

# Learning with Auxiliary Activation for Memory-Efficient Training

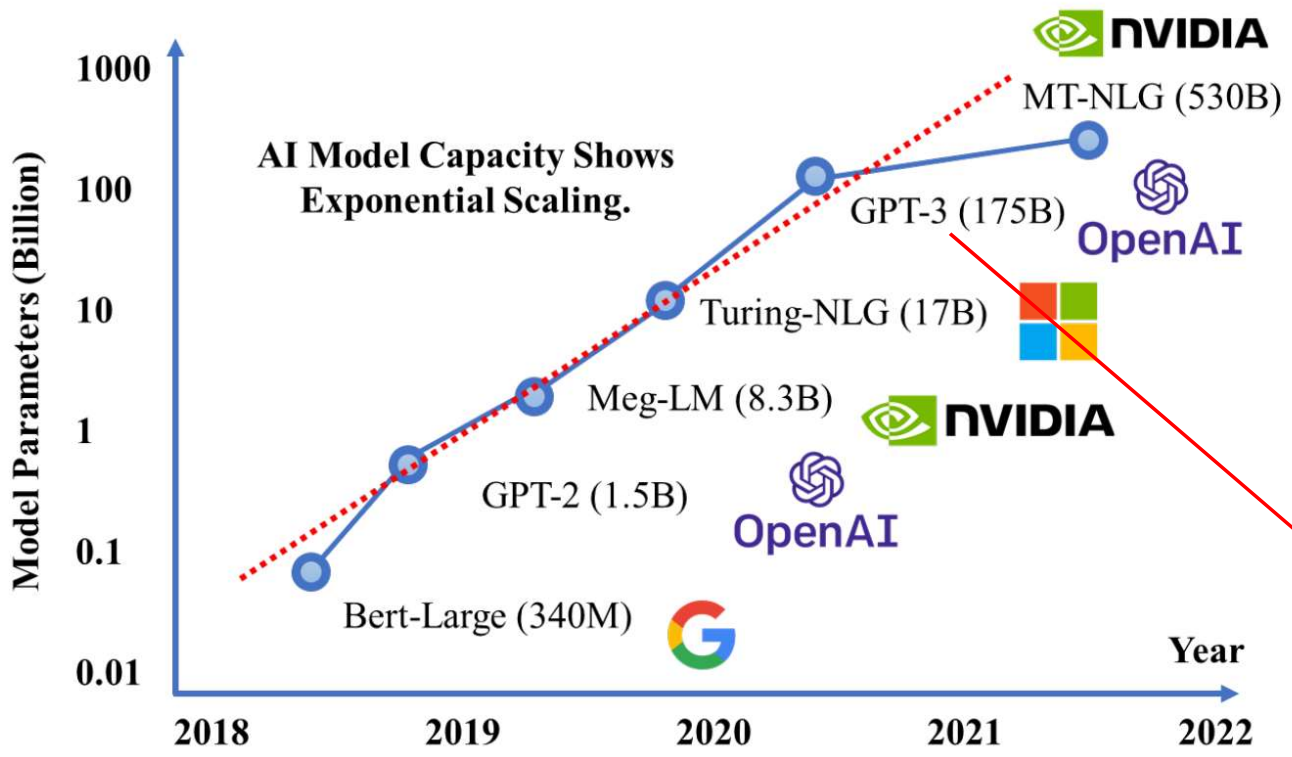
Sunghyeon Woo, Dongsuk Jeon

{wsh0917, djeon1@snu.ac.kr}

Seoul National University

 **Mobile Multimedia Systems Group**

# Introduction



< Rapid growth of model parameters (Yu et al, 2021) >

Model size ↑

Accuracy ↑

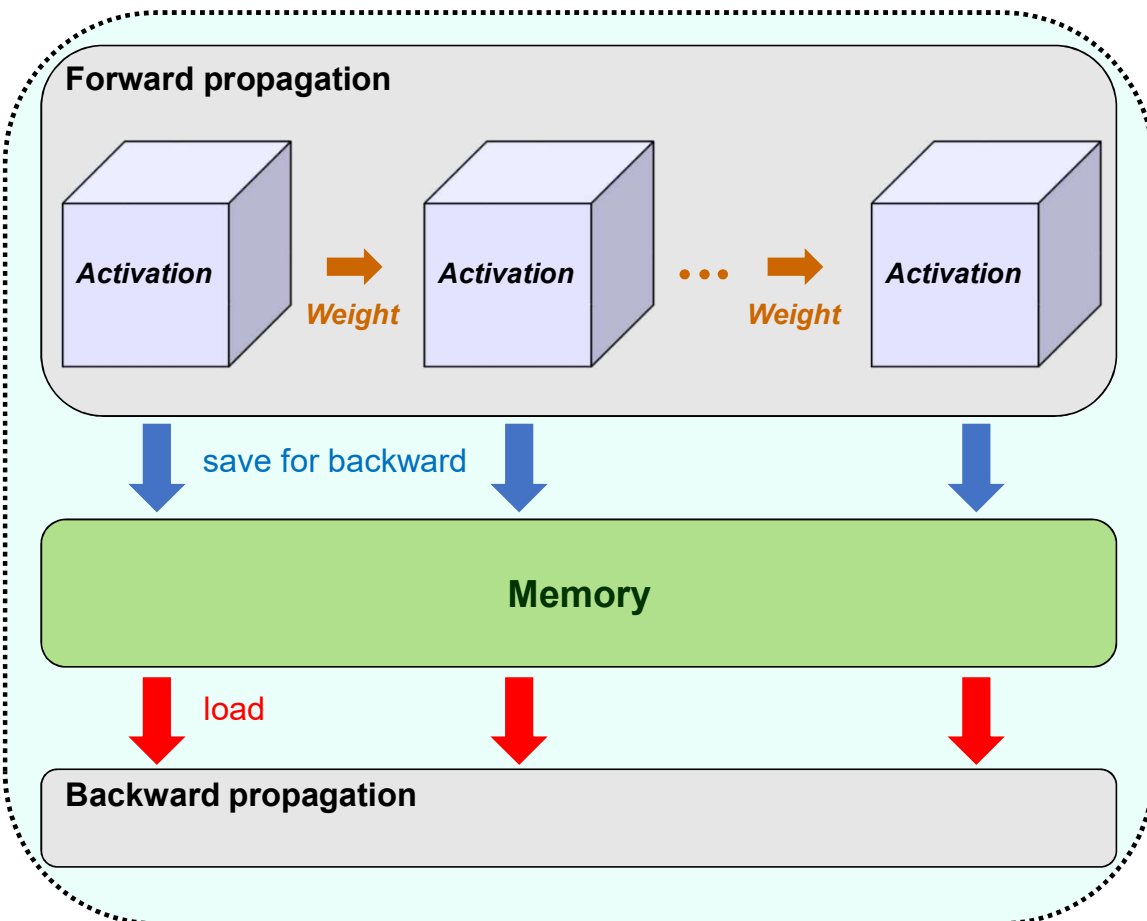
Training memory ↑



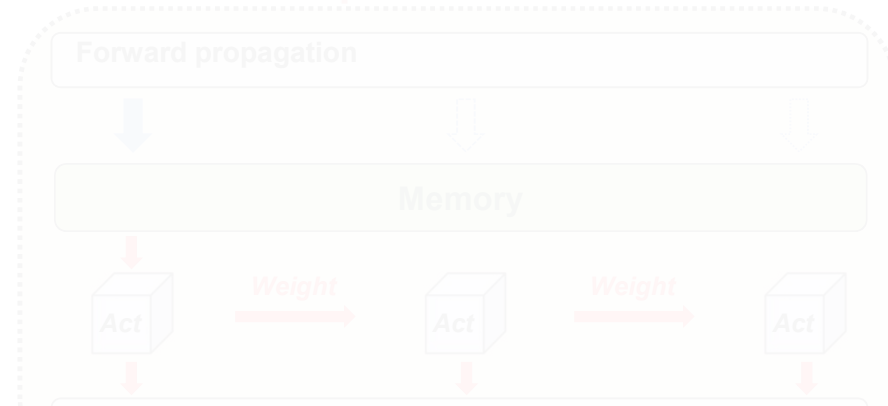
1024 x A100 Requires ...

