

## The Problem

A fast, calibration-free quantizer that closes the gap to calibrated & non-uniform methods, by adding one extra scale vector.

**Post-training quantization (PTQ)** is the standard way to deploy LLMs cheaply. But at  $\leq 4$  bits perplexity degrades: a few *outlier* weights share a scale with their whole row/column, wrecking precision for everyone else.

The cleanest, fastest recipe (**calibration-free, uniform** quantization) suffers most. Calibrated (AWQ/GPTQ) and non-uniform (NF4) methods help, but add quantization-time cost, overfitting risk, and inference look-ups.

### Contributions

- Column-wise weight std. dev. *predicts* input activation scales  $\Rightarrow$  **calibration-free pseudo-activation-awareness**.
- A **dual-scale** parameterization + a fast **Sinkhorn-Knopp** algorithm that balances rows & columns.
- Beats SoTA calibration-free baselines; composes with calibration (AWQ) and non-uniform levels (NF4); **negligible overhead**, drop-in for any linear layer.

## Key Idea: Dual-Scale Quantization

A single scale vector cannot isolate a 2-D outlier. With scales along **both** axes we trade off error between the outlier's row and column:



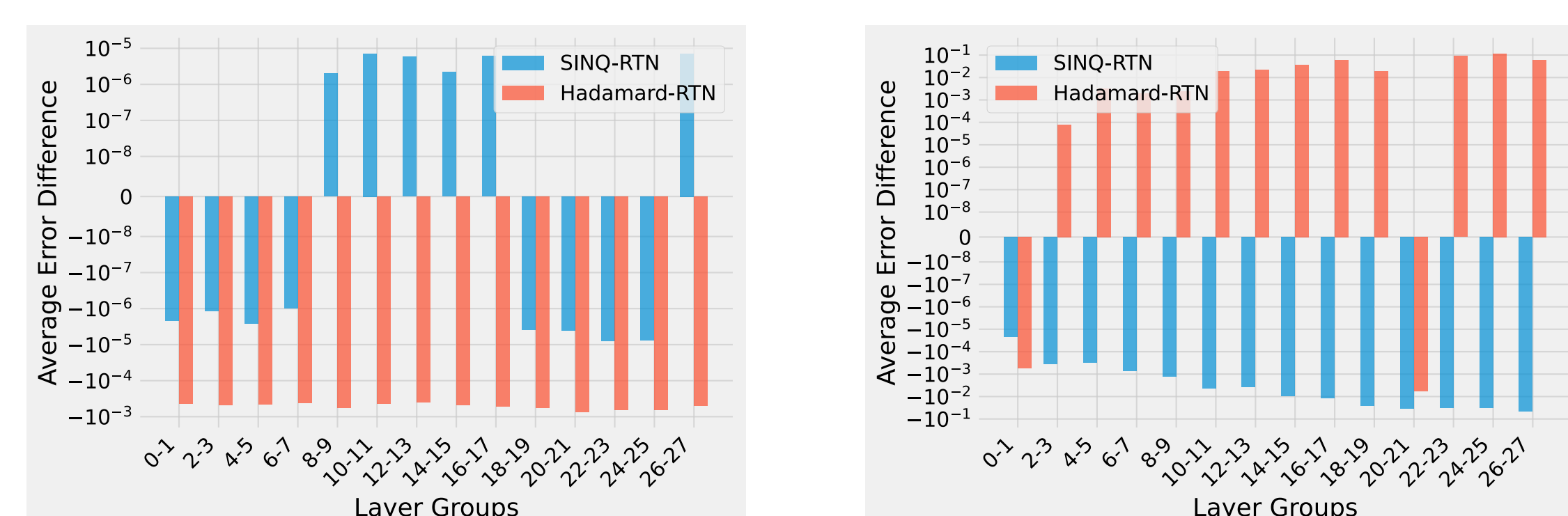
Single scale:  $\mathbf{W} \approx \vec{s} \odot (\mathbf{Q} + \vec{z})$      **Dual scale (ours):**  $\mathbf{W} \approx \vec{s} \odot (\mathbf{Q} + \vec{z}) \odot \vec{t}$

Concretely, for a large outlier  $W_{ij}$ , raising the row scale  $s_i$  and lowering the column scale  $t_j$  shrinks its quantization error and **redistributes the slack** across row  $i$  and column  $j$ , a degree of freedom that single-scale quantization simply does not have.

