

# Network Binarization toward Hardware-friendly Deep Learning

Haotong Qin

ETH Zürich CVL & Beihang University

# Background

## □ Vision

- Classification
- Detection
- Localization
- Segmentation

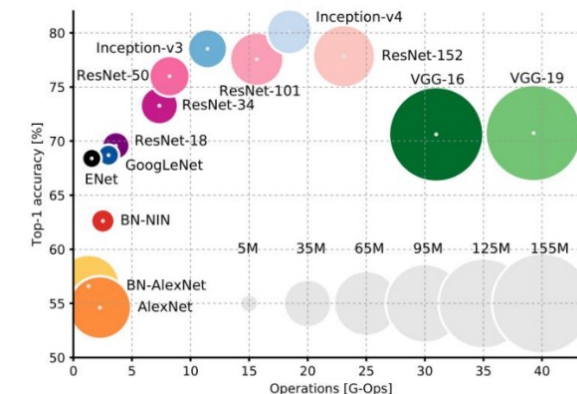
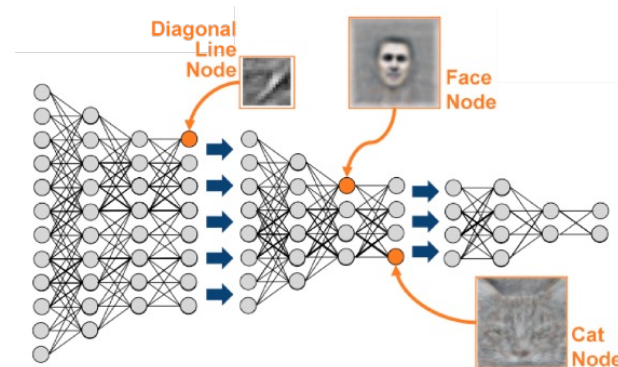
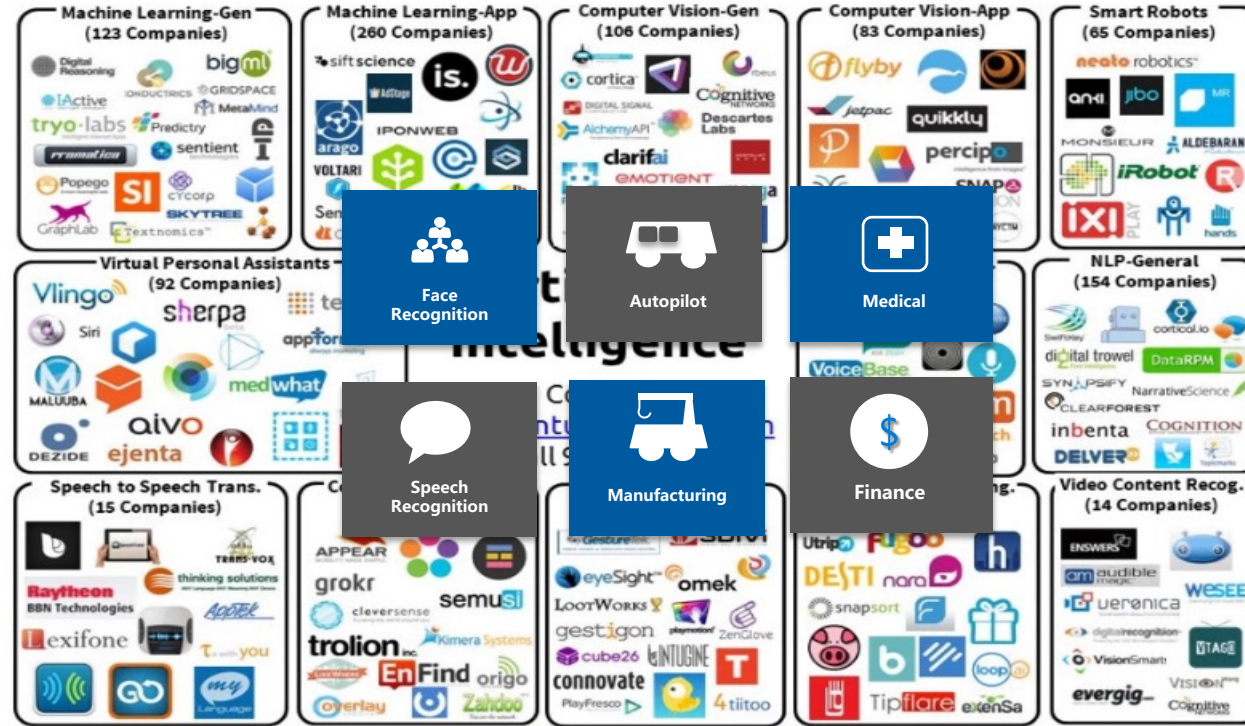
## □ Language

- Information retrieval
- Relation extraction
- Machine translation

## □ Speech

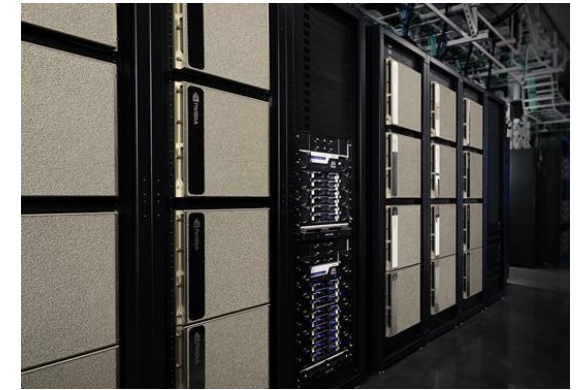
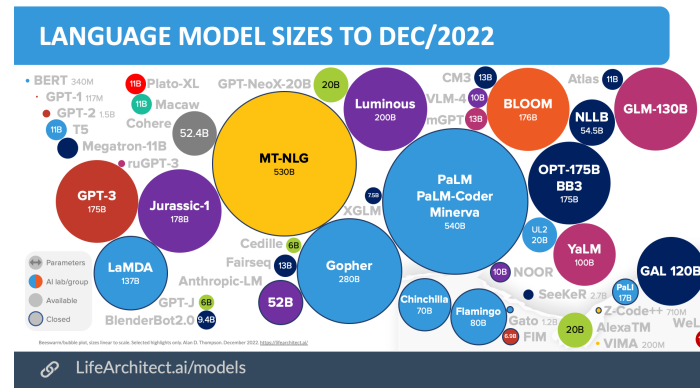
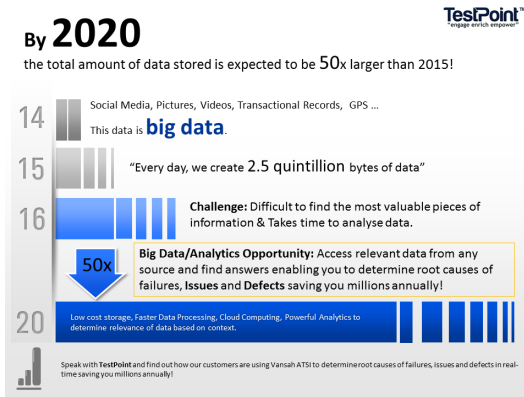
- Language understanding
- Speech recognition

...

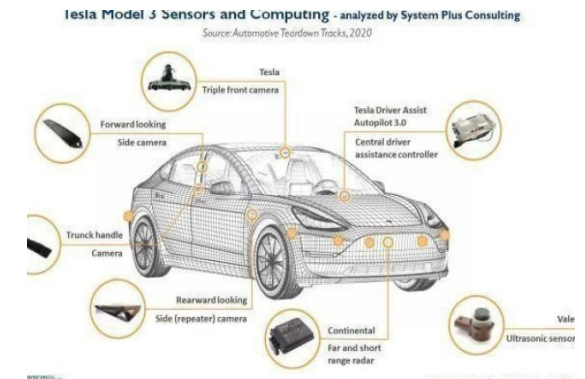
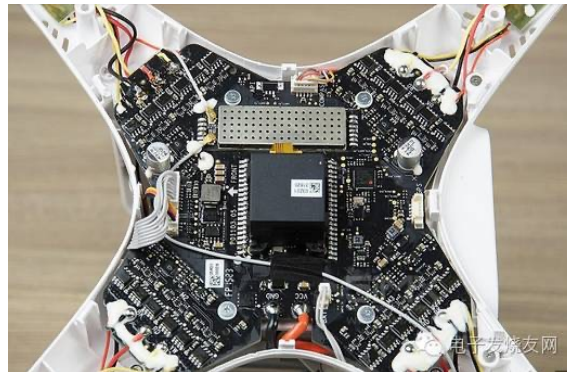


# Background

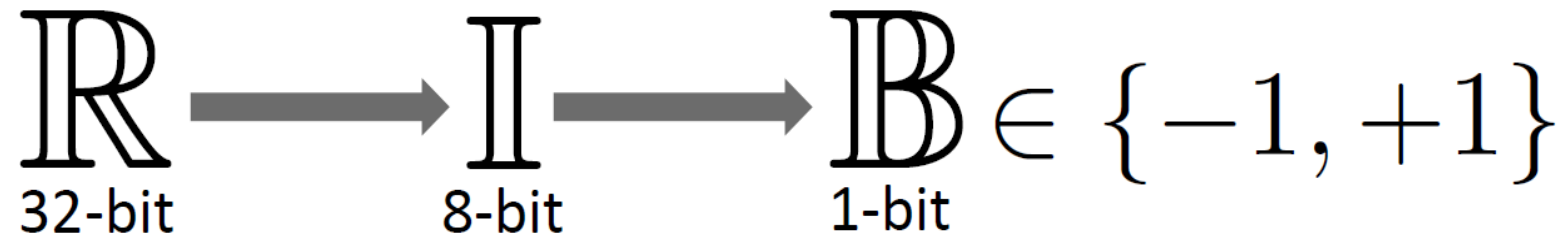
**bigger data**  
and  
**larger model**



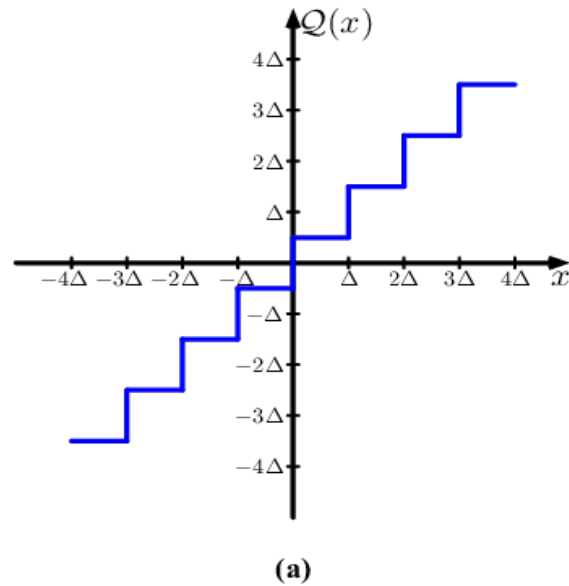
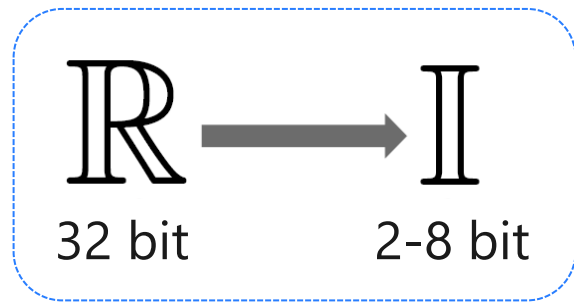
**diverse usage**  
and  
**limited resources**



# Network Quantization and Binarization



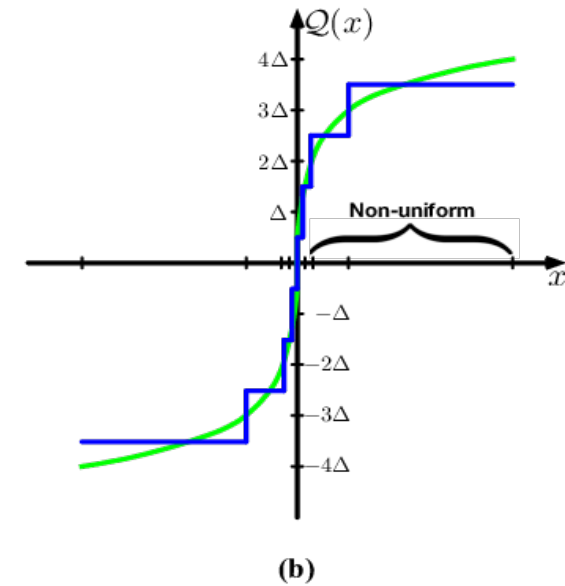
# Network Quantization: 2-8 bit



**Uniform Quantization**

$$Q_U(x) = \text{round}\left(\frac{x}{\Delta}\right) \Delta$$

$$\Delta = \frac{u - l}{2^b - 1}$$

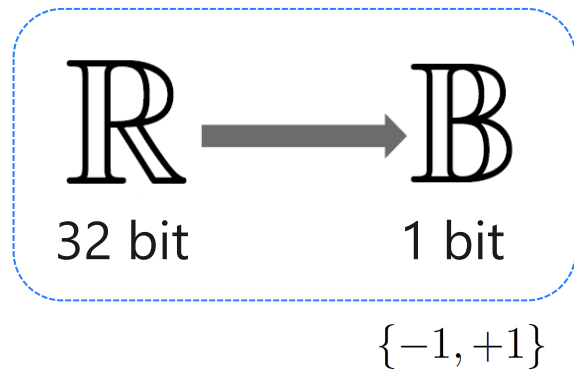


**Non-Uniform Quantization**

$$Q_U(x) = \text{round}\left(\frac{\log_2 x}{\Delta}\right) \Delta$$

$$\Delta = \frac{u - l}{2^b - 1}$$

# Network Binarization: 1-bit



1-Bit Parameters:

$$\mathbf{B}_x = \text{sign}(x) = \begin{cases} +1, & \text{if } x \geq 0 \\ -1, & \text{otherwise.} \end{cases} \quad Q_x(\mathbf{x}) = \alpha \mathbf{B}_x,$$