# Introduction

## Backpropagation

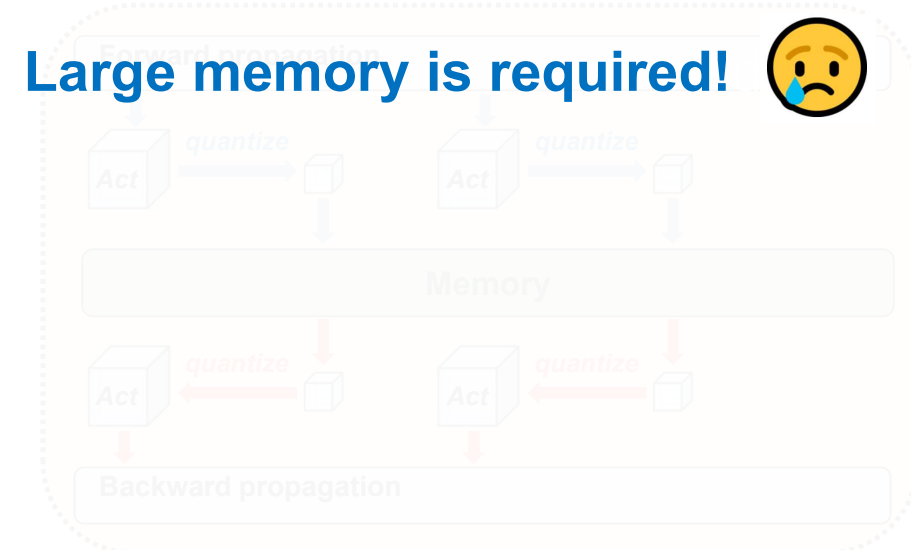


## Activation Recomputation



**Activations are stored in memory for backward propagation...**

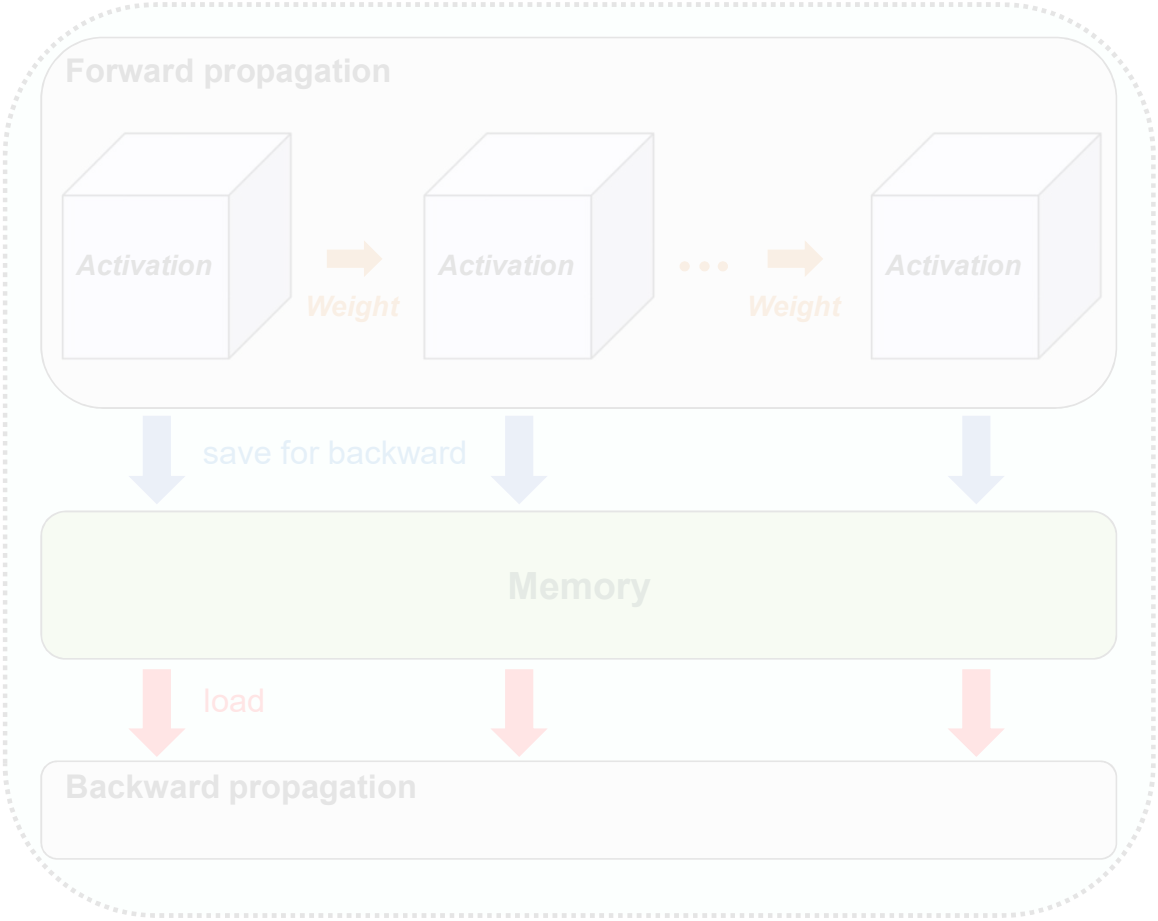
## Activation Compression Training



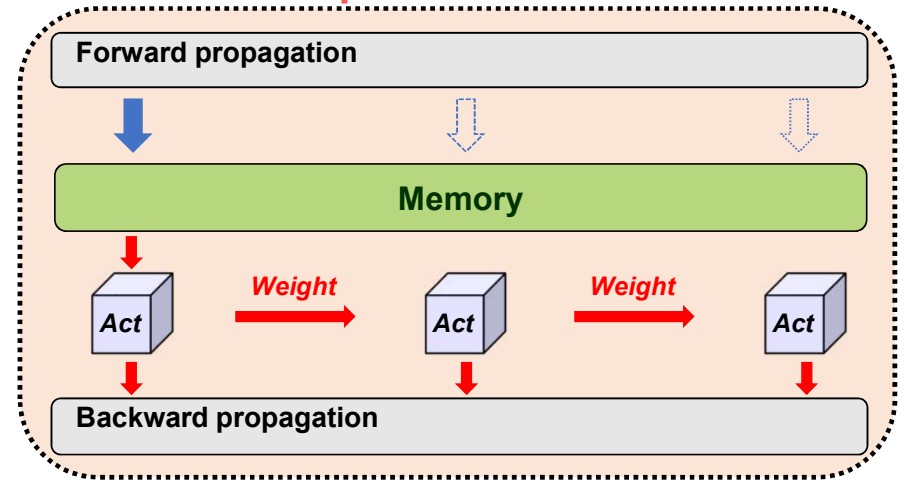
**Large memory is required!** 😞

# Introduction

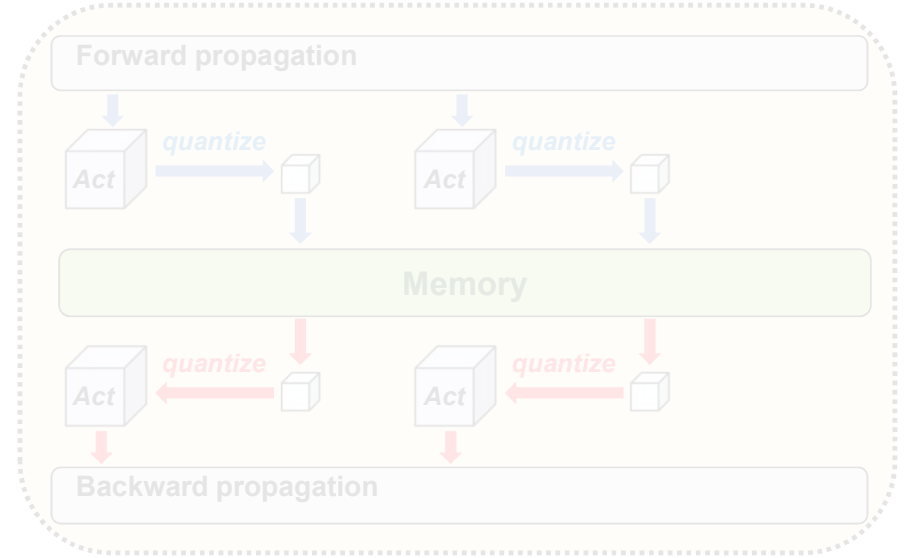
