

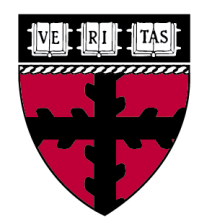


μ -TWO: 3x Faster Multi-model Training with Orchestration and Memory Optimization

Sanket Purandare
Abdul Wasay
Animesh Jain
Stratos Idreos



PyTorch

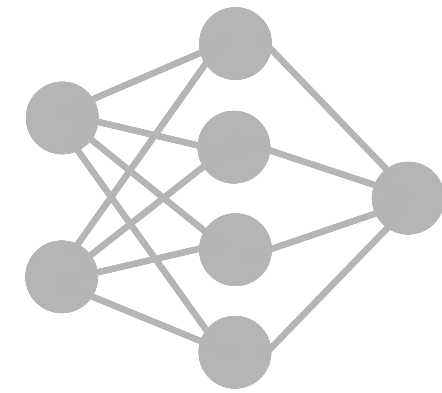


DAS lab
@ Harvard SEAS

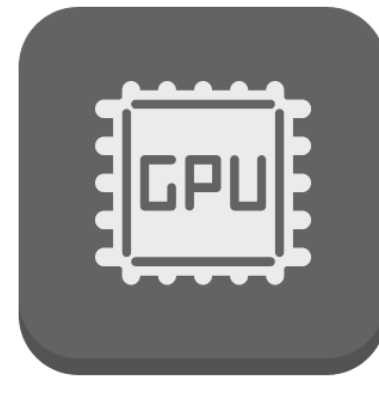
Deep Learning Training



High Computational Cost



200M



64 GPUs



79 hrs



USD 12000



1438 lbs



BERT

Natural Language Processing

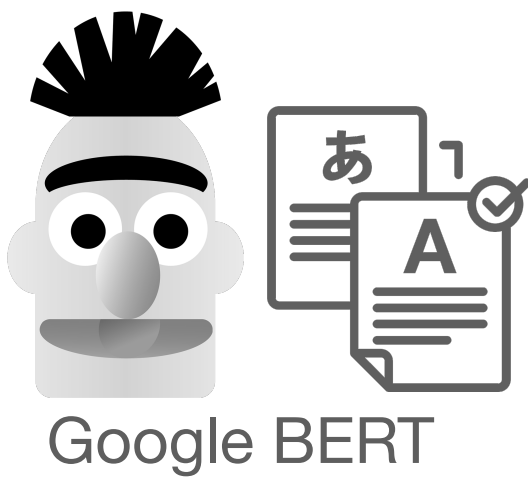
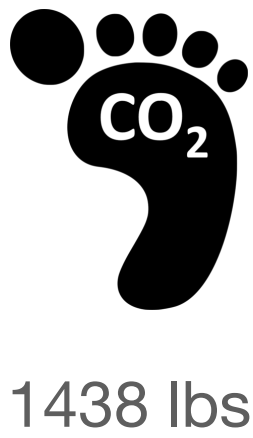
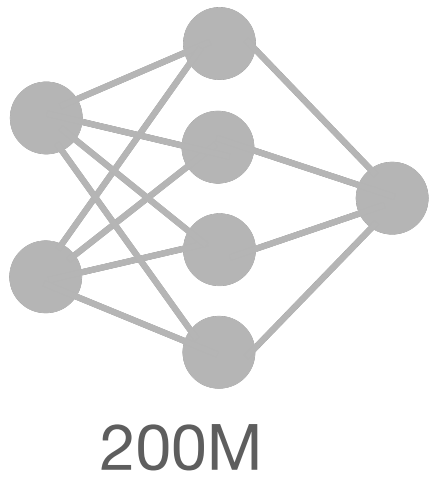
Strubell et al., Energy and policy considerations for deep learning in NLP, 2019



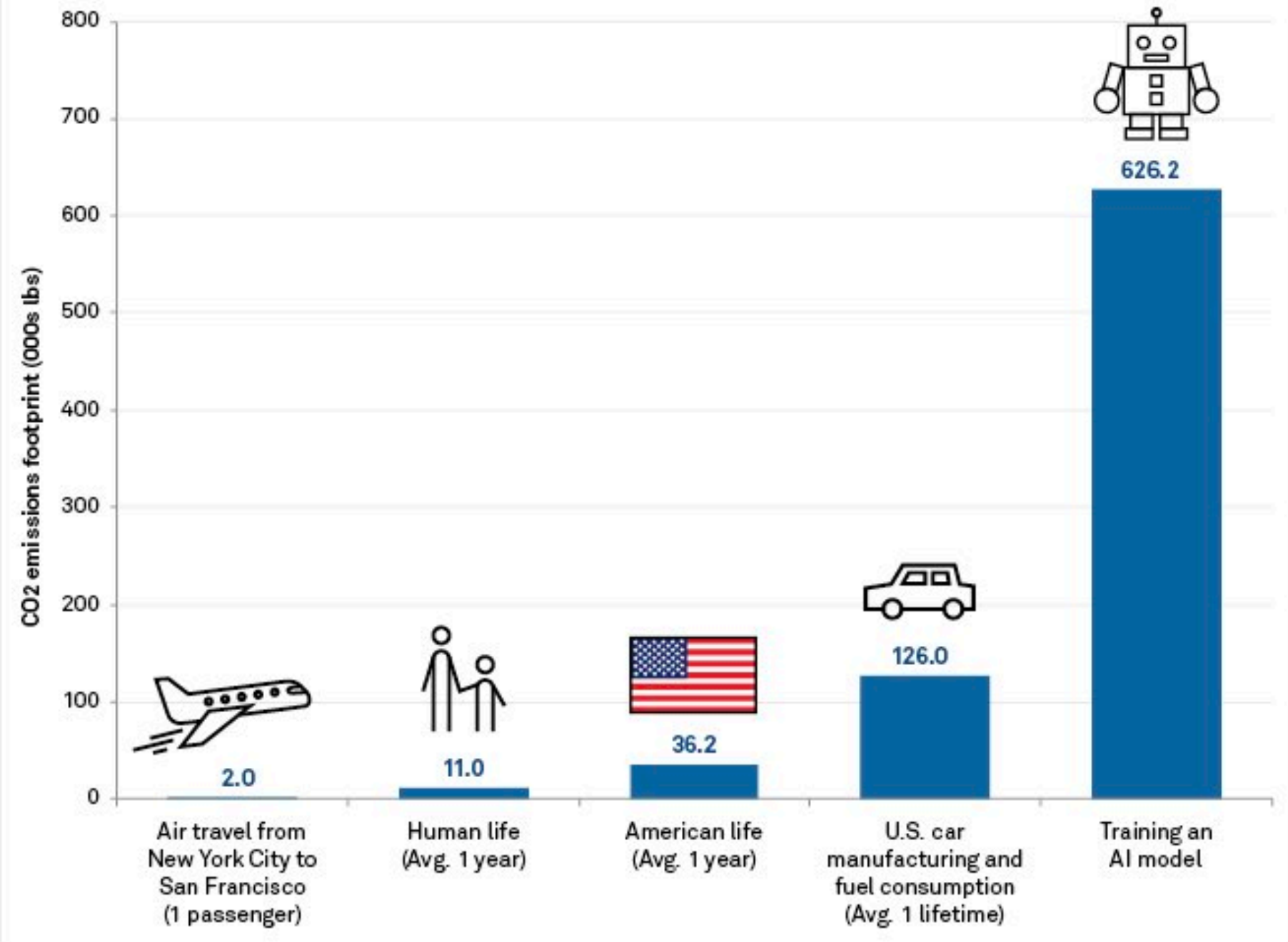
Deep Learning Training



High Computational Cost



CO2 emission benchmarks



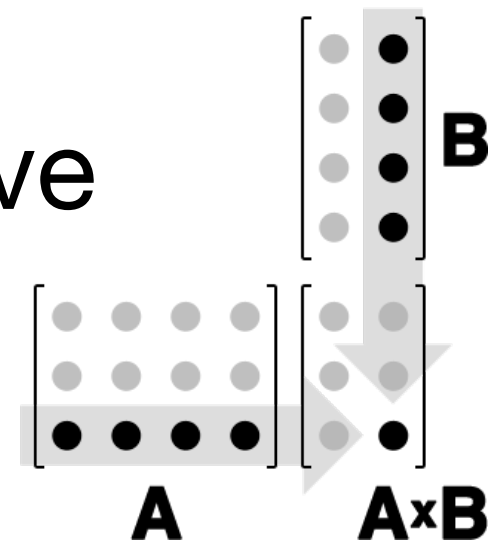
Data compiled Oct. 9, 2019.
An "American life" has a larger carbon footprint than a "Human life" because the U.S. is widely regarded as one of the top carbon dioxide emitters in the world.
Source: College of Information and Computer Sciences at University of Massachusetts Amherst

Strubell et al., Energy and policy considerations for deep learning in NLP, 2019



Sources?

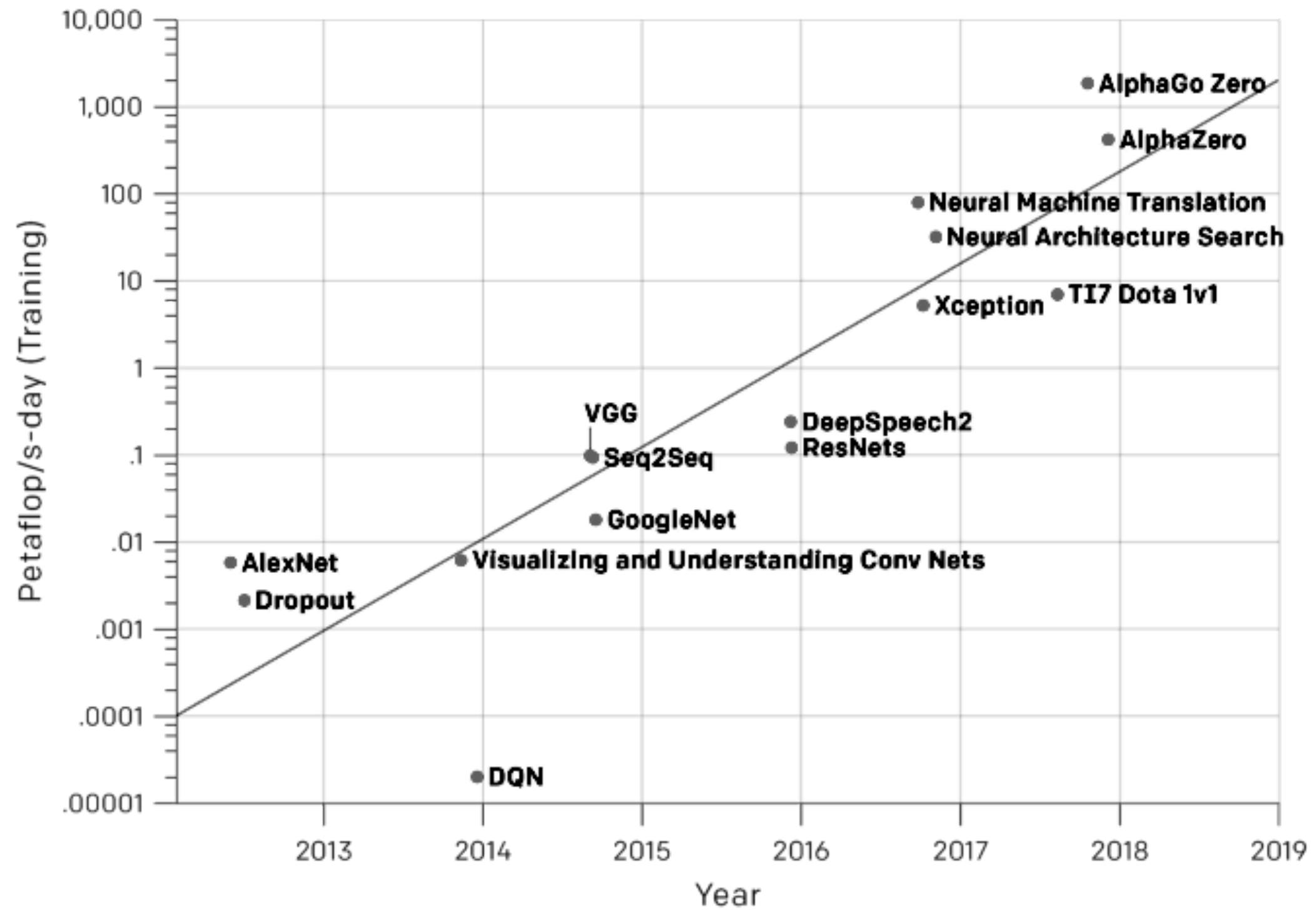
Compute Intensive



Deep Learning Training

High Computational Cost

300,000x Increase in compute



Open AI: <https://openai.com/blog/ai-and-compute/>

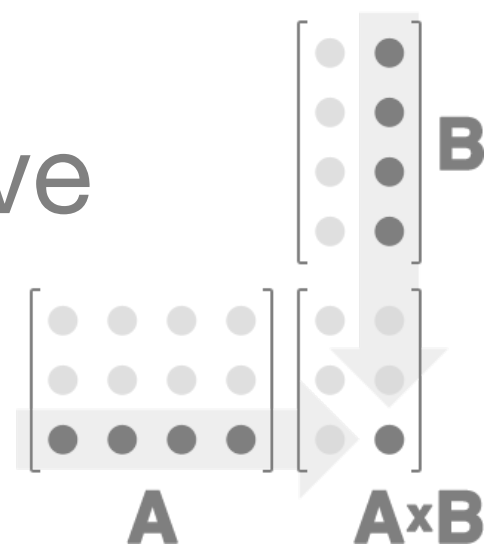


Deep Learning Training

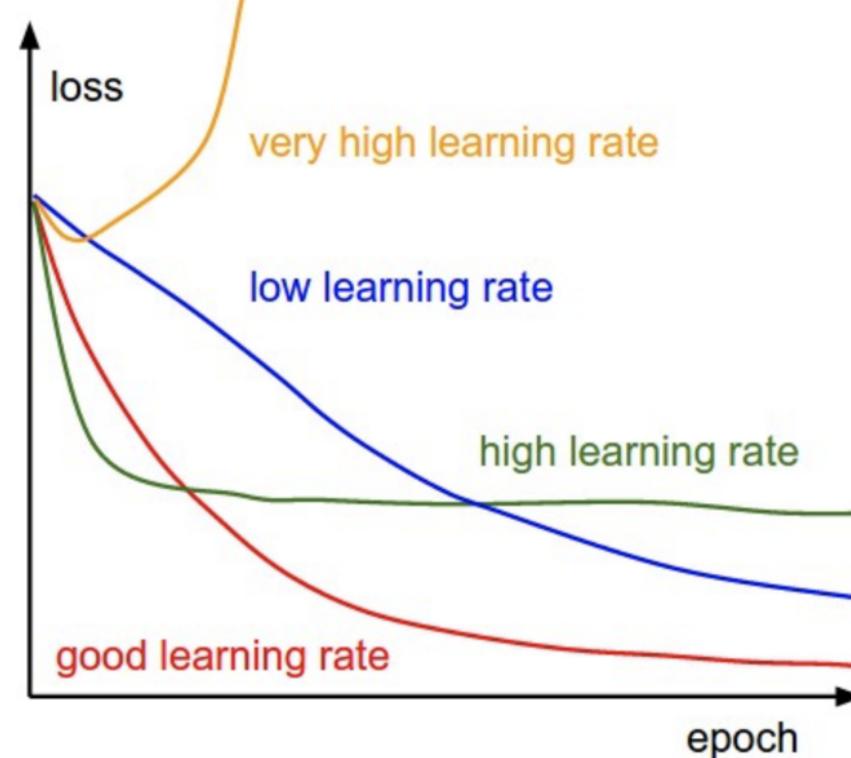
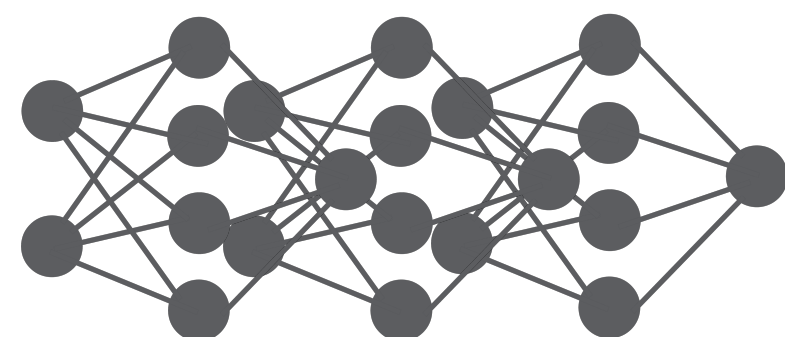
High Computational Cost

Sources?

Compute Intensive

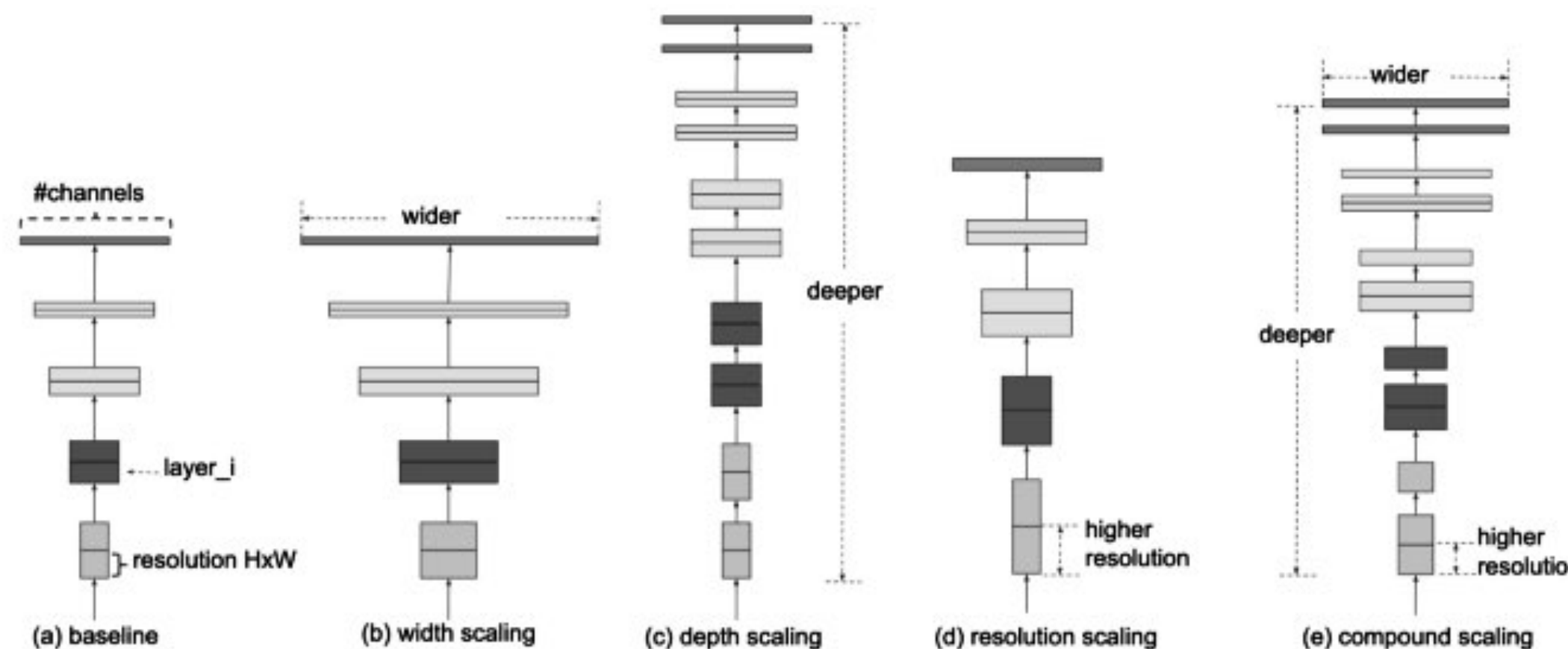


Training Several Networks



Hyperparameter Tuning

Bergstra, J., Algorithms for hyper-parameter optimization, NeurIPS 2011.



Neural Architecture Search

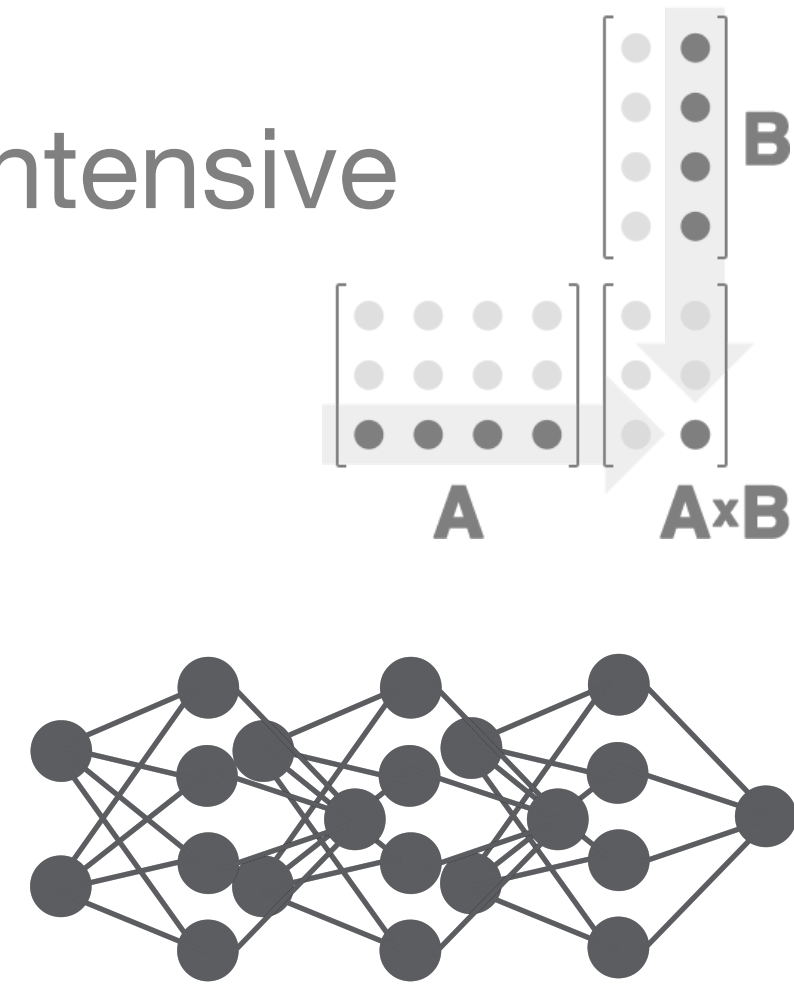
T. Elsken et al, Neural Architecture Search: A Survey, JMLR 2019



Sources?

Compute Intensive

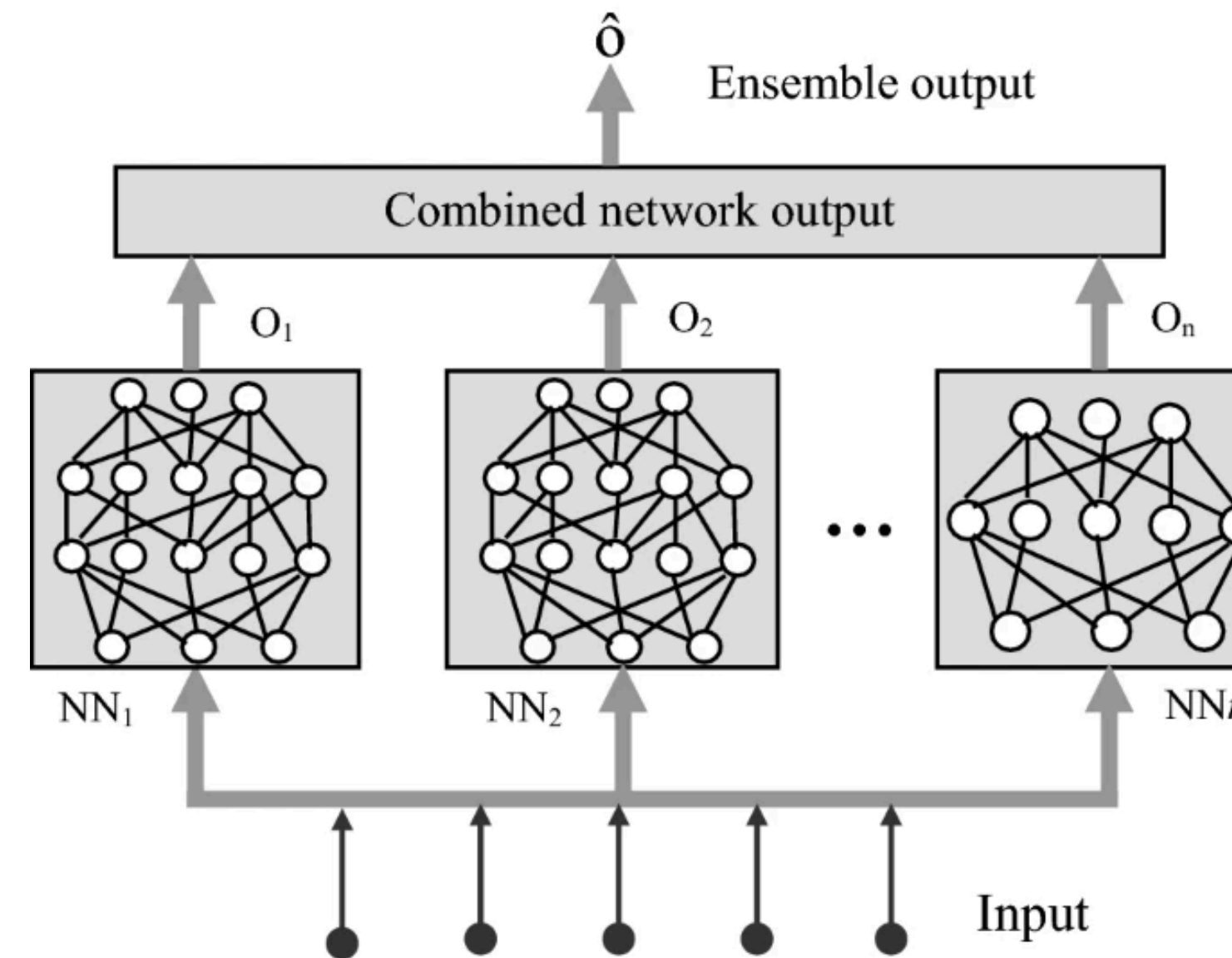
Training
Several
Networks



Deep Learning Training

High Computational Cost

Ensemble Training



Ganaie MA., Ensemble deep learning: A review, 2021



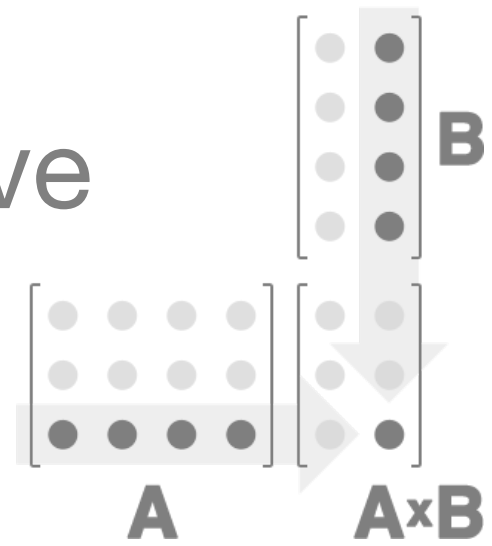
Deep Learning Training



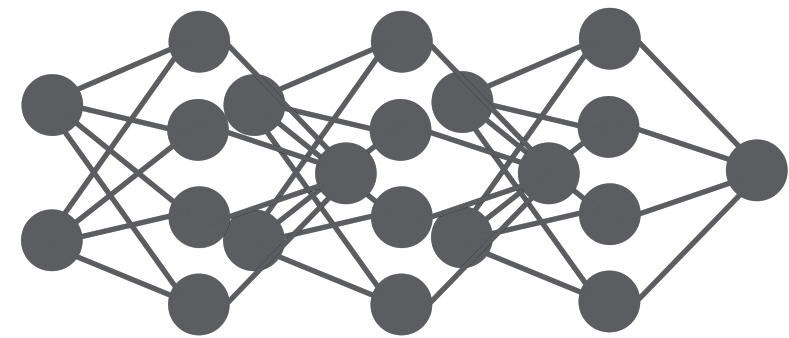
High Computational Cost

Sources?

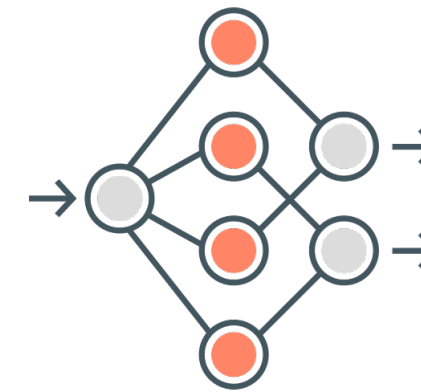
Compute Intensive



Training Several Networks



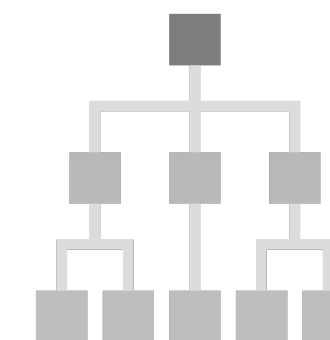
Interpretability Studies: Class Representations



Adversarial ML: Data poisoning



ML for Systems: Indexing



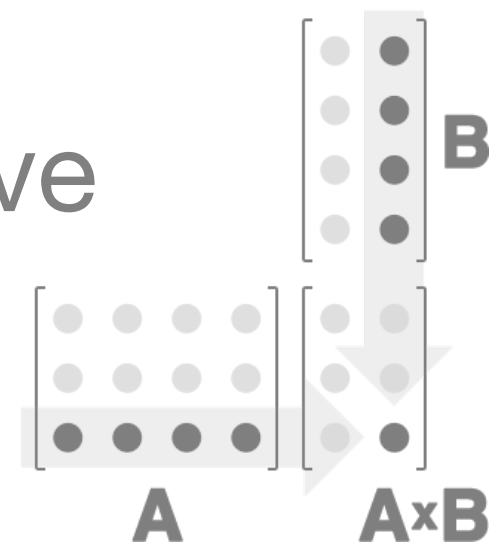
Deep Learning Training



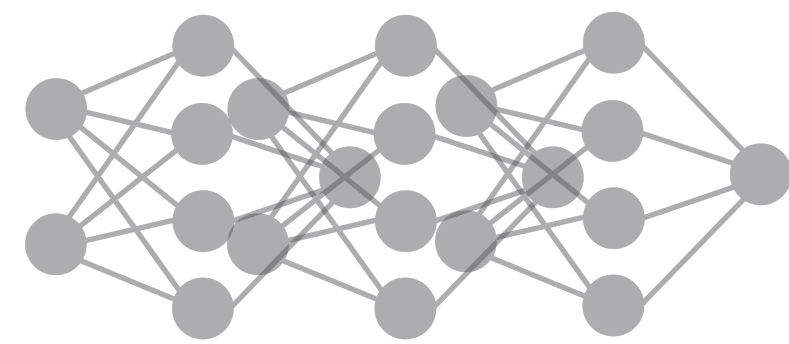
High Computational Cost

Sources?