The second scale  $\vec{t}$  costs only  $M$  extra parameters and, at inference, folds into the **activations** ( $\vec{x} \odot \vec{t}$ ).

## Weights vs. Activations

SINQ minimizes *activation* error, not *weight* error: the opposite of Hadamard rotation, and what matters end-to-end.

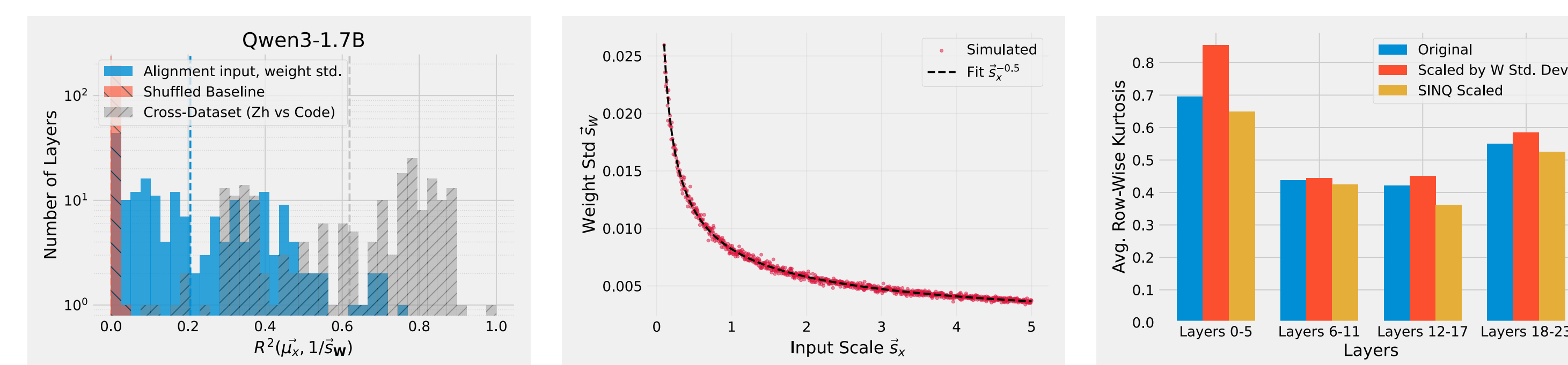


(a) Matrix reconstruction: Hadamard wins

(b) Activation reconstruction: SINQ wins

## Why It Works: Pseudo-Activation-Awareness

The reciprocal column-wise weight std. dev. is a strong predictor of the *average input magnitude*  $\bar{\mu}_x$ , exactly the quantity AWQ needs calibration data to estimate. We recover it **for free** from the weights.



(a) High  $R^2$  between  $\bar{\mu}_x$  and  $1/\bar{s}_w$

(b) Adam training drives  $\bar{s}_w \propto |\bar{\mu}_x|^{-1/2}$

(c) SINQ avoids new row outliers

$$\bar{s}_w \propto \frac{1}{\sqrt{|\bar{\mu}_x|}} \quad (\text{weight updates} \sim \text{outer product of inputs \& gradients})$$

Naively scaling by  $1/\bar{s}_w$  creates *new* row outliers (panel c, red). SINQ scales columns toward activation-awareness while keeping rows tame.

## The SINQ Algorithm

A dampened **Sinkhorn-Knopp** iteration rescales rows and columns toward a common target std. dev.  $\tau$ , minimizing the **matrix imbalance**  $I$  before round-to-nearest:

$$I(\mathbf{W}) = \frac{\max(\max_i \sigma_i^{\text{row}}, \max_j \sigma_j^{\text{col}})}{\min(\min_i \sigma_i^{\text{row}}, \min_j \sigma_j^{\text{col}})}$$

- Multiplicative scales  $\vec{s}, \vec{t}$ ; the clamp damps each step to avoid overshoot.
- Snapshot the scales achieving the **lowest imbalance**.
- Apply any standard quantizer (INT4 / NF4 / ...) to the balanced matrix.
- Recover  $\vec{s}, \vec{t}$ ;  $\sim 1.1 \times$  RTN runtime (vs. HQQ  $> 2 \times$ , AWQ  $> 30 \times$ ).