## Backpropagation



## Activation Recomputation

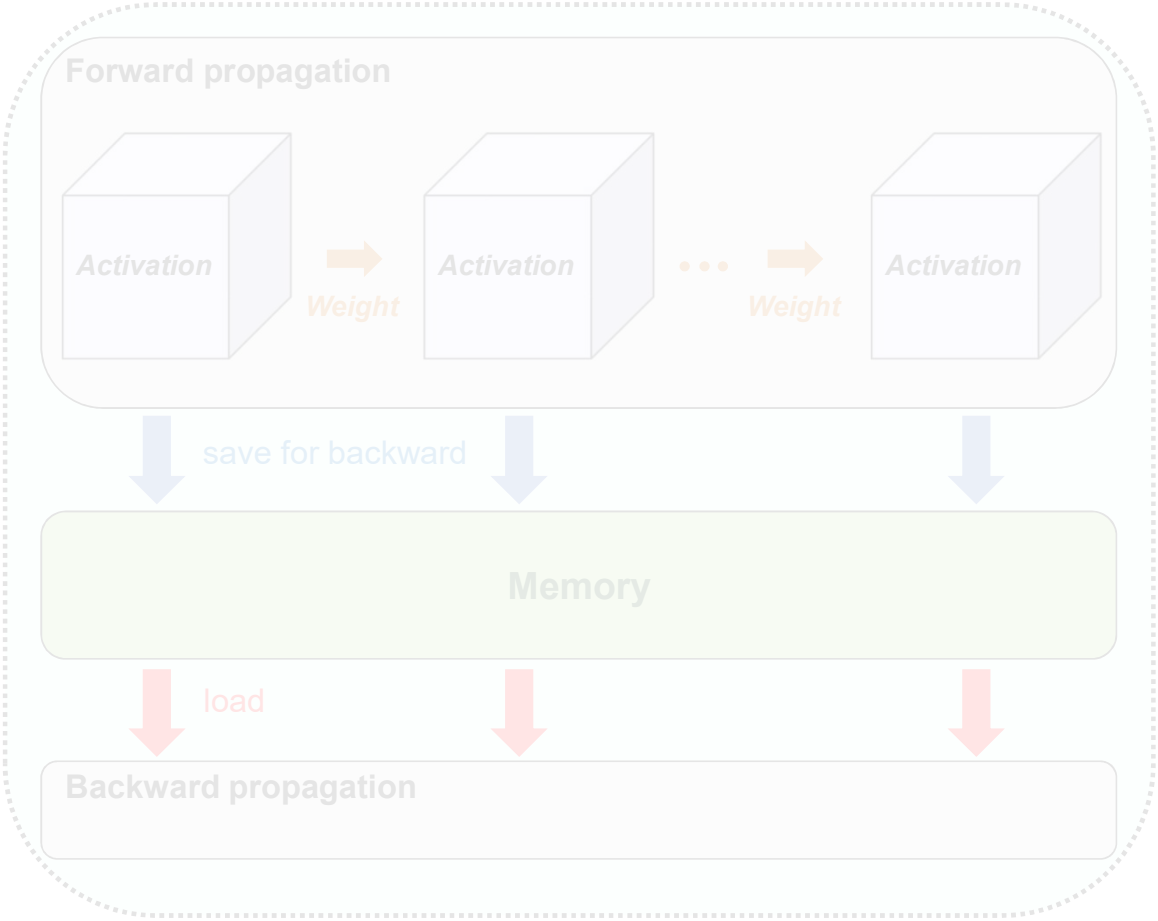


## Activation Compression Training

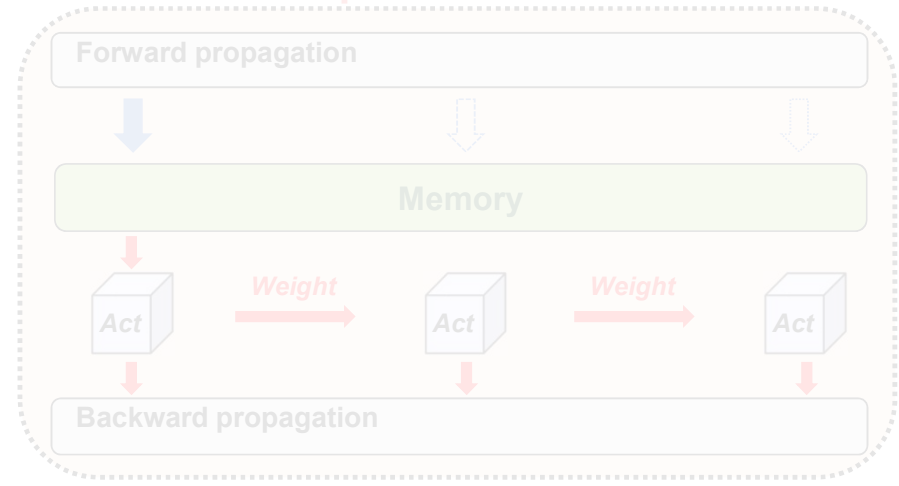


# Introduction

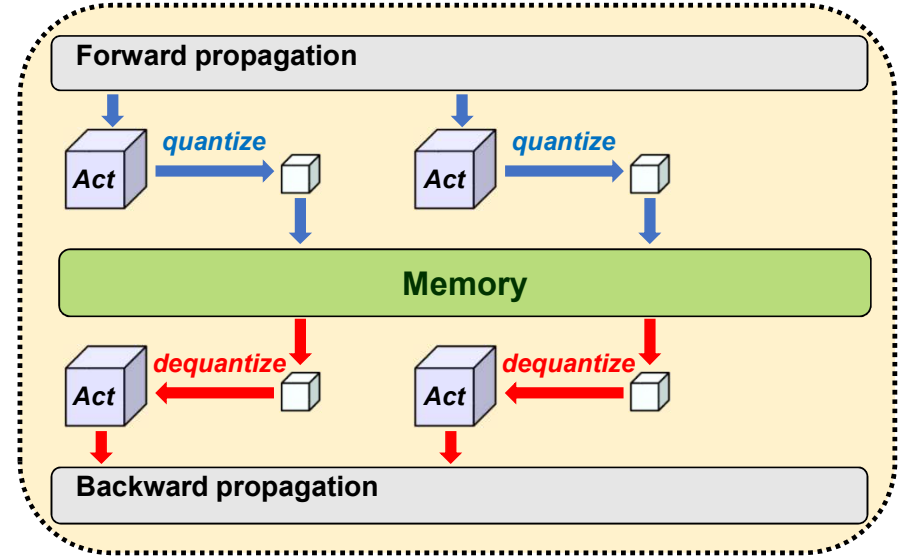
## Backpropagation



## Activation Recomputation



## Activation Compression Training

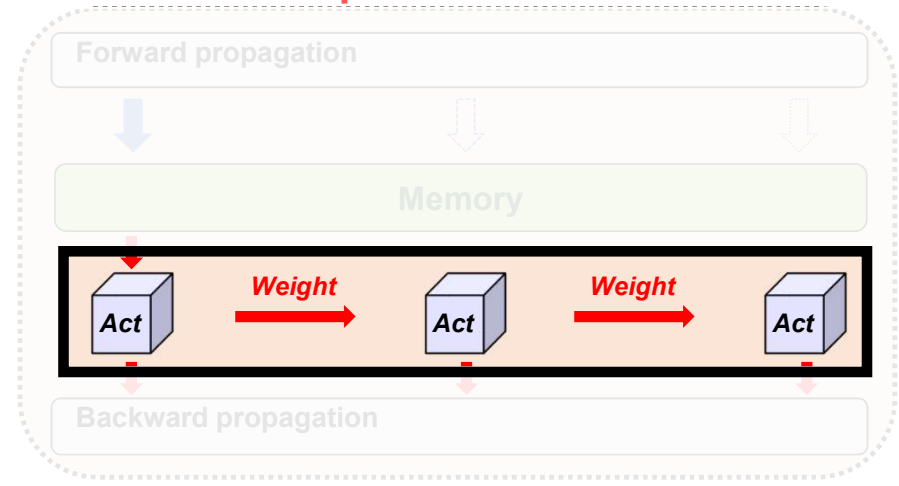


# Introduction

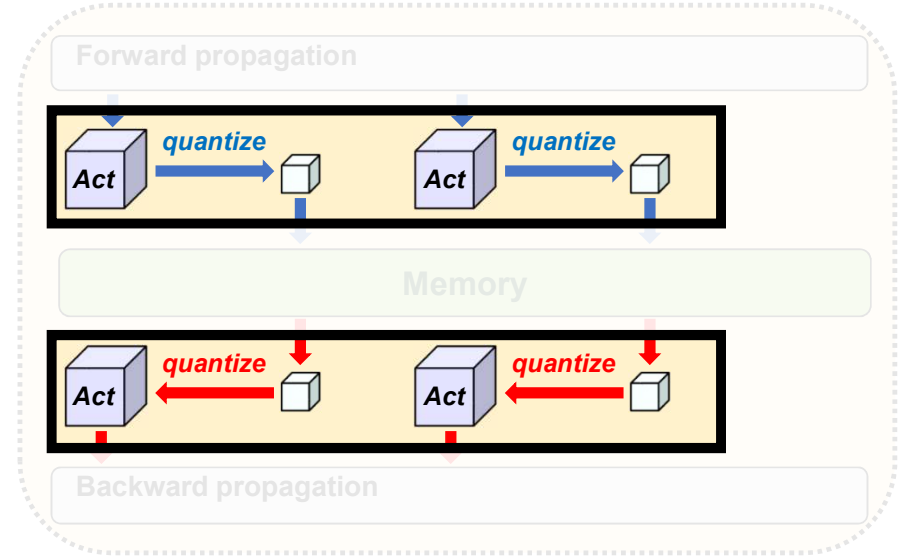
## Backpropagation



## Activation Recomputation

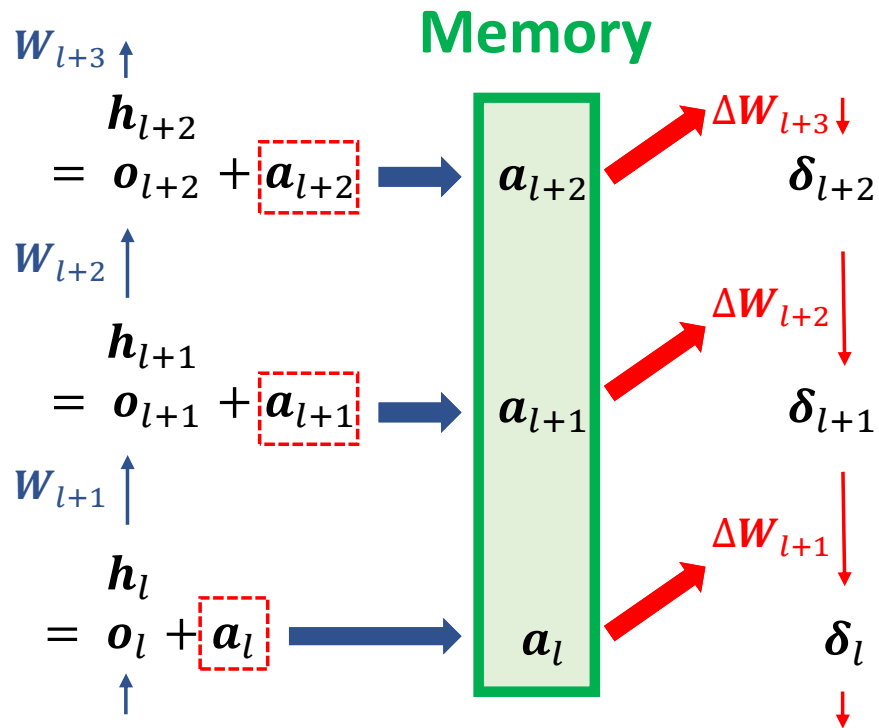