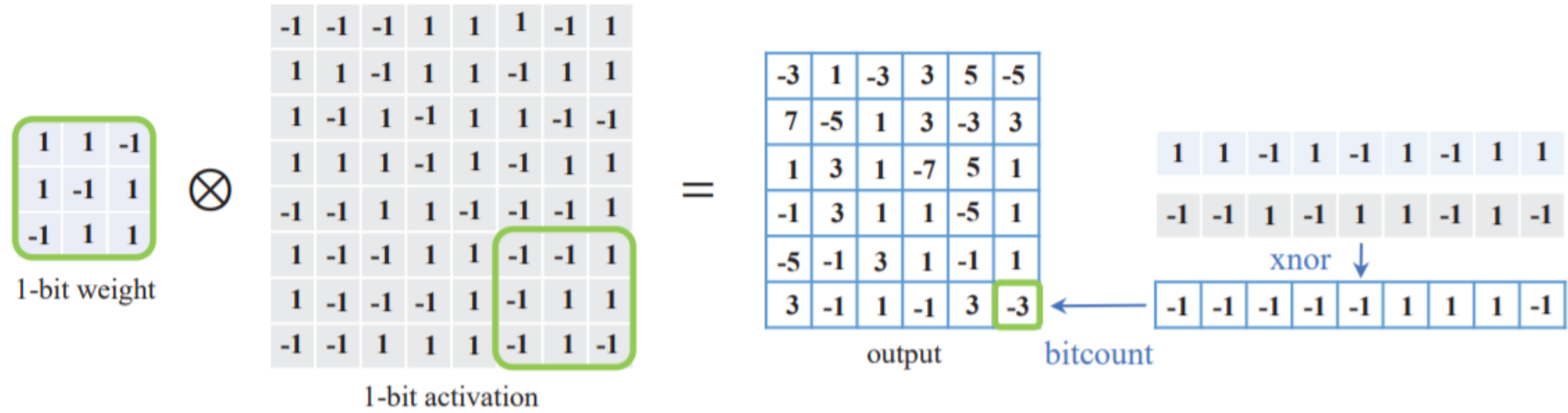
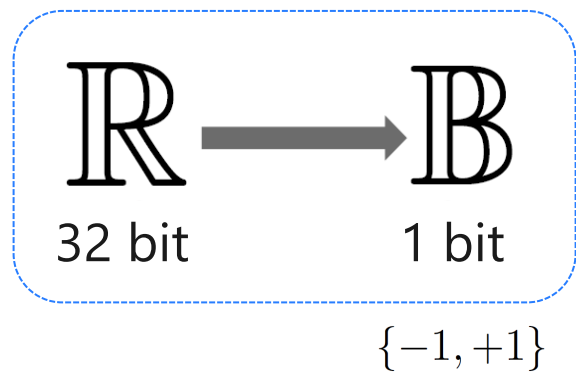


Bitwise Operations:

$$\mathbf{z} = \sigma(Q_w(\mathbf{w}) \otimes Q_a(\mathbf{a})) = \sigma(\alpha\beta(\mathbf{b}_w \odot \mathbf{b}_a))$$



# Network Binarization: 1-bit



# Network Binarization

Full-Precision  
Neural Networks



Massive  
Parameters



Complex  
Computation



High Power  
Consumption



# Network Binarization

Full-Precision  
Neural Networks



Massive  
Parameters



Complex  
Computation



High Power  
Consumption

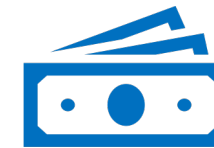
Binarized  
Neural Networks



Binarized  
Parameters



Efficient  
Instructions



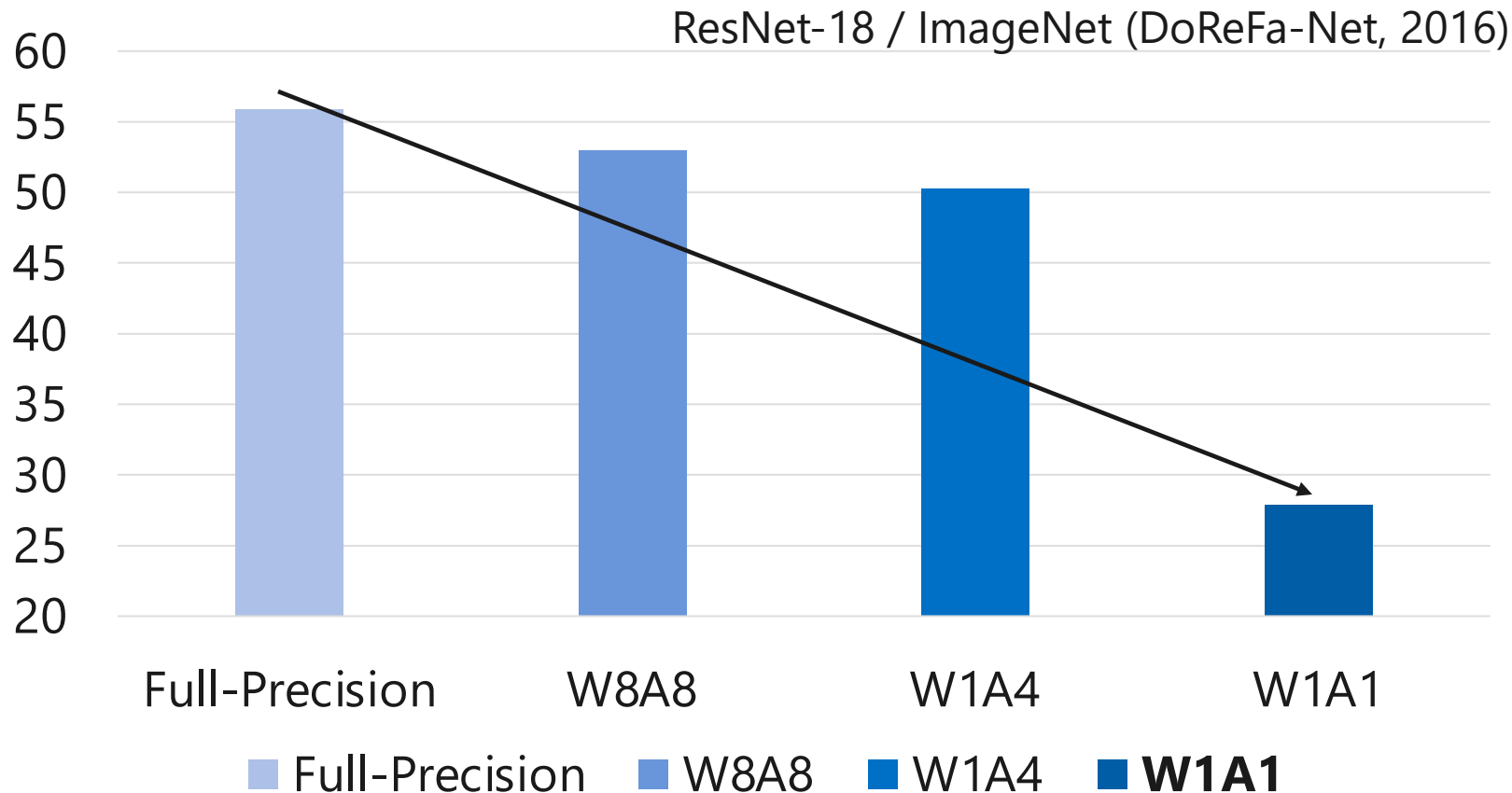
Low Power  
Consumption



# Network Binarization: challenges

**Goal:** accurate extreme-low bit quantization (binarization)

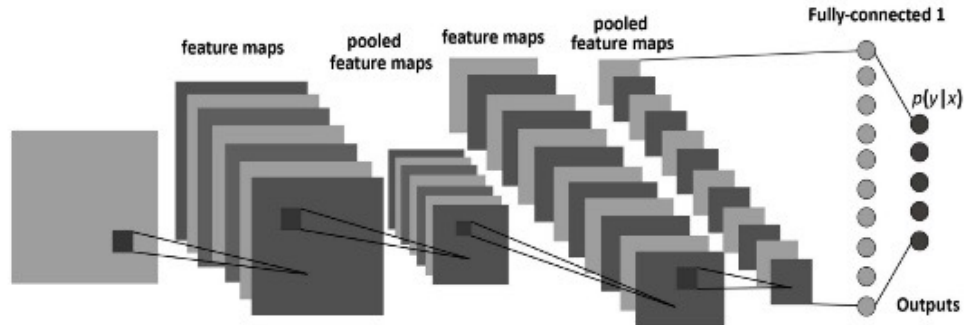
1. the **accuracy** of the binary network has dropped seriously



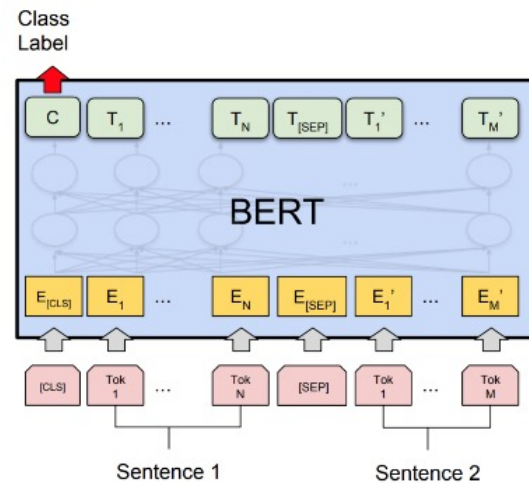
# Network Binarization: challenges

**Goal:** accurate extreme-low bit quantization (binarization)

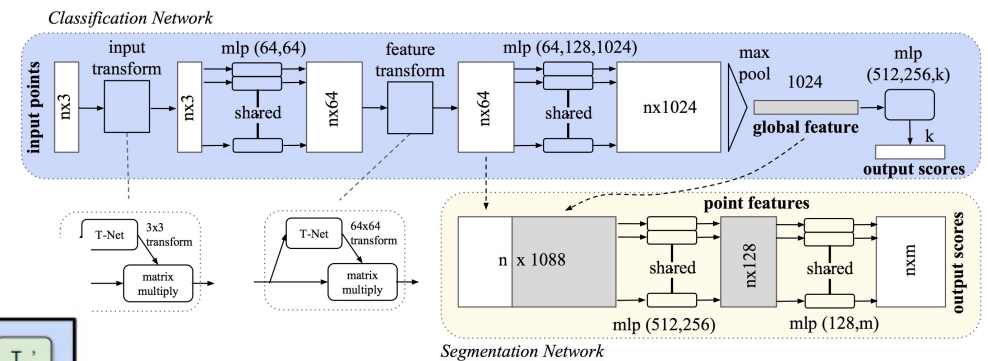
2. binarization methods are not generic across different **architecture**



CNN (ResNet, VGG, ...)



Transformer (BERT, ViT,...)

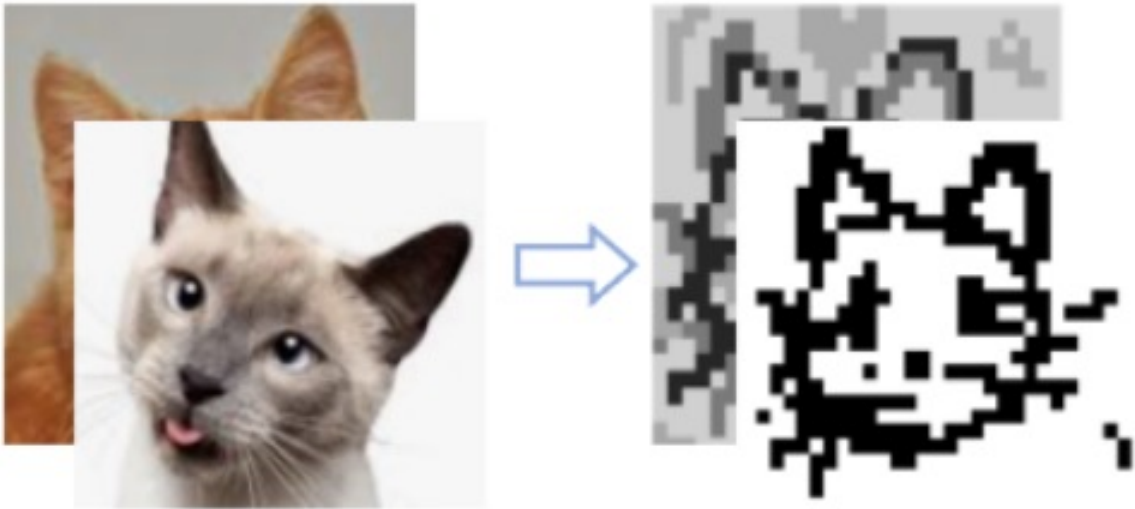


MLP (PointNet, MLP-Mixer, ...)

