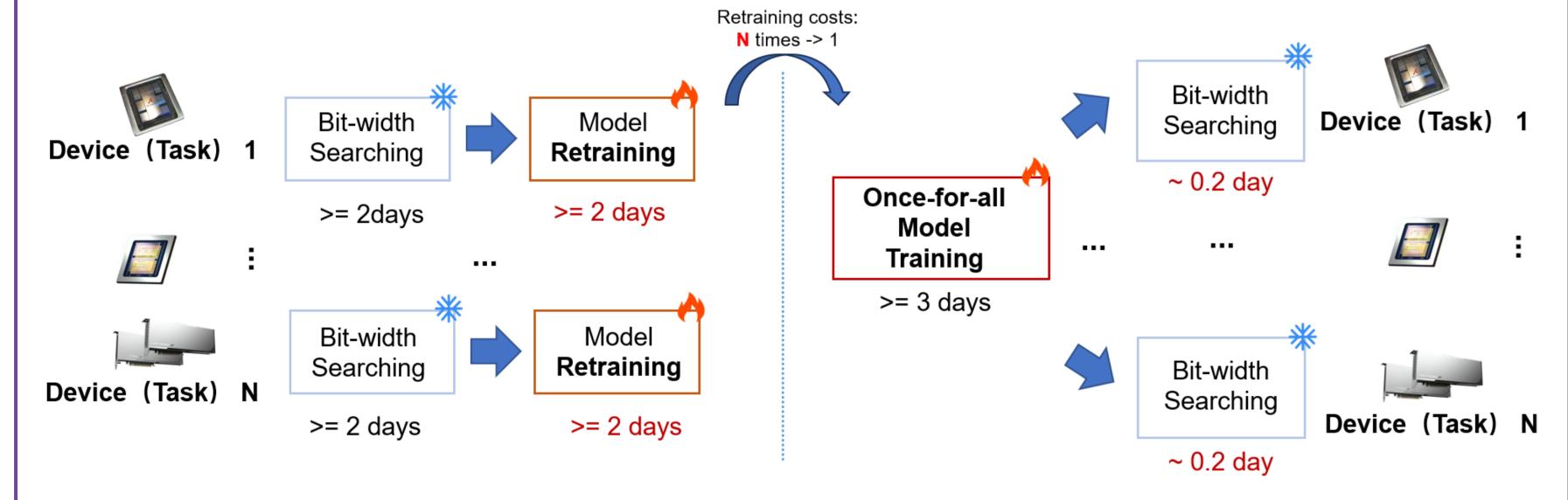


- optimal bit-width for each layer.
- well-performing ones.
- without introducing any additional training costs.

