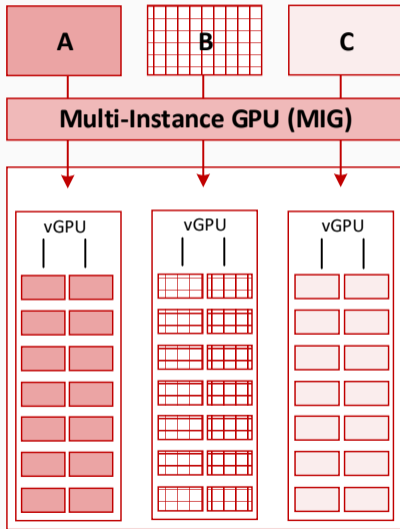
- ▶ CUDA-stream per process
- ▶ Convenient but limited efficiency

# Collocation on GPUs: MPS



- ▶ Convenient but single user
- ▶ High efficiency but no isolation

# Collocation on GPUs: MIG



- ▶ Rigid partitioning of GPU
- ▶ High isolation but overhead-prone

## Nvidia A100

compute	1	2	3	4	5	6	7	X	
memory	5gb	5gb	5gb	5gb	5gb	5gb	5gb	5gb	

For the A100 GPU:

- ▶ 7 compute instances of 14 streaming multiprocessors + overhead instance
- ▶ 8 memory instances of 5gb



# MIG 1g.5gb

vgpu	1	2	3	4	5	6	7	
compute	1	2	3	4	5	6	7	X
memory	5gb	5gb	5gb	5gb	5gb	5gb	5gb	5gb

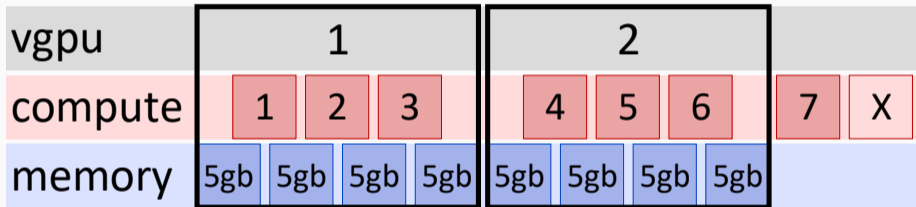
e.g. training many variations at the same time

# MIG Mixed Instances

vgpu	1				2		3	
compute	1	2	3	4	5	6	7	X
memory	5gb	5gb	5gb	5gb	5gb	5gb	5gb	5gb

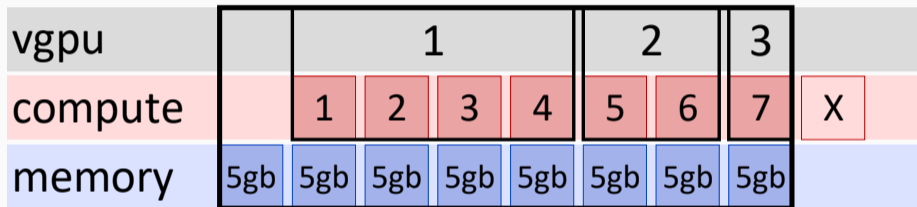
e.g. training models of different sizes

# MIG 3g.20gb



e.g. two larger models

# MIG Shared Memory



e.g. one of the models requires low compute but high memory

# Homogeneous Collocation

---

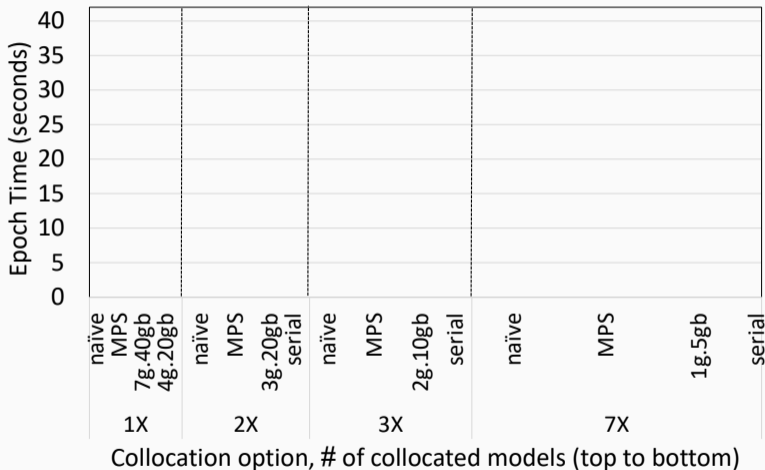
## Scales

- ▶ Small: ResNet26 on Cifar10
- ▶ Medium: EfficientNetv2 s on ImageNet64
- ▶ Large: CaiT on ImageNet2012