# CNN Binarization: information loss and retention

## Effects of BNN in the Forward and Backward Propagation

limited representation

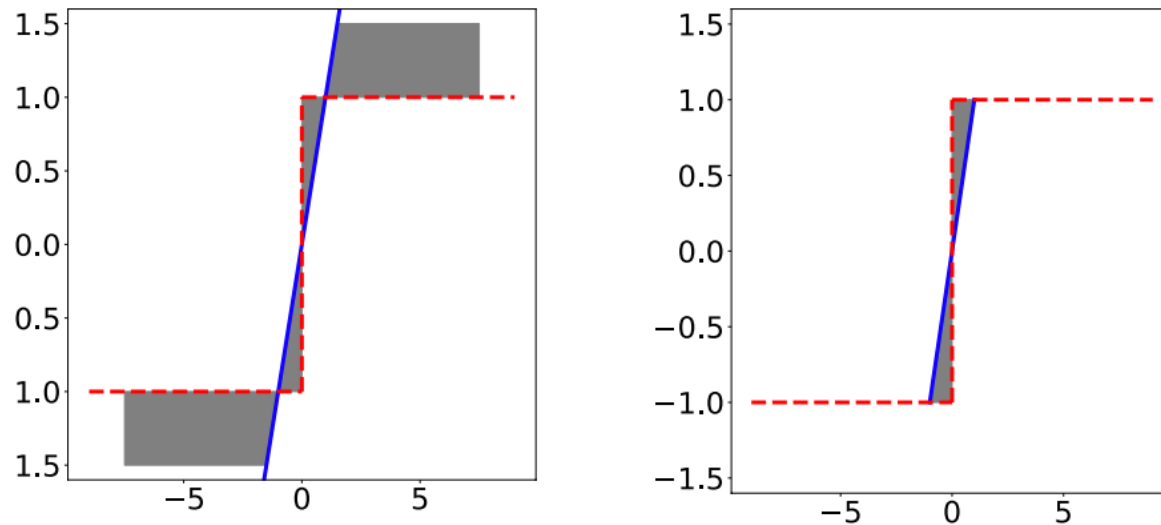


$$\mathbf{B}_x = \text{sign}(\mathbf{x}) = \begin{cases} +1, & \text{if } \mathbf{x} \geq 0 \\ -1, & \text{otherwise.} \end{cases}$$

# CNN Binarization: information loss and retention

## Effects of BNN in the Forward and Backward Propagation

gradient mismatch

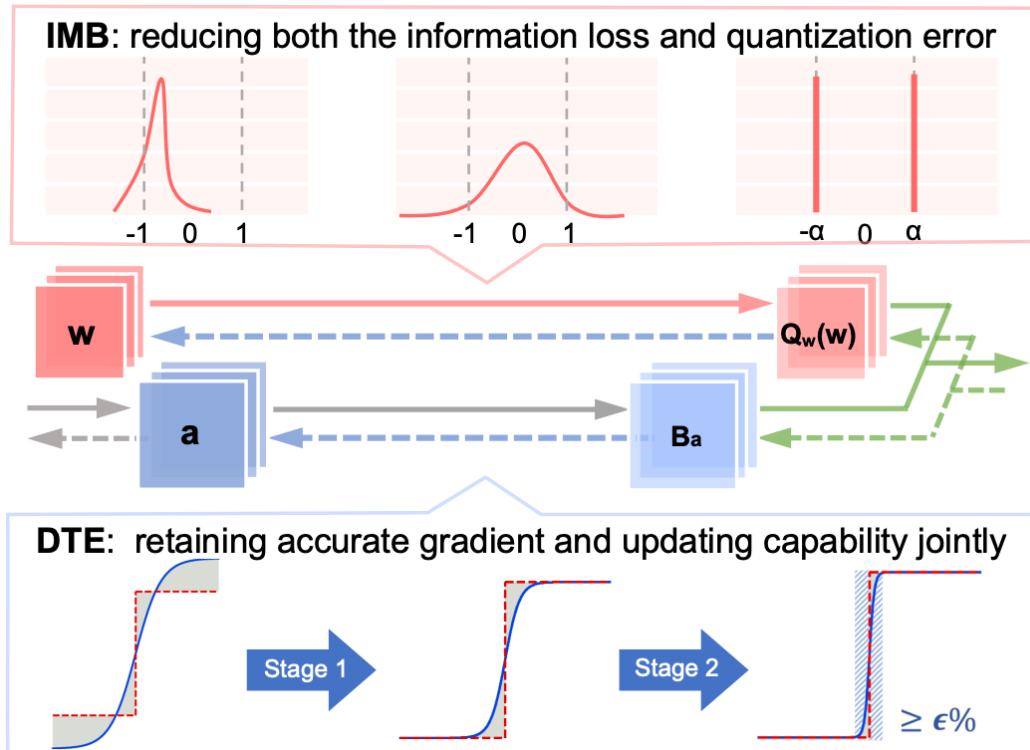


Identity :  $y = x$

Clip :  $y = \text{Hardtanh}(x)$

# CNN Binarization: information loss and retention

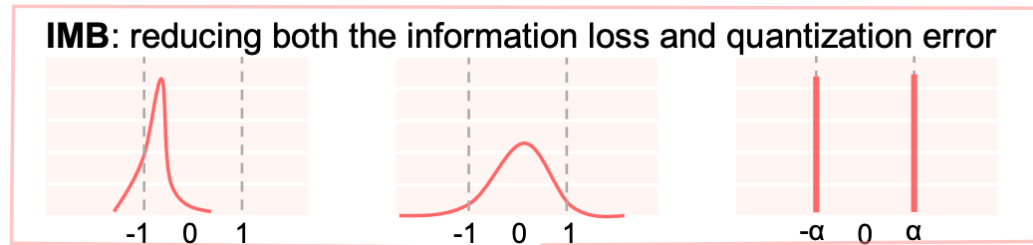
## Distribution-sensitive Information Retention (DIR-Net)



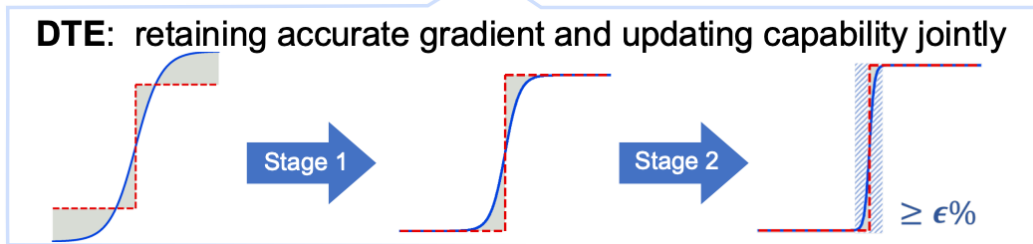
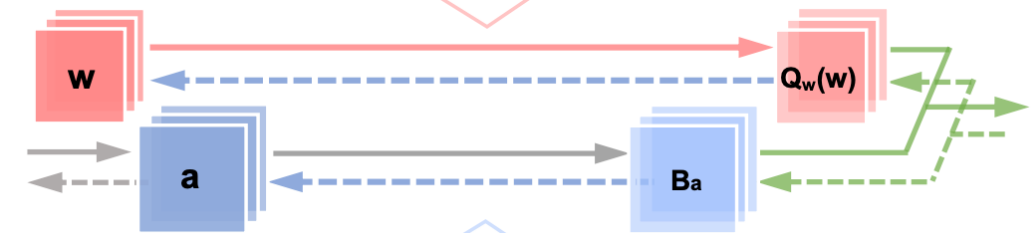
# CNN Binarization: information loss and retention

## Distribution-sensitive Information Retention (DIR-Net)

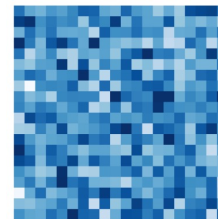
Maximizing the information entropy



Forward



$\mathcal{H} = 0 (min)$



original

**Binarize**



$\mathcal{H} = 0.5$



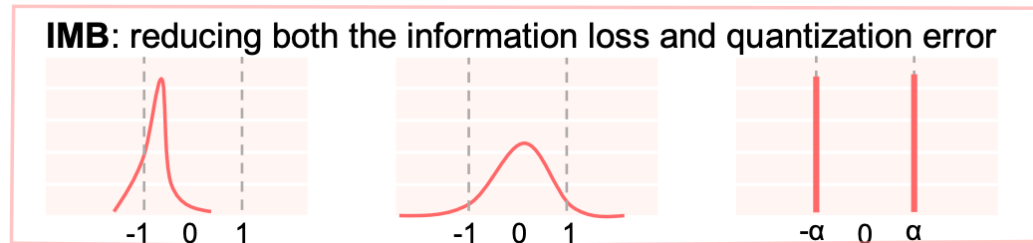
$\mathcal{H} = 0.69 (max)$



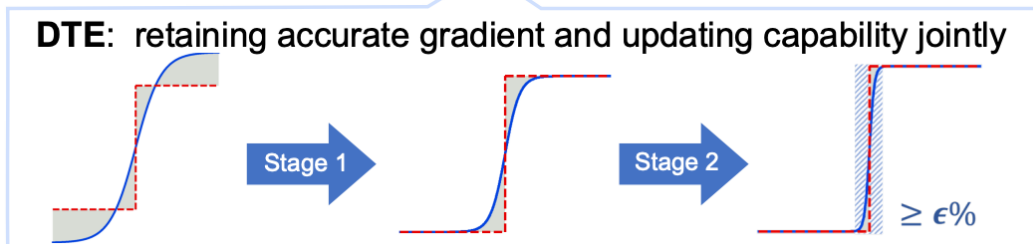
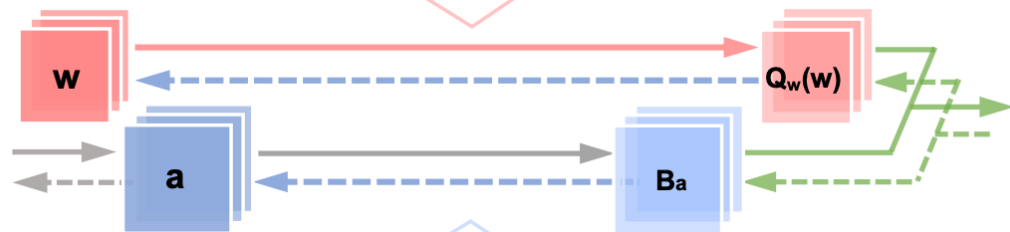
# CNN Binarization: information loss and retention

## Distribution-sensitive Information Retention (DIR-Net)

Maximizing the information entropy

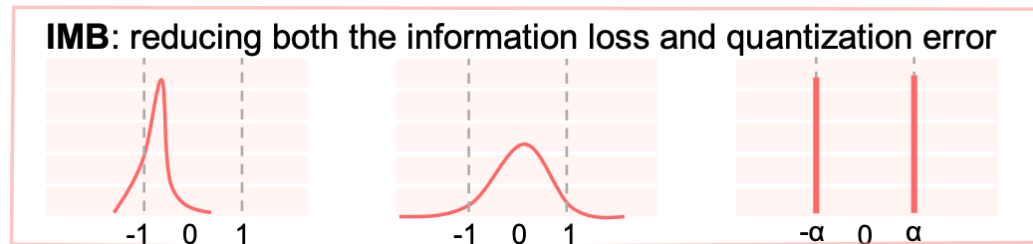


Forward

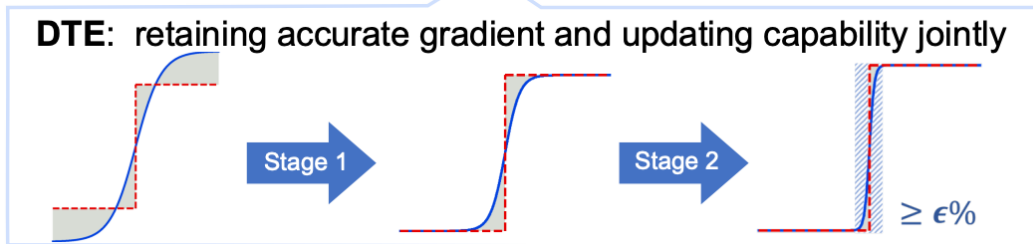
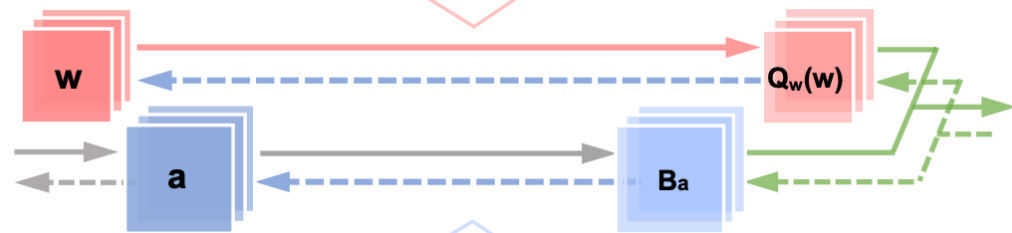