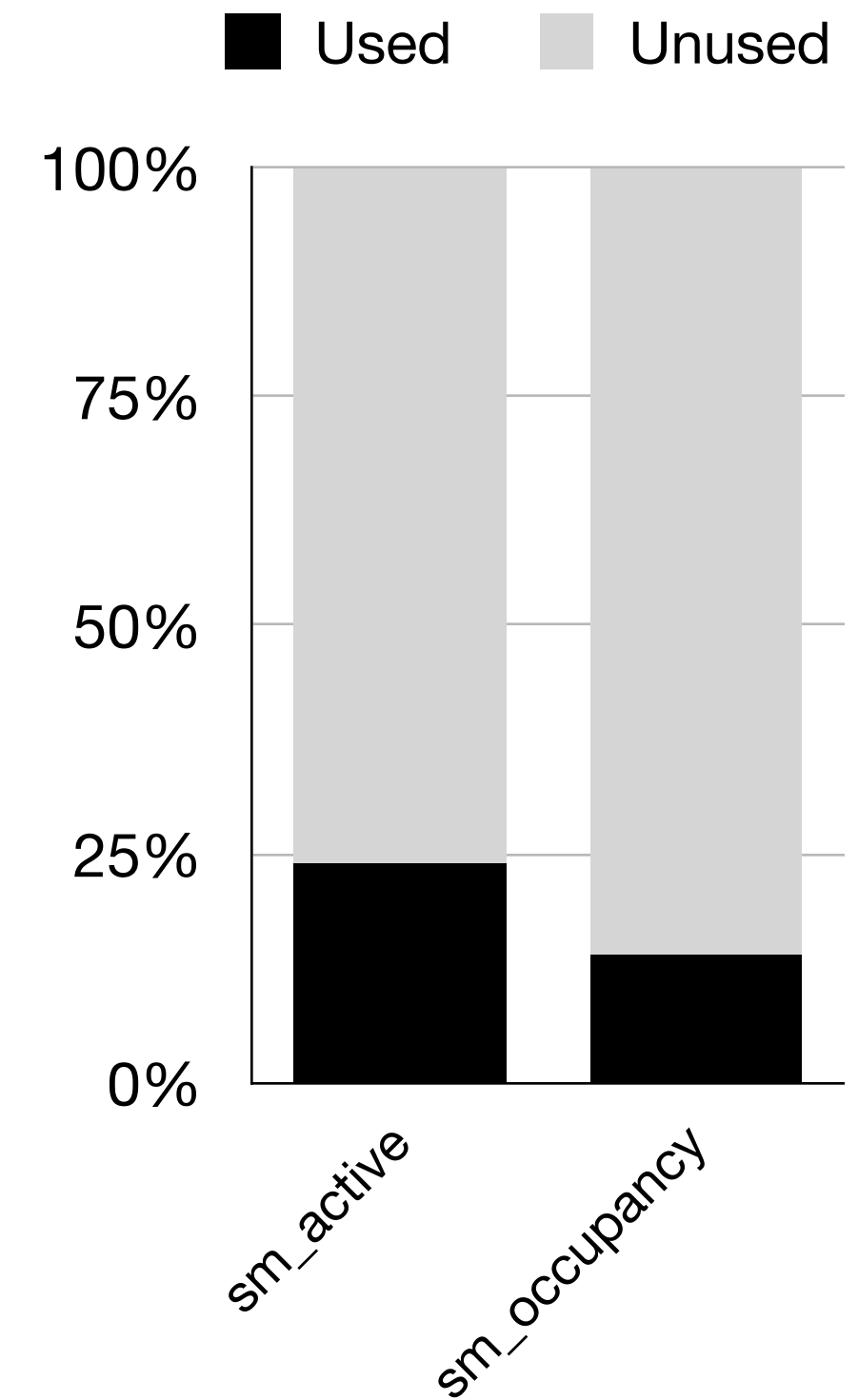
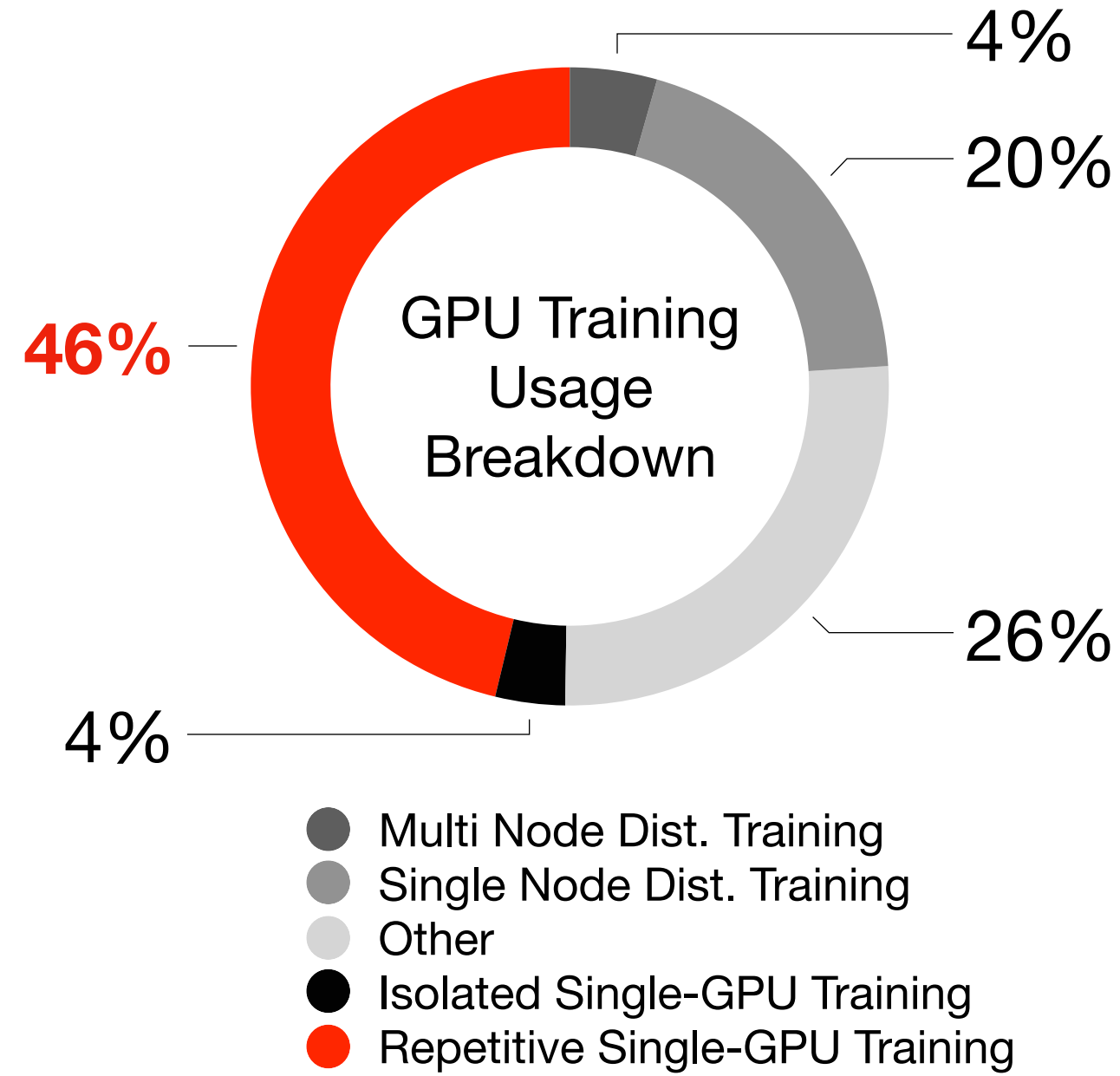
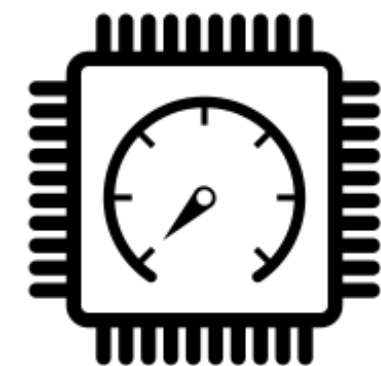
Compute Intensive



Training Several Networks



Sub-optimal Hardware Utilization



GPU Performance Counters



Deep Learning Training



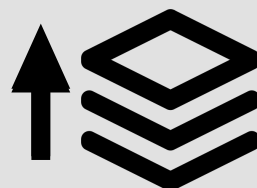
High Computational Cost

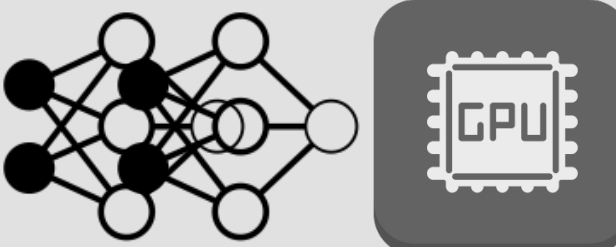


Sub-optimal
Hardware Utilization

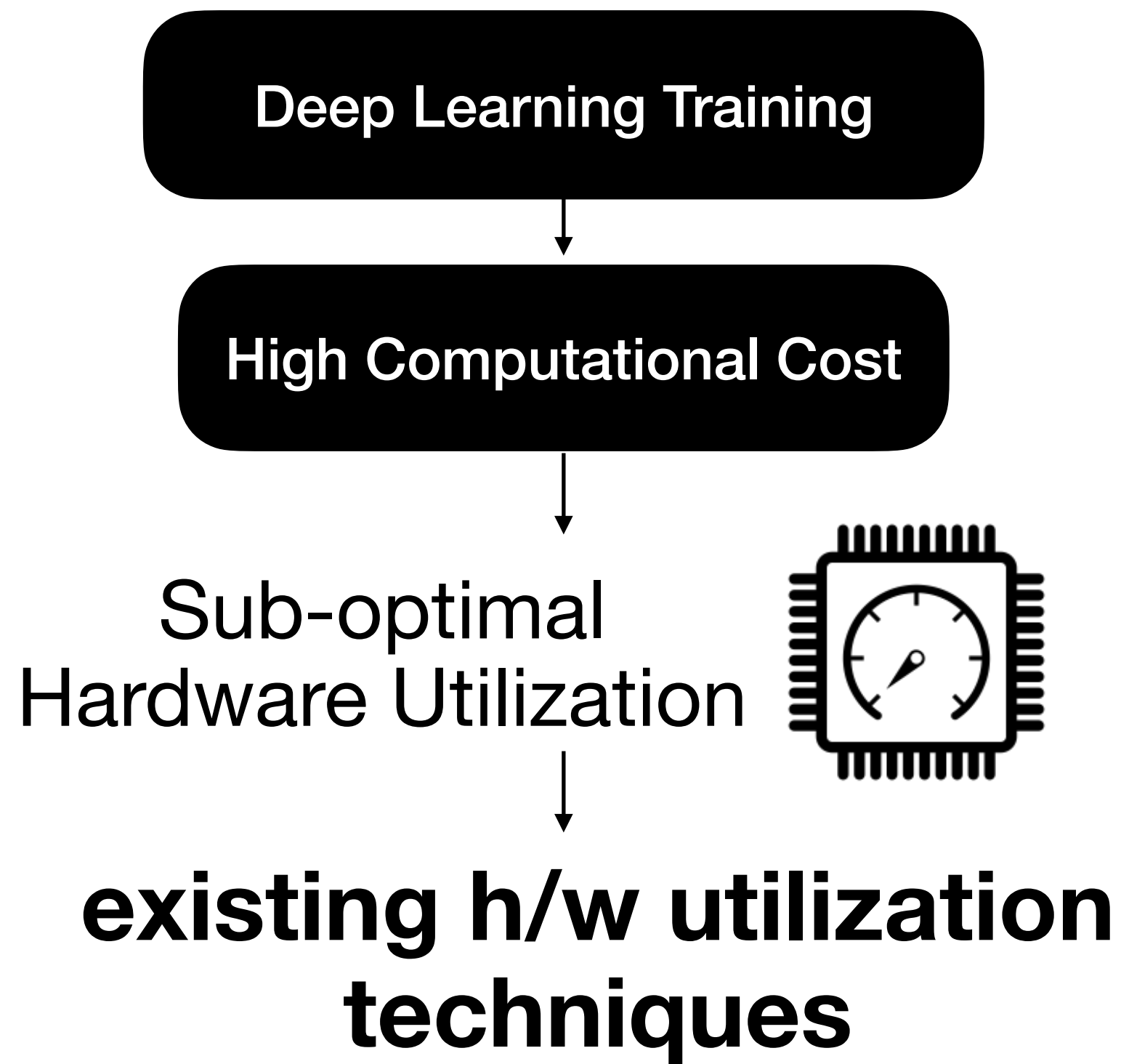


existing h/w utilization techniques

1  Increasing mini-batch size
(*data parallelism*)

2  Concurrent-training
(*resource sharing*)

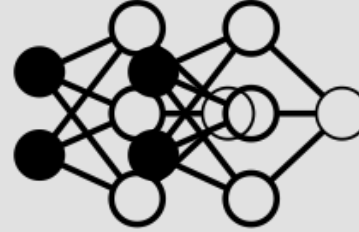





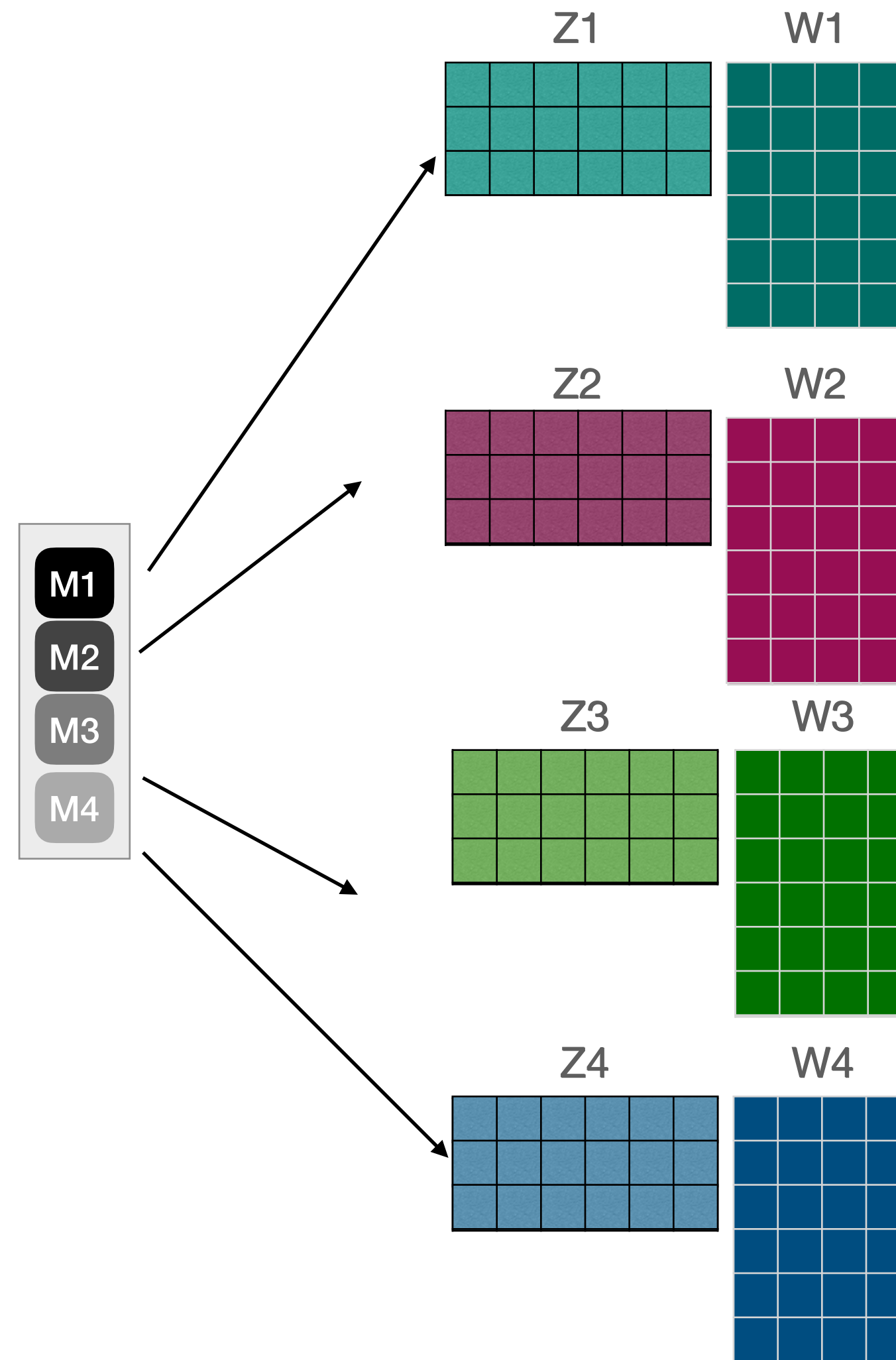
Horizontal Fusion

1  ↑ 

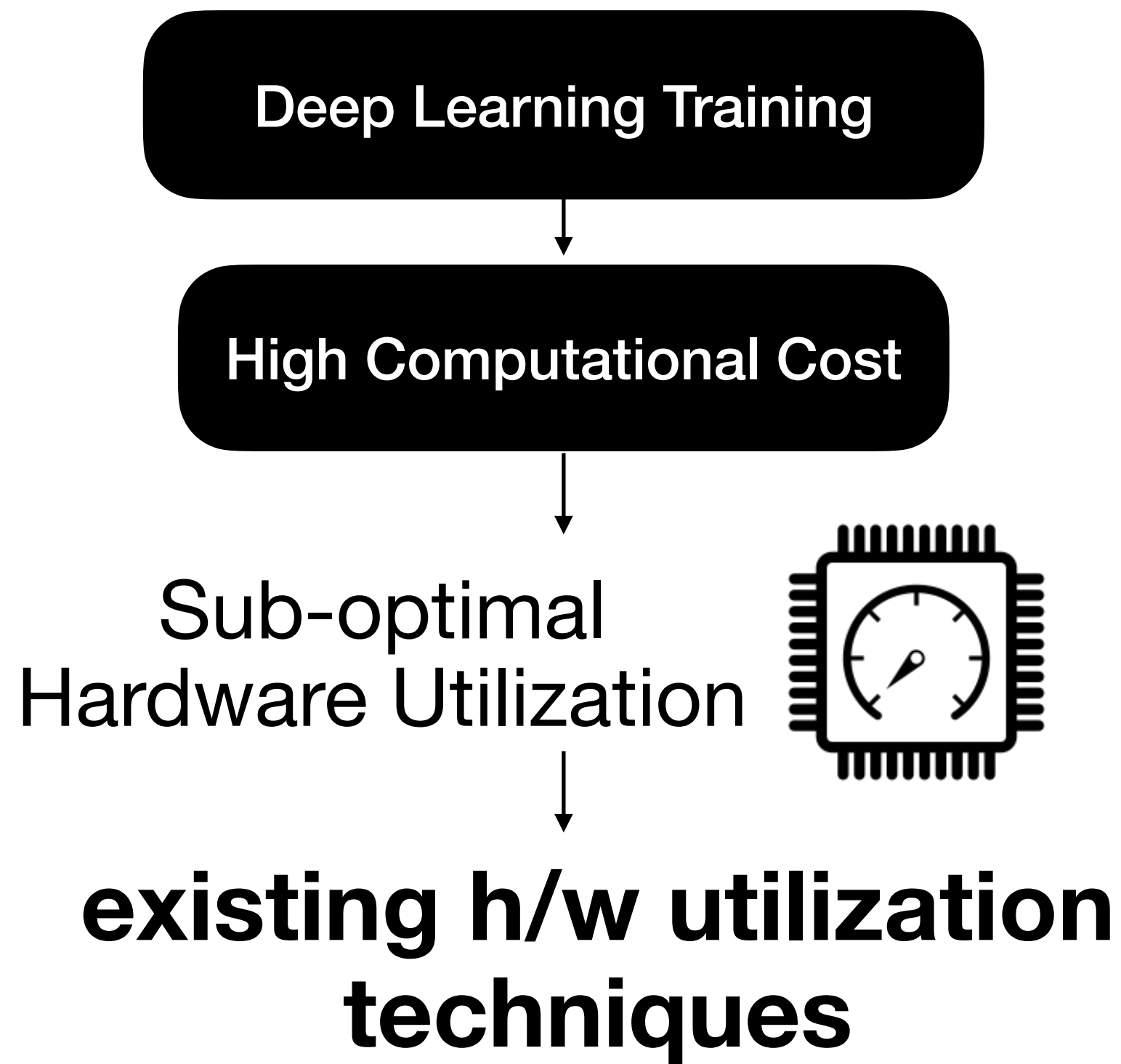
Increasing mini-batch size
(*data parallelism*)

2  

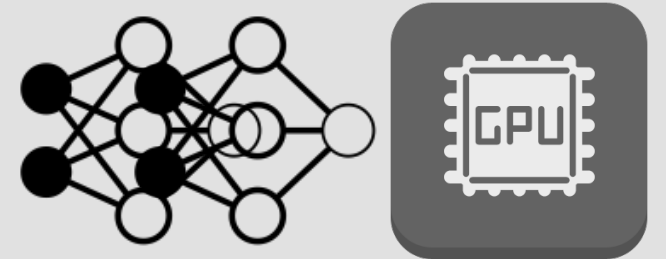
Concurrent-training
(*resource sharing*)



Matrix Multiply Operator

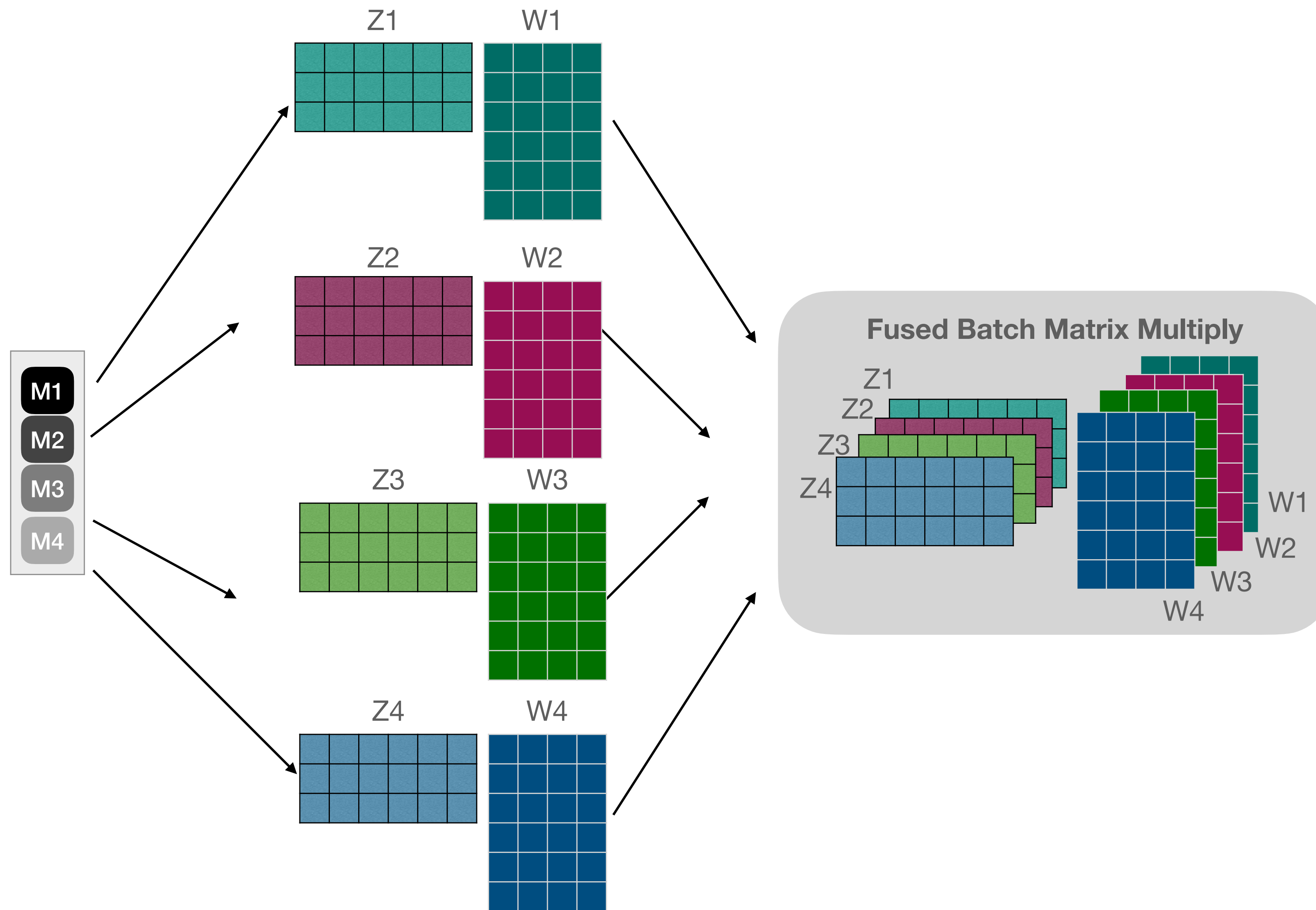


1  Increasing mini-batch size
(*data parallelism*)

2  Concurrent-training
(*resource sharing*)



Horizontal Fusion



Matrix Multiply Operator

Deep Learning Training



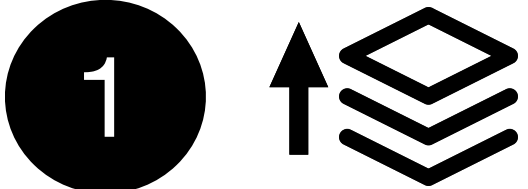
High Computational Cost



Sub-optimal Hardware Utilization



existing h/w utilization techniques



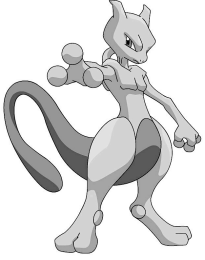
Increasing mini-batch size
(data parallelism)

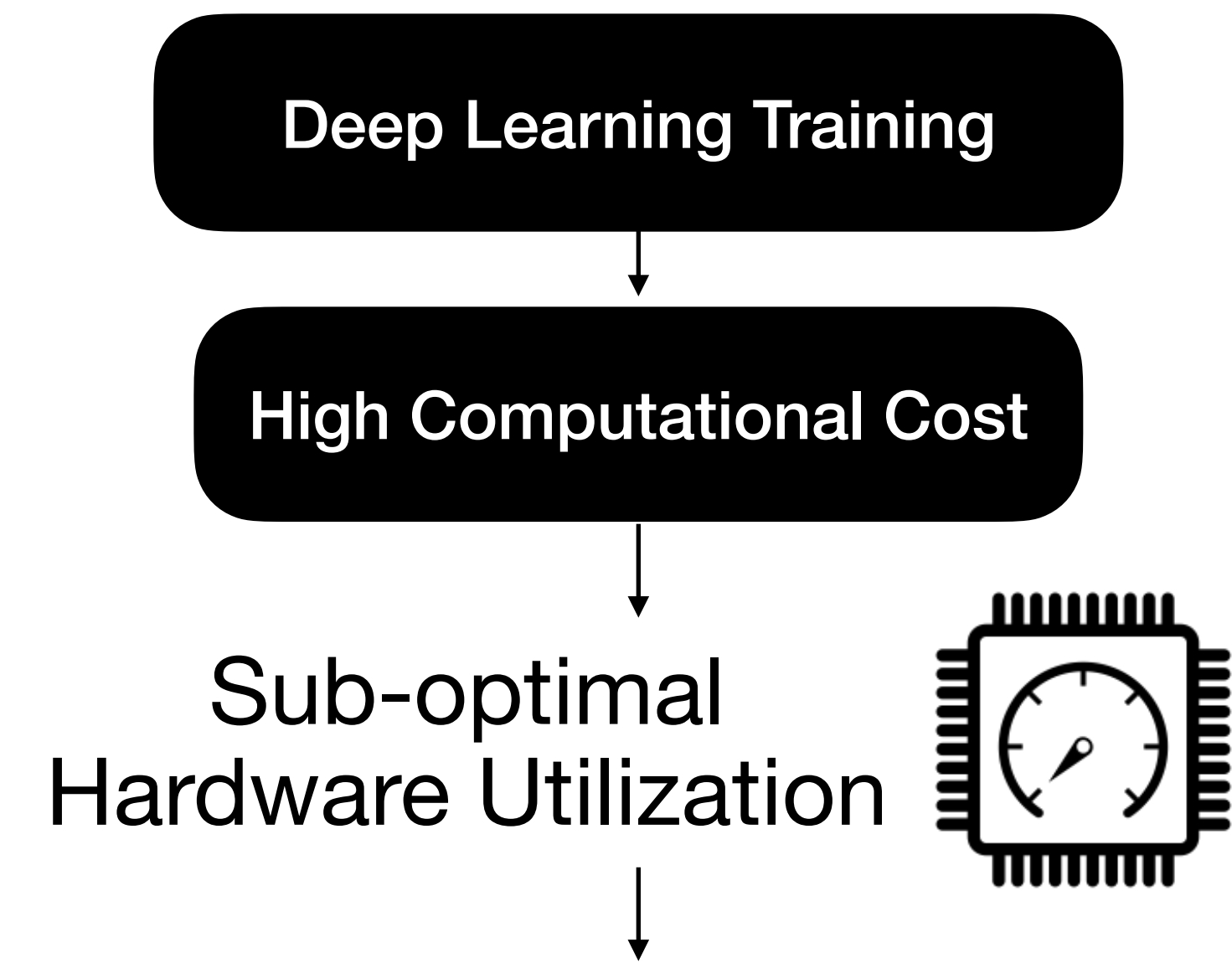


Concurrent-training
(resource sharing)

