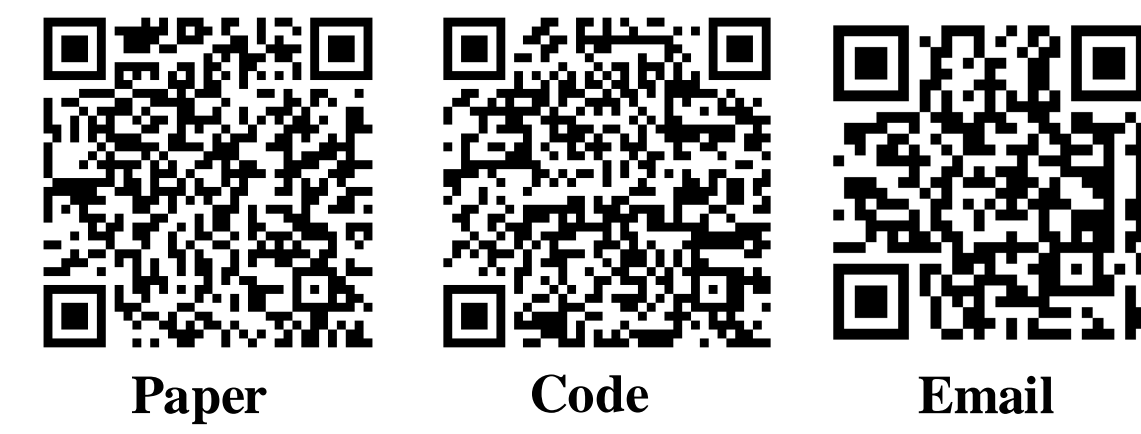


OSTQuant: Refining Large Language Model Quantization with Orthogonal and Scaling Transformation for Better Distribution Fitting

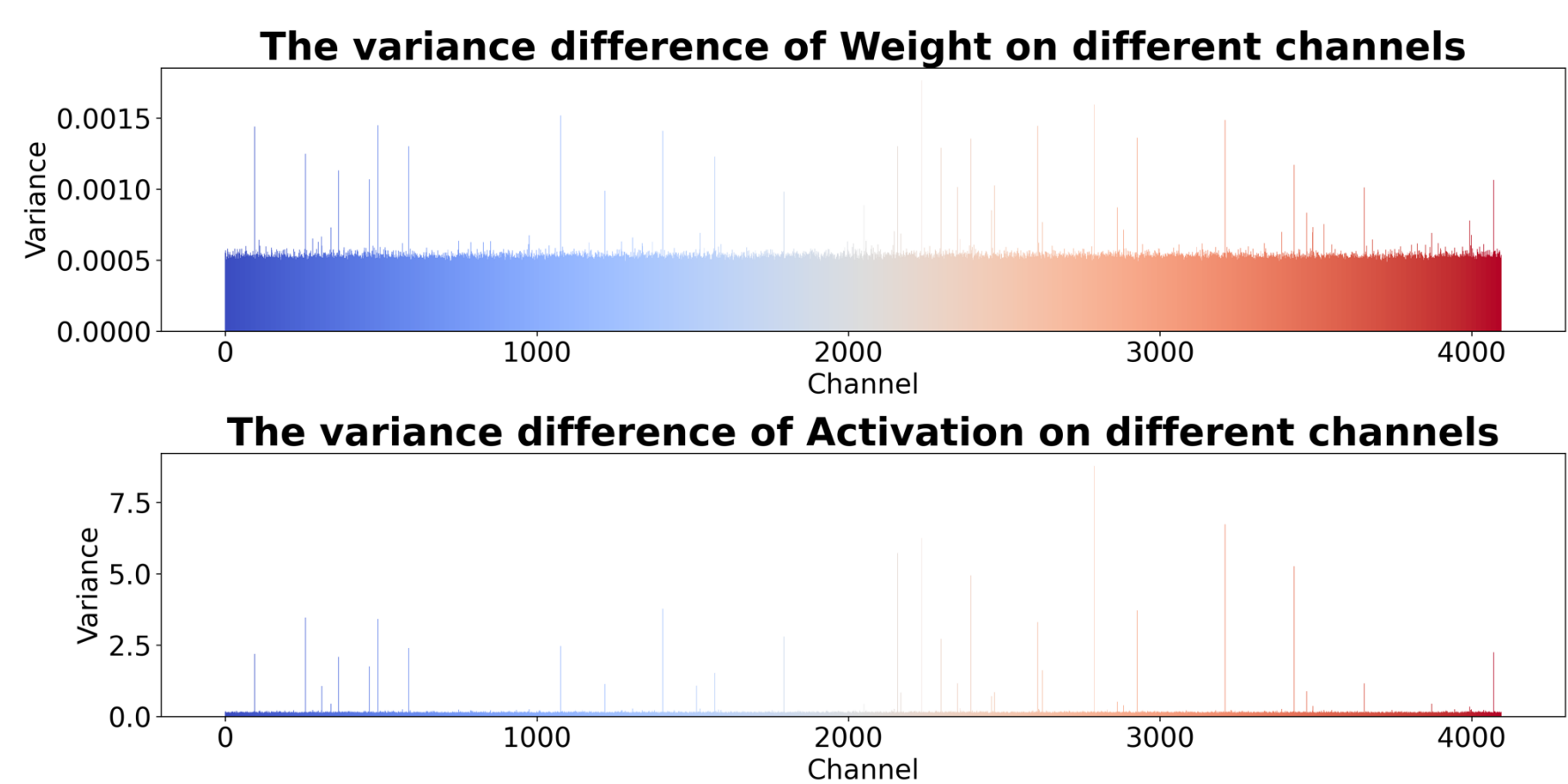
Xing Hu^{1*}, Yuan Cheng^{12*}, Dawei Yang^{1*}✉, Zukang Xu¹, Sifan Zhou³, Jiangyong Yu¹, Chen Xu¹, Zhe Jiang³, Zhihang Yuan¹✉

