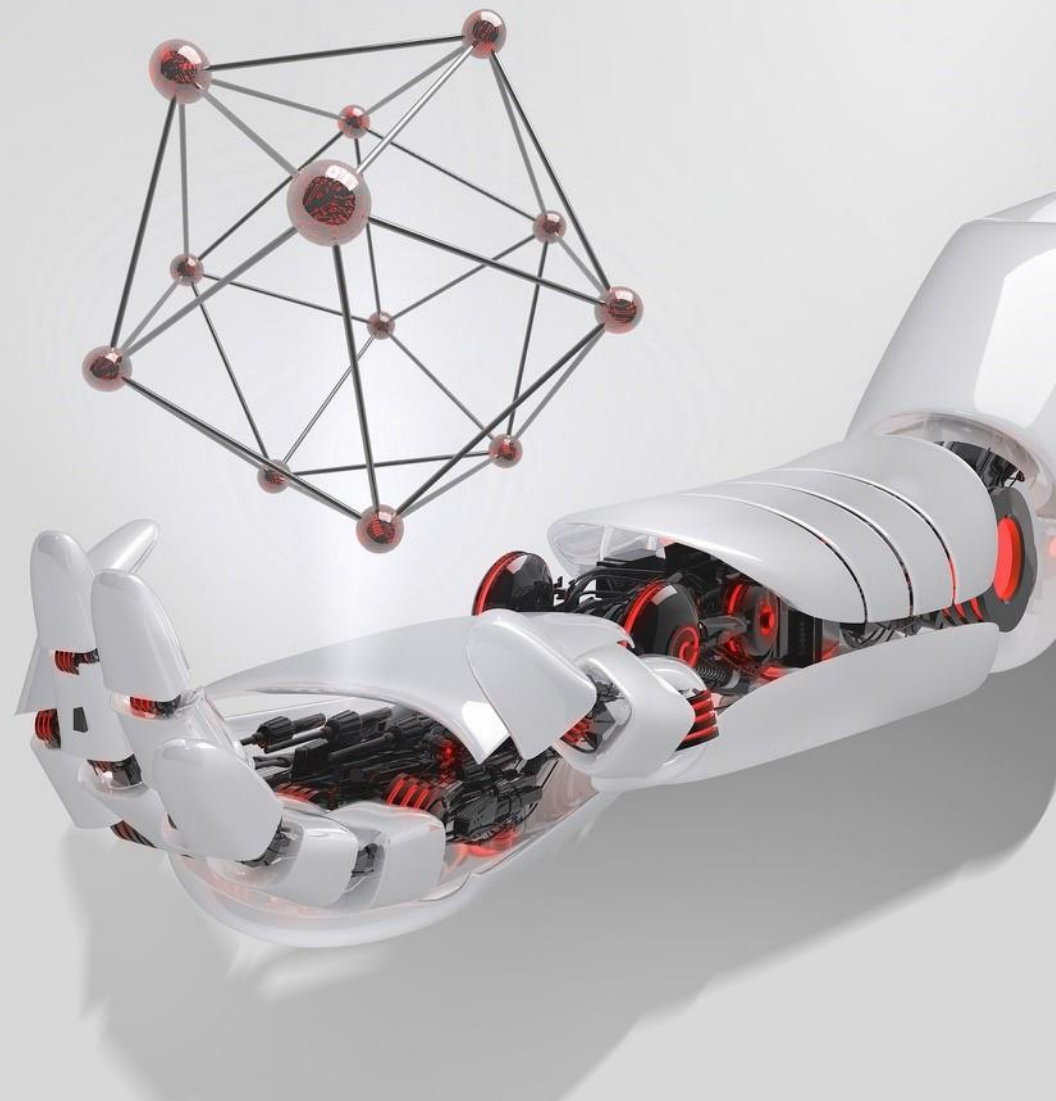


# Highly Efficient and Effective LLMs with Multi-Boolean Architectures

Ba-Hien Tran

[ba.hien.tran@huawei.com](mailto:ba.hien.tran@huawei.com)