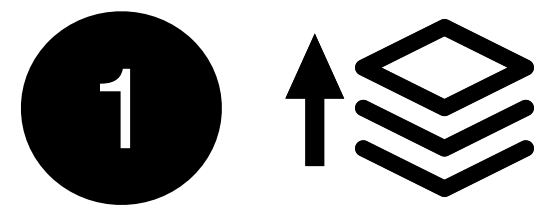


Added Memory Pressure

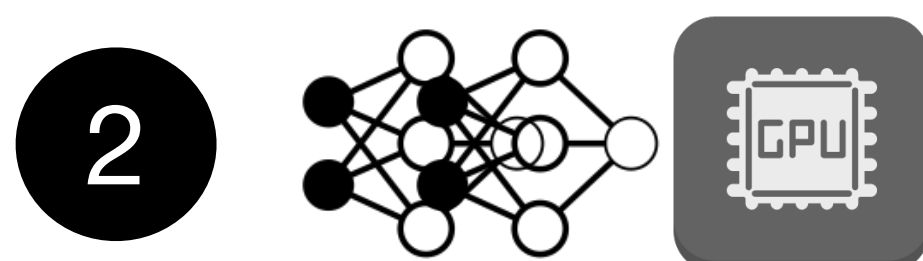




existing h/w utilization techniques



Increasing mini-batch size
(*data parallelism*)

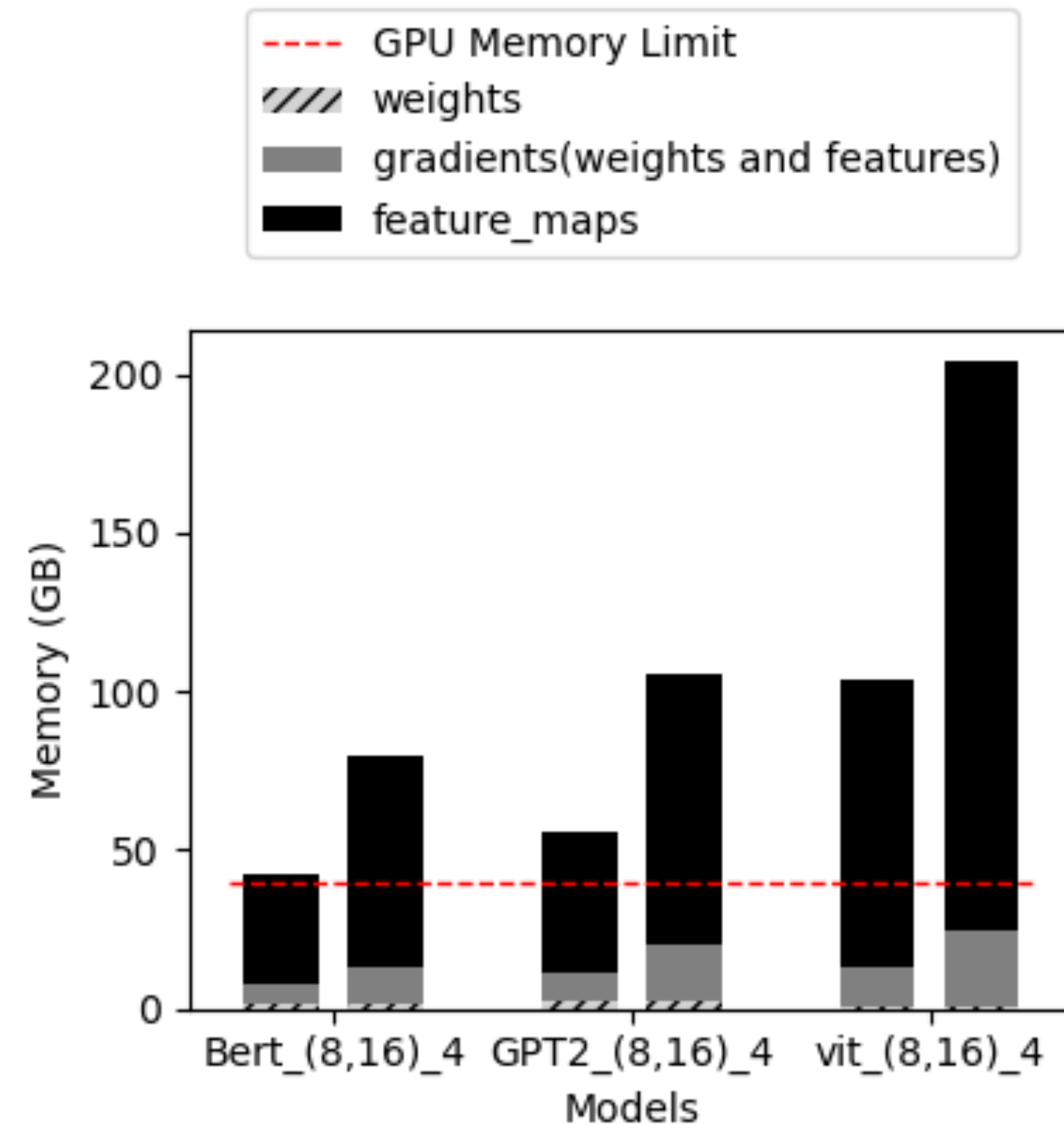


Concurrent-training
(*resource sharing*)




Memory Oversubscription -> Limits Scaling

- A Large model sizes
- B Greater number of models
- C Large training memory footprint
- D Limited GPU memory capacity

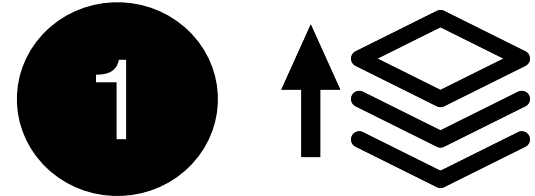


Deep Learning Training

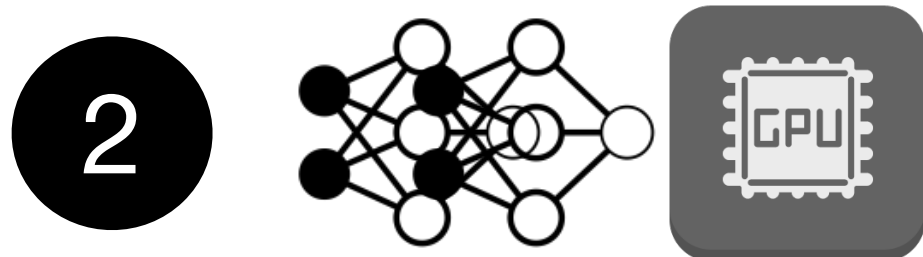
High Computational Cost

Sub-optimal Hardware Utilization 

existing h/w utilization techniques



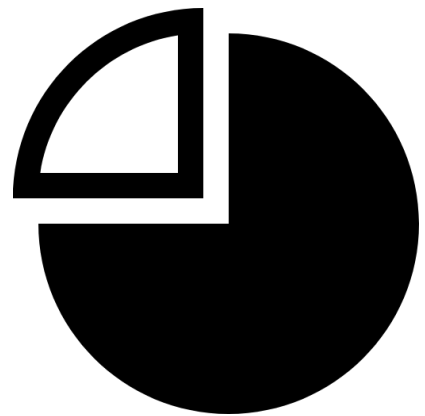
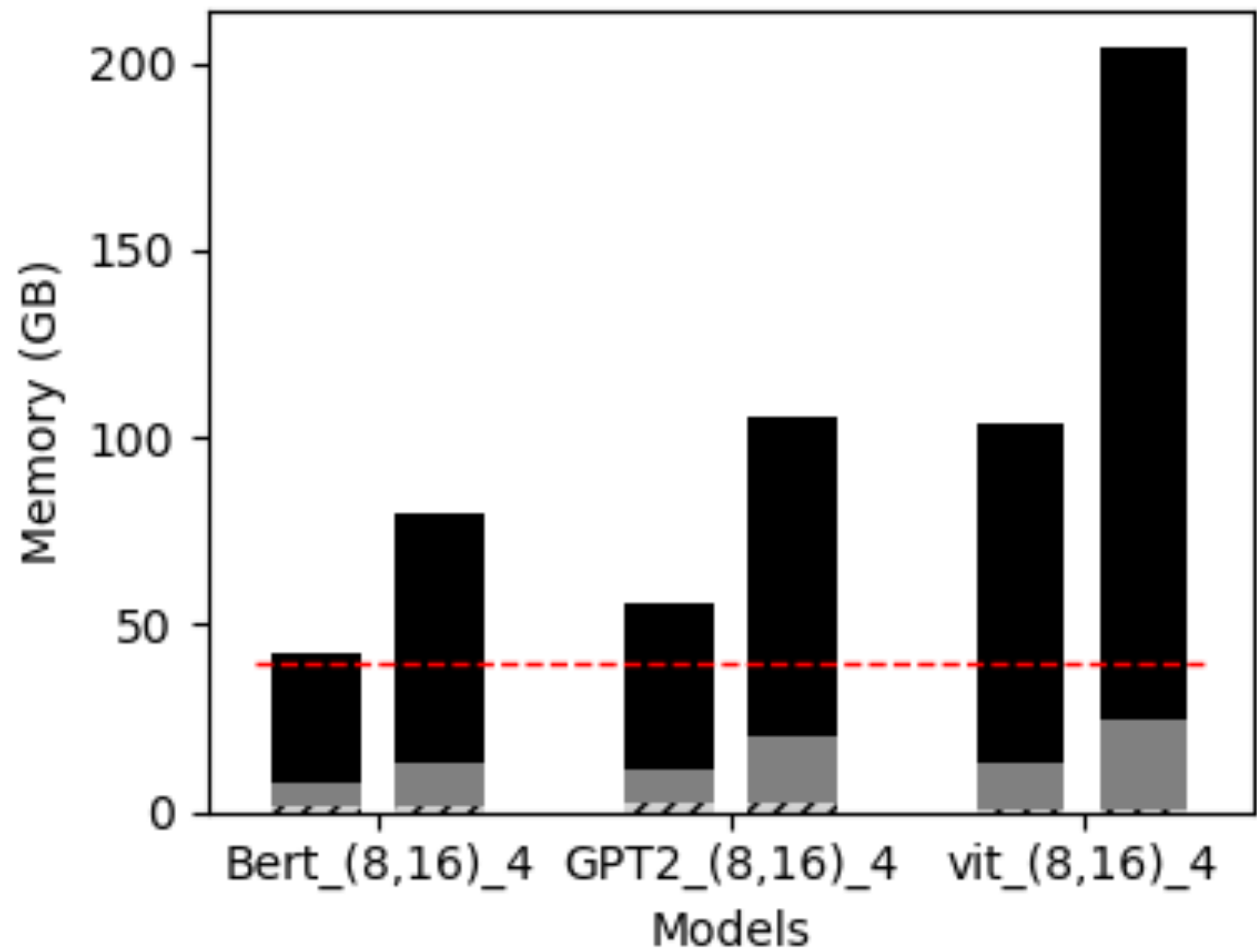
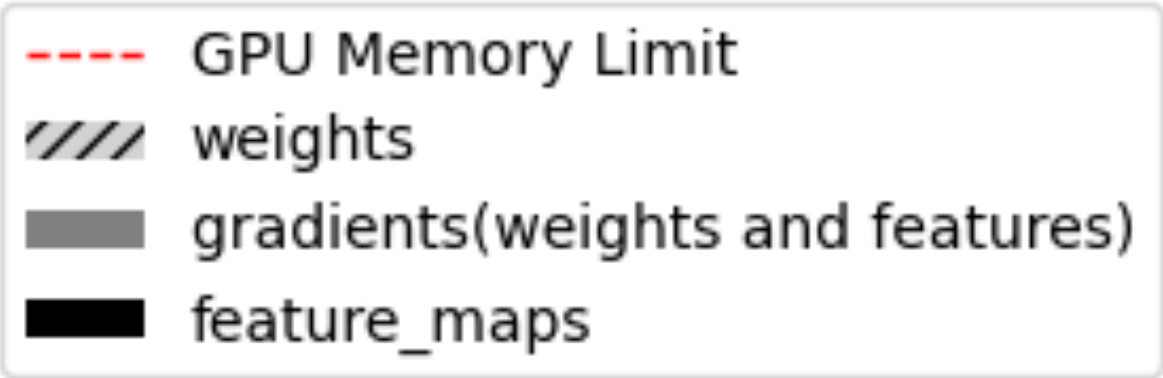
Increasing mini-batch size
(data parallelism)



Concurrent-training
(resource sharing)




Addressing Memory Oversubscription



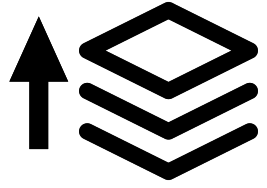
Feature Maps have the largest and most significant share in memory consumption.

Deep Learning Training

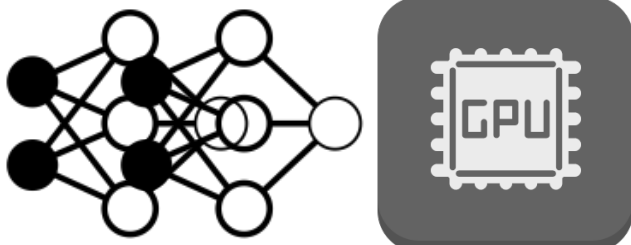
High Computational Cost

Sub-optimal Hardware Utilization 

existing h/w utilization techniques

1 

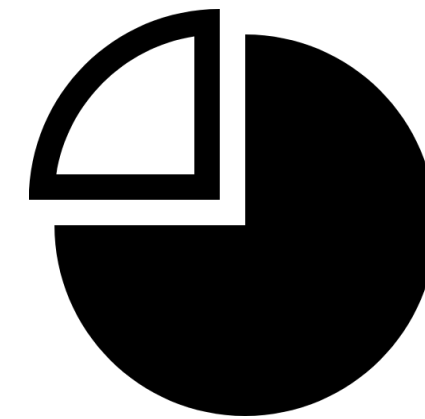
Increasing mini-batch size
(data parallelism)

2 

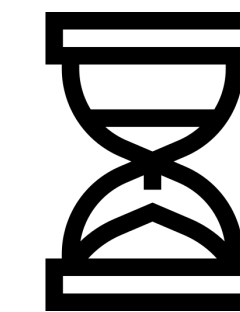
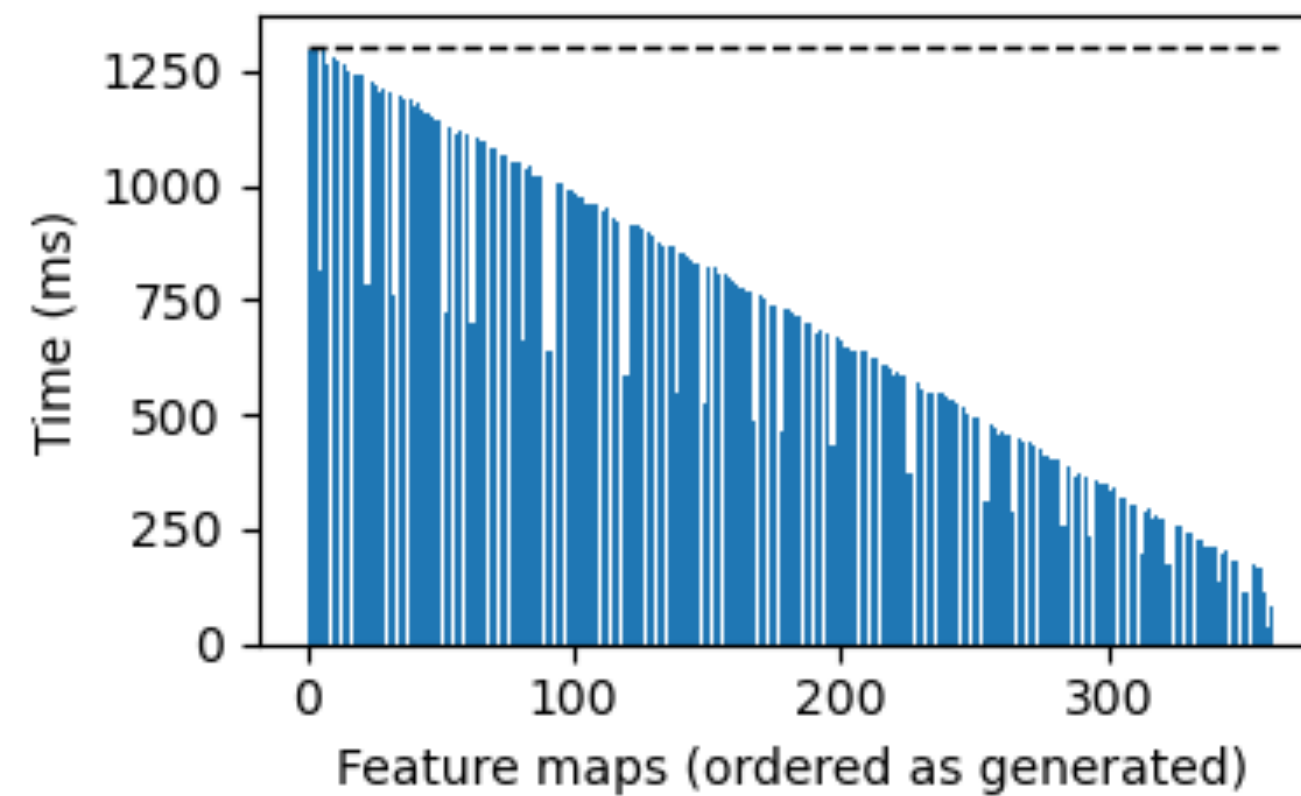
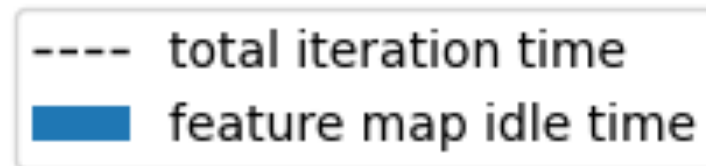
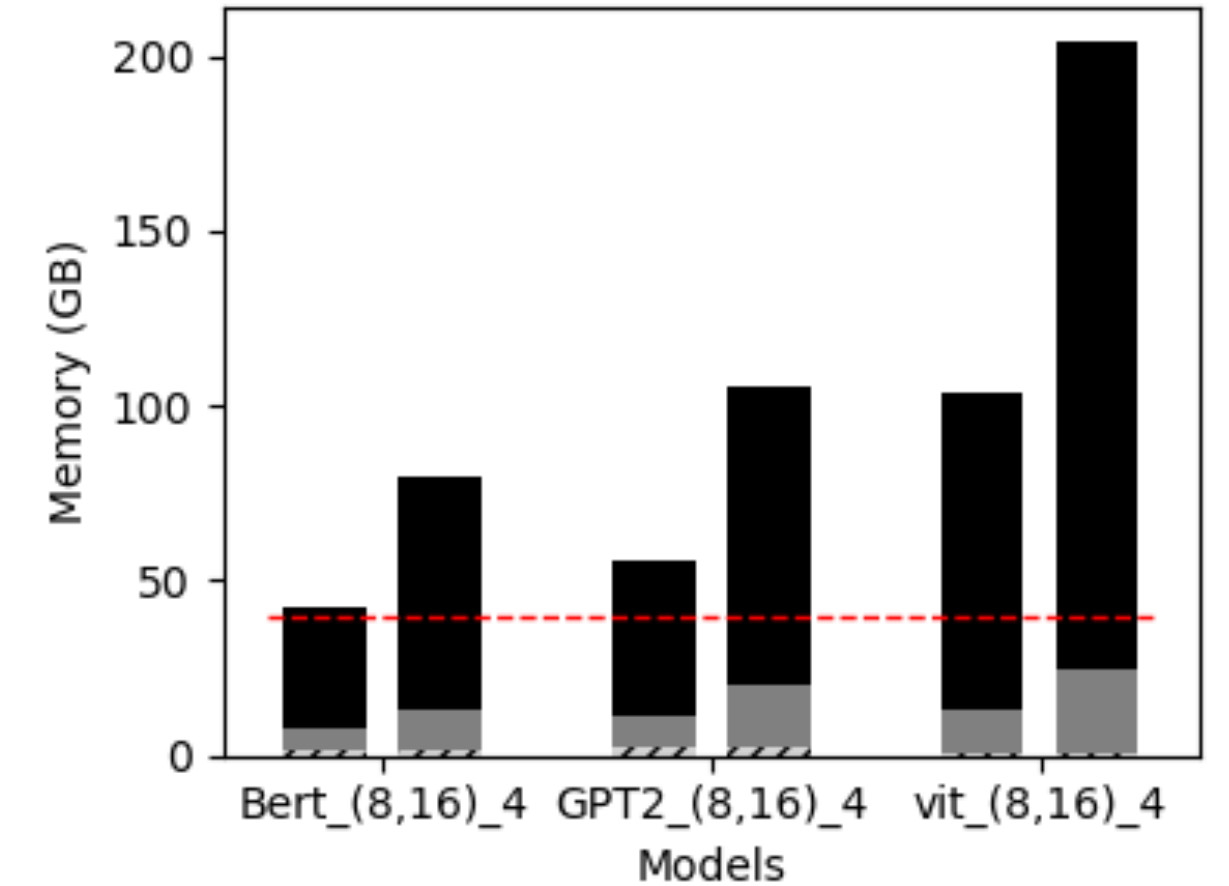
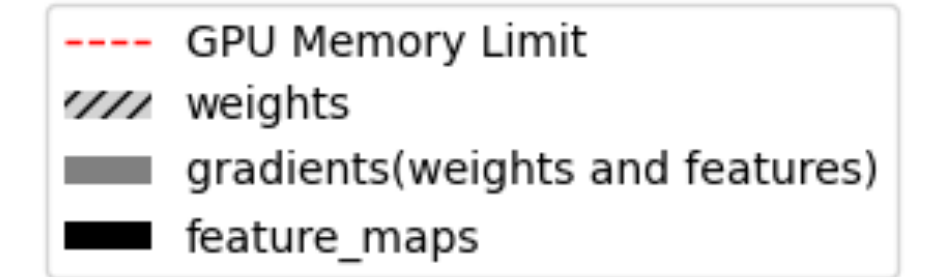
Concurrent-training
(resource sharing)



Addressing Memory Oversubscription




Feature Maps have the largest and most significant share in memory consumption.



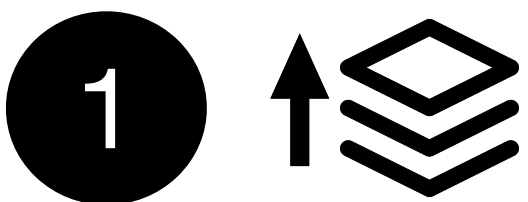
Feature Maps lie idle in GPU memory between forward and backward pass

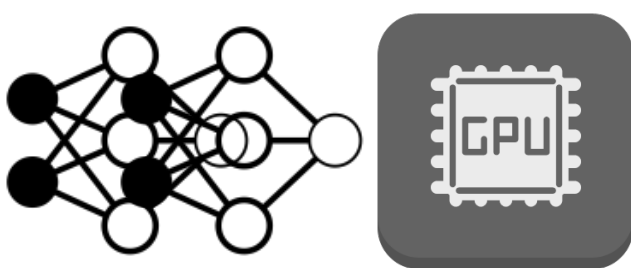
Deep Learning Training

High Computational Cost

Sub-optimal Hardware Utilization 

existing h/w utilization techniques

1 
Increasing mini-batch size
(*data parallelism*)

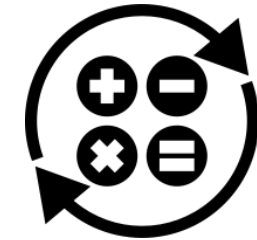
2 
Concurrent-training
(*resource sharing*)

Addressing Memory Oversubscription

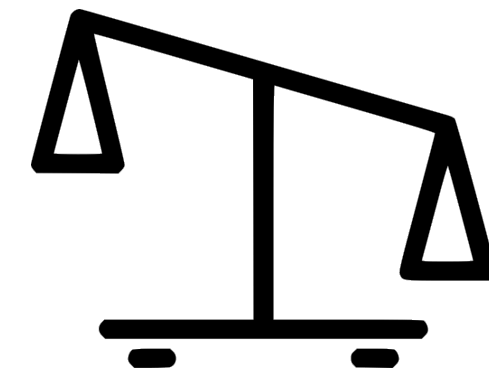
Recomputation



Discard feature maps after use in forward pass



Recompute when needed during backward pass



Tradeoff Compute for Memory



Superfluous Compute Overhead

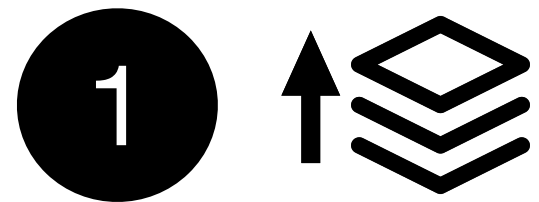


Deep Learning Training

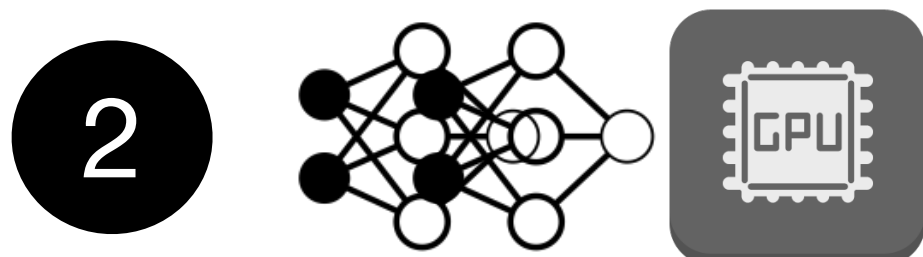
High Computational Cost

Sub-optimal Hardware Utilization 

existing h/w utilization techniques



Increasing mini-batch size
(data parallelism)

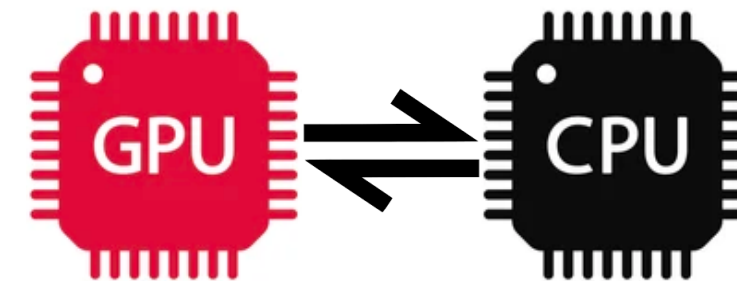


Concurrent-training
(resource sharing)



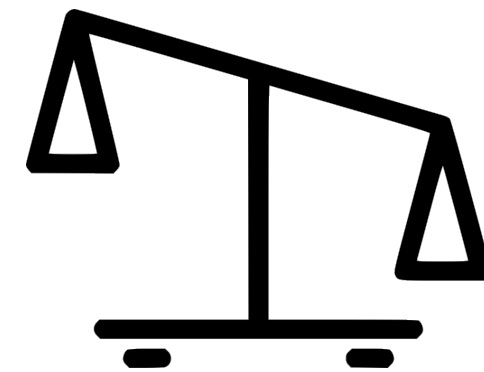
Addressing Memory Oversubscription

Swapping



Offload feature maps after use in forward pass to larger host memory

Fetch feature maps to GPU memory when needed during backward pass



Tradeoff stall time for Memory



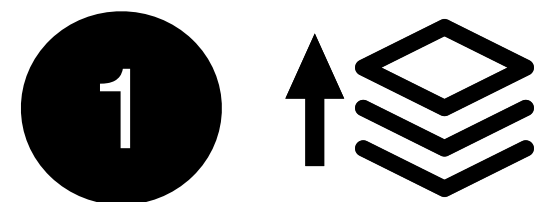
Stalling Overhead

Deep Learning Training

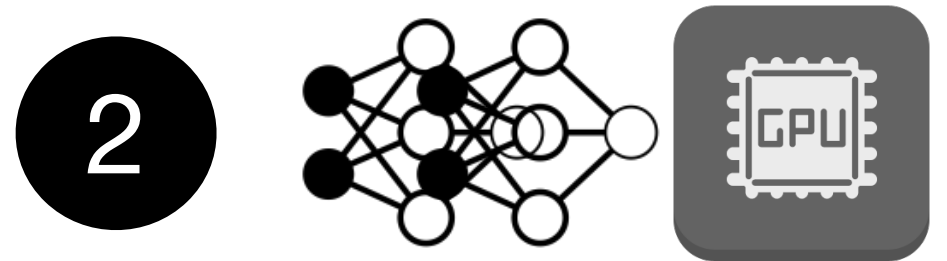
High Computational Cost

Sub-optimal Hardware Utilization 

existing h/w utilization techniques



Increasing mini-batch size
(data parallelism)

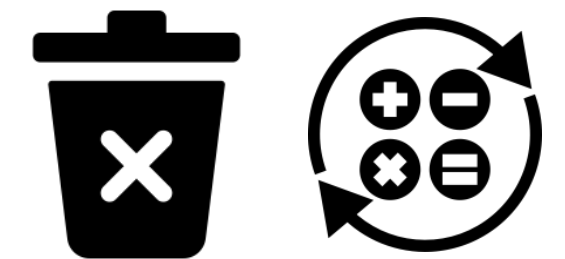


Concurrent-training
(resource sharing)

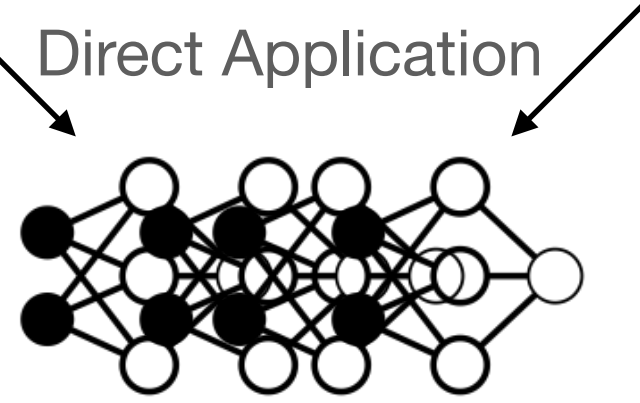
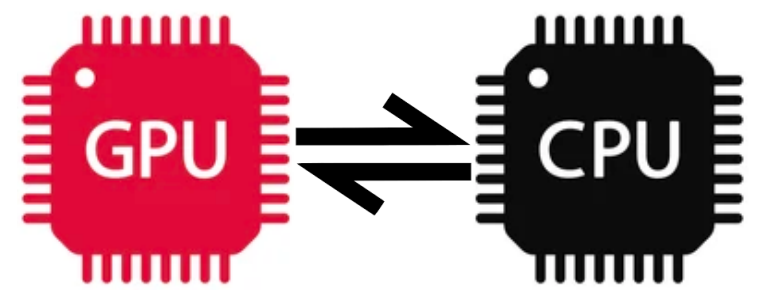


