

Can Less Precise Be More Reliable? A Systematic Evaluation of Quantization’s Impact on CLIP Beyond Accuracy

Aymen Bouguerra,¹ Daniel Montoya,¹ Alexandra Gomez-Villa,² Fabio Arnez¹ Chokri Mraidha¹

¹ Université Paris-Saclay, CEA, List, F-91120, Palaiseau, France

² Computer Vision Center, Barcelona, Spain

{aymen.bouguerra, daniel-alfonso.montoyavasquez, fabio.arnez}@cea.fr, agomezvi@cvc.uab.es

Abstract

The powerful zero-shot generalization capabilities of vision-language models (VLMs) like CLIP have enabled new paradigms for safety-related tasks such as out-of-distribution (OOD) detection. However, additional aspects crucial for the computationally efficient and reliable deployment of CLIP are still overlooked. In particular, the impact of quantization on CLIP’s performance beyond accuracy remains underexplored. This work presents a large-scale evaluation of quantization on CLIP models, assessing not only in-distribution accuracy but a comprehensive suite of reliability metrics and revealing counterintuitive results driven by pre-training source. We demonstrate that quantization consistently improves calibration for typically underconfident pre-trained models, while often degrading it for overconfident variants. Intriguingly, this degradation in calibration does not preclude gains in other reliability metrics; we find that OOD detection can still improve for these same poorly calibrated models. Furthermore, we identify specific quantization-aware training (QAT) methods that yield simultaneous gains in zero-shot accuracy, calibration, and OOD robustness, challenging the view of a strict efficiency-performance trade-off. These findings offer critical insights for navigating the multi-objective problem of deploying efficient, reliable, and robust VLMs by utilizing quantization beyond its conventional role.

1 Introduction

Vision-Language Models (VLMs), particularly CLIP (Radford et al. 2021), have revolutionized computer vision through their remarkable generalization capabilities. Their powerful zero-shot performance has made them a go-to model for safety-related tasks, particularly out-of-distribution (OOD) detection. Consequently, a significant body of work has emerged to rigorously benchmark their reliability, with large-scale evaluations like OpenOOD and ImageNet-X specifically designed to probe the OOD robustness of foundation models (Zhang et al. 2024; Noda et al. 2025; Miyai et al. 2025; Mayilvahanan et al. 2023; Tu, Deng, and Gedeon 2023).

However, the substantial computational requirements of these models pose a significant barrier to real-world deployment. To mitigate this, model quantization has become the

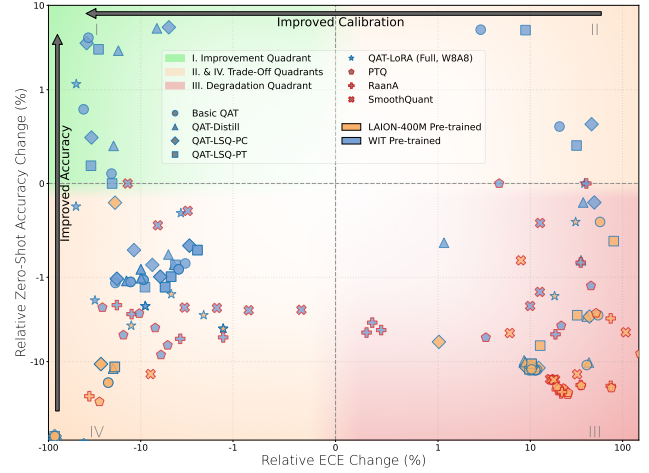


Figure 1: The dichotomous impact of quantization on zero-shot Performance. WIT models (blue) consistently improve in calibration (left), with several QAT methods achieving simultaneous accuracy gains. In contrast, LAION models (orange) show systematic degradation in calibration (right). The lack of points near the origin suggests that quantization is always impactful.

standard compression method, dramatically reducing memory and computational overhead by lowering the precision of the model’s weights (Courbariaux, Bengio, and David 2015; Esser et al. 2020). These two domains (rigorous reliability benchmarking and efficiency-driven quantization) have evolved in parallel, creating a critical blind spot. While the community has extensively characterized the intrinsic robustness of full-precision CLIP, it remains largely unknown if these safety-critical properties survive the aggressive, lossy compression of quantization. Our work directly addresses this gap.

Despite quantization’s popularity, current evaluation practices suffer from a critical limitation: an overwhelming focus on accuracy-based metrics while neglecting essential reliability considerations. This narrow paradigm is particularly problematic for safety-critical applications, where model trustworthiness extends far beyond classification accuracy. The oversight becomes more concerning when compressed

models are deployed in real-time systems, where reliability failures would have immediate consequences.

To address this gap, we move beyond conventional accuracy-focused evaluation and ask:

How does quantization impact VLMs’ reliability?

To answer this question, we systematically evaluate quantization’s impact across four critical reliability dimensions: (1) Robustness to quantization noise, examining stability across bit-widths and methods; (2) Uncertainty quality and calibration, evaluating whether models provide reliable confidence estimates; (3) Out-of-distribution detection, evaluating whether the model can distinguish between two to more semantic distribution of samples; and (4) Distribution shift robustness, assessing performance under realistic data variations.