## Activation Compression Training



# Auxiliary Activation Learning

$a_l$  : Auxiliary Activation

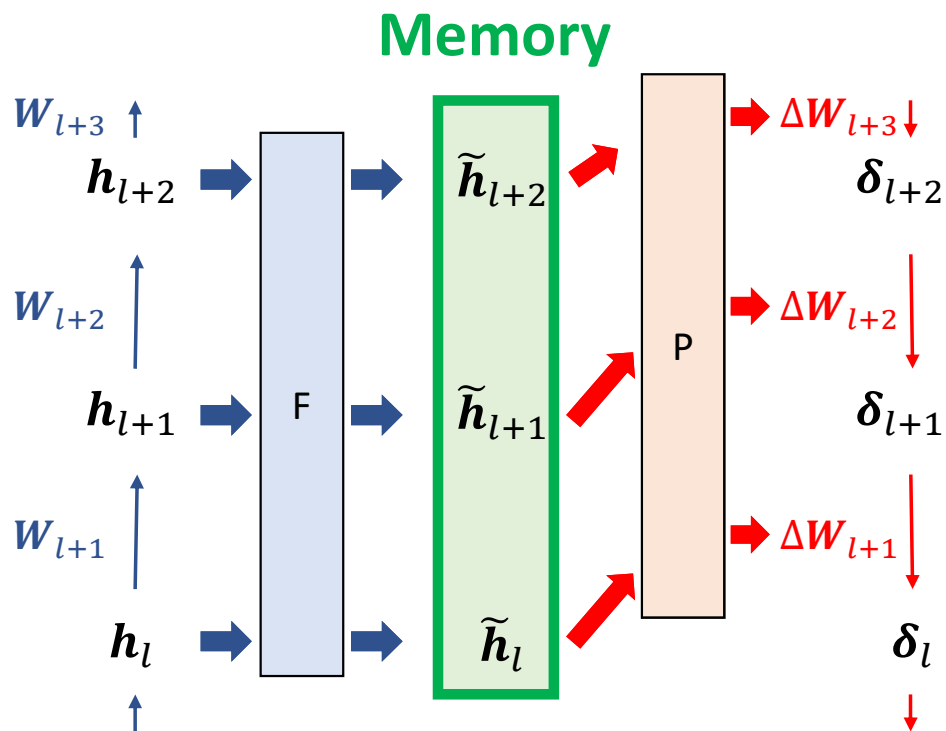


Auxiliary Activation Learning

- **Forward Propagation**  
: Add auxiliary activation to output activation of layer and store auxiliary activation instead of actual activation
- **Backward Propagation**  
: Use auxiliary activation for updating weights

