

# Compacting Binary Neural Networks by Sparse Kernel Selection

Yikai Wang<sup>1</sup>, Wenbing Huang<sup>2</sup>, Yinpeng Dong<sup>1,3</sup>, Fuchun Sun<sup>1</sup>, Anbang Yao<sup>4</sup>

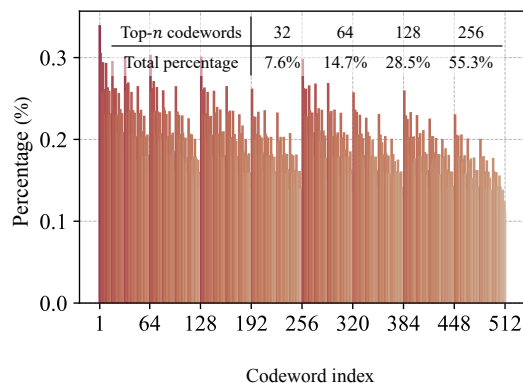
<sup>1</sup>BNRist Center, State Key Lab on Intelligent Technology and Systems,  
Department of Computer Science and Technology, Tsinghua University

<sup>2</sup>Gaoling School of Artificial Intelligence, Renmin University of China   <sup>3</sup>RealAI   <sup>4</sup>Intel Labs China

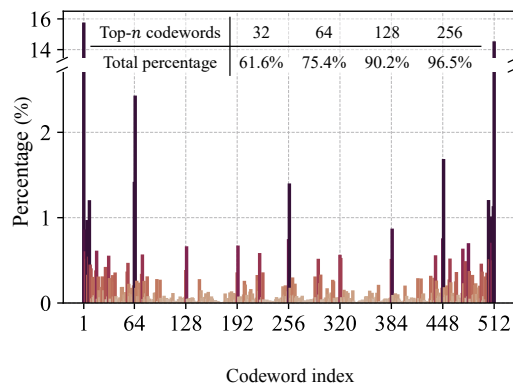
Primary Contact: Yikai Wang (yikaiw@outlook.com)

# Compacting Binary Neural Networks by Sparse Kernel Selection

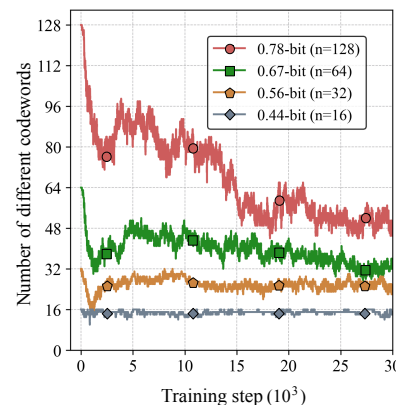
- We introduce how to compact and accelerate BNN further by Sparse Kernel Selection, abbreviated as **Sparks**.
- Our work is built based on a previously revealed phenomenon (by SNN<sup>[1]</sup>) that the  $3\times 3$  binary kernels in successful BNNs are nearly power-law distributed, **their values being mostly clustered into a small portion of codewords**. See the difference between Figure (a) and (b).
- In SNN, we observe that the sub-codebook is easy to degenerate during training (see Figure (c)), since codewords tend to be repetitive when being updated independently.
- While in our Sparks (Figure (d)), the diversity of codewords preserves by **selection-based learning**.



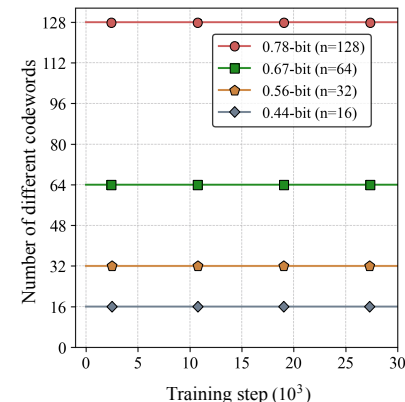
(a) Codebook constructed with sub-vectors (multi-channel codewords)