Addressing Memory Oversubscription

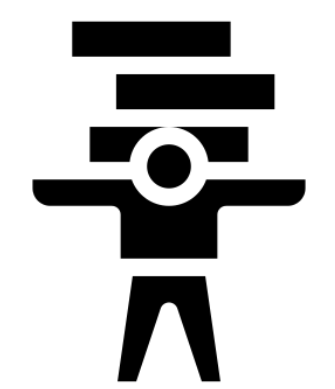
Recomputation



Swapping



Multi-Model Training

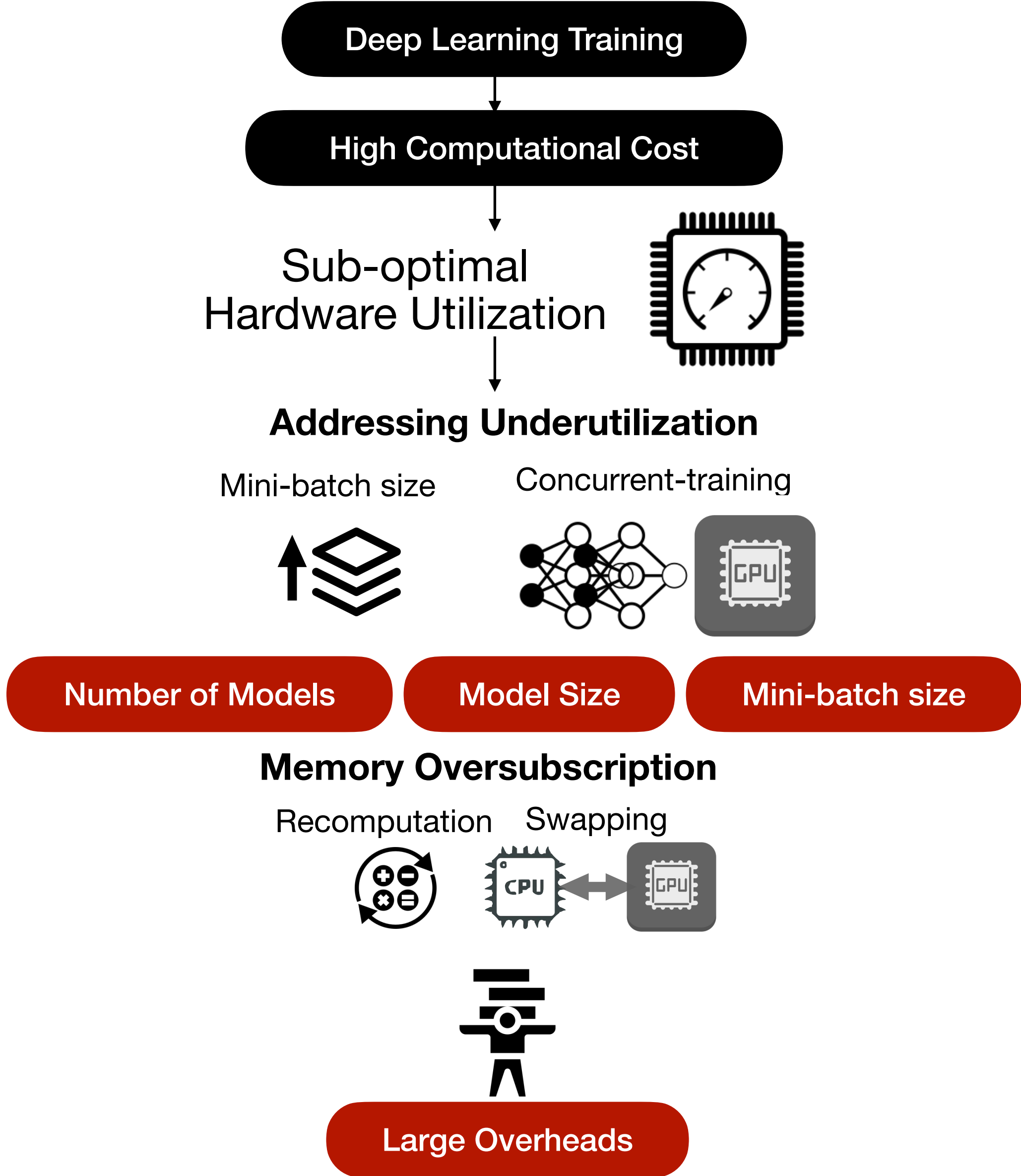


Significant Overheads

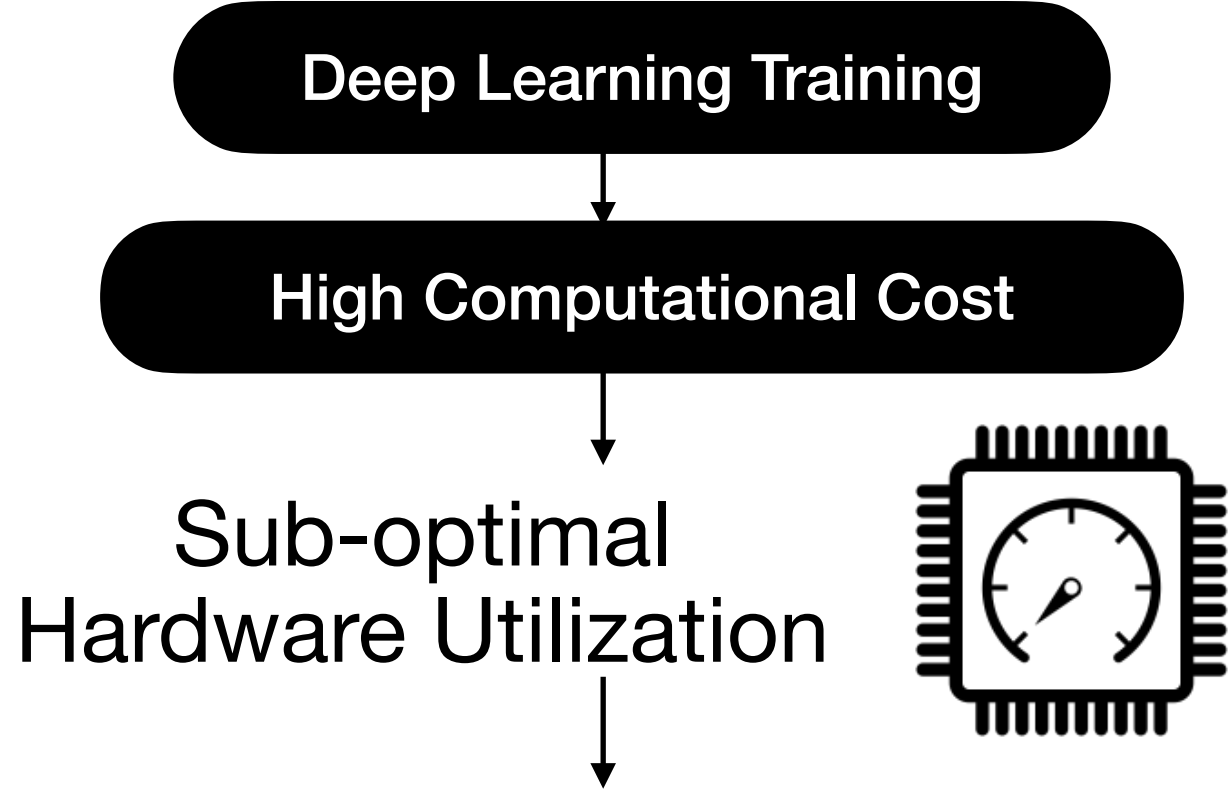


50% Slowdown

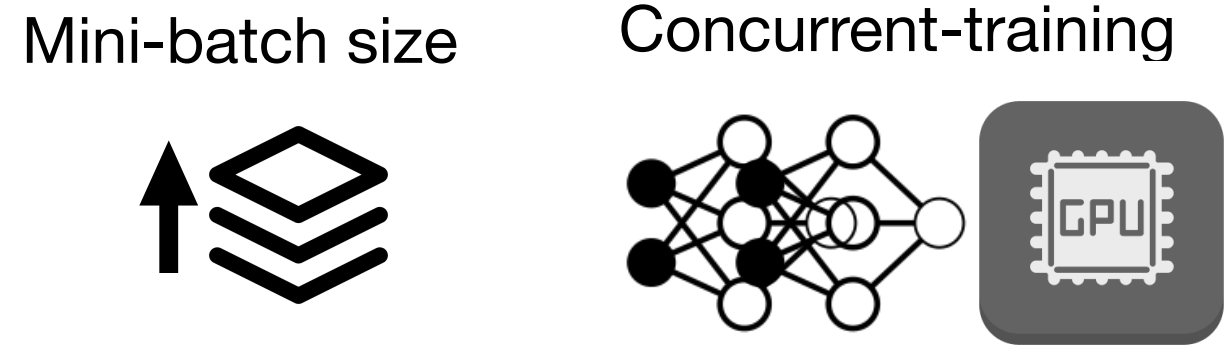
Problem Space



Problem Definition



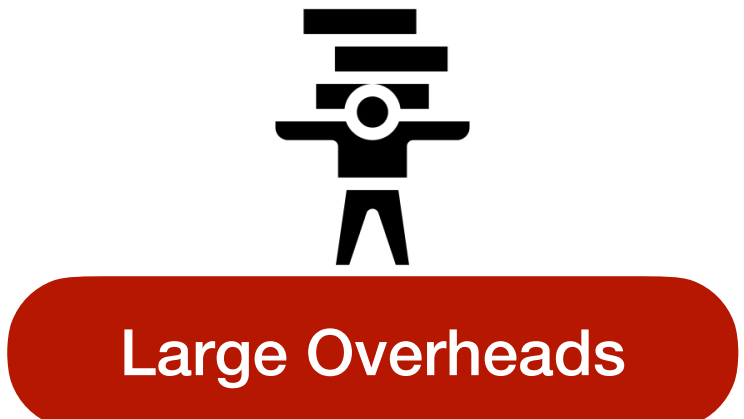
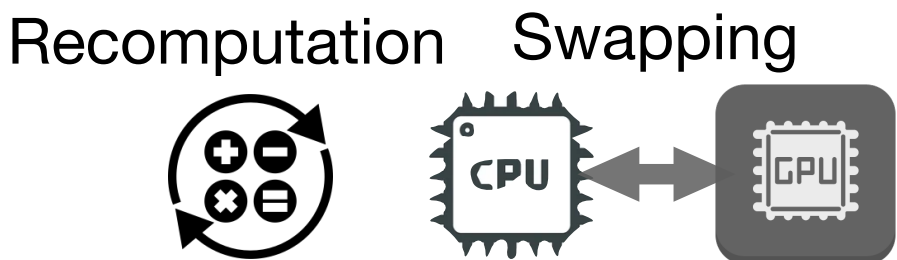
Addressing Underutilization



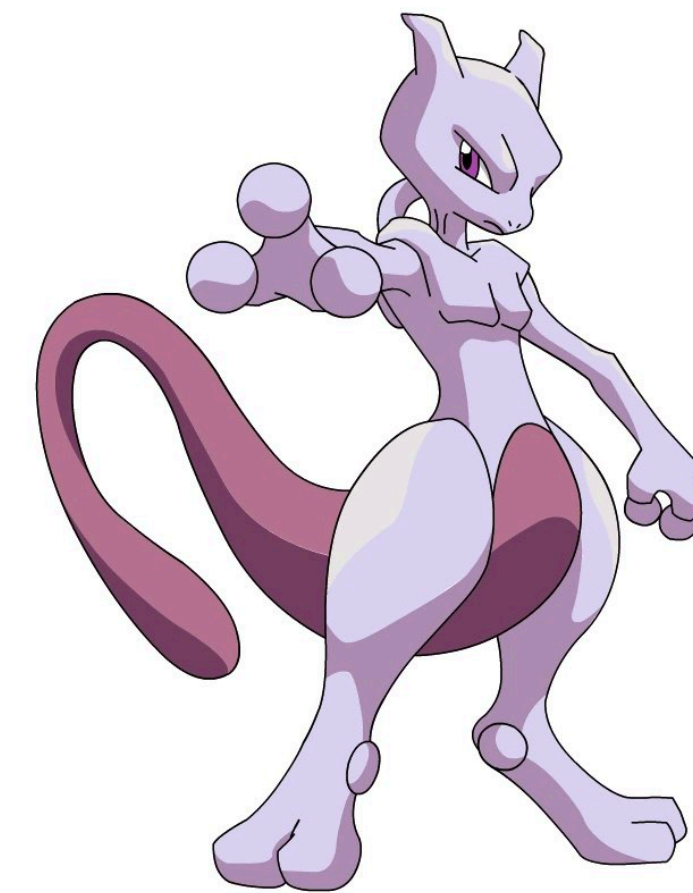
How to scale concurrent multi-model training as models grow and peak memory requirement surpasses the available GPU memory capacity?



Memory Oversubscription



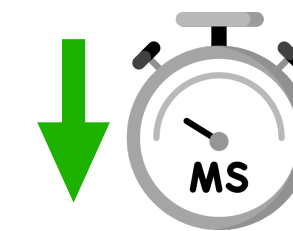
Solution



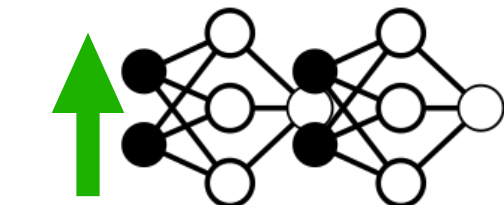
μ -TWO

Multi-model training compiler

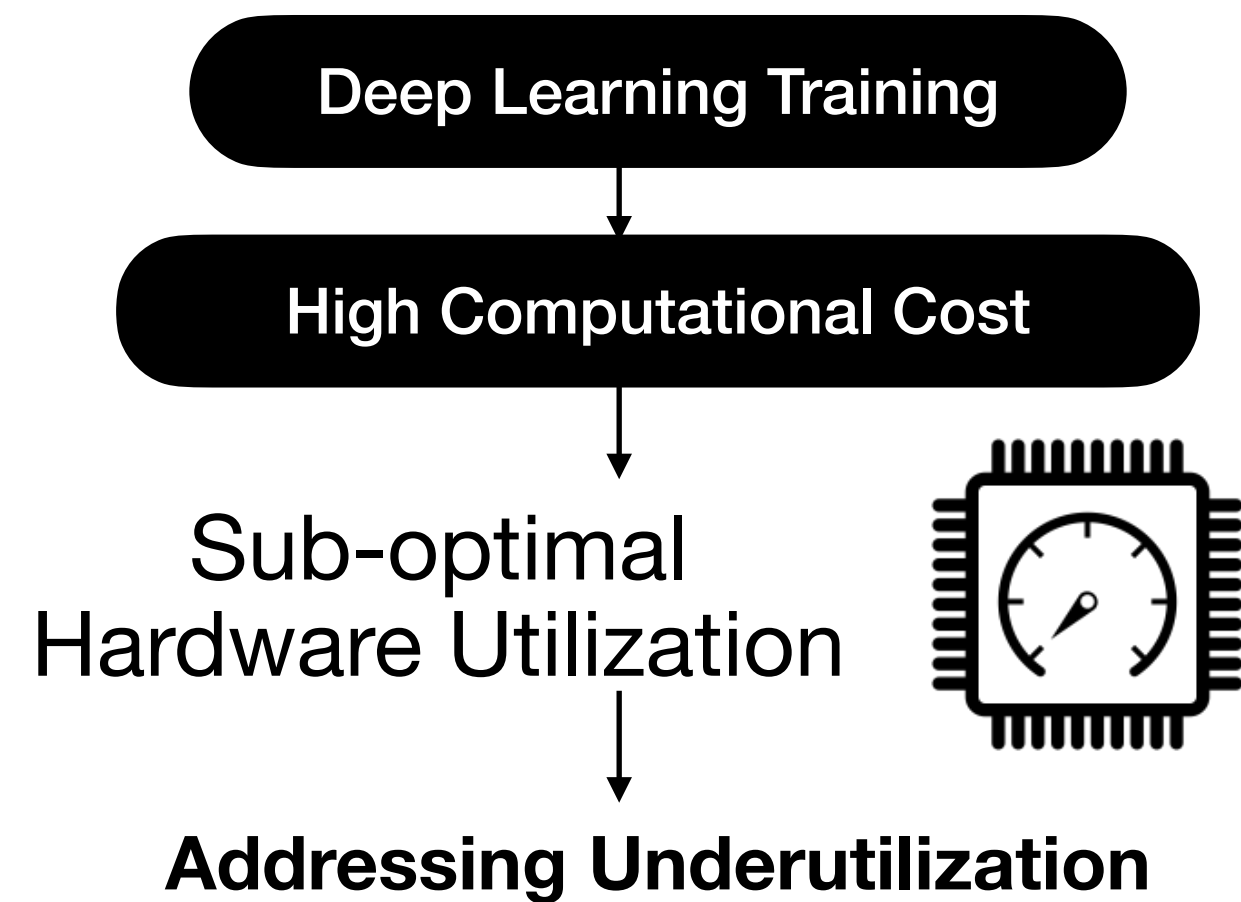
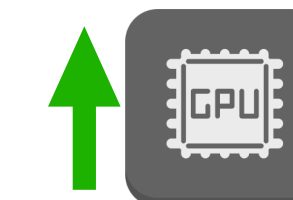
3x latency speed-up



3-5x more models

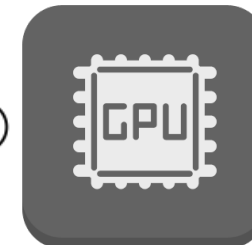
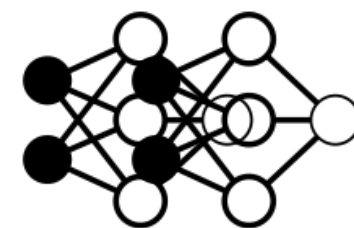
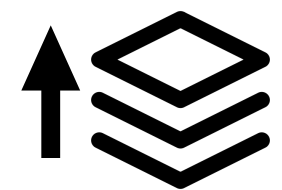


6x the GPU memory size



Mini-batch size

Concurrent-training



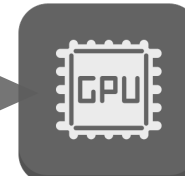
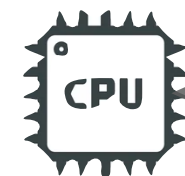
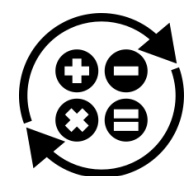
Number of Models

Model Size

Mini-batch size

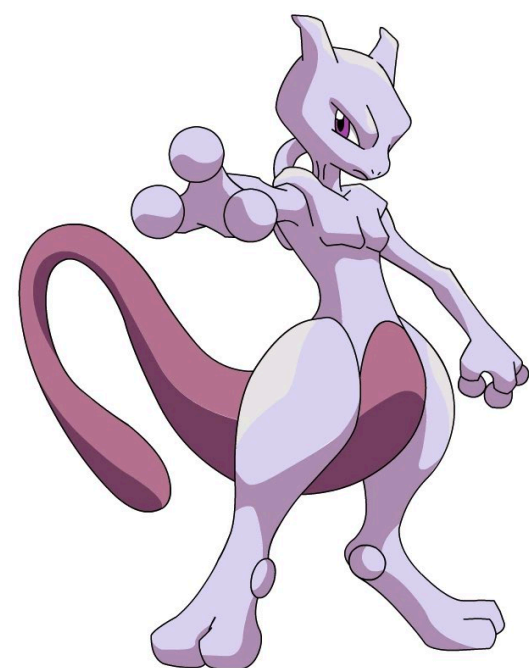
Memory Oversubscription

Recomputation Swapping



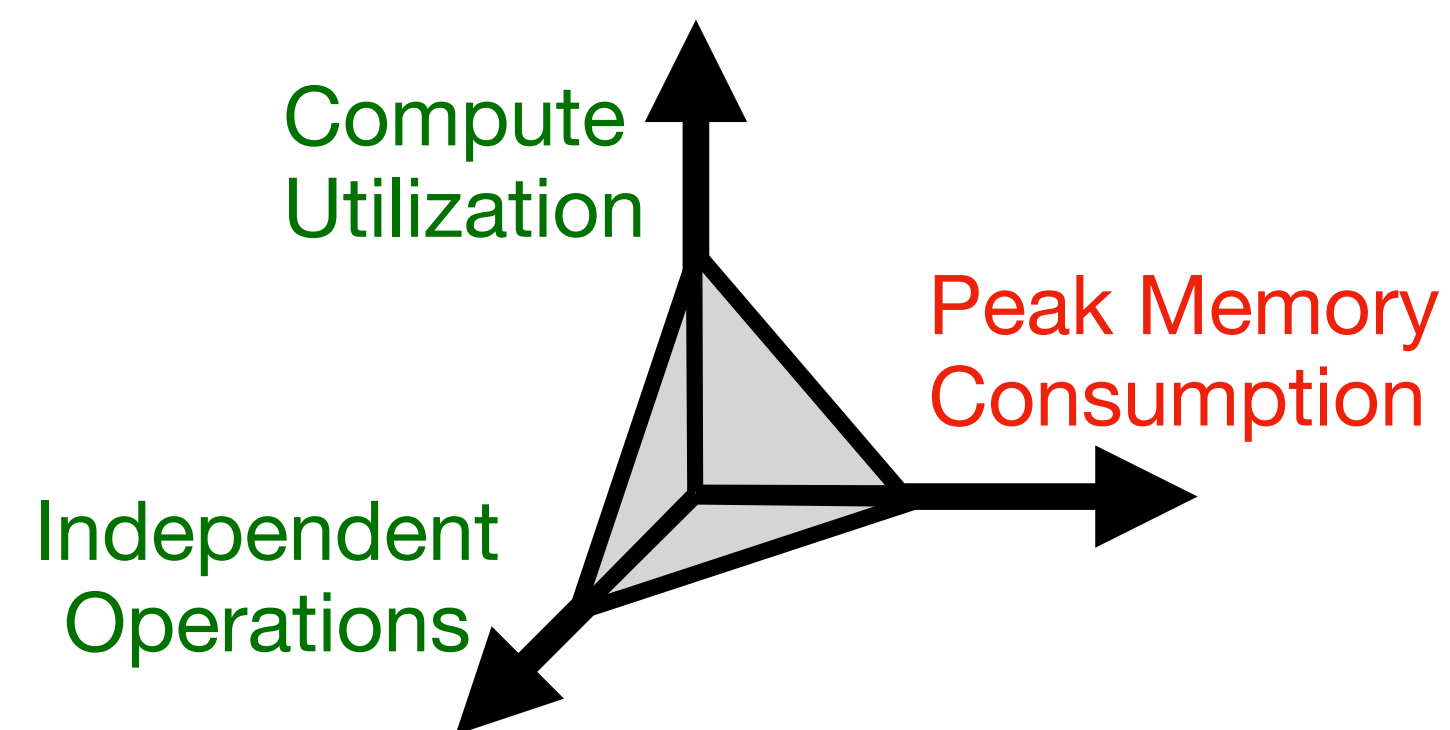
Large Overheads

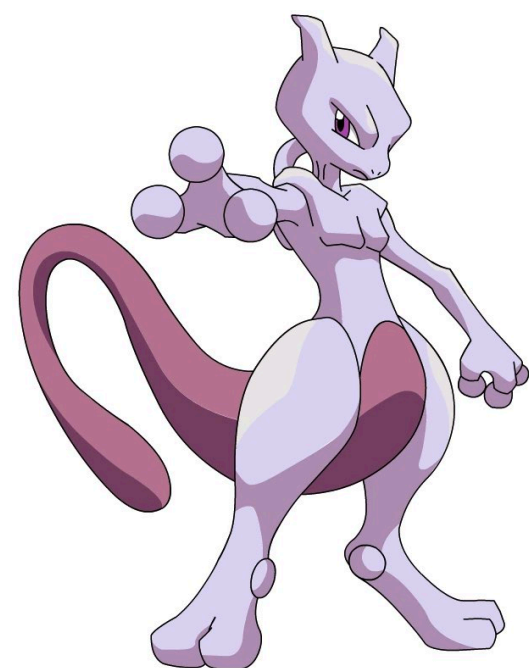




μ -TWO: Multi-Model Training with Orchestration

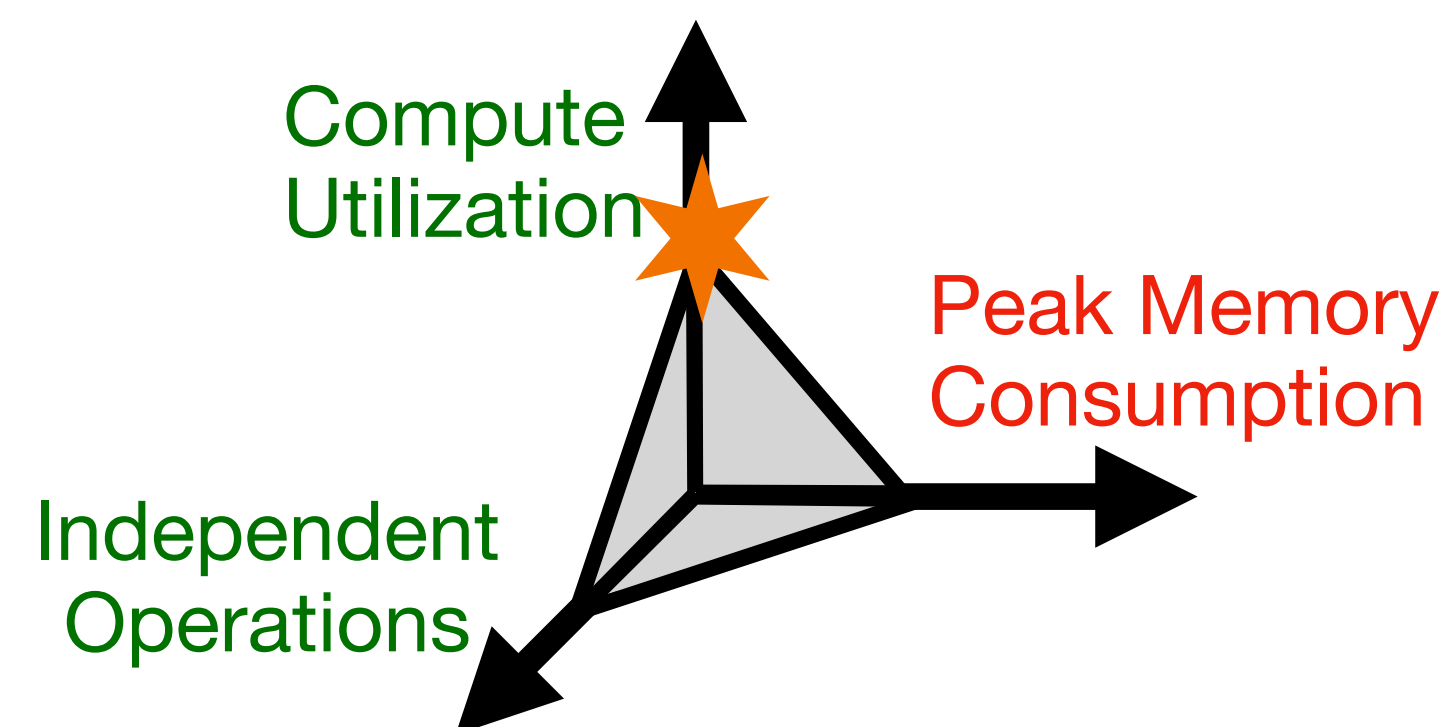
Latency of a multi-model training schedule





μ -TWO: Multi-Model Training with Orchestration

Latency of a multi-model training schedule



Efficiently navigates the trade-off for a given set of models and target GPU and finds a sweet spot

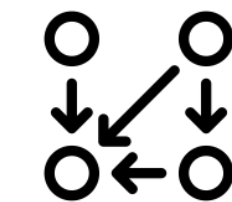


μ -TWO: Multi-Model Training with Orchestration

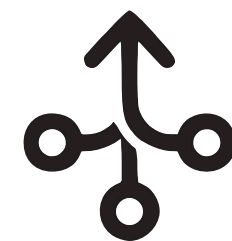
For a given set of models and target GPU