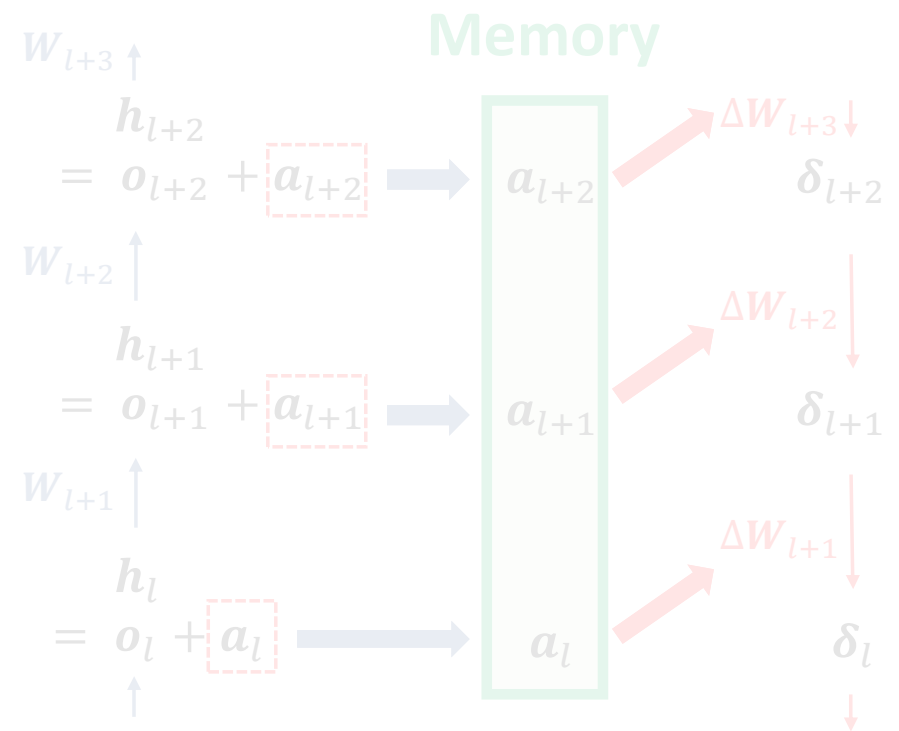
# Auxiliary Activation Learning

**F** : forward process    **P** : backward process



Conventional memory-saving algorithms

$a_l$  : Auxiliary Activation

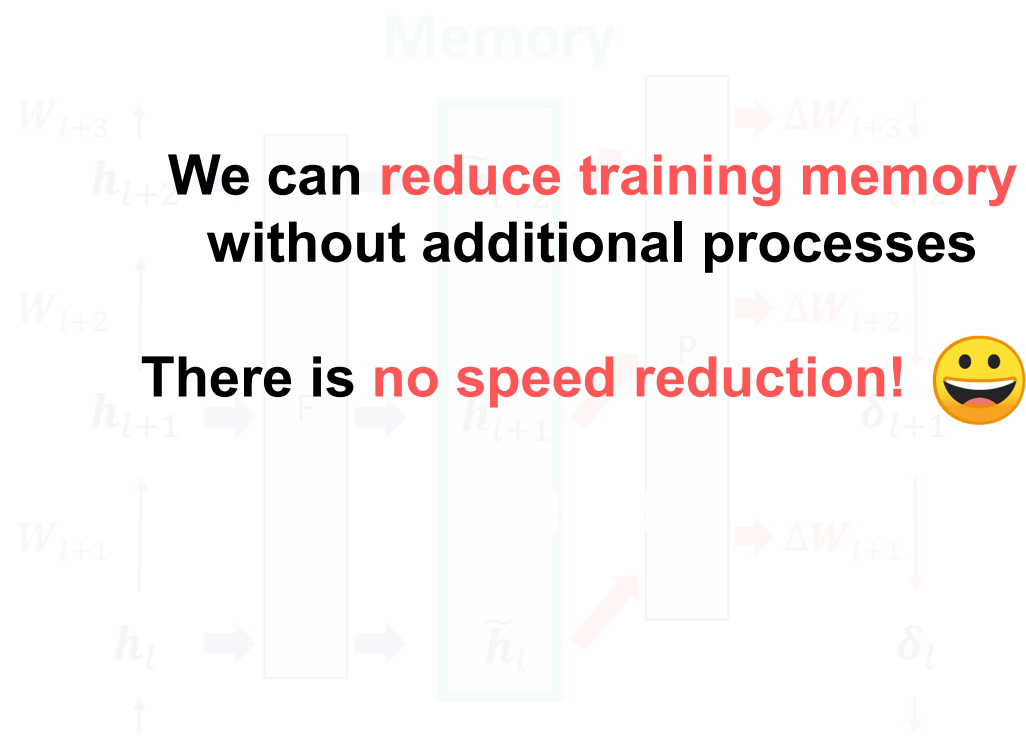


Auxiliary Activation Learning



# Auxiliary Activation Learning

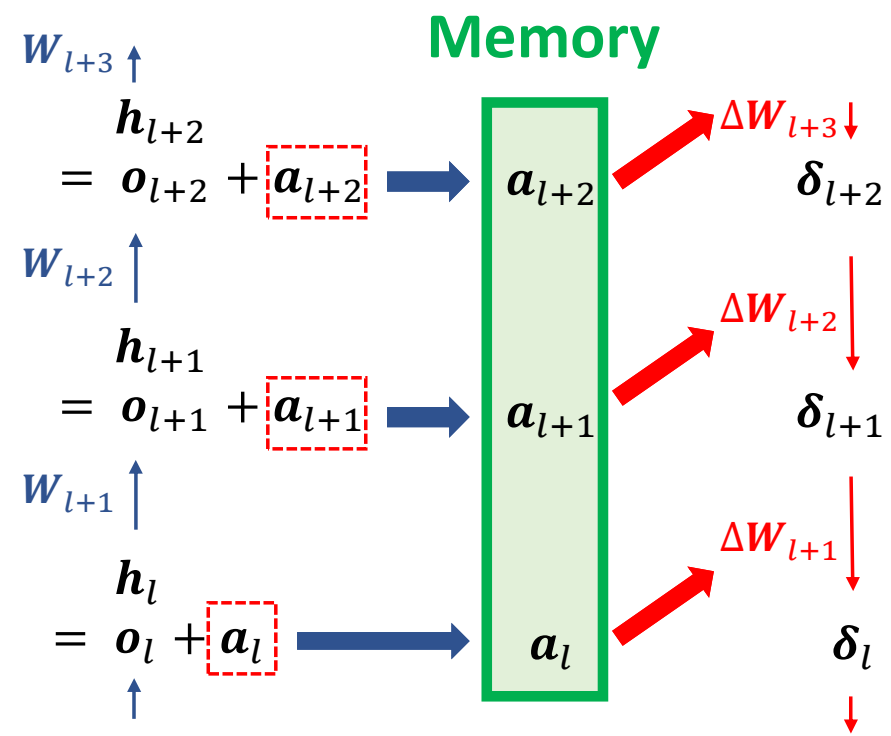
[F] : forward process    [P] : backward process



We can **reduce training memory** without additional processes

There is **no speed reduction!** 😊

$a_l$  : Auxiliary Activation



Auxiliary Activation Learning

## Learning criterion for alternative activation

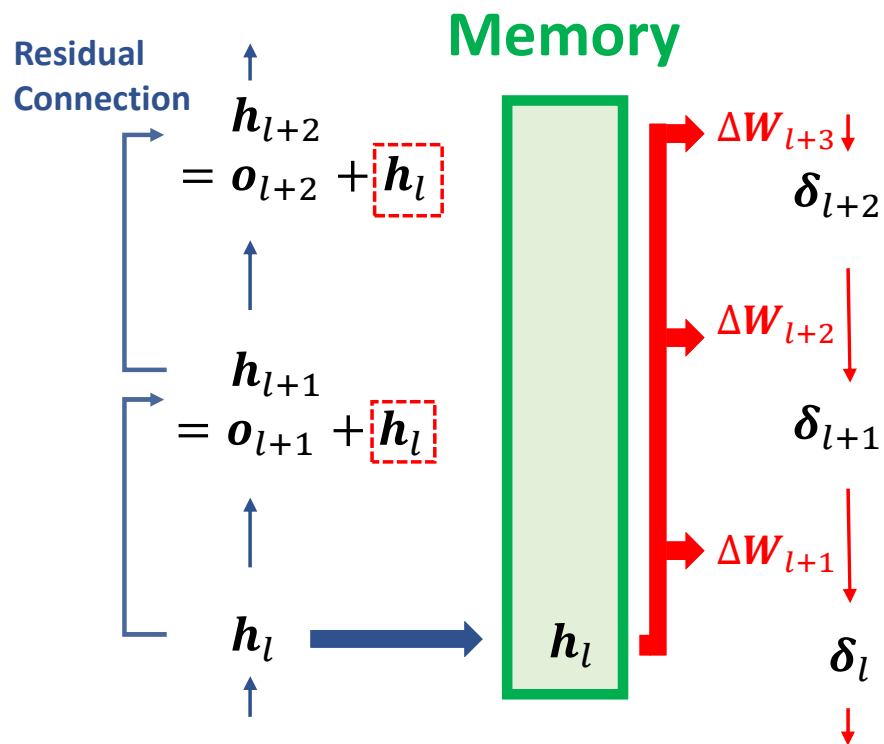
**Theorem 1.** Let weight updates  $W_{l+1}$  is calculated by alternative activation  $r_l$  instead of actual activation  $h_l$ . In this case, if the gradient of loss function  $f(W)$  is  $L$ -Lipschitz continuous, learning rate  $\eta$  satisfies  $0 \leq \eta \leq \frac{1}{L}$ , and  $\mathbf{r}_l^T (2\mathbf{h}_l - \mathbf{r}_l) \geq 0$ , then loss function  $f(W)$  is converged.