We evaluate these four dimensions by applying a suite of quantization techniques to CLIP models pre-trained on WIT and LAION, and testing them on rigorous benchmarks designed to assess calibration, OOD detection, and data-shift robustness. Our analysis reveals interesting findings that are partially illustrated in Figure 19, showing that quantization always has a measurable impact beyond accuracy, affecting all reliability attributes, indicating a complex interplay between pre-training source and compression strategy. Our findings challenge the view of quantization as a simple compression or regularization tool, and suggest it can be used as a complex operator that could predictably improve both model efficiency and reliability. In summary, our contributions are threefold:

1. We reveal a fundamental dichotomy where the direct impact of quantization is critically dependent on the model’s pre-training source. The same quantization procedure affects the same VLMs very differently depending on their training data. We also highlight a post-quantization adaptation that drastically improves calibration.
2. We demonstrate a surprising decoupling of reliability metrics, showing that quantization can degrade a model’s accuracy and calibration while simultaneously preserving or even improving its out-of-distribution detection capabilities.
3. We uncover how quantization, and more specifically, quantization-aware training, enables the model to retain coarse-grained features while suppressing the fine-grained ones, increasing robustness to detail-suppressing noise like blur while decreasing it to covariate shift and spurious correlations.

2 Related Work

Quantization of Large Pre-trained Models. Quantizing large models involves a trade-off between data-free Post-Training Quantization (PTQ) (Jacob et al. 2018) and the more accurate but data-dependent Quantization-Aware Training (QAT) (Courbariaux, Bengio, and David 2015). Applying QAT to foundational models requires fine-tuning on training or proxy datasets, creating a significant risk of catastrophic forgetting. While parameter-efficient methods

like LoRA (Hu et al. 2021) help, the core tension between adaptation and forgetting remains an open problem.

Quantization as a Regularizer. Quantization is increasingly understood as an implicit regularizer that can improve generalization by finding flatter, more robust minima in the loss landscape (Hochreiter and Schmidhuber 1997; Tallec, Blier, and Ollivier 2023; Saqib, Hieu, and Mathieu 2025). However, this effect has only been studied on weight-only quantization and exclusively through the lens of accuracy and domain generalization. Its impact on a wider range of reliability metrics, and how it interacts with a model’s pre-training source (e.g., WIT (Radford et al. 2021) vs. LAION (Schuhmann et al. 2022)), remains largely unexplored. In this work, we will consider quantization as a complex operation that can act as a regularizer, and study its impact on several reliability detentions, where each benchmark aims to answer two critical questions: what does quantizing a model do to its attributes, and how can these changes help us understand quantization to the fullest

Benchmarking VLM Reliability. The evaluation of VLM robustness and reliability has matured, with sophisticated benchmarks for OOD detection (OpenOOD) (Yang, Zhou, and Liu 2022; Zhang et al. 2024; Wang et al. 2024) and spurious correlations (CounterAnimal) (Hochlehnert et al. 2025). While these benchmarks have been instrumental in characterizing full-precision models, their application to systematically study the impact of architectural interventions like quantization is nascent. Closely related to our work, Tu, Deng, and Gedeon (2023) evaluate CLIP’s robustness to distribution shifts, OOD detection, and predictive uncertainty. In contrast, in this work, we go one step further, considering the impact of CLIP quantization, a key requirement for computationally efficient deployment.

3 Experiments and Results

Trustworthy AI systems require deep learning models to be developed with a focus on robustness and resilience as a means to ensure a reliable and safe deployment (European Parliament and Council of the European Union 2024; Arnez Yagualca 2023; Hendrycks et al. 2021b). Consequently, we assess the impact of quantization on CLIP’s reliability in zero-shot settings and across four key dimensions: (1) Robustness to quantization noise, (2) Impact on predictive uncertainty, (3) OoD detection, and (4) Robustness to distribution shift. In the first dimension, we examine performance stability across bit-widths and quantization schedules. Next, we assess the CLIP’s capacity to provide reliable confidence estimates by observing the calibration of predicted probabilities. Then we evaluate CLIP’s capacity to detect distribution shifts using its outputs and joint vision-language features. Finally, we examine accuracy and shift-detection under realistic data variations & perturbations.

In our experiments, we quantize only the visual encoder, using the frozen text encoder as a stable semantic anchor to mitigate semantic drift during adaptation (Li et al. 2022b). This allows us to reduce the computational requirements as the text embeddings are only computed once. Crucially, all

downstream classification and OOD evaluation datasets are used exclusively for testing. They are never seen or used during any quantization or adaptation phase, ensuring a strict zero-shot evaluation of all quantization and discrimination methods. We considered two pretrained sources, two visual backbones, 10 quantization methods, and 4 types of weight/activation bit-width across 10 datasets. In total, we conducted 4560 evaluations.

3.1 Experimental Setup

Model Architectures. We conduct our experiments on various pre-trained Open-CLIP models (e.g., ViT-B/32, ViT-L/14) stemming from different pre-training sources (e.g., WIT, LAION).

Quantization strategies. Our evaluation covers multiple quantization techniques from data-free post-training quantization (PTQ) (Jacob et al. 2018; Xiao et al. 2023) to quantization-aware training (QAT) (Courbariaux, Bengio, and David 2015). QAT is performed on the CC3M proxy dataset (Sharma et al. 2018)—distinct from the models’ original WIT or LAION pre-training data—using advanced methods like LSQ (Esser et al. 2020) and LoRA-accelerated distillation (Polino, Pascanu, and Alistarh 2018; Hu et al. 2021). Our primary focus is Light QAT, which uses minimal iterations and samples to adapt to the quantization grid while mitigating overfitting to the proxy data (Kirkpatrick et al. 2017) or forgetting previously acquired knowledge.