# CNN Binarization: information loss and retention

## Distribution-sensitive Information Retention (DIR-Net)

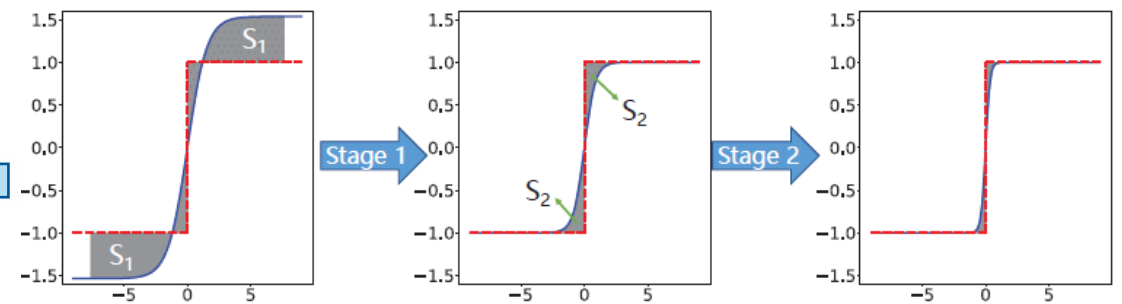


Forward



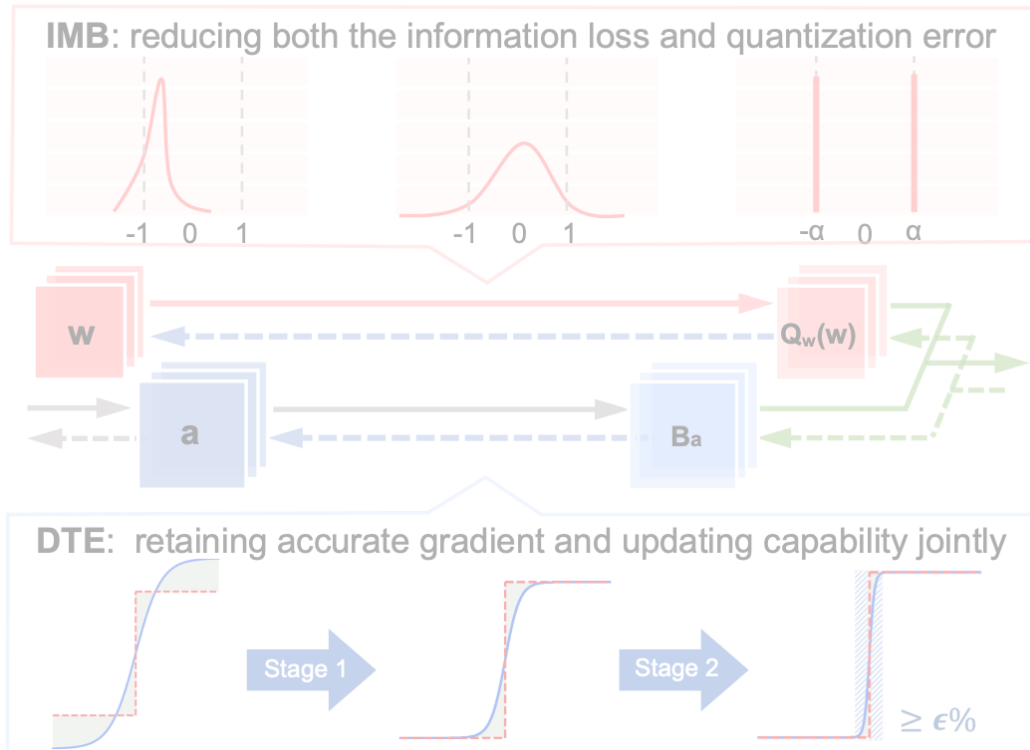
Backward

## Changing the shape of estimator

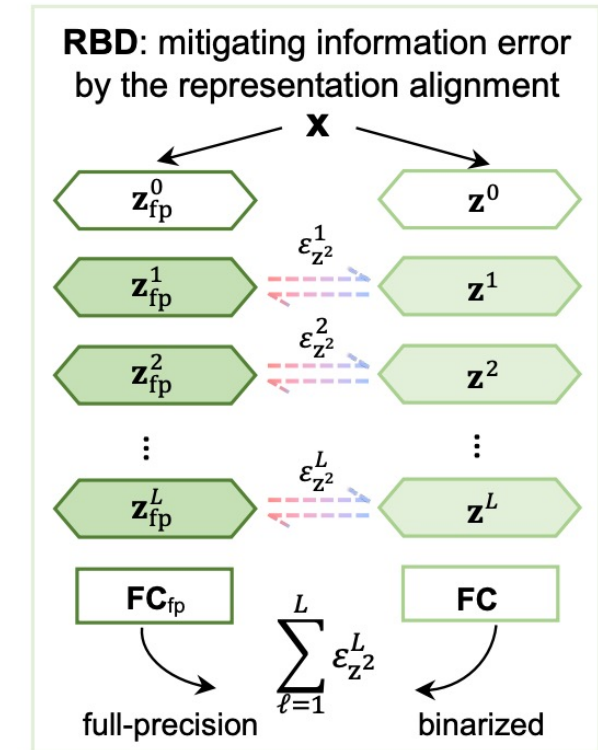


# CNN Binarization: information loss and retention

## Distribution-sensitive Information Retention (DIR-Net)



**external:**  
binarization-aware  
distillation



# CNN Binarization: information loss and retention

## Performance

ResNet-34	Full-Precision	32/32	73.3	91.3
	ABC-Net	1/1	52.4	76.5
	Bi-Real	1/1	62.2	83.9
	IR-Net	1/1	62.9	84.1
	Si-BNN	1/1	63.3	84.4
	ReActNet	1/1	67.3	87.9
	DIR-Net <sup>1</sup> (ours)	1/1	64.1	85.3
	DIR-Net <sup>2</sup> (ours)	1/1	<b>67.9<math>\pm</math>0.09</b>	<b>88.2</b>



The accuracy reached **90%** of the full precision ResNet

Full-Precision	32/32	73.3	91.3
ABC-Net	1/32	68.8	86.1
Bi-Real	1/32	69.7	88.9
Si-BNN	1/32	70.1	89.7
IR-Net	1/32	70.4	89.5
DIR-Net (ours)	1/32	<b>71.1<math>\pm</math>0.03</b>	<b>90.4</b>

DARTS	Full-Precision	32/32	73.3	91.3
	BNN	1/1	52.2	76.6
	Bi-Real	1/1	61.5	83.8
	IR-Net	1/1	62.1	84.2
	ReActNet	1/1	65.1	86.4
	DIR-Net <sup>1</sup> (ours)	1/1	63.3	85.1
	DIR-Net <sup>2</sup> (ours)	1/1	<b>65.6<math>\pm</math>0.12</b>	<b>87.2</b>

# CNN Binarization: information loss and retention

## Performance

ResNet-34	Full-Precision	32/32	73.3	91.3
	ABC-Net	1/1	52.4	76.5
	Bi-Real	1/1	62.2	83.9
	IR-Net	1/1	62.9	84.1
	Si-BNN	1/1	63.3	84.4
	ReActNet	1/1	67.3	87.9
	DIR-Net <sup>1</sup> (ours)	1/1	64.1	85.3
	DIR-Net <sup>2</sup> (ours)	1/1	<b>67.9<math>\pm</math>0.09</b>	<b>88.2</b>



The accuracy reached **90%** of the full precision ResNet

	Full-Precision	32/32	73.3	91.3
	ABC-Net	1/32	68.8	86.1
	Bi-Real	1/32	69.7	88.9
	Si-BNN	1/32	70.1	89.7
	IR-Net	1/32	70.4	89.5
	DIR-Net (ours)	1/32	<b>71.1<math>\pm</math>0.03</b>	<b>90.4</b>



Accurate binarization on **lightweight** architectures

DARTS	Full-Precision	32/32	73.3	91.3
	BNN	1/1	52.2	76.6
	Bi-Real	1/1	61.5	83.8
	IR-Net	1/1	62.1	84.2
	ReActNet	1/1	65.1	86.4
	DIR-Net <sup>1</sup> (ours)	1/1	63.3	85.1
	DIR-Net <sup>2</sup> (ours)	1/1	<b>65.6<math>\pm</math>0.12</b>	<b>87.2</b>

# CNN Binarization: information loss and retention

## Performance

ResNet-34	Full-Precision	32/32	73.3	91.3
	ABC-Net	1/1	52.4	76.5
	Bi-Real	1/1	62.2	83.9
	IR-Net	1/1	62.9	84.1
	Si-BNN	1/1	63.3	84.4
	ReActNet	1/1	67.3	87.9
	DIR-Net <sup>1</sup> (ours)	1/1	64.1	85.3
	DIR-Net <sup>2</sup> (ours)	1/1	<b>67.9<math>\pm</math>0.09</b>	<b>88.2</b>

ResNet-34	Full-Precision	32/32	73.3	91.3
	ABC-Net	1/32	68.8	86.1
	Bi-Real	1/32	69.7	88.9
	Si-BNN	1/32	70.1	89.7
	IR-Net	1/32	70.4	89.5
	DIR-Net (ours)	1/32	<b>71.1<math>\pm</math>0.03</b>	<b>90.4</b>

DARTS	Full-Precision	32/32	73.3	91.3
	BNN	1/1	52.2	76.6
	Bi-Real	1/1	61.5	83.8
	IR-Net	1/1	62.1	84.2
	ReActNet	1/1	65.1	86.4
	DIR-Net <sup>1</sup> (ours)	1/1	63.3	85.1
	DIR-Net <sup>2</sup> (ours)	1/1	<b>65.6<math>\pm</math>0.12</b>	<b>87.2</b>

Table 5: Comparison of time cost of ResNet-18 with different bits (single thread).

Method	Bit-width (W/A)	Size (Mb)	Time (ms)
FP	32/32	46.77	1418.94
NCNN	8/8	–	935.51
DSQ	2/2	–	551.22
Ours (without bit-shift scales)	1/1	<b>4.20</b>	<b>252.16</b>
Ours	1/1	<b>4.21</b>	<b>261.98</b>

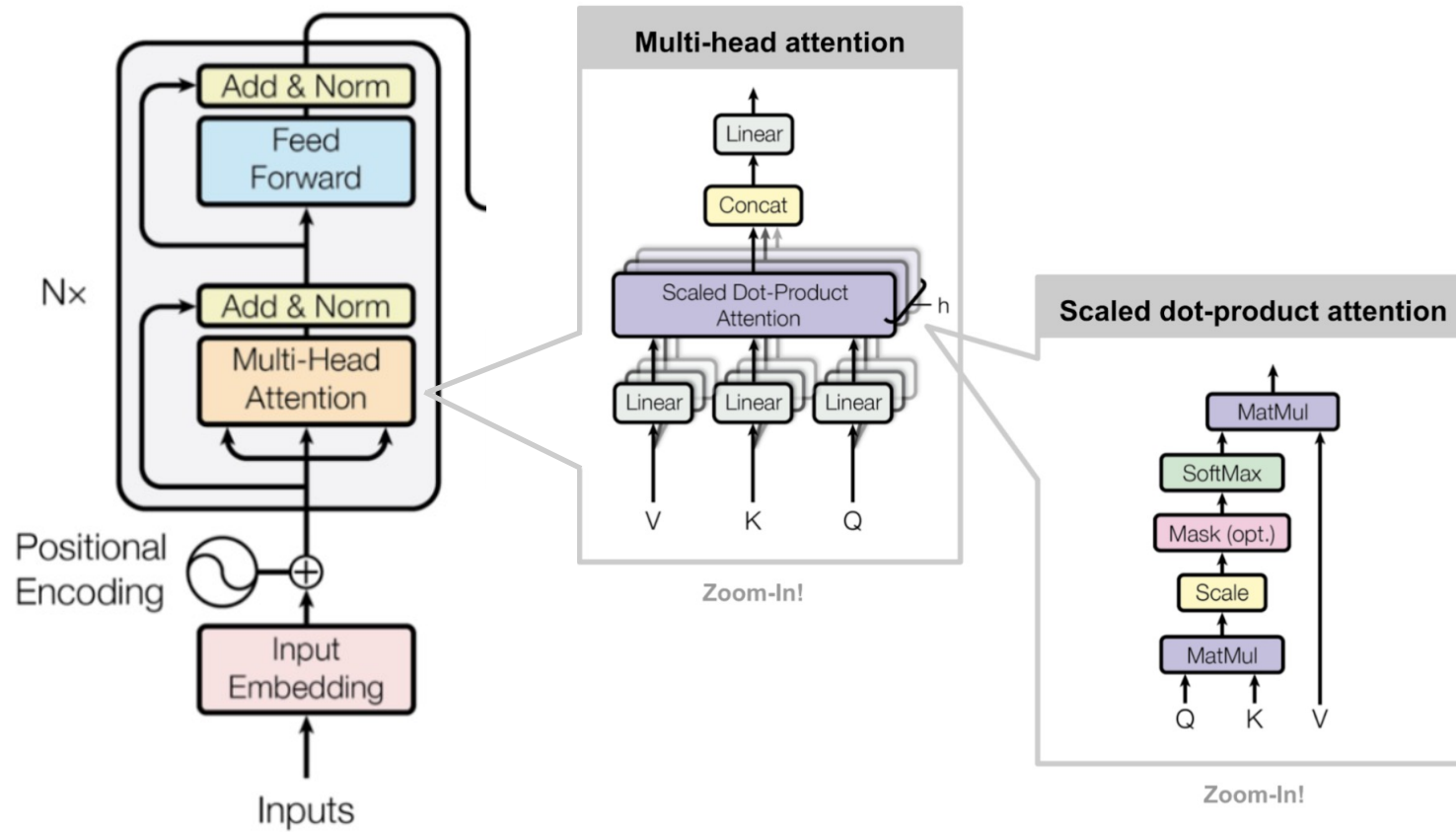


**11.1x** storage saving

**5.4x** speedup

# Transformer Binarization: attention crash and recovery

## Bottlenecks of Fully Binarized BERT Baseline

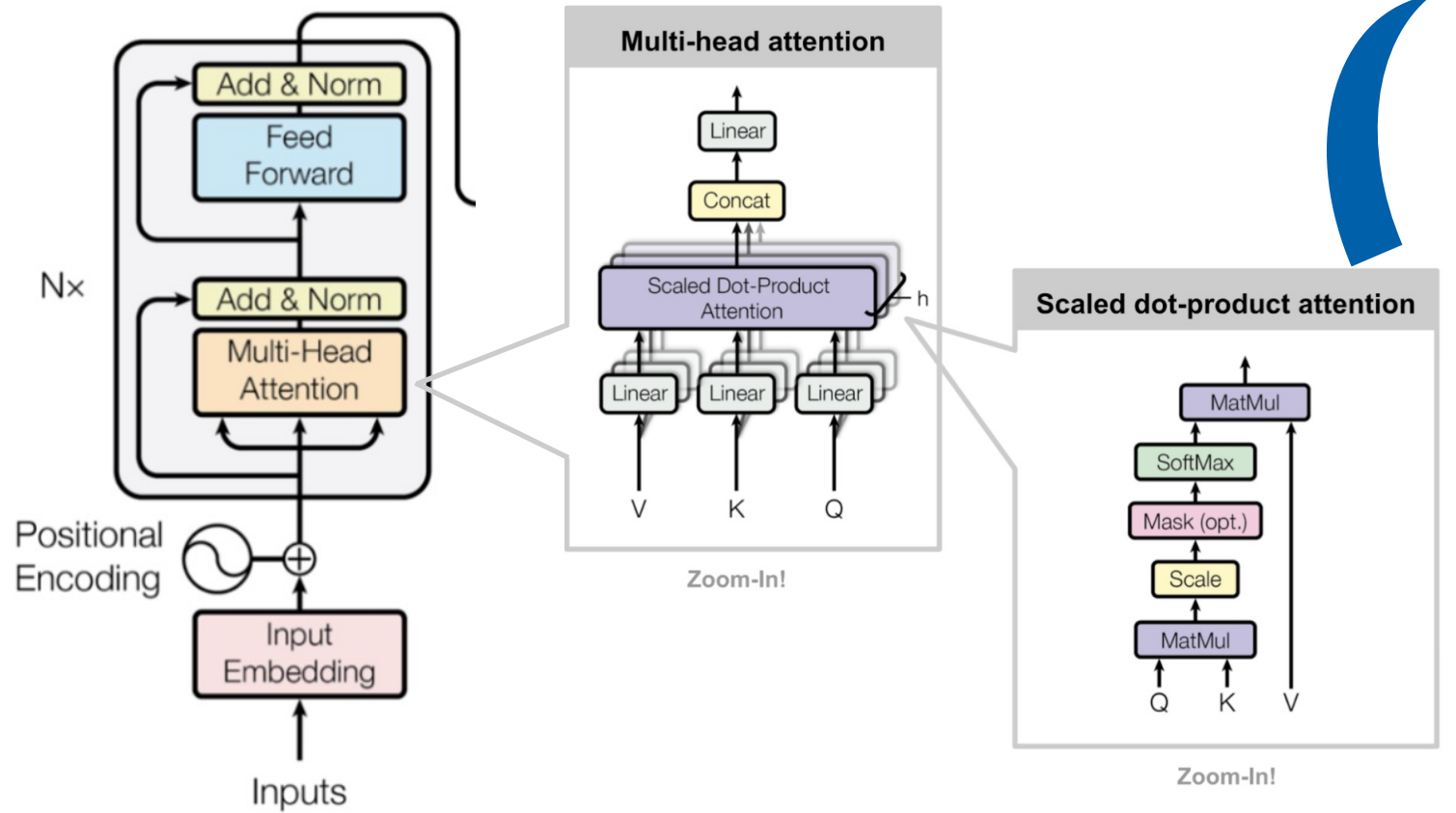


<https://deepfrench.gitlab.io/deep-learning-project/>



# Transformer Binarization: attention crash and recovery

## Bottlenecks of Fully Binarized BERT Baseline



**Binarize  
(directly)**

$$\mathbf{B}_Q = \text{sign}(\mathbf{Q}), \quad \mathbf{B}_K = \text{sign}(\mathbf{K})$$