Lightweight
Profiling



Static
Analysis



Operators
to fuse

saturate compute

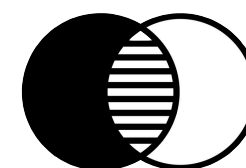


Tailored strategy



Swap/
recompute

optimize memory usage



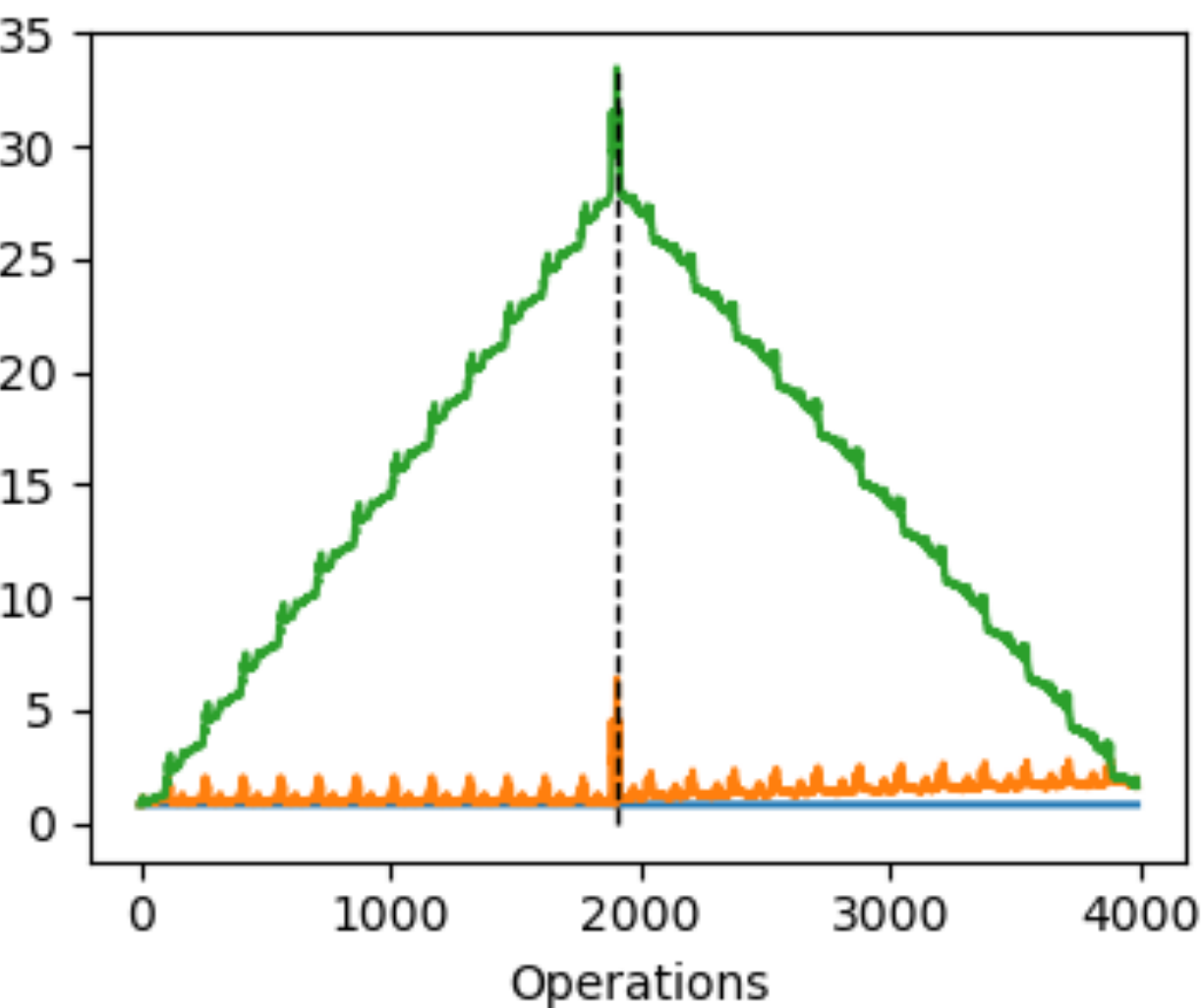
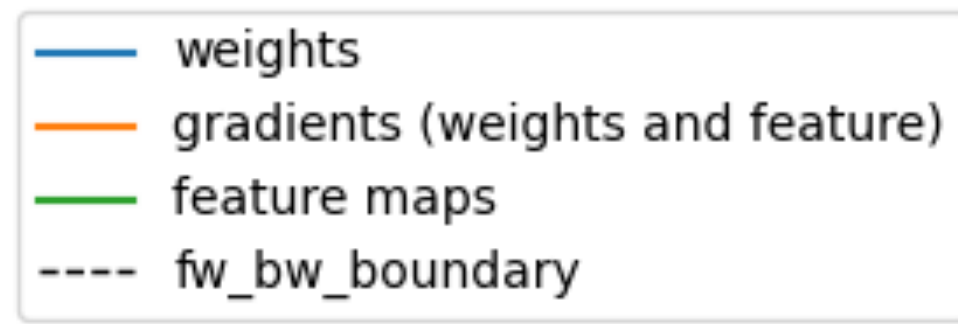
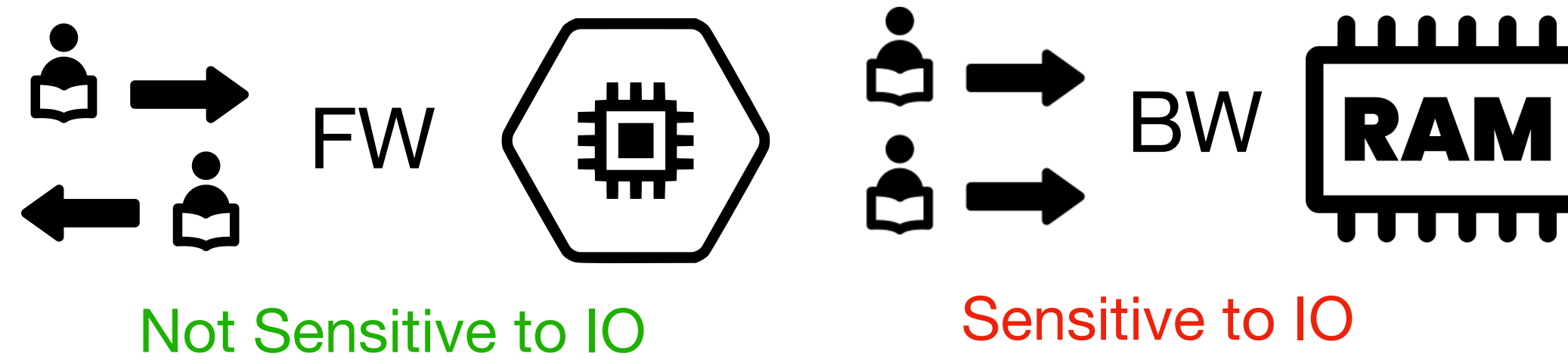
Maximally overlapped
swapping

eliminate stalling

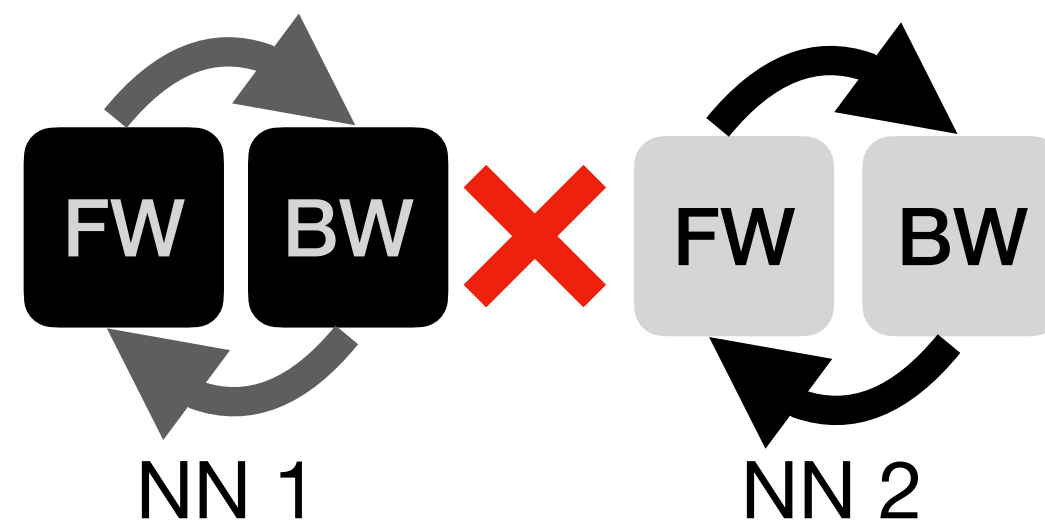


Insights for μ -TWO

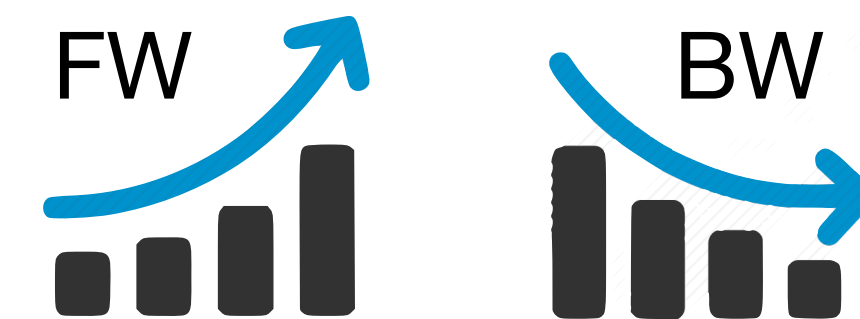
1 CHANGED NATURE OF FORWARD AND BACKWARD PASS DURING SWAPPING



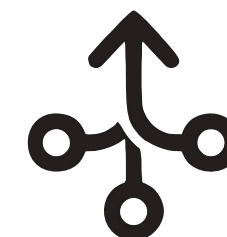
2 NO INTER-NETWORK DEPENDENCY



3 PEAK MEMORY OCCUPANCY

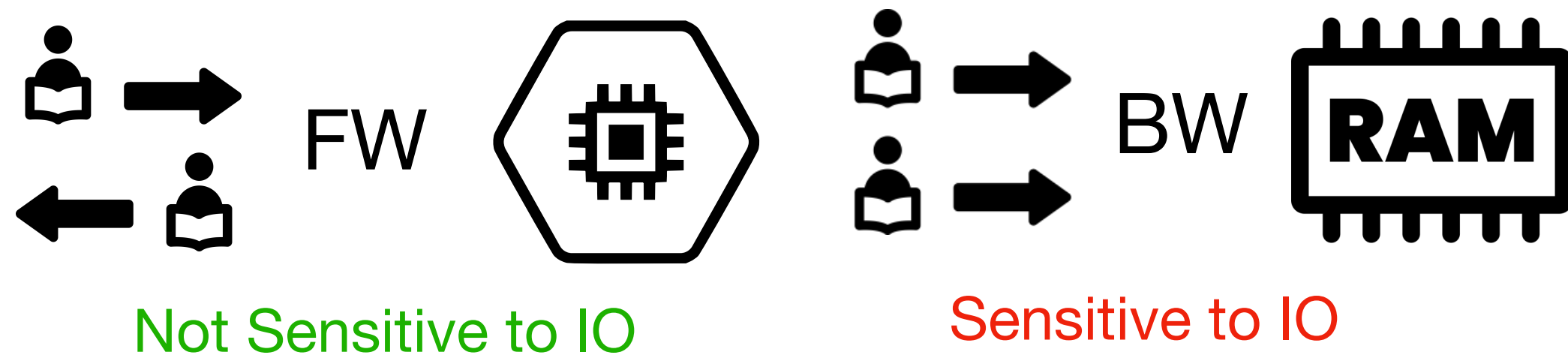


4 HORIZONTAL FUSION IS ESSENTIAL FOR COMPUTE UTILIZATION



Design Implications for μ -TWO

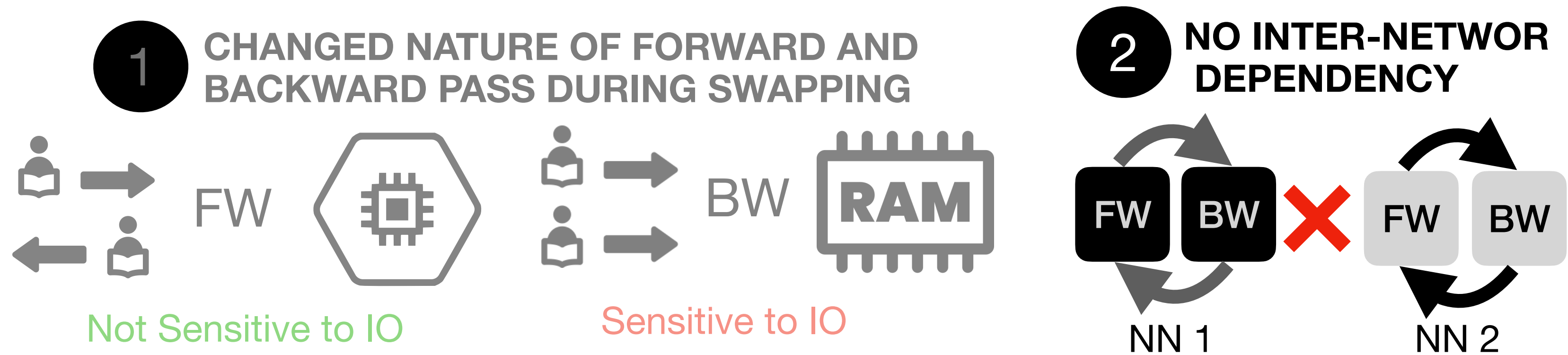
1 CHANGED NATURE OF FORWARD AND BACKWARD PASS DURING SWAPPING



1. Conservatively schedule Swapping, overlap with compute as much as possible



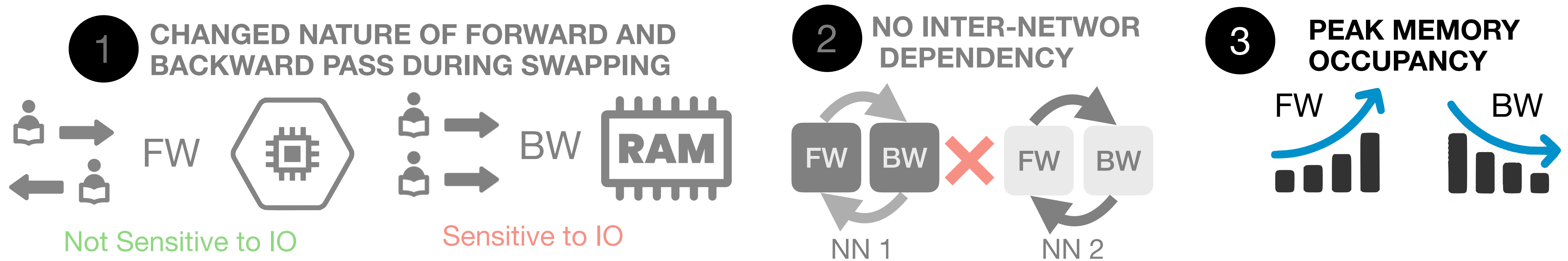
Design Implications for μ -TWO



1. Conservatively schedule Swapping, overlap with compute as much as possible
2. Operations from one model can be used to overlap IO operations of other models



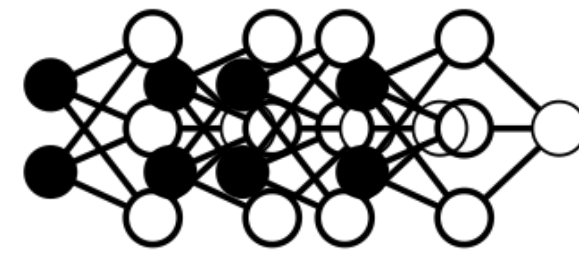
Design Implications for μ -TWO



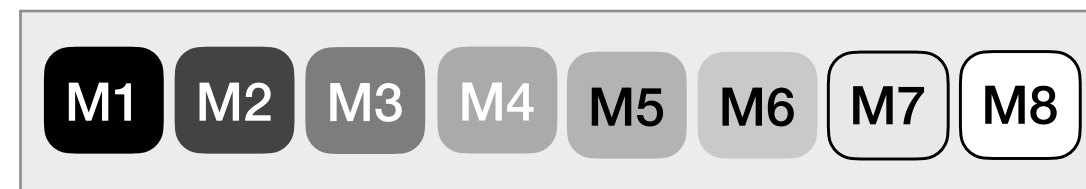
1. Conservatively schedule Swapping, overlap with compute as much as possible
2. Operations from one model can be used to overlap IO operations of other models
3. Only forward pass operations should be used for overlapping backward pass IO operations



Design Trade-Offs for μ -TWO



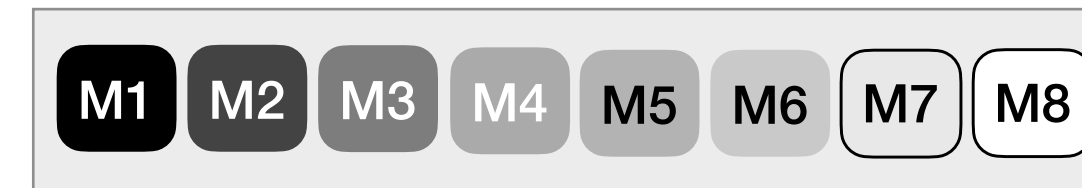
Multiple Models



Monolithic inseparable forward and backward pass operations
(Minimal opportunity to overlap)

High Compute Utilization
due to maximal fusion

High Peak memory consumption
due to fusion



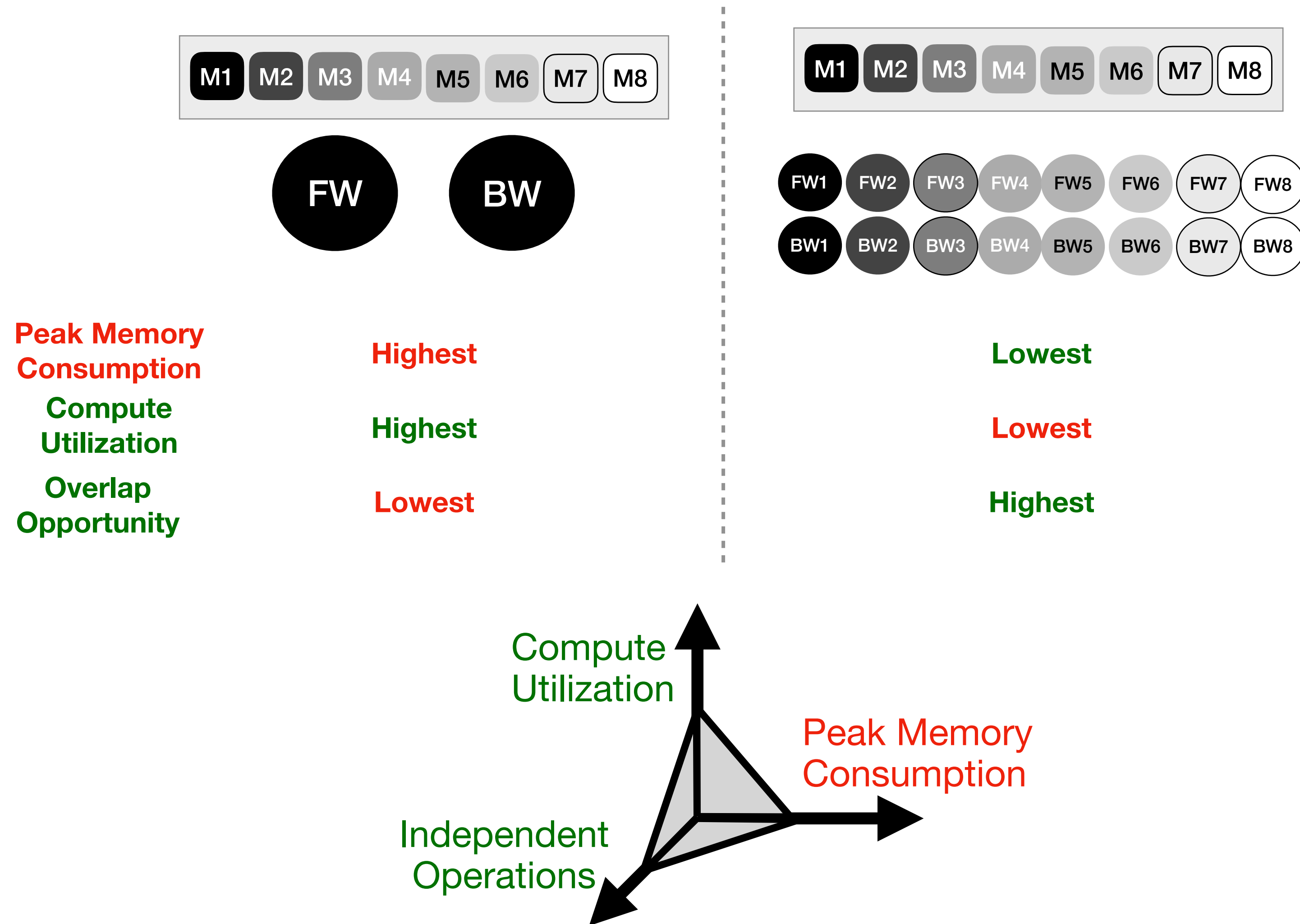
Separate forward and backward pass operations
(Maximum opportunity to overlap)

Low Compute Utilization
due to no fusion

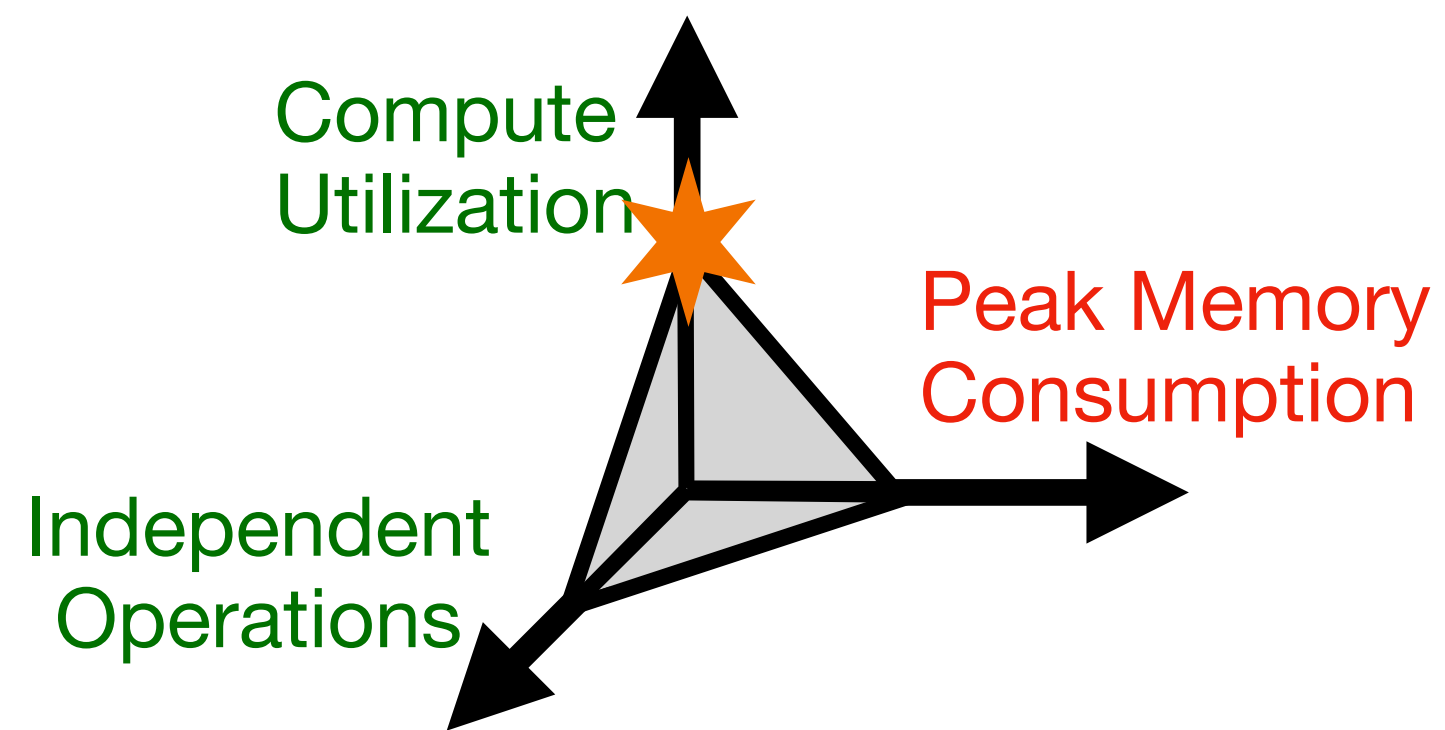
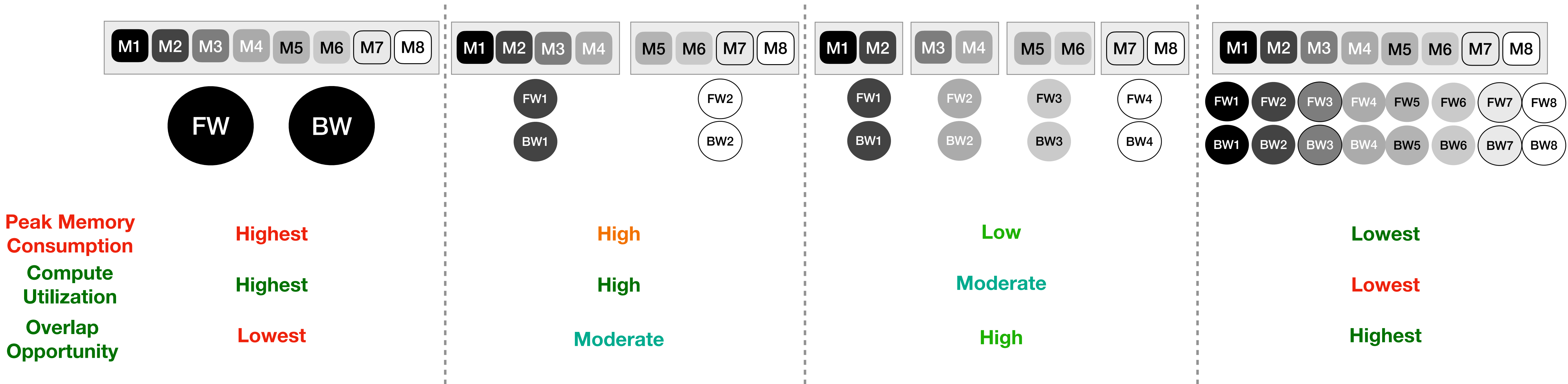
Lowest Peak memory consumption
due to absence of fusion



Design Trade-Offs for μ -TWO



Design Trade-Offs for μ -TWO



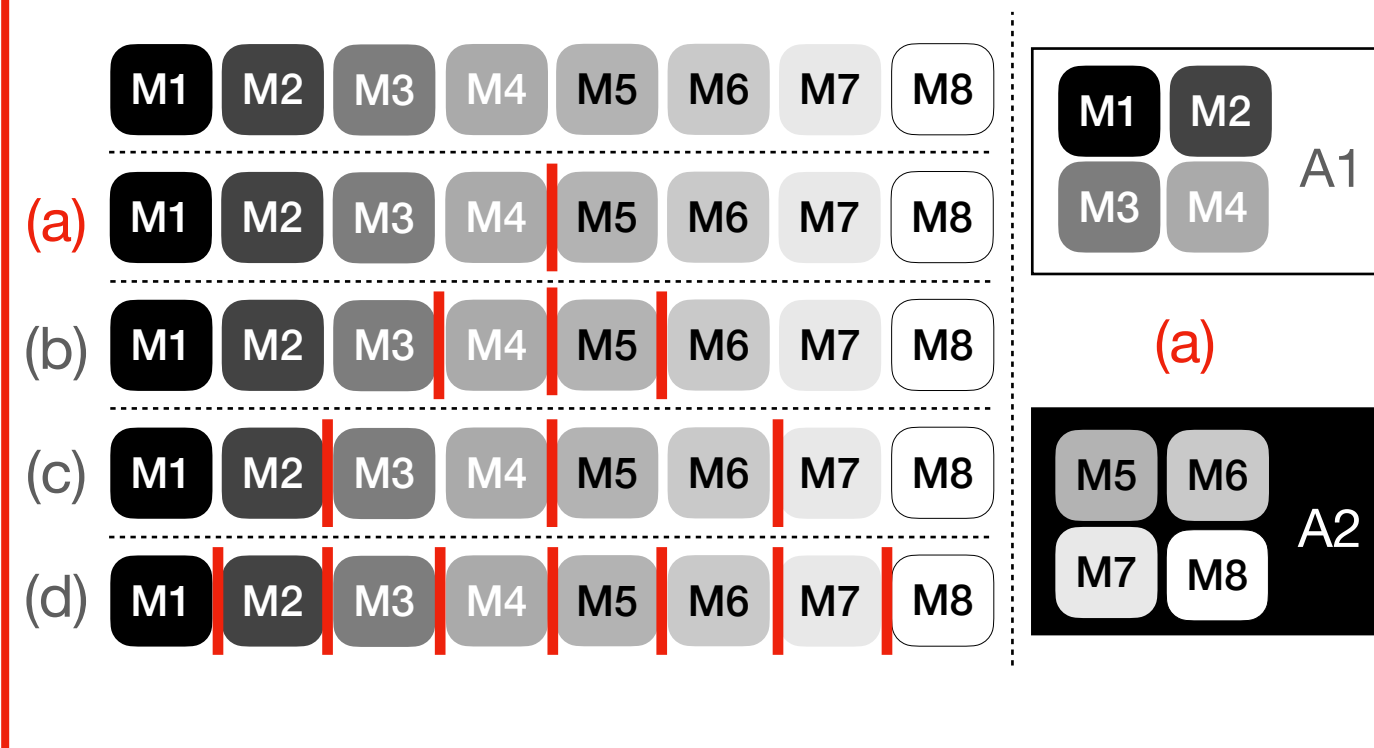
Use fusion granularity to navigate the trade-off



μ -TWO System Overview

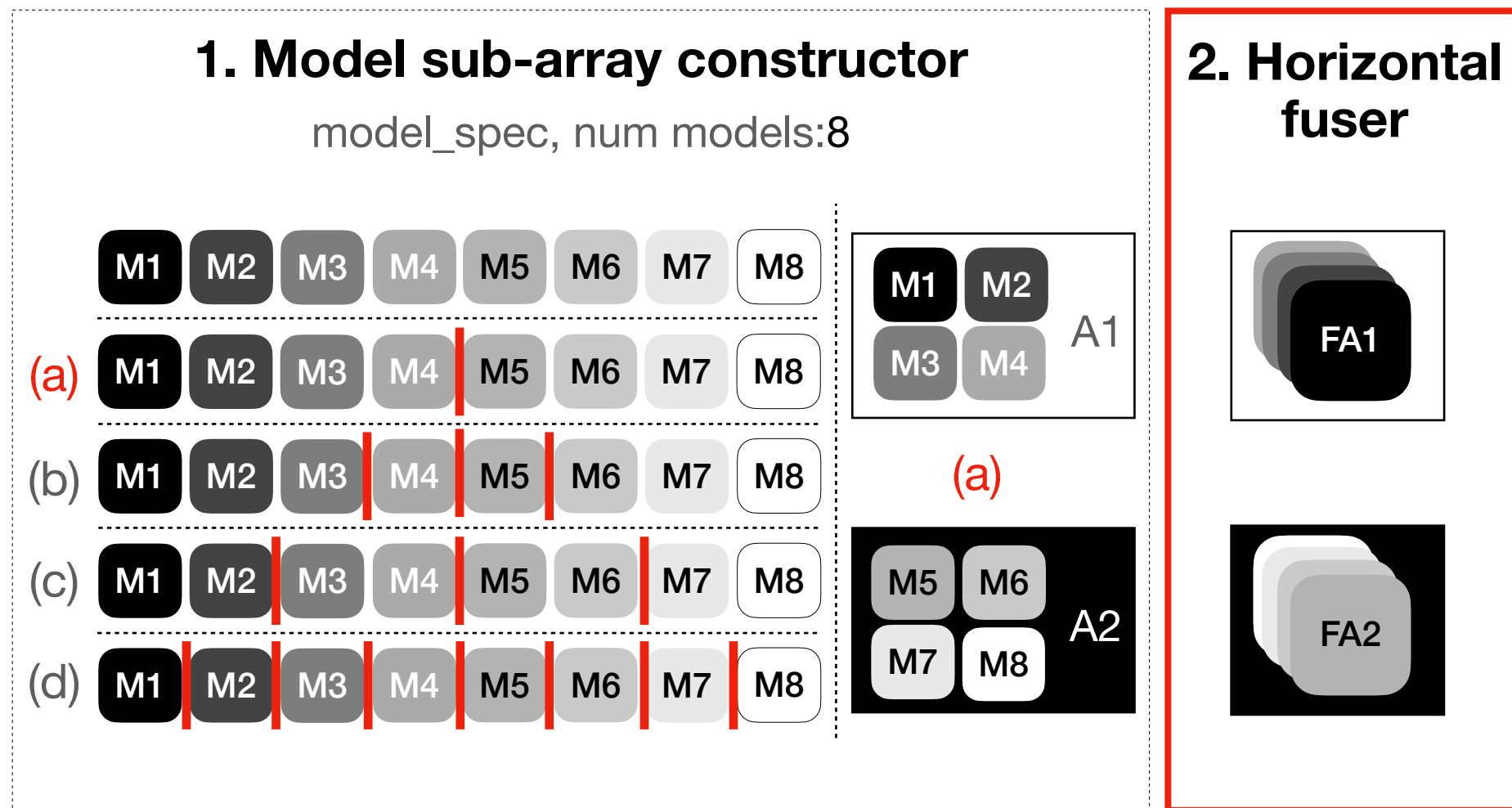
1. Model sub-array constructor

model_spec, num models:8



Enumerate possible sub-array partitions of models

μ -TWO System Overview



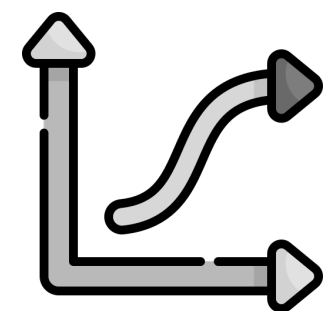
Horizontally fuse models within each sub-array

