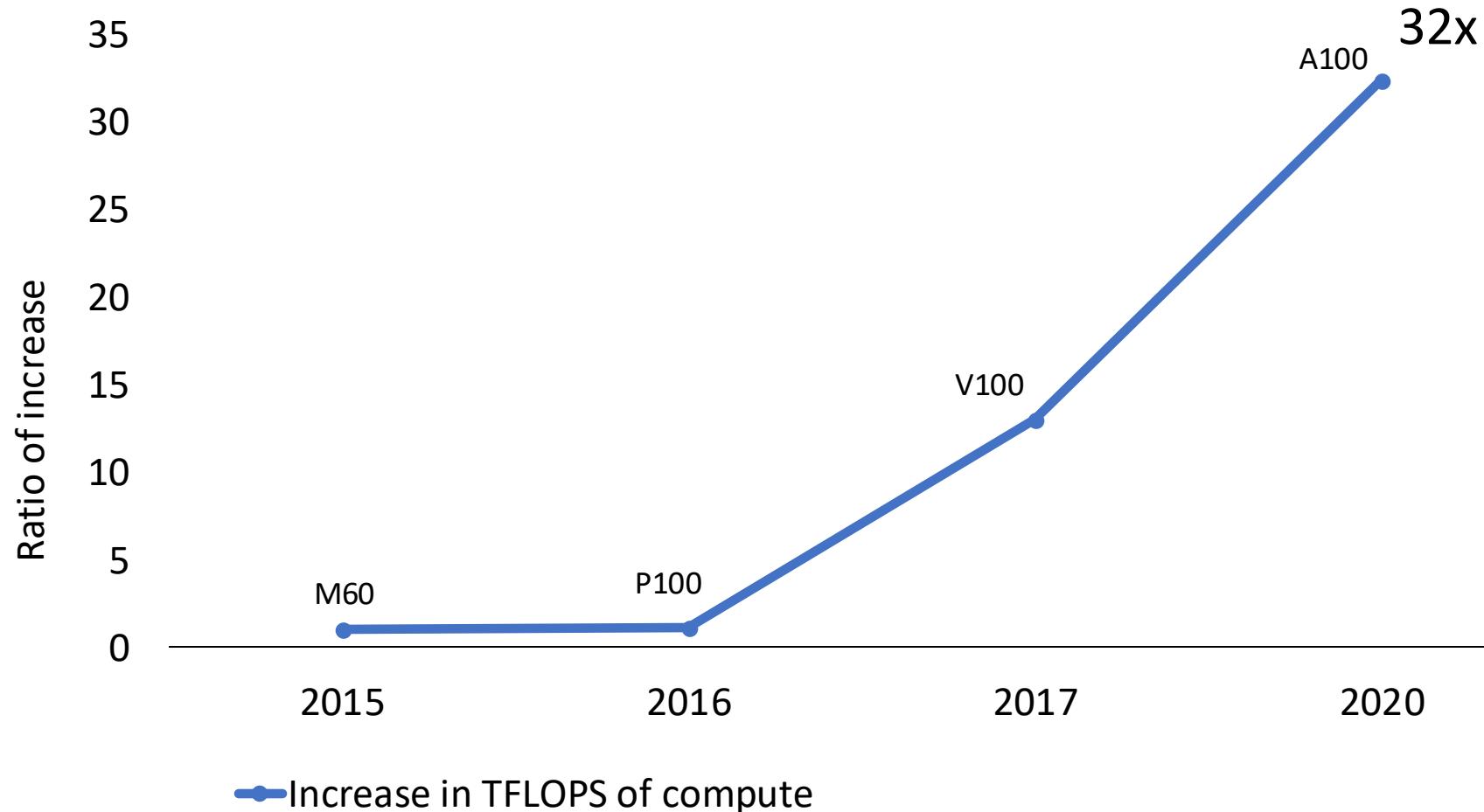


Memory Optimization for Deep Networks

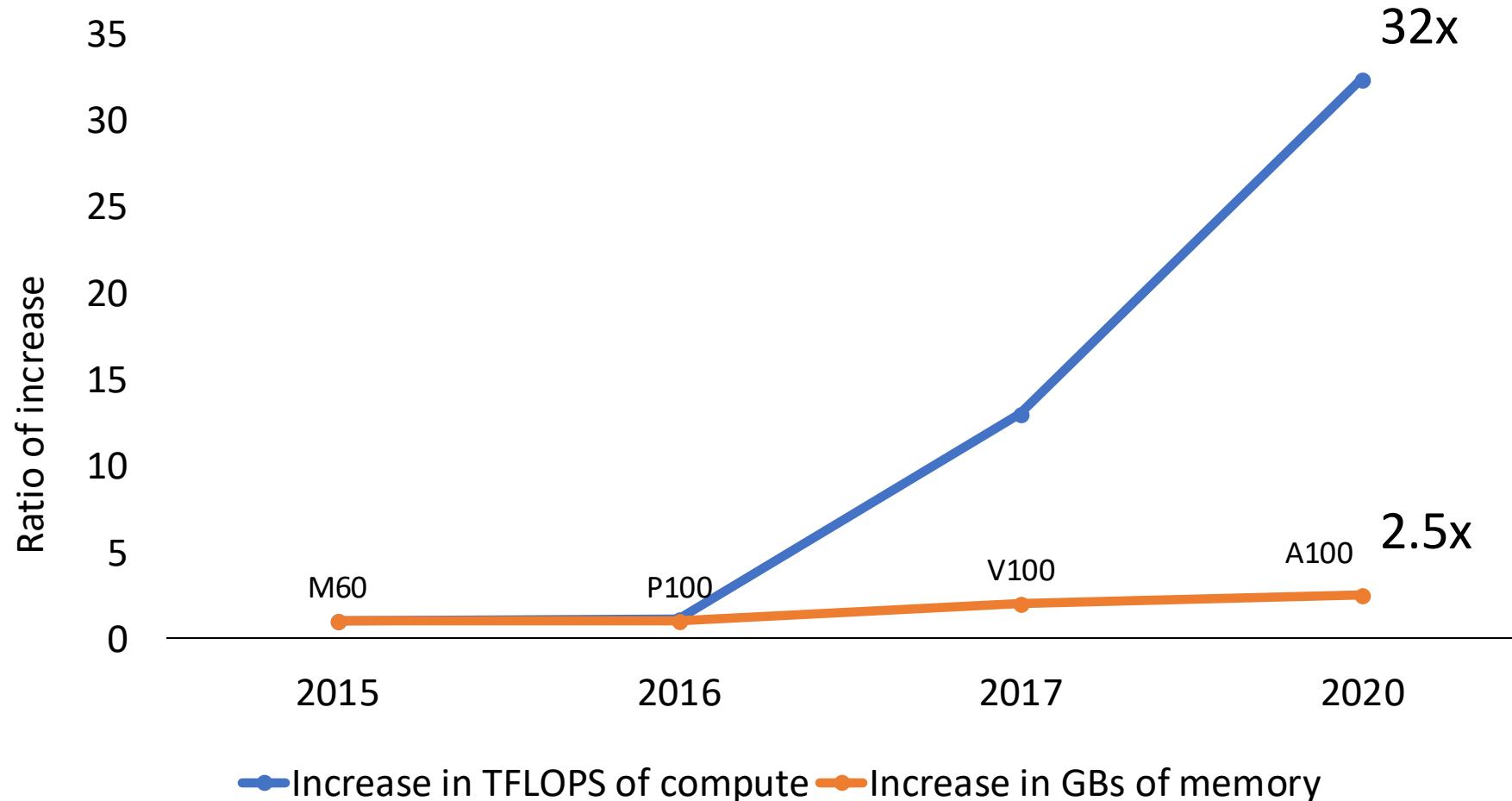
Aashaka Shah¹, Chao-Yuan Wu¹, Jayashree Mohan¹,
Vijay Chidambaram^{1,2}, Philipp Krähenbühl¹



Memory and Compute Trends in GPUs



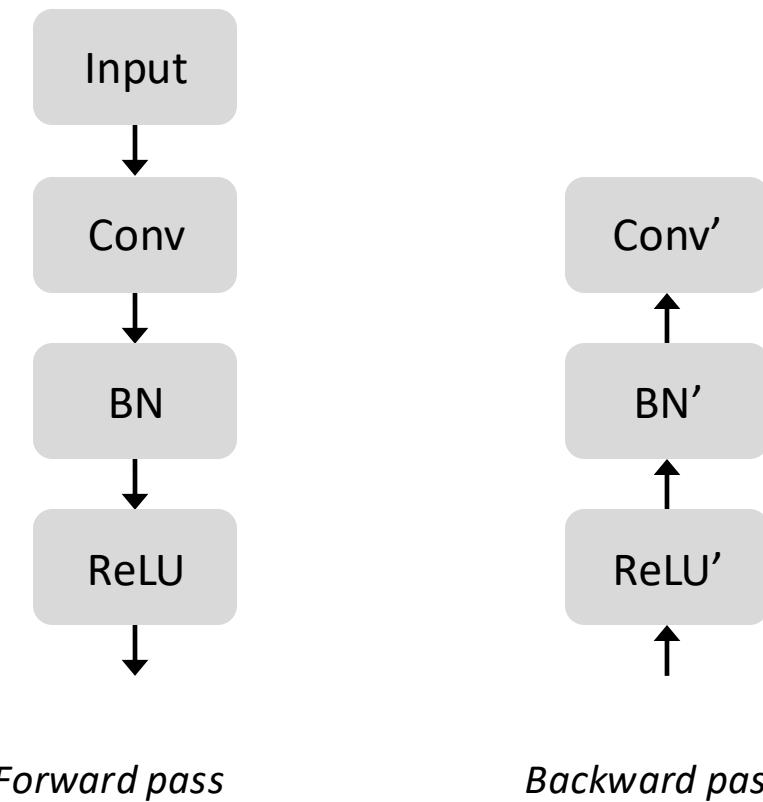
Memory and Compute Trends in GPUs



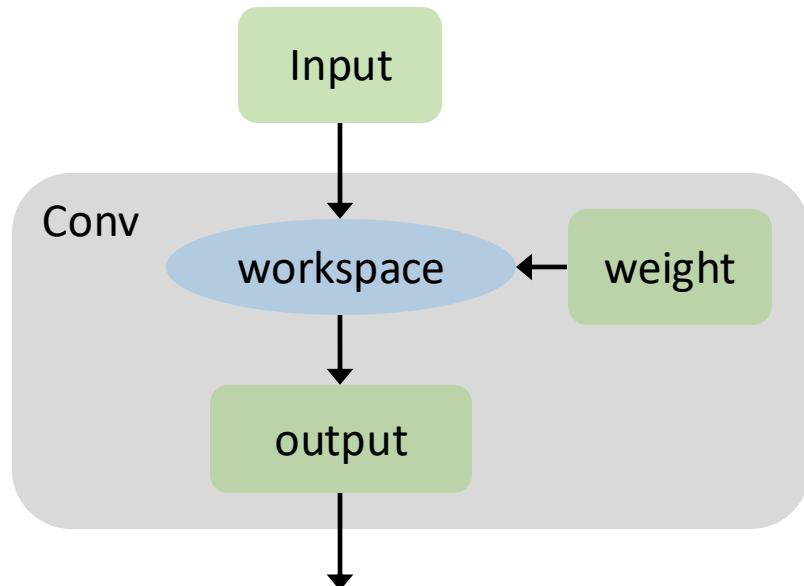
(pytorch_env) cc@gpu:~\$

]