### SINQ Normalization

```
In:  $\mathbf{W} \in \mathbb{R}^{m \times n}$ , iters  $K$ , bits  $b$ , clamp  $[c_-, c_+]$ 
Out:  $\mathbf{Q}, \vec{s}, \vec{t}$ ; scales  $\vec{s} \in \mathbb{R}^m, \vec{t} \in \mathbb{R}^n$ 

 $\tau \leftarrow \min(\min_i \sigma_i^{\text{row}}, \min_j \sigma_j^{\text{col}})$  target
 $\vec{s} \leftarrow \mathbf{1}_m, \vec{t} \leftarrow \mathbf{1}_n, I^* \leftarrow \infty$ 
for  $k = 1, \dots, K$ :
   $\mathbf{W} \leftarrow \mathbf{W} \odot \vec{s} \odot \vec{t}$  apply scales
  if  $I(\mathbf{W}) < I^*$ :  $I^* \leftarrow I(\mathbf{W}), \vec{s}, \vec{t} \leftarrow \vec{s}, \vec{t}$ 
   $\vec{t} \leftarrow \vec{t} \cdot \text{clamp}_{[c_-, c_+]}(\sigma^{\text{col}}(\mathbf{W})/\tau)$ 
   $\vec{s} \leftarrow \vec{s} \cdot \text{clamp}_{[c_-, c_+]}(\sigma^{\text{row}}(\mathbf{W})/\tau)$ 
 $\mathbf{W} \leftarrow \mathbf{W} \odot \vec{s} \odot \vec{t}$ 
 $\mathbf{Q}, \vec{s}, \vec{t} \leftarrow \text{RoundToNearest}(\mathbf{W}, b)$ 
return  $\mathbf{Q}, \vec{s}, \vec{t}$ 
```

## Fast & (Almost) Free at Inference

**End-to-end decode** (SGLang, W4A16 speedup over FP16):

	Qwen3-14B	Qwen3-32B
FP16	48 tps	21 tps
AWQ	2.4 $\times$	2.9 $\times$
<b>SINQ</b>	2.4 $\times$	2.8 $\times$

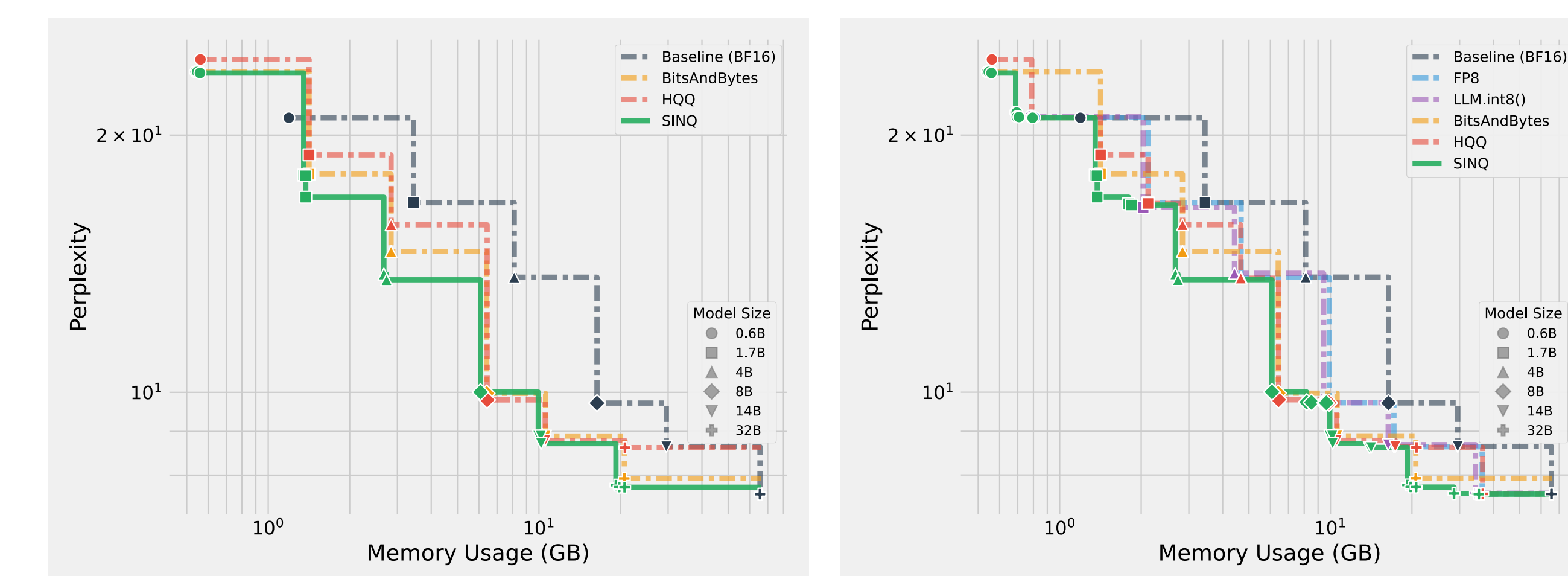
**Kernel overhead** of the 2nd scale (gemlite W4A16):