Joint work with Van Minh Nguyen



Huawei Paris Research Center

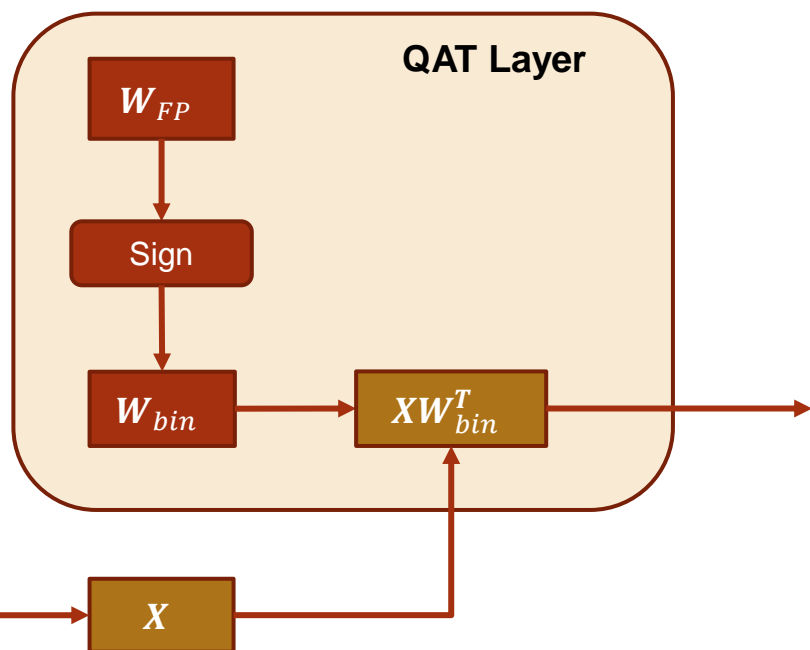


# Pitfalls of Binarized Neural Networks

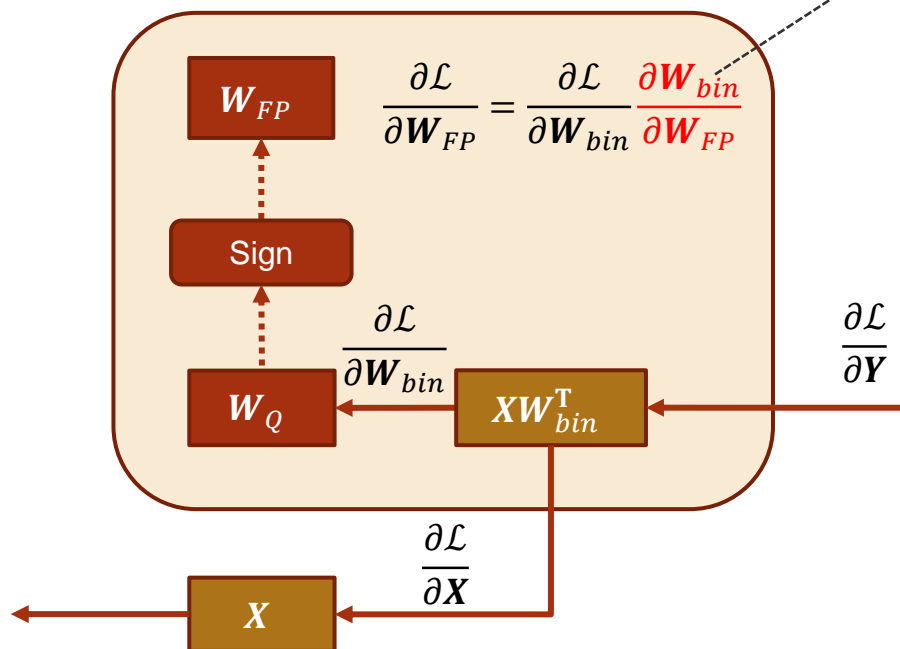
- Binarized linear layer:  $Y = \alpha \cdot XW_{bin}^T + \mathbf{b}$ , with  $W_{bin} = \text{Sign}(W_{FP})$
- Binarized weights are updated as :  $W_{bin} = \text{Sign}(W_{FP} - \eta \cdot \frac{\partial \mathcal{L}}{\partial W_{FP}})$