## 「有有人」 **Retraining-free Model Quantization via One-Shot Weight-Coupling Learning** Chen Tang\*, Yuan Meng\*, Jiacheng Jiang, Shuzhao Xie, Rongwei Lu, Xinzhu Ma, Zhi Wang, Wenwu Zhu **Contributions of Retraining-free Quantization** Weight-sharing for Quantization **Retraining-free Quantization (ReFQ) devises a one-**To jointly learn the quantized weights, we approximate the expectation term Bit-width Top-1 Method (W/A) shot training-then-searching paradigm for mixedover the whole search space of bit-width by randomly sampling the bit-width Baseline 32/32 70. PACT [7] 3/3 68. precision model compression. In the first stage, all configurations at each training step using a set of shared weights for different potential bit-width configurations are coupled and bit-width, the forward propagation is defined as follow: LSQ\* [12] 3/3 EWGS [28] thus optimized simultaneously within a set of $\underset{\mathbf{W}}{\operatorname{arg\,min}} \quad \mathbb{E}_{\mathcal{S}\sim\mathcal{A}}\left[\mathcal{L}(f(\mathbf{x};\mathcal{S},w^{(\mathcal{S})}),\mathbf{y})\right] \approx \underset{\mathbf{W}}{\operatorname{arg\,min}} \frac{1}{K} \quad \sum_{\mathbf{W}} \quad \left[\mathcal{L}(f(\mathbf{x};\mathcal{S}_{k},\hat{\mathbf{W}}^{(\mathcal{S}_{k})}),\mathbf{y})\right],$ EdMIPS [12] 3<sub>MP</sub> / 3<sub>MP</sub> 68.2 shared weights. In the second stage, a lightweight $GMPQ^*$ [60] $3_{MP} / 3_{MP}$ 68.6 DNAS [61] 3<sub>MP</sub> / 3<sub>MP</sub> 68. inference algorithm is applied to determine the FracBits [64] 3<sub>MP</sub> / 3<sub>MP</sub> 69.4 Bit-width Interference among Highly Coupled Weights during Optimization LIMPQ [47] 3<sub>MP</sub> / 3<sub>MP</sub> 69. SEAM [49] 3<sub>MP</sub> / 3<sub>MP</sub> 70.0 To solve the bit-width interference problem in the Ours 2<sub>MP</sub> / 3<sub>MP</sub> 67. — w 2bits 🔶 target weight 🧼 —— 2-bits gradient –— 4-bits only gradient first stage, we design a bit-width scheduler to dynamically freeze the most turbulent bit-width of 0.034 0.58 -0.58 -A 0.0975 0.56 -0.033 -0.56 layers during training, to ensure the rest bit-widths w 2bits converged properly. We then present an informa-0.05 -0.52 -0.52 tion distortion mitigation technique to align the 0.04 behavior of the bad-performing bit-widths to the training step 200 400 **<u>Bit-width interference</u>:** While training becomes possible with weight-sharing Fig. 3 In second stage, an inference-only greedy search quantization, we have observed it exhibits training instability and the weight moves closer to scheme is devised to evaluate configurations of retraining. the quantization bound more frequently (Fig.2), introducing extra small bit-widths induces significant random oscillations for the learning process (Fig.3), signifying hindering model **Results: Transfer Learning** convergence. **Overview of Retraining-free Quantization** Training and Search Techniques of One-shot Retraining-free Quantization To effectively train the weight-sharing model, we develop: Current pipeline: search-retraining Proposed pipeline: training-then-search Make use of quantization results with weight-sharin 32/32 Retraining are repetitive and costly 79.4



- One-shot Training for All Deployments: unlike previous searching-then-retraining pipelines, our one-shot method shifts the costly training process to the front end of the pipeline, forming a training-then-searching pipeline, thus the training <mark>only performs once</mark> to support diverse deployment requirements.
- Interference-aware Training: we identify the bitwidth interference problem and design several key techniques to improve the convergence of one-shot model.

**Dynamic Bit-width Schedule** selectively freezes the bit-width causing weight **interference** from unstable bit-width set  $\Omega$  to ensure proper convergence for remaining bit-widths during training:

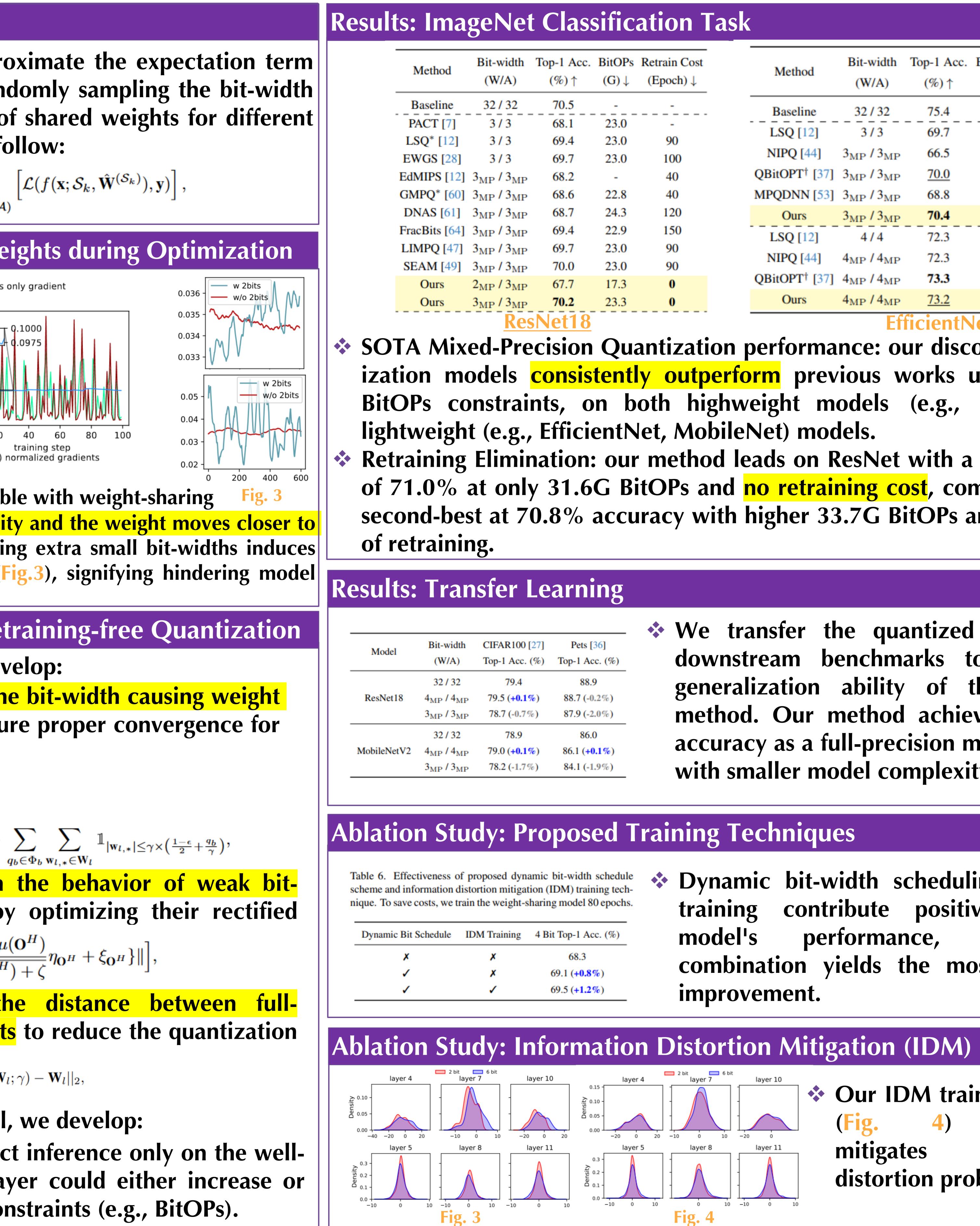
 $\Omega \leftarrow \textbf{TopKToFreeze}(\hat{\Delta} \mathbf{W}^{\text{unstable}}; \mathcal{K}), \text{ where}$ 

 $\hat{\Delta}\mathbf{W}^{\text{unstable}} \triangleq \{\hat{\Delta}\mathbf{W}_{l}^{\text{unstable}}\}_{l=0}^{L-1} \quad \text{where} \quad \hat{\Delta}\mathbf{W}_{l}^{\text{unstable}} = \sum_{b \in B^{(w)}} \quad \frac{1}{2^{b}} \frac{1}{\|\mathbf{W}_{l}\|_{0}} \cdot \sum_{q_{b} \in \Phi_{b}} \sum_{\mathbf{W}_{l,*} \in \mathbf{W}_{l}} \mathbb{1}_{|\mathbf{W}_{l,*}| \le \gamma \times \left(\frac{1-\epsilon}{2} + \frac{q_{b}}{\gamma}\right)},$ Information Distortion Mitigation aims to align the behavior of weak bitwidths to their high-performing counterparts by optimizing their rectified Euclidea  $\mathbb{E}\Big[\big\|\max\{Q, \frac{\mathbf{O}^S - \mu(\mathbf{O}^S)}{\sqrt{\sigma(\mathbf{O}^S) + \zeta}}\eta_{\mathbf{O}^S} + \xi_{\mathbf{O}^S}\} - \max\{Q, \frac{\mathbf{O}^H - \mu(\mathbf{O}^H)}{\sqrt{\sigma(\mathbf{O}^H) + \zeta}}\eta_{\mathbf{O}^H} + \xi_{\mathbf{O}^H}\}\big\|\Big],$ 