1. Houmo AI 2. Nanjing University 3. Southeast University



Challenges

Uneven & heavy-tailed distributions in LLMs expands the quantization range, thereby making the quantization for LLMs challenging.



Previous approaches **remain heuristic** and do not optimize the distribution across the entire quantization space.

QSUR

The Quantization Space Utilization Rate(**QSUR**) of X is defined as the ratio of the quantization space V_{SX} to the hypervolume occupied by X :

$$QSUR_X = \frac{V_X}{V_{SX}} \quad (1)$$

For $X \in N(\mu, \Sigma)$, We get:

$$QSUR_X = \frac{\pi^{d/2}}{\Gamma(d/2+1)} \cdot (\chi_d^2(\alpha))^{d/2} \cdot \sqrt{\det(\Lambda)} = \frac{\pi^{d/2}}{\Gamma(d/2+1)} \cdot \sqrt{\prod_{i=1}^d \lambda_i} = \frac{2^d \left(\max \left(\sqrt{\chi_d^2(\alpha)} \cdot \lambda_1 \cdot q_1 \right) \right)^d}{2^d \left(\max(\sqrt{\lambda_1} \cdot q_1) \right)^d}$$

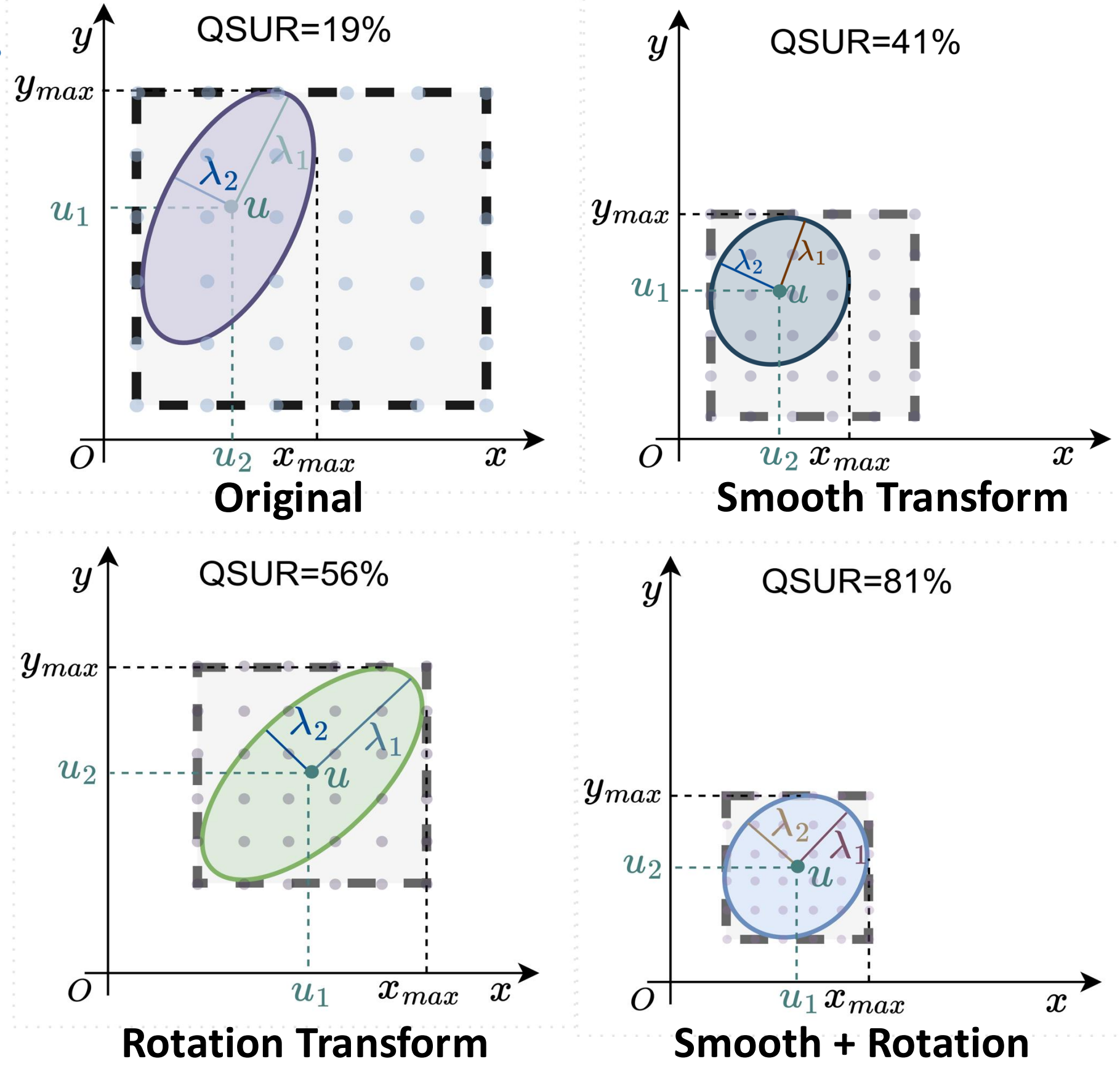
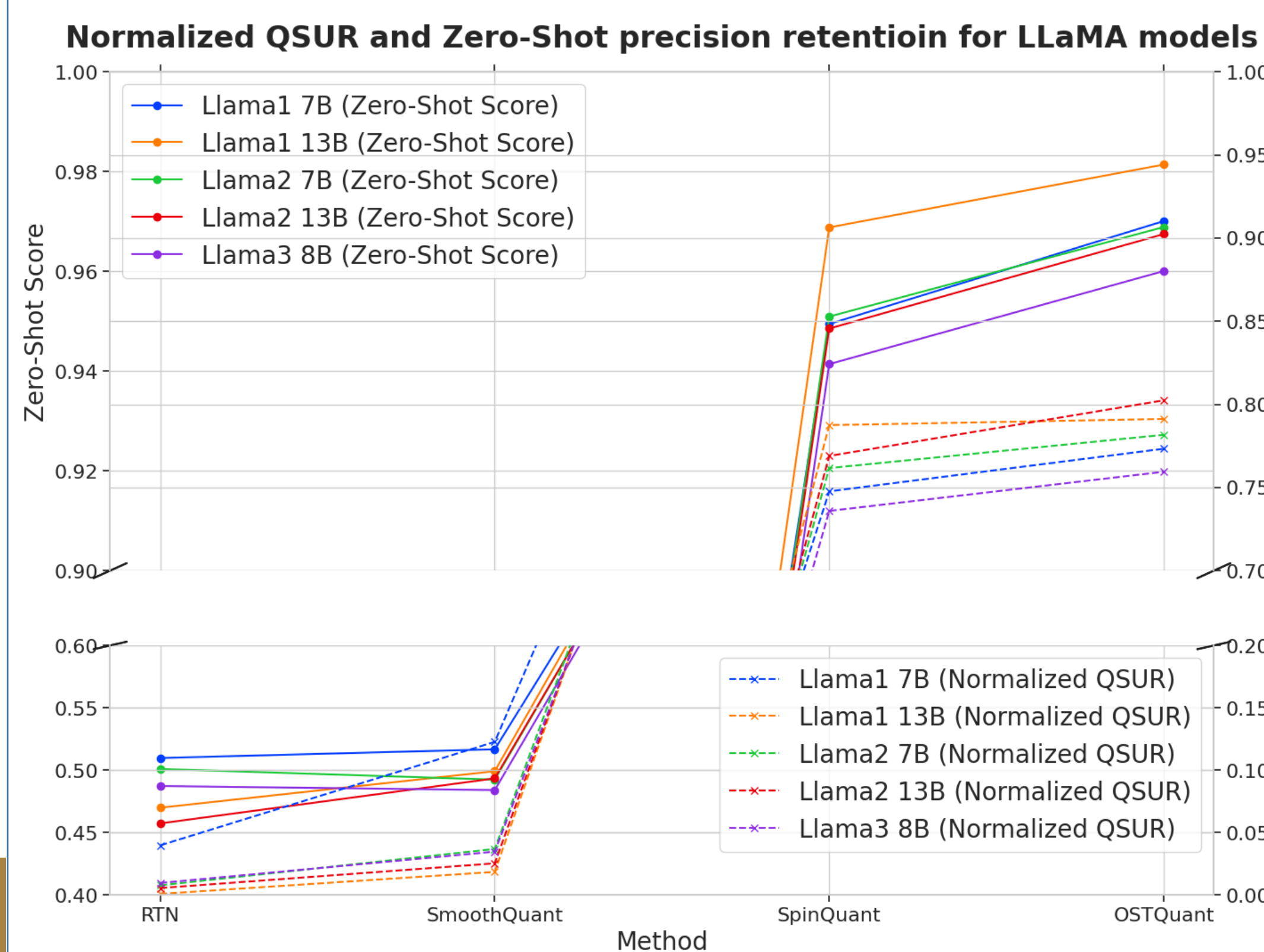
From equation above, We find that:

1) QSUR is proportional to the product of the ratios of each eigenvalue λ_i to the largest eigenvalue λ_1

2) The maximum component of the eigenvector q_1 is inversely proportional to QSUR.

Motivation

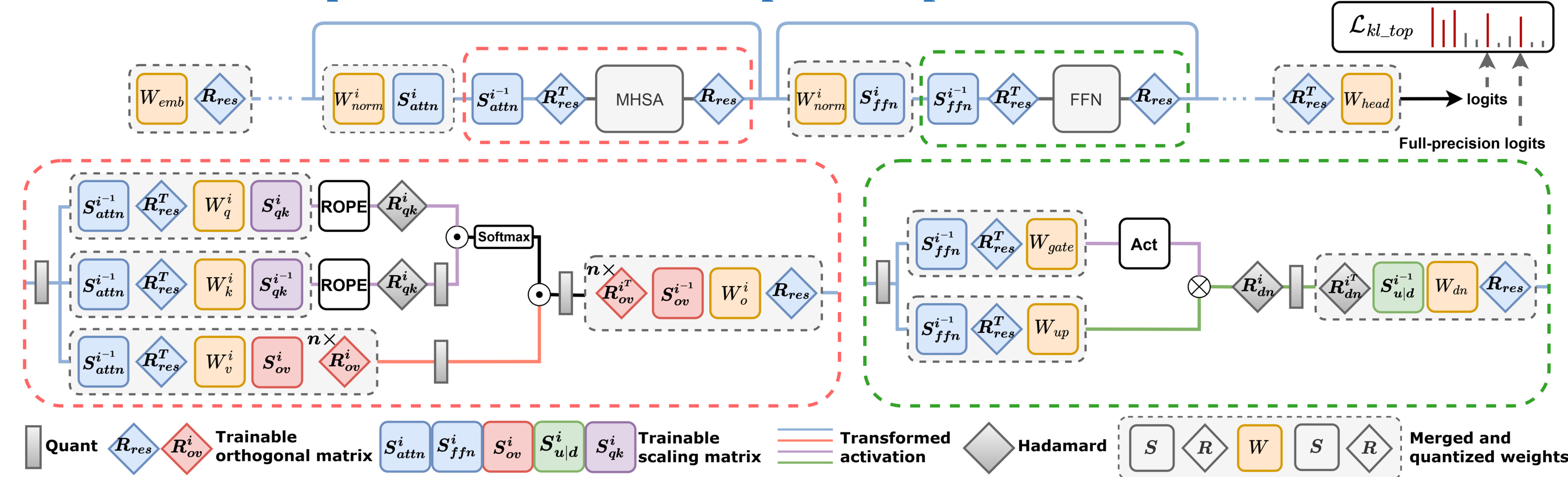
➤ **QSUR exhibits a positive correlation with accuracy.**



➤ **Orthogonal & Scaling enhance QSUR.**

Method

➤ **Use learnable OS equivalent transformation pairs to optimize the distributions.**



Equivalent Transformation Pair $T = A Q^T$, forward like:

$$y = Q(x W_1 O \Lambda) Q(\Lambda^{-1} O^T W_2)$$

Optimization objective:

$$\arg \min_{A_i, O_i} \mathcal{L}(\hat{y}, y; A_i, O_i, \theta)$$

Initialization:

$$O = E_W H \quad A = I$$

Experiments

➤ **Performance across different tasks and models.**

#Bits	Method	LLaMA-3 8B	LLaMA-3 70B	LLaMA-2 7B	LLaMA-2 13B	LLaMA-2 70B	LLaMA 7B	LLaMA 13B	LLaMA 30B
16-16-16	0-shot ⁰	68.09	6.14	73.81	2.86	65.21	5.47	67.61	4.88
	Wiki Avg.(↑)	68.09	6.14	73.81	2.86	65.21	5.47	67.61	4.88
4-16-16	RTN	63.70	8.13	31.15	1e5	61.27	7.02	60.24	6.39
	SmoothQuant	62.79	8.12	67.94	6.70	58.88	8.03	62.03	5.86
	GPQ	61.03	7.43	31.45	9e3	60.86	9.84	64.71	5.79
	OmniQuant	65.66	7.19	-	-	63.19	5.74	66.38	5.02
	AWQ	67.03	7.36	68.92	5.92	63.89	5.83	66.25	5.07
	QuaRot	67.27	6.53	72.93	3.53	64.30	5.62	66.95	5.00