→ We train the weights in full-precision not binary domain

Forward Pass

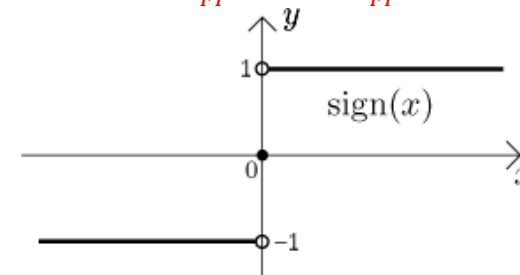


Backward Pass



- Not differentiable

$$\frac{\partial W_{bin}}{\partial W_{FP}} = \frac{\partial \text{Sign}(W_{FP})}{\partial W_{FP}}$$



- Derivative is 0 almost everywhere

- Model will learn nothing

- Gradient Approximation by Straight-through-estimator (STE)

$$\frac{\partial W_{bin}}{\partial W_{FP}} \approx 1$$

# Native Boolean Neural Networks

- **Bool linear layer:**  $Y_{[k,j]}^{(l)} = \sum_{i=1}^n L(X_{[k,i]}^{(l)}, W_{[i,j]}^{(l)}) + b_{[j]}^{(l)}$ , where  $L$  is a logic gate (**xnor**), and  $W_{[i,j]}^{(l)}$  are Boolean
- Logic gate can be extended to handle mixed-type data, such as real input and Boolean weights
- **Backward.** We can compute the loss signal for the weights as

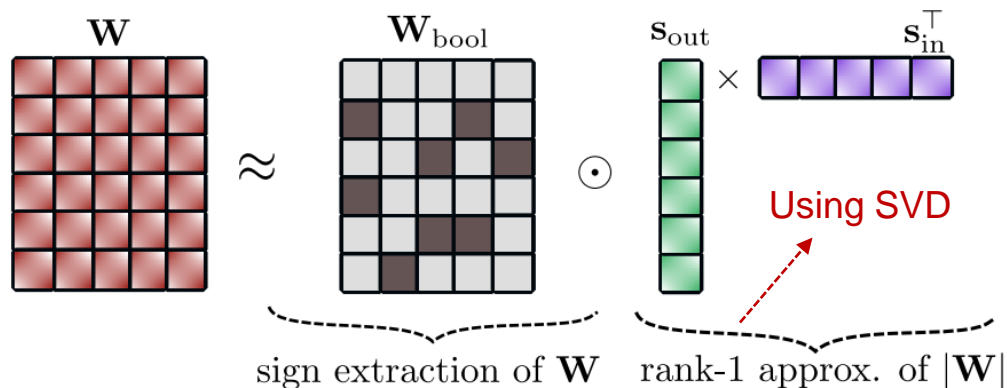
$$\frac{\delta \mathcal{L}}{\delta W_{[i,j]}^{(l)}} =: \mathbf{q}_{[i,j]}^{(l)} = \sum_{k=1}^b \mathbf{1}(q_{[k,i,j]}^{(l)} = \text{True}) |q_{[k,i,j]}^{(l)}| - \sum_{k=1}^b \mathbf{1}(q_{[k,i,j]}^{(l)} = \text{False}) |q_{[k,i,j]}^{(l)}|$$

where  $q_{[k,i,j]}^{(l)} = \text{xnor}(Z_{[k,j]}^{(l)}, X_{[k,i]}^{(l)})$

- **Update rule.**  $w_{[i,j]}^{(l)} = \neg w_{[i,j]}^{(l)}$  if  $\text{xnor}(q_{[k,j]}^{(l)}, w_{[i,j]}^{(l)}) = \text{True}$
- **Accumulator.**  $M_{[i,j]}^{(l),(t+1)} \leftarrow \beta^t M_{[i,j]}^{(l),t} + \eta \mathbf{q}_{[i,j]}^{(l),t}$
- We train the weights directly in Boolean domain
- $\text{xnor}(s, w) = s \times w$ , enabling direct use of existing linear algebra operations

# Boolean Reformulation for FP Linear Layers

## Sign-value-independent Decomposition (SVID)



We can prove that this initialization is optimal for preserving information from FP weights

**Proposition 4.1** (Xu et al. [59]). For  $\mathbf{W} \in \mathbb{R}^{m \times n}$ , write  $\mathbf{W} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^{\top}$  its SVD. Let  $\mathbf{a} = \sqrt{\tilde{\sigma}_1}\tilde{\mathbf{U}}_{[:,1]}$ , and  $\mathbf{b} = \sqrt{\tilde{\sigma}_1}\tilde{\mathbf{V}}_{[:,1]}$ . With the notations as described above, we have:

$$\|\mathbf{W} - \mathbf{W}_{\text{bool}} \odot \mathbf{s}_{\text{out}}\mathbf{s}_{\text{in}}^{\top}\|_F^2 \leq \|\mathbf{W} - \mathbf{a}\mathbf{b}^{\top}\|_F^2. \quad (4)$$