(b) Codebook constructed with kernels (single-channel codewords) (**ours**)



(c) Product quantization-based optimization for binary codewords



(d) Selection-based optimization for binary codewords (**ours**)

# Compacting Binary Neural Networks by Sparse Kernel Selection

( $K = 3$  for  $3 \times 3$  binary kernels)

**Property 1** We denote  $\mathbb{B} = \{-1, +1\}^{K \times K}$  as the codebook of binary kernels. For each  $\mathbf{w} \in \mathbb{R}^{K \times K}$ , the binary kernel  $\hat{\mathbf{w}}$  can be derived by a grouping process:

$$\hat{\mathbf{w}} = \text{sign}(\mathbf{w}) = \arg \min_{\mathbf{u} \in \mathbb{B}} \|\mathbf{u} - \mathbf{w}\|_2. \quad (1)$$

We compact BNNs by recasting the grouping as  $\hat{\mathbf{w}} = \arg \min_{\mathbf{u} \in \mathbb{U}} \|\mathbf{u} - \mathbf{w}\|_2$ , s.t.  $\mathbb{U} \subseteq \mathbb{B}$ .

Matrix representation, where  $\mathbf{P}$  is a permutation matrix and  $\mathbf{V}$  is fixed as a certain initial selection,

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{u} \in \mathbb{U}} \|\mathbf{u} - \mathbf{w}\|_2, \text{ s.t. } \mathbf{U} = \mathbf{B}\mathbf{P}\mathbf{V}, \mathbf{P} \in \mathbb{P}_N,$$

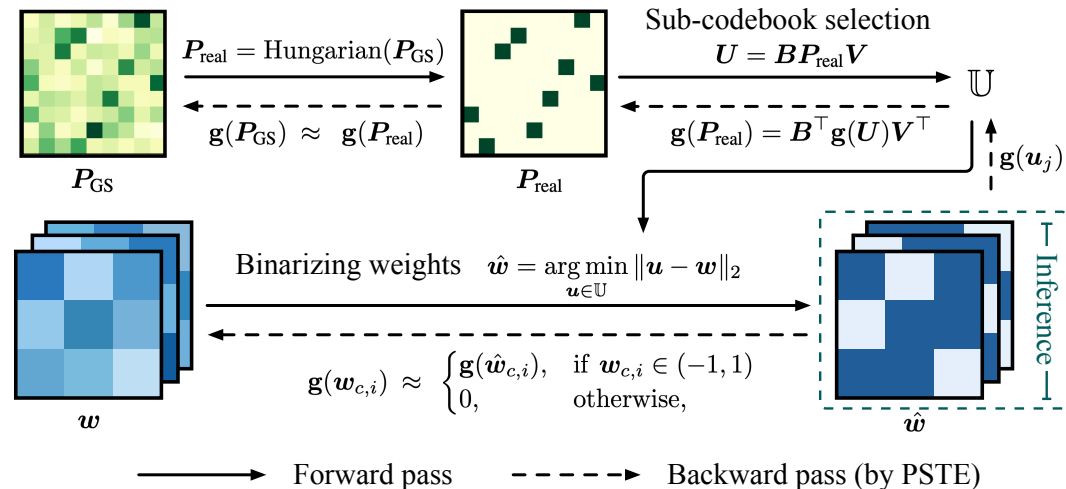
We learn the permutation matrix  $\mathbf{P}$  by Gumbel-Sinkhorn, denoted as  $\mathbf{P}_{\text{GS}}$ .