	SpinQuant	66.54	6.49	72.90	3.49	63.59	5.58	67.14	5.00
	OSTQuant	67.80	6.53	73.69	3.19	64.37	5.64	67.31	4.94
	OSTQuant	67.80	6.53	73.69	3.19	64.37	5.64	67.31	4.94
	OSTQuant	67.80	6.53	73.69	3.19	64.37	5.64	67.31	4.94
4-4-16	RTN	33.42	6e2	31.21	8e3	32.44	nan	30.86	8e3
	SmoothQuant	33.04	1e3	34.67	2e2	32.13	nan	34.26	1e3
	GPQ	32.98	5e2	31.47	4e4	32.72	nan	30.11	4e3
	OmniQuant	61.69	8.02	65.56	6.35	61.87	6.05	65.13	5.35
	QuaRot	64.11	7.28	66.99	6.10	67.37	6.78	63.23	5.24
	SpinQuant	65.14	7.24	72.21	3.97	63.90	6.68	66.24	5.14
	OSTQuant	65.14	7.24	72.21	3.97	63.90	6.68	66.24	5.14
	OSTQuant	65.14	7.24	72.21	3.97	63.90	6.68	66.24	5.14
	OSTQuant	65.14	7.24	72.21	3.97	63.90	6.68	66.24	5.14
	OSTQuant	65.14	7.24	72.21	3.97	63.90	6.68	66.24	5.14
4-4-4	RTN	33.18	7e2	30.82	8e3	32.67	nan	30.93	7e3
	SmoothQuant	32.96	1e3	33.76	3e2	32.12	nan	33.36	1e3
	GPQ	33.71	6e2	31.20	4e4	33.52	nan	27.85	5e3
	OmniQuant	32.33	4e2	-	-	48.40	14.26	50.35	12.30
	QuaRot	61.38	8.18	65.33	6.6	61.48	6.11	65.16	5.39
	SpinQuant	64.10	7.35	66.31	6.24	62.01	5.96	64.13	5.74
	OSTQuant	65.37	7.29	71.69	4.01	63.18	5.91	65.41	5.25
	OSTQuant	65.37	7.29	71.69	4.01	63.18	5.91	65.41	5.25
	OSTQuant	65.37	7.29	71.69	4.01	63.18	5.91	65.41	5.25
	OSTQuant	65.37	7.29	71.69	4.01	63.18	5.91	65.41	5.25