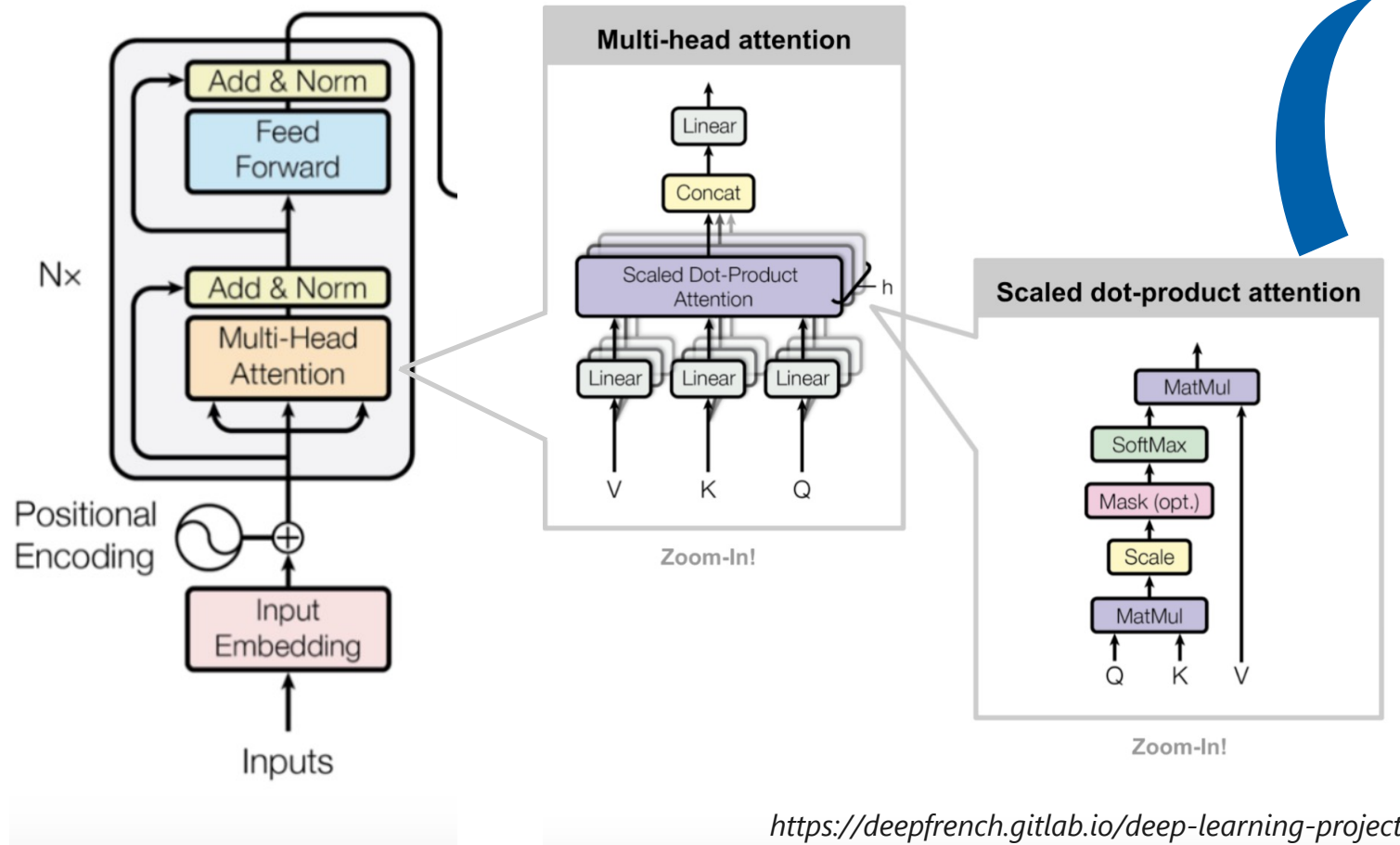
$$\mathbf{A} = \frac{1}{\sqrt{D}} \left( \mathbf{B}_Q \otimes \mathbf{B}_K^T \right)$$

$$\mathbf{B}_A^s = \text{sign}(\text{softmax}(\mathbf{A}))$$

<https://deepfrench.gitlab.io/deep-learning-project/>

# Transformer Binarization: attention crash and recovery

## Bottlenecks of Fully Binarized BERT Baseline



**Binarize  
(directly)**

$$\mathbf{B}_Q = \text{sign}(\mathbf{Q}), \quad \mathbf{B}_K = \text{sign}(\mathbf{K})$$

$$\mathbf{A} = \frac{1}{\sqrt{D}} \left( \mathbf{B}_Q \otimes \mathbf{B}_K^T \right)$$

$$\mathbf{B}_A^s = \text{sign}(\text{softmax}(\mathbf{A}))$$

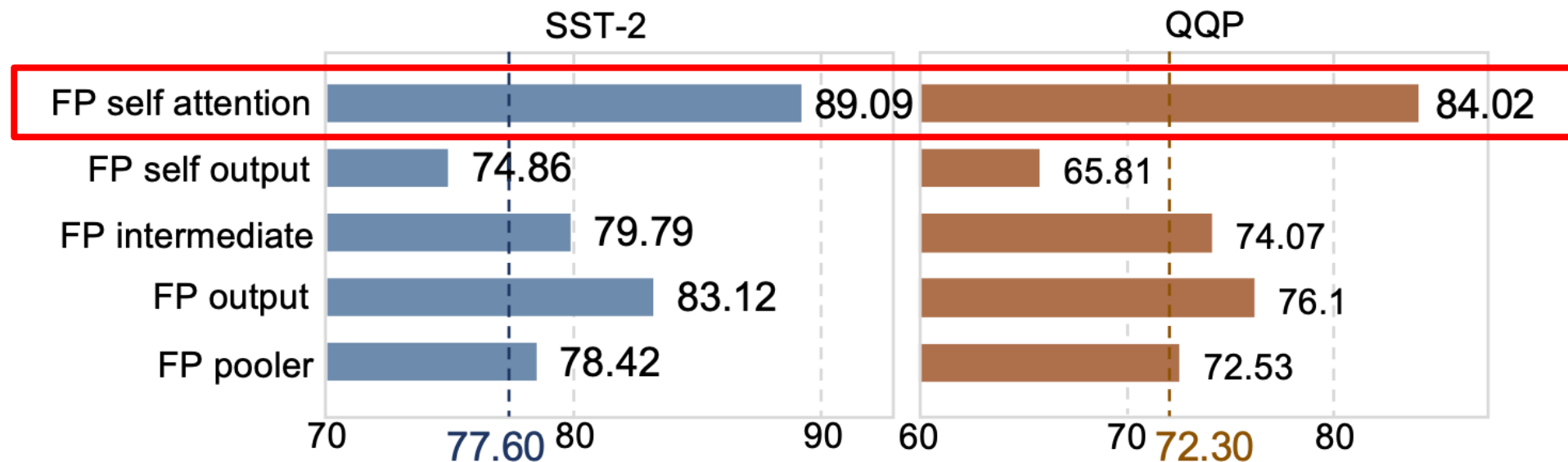
**Severely dropped!**  
(Avg: 83.9% -> **50.4%**)

<https://deepfrench.gitlab.io/deep-learning-project/>

# Transformer Binarization: attention crash and recovery

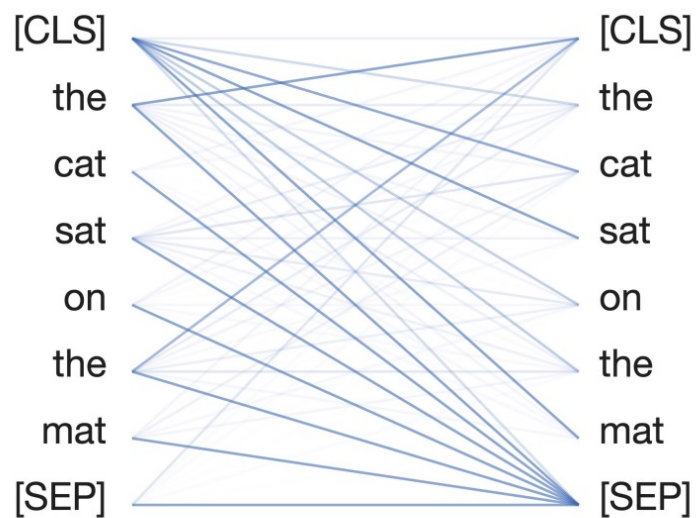
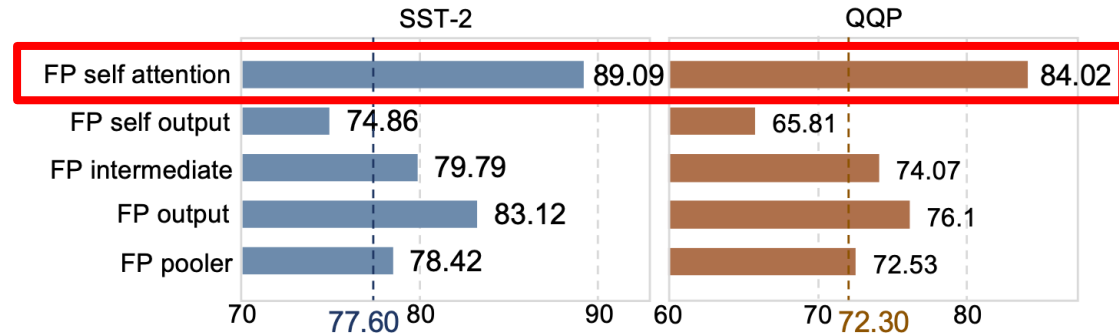
## Bottlenecks of Fully Binarized BERT Baseline

Which part caused the **biggest drop**?