Out-of-distribution Detection Methods. We employ a diverse set of OOD scoring functions, categorized into two groups. The first includes classic, output-based methods that operate on the final logits: maximum softmax probability (MSP) (Hendrycks and Gimpel 2017) and the Energy Score (Liu et al. 2020). The second, more sophisticated group consists of VLM-specific methods that operate directly within the joint embedding space. This includes maximum concept matching (MCM), which calculates the maximum similarity score between an image and the set of text embeddings for all in-distribution class names; a low score indicates the image is semantically distant from all known concepts and thus likely OOD (Ming and Li 2022). We also use methods based on negative text prompts, such as a Negative Label (Li et al. 2023) score and a Generic Negative score. These methods gauge novelty by measuring an image’s similarity to explicitly out-of-distribution or generic negative concepts (e.g., “a photo of something”), a technique explored for its robustness in large-scale models

Evaluation Metrics. We employ a suite of standard metrics to assess model performance and reliability:

- **Accuracy:** Standard top-1 classification accuracy.
- **ECE:** Expected Calibration Error, to measure the alignment of confidence with accuracy (Guo et al. 2017).
- **NLL:** Negative Log-Likelihood evaluates the probabilistic correctness of a model’s predictions, heavily penalizing overconfident errors (Bishop 2006). Please refer to the appendix for NLL results.

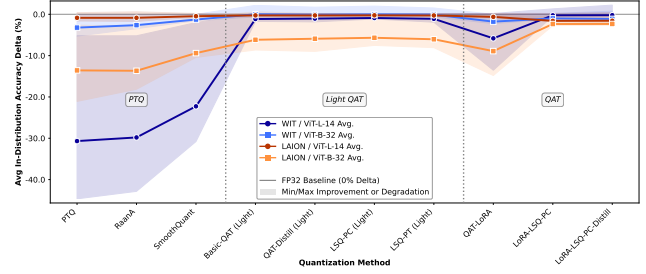


Figure 2: Average In-distribution accuracy change for WIT (blue) and LAION (orange) sources for both ViT/B-32 and ViT/L-14 backbones under various quantization methods, relative to the FP32 baseline (0%).

- **AUROC:** Area Under the ROC Curve, for threshold-independent OOD detection performance (Fawcett 2006).
- **FPR@95TPR:** False positive rate at 95% true positive rate, to evaluate a practical OOD detection operating point (Hendrycks and Gimpel 2017).

Benchmark Datasets Our benchmark datasets, inspired by recent surveys on OOD detection (Yang et al. 2022; Zhang et al. 2024) in VLMs (Yang, Zhou, and Liu 2022; Miyai et al. 2025), categorize datasets in 4 groups along two axes: *covariate Shift* and *semantic Shift*. Please refer to the appendix for an illustration of our experimental protocol.

1. **In-Distribution (ID):** Standard sets including ImageNet-1k (Deng et al. 2009) and CIFAR (Krizhevsky and Hinton 2009).
2. **Data-shifted ID (Covariate Shift):** Datasets with input perturbations, including CIFAR100-C, ImageNet-C (Hendrycks and Dietterich 2019), -R, -A (Hendrycks et al. 2021a), -Sketch (Wang et al. 2019), and -V2 (Recht et al. 2019).
3. **Near-OOD (Semantic Shift):** Datasets with moderate semantic shifts and subsets of Imagenet1k, including ImageNet-X, a challenging benchmark for real-world OOD evaluation (Noda et al. 2025).
4. **Far-OOD (Semantic Shift):** Datasets with distinct categories, including SUN-397 (Xiao et al. 2010), Places365 (Zhou et al. 2017), iNaturalist (Van Horn et al. 2018), and DTD (Cimpoi et al. 2014).

3.2 Robustness to Quantization Noise

Impact on Zero-Shot Accuracy: Model scale inverts the relationship between pre-training noise and quantization robustness. Our results reveal a complex, scale-dependent interplay between the regularization from noisy pre-training data and the subsequent pressure of quantization, a phenomenon reminiscent of the “double descent” behavior where more capacity can paradoxically improve generalization (Nakkiran et al. 2021). At the scale of ViT-B/32, as shown in Figure 2, the LAION-pretrained model (orange) is more fragile to quantization than its WIT counterpart (blue).

We hypothesize that for this smaller model, the regularization from noisy data consumes its limited parameter redundancy. The subsequent application of quantization (Courbariaux et al. 2016) then acts as a second injection of noise, forcing the model past a critical point in the information bottleneck and causing it to lose essential, task-relevant information (Tishby, Pereira, and Bialek 2000).

However, this trend is not universal and, counter-intuitively, reverses in larger models. For a ViT-L/14 model, we observe that the LAION pre-training yields a model **more** robust to quantization. In this high-capacity regime, large neural networks are known to be highly over-parameterized, containing significant redundancy in their weights (Frankle and Carbin 2019). Here, the noisy data acts as a beneficial regularizer, and there remains ample superfluous capacity for quantization to compress without harming the core learned representations. Therefore, the impact of pre-training data on a model’s “quantizability” is not a fixed characteristic but is critically mediated by model capacity.

Less precise can also be more accurate. As illustrated in Figure 2, several quantized models have surpassed their full precision counterparts, improving accuracy beyond the baseline by acting as an implicit sharpness-aware optimizer. Please refer to the appendix for further details about this phenomenon.