$$\text{Learning Indicator (LI)} \equiv \frac{\mathbf{r}_l^T (2\mathbf{h}_l - \mathbf{r}_l)}{\|\mathbf{h}_l\|^2}$$

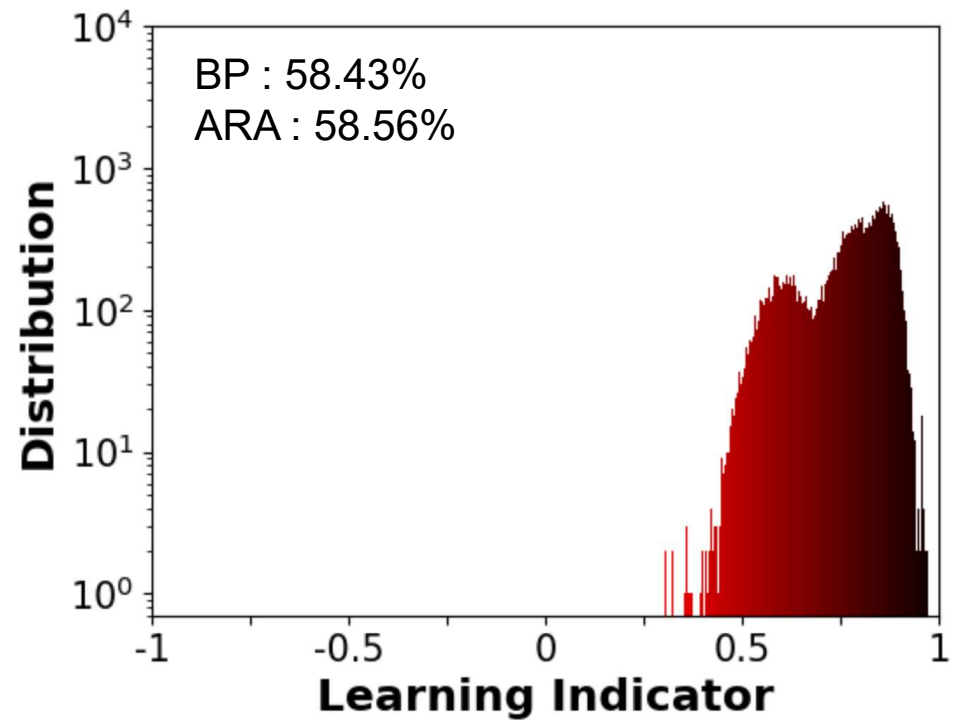
- $LI > 0$  : The loss is converged. The larger Learning indicator is, the better loss converges.
- $LI \leq 0$  : The loss is diverged.

# Two candidates of Auxiliary Activation

$h_l$  : Auxiliary Residual Activation (ARA)