➤ **Speedup and Memory saving.**

Model	Decoder Speed (tokens/sec)	Quantized	Speed up	Memory Use (GB)	Quantized	Memory Saving
LLaMA-2-7B	47.32	89.4	1.89x	13.94	4.32	3.23x
LLaMA-3-8B	38.33	77.71	2.03x	15.83	5.88	2.69x
LLaMA-2-13B	23.7	55.35	2.34x	23.7	8.5	2.79x
LLaMA-30B	OOM	30.49	-	OOM	18.19	-
LLaMA-3-70B	OOM	14.68	-	OOM	38.41	-

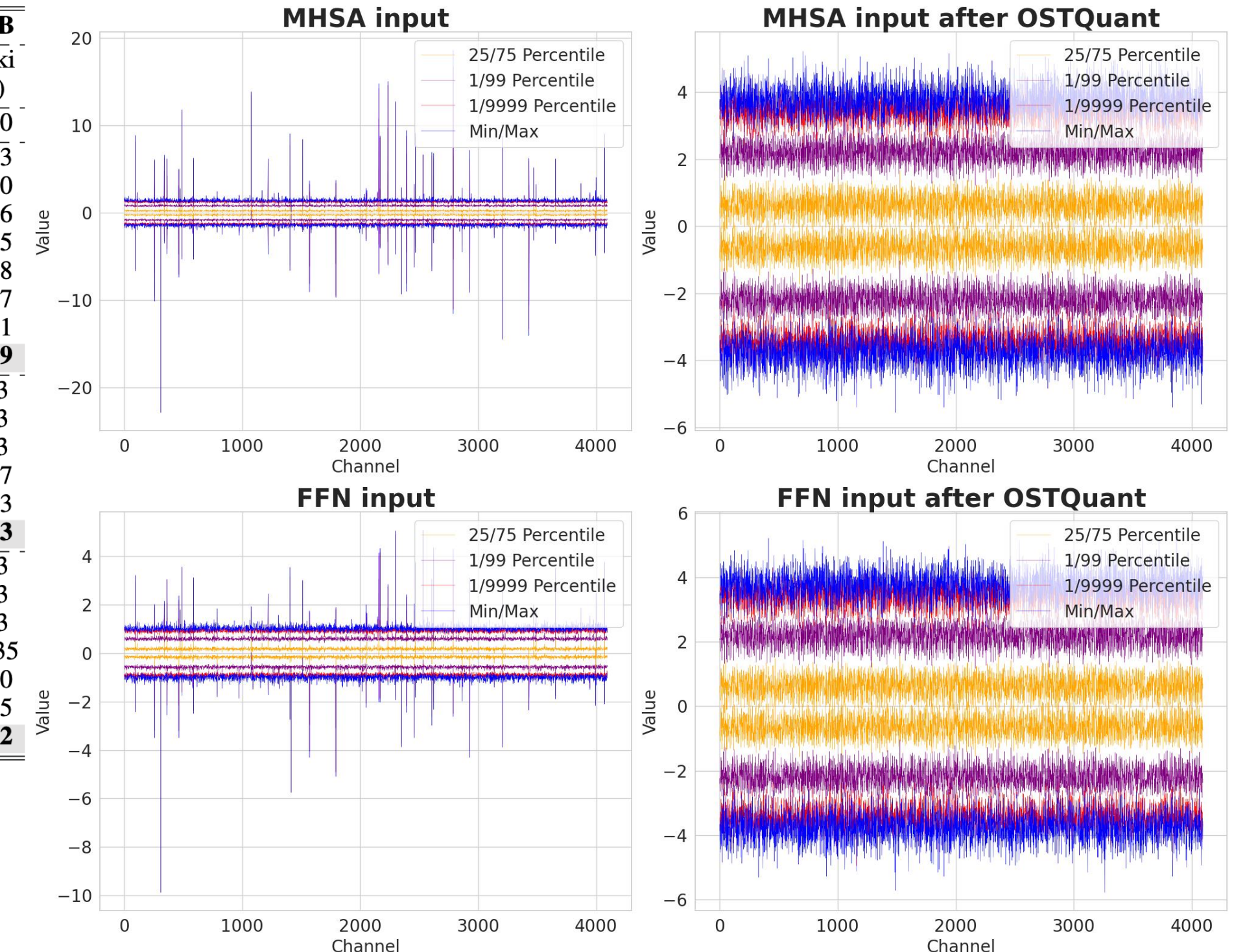
➤ **Ablation Study.**

Ablation on loss function

Model	Loss Type	Wiki PPL	Arc-Easy Score	Arc-Challenge Score
LLaMA-2-7B	Origin	5.38	69.87	42.41
	KL-Top	5.94	72.69	44.62
LLaMA-2 13B	Origin	5.12	75.09	46.08
	KL-Top	5.25	75.29	47.10
LLaMA-3 8B	Origin	6.80	76.68	49.26
	KL-Top	7.29	76.73	49.32

Setting	Metric	k=5	k=50	k=100	k=500	k=1000	k=5000	k=10000
W3 Only	Zero-Shot ⁰ Score	61.87	61.88	61.75	62.18	62.30	61.25	61.21
	Wiki PPL	6.06	6.116	6.13	6.07	6.06	6.06	6.12
W4A4KV4	Zero-Shot ⁰ Score	62.4	62.13	62.38	62.34	63.18	62.44	62.11
	Wiki PPL	5.99	5.96	5.95	5.96	5.96	5.93	5.94