## Hardware

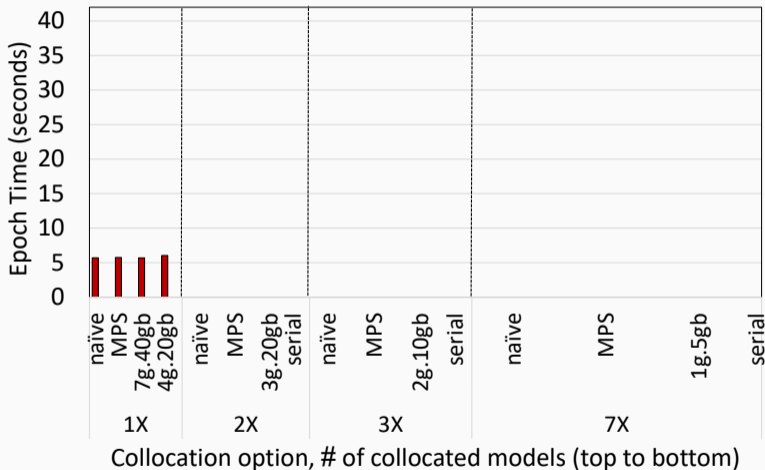
DGX A100 Station  
64-Core AMD Epyc  
NVIDIA A100 40GB GPU