Forward pass

$$\begin{aligned} \mathbf{P}_{\text{real}} &= \text{Hungarian}(\mathbf{P}_{\text{GS}}), \\ \mathbf{U} &= \mathbf{B}\mathbf{P}_{\text{real}}\mathbf{V}, \\ \hat{\mathbf{w}}_c &= \arg \min_{\mathbf{u} \in \mathbb{U}} \|\mathbf{u} - \mathbf{w}_c\|_2, \end{aligned}$$

Backward pass

$$\begin{aligned} \mathbf{g}(\mathbf{w}_{c,i}) &\approx \begin{cases} \mathbf{g}(\hat{\mathbf{w}}_{c,i}), & \text{if } \mathbf{w}_{c,i} \in (-1, 1), \\ 0, & \text{otherwise,} \end{cases} \\ \mathbf{g}(\mathbf{u}_j) &= \sum_{c=1}^{C_{\text{in}} \times C_{\text{out}}} \mathbf{g}(\hat{\mathbf{w}}_c) \cdot \mathbb{I}_{\mathbf{u}_j = \arg \min_{\mathbf{u} \in \mathbb{U}} \|\mathbf{u} - \mathbf{w}_c\|_2}, \\ \mathbf{g}(\mathbf{P}_{\text{real}}) &= \mathbf{B}^\top \mathbf{g}(\mathbf{U})\mathbf{V}^\top, \\ \mathbf{g}(\mathbf{P}_{\text{GS}}) &\approx \mathbf{g}(\mathbf{P}_{\text{real}}), \text{ (our PSTE, will be introduced)} \end{aligned}$$



# Compacting Binary Neural Networks by Sparse Kernel Selection

How does Gumbel-Sinkhorn in our setting work?

Given a matrix  $\mathbf{X} \in \mathbb{R}^{N \times N}$  ( $N = |\mathbb{B}|$ ), the Sinkhorn operator over  $\mathcal{S}(\mathbf{X})$  is proceeded as follow,

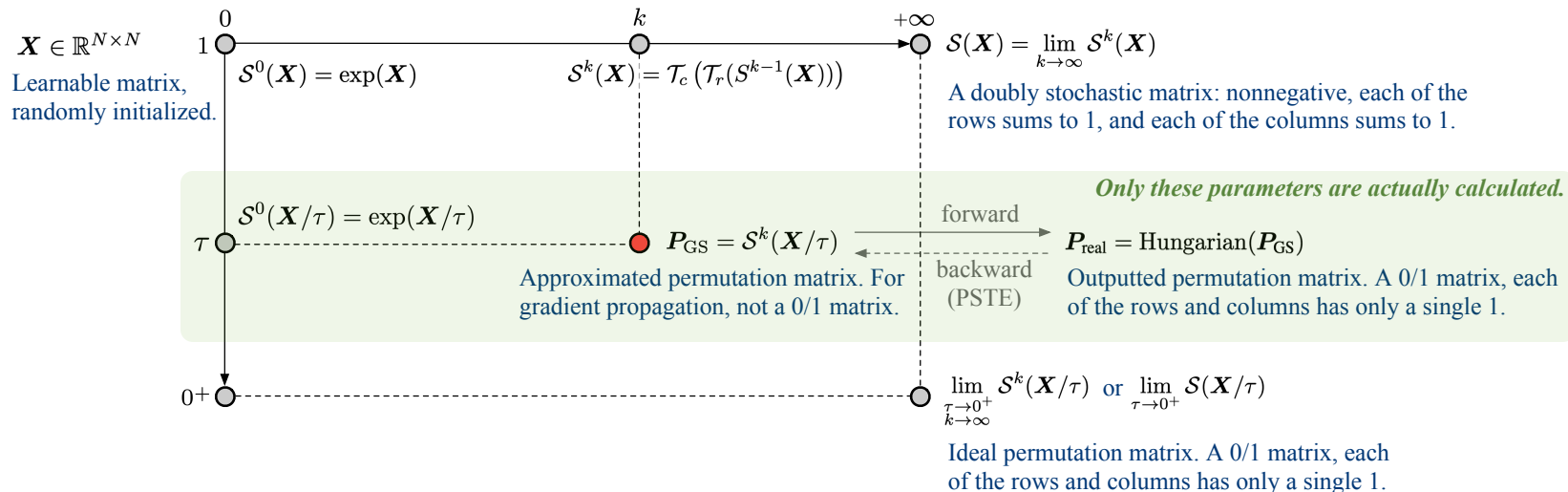
$$\mathcal{S}^0(\mathbf{X}) = \exp(\mathbf{X}), \tag{5}$$

$$\mathcal{S}^k(\mathbf{X}) = \mathcal{T}_c(\mathcal{T}_r(\mathcal{S}^{k-1}(\mathbf{X}))), \tag{6}$$