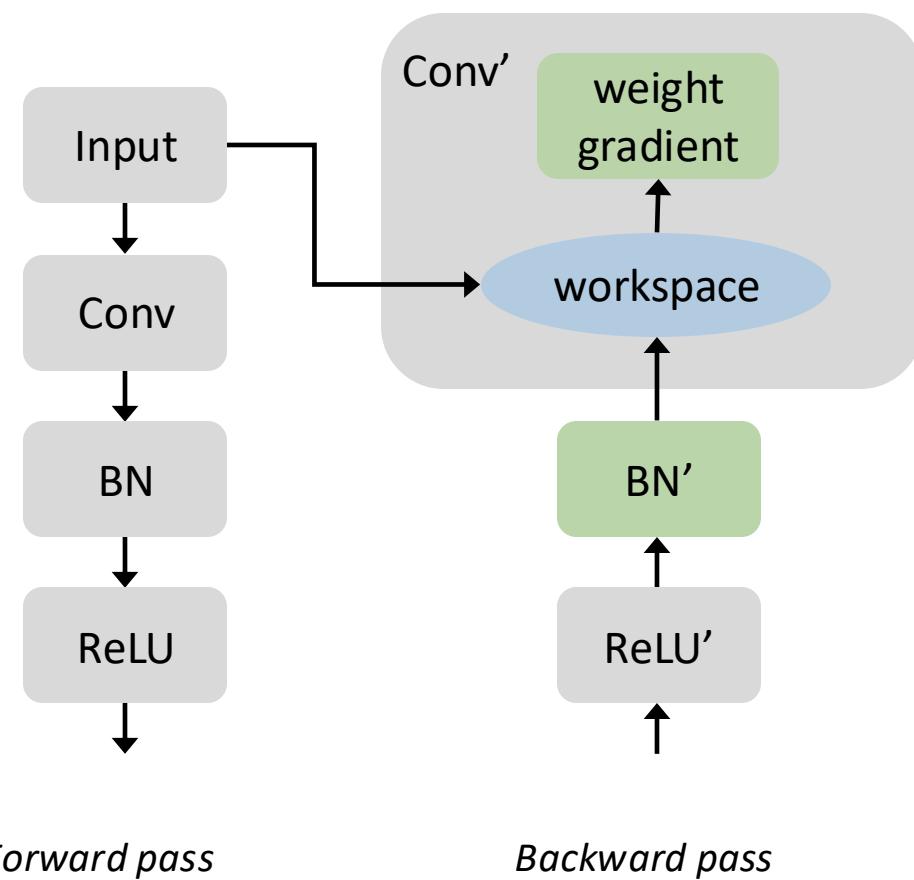
Memory usage in deep networks



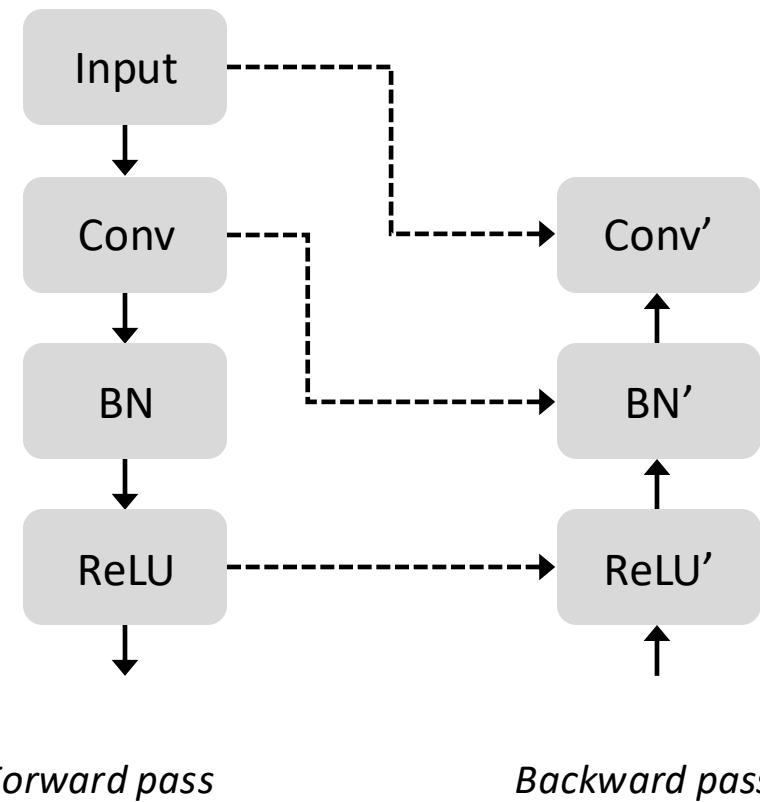
Memory usage in deep networks



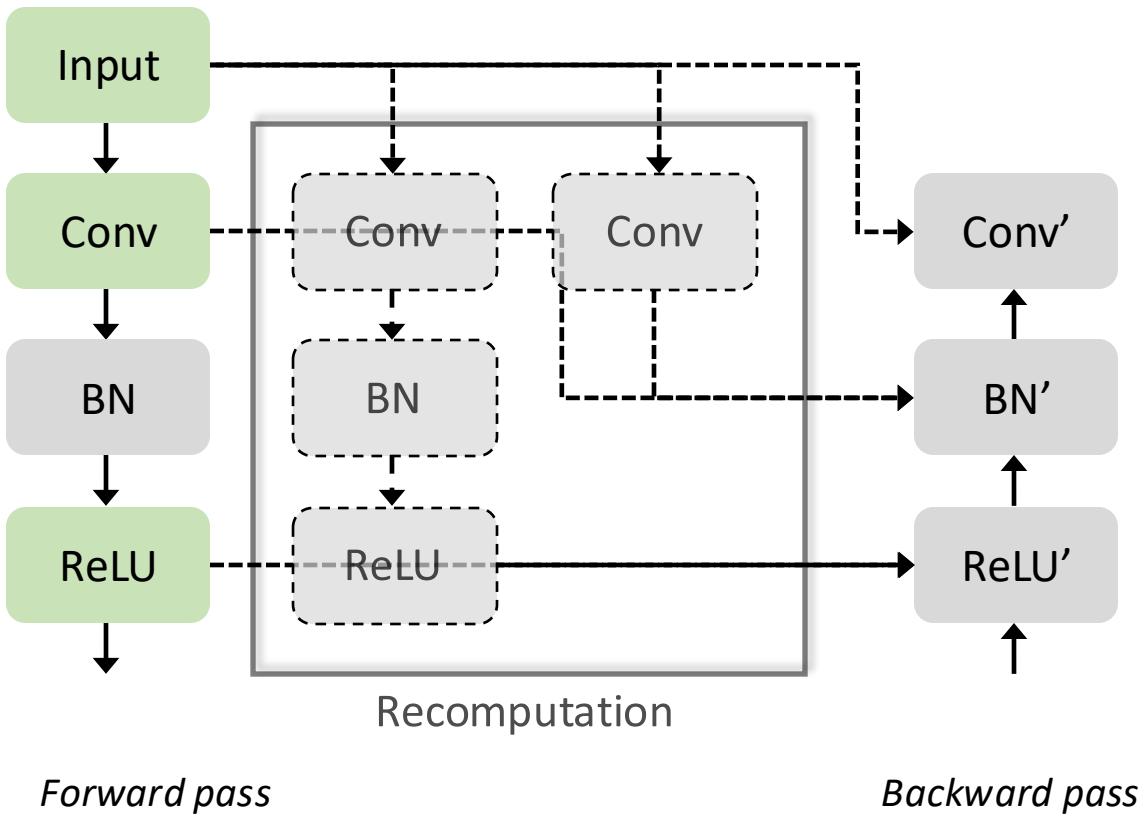
Memory usage in deep networks



Memory usage in deep networks



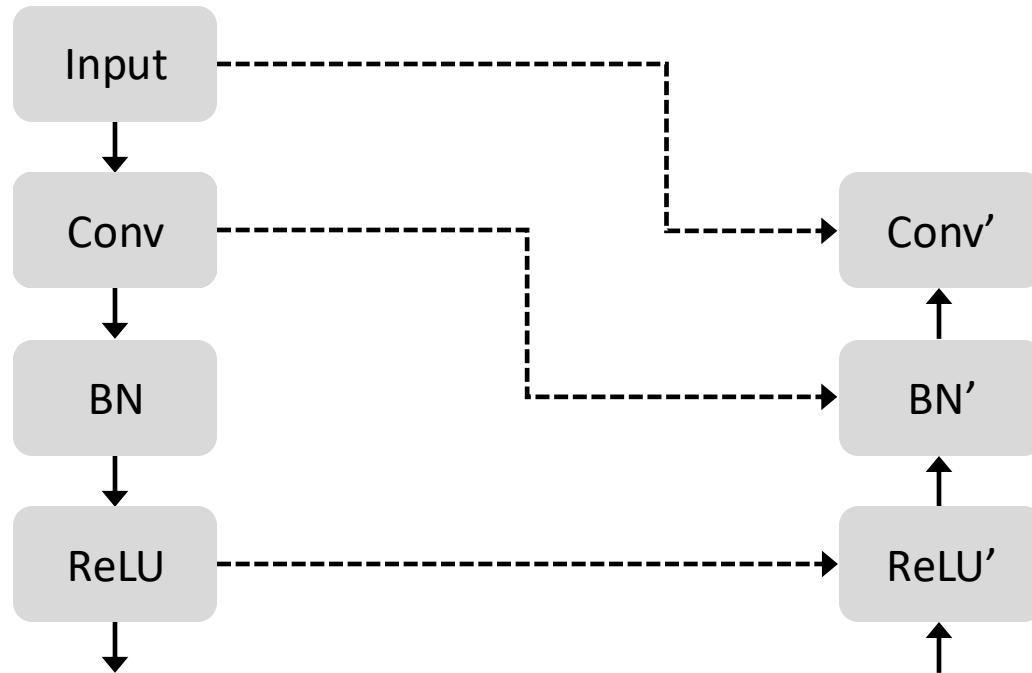
Checkpointing



- [1] Jain, Paras, et al. "Checkmate: Breaking the Memory Wall with Optimal Tensor Rematerialization." *MLSys*. 2020.
- [2] Chen, Tianqi, et al. "Training deep nets with sublinear memory cost." *arXiv preprint arXiv:1604.06174* (2016).