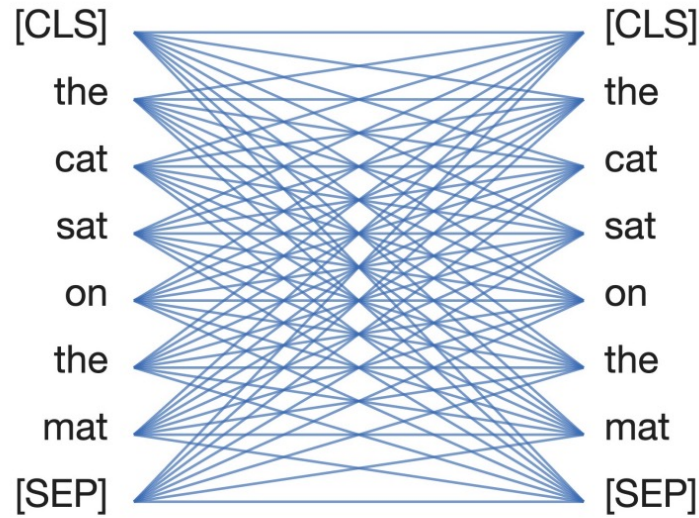


# Transformer Binarization: attention crash and recovery

## Bottlenecks of Fully Binarized BERT Baseline



(a) Full-precision

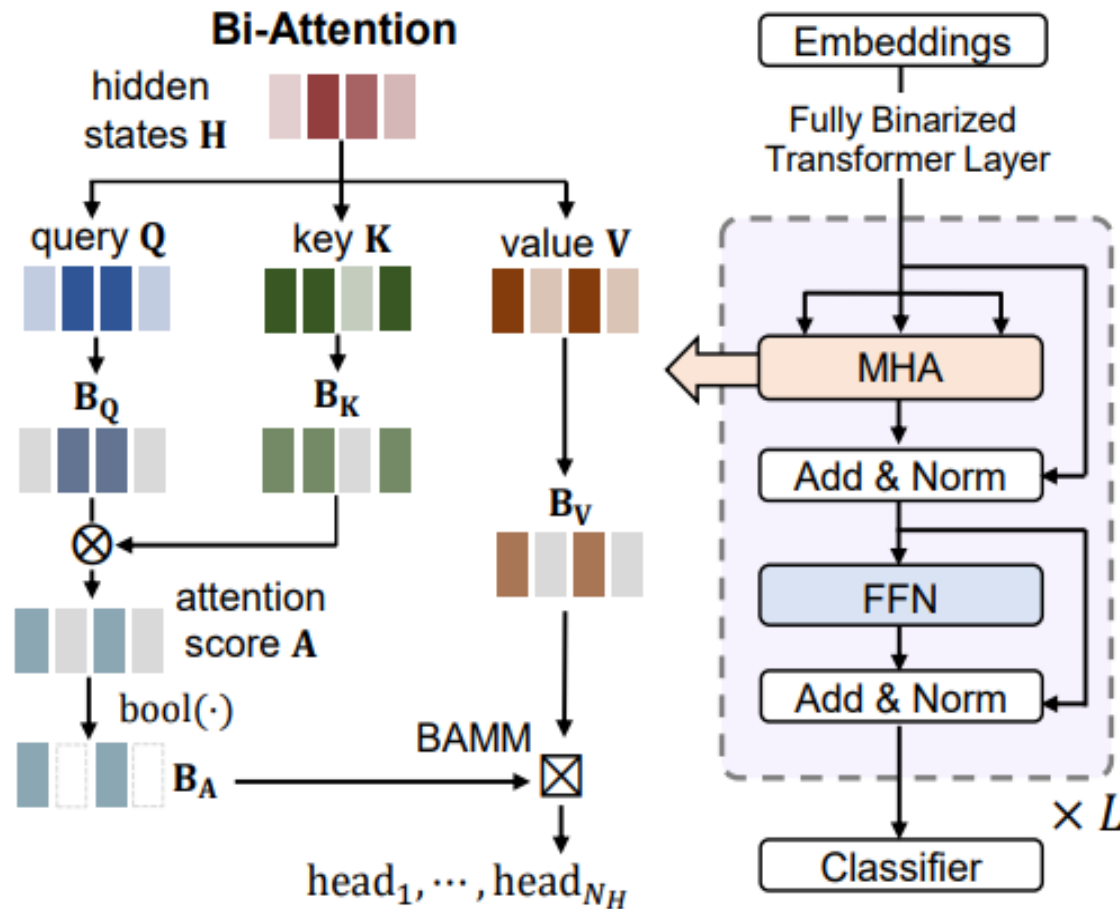


(b) Fully binarized BERT baseline

**attention  
mechanism  
crashed**

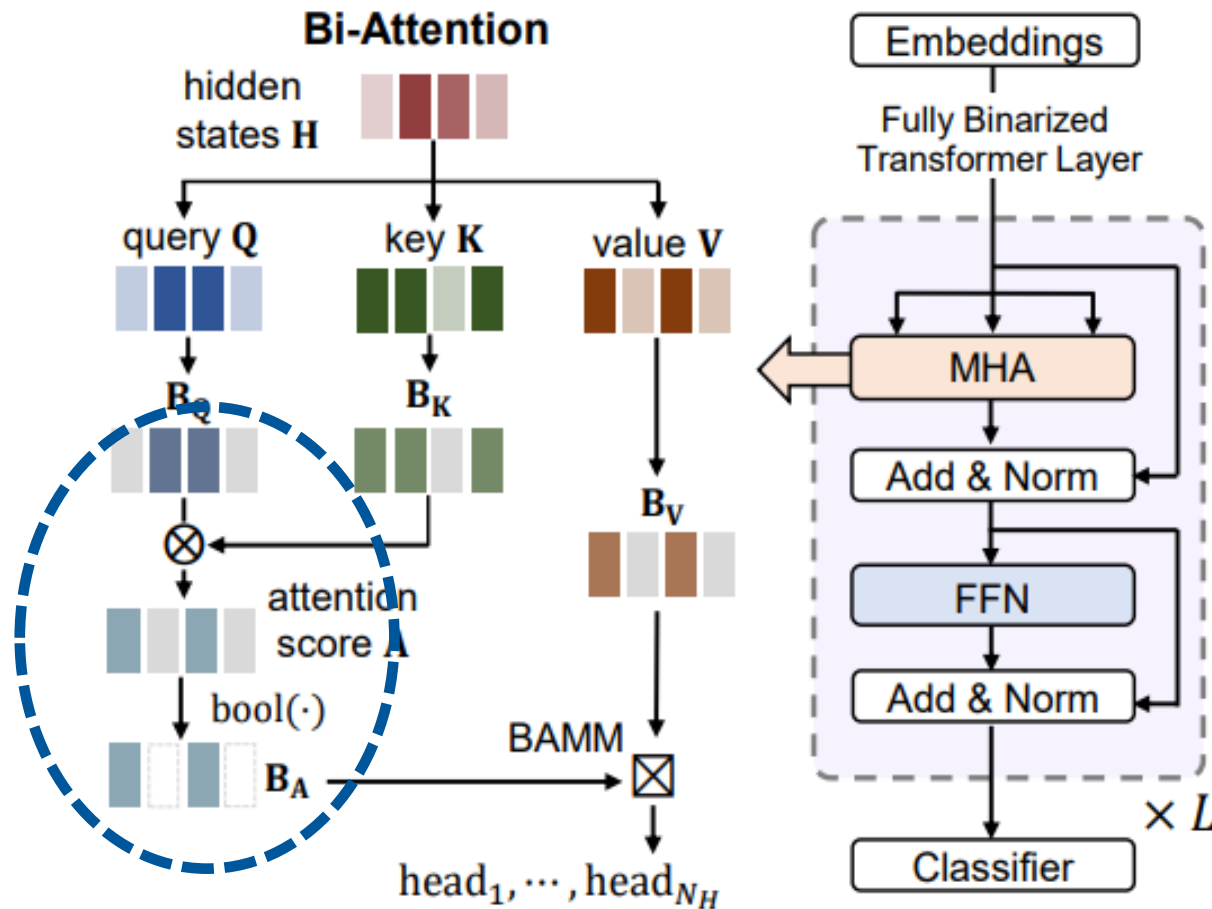
# Transformer Binarization: attention crash and recovery

## Accurate Fully Binarized BERT (BiBERT)



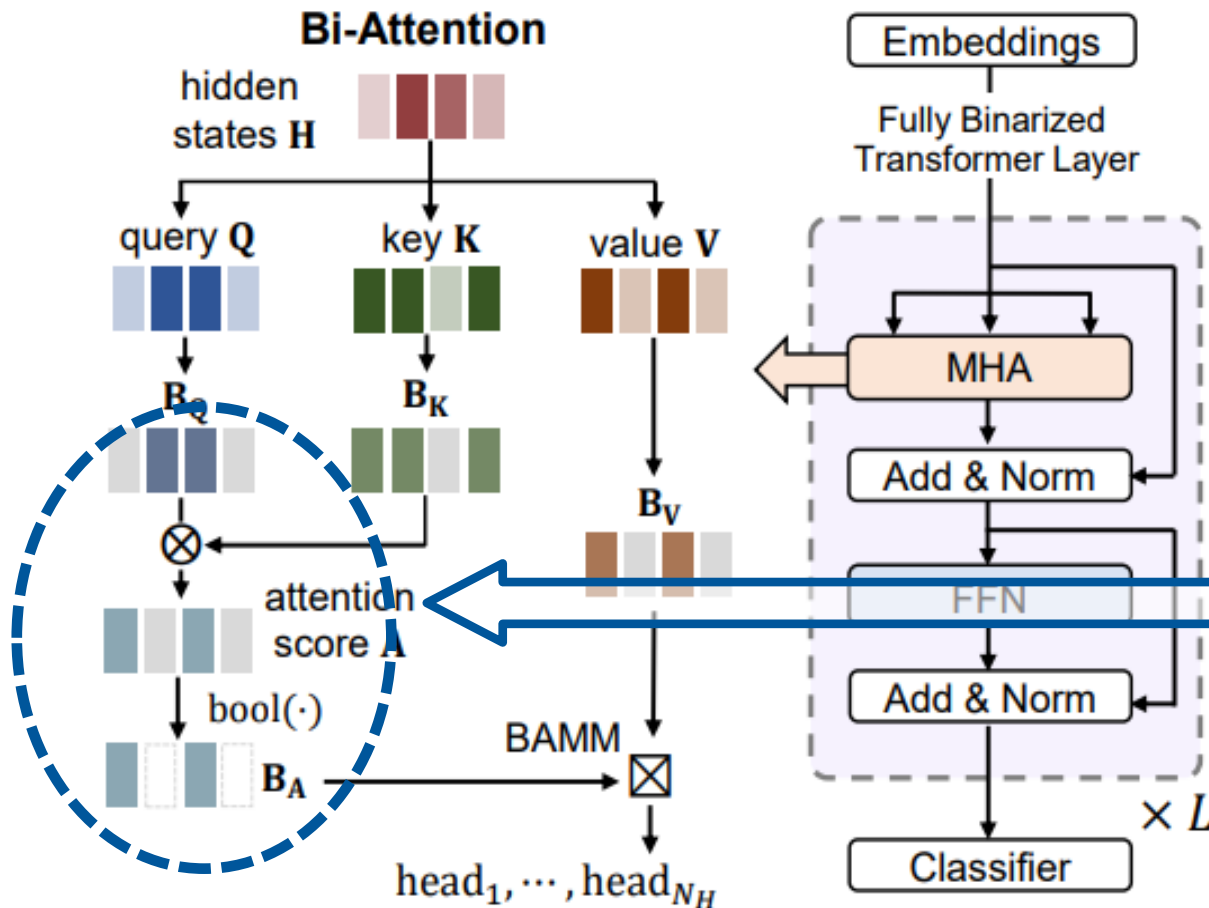
# Transformer Binarization: attention crash and recovery

## Accurate Fully Binarized BERT (BiBERT)



# Transformer Binarization: attention crash and recovery

## Accurate Fully Binarized BERT (BiBERT)



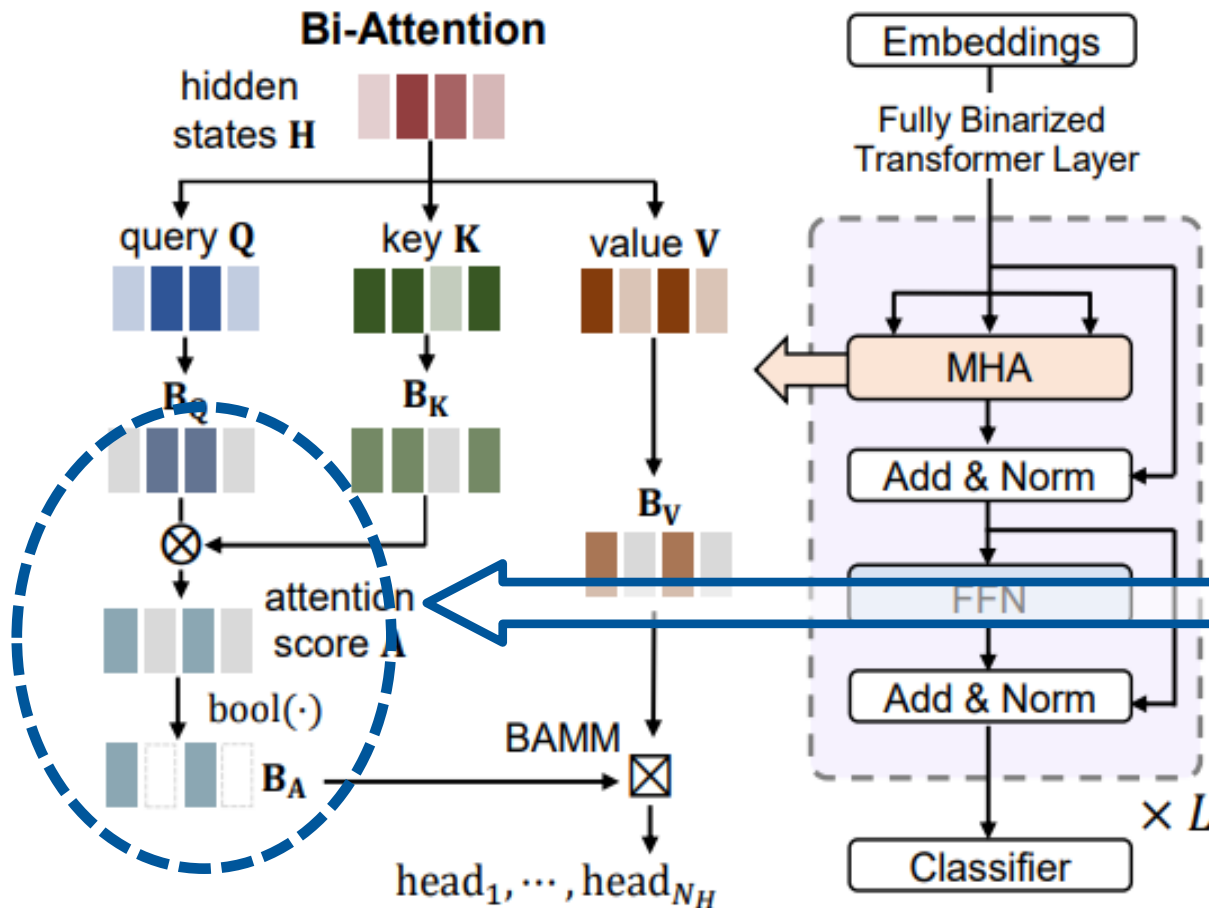
$$\text{bool}(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\frac{\partial \text{bool}(x)}{\partial x} = \begin{cases} 1, & \text{if } |x| \leq 1 \\ 0, & \text{otherwise} \end{cases}$$