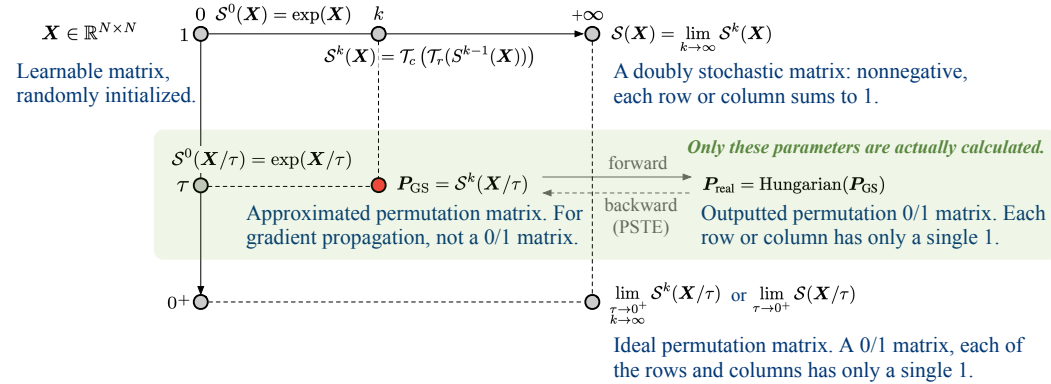
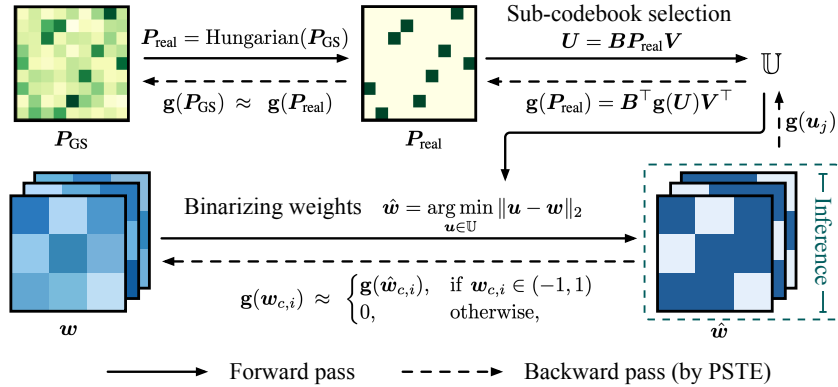
$$\mathcal{S}(\mathbf{X}) = \lim_{k \rightarrow \infty} \mathcal{S}^k(\mathbf{X}), \tag{7}$$

where  $\mathcal{T}_r(\mathbf{X}) = \mathbf{X} \oslash (\mathbf{X} \mathbf{1}_N \mathbf{1}_N^\top)$  and  $\mathcal{T}_c(\mathbf{X}) = \mathbf{X} \oslash (\mathbf{1}_N \mathbf{1}_N^\top \mathbf{X})$  are the row-wise and column-wise normalization operators, and  $\oslash$  denotes the element-wise division. For stability purpose, both normalization operators are calculated in the log domain in practice. The work by [41] proved that  $\mathcal{S}(\mathbf{X})$  belongs to the Birkhoff polytope—the set of doubly stochastic matrices.