Robustness to Quantization Aggressiveness: Quantization methods are not made equal. While our previous results show that accuracy can be recovered by most methods with ease for 8-bit quantization, we find it extremely hard to replicate using lower precisions. As shown in Figure 3, all models degrade as precision decreases, but methods that optimize the quantized representation range are highly more resilient to extremely low precision. Simpler QAT techniques, like basic ‘QAT-LoRA’, are brittle and fall off a “quantization cliff” at 4-bit precision, with performance collapsing to near-zero. The performance inversion when using advanced methods is puzzling: the LAION model, which was less robust at 8-bits, now significantly outperforms the WIT model when quantized to 4-bits. This suggests that while the WIT model is easier to compress naively, the LAION model’s representations, when guided by a powerful regularizer like distillation, can be molded into a more fundamentally robust low-bit state, overcoming its initial fragility.

Robustness to Catastrophic Forgetting during QAT: The optimal QAT duration is a trade-off between adapting to quantization and catastrophic forgetting on the proxy dataset. A core challenge in deploying large pre-trained models is that the original training data is often inaccessible, forcing any adaptation to be performed on a smaller proxy dataset. This turns QAT into a high-stakes balancing act. As shown in Figure 4, the optimal number of training iterations is dictated by the tension between two opposing forces: adapting to quantization noise and forgetting the model’s original general-purpose knowledge by overfitting to the narrow proxy dataset. Although we achieve better performance on ‘w8a8’ using light QAT, adaptation using more diverse data and for longer iterations becomes required as

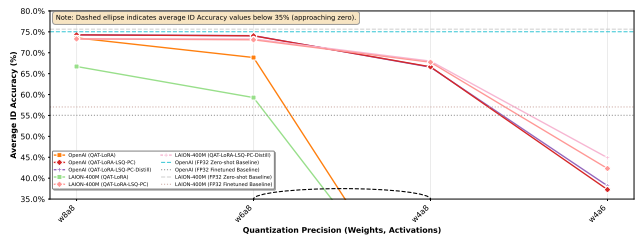


Figure 3: Robustness to Decreasing Quantization Precision. Average ID accuracy vs. bit-width. While simpler QAT methods collapse at 4-bit precision, advanced methods are more robust.

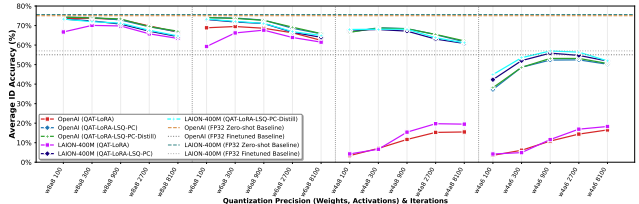


Figure 4: Accuracy evolution of quantized model accuracy relative to iteration steps; higher precision models exhibit catastrophic forgetting behavior while lower precision models struggle to recuperate.

we reduce bit-width precision.

At higher precision (e.g., ‘w8a8’), the quantization noise is minimal. The model requires only a short QAT schedule to adapt its weights. Further training provides diminishing returns and increases the risk of catastrophic forgetting, explaining the quick performance plateau, then degradation.

In contrast, at ultra-low precision (e.g., ‘w4a6’), the aggressive quantization effectively destroys the information encoded in the weights, inducing a state of precision-based amnesia. The model has insufficient parameter redundancy to function. In this regime, a long QAT schedule is mandatory for recovery. The model must use the proxy dataset to re-learn its function from scratch within the highly constrained low-bit space. This difficult search through a complex optimization landscape requires thousands of iterations to succeed (Gong et al. 2019). The upward curve shows that recovering from precision-induced forgetting is the dominant and necessary effect, even at the cost of specializing to the proxy dataset.

3.3 Uncertainty quality and calibration

Direct quantization’s impact on calibration is directly tied to its initial, pre-quantization calibration state Our findings reveal a stark dichotomy in how quantization impacts model calibration, a key measure of reliability. As shown in Figure 5, QAT consistently **improves** calibration for WIT-trained models (15% ECE improvement on average when using basic QAT), where its regularizing effect (Tallec, Blier, and Ollivier 2023) appears to correct for under-confidence. Conversely, the same methods systematically **degrade** calibration for LAION-trained models, suggesting

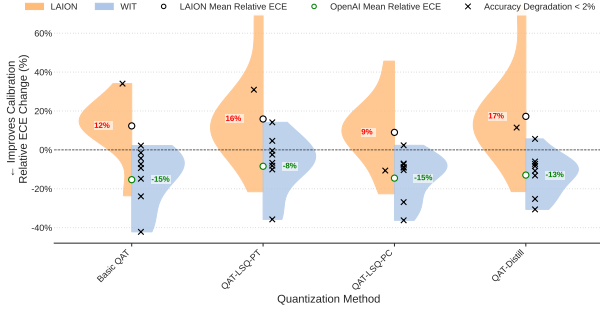


Figure 5: Impact of QAT Methods on CLIP Model Calibration on ViT-B/32. Violin plots show Relative ECE Change (%), comparing LAION (left, blue) and WIT (right, orange) pre-training. Negative values signify calibration improvement. Black crosses represent runs with less than 2% accuracy degradation. Please refer to our appendix for ViT-B/16 and ViT-L/14 results.

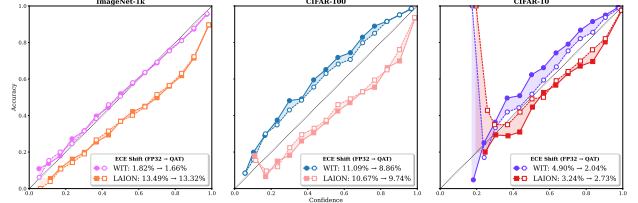
the added quantization noise exacerbates poor priors learned from web-scale data. This effect holds even for models with no accuracy loss (black crosses), exposing a critical reliability trade-off for safety-critical applications where predictive uncertainty is paramount.

