

## Understanding the Impact of Compression Order in Joint Model Compression

Minjun Kim, Jaehyeon Choi, Hyunwoo Yang, Jongjin Kim, Jinho Song, and U Kang\*  
 Seoul National University, Seoul, South Korea, ✉ {minjun.kim, ukang}@snu.ac.kr, \*: Corresponding Author

### Summary

#### The Progressive Intensity Hypothesis

- **TL;DR:** Weaker perturbations should precede stronger ones for better performance in joint model compression
- **GitHub:** <https://github.com/snudatalab/PQQP>

### Joint Model Compression

#### Joint Model Compression

- Combining two or more compression techniques
- Two categories: 1) co-designed or 2) **post-hoc** methods
- Sequential application of  $f_1(\cdot)$  and  $f_2(\cdot)$ :  $f_1 \rightarrow f_2$  or  $(f_2 \circ f_1)(\cdot)$
- Compression rate  $C$ : memory usage ratio of original  $\phi$  and compressed  $\phi'$  models ( $C_P = 1/(1-p)$ ,  $C_Q = B_{orig}/B_Q$ )

#### Motivation. Compression Order in Joint Model Compression

- Identifying an optimal order can yield a **“free lunch”!**
- Error rates deviates up to **3.41x** (at compression rate  $C = 4$ )
- However, the role of compression order has been largely overlooked

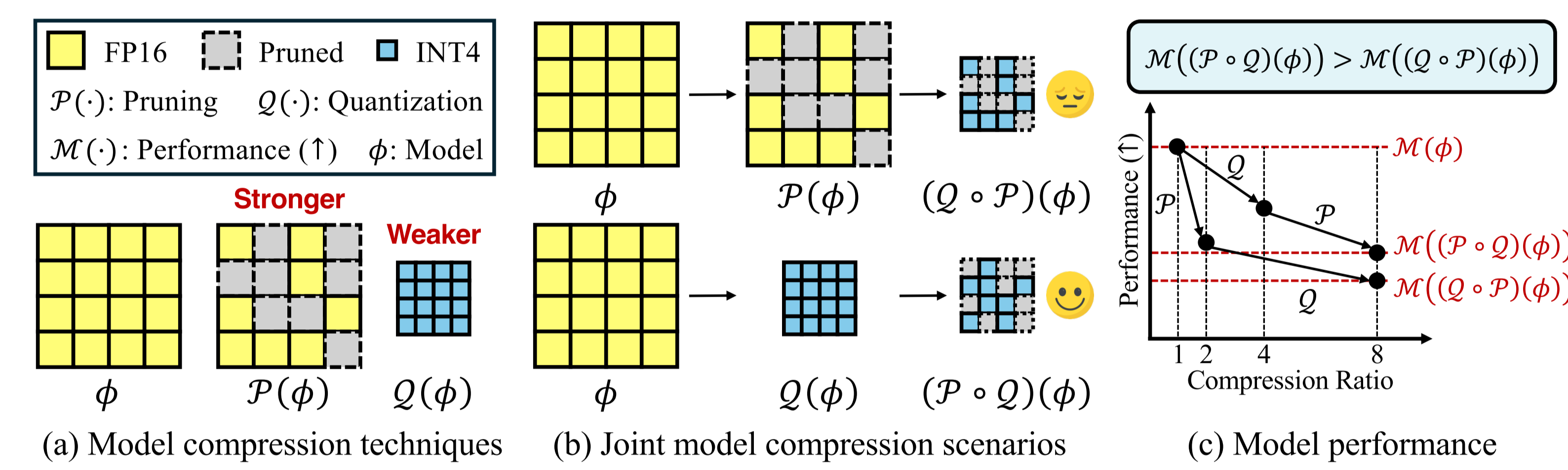
Pruning Type	Order	FP32	INT8	WikiText-2 PPL (↓) of OPT-125M model
Dense model (0% Pruning)		27.65	28.06	
50% (Unstructured)	$\mathcal{P} \rightarrow \mathcal{Q}$		<b>30.22</b>	Harma et al., ICLR 2025
	$\mathcal{Q} \rightarrow \mathcal{P}$	29.94	34.71	
2:4 (Semi-structured)	$\mathcal{P} \rightarrow \mathcal{Q}$	31.89	<b>32.76</b>	Error rate: 18.48%
	$\mathcal{Q} \rightarrow \mathcal{P}$		<b>45.06</b>	