By substituting the Gumbel-Sinkhorn matrix, we characterize the sub-codebook selection as  $\mathbf{U} = \mathbf{B} \mathcal{S}^k((\mathbf{X} + \epsilon)/\tau) \mathbf{V}$ ,



# Compacting Binary Neural Networks by Sparse Kernel Selection



**PSTE:** Approximate the gradient of the Gumbel-Sinkhorn matrix  $P_{GS}$  with  $P_{real}$ . We have the following theorem to guarantee the convergence for sufficiently large  $k$  and small  $\tau$ .

**Lemma 1** For sufficiently large  $k$  and small  $\tau$ , we define the entropy of a doubly-stochastic matrix  $P$  as  $h(P) = -\sum_{i,j} P_{i,j} \log P_{i,j}$ , and denote the rate of convergence for the Sinkhorn operator as  $r$  ( $0 < r < 1$ )<sup>3</sup>. There exists a convergence series  $s_\tau$  ( $s_\tau \rightarrow 0$  when  $\tau \rightarrow 0^+$ ) that satisfies

$$\|P_{real} - P_{GS}\|_2^2 = \mathcal{O}(s_\tau^2 + r^{2k}). \quad (18)$$

**Theorem 1** Assume that the training objective  $f$  w.r.t.  $P_{GS}$  is  $L$ -smooth, and the stochastic gradient of  $P_{real}$  is bounded by  $\mathbb{E}\|g(P_{real})\|_2^2 \leq \sigma^2$ . Denote the rate of convergence for the Sinkhorn operator as  $r$  ( $0 < r < 1$ ) and the stationary point as  $P_{GS}^*$ . Let the learning rate of PSTE be  $\eta = \frac{c}{\sqrt{T}}$  with  $c = \sqrt{\frac{f(P_{GS}^0) - f(P_{GS}^*)}{L\sigma^2}}$ . For a uniformly chosen  $\mathbf{u}$  from the iterates  $\{P_{real}^0, \dots, P_{real}^T\}$ , concretely  $\mathbf{u} = P_{real}^t$  with the probability  $p_t = \frac{1}{T+1}$ , it holds in expectation over the stochasticity and the selection of  $\mathbf{u}$ :

$$\mathbb{E}\|\nabla f(\mathbf{u})\|_2^2 = \mathcal{O}\left(\sigma \sqrt{\frac{f(P_{GS}^0) - f(P_{GS}^*)}{T/L} + L^2(s_\tau^2 + r^{2k})}\right). \quad (19)$$



- Comparisons of top-1 and top-5 accuracies with state-of-the-art methods on ImageNet based on ResNet-18.