# Small: ResNet26 + Cifar10



Small model training can greatly benefit from collocation

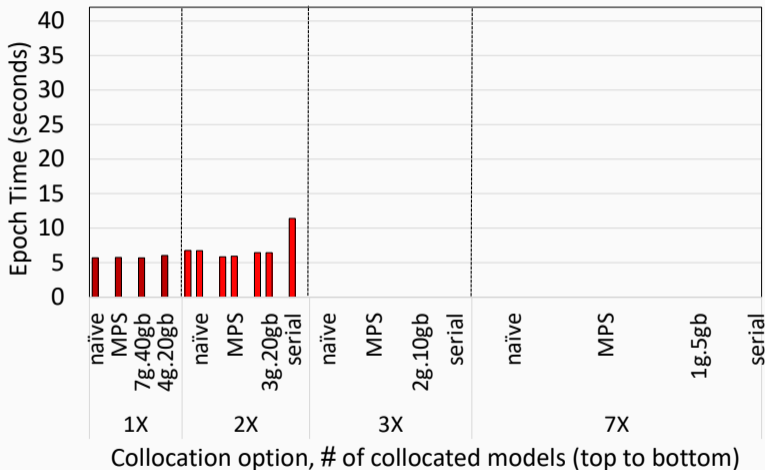
# Small: ResNet26 + Cifar10



Small model training can greatly benefit from collocation

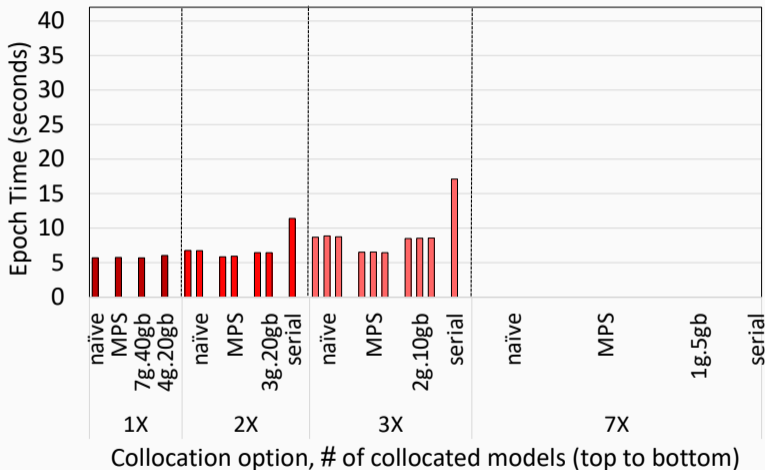


# Small: ResNet26 + Cifar10



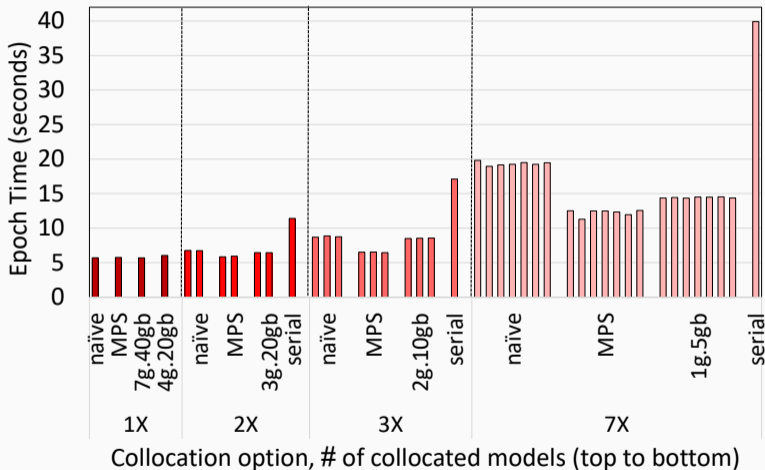
Small model training can greatly benefit from collocation

# Small: ResNet26 + Cifar10



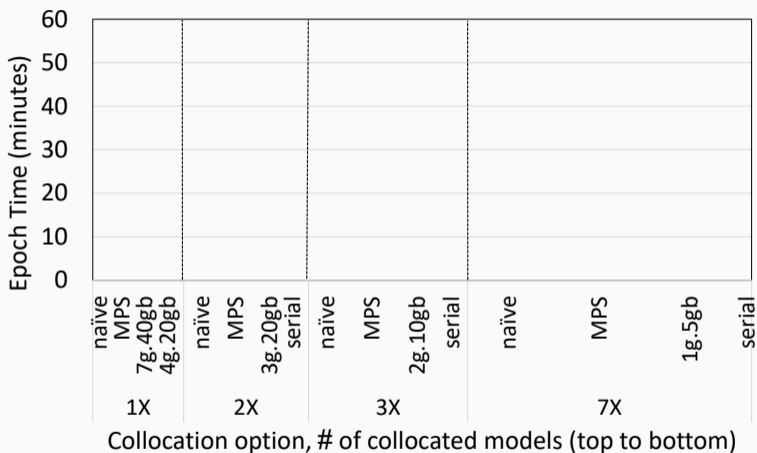
Small model training can greatly benefit from collocation

## Small: ResNet26 + Cifar10



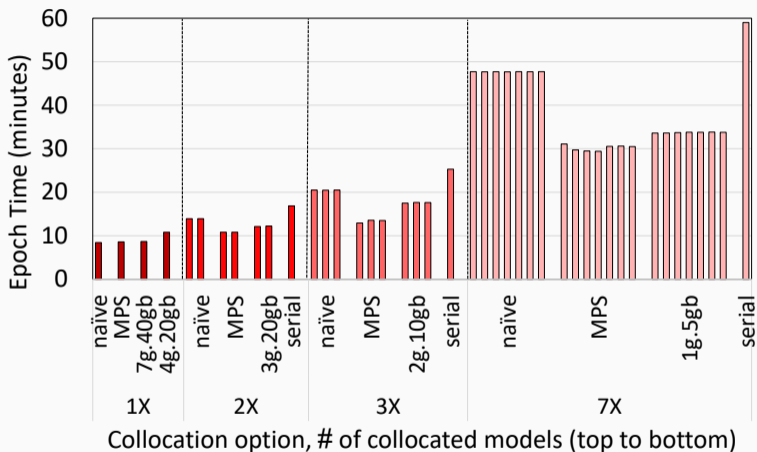
Small model training can greatly benefit from collocation