## Approximate the full-precision weights

$$\mathbf{W} = \mathbf{W}_{\text{bool}} \odot |\mathbf{W}| \approx \mathbf{W}_{\text{bool}} \odot (\mathbf{s}_{\text{out}}\mathbf{s}_{\text{in}}^{\top})$$

**Proposition 4.3.** For  $\mathbf{W} \in \mathbb{R}^{m \times n}$  and the notations as described above, we have:

$$\|\mathbf{W} - \mathbf{W}_{\text{bool}} \odot \mathbf{s}_{\text{out}}\mathbf{s}_{\text{in}}^{\top}\|_F^2 \leq \|\mathbf{W} - \mathbf{W}_{\text{bool}} \odot \mathbf{c}\mathbf{d}^{\top}\|_F^2, \quad \forall \mathbf{c} \in \mathbb{R}^{m \times 1}, \forall \mathbf{d} \in \mathbb{R}^{n \times 1}. \quad (5)$$

## The linear layer becomes

$$\mathbf{X}\mathbf{W}_{\text{FP}}^{\top} \approx [(\mathbf{X} \odot \mathbf{s}_{\text{in}}^{\top}) \mathbf{W}_{\text{bool}}] \odot \mathbf{s}_{\text{out}}^{\top}$$

# Enhanced Expressivity with Multiple Boolean Kernels

SVID is good but still not expressive enough to fully capture FP weights

- We employ multiple kernels, where each kernel utilizes distinct Boolean weights and scaling factors

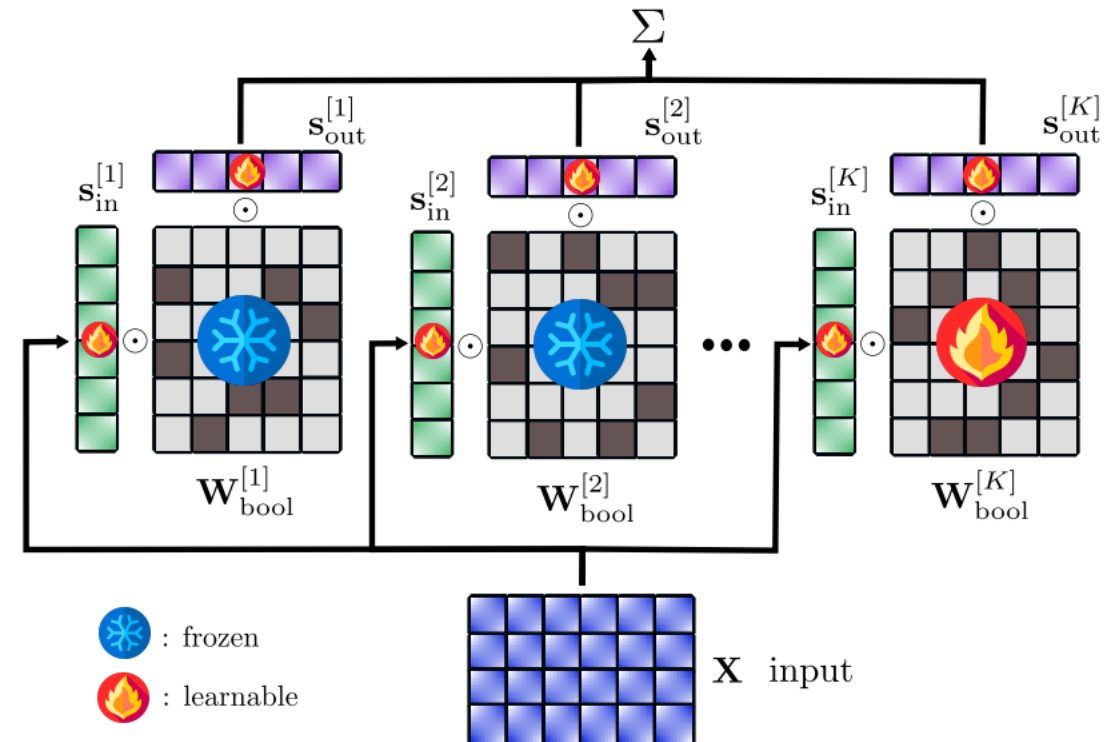
$$\mathbf{W}_{\text{FP}} \approx \mathbf{W}_{\text{approx}} \triangleq \sum_{k=1}^K \mathbf{W}_{\text{bool}}^{[k]} \odot (\mathbf{s}_{\text{out}}^{[k]} \mathbf{s}_{\text{in}}^{[k] \top})$$