QAT’s optimization in the quantized domain directly reshapes the final logit distribution, but the outcome is fundamentally dependent on the model’s initial state. TQAT minimizes task loss under quantization noise, directly altering the final logits that determine confidence. For underconfident WIT models, this optimization beneficially sharpens the output distribution, increasing confidence to match empirical accuracy, as seen in Figure 6.

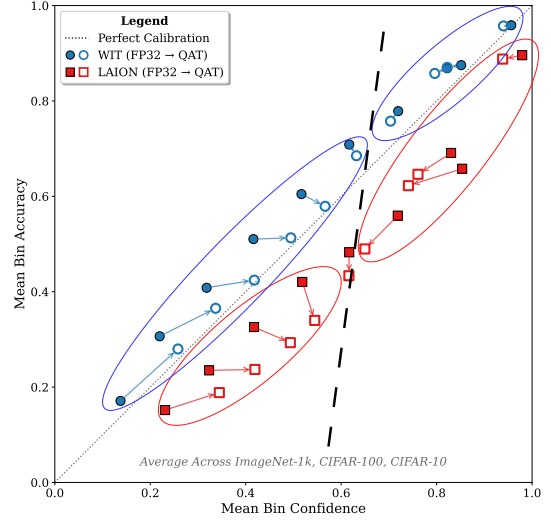
Conversely, this can be detrimental for LAION models. While quantization correctly regulates the highest-confidence logits, it perversely boosts the logits of already overconfident low-to-mid certainty predictions. This confidence distortion is not an artifact of accuracy loss, as it persists even when accuracy is maintained at larger scales. This demonstrates that the regularization from quantization is not a simple smoothing but a complex reshaping of the logit structure, contingent on the model’s initial properties.

A standard reliability diagram is insufficient to reveal these dynamics, as it compares bin populations that are not composed of the same samples before and after QAT. The bin-wise evolution in Figure 6(b) is thus essential, as it explicitly follows the same initial groups of samples to deconstruct these internal, underlying transformations.

A re-adaptation of the logit scale is crucial to correct the logit magnitude distortion introduced by QAT. While QAT adapts the weights, it simultaneously alters the dynamic range of the final logit vectors, rendering the pre-trained logit scale suboptimal. This mismatch is a known issue in large pre-trained models, which often require recalibration before deployment (Desai and Durrett 2020). A second recalibration step is therefore necessary. This step is particularly effective for the LAION model. As QAT distorts



(a) Reliability Curves

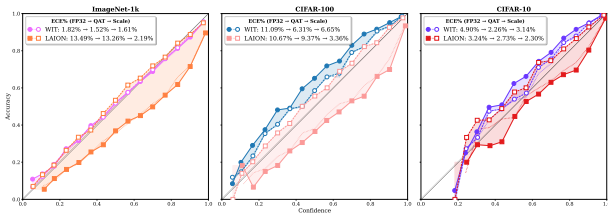


(b) Aggregated Bin-Wise Shift

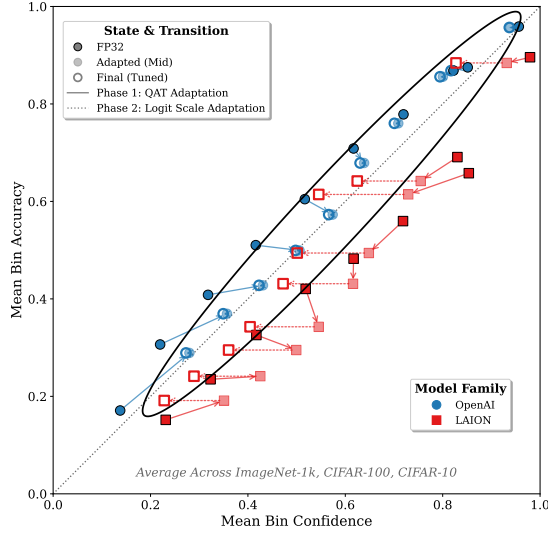
Figure 6: Direct Impact of QAT on Calibration. These plots show the state after only QAT has been applied. The process performs a coarse correction, improving calibration by squeezing the confidences towards the center axis. Please refer to our appendix for the dataset-specific bin-wise shift.

their logit magnitudes, exacerbating overconfidence without necessarily altering the ‘argmax’ or the relative ordering of the logits. Logit scale adaptation directly targets and corrects this magnitude distortion. By learning a single temperature parameter (Kull et al. 2019), we drastically reduce both the inherent overconfidence as well as the one introduced by QAT and pull the bins back towards the diagonal. Figure 7 shows the final, superior state after this adaptation, demonstrating its necessity for achieving a well-calibrated model. Note that we only use ID data or classification datasets for evaluations.

QAT squeezes confidence, while logit scale tuning adapts them to the quantized model’s accuracy The aggregated plots in Figure 6.b summarize the two-phase re-calibration. QAT’s regularization acts as a structural operator on the logit vector. The optimization under quantization noise penalizes extreme logit values; large positive logits, which create high confidence, are brittle to weight perturbations, while near-zero logits, creating low confidence, are unstable. QAT finds a robust middle ground by attenuating large logits and amplifying small ones, resulting in the observed confidence-



(a) Reliability Curves



(b) Aggregated Bin-Wise Shift

Figure 7: Final State after QAT & Logit Scale Re-Adaptation. After adapting the logit scale to the new quantized model, calibration is further improved. Please refer to our appendix for the dataset-specific bin-wise shift.

squeezing effect towards the dashed line in Figure 6(b).

