



SynQ: Accurate Zero-shot Quantization by Synthesis-aware Fine-tuning





Summary

SYNQ (<u>Syn</u>thesis-aware Fine-tuning for Zero-shot <u>Quantization</u>)

- **TL;DR**: We clearly illustrate and address the three major challenges in Zero-shot Quantization (ZSQ)
- **Github**: <u>https://github.com/snudm-starlab/SynQ</u>

Challenges

Challenge 1. Noise in the Synthetic Dataset

Synthetic samples contain **high-frequency noise** unlike real images that concentrate on low frequencies



Magnitude spectrum (log scale) (a) ImageNet dataset





Challenge 2. Predictions based on Off-target Patterns

Quantized models rely on incorrect image patterns for predictions









Challenge 3. Misguidance by Erroneous Hard Labels

• Erroneous hard labels of **difficult samples** lead the quantized model into inaccuracy

Minjun Kim

minjun.kim@snu.ac.kr

Jongjin Kim j2kim@snu.ac.kr

U Kang

ukang@snu.ac.kr

Data Mining _ABORATORY

SEOUL NATIONAL UNIVERSITY

lethod	R-20 (CIFAR-10)		R-20 (CIFAR-100)		R-18 (ImageNet)		R-50 (ImageNet)		MV2 (ImageNet)	
	W4A4	W3A3	W4A4	W3A3	W4A4	W3A3	W4A4	W3A3	W4A4	W3A3
ull Precision (W32A32)	93.89		70.33		71.47		77.73		73.03	
daDFQ (Qian et al., 2023a)	92.31	84.89	66.81	52.74	66.53	38.10	68.38	17.63	65.41	28.99
IAST (Li et al., 2023a)	92.36	86.34	66.68	55.67	66.91	42.58	-	-	65.60	-
exQ (Chen et al., 2023)	<u>92.68</u>	86.47	<u>67.18</u>	55.87	67.73	<u>50.28</u>	70.72	25.27	<u>67.07</u>	<u>32.80</u>
LF (Fan et al., 2024)	92.47	88.04	66.94	57.03	67.02	-	68.97	-	-	-
YNQ (Proposed)	92.76	88.11	67.34	57.28	67.90	52.02	71.05	26.89	67.27	34.21

+ SYNQ (Proposed)	$\textbf{54.97} \pm \textbf{0.35}$	$\textbf{65.88} \pm \textbf{0.27}$	$\textbf{67.42} \pm \textbf{0.21}$	$\textbf{69.88} \pm \textbf{0.19}$	64.54
Genie (Jeon et al., 2023b)	54.01	65.10	66.84	69.66	63.90
Iethod	W2A2	W2A4	W3A3	W4A4	Average

Method	DeiT-Tiny	DeiT-Small	Swin-Tiny	Swin-Small	Average	
Full Precision	72.21	79.85	81.35	83.20	79.15	
PSAQ-ViT (Li et al., 2022) SynQ (Proposed)	$\begin{array}{c} 65.57 \pm 0.10\\ 65.90 \pm 0.07 \end{array}$	$\begin{array}{c} 72.04\pm0.19\\ \textbf{72.28}\pm\textbf{0.34}\end{array}$	69.78 ± 1.67 70.76 ± 1.61	$\begin{array}{c} 75.03 \pm 0.63 \\ \textbf{75.82} \pm \textbf{0.54} \end{array}$	70.61 71.19	
PSAQ-ViT (Li et al., 2022) SYNQ (Proposed)	$\begin{array}{c} 71.56\pm0.03\\ \textbf{71.74}\pm\textbf{0.03}\end{array}$	$\begin{array}{c} 75.97 \pm 0.20 \\ \textbf{76.16} \pm \textbf{0.29} \end{array}$	$\begin{array}{c} 73.54 \pm 1.61 \\ \textbf{74.11} \pm \textbf{1.82} \end{array}$	$\begin{array}{c} 76.68 \pm 0.53 \\ \textbf{77.32} \pm \textbf{0.59} \end{array}$	74.44 74.83	
	Method Full Precision PSAQ-ViT (Li et al., 2022) SYNQ (Proposed) PSAQ-ViT (Li et al., 2022) SYNQ (Proposed)	Method DeiT-Tiny Full Precision 72.21 PSAQ-ViT (Li et al., 2022) 65.57 ± 0.10 SYNQ (Proposed) 65.90 ± 0.07 PSAQ-ViT (Li et al., 2022) 71.56 ± 0.03 SYNQ (Proposed) 71.74 ± 0.03	MethodDeiT-TinyDeiT-SmallFull Precision 72.21 79.85 PSAQ-ViT (Li et al., 2022) 65.57 ± 0.10 72.04 ± 0.19 SYNQ (Proposed) 65.90 ± 0.07 72.28 ± 0.34 PSAQ-ViT (Li et al., 2022) 71.56 ± 0.03 75.97 ± 0.20 SYNQ (Proposed) 71.74 ± 0.03 76.16 ± 0.29	MethodDeiT-TinyDeiT-SmallSwin-TinyFull Precision72.2179.8581.35PSAQ-ViT (Li et al., 2022) 65.57 ± 0.10 72.04 ± 0.19 69.78 ± 1.67 SYNQ (Proposed) 65.90 ± 0.07 72.28 ± 0.34 70.76 ± 1.61 PSAQ-ViT (Li et al., 2022) 71.56 ± 0.03 75.97 ± 0.20 73.54 ± 1.61 SYNQ (Proposed) 71.74 ± 0.03 76.16 ± 0.29 74.11 ± 1.82	MethodDeiT-TinyDeiT-SmallSwin-TinySwin-SmallFull Precision72.2179.85 81.35 83.20 PSAQ-ViT (Li et al., 2022) 65.57 ± 0.10 72.04 ± 0.19 69.78 ± 1.67 75.03 ± 0.63 SYNQ (Proposed) 65.90 ± 0.07 72.28 ± 0.34 70.76 ± 1.61 75.82 ± 0.54 PSAQ-ViT (Li et al., 2022) 71.56 ± 0.03 75.97 ± 0.20 73.54 ± 1.61 76.68 ± 0.53 SYNQ (Proposed) 71.74 ± 0.03 76.16 ± 0.29 74.11 ± 1.82 77.32 ± 0.59	