- The linear layers become

$$\mathbf{X} \mathbf{W}_{\text{FP}}^{\top} \approx \sum_{k=1}^K \left[ \left( \mathbf{X} \odot \mathbf{s}_{\text{in}}^{[k] \top} \right) \mathbf{W}_{\text{bool}}^{[k]} \right] \odot \mathbf{s}_{\text{out}}^{[k] \top}$$

- The dominant computational cost is from the matrix multiplication between the **Boolean weights** and **scaled inputs**

- We showed that training the last kernel is enough!



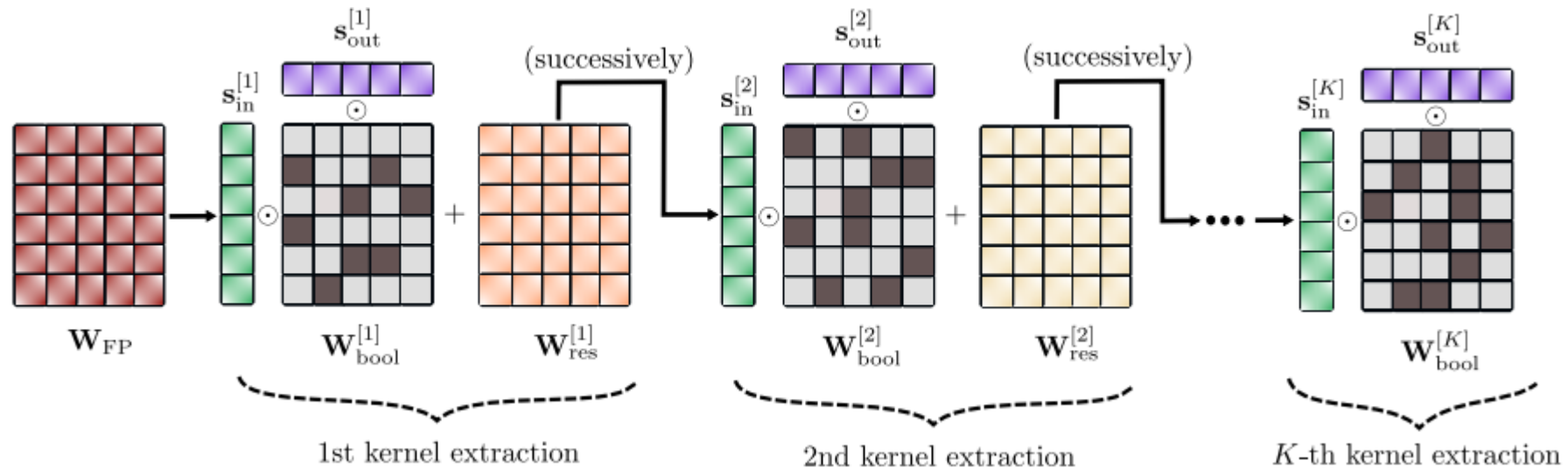
# Effective Knowledge Transfer into Boolean Models

## Successive Extraction using SVID

- Further proceed to SVID process to approximate **residual error** introduced by the previous step