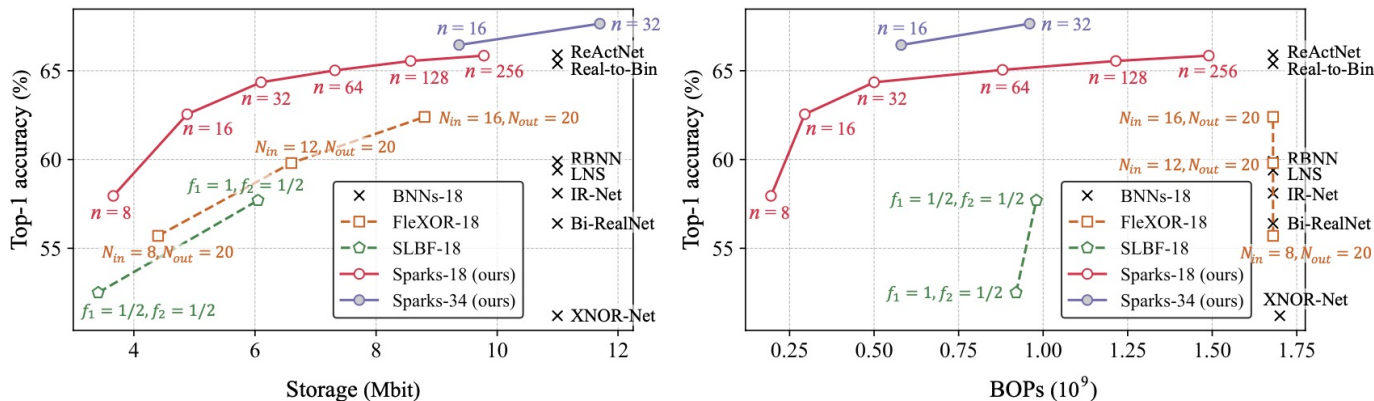
Method	Bit-width (W/A)	Accuracy (%)		Storage (Mbit)	BOPs ( $\times 10^9$ )
		Top-1	Top-5		
Full-precision	32/32	69.6	89.2	351.5	107.2 (1 $\times$ )
BNN [16]	1/1	42.2	69.2	11.0 (32 $\times$ )	1.70 (63 $\times$ )
XNOR-Net [37]	1/1	51.2	73.2	11.0 (32 $\times$ )	1.70 (63 $\times$ )
Bi-RealNet [31]	1/1	56.4	79.5	11.0 (32 $\times$ )	1.68 (64 $\times$ )
IR-Net [36]	1/1	58.1	80.0	11.0 (32 $\times$ )	1.68 (64 $\times$ )
LNS [10]	1/1	59.4	81.7	11.0 (32 $\times$ )	1.68 (64 $\times$ )
RBNN [26]	1/1	59.9	81.9	11.0 (32 $\times$ )	1.68 (64 $\times$ )
Ensemble-BNN [52]	(1/1) $\times$ 6	61.0	-	65.9 (5.3 $\times$ )	10.6 (10 $\times$ )
ABC-Net [28]	(1/1) $\times$ 5 <sup>2</sup>	65.0	85.9	274.5 (1.3 $\times$ )	42.5 (2.5 $\times$ )
Real-to-Bin [33]	1/1	65.4	86.2	11.0 (32 $\times$ )	1.68 (64 $\times$ )
ReActNet [32]	1/1	65.9	86.4	11.0 (32 $\times$ )	1.68 (64 $\times$ )
SLBF [24]	0.55/1	57.7	80.2	6.05 (58 $\times$ )	0.92 (117 $\times$ )
SLBF [24]	0.31/1	52.5	76.1	3.41 (103 $\times$ )	0.98 (110 $\times$ )
FleXOR [25]	0.80/1	62.4	83.0	8.80 (40 $\times$ )	1.68 (64 $\times$ )
FleXOR [25]	0.60/1	59.8	81.9	6.60 (53 $\times$ )	1.68 (64 $\times$ )
Sparks (ours)	0.78/1	65.5	86.2	8.57 (41 $\times$ )	1.22 (88 $\times$ )
Sparks (ours)	0.67/1	65.0	86.0	7.32 (48 $\times$ )	0.88 (122 $\times$ )
Sparks (ours)	0.56/1	64.3	85.6	6.10 (58 $\times$ )	0.50 (214 $\times$ )

- Results when extending our Sparks to wider or deeper models.