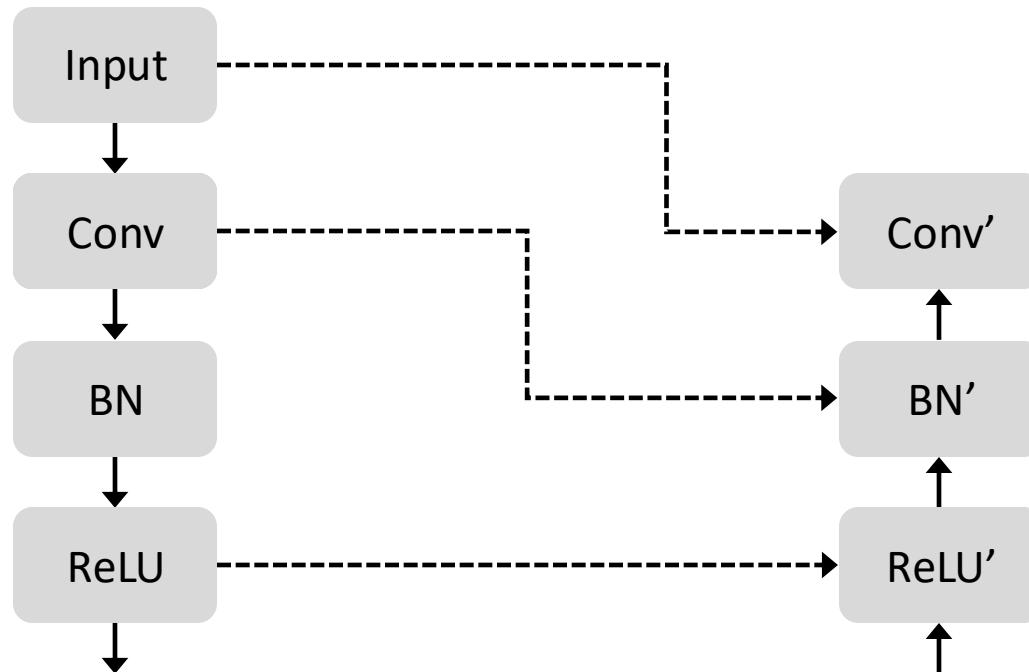
Operator optimizations

Convolution algorithm selection



Operator optimizations

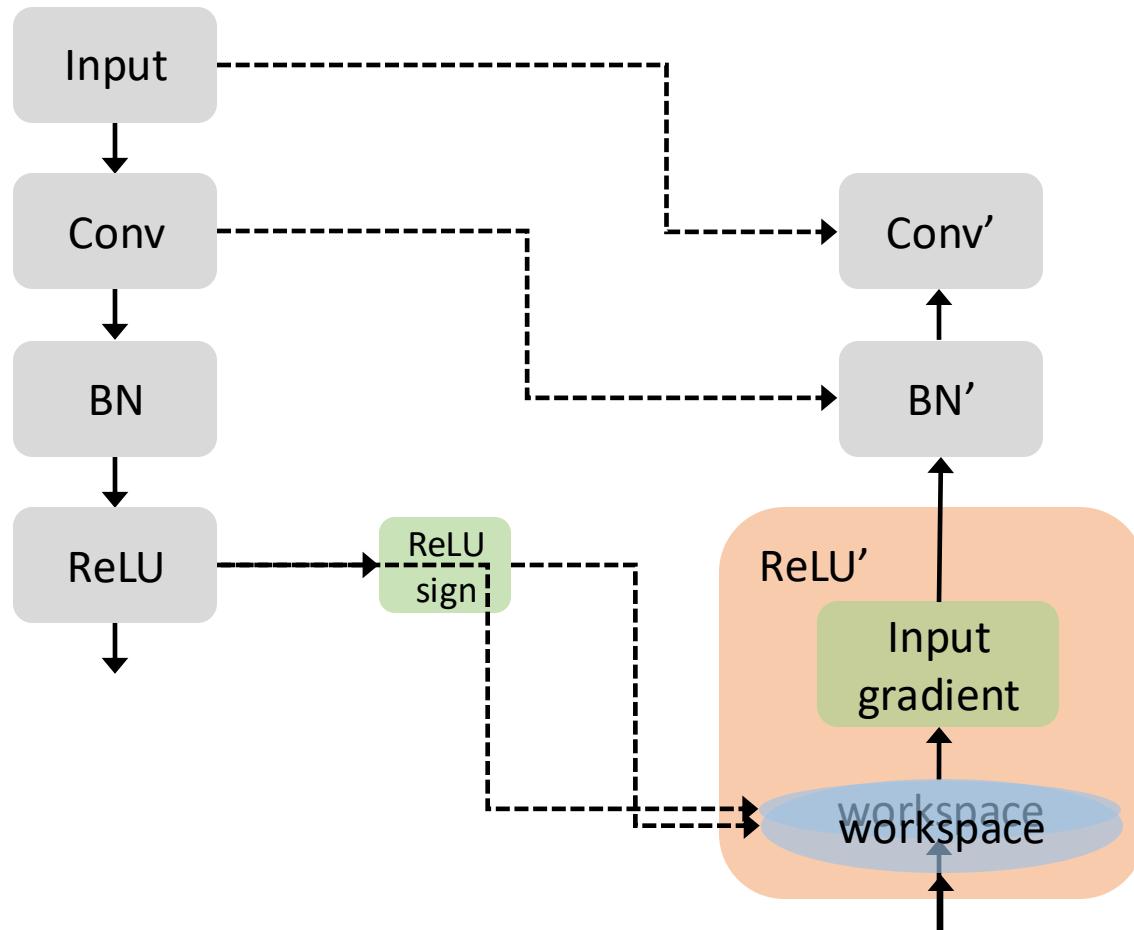
Intermediate-activated ReLU



[1] Jain, Animesh, et al. "Gist: Efficient data encoding for deep neural network training." *ISCA 2018*.

Operator optimizations

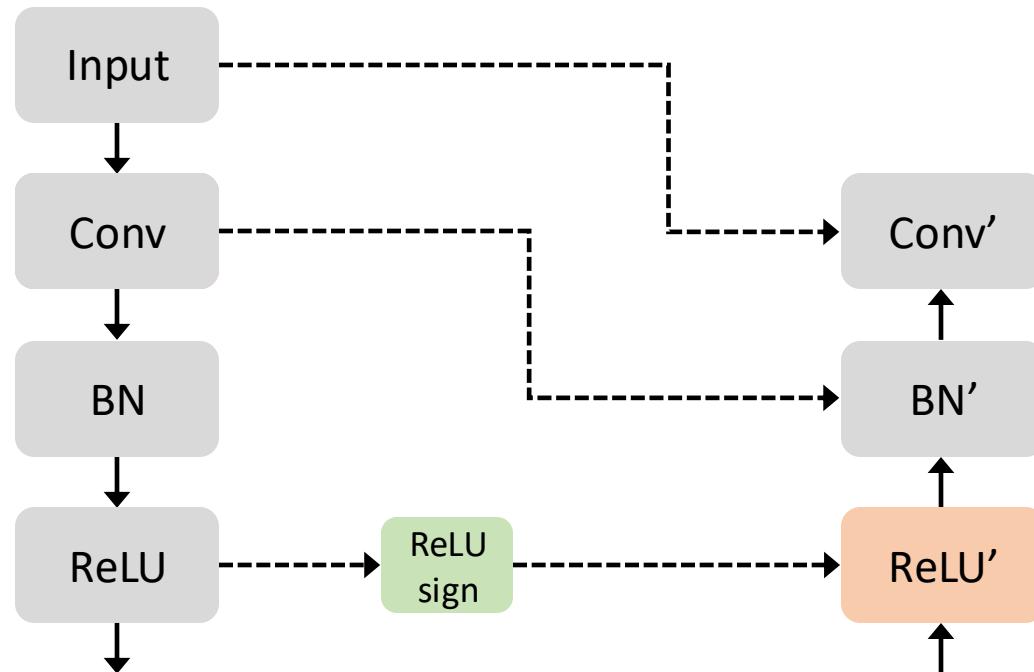
Intermediate-activated ReLU



[1] Jain, Animesh, et al. "Gist: Efficient data encoding for deep neural network training." *ISCA 2018*.

Operator optimizations

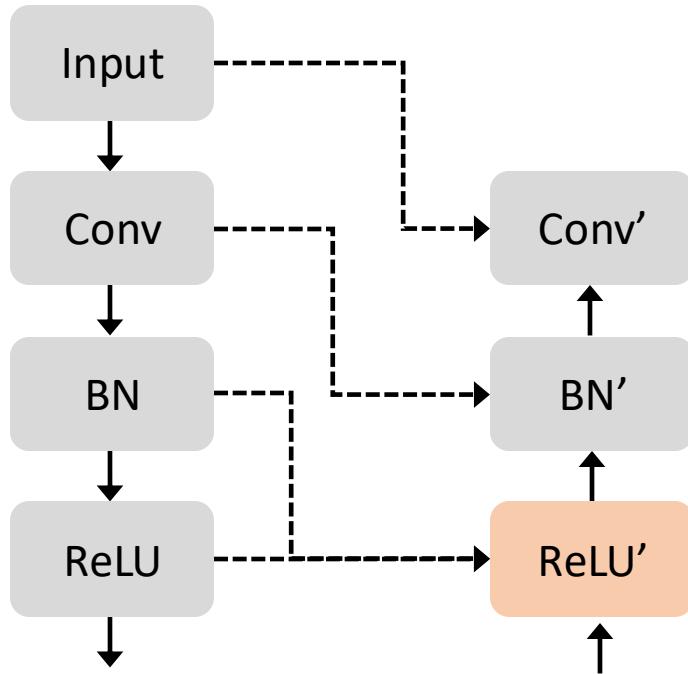
Intermediate-activated ReLU



[1] Jain, Animesh, et al. "Gist: Efficient data encoding for deep neural network training." *ISCA 2018*.