Models in each sub-array must have same architecture to fuse kernels

Can have different hyper parameters

Loss function

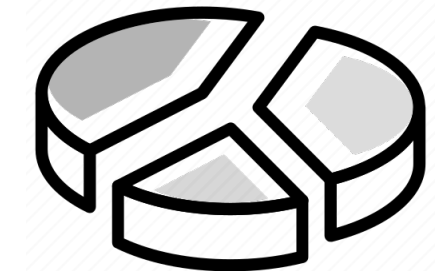
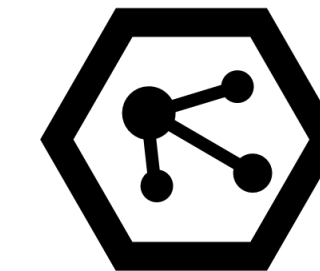
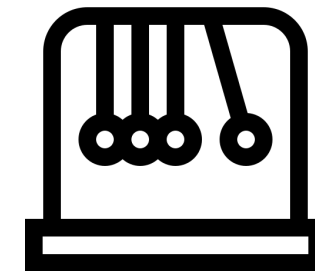
Weight initialization



Learning rate

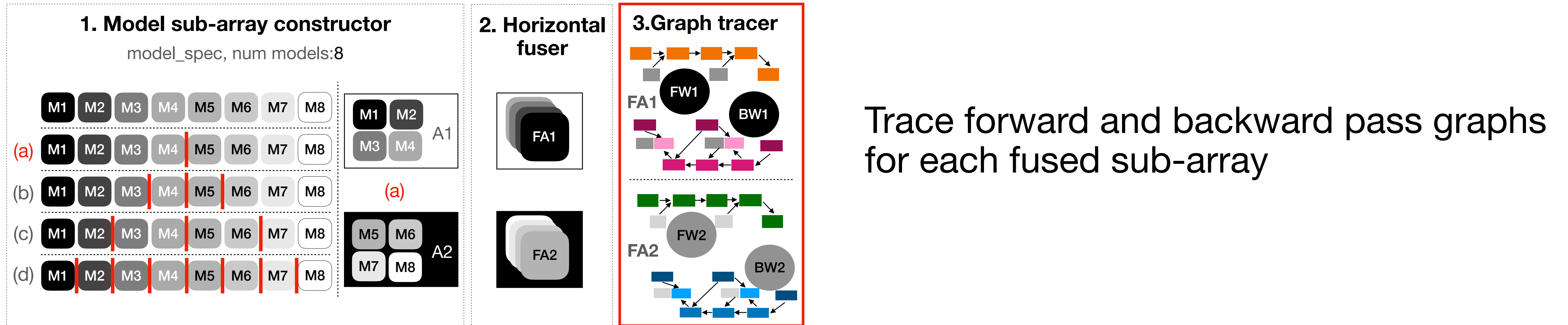


Momentum

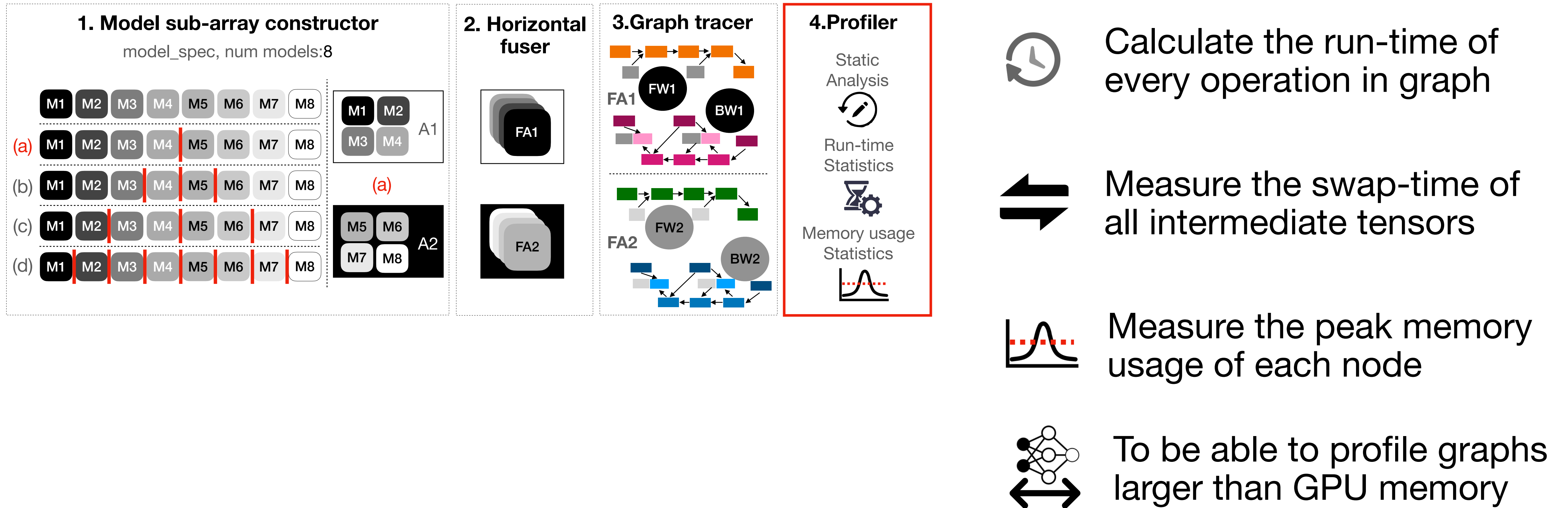


Can train on different data set/partitions

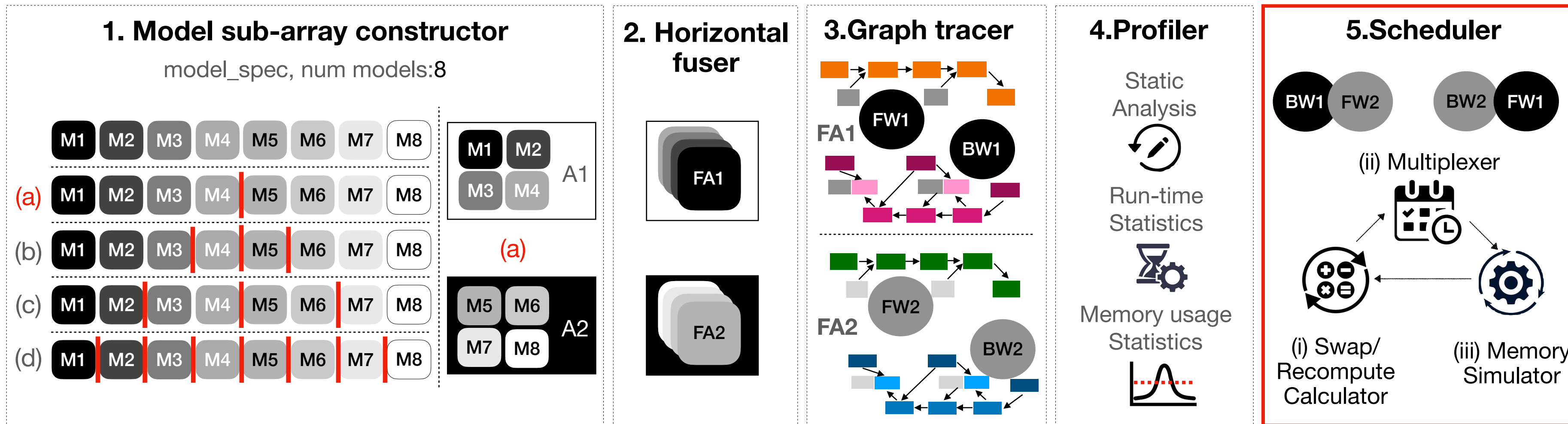
μ -TWO System Overview



μ -TWO System Overview

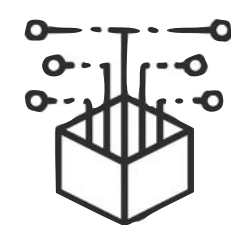


μ-TWO System Overview



Choose whether an intermediate tensor should be

Recomputed



Sources to recompute the tensor?



Recompute overhead

Swapped

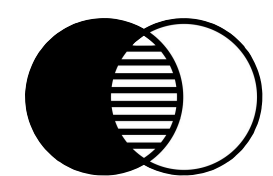
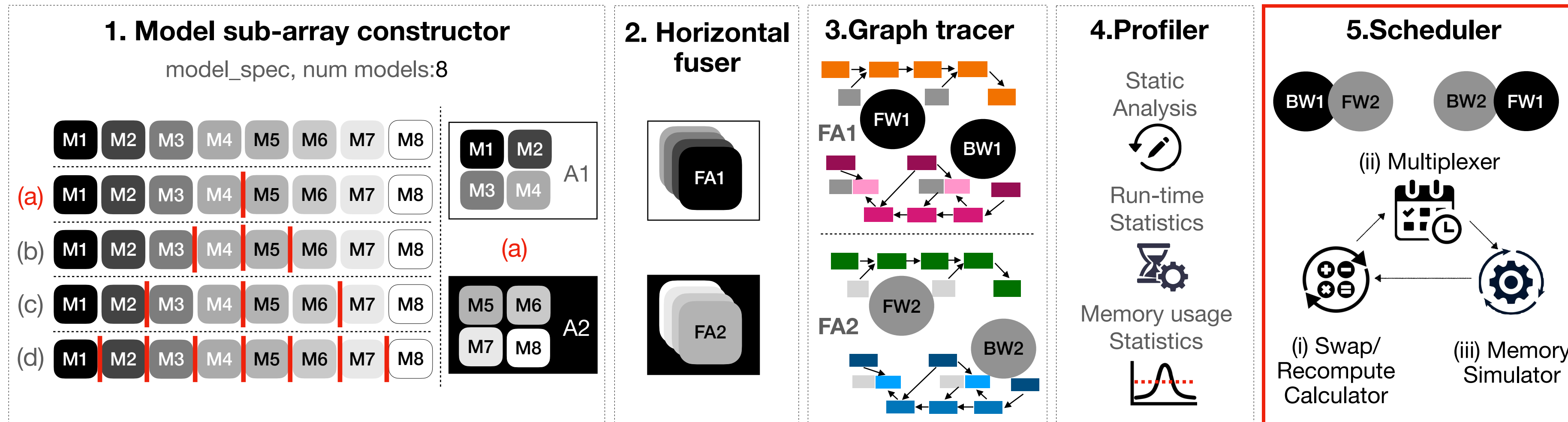


When to prefetch/offload



Swap overhead

μ -TWO System Overview

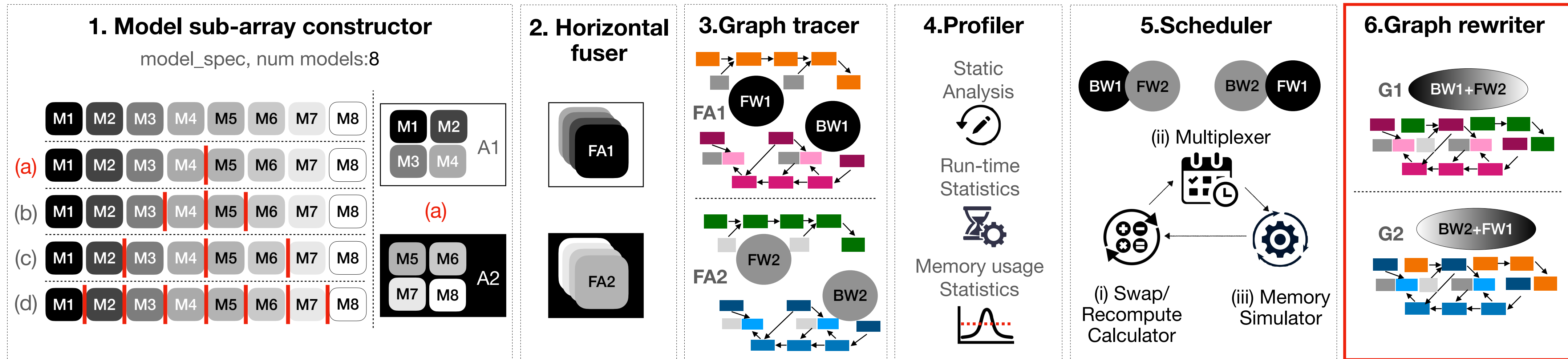


Multiplexer overlaps all the swaps in the backward pass of one fused sub-array with forward pass of another



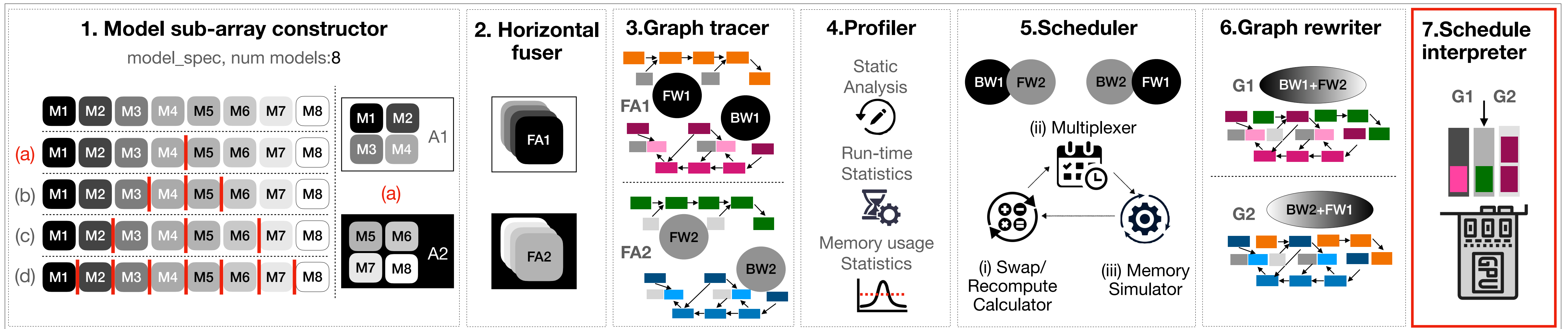
Memory simulator ensures that the peak memory consumption post the reordering pass is within the GPU memory limit

μ -TWO System Overview



Embed the scheduling information in the graphs

μ -TWO System Overview



Enqueue the appropriate operations across different execution queues

Add synchronization markers for coordination across queues

Experimental Setup

Instance	Nvidia GPU Version	GPU Mem (GB)	Tensor Cores	CPU-GPU Link	CPUs	CPU Mem
AWS p4d24-large	A-100	40	Yes	PCI-e Gen 4 x16 (32GB/s)	16	1152
Dell Claudron DSS 8440	Tesla V-100	32	Yes	PCI-e Gen 4 x16 (32GB/s)	16	384



Workload

Application	Model Name	Functionality	Architectural Features	Params	Batch Sizes
Vision	Vision Transformer	Image Classification, Image Segmentation, Action Recognition	Positional image embeddings, transformers	60M	8 16
	Mobilenet v3 large		Depthwise separable convolutions	5.4M	64 128
	Resnet101		Convolutions, Skip Connections	44.5M	48 64
Natural Language Processing	Bert	Predict Next Sentence	Transformer Encoders	100M	16 24
	GPT2	Predict Next Token	Transformer Decoders	124M	8 16
Recommender Systems	NVIDIA DLRM	Item Recommendation	Encoders, Decoders, sparse embeddings	40M	512 1024

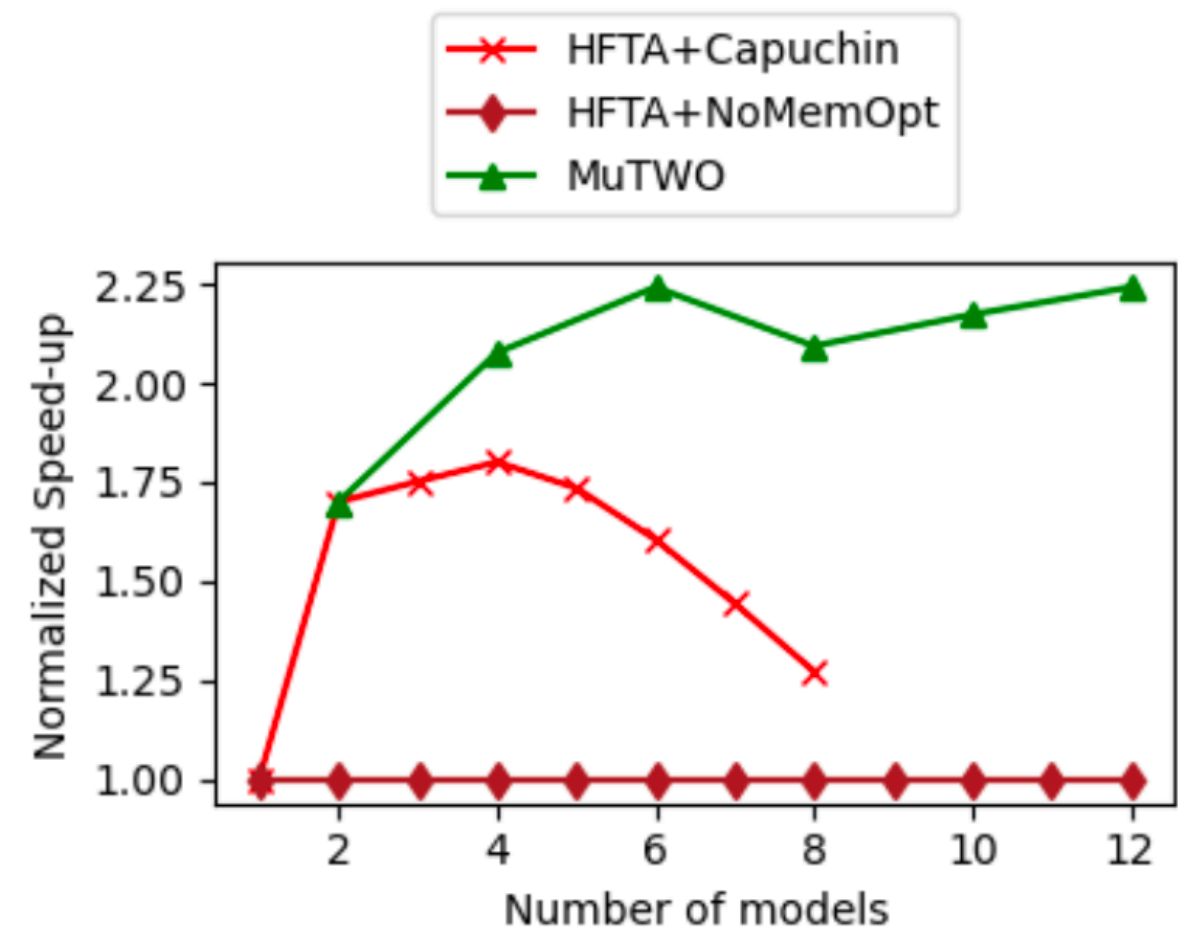


Baselines

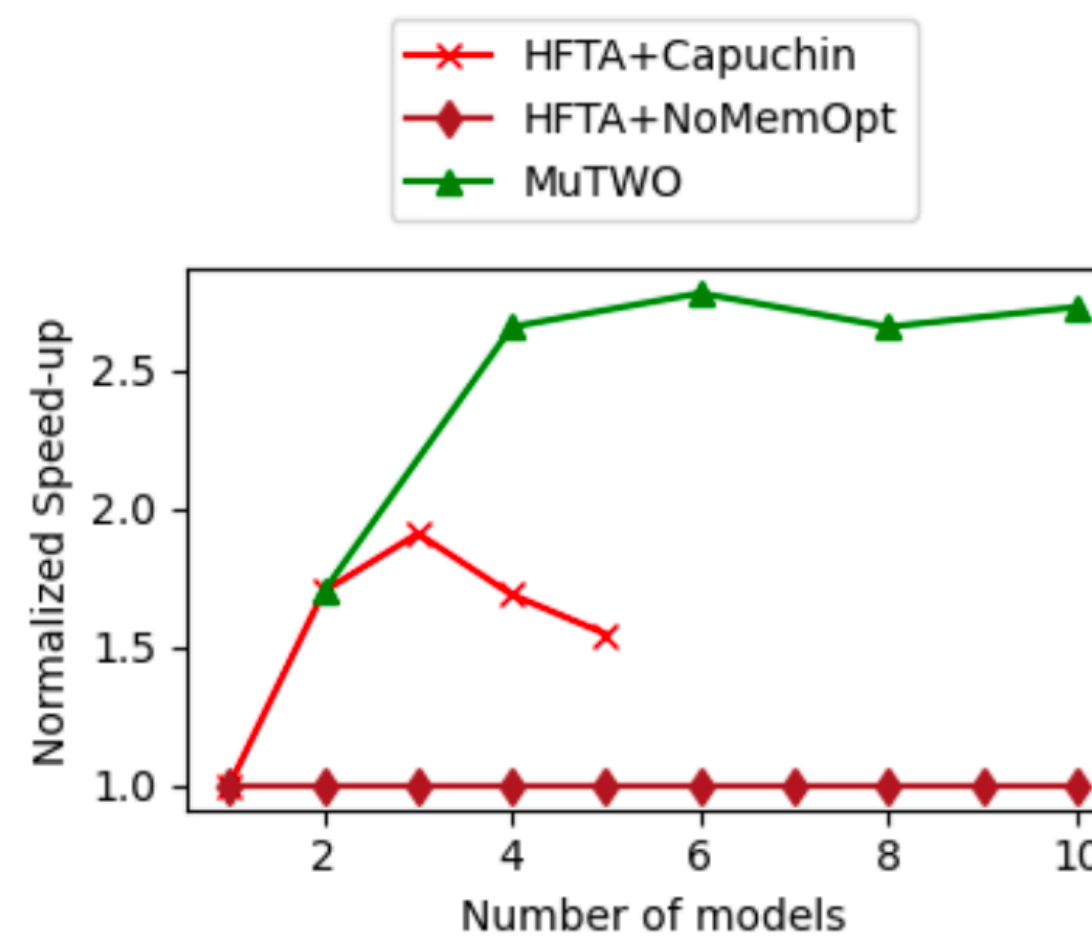
- HFTA-NoMemOpt - Horizontal Fusion only with no memory optimization
- HFTA-Capuchin - HFTA with Capuchin Algorithm applied directly
- μ -TWO - Multi-model training with orchestration and memory optimization



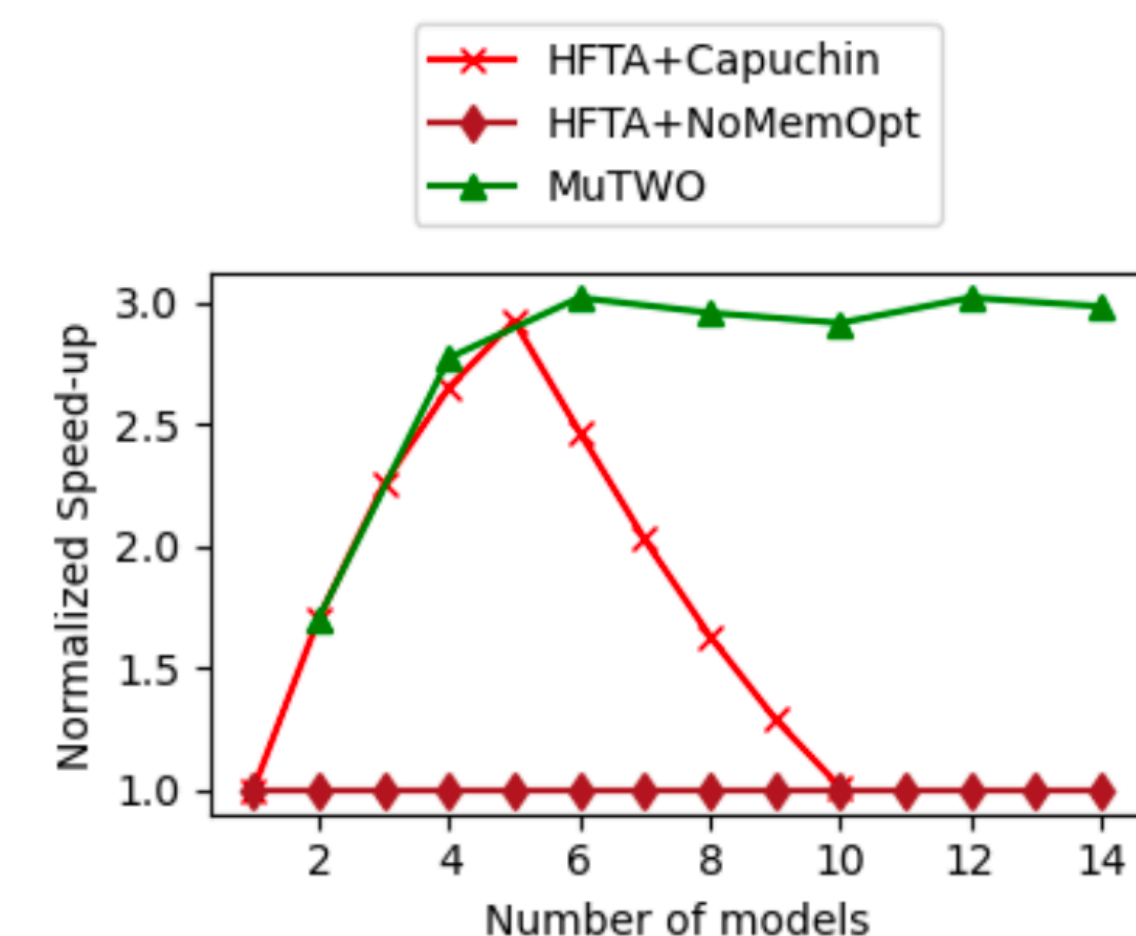
μ -TWO achieves upto 3x Speed-up



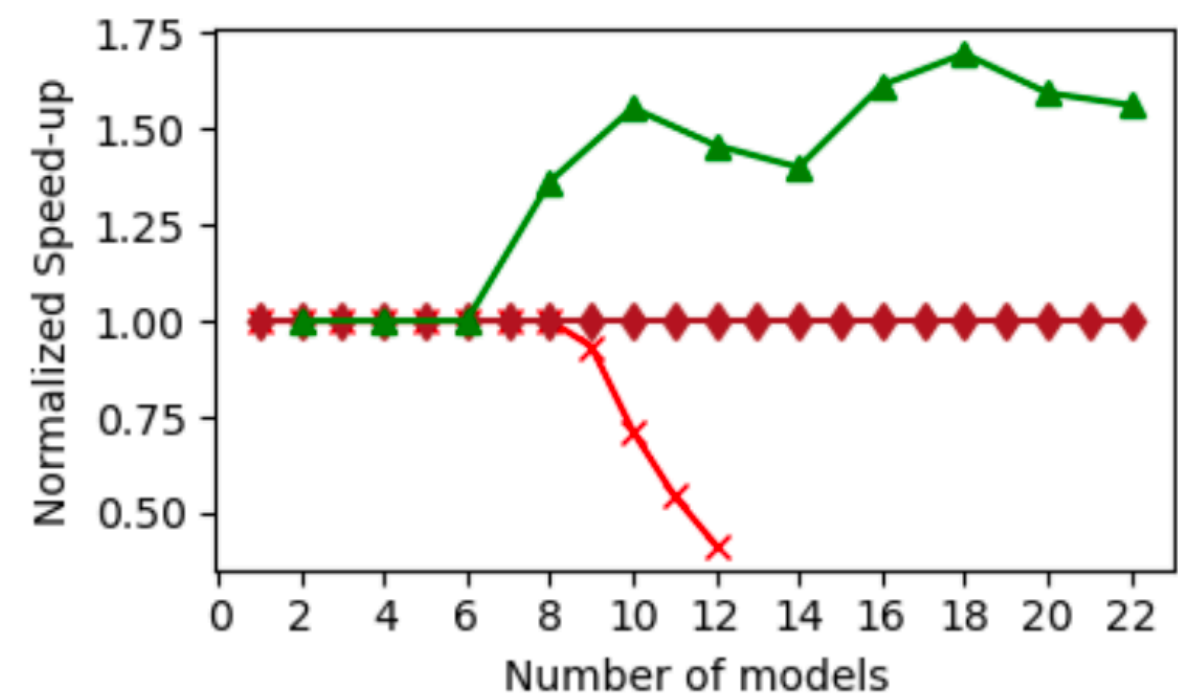
(a) Bert (Batch size: 24)