However, this mechanical confidence-squeezing is agnostic to the model’s actual accuracy, leading to a new, often mismatched, calibration state as seen in Figure 6(b). This is where logit scale adaptation becomes critical. It is this final, global correction that aligns confidence with accuracy, moving the model to the superior state shown in Figure 7(b).

3.4 Impact on Out-of-Distribution Detection

ID Accuracy does not necessarily correlate with better OOD detection. The apparent decoupling of OOD performance from accuracy is explained by where in the model the OOD score is calculated. As shown in Figure 9, methods that rely on the final output confidence, like MSP and Energy (square marker), fail catastrophically for the quantized LAION model. This is a direct consequence of our earlier findings: because QAT severely distorts the LAION model’s output logits, the core assumption of these methods—that OOD samples will have lower confidences (Hendrycks and Gimpel 2017).

In stark contrast, methods that operate on the model’s

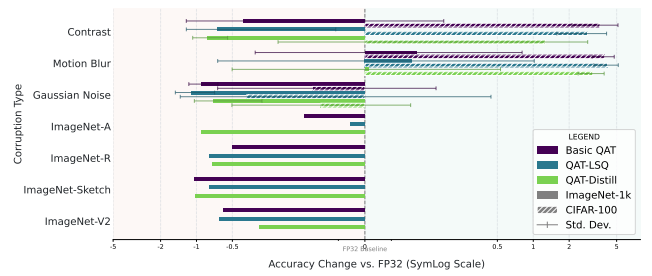


Figure 8: Divergent Impact of QAT on Robustness to Covariate Shift. QAT improves robustness on low-level shifts but degrades it on complex ones like ImageNet-V2, -S, -R and -R.

deeper representations remain effective. The resilience of the VLM methods strongly suggests that even if quantization degrades the quality of the max logit magnitudes, it largely preserves the information in deeper feature space. We also note that quantization can improve out-of-distribution detection, illustrated by the colored vertical lines exceeding the full precision baseline. Please refer to the appendix for FPR@95 results.

3.5 Robustness to Covariate Shift

QAT enhances robustness to low-level, information-suppressing corruptions but degrades generalization on more complex covariate shifts. To evaluate how quantization affects robustness to covariate shifts. The results, shown in Figure 8, vary based on the nature of the distribution shift.

For common, low-level corruptions like blur and noise (top three rows) on CIFAR, QAT acts as a powerful regularizer. It not only improves the model’s accuracy on these corrupted images (Figure 8) but also substantially boosts its ability to detect them as out-of-distribution (please refer to the appendix for the data-shift detection results). This “rare double win” as these goals are contradictory. These results suggest that the QAT process makes the model’s representations less sensitive to simple input perturbations, a known benefit of certain regularization techniques.

However, this trend reverses for more complex and challenging datasets like ImageNet, and on semantic-style covariate shifts like ImageNet-R (renditions) and ImageNet-A (adversarial examples). For these datasets (bottom four rows), QAT consistently slightly degrades the model’s accuracy. This can be explained by its tendency to ignore fine-grained detail, whether they are meaningful or not, and focus on coarse-grained features that could penalize generalization, which is not a universal robustness enhancer, and its benefits are specific to the type of distribution shift encountered.

Robustness to Spurious Correlations: Quantization amplifies reliance on spurious correlations by degrading fine-grained shape features, forcing the model to rely on coarse-grained textures. This occurs because quantization’s limited bit-depth creates an information bottleneck.

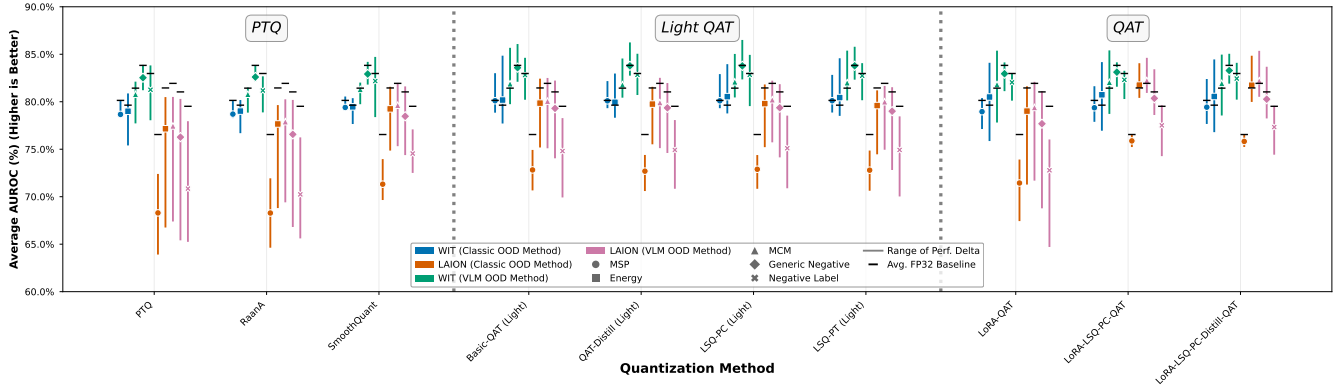


Figure 9: Impact of quantization on OOD Detection (AUROC). Average AUROC across quantization methods (higher is better). QAT methods (center, right) maintain OOD performance for the LAION model, despite this model suffering from significant accuracy and calibration degradation. VLM-specific OOD methods consistently outperform classic methods. Vertical lines represent the maximum improvement and degradation relative to the full precision baseline of that experiment.