Operator optimizations

Changing backward dependencies

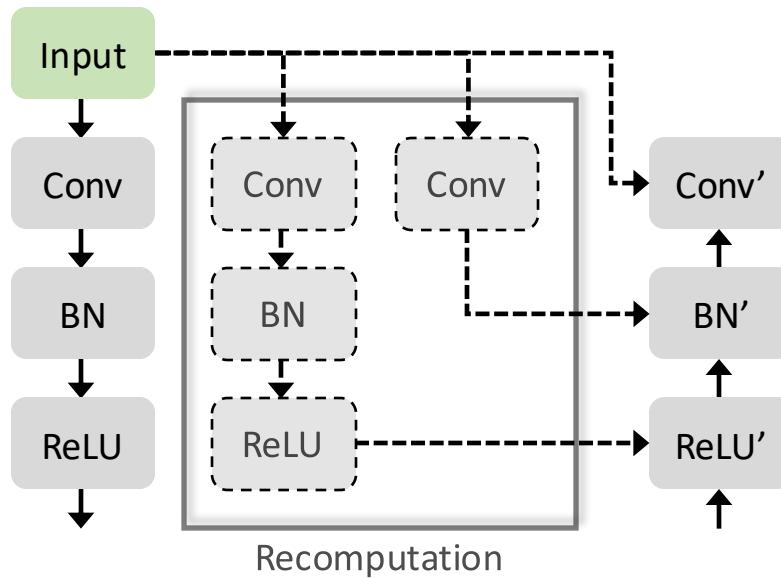


[1] Bulo, Samuel Rota, Lorenzo Porzi, and Peter Kortscheder. "In-place activated batchnorm for memory-optimized training of dnns." *CVPR 2018*.

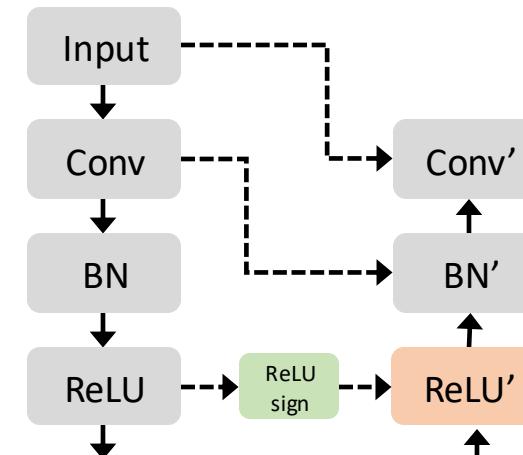
MONeT

Optimizes for minimal memory v/s compute tradeoff

Reduces memory usage by 3x while incurring 9–16% compute overhead over PyTorch for ResNet-50, GoogleNet, etc.