#### Problem. Joint Compression Order Optimization

- **Input:** a set of compression methods  $\mathbb{F} = \{f_1(\cdot), f_2(\cdot), \dots, f_n(\cdot)\}$ , a pre-trained model  $\phi$ , and performance metric  $\mathcal{M}(\cdot) \uparrow$
- **Output:** the optimal permutation  $\pi^* = \operatorname{argmax}_{\pi \in \Pi} \mathcal{M}(\pi(\phi))$  for set  $\Pi = \{\pi: \mathbb{F} \rightarrow \mathbb{F} \mid \pi \text{ is bijective}\}$  of all permutations over  $\mathbb{F}$

### The Progressive Intensity Hypothesis

**The Progressive Intensity Hypothesis.** Neural networks compressed by multiple methods perform better when **weaker** perturbations are applied first and **stronger** ones later.



### Theoretical Analysis

#### Intensity of Compression

- We assess intensity through model performance
- Performance gap  $\mathcal{G}(f_1, f_2) := \mathcal{M}(f_1(\phi)) - \mathcal{M}(f_2(\phi))$
- Compression Equivalence Ratio (CER)  $C_f^*: B_Q$ -axis

$$\mathcal{M}(f(\phi; C)) = \mathcal{M}(Q(\phi; C_f^*))$$

- Compression order advantage  $\mathcal{A}(f_1 \rightarrow f_2)$ : order effect

$$\mathcal{A}(f_1 \rightarrow f_2) := \mathcal{G}(f_1 \rightarrow f_2, f_1 \rightarrow f_2)$$

#### Pruning and Quantization

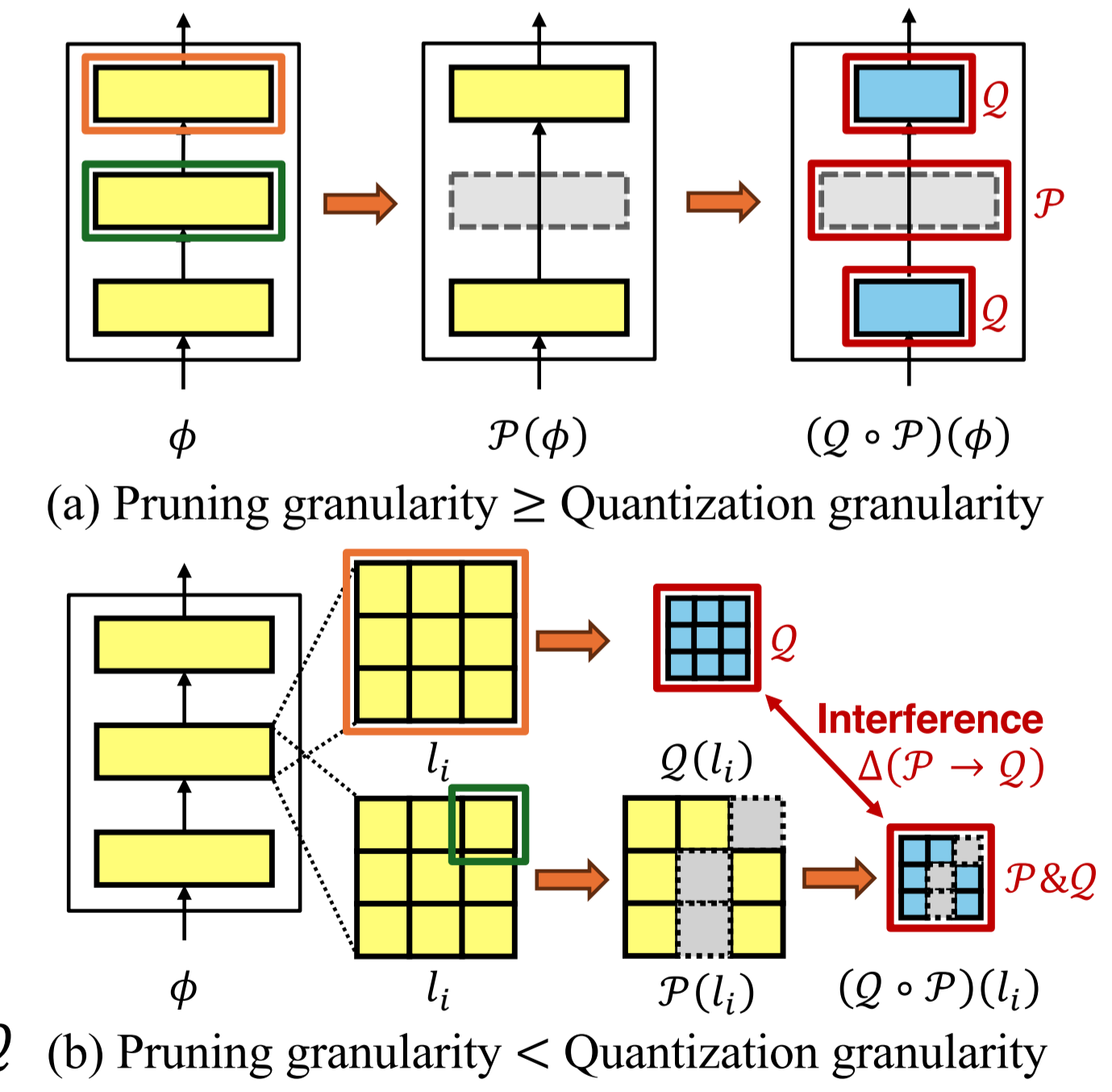
- $\mathcal{A}(Q \rightarrow P)$  grows monotonically with  $\mathcal{G}(Q, P)$  or  $C_P^* - C_Q$