## Medium: EfficientNetv2 s + ImageNet64



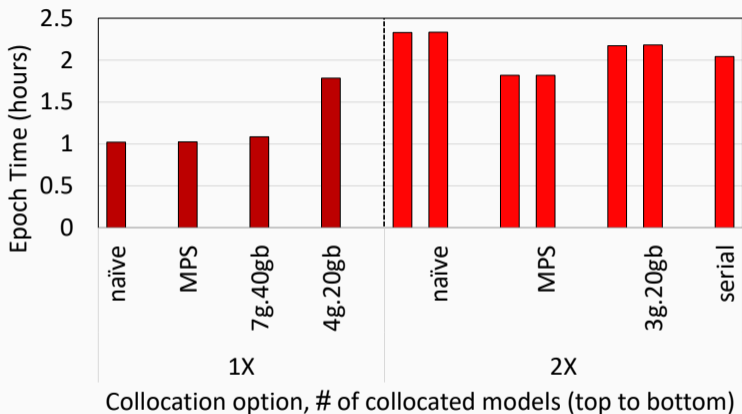
MPS and MIG can fully utilize the GPU

## Medium: EfficientNetv2 s + ImageNet64



MPS and MIG can fully utilize the GPU

# Large: CaiT + ImageNet

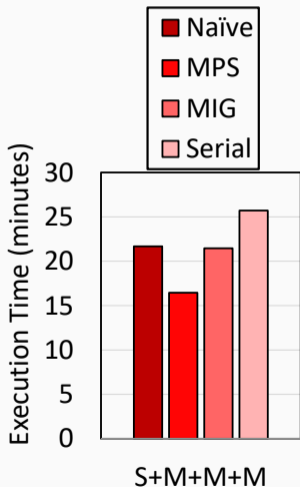


Collocation's benefit is diminished when the GPU is over-allocated

# Mixed Collocation

---

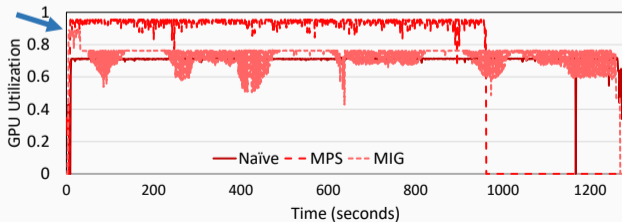
# Mixed vision workloads



Mixing ResNets:

- ▶ S: ResNet26 on Cifar10
- ▶ M: ResNet50 on ImageNet64

MIG cannot reallocate resources





# DLRM (Recommender) + ResNet152

Collocation	GPU Util.	Memory (GB)	Time (h)		
			ResNet	DLRM	Total
DLRM (no-colloc)					
ResNet152 (no-colloc)					
Naïve					
MPS					
MIG (2x 3g compute, shared memory)					

Models that stress different parts of the GPU are ideal candidates for collocation

# DLRM (Recommender) + ResNet152

Collocation	GPU Util.	Memory (GB)	Time (h)		
			ResNet	DLRM	Total
DLRM (no-colloc)	5%	29.14			
ResNet152 (no-colloc)	82%	8.47			
Naïve	81%	37.75			
MPS	81%	37.62			
MIG (2x 3g compute, shared memory)	39%	37.86			

Models that stress different parts of the GPU are ideal candidates for collocation

# DLRM (Recommender) + ResNet152

Collocation	GPU Util.	Memory (GB)	Time (h)		
			ResNet	DLRM	Total
DLRM (no-colloc)	5%	29.14	-	5.36	<b>6.41</b>
ResNet152 (no-colloc)	82%	8.47	1.05	-	
Naïve	81%	37.75			
MPS	81%	37.62			
MIG (2x 3g compute, shared memory)	39%	37.86			

Models that stress different parts of the GPU are ideal candidates for collocation

# DLRM (Recommender) + ResNet152

Collocation	GPU Util.	Memory (GB)	Time (h)		
			ResNet	DLRM	Total
DLRM (no-colloc)	5%	29.14	-	5.36	<b>6.41</b>
ResNet152 (no-colloc)	82%	8.47	1.05	-	
Naïve	81%	37.75	1.11	6.09	<b>6.09 (-5%)</b>
MPS	81%	37.62	1.10	5.57	<b>5.57 (-13%)</b>
MIG (2x 3g compute, shared memory)	39%	37.86	1.40	5.60	<b>5.60 (-13%)</b>

Models that stress different parts of the GPU are ideal candidates for collocation

# GPU Collocation

- ▶ Highly beneficial for small- and medium-sized workloads
- ▶ Mix models with varying resource requirements for high utilization
- ▶ MIG for strict separation, MPS for general use

Thank you!



[dasya.itu.dk](http://dasya.itu.dk)



[rad.itu.dk](http://rad.itu.dk)



[dff.dk](http://dff.dk)

IT UNIVERSITY OF COPENHAGEN

[itu.dk](http://itu.dk)

# GPU Collocation

- ▶ Highly beneficial for small- and medium-sized workloads
- ▶ Mix models with varying resource requirements for high utilization
- ▶ MIG for strict separation, MPS for general use

Thank you!



[dasya.itu.dk](http://dasya.itu.dk)



[rad.itu.dk](http://rad.itu.dk)



[dff.dk](http://dff.dk)

IT UNIVERSITY OF COPENHAGEN

[itu.dk](http://itu.dk)