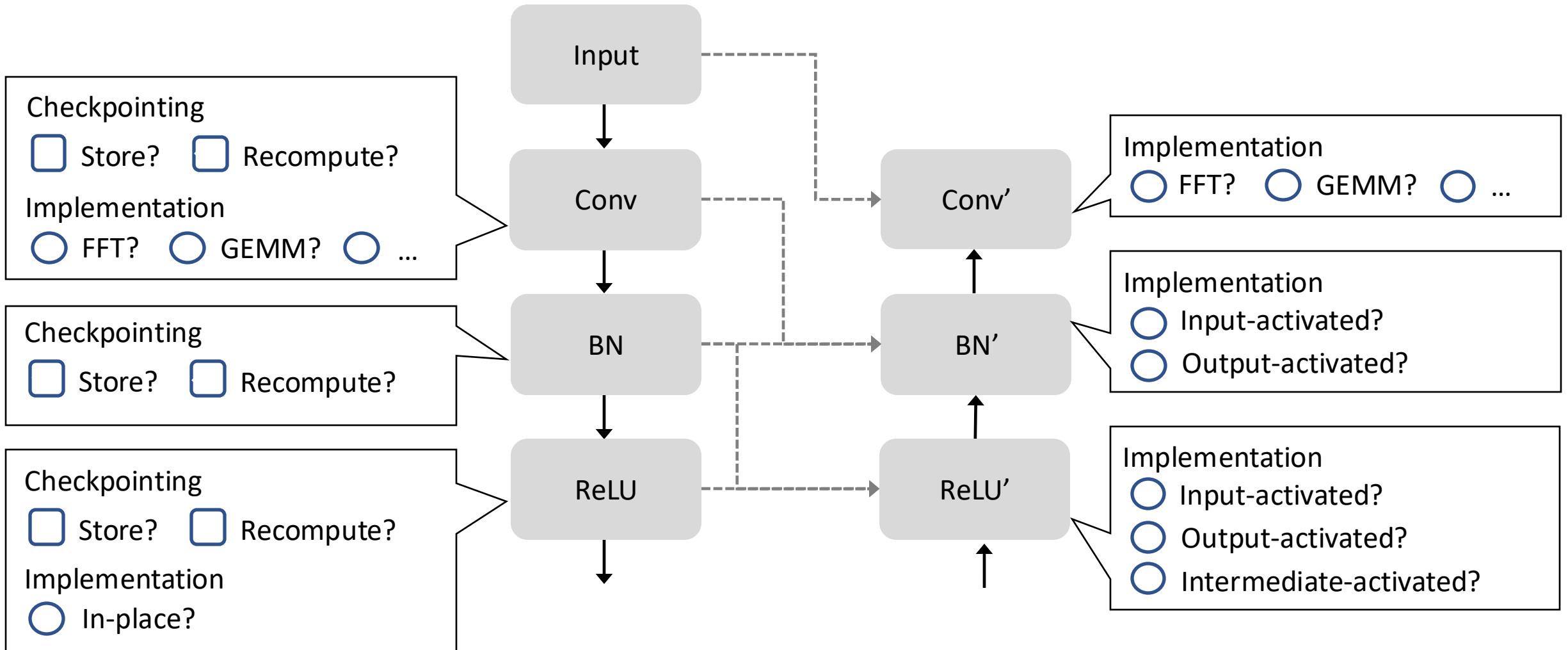


Checkpointing



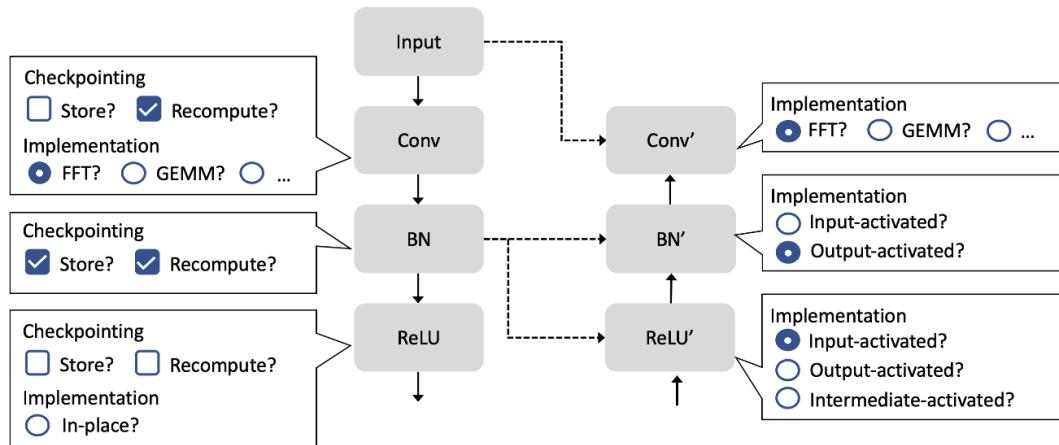
Operator optimizations

MONeT



MONeT

MONeT formulation



Optimization subproblem

Off-the-shelf solver

Boolean linear programming problem

Objective: Minimize computation time

Memory constraints

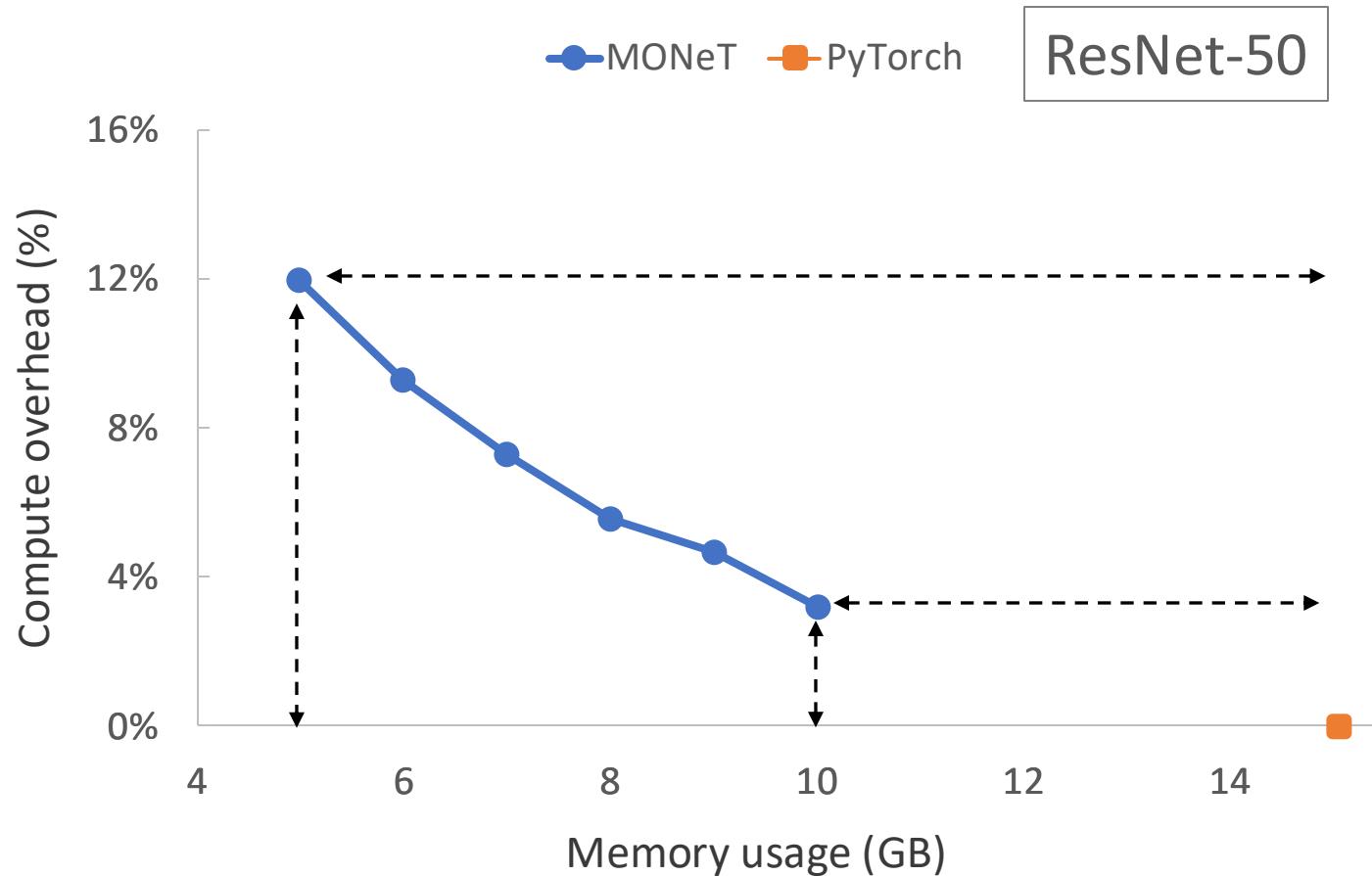
Peak backward memory ≤ Memory budget

Peak recomputation memory ≤ Memory budget

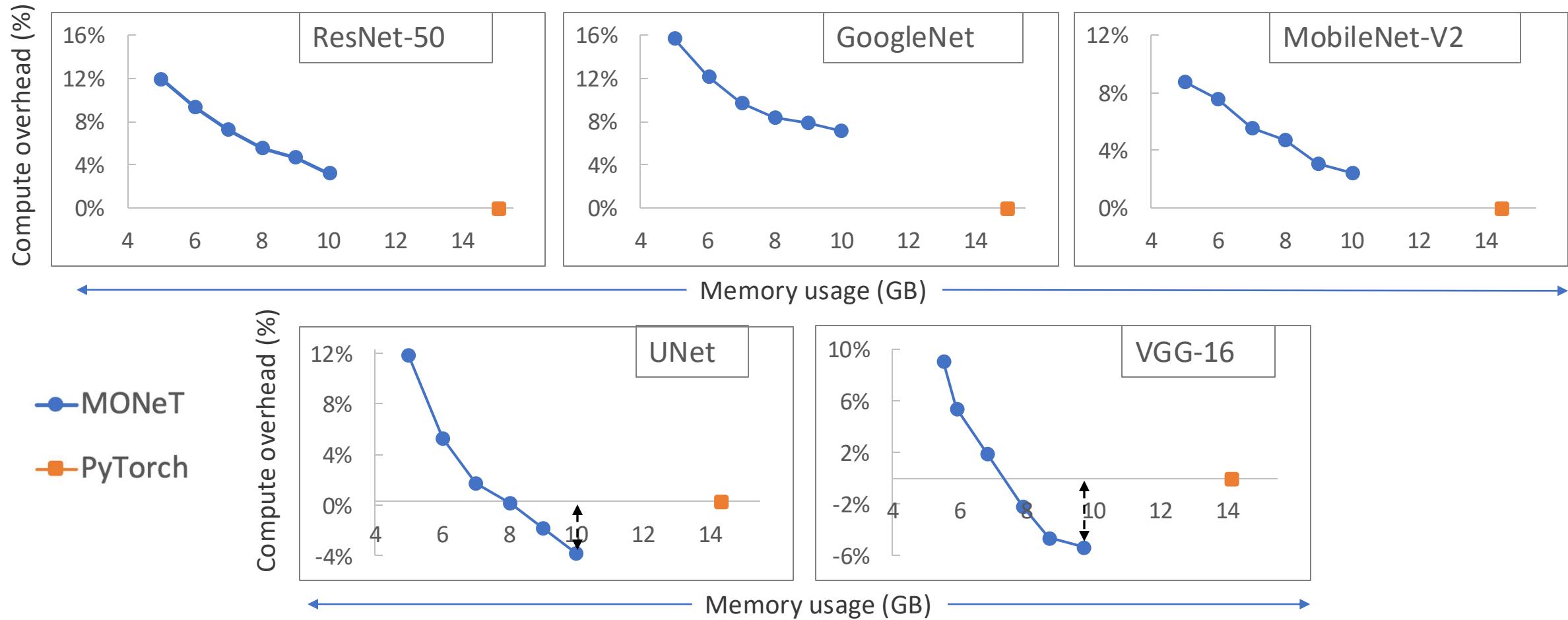
Checkpointing constraints

Only use operator output if present in memory

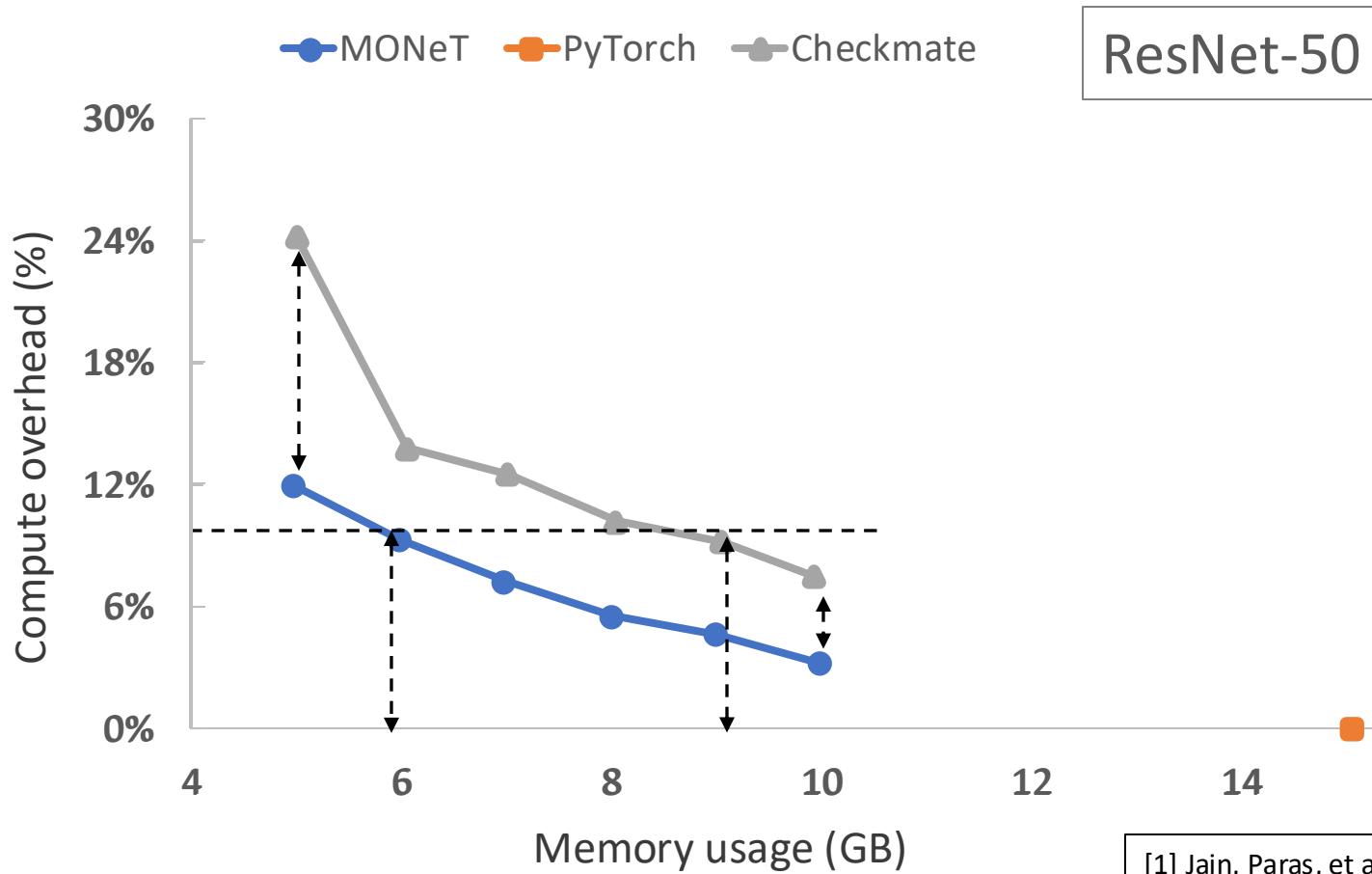
Evaluation



Evaluation

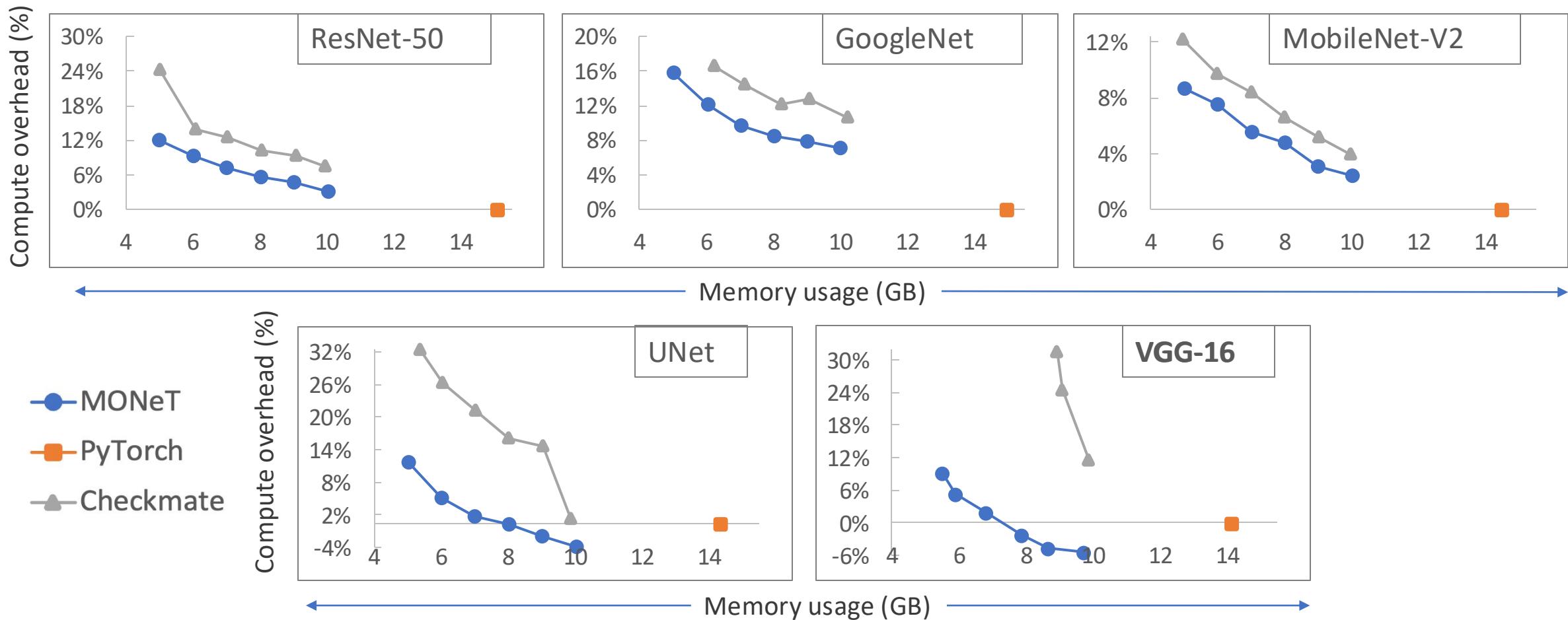


Checkmate

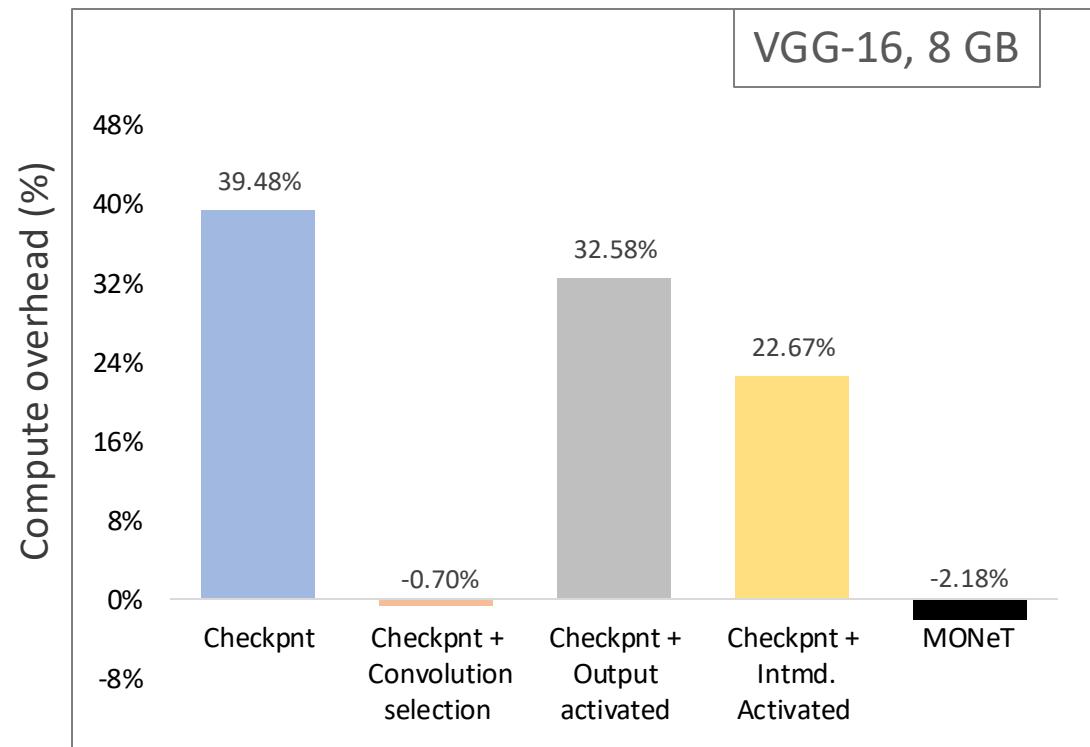
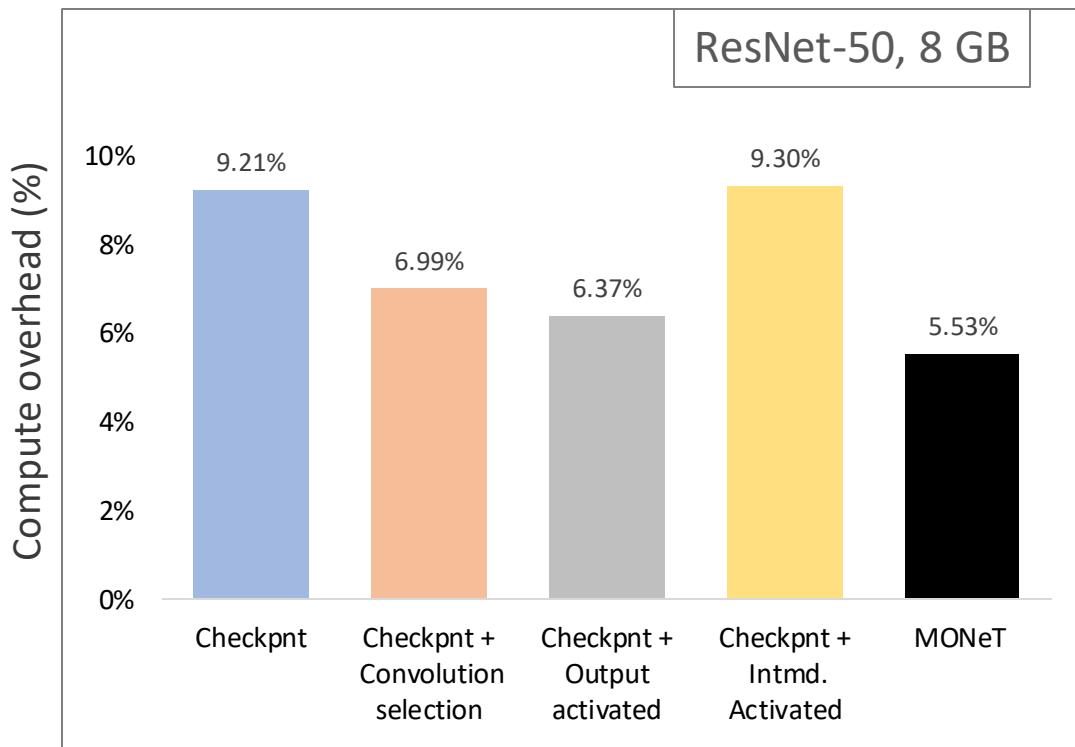


[1] Jain, Paras, et al. "Checkmate: Breaking the Memory Wall with Optimal Tensor Rematerialization." *MLSys*. 2020.

Checkmate

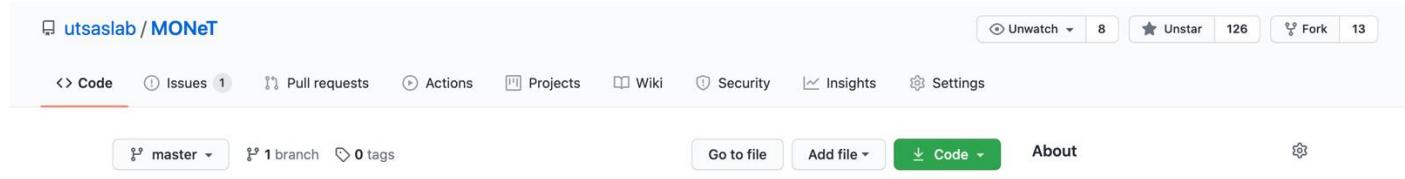


Ablation