It disproportionately degrades the high-frequency signals corresponding to fine-grained object shapes while preserving the more robust, low-frequency signals of coarse background textures. When QAT optimizes the model under this constraint, it finds the path of least resistance: it learns to rely more heavily on the now-dominant texture cues to minimize loss on the proxy dataset. This active re-weighting of feature importance towards texture bias explains why the model’s performance collapses when the spurious background cue is removed (Teney, Abbasi, and van den Hengel 2022). We test this vulnerability using datasets with spurious animal-background correlations, and then with uncommon backgrounds (e.g., a polar bear in the jungle), CounterAnimals (Hochlehnert et al. 2025). As shown in Table 1, QAT consistently increases the model’s vulnerability, making it more easily fooled.

To verify these claims, we analyze the 2D Fourier spectrum of the model’s feature maps. As shown in Figure 10, the full-precision (FP32) spectrum contains significant energy in both the low-frequency (coarse-grained, center) and high-frequency (fine-grained, outer) bands. Both PTQ and QAT act as low-pass filters, suppressing the magnitude of the high-frequency components, which correspond to fine-grained textures and sharp edges. The relative error plots confirm this, showing that the quantization error is highest in these high-frequency bands. This provides an explanation for our earlier finding on spurious correlations and data-shift robustness: Quantization forces the model to rely on the more general, less-changing features.

4 Conclusion

This work reframes quantization not as a simple compression tool, but as a complex operator with profound and often multidimensional impacts on model reliability. Our primary finding is that a model’s pre-training source is a fundamental determinant of its response to quantization, creating a stark dichotomy where the same technique can improve the reliability of one model while severely degrading another.

Model	FP32 Vuln.	QAT Vuln.	Added Vuln.
WIT	16.7%	17.9%	+1.2%
(Normal → Counter)	(83.1 → 66.4)	(82.5 → 64.6)	
LAION	23.1%	24.8%	+1.7%
(Normal → Counter)	(84.0 → 60.9)	(82.5 → 57.7)	

Table 1: Quantization exacerbates vulnerability to spurious correlations. The final column isolates the additional accuracy drop due to QAT, a direct measure of increased vulnerability.

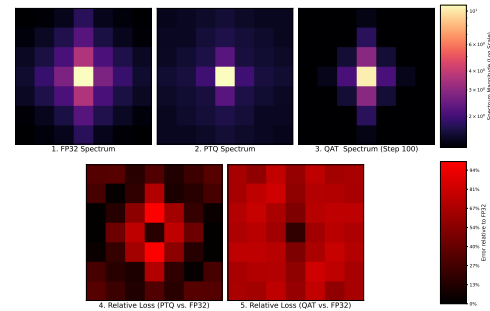


Figure 10: Frequency-domain impact of quantization. Top: The Fourier spectrum of feature maps for FP32, PTQ, and QAT. Bottom: The relative error introduced by PTQ and QAT, respectively.

We demonstrate that this manifests as a critical trade-off: by depressing fine-grained learned features, quantization can improve robustness to high-frequency noise while simultaneously amplifying reliance on spurious correlations and degrading generalization on complex covariate shifts. Furthermore, we reveal a surprising decoupling of reliability metrics, where a model’s calibration can be severely degraded even as its OOD detection capabilities remain robust.

This work encourages the community to advance quantization beyond a simple compression method and toward a multi-objective optimization tool. Our findings reveal that quantization’s impact is not a fixed trade-off but a malleable outcome dependent on the method and model. This presents an opportunity to develop new quantization strategies that are not merely designed to preserve accuracy but are explicitly optimized to simultaneously enhance both model efficiency and reliability.