# Transformer Binarization: attention crash and recovery

## Accurate Fully Binarized BERT (BiBERT)



### 1. ~~SoftMax~~

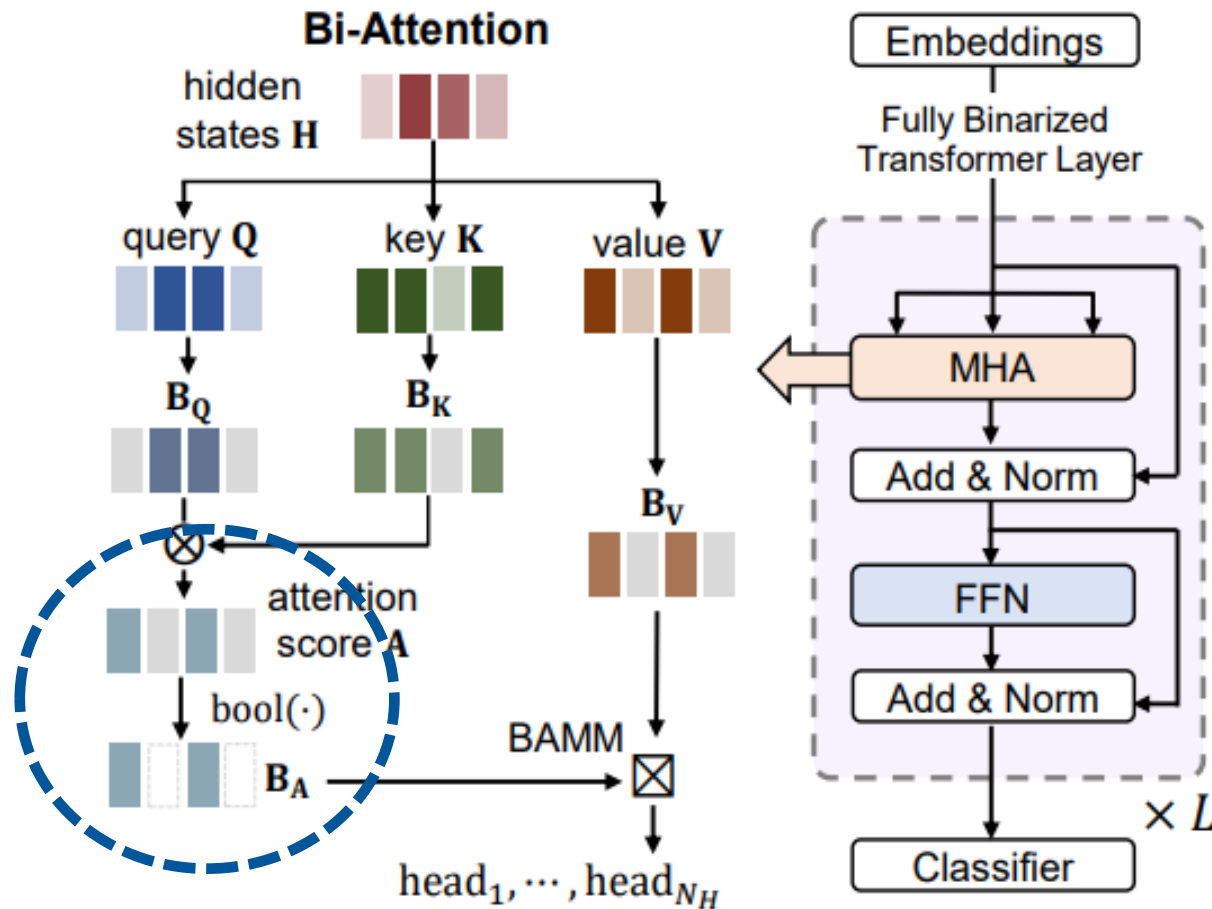


$$\text{bool}(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

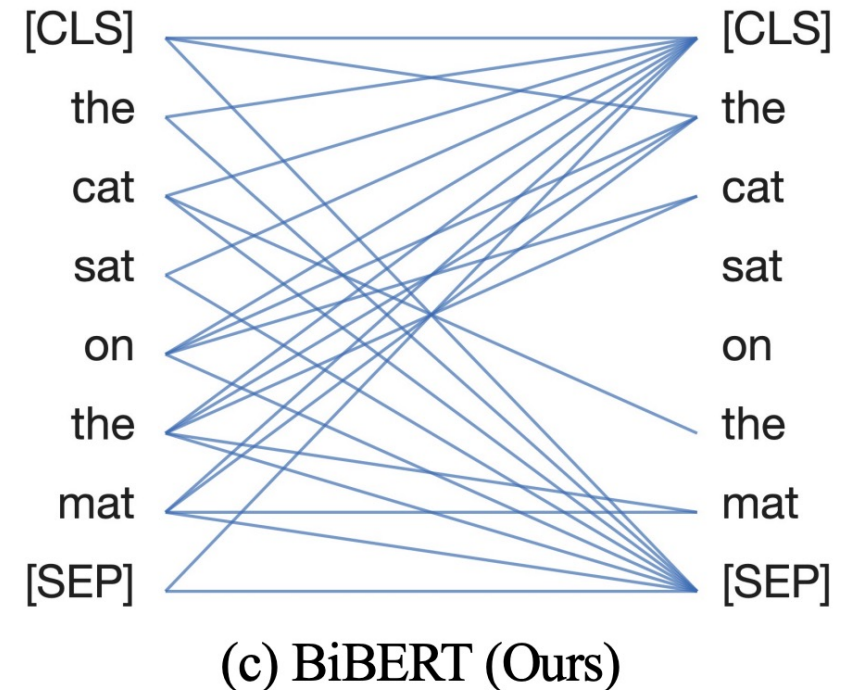
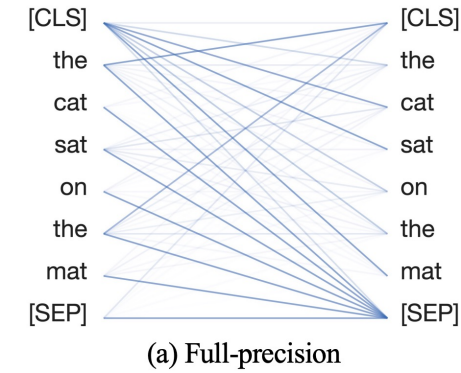
$$\frac{\partial \text{bool}(x)}{\partial x} = \begin{cases} 1, & \text{if } |x| \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

# Transformer Binarization: attention crash and recovery

## Accurate Fully Binarized BERT (BiBERT)

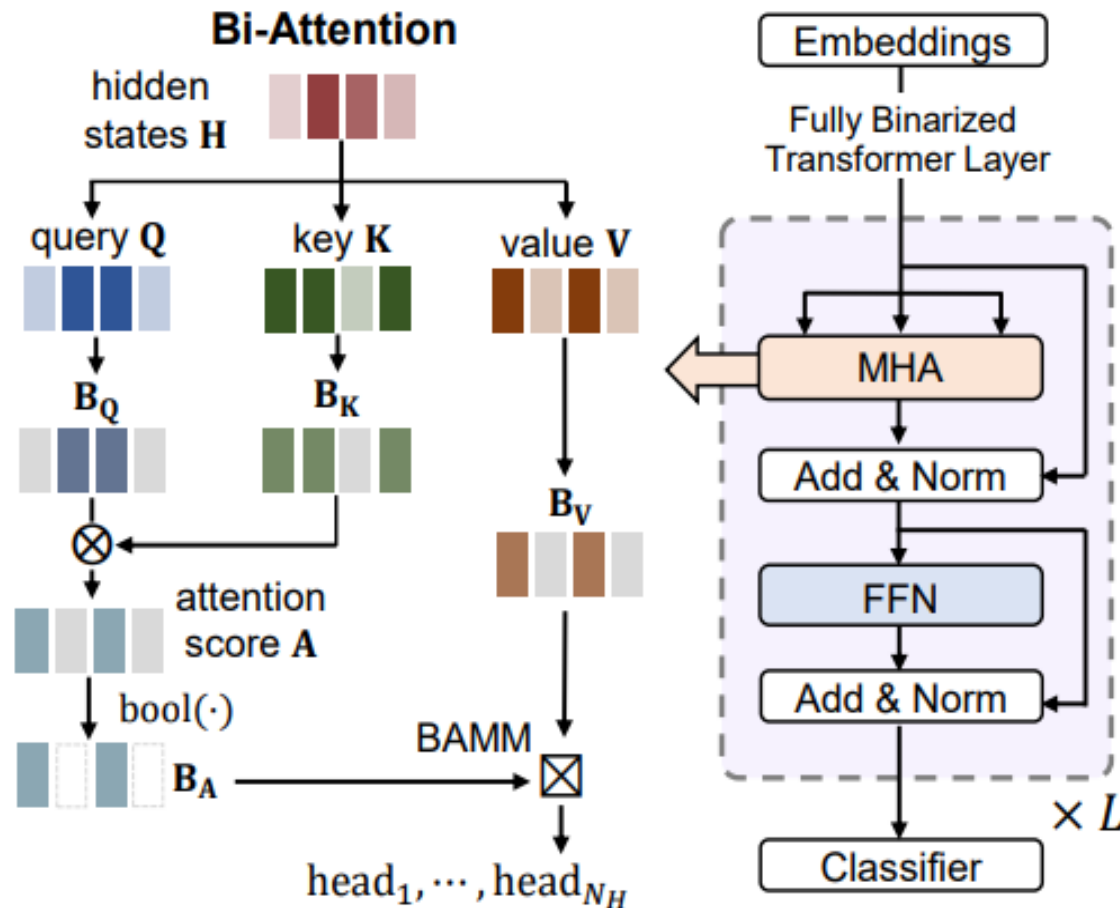


2.



# Transformer Binarization: attention crash and recovery

## Accurate Fully Binarized BERT (BiBERT)

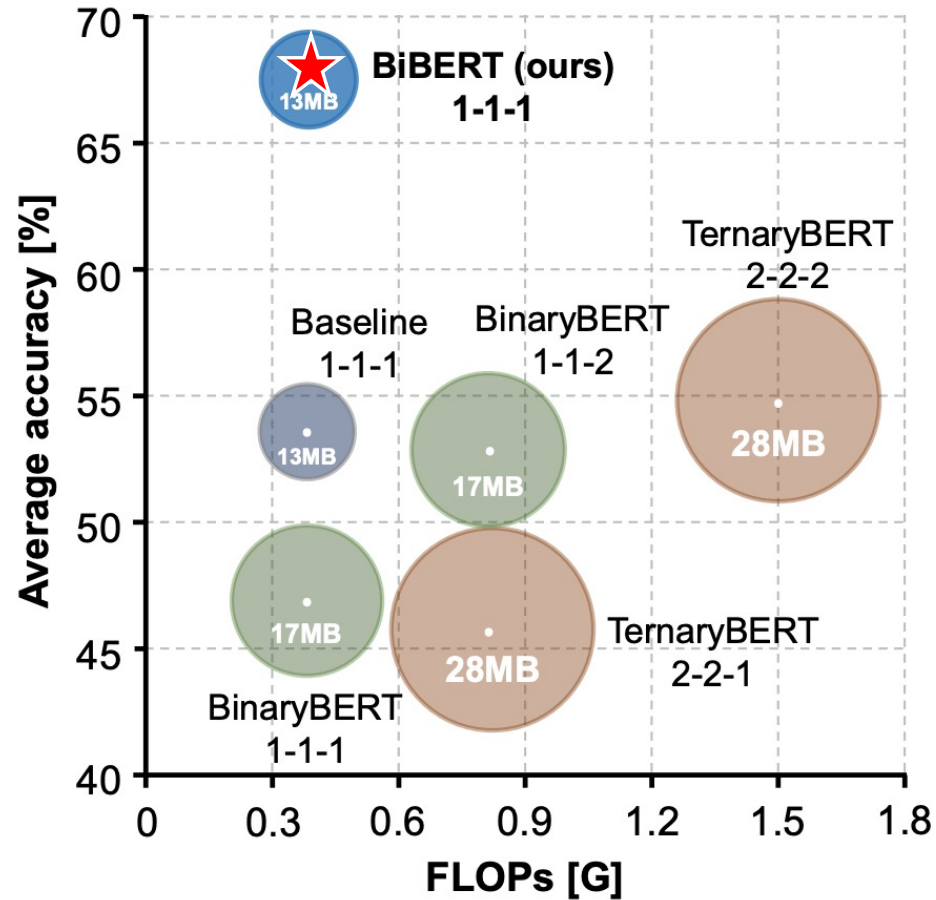


$$\mathbf{B}_A = \text{bool}(\mathbf{A}) = \text{bool}\left(\frac{1}{\sqrt{D}}\left(\mathbf{B}_Q \otimes \mathbf{B}_K^T\right)\right)$$

$$\text{Bi-Attention}(\mathbf{B}_Q, \mathbf{B}_K, \mathbf{B}_V) = \mathbf{B}_A \boxtimes \mathbf{B}_V$$

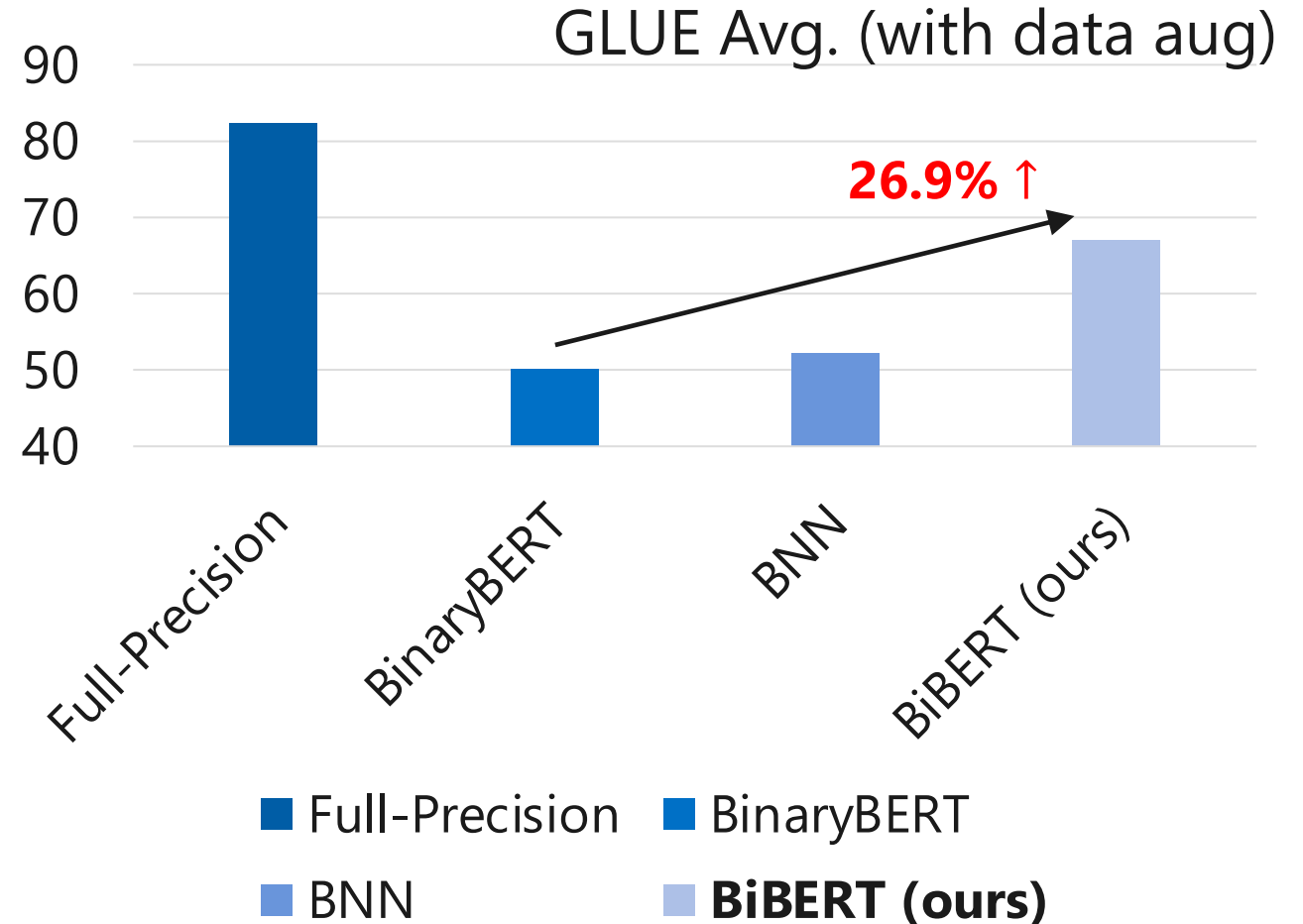
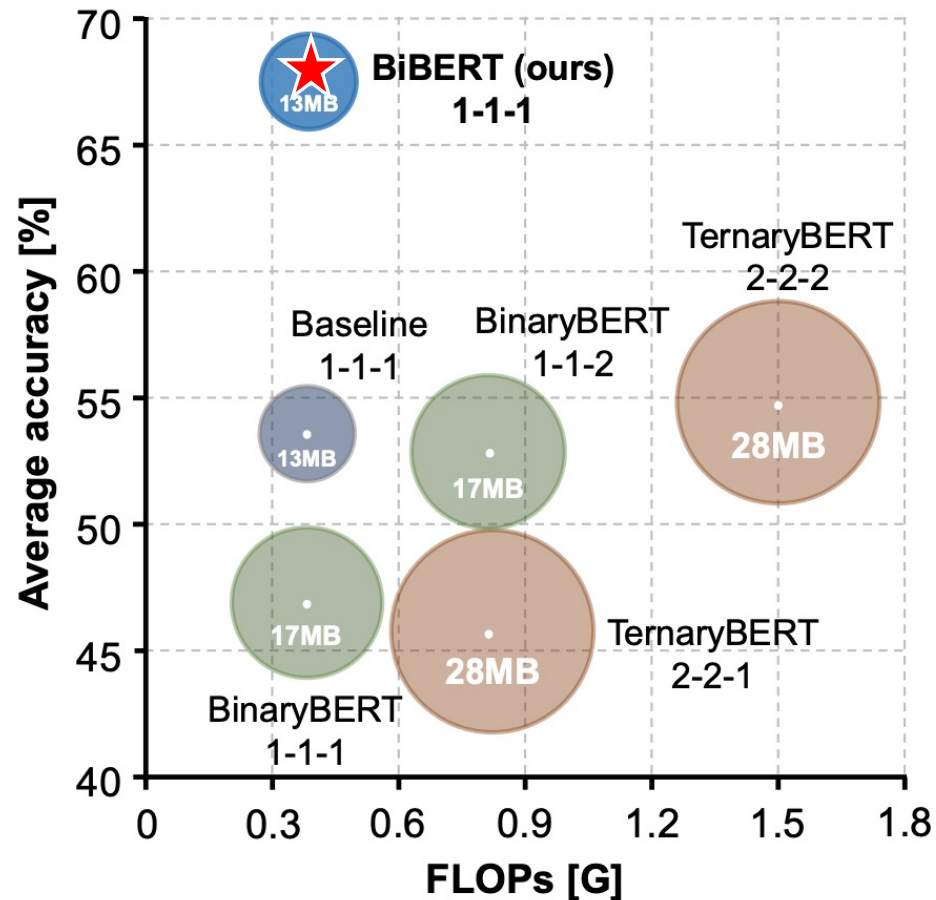
# Transformer Binarization: attention crash and recovery