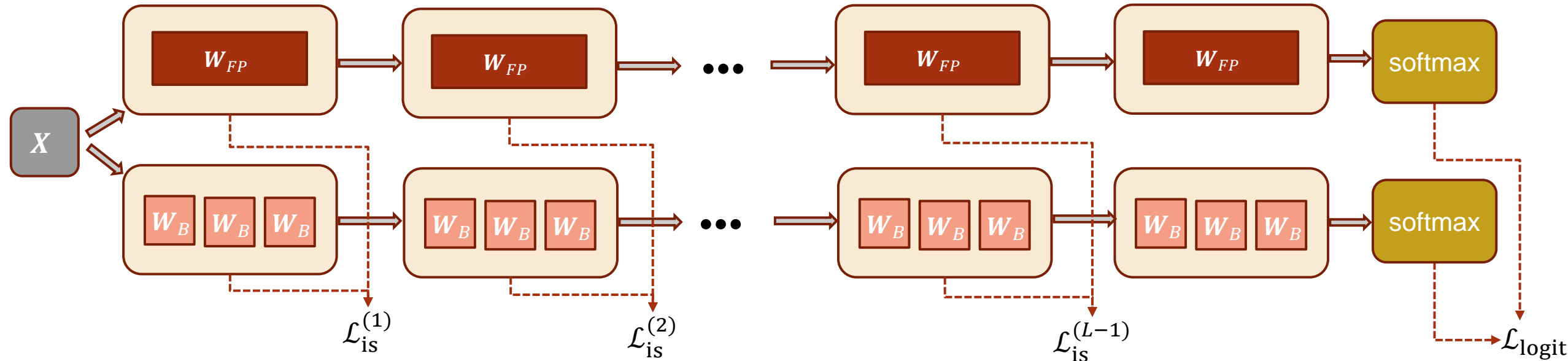
$$\mathbf{W}_{\text{res}}^{[k]} = \mathbf{W}_{\text{input}}^{[k]} - \mathbf{W}_{\text{bool}}^{[k]} \odot \left( \mathbf{s}_{\text{out}}^{[k]} \mathbf{s}_{\text{in}}^{[k] \top} \right)$$

Residual matrix



# Effective Knowledge Transfer into Boolean Models

## Finetuning with Knowledge Distillation

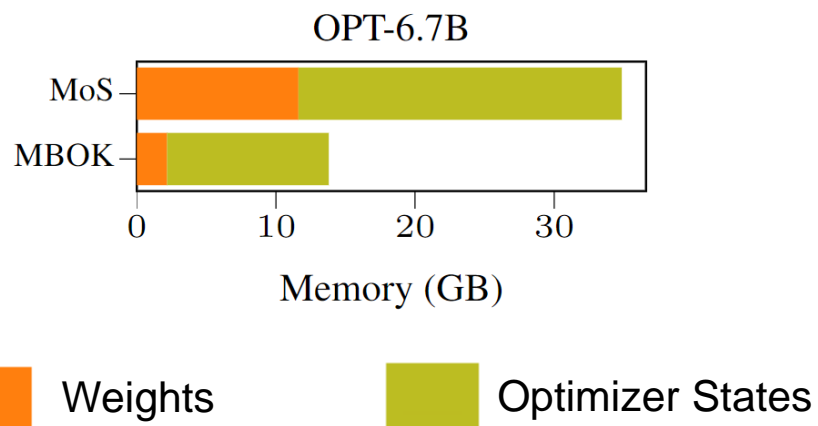


- Using multiple kernels effectively captures the original weight matrix, very small residual errors may still remain
- These errors can accumulate as they propagate through the layers
- Knowledge distillation for calibrating the whole model using very efficient Boolean optimizer

$$\mathcal{L} = \mathcal{L}_{logits} + \gamma \mathcal{L}_{is}$$
$$\mathcal{L}_{logits} = \frac{1}{L} \sum_{j=1}^L D_{logits} (p_{FP}(\mathbf{X}_{[j]}; \tau), p_{bool}(\mathbf{X}_{[j]}; \tau))$$
$$\mathcal{L}_{is} = \frac{1}{L} \sum_{h \in H} \sum_{j=1}^L \left\| \mathbf{Q}_{FP}^{j,h} - \mathbf{Q}_{bool}^{j,h} \right\|_2^2$$