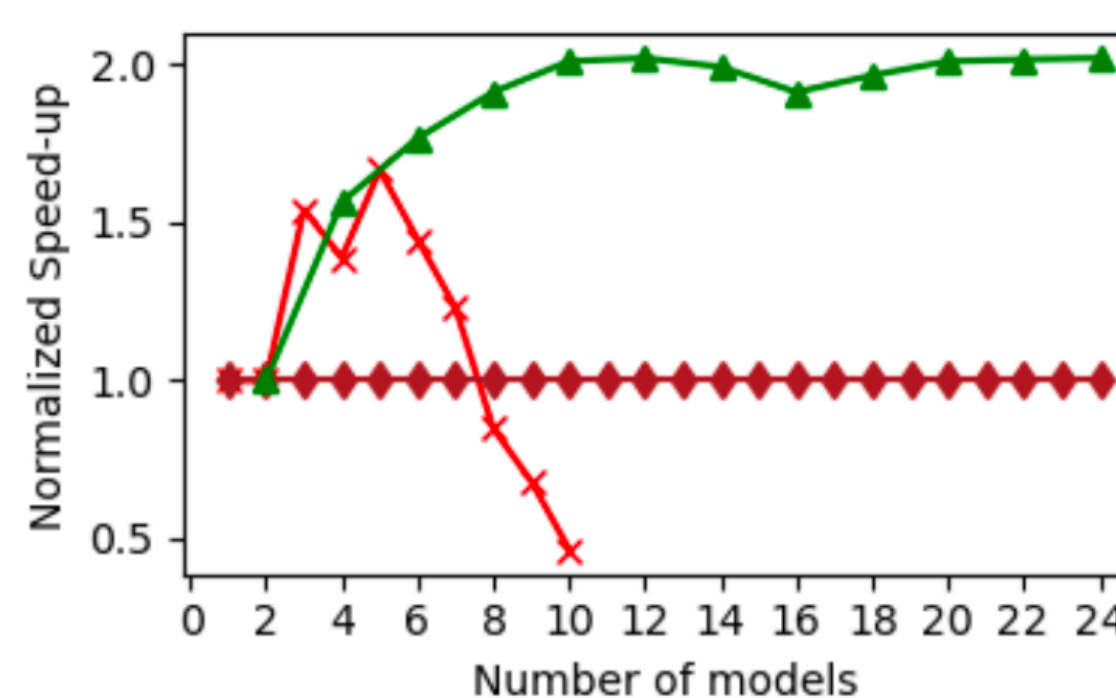
(b) GPT2 (Batch size: 16)



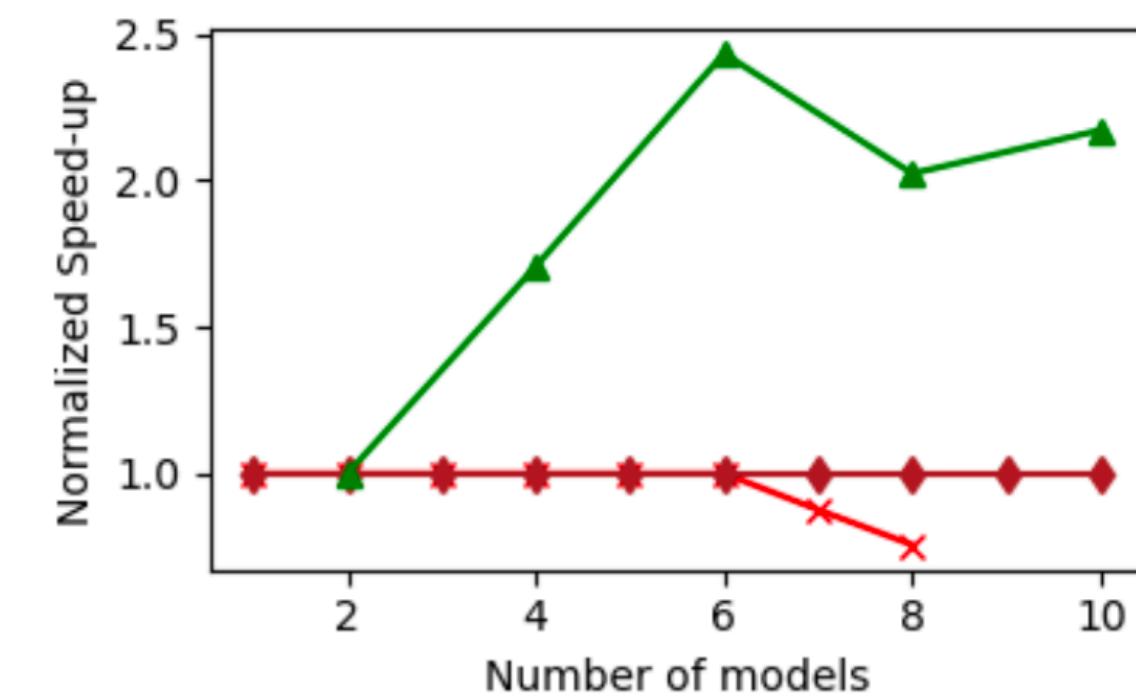
(c) Vision Transformer (Batch size: 8)



(d) Mobilenet v3 large (Batch size: 64)



(e) Resnet101 (Batch size: 48)



(f) NV DLRM (Batch size: 1024)



Performance Breakdown

Useful Compute: Computation time spent in necessary operations

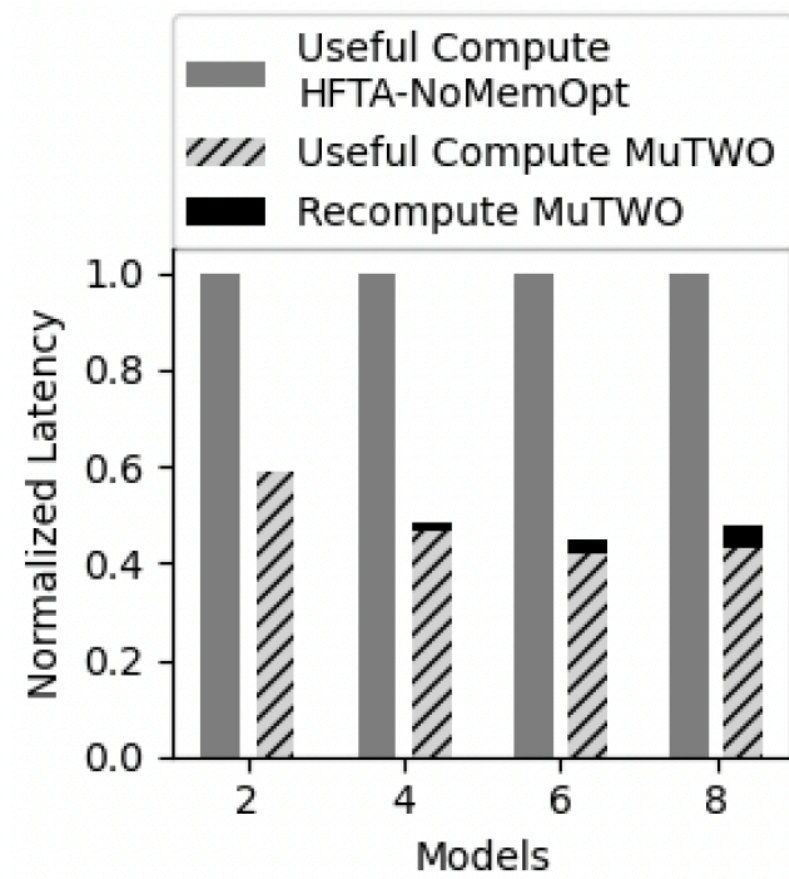
Recomputation: Computation time spent in recompute operations

Swap Overlap: Successful overlap with compute operations

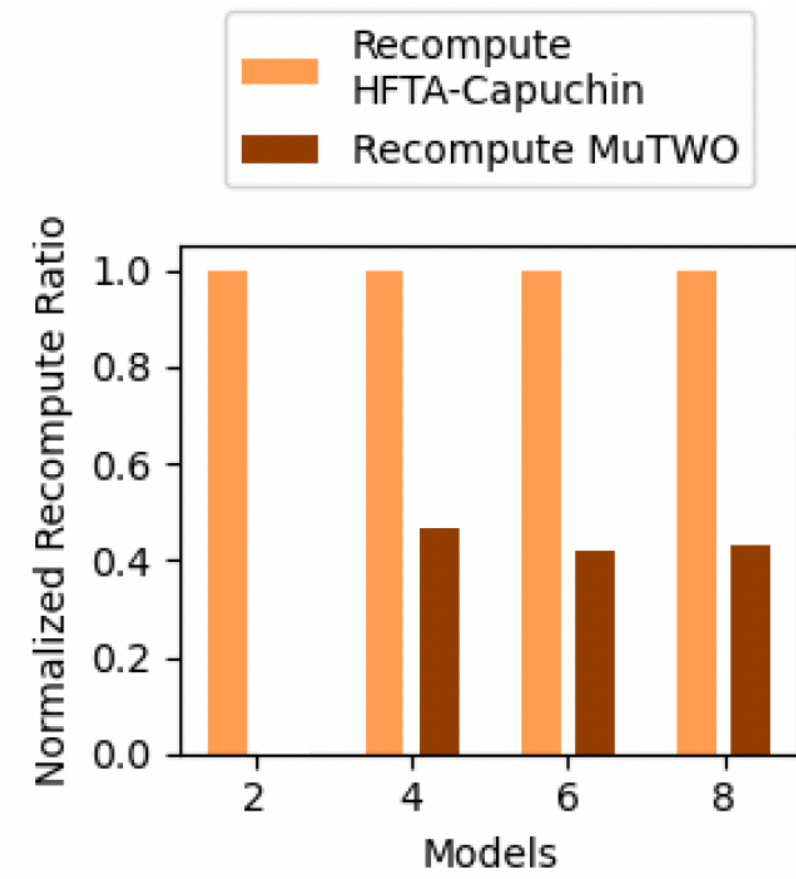
Peak memory Consumption: Maximum memory consumed at any point during the entire iteration