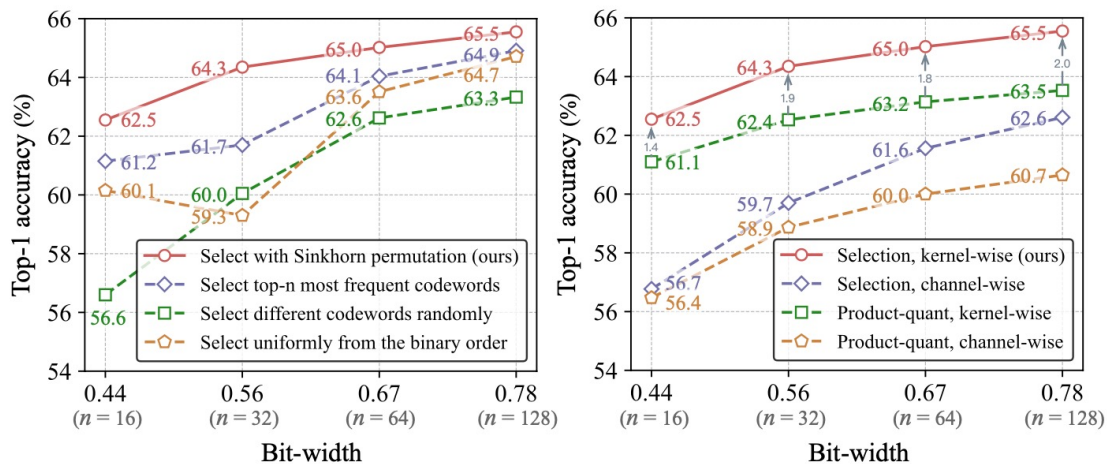
Method	Backbone	Bit-width (W/A)	Accuracy (%)		Storage (Mbit)	BOPs ( $\times 10^9$ )
			Top-1	Top-5		
ReActNet [32]	ResNet-18	1/1	65.9	86.4	11.0	1.68
Sparks-wide	ResNet-18 (+ABC-Net [28])	(0.56/1) $\times$ 3	<b>66.7</b>	<b>86.9</b>	18.3	<b>1.50</b>
Sparks-deep	ResNet-34	0.56/1	<b>67.6</b>	<b>87.5</b>	11.7	<b>0.96</b>
Sparks-deep	ResNet-34	0.44/1	<b>66.4</b>	<b>86.7</b>	<b>9.4</b>	<b>0.58</b>

# Compacting Binary Neural Networks by Sparse Kernel Selection

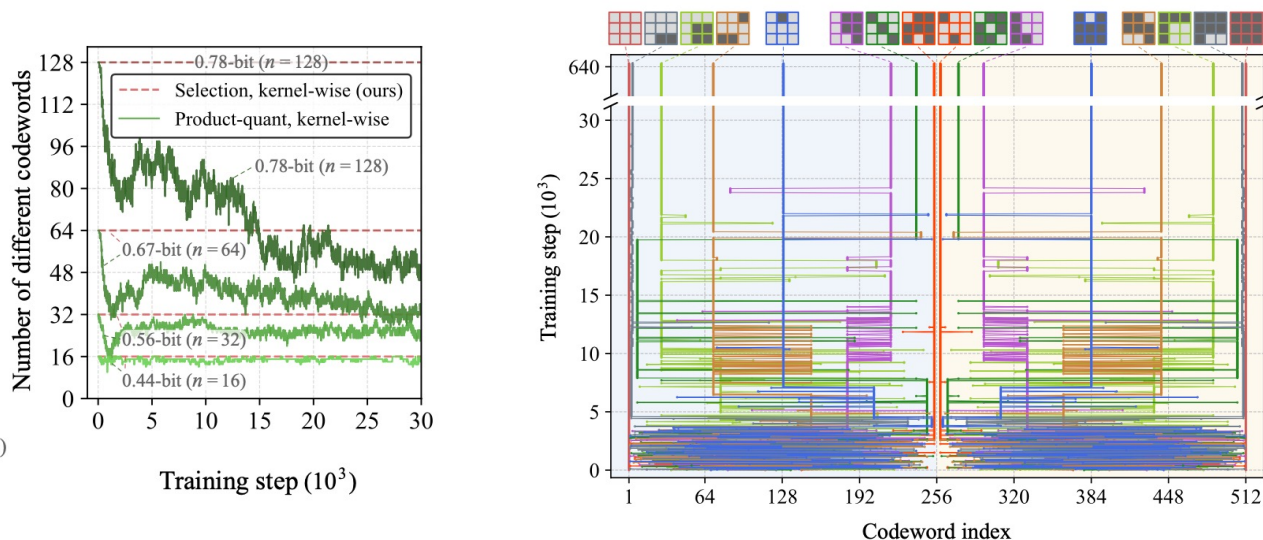
● Trade-off between performance and complexity on ImageNet,



● Ablation studies on ImageNet with ResNet-18,



● Codewords selection during training,





**Thanks**