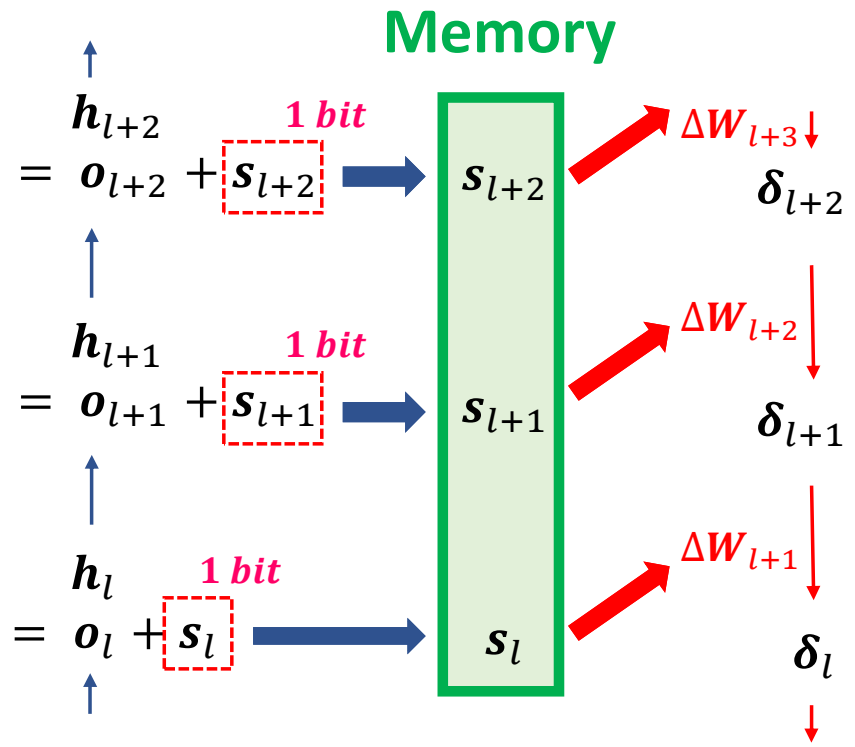


Train ResNet-18 on Tiny ImageNet

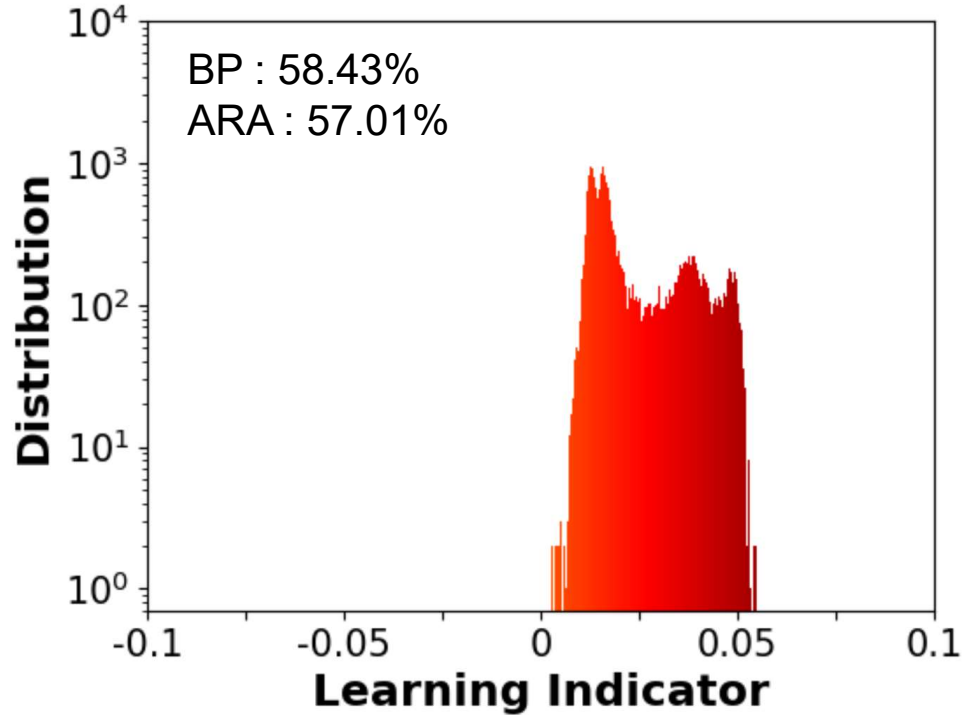


# Two candidates of Auxiliary Activation

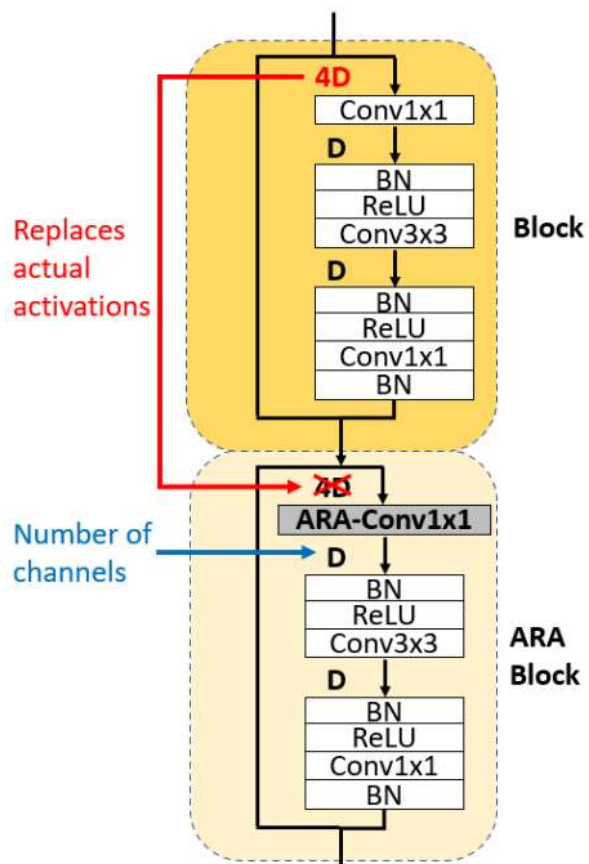
$s_l$  : Auxiliary Sign Activation (**ASA**)



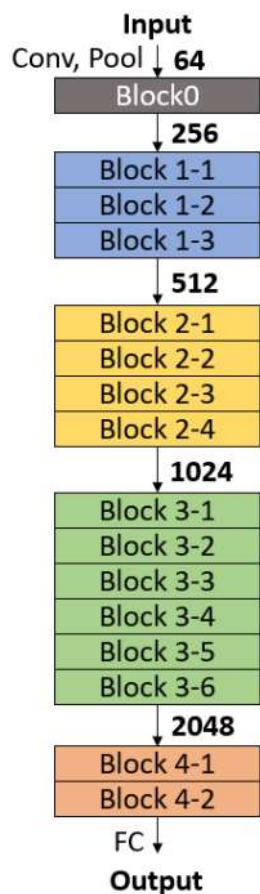
Train ResNet-18 on Tiny ImageNet



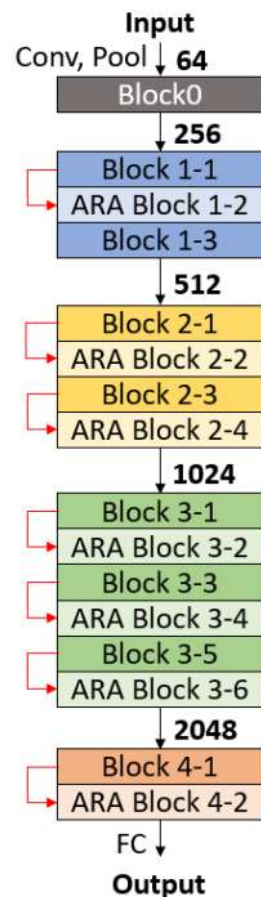
# Results



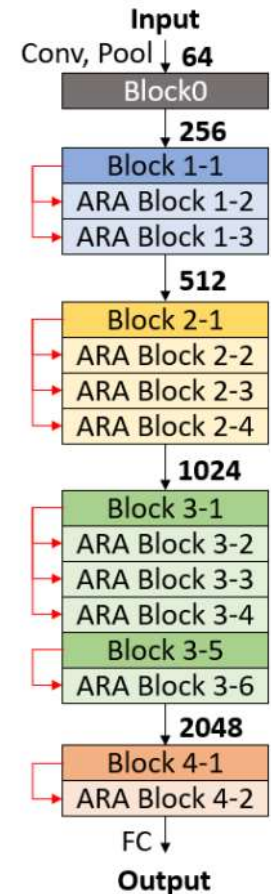
(a) ARA Block



(b) BP



(c) ARA(2,2,2,2)



(d) ARA(3,4,4,2)

Figure 3: Training ResNet-50 using Auxiliary Residual Activation (ARA).

# Results

Table 2: Test accuracy, training memory with compression rate (bracketed), and training time (italic) of one epoch in ResNet training on ImageNet with 512 batch size and six RTX-3090 GPUs.

Models	Baseline	-	ARA (2,2,2,2)	ARA (3,4,2,2)	ARA (3,4,4,2)	ARA (3,4,6,2)	Save memory without speed reduction
ResNet-50	BP	<b>76.01</b> 44.6 GB <i>17m 35s</i>	75.97 (1.12x) <i>17m 59s</i>	75.89 (1.2x) <i>17m 54s</i>	75.62 (1.21x) <i>17m 55s</i>	75.23 (1.22x) <i>18m 3s</i>	→
	GCP	<b>76.01</b> (2.2x) <i>36m 42s</i>	75.97 (2.88x) <i>37m 2s</i>	75.89 (3.44x) <i>37m 9s</i>	75.62 (3.54x) <i>36m 58s</i>	75.23 (3.66x) <i>37m 14s</i>	
	ActNN	75.96 (11.8x) <i>48m 48s</i>	75.93 (14x) <i>48m 48s</i>	75.67 (14.6x) <i>48m 10s</i>	75.51 (14.8x) <i>47m 58s</i>	75.12 ( <b>15.01x</b> ) <i>47m 08s</i>	
ResNet-152	BP	77.38 90.2 GB <i>35m 37s</i>	77.14 (1.16x) <i>36m 17s</i>	<b>77.41</b> (1.21x) <i>36m 24s</i>	76.84 (1.27x) <i>36m 11s</i>	76.64 (1.29x) <i>36m 30s</i>	→
	GCP	77.38 (2.1x) <i>1h 18m</i>	77.14 (2.92x) <i>1h 19m</i>	<b>77.41</b> (3.26x) <i>1h 19m</i>	76.84 (3.75x) <i>1h 19m</i>	76.64 ( <b>3.95x</b> ) <i>1h 20m</i>	



# Results

Table 2: Test accuracy, training memory with compression rate (bracketed), and training time (italic) of one epoch in ResNet training on ImageNet with 512 batch size and six RTX-3090 GPUs.

Models	Baseline	-	ARA (2,2,2,2)	ARA (3,4,2,2)	ARA (3,4,4,2)	ARA (3,4,6,2)
ResNet-50	BP	<b>76.01</b> 44.6 GB <i>17m 35s</i>	75.97 (1.12x) <i>17m 59s</i>	75.89 (1.2x) <i>17m 54s</i>	75.62 (1.21x) <i>17m 55s</i>	75.23 (1.22x) <i>18m 3s</i>
	GCP	<b>76.01</b> (2.2x) <i>36m 42s</i>	75.97 (2.88x) <i>37m 2s</i>	75.89 (3.44x) <i>37m 9s</i>	75.62 (3.54x) <i>36m 58s</i>	75.23 (3.66x) <i>37m 14s</i>
	ActNN	75.96 (11.8x) <i>48m 48s</i>	75.93 (14x) <i>48m 48s</i>	75.67 (14.6x) <i>48m 10s</i>	75.51 (14.8x) <i>47m 58s</i>	75.12 <b>(15.01x)</b> <i>47m 08s</i>
ResNet-152	BP	77.38 90.2 GB <i>35m 37s</i>	77.14 (1.16x) <i>36m 17s</i>	<b>77.41</b> (1.21x) <i>36m 24s</i>	76.84 (1.27x) <i>36m 11s</i>	76.64 (1.29x) <i>36m 30s</i>
	GCP	77.38 (2.1x) <i>1h 18m</i>	77.14 (2.92x) <i>1h 19m</i>	<b>77.41</b> (3.26x) <i>1h 19m</i>	76.84 (3.75x) <i>1h 19m</i>	76.64 <b>(3.95x)</b> <i>1h 20m</i>