## Performance



# Transformer Binarization: attention crash and recovery

## Performance

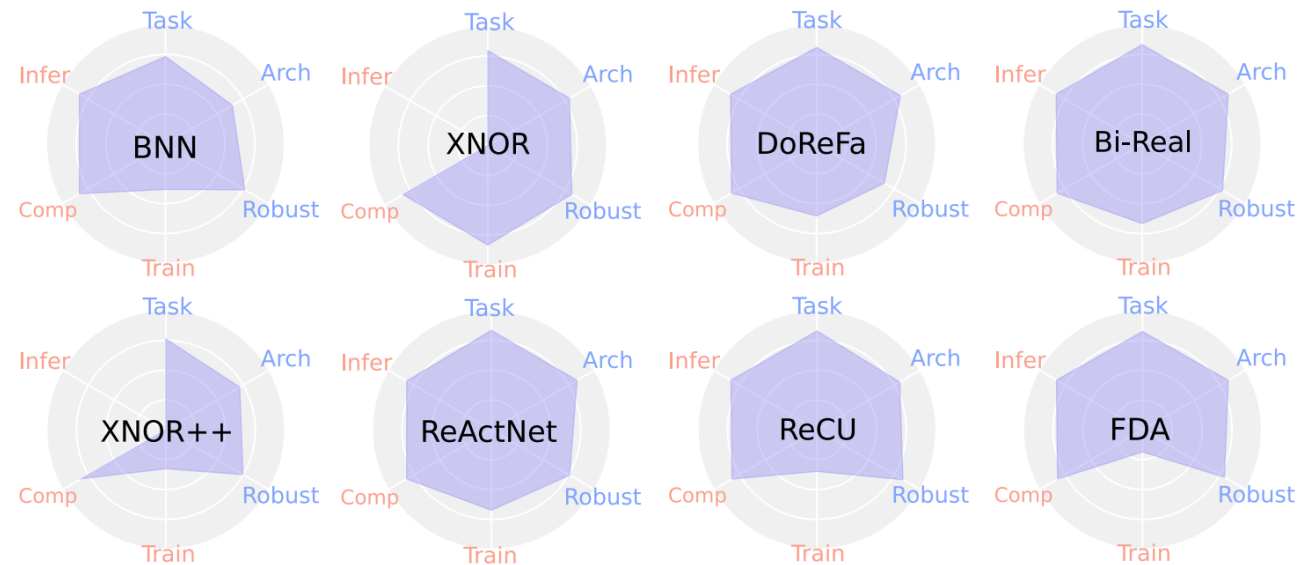
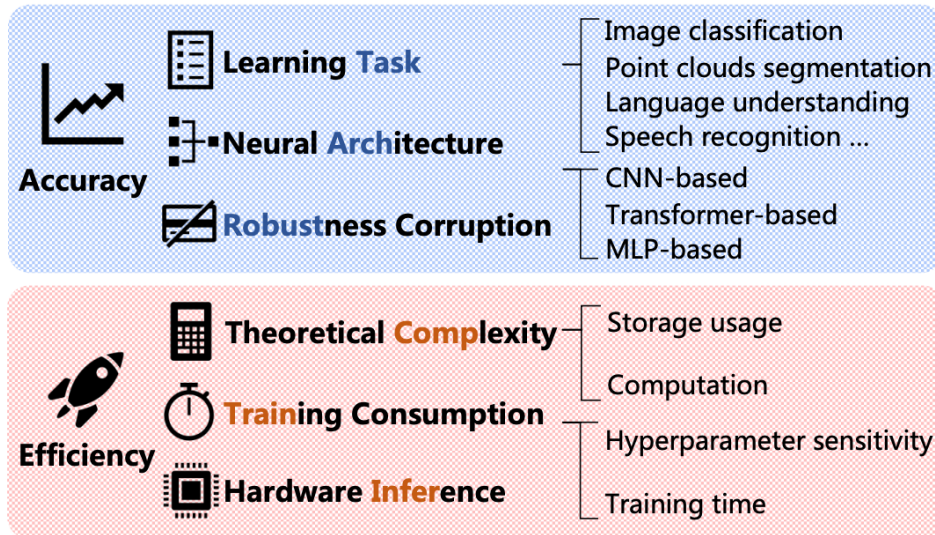


# Network Binarization: benchmark



## *BiBench: Benchmarking and Analyzing Network Binarization*

### Evaluation Tracks for Network Binarization



**6 Evaluation Tracks on Accuracy and Efficiency**

- 8 Binarization Algorithm
- 9 Deep Learning Datasets
- 13 Neural Architectures
- 2 Deployment Libraries
- 14 Hardware Chips



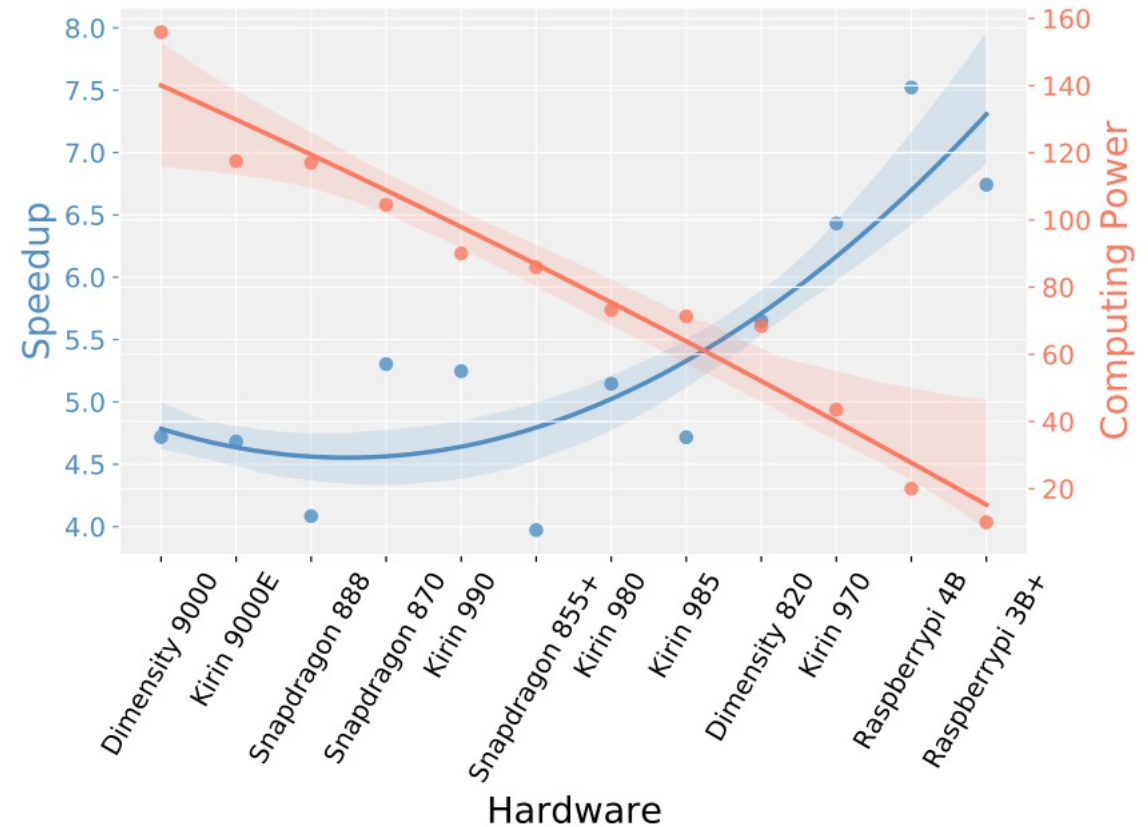
# Network Binarization: benchmark



**BiBench:** Benchmarking and Analyzing Network Binarization

## The 3 Most Effective Techniques for Generic Binarization:

- (1) Soft gradient approximation
- (2) Channel-wise scaling factors
- (3) Prebinarization parameter redistributing



Interesting Finding for Binarization:  
**Born for Edge**



# Network Binarization: survey

## *Binary Neural Networks: A Survey*

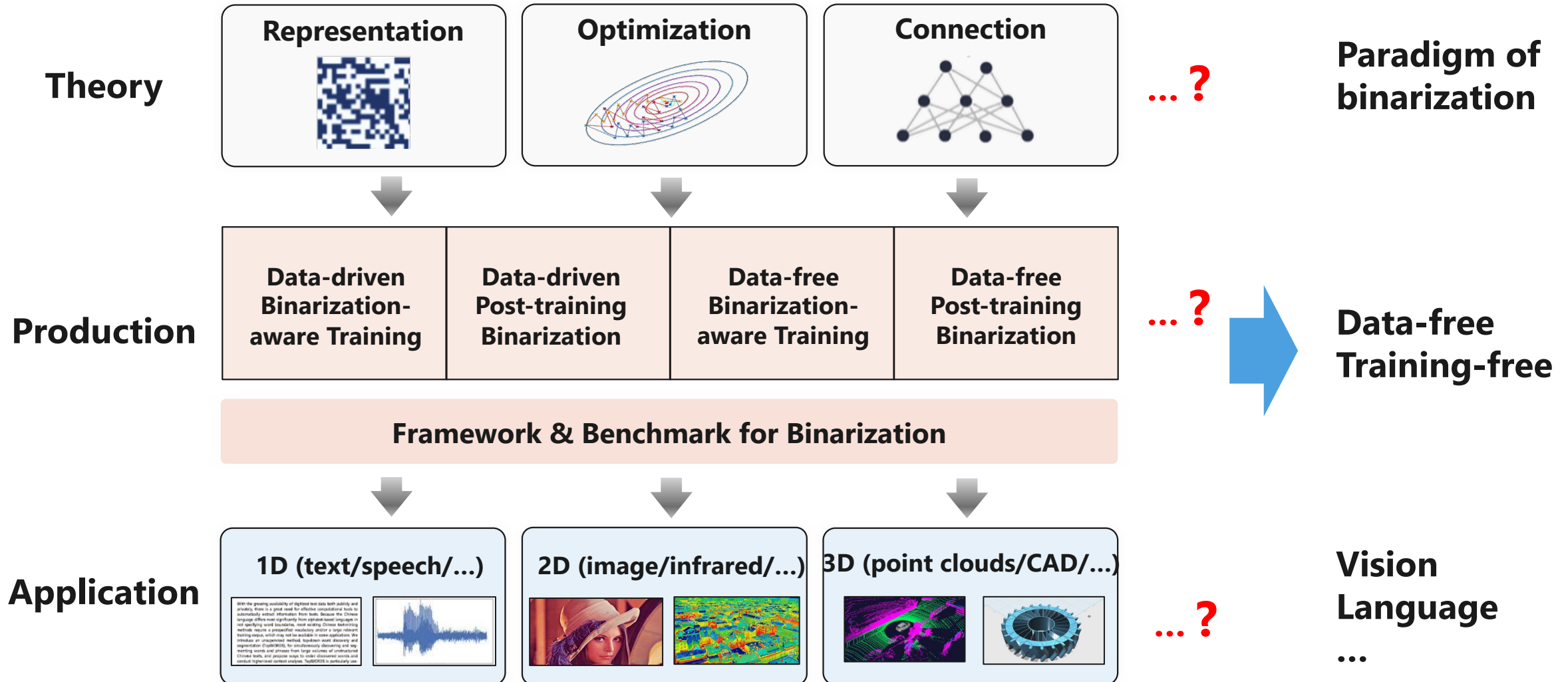
Haotong Qin, et al.  
*Pattern Recognition*  
Volume 105, 2020  
(132 citations)

<b>Naive BNNs</b>		BinaryConnect
		BNN
		Bitwise Neural Network
<b>Optimization Based BNNs</b>	<b>Minimize the Quantization Error</b>	XNOR-Net
		DoReFa-Net
		BWHN
		WRPN
		...
		LAB
	<b>Improve Network Loss Function</b>	RAD
		...
		HWGQ
	<b>Reduce the Gradient Error</b>	Bi-Real Net
		...



<https://github.com/htqin/awesome-model-quantization/>

# Network Binarization: future



**Thank you!**

**Q&A**

Haotong Qin

ETH Zürich CVL & Beihang University