#### Case (a). Disjoint Selectivity

- All units are exclusively assigned to one  $f(\cdot)$
- Layer-wise error analysis: showing that weak  $\rightarrow$  strong ordering consistently yields lower layer-wise degradation (under fixed  $C_P$ )

#### Case (b). Interference

- Additional error from mutual interaction
- The magnitude of interference is determined solely by pruning ratio  $p$  (preserving monotonic trend)



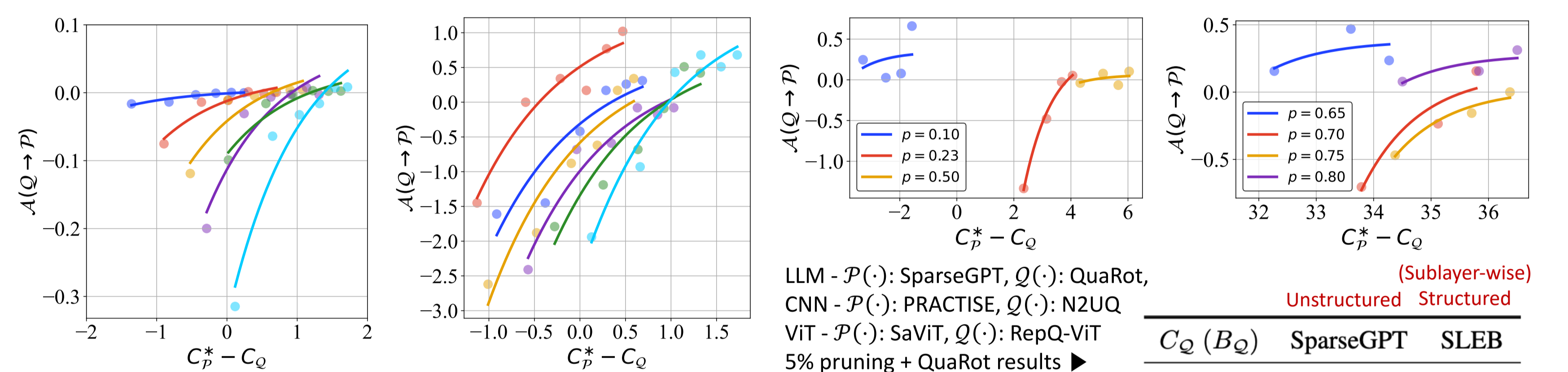
### Experimental Findings

#### 1. The hypothesis holds across diverse scenarios

- Regardless of modality, model architecture, task, intensity, method, and technique
- **2. Pruning granularity defines interference**
- Toward precise sign / value prediction of compression-order advantage  $\mathcal{A}(f_1 \rightarrow f_2)$