Try MONeT!

<https://github.com/utsaslab/MONeT>

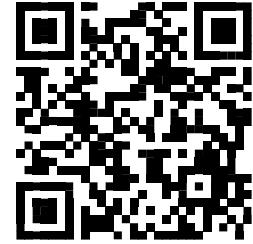


Works out-of-box with PyTorch 1.5.1

```
from monet=cvxpy_solver import Solution  
from monet=monet_wrapper import MONeTWrapper  
  
monet_model = MONeTWrapper(model, solution_file, input_shape)
```

No need to run a solver, plug in the pre-computed schedules ☺

<https://github.com/utsaslab/monet-schedules>



Conclusion

MONeT is a framework for memory optimization in training that is **automated** and creates schedules which **jointly optimize** global checkpointing and local operator implementations

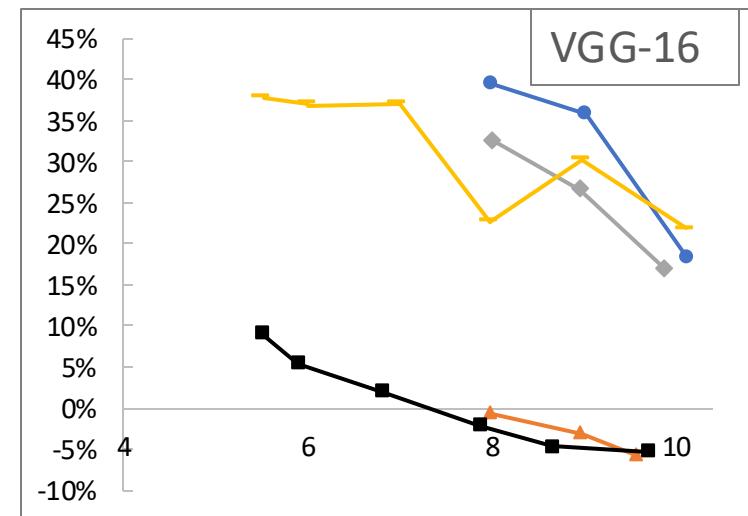
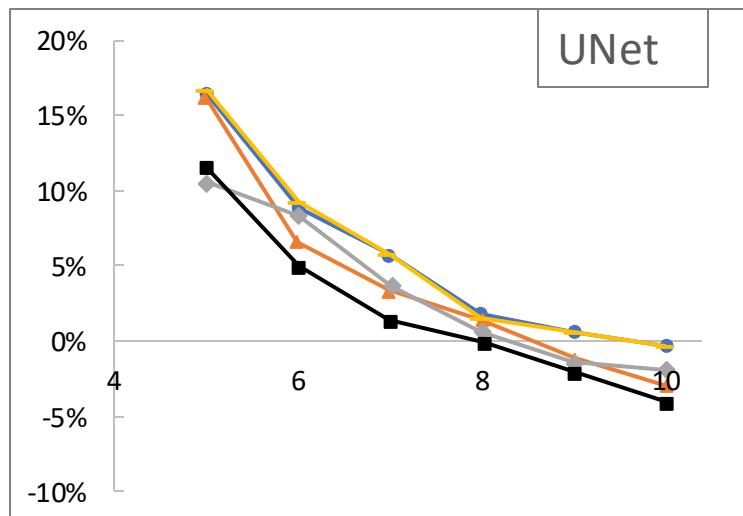
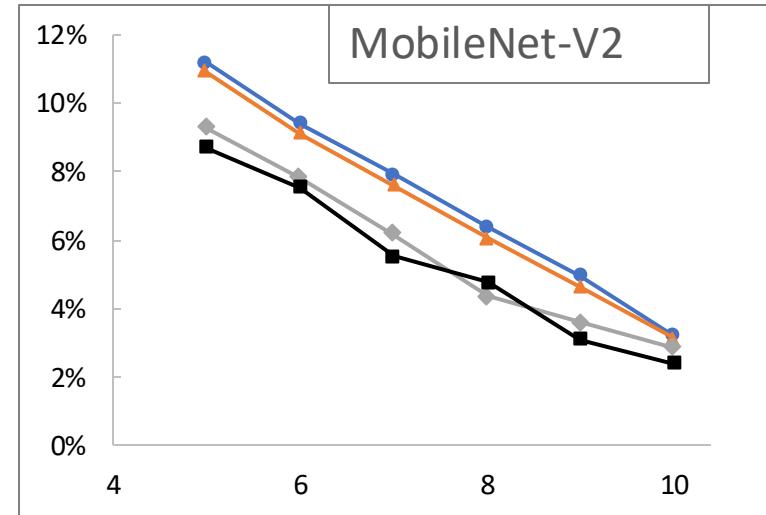
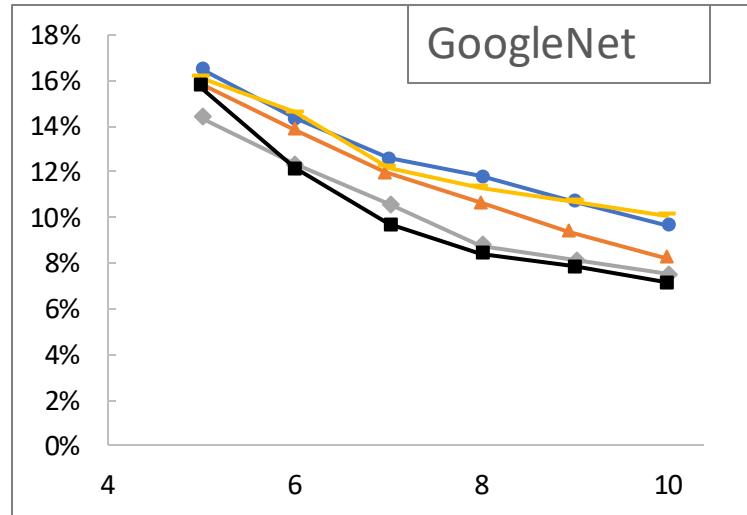
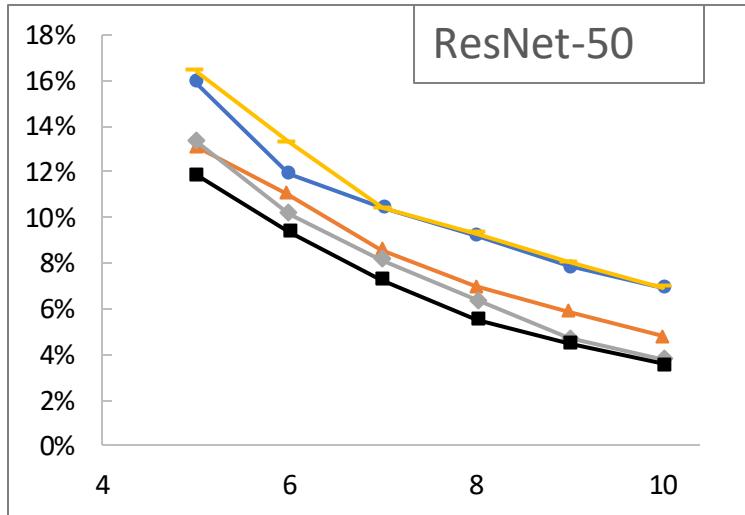
Reduces memory usage by **3x** while incurring **9–16% compute overhead** over PyTorch for ResNet-50, GoogleNet, etc.