Quantization Error Minimization optimizes the distance between fullprecision latent weights and the quantized weights to reduce the quantization error:  $\mathcal{L}_{\text{QE}} = \frac{1}{K} \sum_{k=0}^{K-1} ||\hat{\mathbf{W}}^{(\mathcal{S}_k)} - \mathbf{W}||_2 = \frac{1}{K} \sum_{k=0}^{K-1} \sum_{l=0}^{L-1} ||Q_{b_{l,w}^{(k)} \in \mathcal{S}_k}(\mathbf{W}_l; \gamma) - \mathbf{W}_l||_2,$ 

To efficiently search from the weight-sharing model, we develop: **Bidirectional Greedy Search: We directly conduct inference only on the well**trained one-shot model, at a time, only one layer could either increase or

decrease its bit-width until it fits the expected constraints (e.g., BitOPs).







Acc.	BitOPs $(G) \downarrow$	Retrain Cost (Epoch)↓		Method	Bit-width (W/A)	Top-1 Acc. (%)↑	BitOPs (G)↓	Retrain Cost (Epoch)↓
.5				Baseline	32/32	75.4	-	-
.1	23.0	-		LSQ [12]	3/3	69.7	4.2	
.4	23.0	90			575		4.2	
.7	23.0	100		NIPQ [44]	$3_{\mathrm{MP}}$ / $3_{\mathrm{MP}}$	66.5	-	43
.2	-	40		QBitOPT <sup>†</sup> [37]	$3_{\rm MP}$ / $3_{\rm MP}$	<u>70.0</u>	-	<u>30</u>
.6	22.8	40		MPQDNN [53]	$3_{\rm MP}$ / $3_{\rm MP}$	68.8	-	50
.7	24.3	120		Ours	$3_{\mathrm{MP}}$ / $3_{\mathrm{MP}}$	70.4	4.5	0
.4	22.9	150		LSQ [12]	4/4	72.3	6.8	90
.7	23.0	90					0.0	
.0	23.0	90		NIPQ [44]	$4_{\mathrm{MP}}$ / $4_{\mathrm{MP}}$	72.3	-	43
.7	17.3	0		QBitOPT <sup>†</sup> [37]	$4_{\mathrm{MP}}$ / $4_{\mathrm{MP}}$	73.3	-	<u>30</u>
.2	23.3	0		Ours	$4_{\rm MP}$ / $4_{\rm MP}$	<u>73.2</u>	6.9	0
18		EfficientNet						

SOTA Mixed-Precision Quantization performance: our discovered quantization models consistently outperform previous works under various BitOPs constraints, on both highweight models (e.g., ResNet) and

Retraining Elimination: our method leads on ResNet with a top accuracy of 71.0% at only 31.6G BitOPs and no retraining cost, compared to the second-best at 70.8% accuracy with higher 33.7G BitOPs and 90 epochs

We transfer the quantized weights for downstream benchmarks to verify the generalization ability of the proposed method. Our method achieves the same accuracy as a full-precision model at 4-bits with smaller model complexity.

Dynamic bit-width scheduling and IDM training contribute positively to the performance, their and combination yields the most significant improvement.

• Our IDM training technique significantly Fig. information mitigates distortion problem (Fig. 3).