Orthogonality:  
maximize  
memory saving

# Results

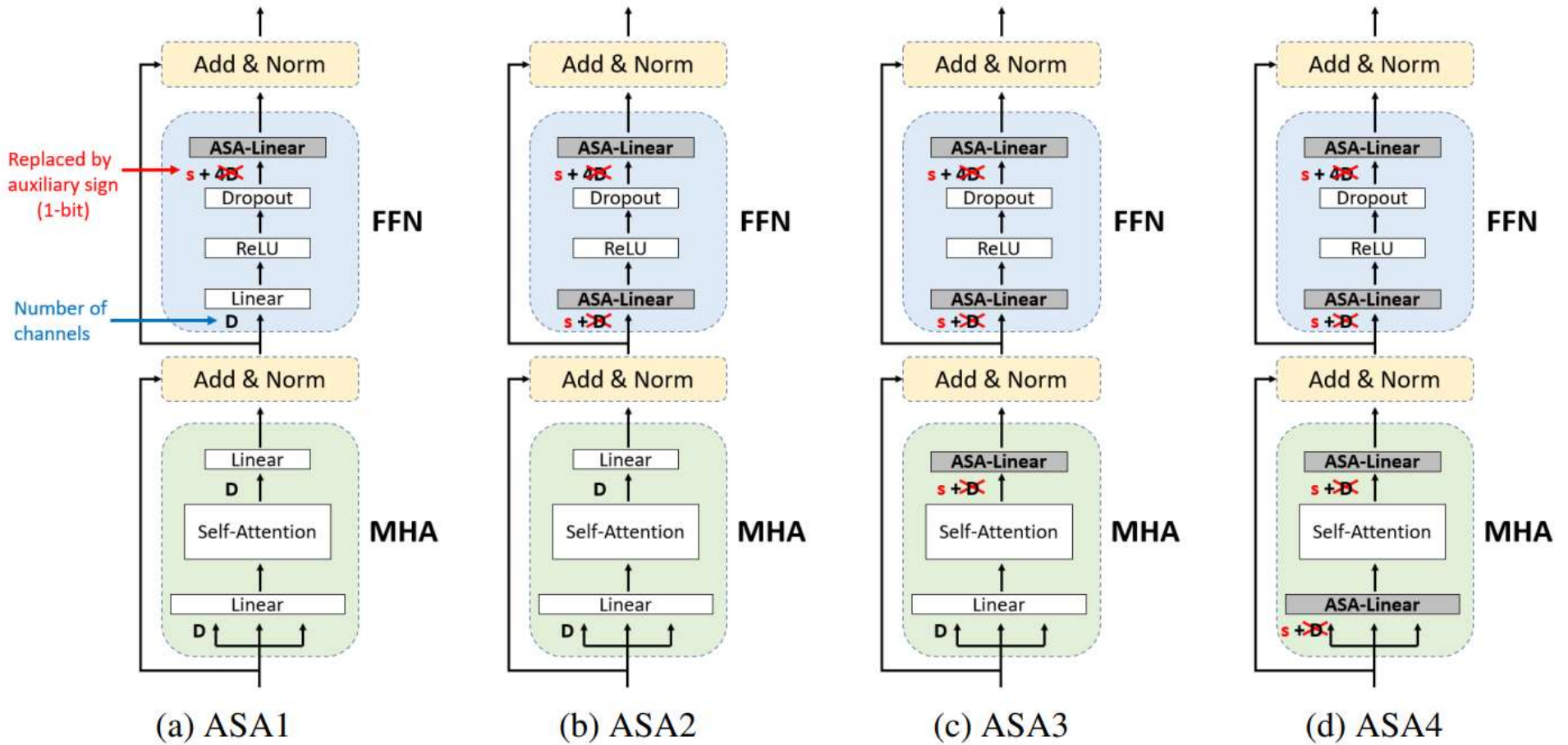


Figure 5: Transformer-like networks trained using Auxiliary Sign Activations (ASA). MHA and FFN represent multihead attention and feedforward network, respectively.



# Results

Table 3: Test/validation scores, compression rate (bracketed), and training time (italic) of one epoch of Transformer, BERT, ViT, and MLP-Mixer training with 4096, 32, 512, and 256 batch sizes, respectively.

Models	Dataset	Baseline	-	ASA1	ASA2	ASA3	ASA4
Trans- former	IWSLT	BP	35.23 3.6 GB <i>3m 15s</i>	35.44 (1.1x) <i>3m 21s</i>	34.87 (1.1x) <i>3m 23s</i>	34.9 (1.2x) <i>3m 26s</i>	35.02 (1.3x) <i>3m 38s</i>
		Mesa	<b>35.45</b> (1.6x) <i>5m 52s</i>	34.74 (1.7x) <i>5m 47s</i>	35.19 (1.7x) <i>5m 44s</i>	35.11 (1.7x) <i>5m 36s</i>	34.84 <b>(1.8x)</b> <i>5m 8s</i>
BERT- Large	MRPC	BP	88.56 10.9 GB <i>1m 28s</i>	88.69 (1.1x) <i>1m 29s</i>	88.23 (1.2x) <i>1m 30s</i>	<b>88.97</b> (1.3x) <i>1m 30s</i>	88.32 (1.3x) <i>1m 31s</i>
		Mesa	88.3 (2.1x) <i>2m 8s</i>	88.35 (2.3x) <i>2m 5s</i>	88.51 (2.4x) <i>2m 3s</i>	88.25 (2.4x) <i>2m 1s</i>	88.23 <b>(2.5x)</b> <i>1m 55s</i>
	MNLI	BP	86.52 10.9 GB <i>2h 37m</i>	<b>86.65</b> (1.1x) <i>2h 40m</i>	86.49 (1.2x) <i>2h 40m</i>	86.42 (1.3x) <i>2h 41m</i>	86.39 (1.3x) <i>2h 41m</i>
		Mesa	86.32 (2.1x) <i>3h 51m</i>	86.37 (2.3x) <i>3h 39m</i>	86.29 (2.4x) <i>3h 35m</i>	86.54 (2.4x) <i>3h 32m</i>	86.17 <b>(2.5x)</b> <i>3h 24m</i>
ViT- Large	CIFAR-100	BP	92.93 48 GB <i>2m 46s</i>	92.81 (1.3x) <i>2m 48s</i>	92.84 (1.5x) <i>2m 50s</i>	<b>92.97</b> (1.5x) <i>2m 51s</i>	92.65 (1.6x) <i>2m 52s</i>
		Mesa	92.89 (3.0x) <i>4m 46s</i>	92.75 (3.5x) <i>4m 15s</i>	92.72 (3.6x) <i>3m 57s</i>	92.95 (3.7x) <i>3m 49s</i>	92.84 <b>(4.3x)</b> <i>3m 33s</i>
Mixer- Large	CIFAR-100	BP	91.62 86.8 GB <i>3m 39s</i>	91.45 (1.3x) <i>3m 47s</i>	91.72 (1.4x) <i>3m 50s</i>	91.66 (1.7x) <i>3m 57s</i>	<b>91.91</b> <b>(2.0x)</b> <i>4m 1s</i>

## Results

Table 4: Largest models that can be trained using a single GPU with 24GB memory. (ResNet: depth = number of layers, width = width of the first bottleneck block. BERT-Large: depth = number of transformer blocks, width = hidden size.)

Models	ResNet				BERT-Large			
	BP	ARA	ActNN	ARA+ActNN	BP	ASA4	Mesa	Mesa+ActNN
Depth	146	165	622	718	50	60	64	70
Width	62	76	214	238	1600	1728	1792	1856

**Thank you very much!**