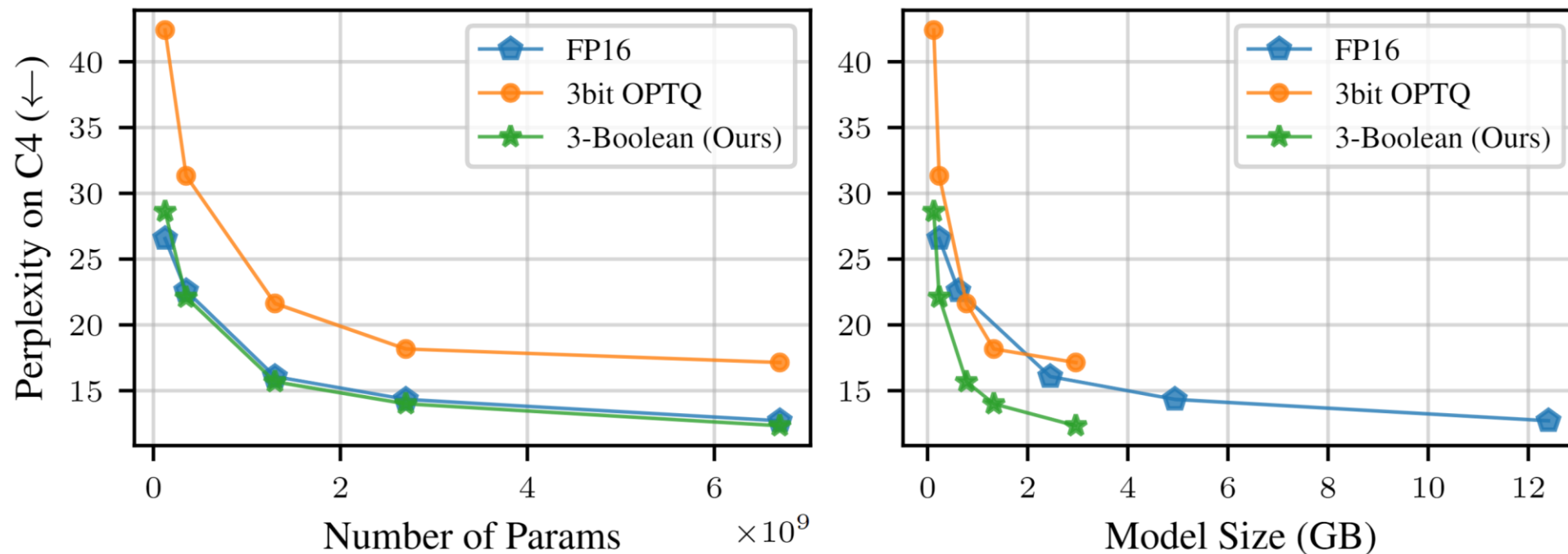
# State-of-the-art Ultra Low-Bit Performance

- We evaluate our method on LLMs
- Ultra low-bit compression is very challenging
- Significantly outperforms SOTA quantization and binarization methods in terms of performance
- Closely matches the performance of FP model
- Efficient in both training and inference