<https://github.com/utsaslab/MONeT>

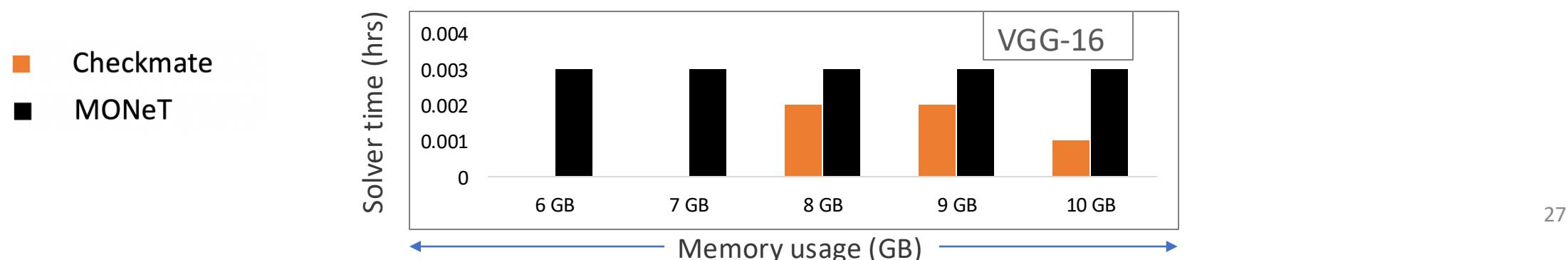
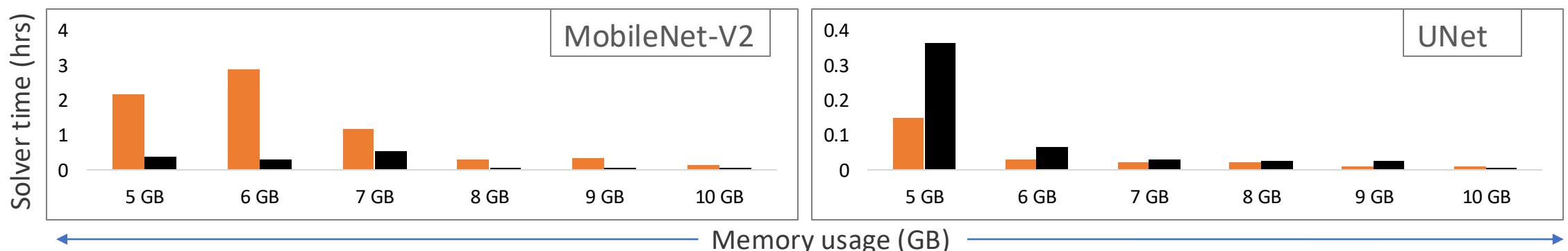
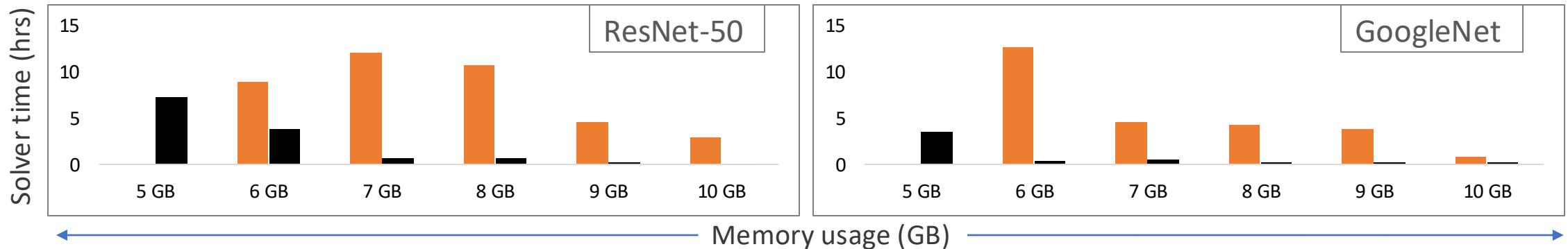
<https://github.com/utsaslab/monet-schedules>

Backup slides

Ablation



Solver time



- Checkmate
- MONeT

Directly adding operator optimizations to current frameworks won't scale

- Directly adding operator optimization in Checkmate formulation will lead to a linear blowup in the number of constraints, and number of auxiliary variables, leading to an at least quadratic expansion on computational costs.

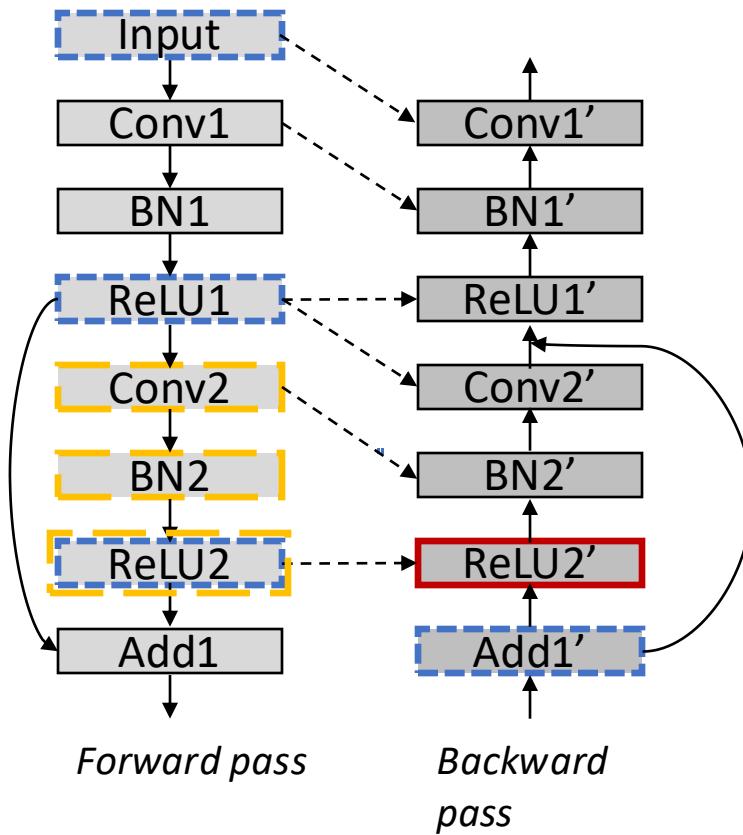
Scope for operator optimization in memory saving

Certain neural network operators can have multiple implementations

1. A convolution can be implemented using at least 8 different algorithms like GEMM, Winograd, FFT, etc. Each algorithm has a different workspace memory requirement and compute usage