➤ **Distribution pre&after OSTQuant**



Ablation on methods of initialization

Model	Quant Setting	Method	Zero-Shot ⁰	Wiki PPL
LLaMA-2-7B	Full-Precision	-	65.21	5.47
	W4A16KV16	Hadarnard	63.32	5.62
	W4A16KV16	WOMI	63.45	5.59
	W4A4KV4	Hadarnard	61.47	6.11
LLaMA-3-8B	Full-Precision	-	68.09	6.14
	W4A16KV16	Hadarnard	67.27	6.53
	W4A16KV16	WOMI	67.41	6.48
	W4A4KV4	Hadarnard	61.38	8.18
	W4A4KV4	WOMI	61.40	8.17

Ablation on different transformation matrices

Metric	Baseline	+R _{res}	+S _{res}	+R _{dn}	+S _{u/d}	+R _{qk}	+S _{qk}	+R _{ov}	+S _{ov}
Wiki PPL	nan	9.70	9.46	6.16	6.00	5.92	5.92	5.94	5.91
Zero-shot ⁰	33.51	54.33	53.74	61.75	61.79	62.35	62.56	63.11	63.18

Conclusion

➤ **Introduce QSUR as an effective metric and support it to guide optimization and method design.**

➤ **OSTQuant: a fast and effective PTQ method helps quantization by optimizing distributions.**

➤ **SOTA: shows strong performance, maintaining high accuracy even at extremely low bitwidths.**