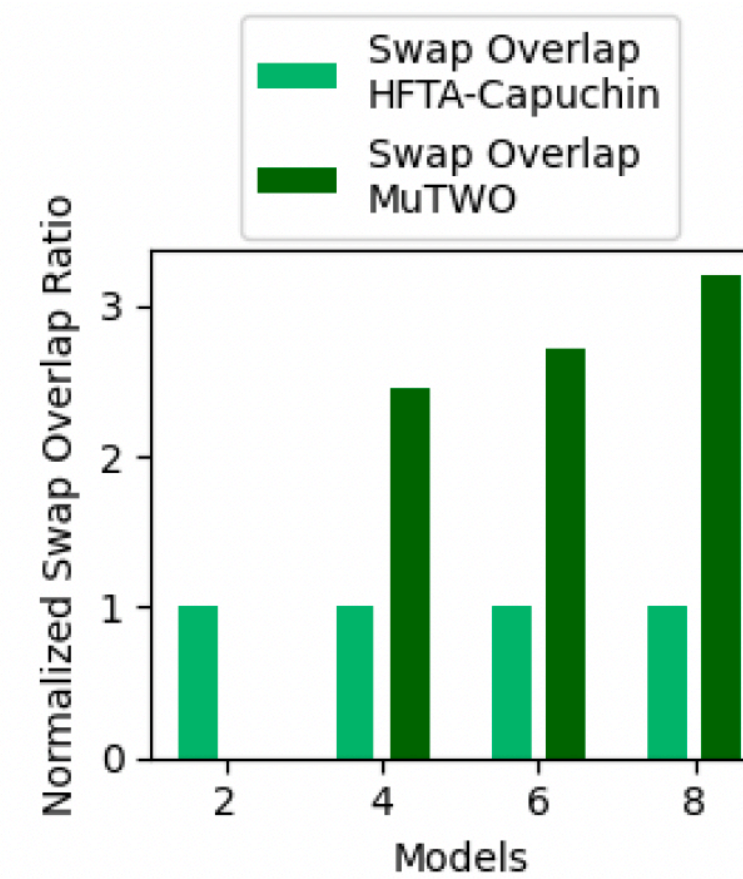
Performance Breakdown



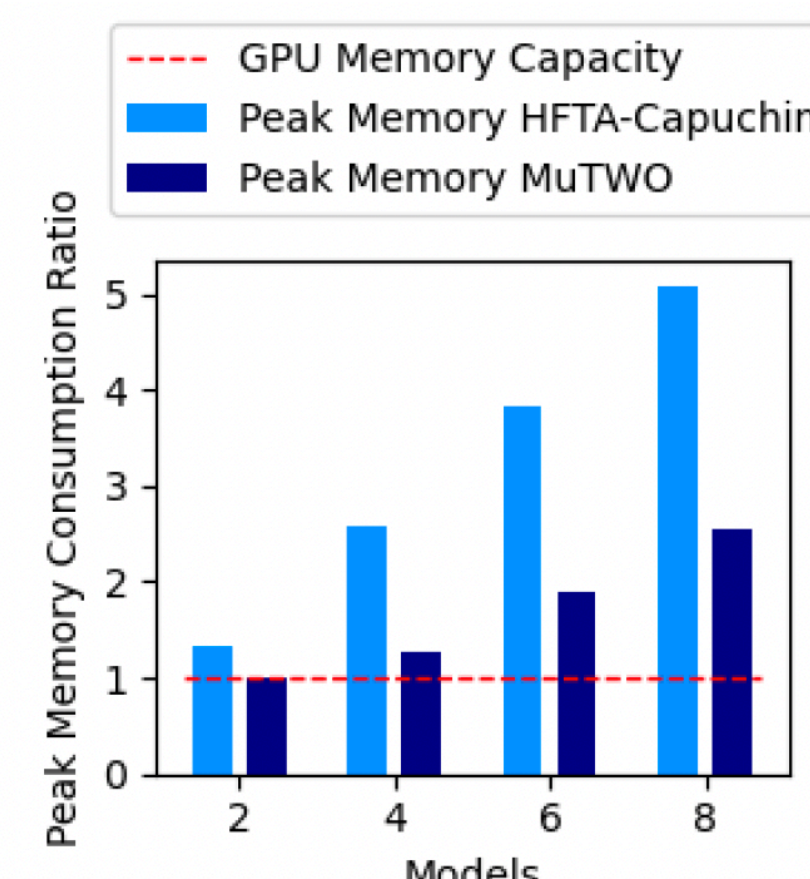
(a) Bert Latency Breakdown



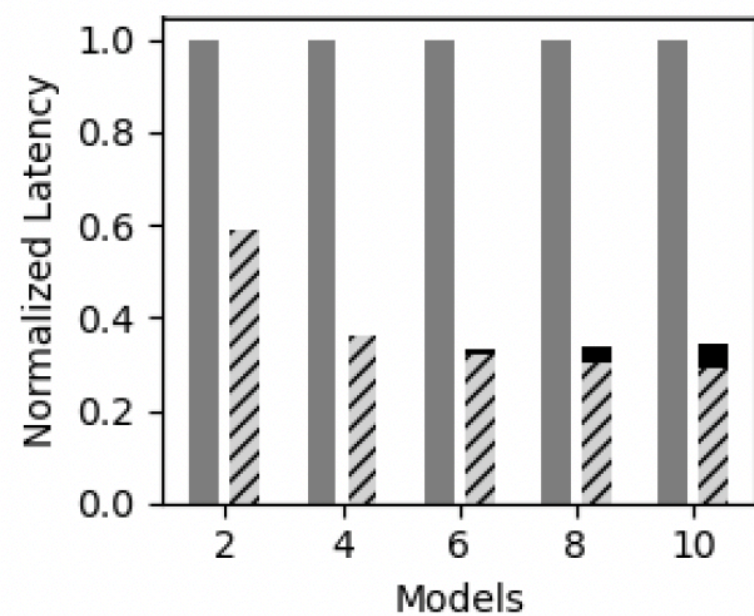
(b) Bert Recomputation Ratio



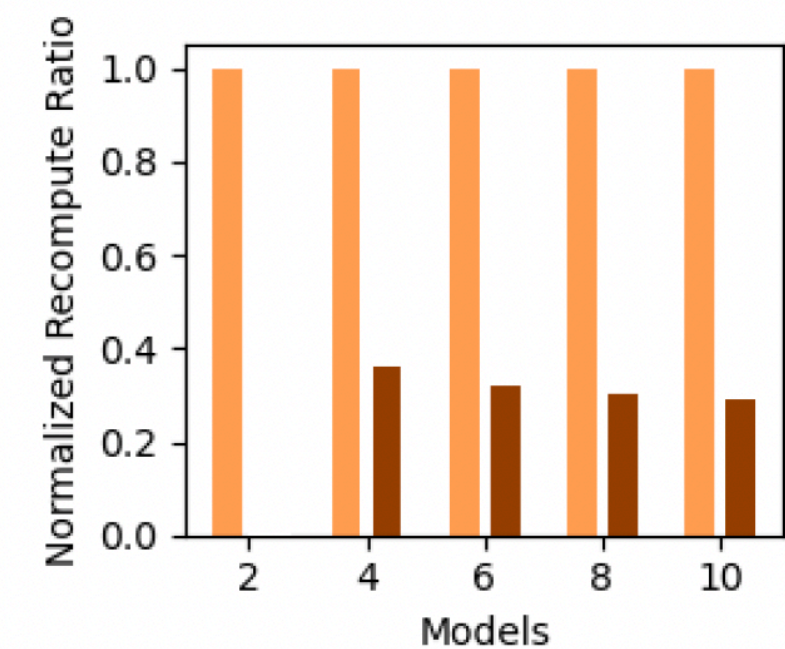
(c) Bert Swap Overlap Ratio



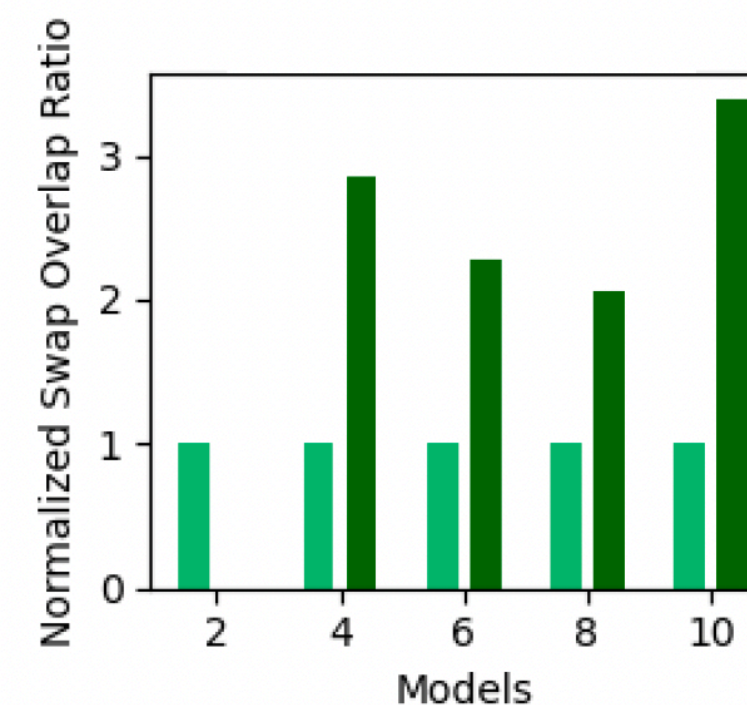
(d) Bert Peak Mem. Cons. Ratio



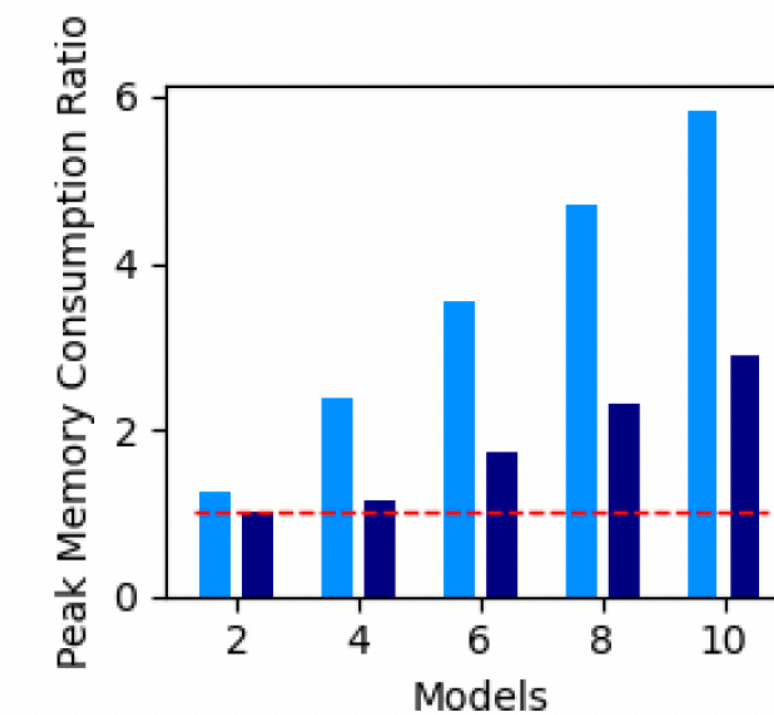
(e) ViT Latency Breakdown



(f) ViT Recomputation Ratio



(g) ViT Swap Overlap Ratio

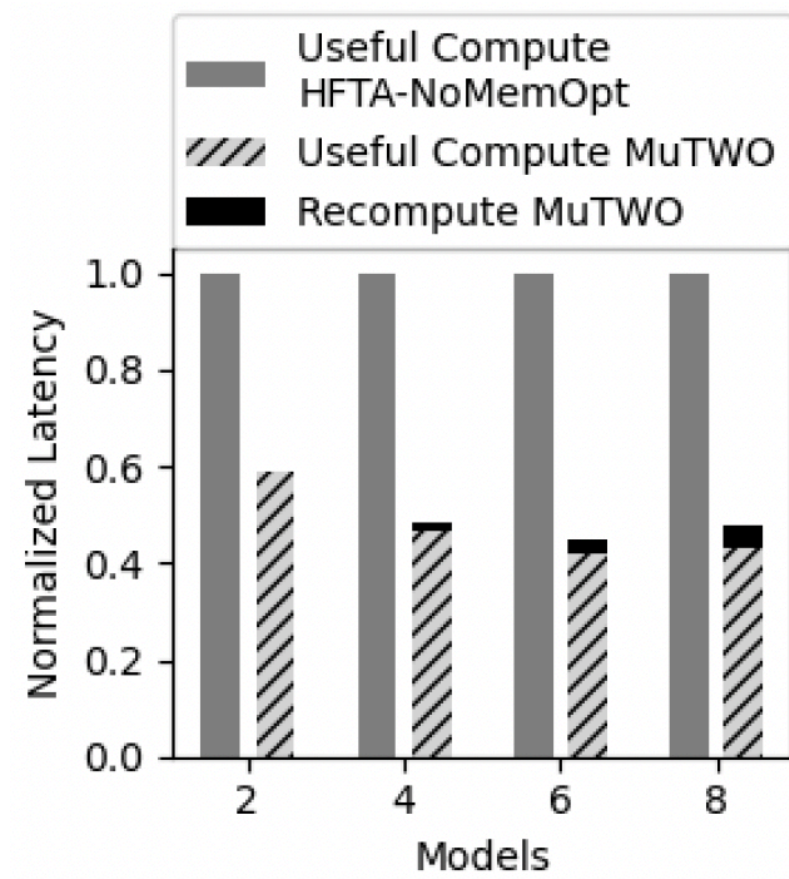


(h) ViT Peak Mem. Cons. Ratio

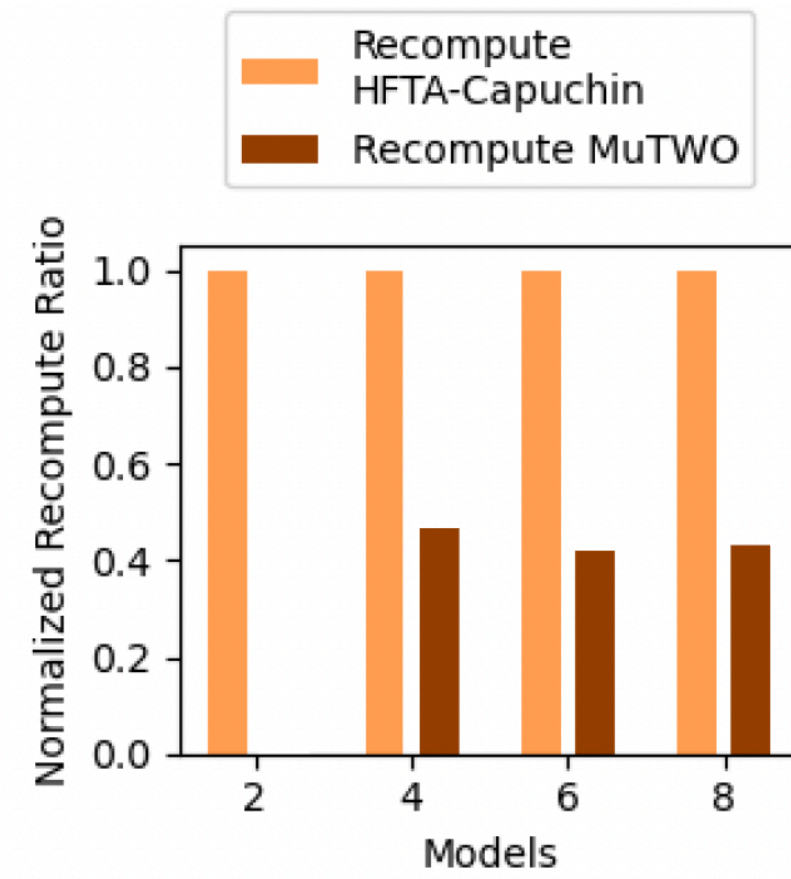
Less than 50% Recomputation



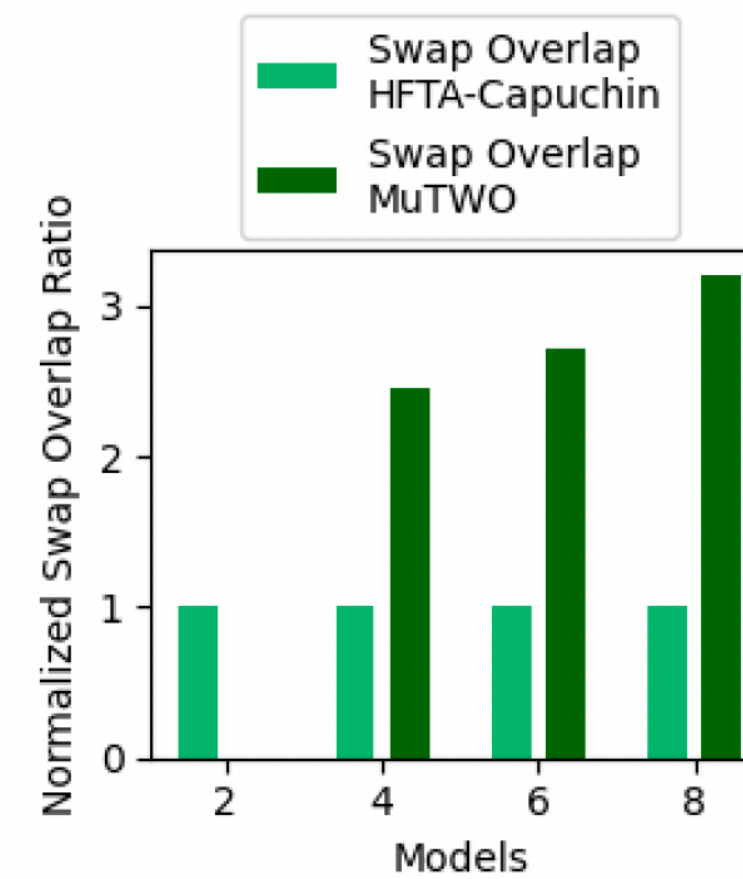
Performance Breakdown



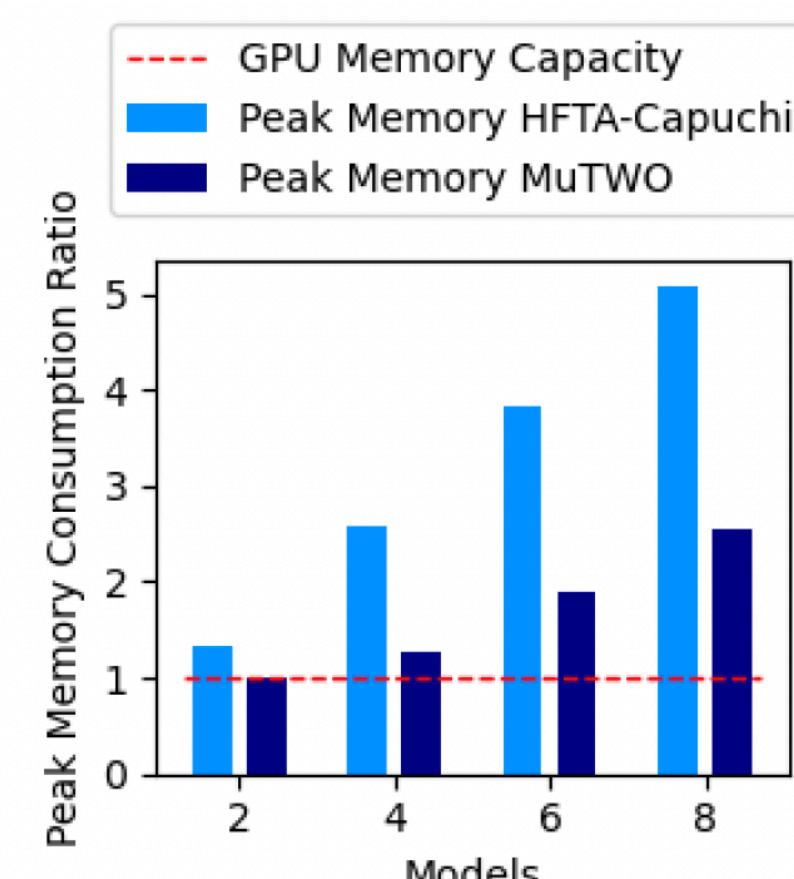
(a) Bert Latency Breakdown



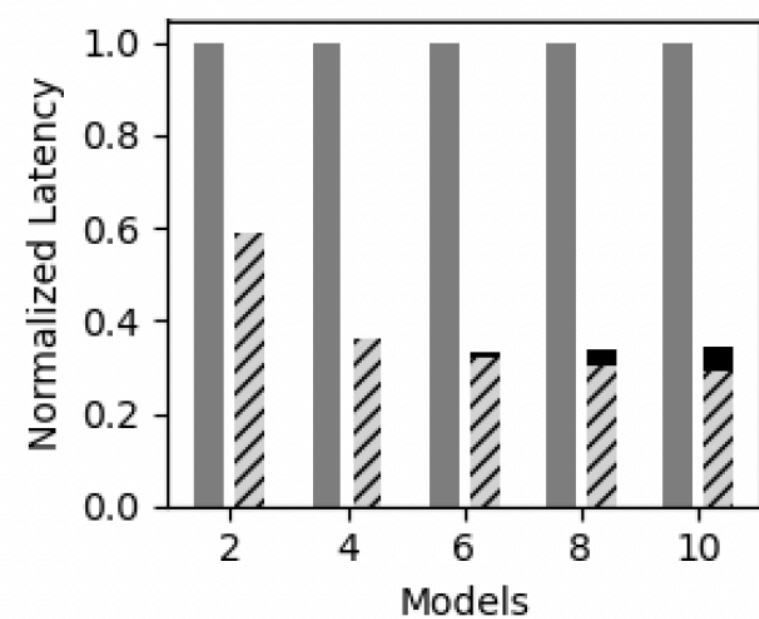
(b) Bert Recomputation Ratio



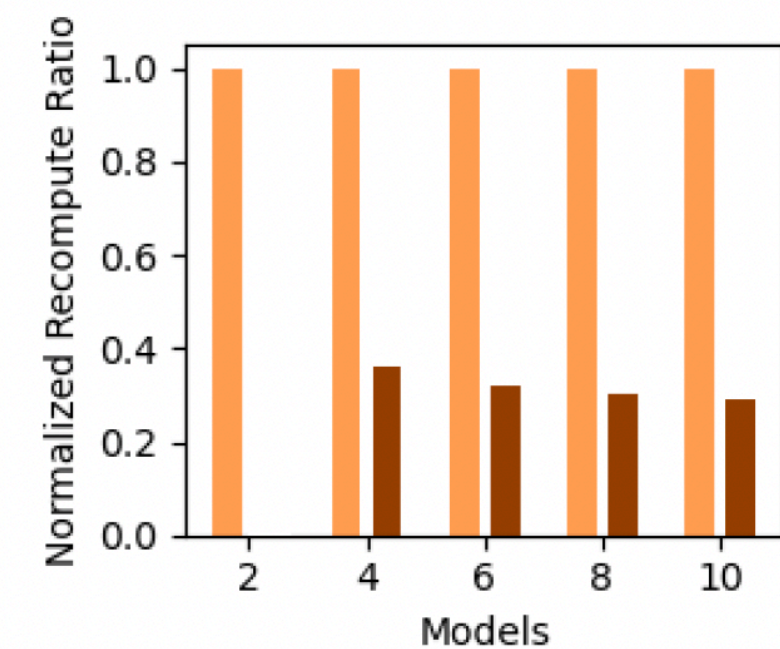
(c) Bert Swap Overlap Ratio



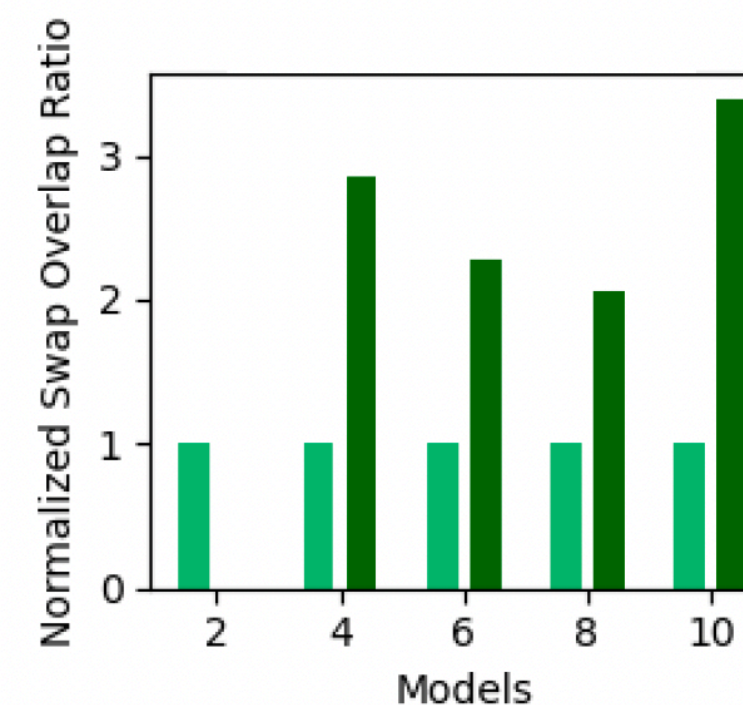
(d) Bert Peak Mem. Cons. Ratio



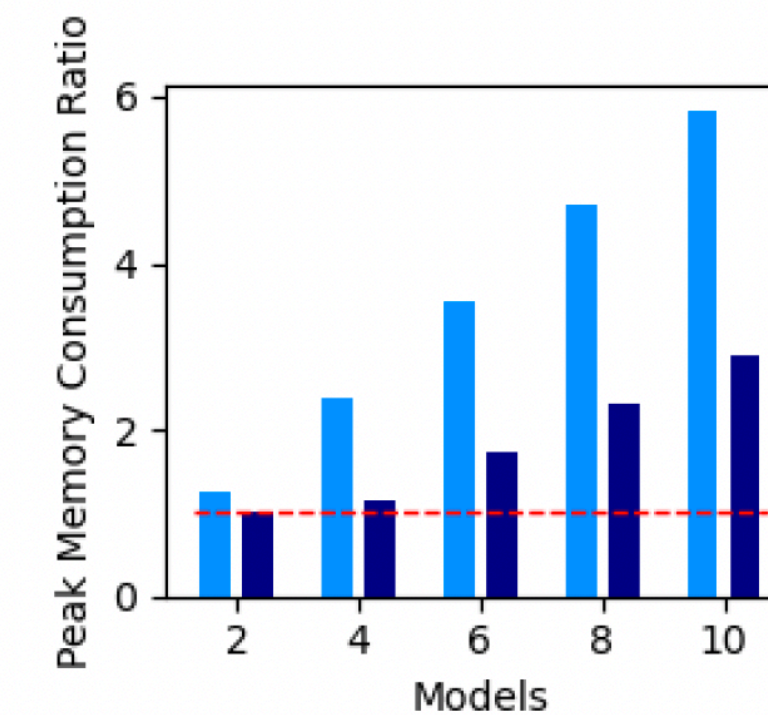
(e) ViT Latency Breakdown



(f) ViT Recomputation Ratio



(g) ViT Swap Overlap Ratio

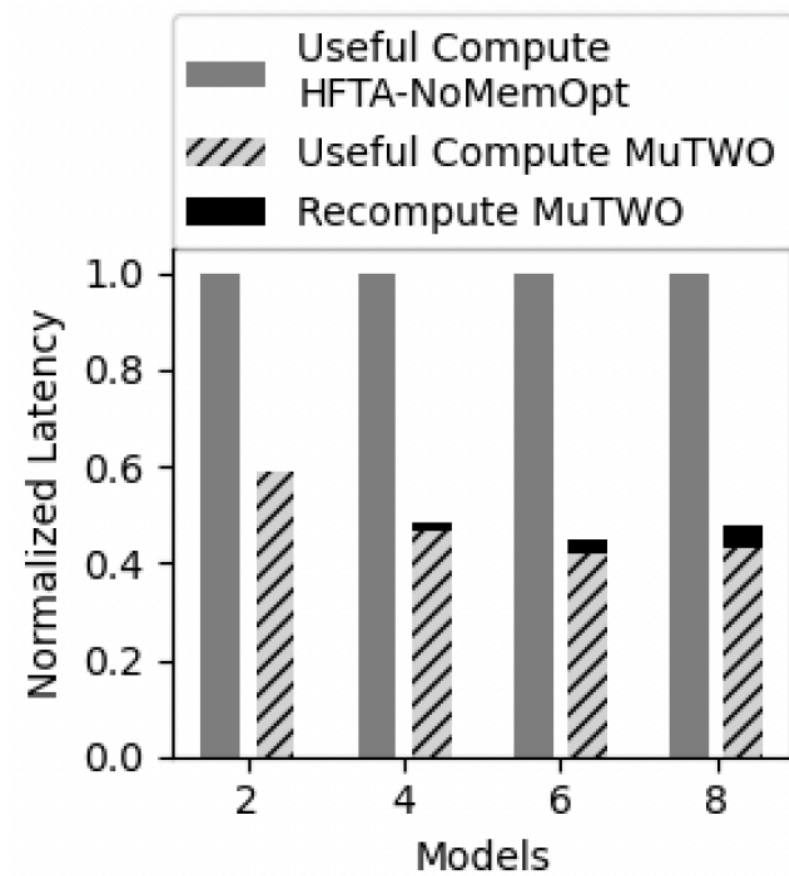


(h) ViT Peak Mem. Cons. Ratio

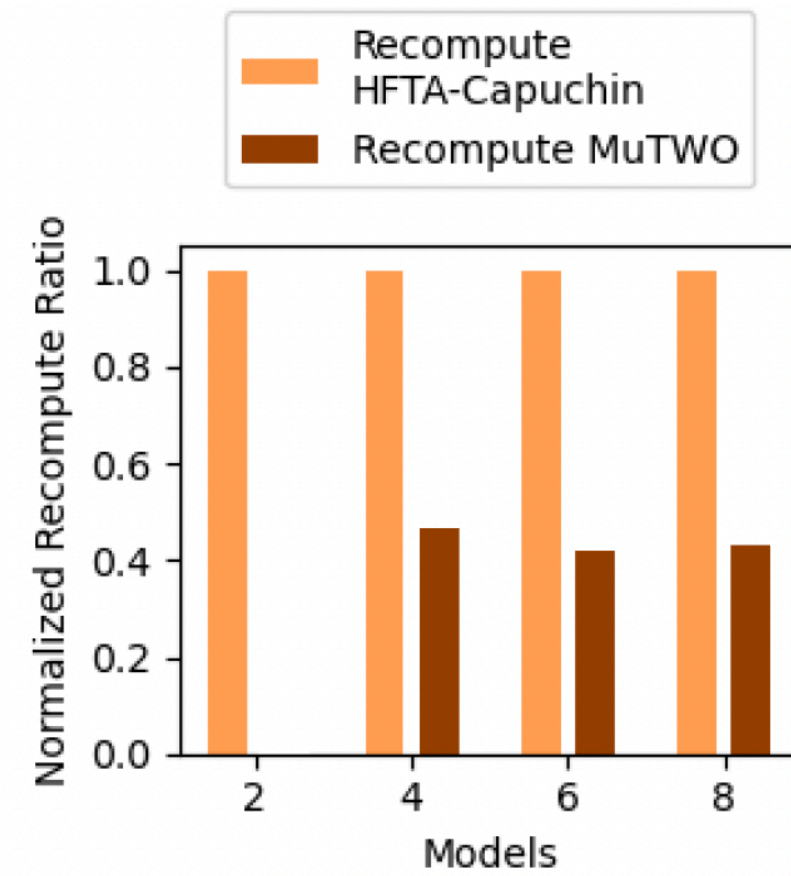
3x more Swap Overlap



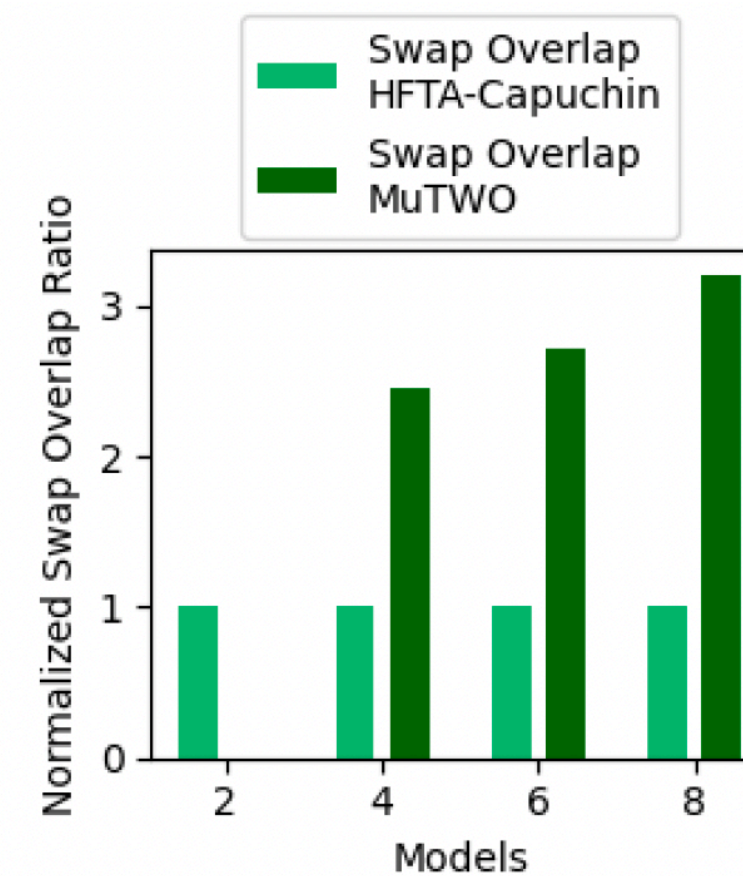
Performance Breakdown



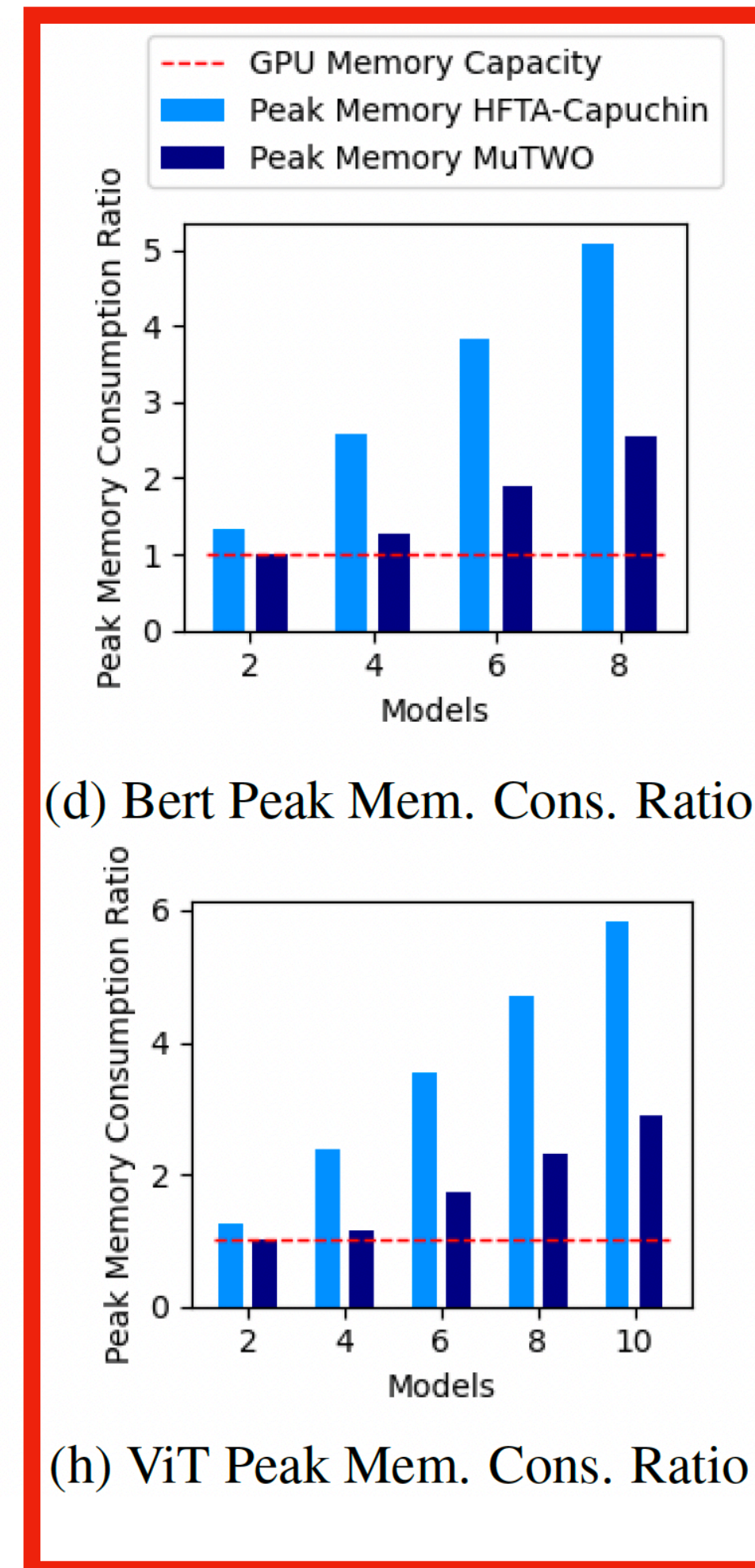
(a) Bert Latency Breakdown



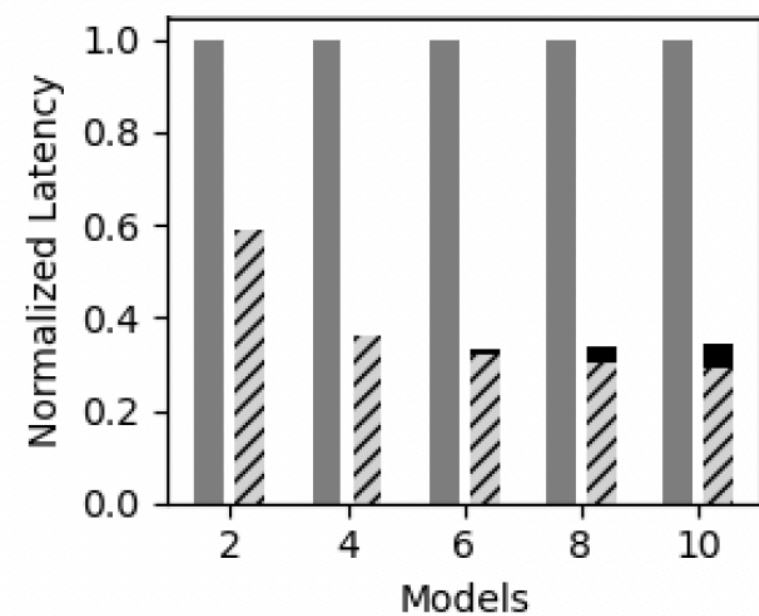
(b) Bert Recomputation Ratio



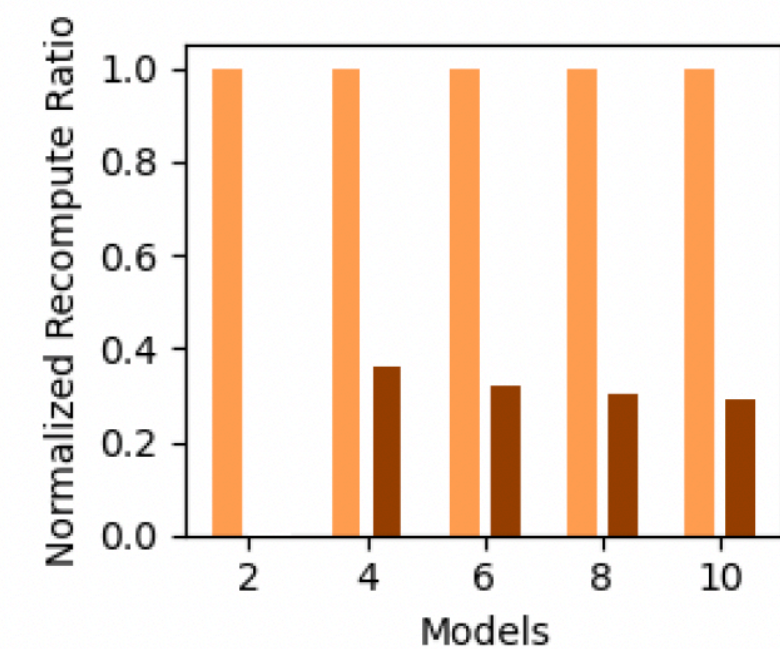
(c) Bert Swap Overlap Ratio



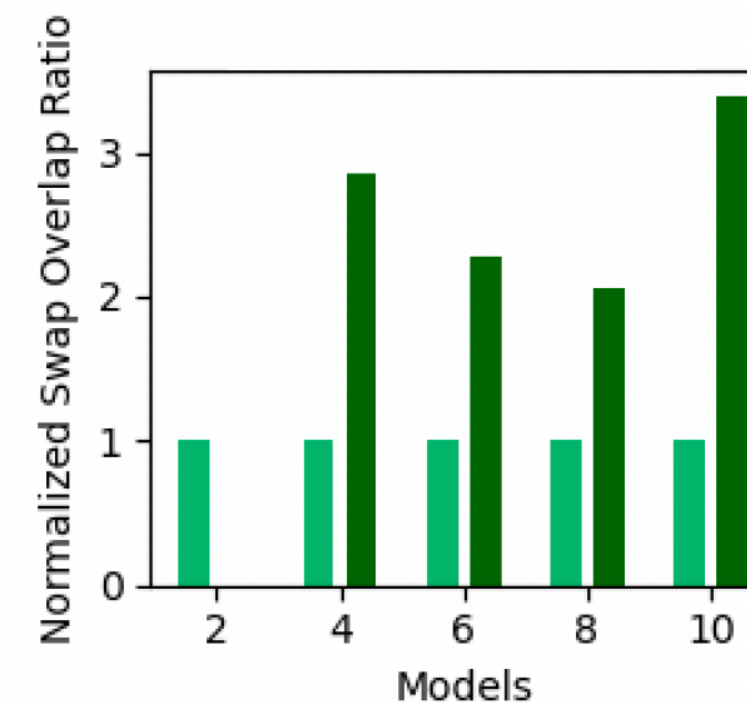
(d) Bert Peak Mem. Cons. Ratio



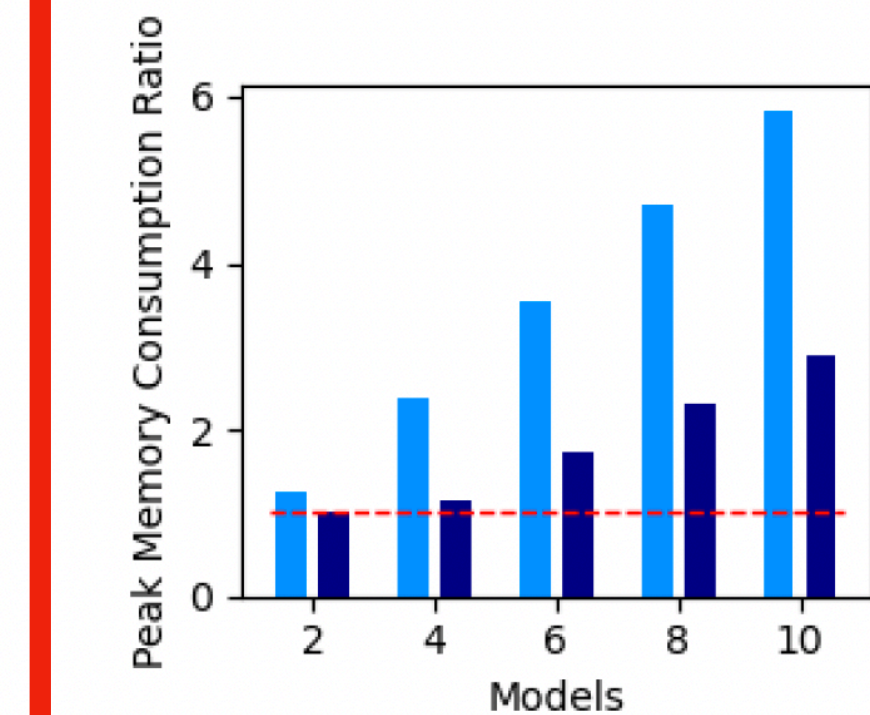
(e) ViT Latency Breakdown



(f) ViT Recomputation Ratio



(g) ViT Swap Overlap Ratio



(h) ViT Peak Mem. Cons. Ratio



Less than 50% peak memory consumption

Implementation Details



- Profiling: PyTorch Profiler
- Parallel Compute and Data Operations: CUDA Streams and Events
- Operator Fusion: PyTorch VMap
- Computational Model Graphs: PyTorch AOT Autograd with FakeTensors
- Runtime Overhead Reduction: CUDA Graphs



State of the Art vs μ -TWO

Parameter	HFTA (MLSys'21)	Checkmate (MLSys'20)	Capuchin (ASPLOS'20)	μ -TWO
High Compute Utilization	Yes	No	No	Yes
High Memory Utilization	No	Yes	Yes	Yes
Large Number of Models	Yes	No	No	Yes
Large Model Size	No	Yes	Yes	Yes
Large Mini-batch Size	No	Yes	Yes	Yes
Stalls	NA	NA	Low	Low
Compute Overhead	NA	High	Moderate	Low



Thank You

Choose whether an intermediate tensor should be

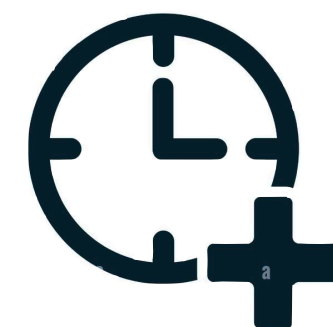
Recomputed

Swapped

Combinatorial in nature

NP-HARD

Jain Paras et al., Checkmate: Breaking the memory wall with optimal tensor rematerialization, *MLSys 2020*

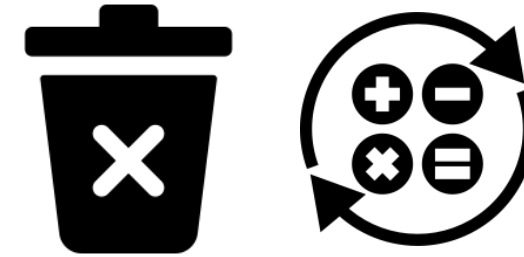


ILP Solvers
Approximation algorithms

GREEDY



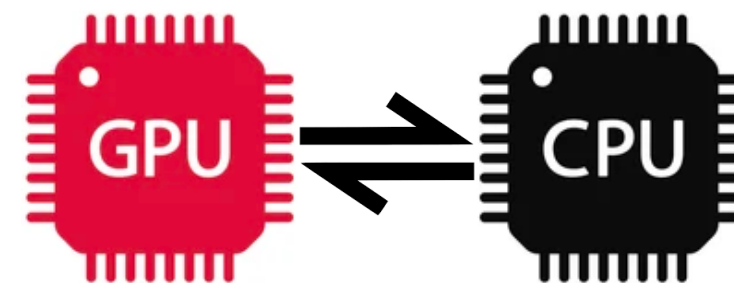
Recomputation



$$\text{recompute_ratio} = \frac{\text{memory_size}}{\text{recomp_time}}$$

Memory Savings Per Second

Swapping

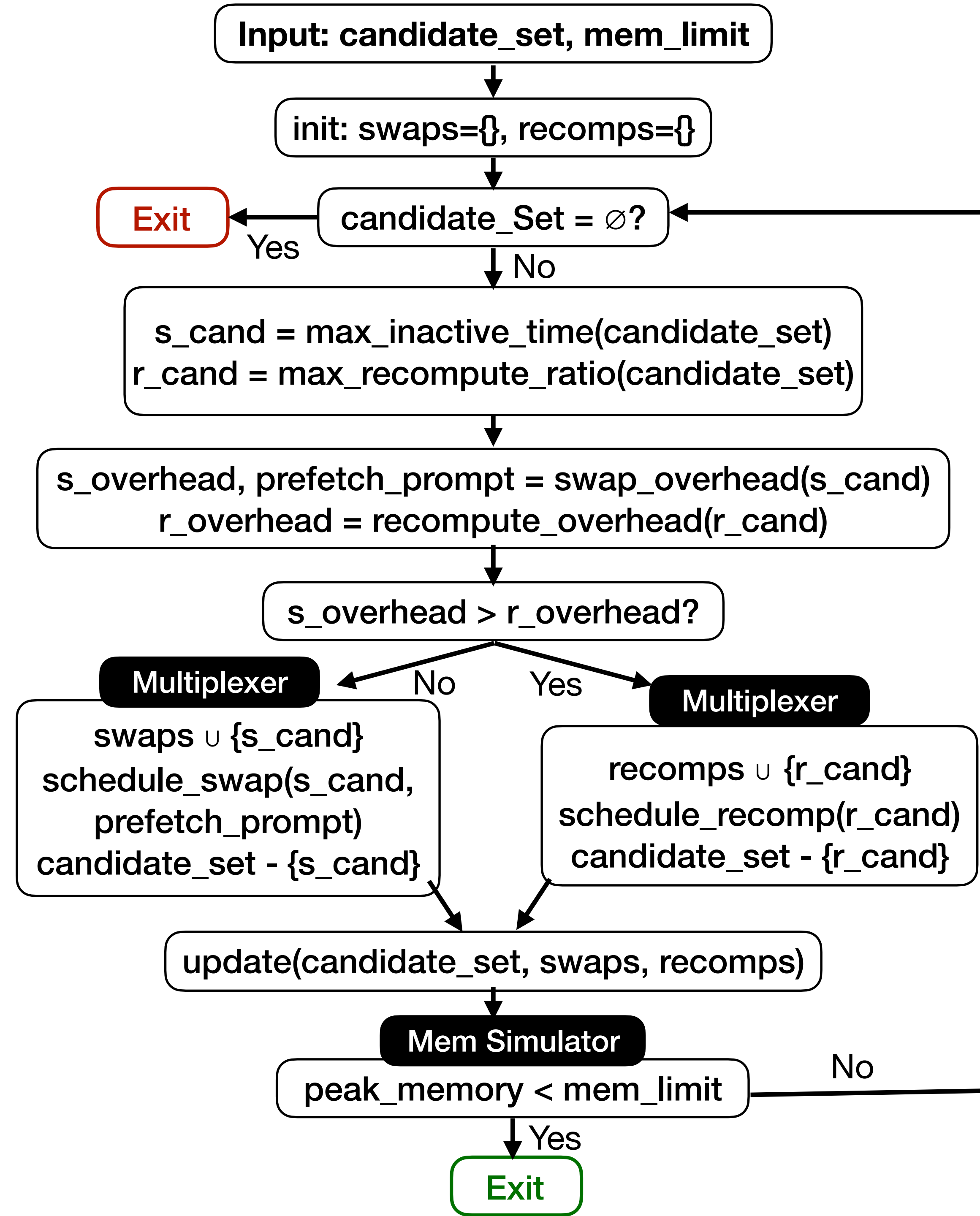


Inactive Time: Last Use (Forward Pass) -> First Use (Backward Pass)

Better opportunity to hide swapping latency



Scheduling Policy



Swap Overhead Calculation