Model	Method	Wbits	Perplexity (↓)				Zero-shot Accuracy (↑)				Average	
			Wiki2	C4	BoolQ	PIQA	Hella.	WinoG.	ARC-e	ARC-c		
OPT-1.3B	FP-16	16	14.62	14.72	57.82	72.42	53.70	59.51	50.97	29.52	53.99	
	PB-LLM	1.7	272.83	175.42	62.17	54.24	27.25	50.27	27.98	23.72	40.94	
	BiLLM	1.11	69.45	63.92	61.92	59.52	33.81	49.32	34.38	22.35	43.55	
	OneBit	1	20.36	20.76	57.85	66.53	39.21	54.61	42.80	23.97	47.50	
	MoS	1	18.45	18.83	60.34	68.66	41.99	53.99	44.87	26.19	49.34	
	GPTQ	2	9.5e3	3.8e3	39.60	52.07	25.57	49.33	26.68	23.63	35.15	
	LLM-QAT	2	4.9e3	2.1e3	37.83	50.05	25.72	49.72	25.76	25.09	34.07	
	OmniQuant	2	42.43	55.64	56.45	60.94	33.39	51.85	38.76	23.38	44.13	
	MBOK [Ours]	2×1		<b>16.13</b>	<b>16.61</b>	58.53	70.67	48.11	56.75	48.19	27.90	<b>51.69</b>
	MBOK [Ours]	3×1		<b>15.30</b>	<b>15.68</b>	60.64	70.78	50.71	56.83	48.82	28.49	<b>52.71</b>
MBOK [Ours]	4×1		<b>14.83</b>	<b>14.92</b>	60.95	70.85	51.02	56.85	49.13	29.24	<b>53.01</b>	
LLaMA-7B	FP-16	16	5.68	7.08	73.21	77.42	72.99	66.85	52.53	41.38	64.06	
	PB-LLM	1.11	41.66	48.15	62.23	58.65	34.64	51.14	33.08	25.68	44.24	
	OneBit	1	8.48	10.49	62.50	70.40	54.03	55.32	41.07	30.88	52.36	
	MoS	1	7.97	9.72	64.59	71.82	58.18	58.88	42.09	31.31	54.48	
	GPTQ	2	1.9e3	7.8e2	43.79	49.95	25.63	49.41	25.84	27.47	37.02	
	LLM-QAT	2	7.1e2	3.0e2	37.83	50.87	24.76	51.78	26.26	25.51	36.17	
	OmniQuant	2	15.34	26.21	58.69	62.79	43.68	52.96	41.54	29.35	48.17	
	MBOK [Ours]	2×1		<b>6.83</b>	<b>8.53</b>	69.20	74.32	64.80	60.30	49.05	34.90	<b>58.76</b>
	MBOK [Ours]	3×1		<b>6.20</b>	<b>7.76</b>	67.89	76.15	68.91	63.30	48.94	37.62	<b>60.47</b>
	MBOK [Ours]	4×1		<b>6.01</b>	<b>7.53</b>	68.16	76.71	69.85	62.09	49.24	38.14	<b>60.70</b>
LLaMA-13B	FP-16	16	5.09	6.61	68.47	79.05	76.24	70.17	59.85	44.54	66.39	
	PB-LLM	1.7	35.83	39.79	62.17	58.70	33.97	52.17	31.86	23.63	43.75	
	BiLLM	1.11	14.56	16.67	62.53	68.17	52.24	59.43	41.91	29.94	52.37	
	OneBit	1	7.65	9.56	63.30	71.98	60.61	59.43	42.85	32.42	55.10	
	MoS	1	7.16	8.81	63.82	73.88	64.05	60.93	44.28	33.11	56.68	
	GPTQ	2	3.2e3	9.9e2	42.39	50.00	25.27	50.67	26.14	27.39	36.98	
	LLM-QAT	2	1.8e3	1.2e3	37.83	50.33	25.40	51.62	27.02	26.87	36.51	
	OmniQuant	2	13.43	19.33	62.20	68.99	54.16	53.83	45.50	30.38	52.51	
	MBOK [Ours]	2×1		<b>6.17</b>	<b>7.88</b>	68.10	76.33	69.88	64.17	52.34	37.88	<b>61.45</b>
	MBOK [Ours]	3×1		<b>5.58</b>	<b>7.15</b>	67.39	77.74	73.37	66.61	54.04	41.21	<b>63.39</b>
MBOK [Ours]	4×1		<b>5.38</b>	<b>6.91</b>	68.69	77.63	74.23	66.53	56.14	41.38	<b>64.10</b>	

# Compression-performance Tradeoff



OPT Model	WBits	Wiki2					C4				
		125M	350M	1.3B	2.7B	6.7B	125M	350M	1.3B	2.7B	6.7B
FULL-PRECISION	16	27.65	22.00	14.63	12.47	10.86	26.56	22.59	16.07	14.34	12.71
RTN [61, 12]	3	37.28	25.94	48.17	16.92	12.10	33.91	26.21	24.51	18.43	14.36
QPTQ [18]	3	31.12	24.24	15.47	12.87	11.39	29.22	24.63	16.97	15.00	13.18
MBOK [Ours]	3×1	<b>29.10</b>	<b>23.12</b>	<b>15.30</b>	<b>13.09</b>	<b>11.03</b>	<b>28.62</b>	<b>22.10</b>	<b>15.68</b>	<b>14.00</b>	<b>12.33</b>

# Low Latency

- Latency (ms) of Linear Layers in LLaMA-13B
- Our approach (MBOK) dramatically speedups the inference

Weight Size	Full-precision 16	MBOK (Ours)	QUIP#	QTIP
5120 x 5120	0.16540	0.05074	0.62260	1.96368
5120 x 13824	0.42830	0.05098	0.62836	5.23681
13824 x 5120	0.43411	0.04987	0.62840	5.21193

- It is even much more significant with a native Boolean accelerator

Thank you.