#### \* Pruning and Quantization: Language and Vision models

- LLaMA-3 8B (WikiText-2 PPL and ARC-C), ResNet-18 and DeiT-Base (ImageNet)



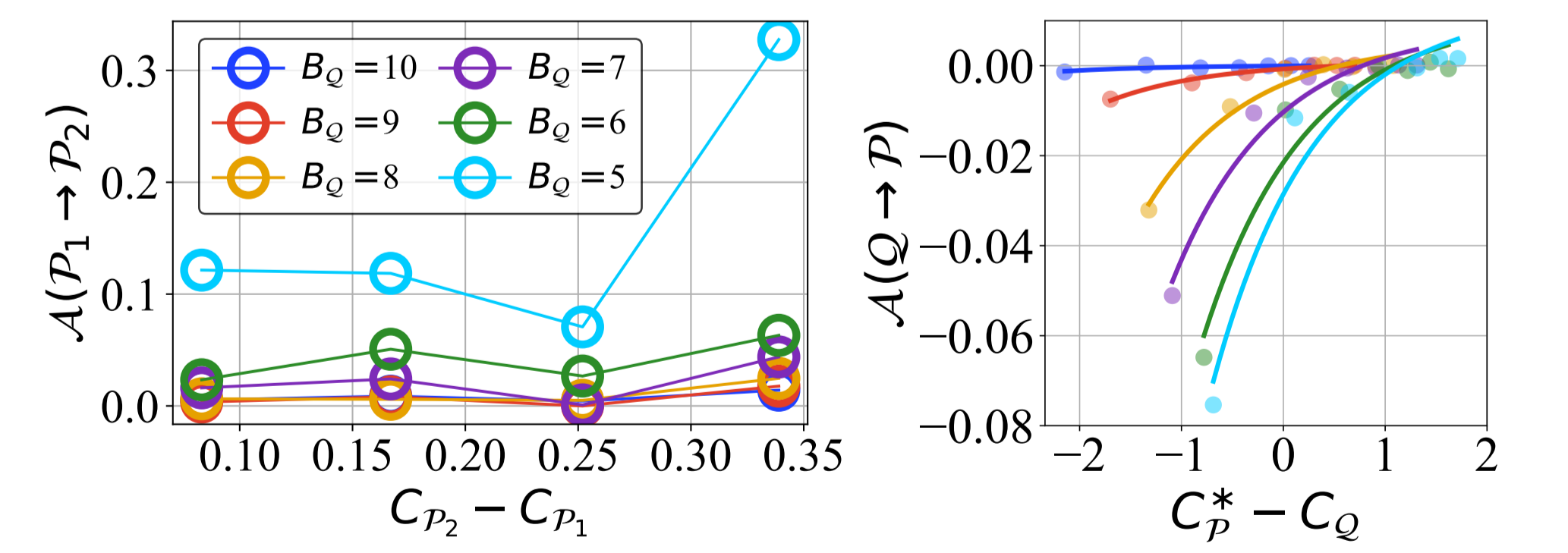
#### \* Pruning granularity and Interference

- Unstructured (SparseGPT): monotonic trend
- Structured (SLEB; sublayer): no interference at early regimes

$C_Q (B_Q)$	SparseGPT	SLEB
1.78 (9)	0.002	0
2.00 (8)	0.001	0
2.28 (7)	-0.003	0
2.68 (6)	-0.013	0
3.20 (5)	-0.017	-0.057
4.00 (4)	-49.899	-9.379

#### \* Toward general pipelines

- **Multi-stage compression.**  $\mathcal{P}_1(\cdot) \rightarrow \mathcal{Q}(\cdot) \rightarrow \mathcal{P}_2(\cdot)$



- **PEFT.**  $\mathcal{P}(\cdot)$  and  $(\mathcal{Q}(\cdot) + \text{LoRA})$

- **Parameter sharing.** Pruning  $\mathcal{P}(\cdot)$  + Basis sharing  $\mathcal{S}(\cdot)$

- **Mixed-precision Quantization.** prog. (8  $\rightarrow$  2 bits) vs. regr. (2  $\rightarrow$  8bits)