	B	D	OH
1	1024	1.8%	
1	2048	1.5%	
64	1024	0.8%	

The no-overhead variant folds  $\vec{t}$  into preceding layers for **identical** inference speed, and also improves GGUF / llama.cpp accuracy at full speed.

## Results: Memory-Perplexity Pareto

Because Qwen3 ships in many sizes, we can show the bf16 Pareto front across the whole memory range, from  $\sim 1.5$  GB to 65 GB. At equal memory **SINQ** matches or beats calibration-free baselines almost everywhere. Panel (b) adds 6- and 8-bit points, with LLM.int8() as a reference.



(a) 4-bit methods

(b) Mixed 4/6/8-bit

## Perplexity (WikiText2) ↓

Qwen3, uniform PTQ. Original BF16: 16.67 / 8.64 / 7.60 (1.7B/14B/32B).

	Calibration-free			Calibrated		
	1.7B	14B	32B	1.7B	14B	32B
<b>3-bit</b>						
RTN	32.43	10.50	30.78	32.21	9.54	9.03
HQQ	32.10	10.73	9.09	24.70	9.61	8.51
<b>SINQ</b>	<b>22.39</b>	<b>9.33</b>	<b>8.79</b>	<b>22.30</b>	<b>9.31</b>	<b>8.45</b>
<b>4-bit</b>						
RTN	18.74	8.95	8.92	19.70	8.81	7.80
HQQ	18.96	8.78	8.62	16.90	8.78	7.79
<b>SINQ</b>	<b>17.14</b>	<b>8.76</b>	<b>7.74</b>	<b>16.66</b>	<b>8.71</b>	<b>7.78</b>

## QA: Flip-Rate & Accuracy (4-bit)

Method	Flip-rate % ↓		Accuracy % ↑	
	14B	32B	14B	32B
RTN	4.13	5.26	73.3	75.6
HQQ	4.11	5.33	72.8	74.6
<b>SINQ</b>	<b>3.46</b>	<b>3.60</b>	<b>72.5</b>	<b>75.2</b>
AWQ	3.20	3.69	72.6	74.9
<b>A-SINQ</b>	<b>3.18</b>	<b>3.60</b>	72.7	75.0

<sup>\*</sup> uses calibration. Flip-rate is a more reliable metric than accuracy.

## Takeaways

- Weight structure alone yields **activation-awareness**: no calibration data.
- One extra scale + a few Sinkhorn iterations **halve the perplexity gap** to BF16.
- Composes with calibration (AWQ) and non-uniform levels (NF4).
- Fast to quantize, fast to run, trivial to apply** to any architecture.
- Validated on Qwen3, Llama, Phi, DeepSeek-V3 and MoE models.

Code: [github.com/huawei-csl/SINQ](https://github.com/huawei-csl/SINQ)