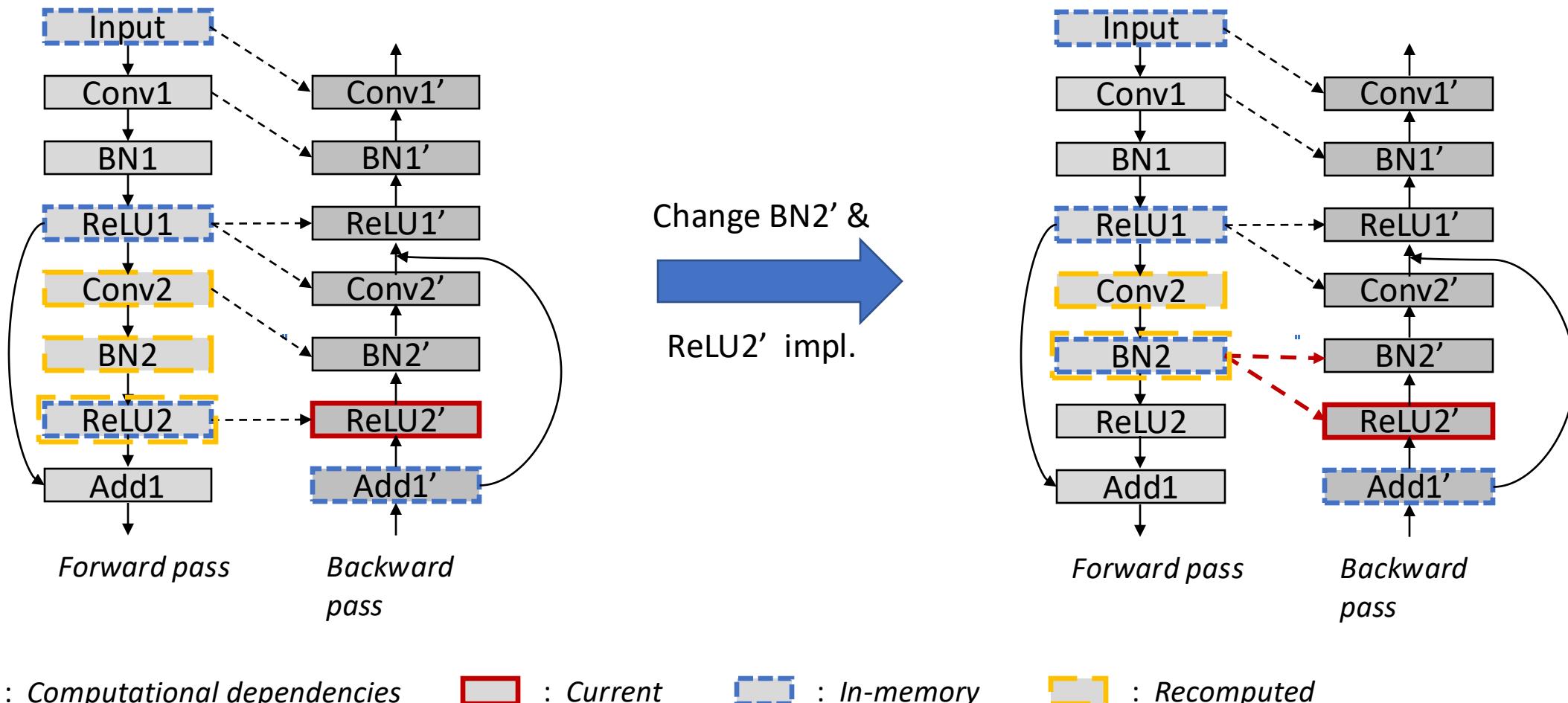


Operator optimizations affect checkpointing

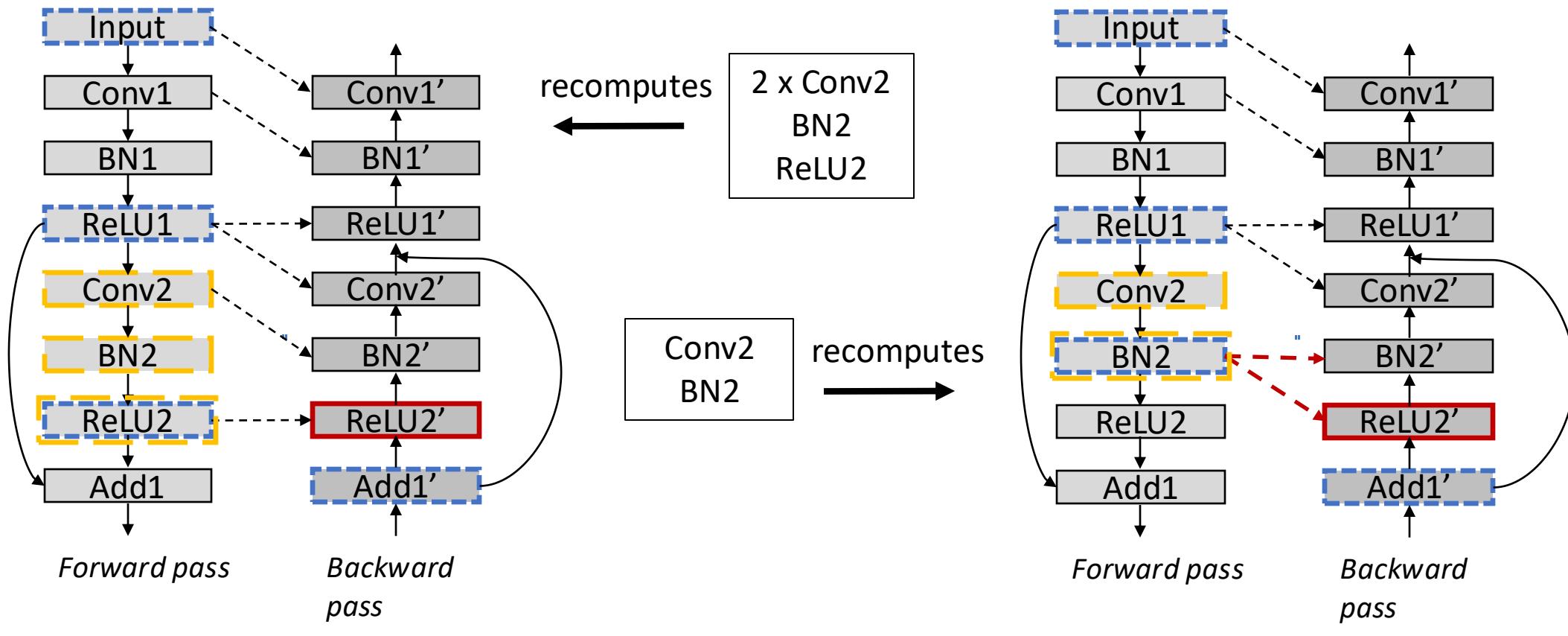


- $\text{ReLU2}'$: recomputes Conv2, BN2, ReLU2
- $\text{BN2}'$: recomputes Conv2

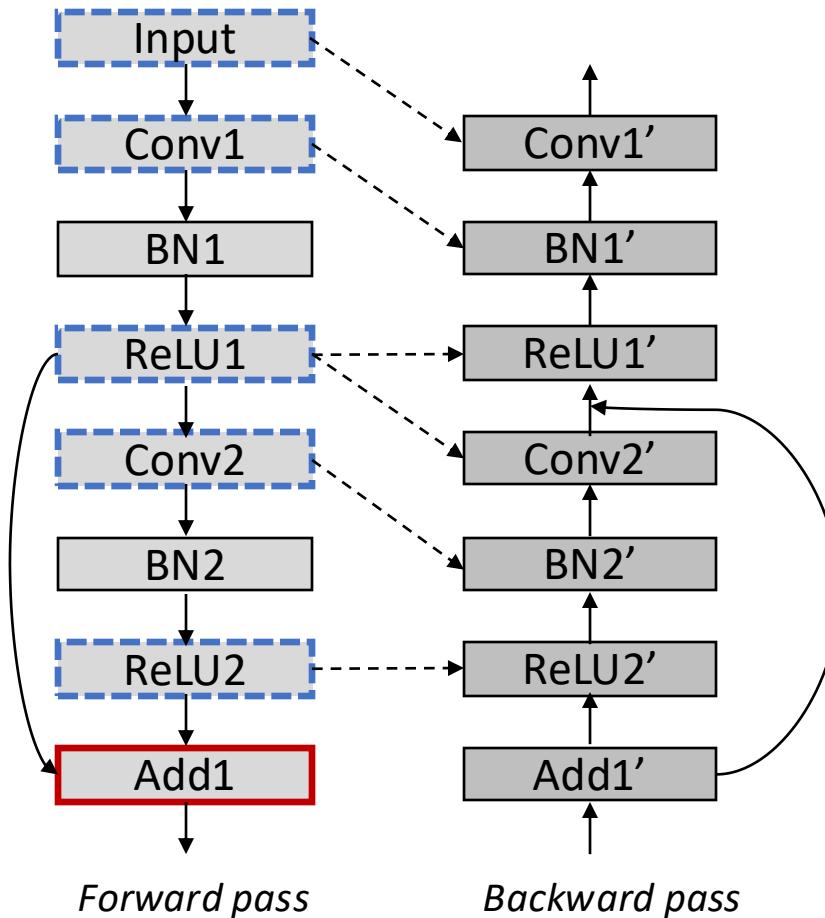
Operator optimizations affect checkpointing



Operator optimizations affect checkpointing



Memory usage during training

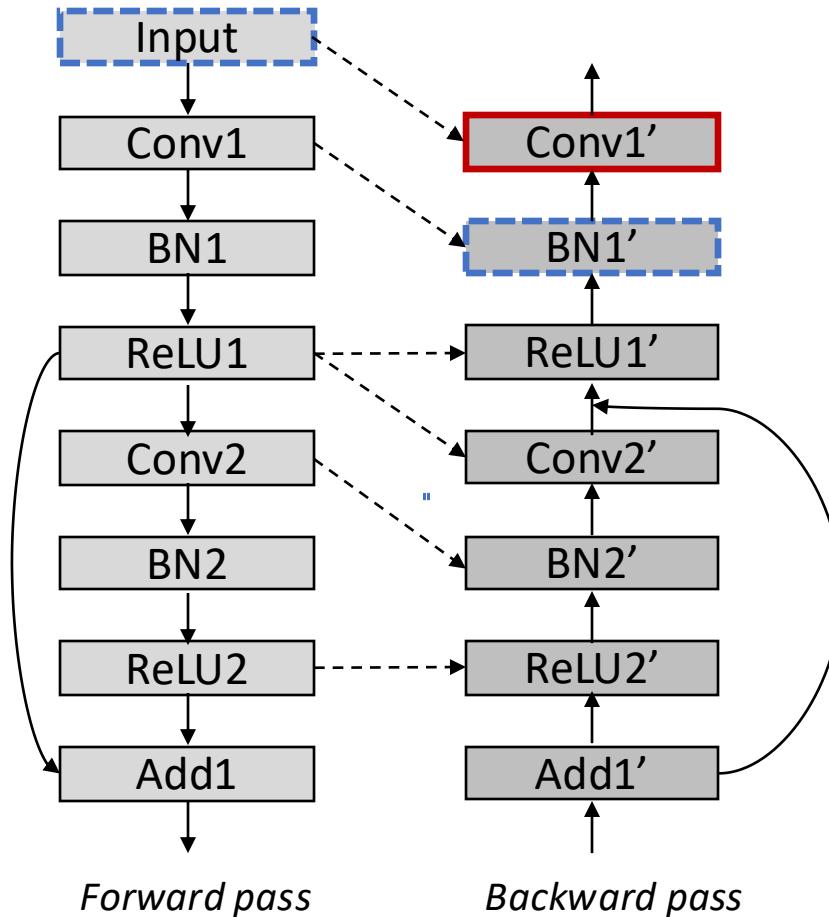


Local tensors (L)	ReLU1, ReLU2
Stored tensors (S)	Input, Conv1, ReLU1, Conv2, ReLU2

$$\text{Memory} = |\text{ReLU1}| + |\text{ReLU2}| + |\text{Input}| + |\text{Conv1}| + |\text{Conv2}| + |\text{Add1}| + \text{workspace}(\text{Add1})$$

Forward pass memory of node i :
 $|L_i \cup S_i| + |x_i| + workspace_i$

Memory usage during training



Local tensors (\hat{L})	BN1'
Stored tensors (S)	Input
Forward deps (D)	Input

$$\text{Memory} = |\text{BN1}'| + |\text{Input}| + |\text{Conv1}'| + \text{workspace}(\text{Conv1}')$$

Backward pass memory of node k :
 $|\widehat{L}_k \cup S_k \cup D_k| + |y_k| + workspace_k$