References

- Arnez Yagualca, F. A. 2023. *Deep neural network uncertainty runtime monitoring for robust and safe AI-based automated navigation*. Theses, Université Paris-Saclay.
- Bishop, C. M. 2006. *Pattern recognition and machine learning*. Springer.
- Bondarenko, Y.; Chiaro, R. D.; and Nagel, M. 2024. Low-Rank Quantization-Aware Training for LLMs. arXiv:2406.06385.
- Cimpoi, M.; Maji, S.; Kokkinos, I.; Mohamed, S.; and Vedaldi, A. 2014. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 3606–3613.
- Courbariaux, M.; Bengio, Y.; and David, J.-P. 2015. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems (NIPS)*, volume 28.
- Courbariaux, M.; Hubara, I.; Soudry, D.; El-Yaniv, R.; and Bengio, Y. 2016. Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1. In *Advances in neural information processing systems (NIPS)*, volume 29.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Desai, S.; and Durrett, G. 2020. Calibration of Pre-trained Transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 3979–3991.
- Esser, S. K.; McKinstry, J. L.; Bablani, D.; Mallya, A.; Appuswamy, R.; and Rath, D. 2020. Learned step size quantization. In *International Conference on Learning Representations (ICLR)*.
- European Parliament and Council of the European Union. 2024. Artificial Intelligence Act. <https://artificialintelligenceact.eu/fr/article/15/>. Regulation (EU) 2024/1689. Specifically referencing Article 15 on ‘Accuracy, robustness and cybersecurity’. Accessed: 2025-08-01.
- Fawcett, T. 2006. An introduction to ROC analysis. *Pattern recognition letters*, 27(8): 861–874.
- Frankle, J.; and Carbin, M. 2019. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations (ICLR)*.
- Gong, R.; Liu, X.; Jiang, S.; Li, T.; Fua, P.; and Yan, S. 2019. Differentiable soft quantization: Bridging full-precision and low-bit neural networks. In *Proceedings of the IEEE/CVF international conference on computer vision (ICCV)*, 4852–4861.
- Guo, C.; Pleiss, G.; Sun, Y.; and Weinberger, K. Q. 2017. On calibration of modern neural networks. In *International conference on machine learning (ICML)*, 1321–1330. PMLR.
- Hendrycks, D.; Basart, S.; Mu, N.; Kadavath, S.; Wang, F.; Dorundo, E.; Desai, R.; Zhu, T.; Parajuli, S.; Hvilshoj, M.; et al. 2021a. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 8340–8349.
- Hendrycks, D.; Carlini, N.; Schulman, J.; and Steinhardt, J. 2021b. Unsolved problems in ml safety. *arXiv preprint arXiv:2109.13916*.
- Hendrycks, D.; and Dietterich, T. 2019. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations (ICLR)*.
- Hendrycks, D.; and Gimpel, K. 2017. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations (ICLR)*.
- Hochlehnert, A.; Bhatnagar, H.; Udandara, V.; Albanie, S.; Prabhu, A.; and Bethge, M. 2025. A Sober Look at Progress in Language Model Reasoning: Pitfalls and Paths to Reproducibility. *arXiv preprint arXiv:2504.07086*.
- Hochreiter, S.; and Schmidhuber, J. 1997. Flat minima. *Neural computation*, 9(1): 1–42.
- Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv preprint arXiv:2106.09685*. Published at the International Conference on Learning Representations (ICLR) 2022.
- Jacob, B.; Kligys, S.; Chen, B.; Zhu, M.; Tang, M.; Howard, A.; Adam, H.; and Kalenichenko, D. 2018. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2704–2713.
- Kar, P.; Arık, S. O.; Choi, D.; Bhattacharjee, B.; Lien, A.-T.; and Pfister, T. 2023. LoCoOp: Few-Shot Out-of-Distribution Detection via Prompt Learning. In *The Eleventh International Conference on Learning Representations (ICLR)*.
- Kirkpatrick, J.; Pascanu, R.; Rabinowitz, N.; Veness, J.; Desjardins, G.; Rusu, A. A.; Milan, K.; Quan, J.; Ramalho, T.; Grabska-Barwinska, A.; et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13): 3521–3526.
- Krizhevsky, A.; and Hinton, G. 2009. Learning multiple layers of features from tiny images. Technical report, University of Toronto.
- Kull, M.; Perello-Nieto, M.; Käng, M.; Filho, T. M.; Song, H.; and Flach, P. 2019. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with Dirichlet calibration. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32.

- Lee, K.; Lee, K.; Lee, H.; and Shin, J. 2018. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Advances in neural information processing systems (NIPS)*, volume 31.
- Li, Y.; Xu, S.; Zhang, B.; Cao, X.; Gao, P.; and Guo, G. 2022a. Q-ViT: Accurate and Fully Quantized Low-bit Vision Transformer. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Li, Z.; Cui, C.; Liu, X.; Zhang, Y.; Chang, S.; Cheng, H.; Cheng, Y.; and Chen, J. 2022b. CLIP-Q: Turning full-precision CLIP into a 4-bit model. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, 24031–24043.
- Li, Z.; Wang, F.; Zhang, Z.; and Li, F. 2023. Neg-CLIP: A Negative-Prompt-based Method for OOD Detection in Vision-Language Models. *arXiv preprint arXiv:2310.03114*.
- Liu, W.; Wang, X.; Owens, J.; and Li, Y. 2020. Energy-based out-of-distribution detection. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 21464–21475.
- Mayilvahanan, P.; Wiedemer, T.; Rusak, E.; Bethge, M.; and Brendel, W. 2023. Does CLIP’s Generalization Performance Mainly Stem from High Train-Test Similarity? *arXiv preprint arXiv:2310.09562*.
- Ming, Y.; and Li, Y. 2022. Delving into the Open-Set World: A Framework for Unsupervised Out-of-Distribution Detection. In *European Conference on Computer Vision (ECCV)*.
- Miyai, A.; Yang, J.; Zhang, J.; Ming, Y.; Lin, Y.; Yu, Q.; Irie, G.; Joty, S.; Li, Y.; Li, H. H.; Liu, Z.; Yamasaki, T.; and Aizawa, K. 2025. Generalized Out-of-Distribution Detection and Beyond in Vision Language Model Era: A Survey. In *Transactions on Machine Learning Research (TMLR)*.
- Nakkiran, P.; Kaplun, G.; Bansal, Y.; Yang, T.; Barak, B.; and Sutskever, I. 2021. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12): 124003.
- Noda, S.; Miyai, A.; Yu, Q.; Irie, G.; and Aizawa, K. 2025. A Benchmark and Evaluation for Real-World Out-of-Distribution Detection using Vision-Language Models. *arXiv preprint arXiv:2501.18463v1*.
- Polino, A.; Pascanu, R.; and Alistarh, D. 2018. Quantization-aware knowledge distillation. In *International Conference on Learning Representations (ICLR) Workshop*.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning (ICML)*, 8748–8763. PMLR.
- Recht, B.; Roelofs, R.; Schmidt, L.; and Shankar, V. 2019. Do ImageNet Classifiers Generalize to ImageNet? In *International Conference on Machine Learning (ICML)*, 5389–5400. PMLR.
- Saqib, J.; Hieu, L.; and Mathieu, S. 2025. QT-DoG: Quantization-aware Training for Domain Generalization. In *International Conference on Learning Representations (ICLR)*.
- Schuhmann, C.; Beaumont, R.; Vencu, R.; Gordon, C.; Wightman, R.; Cherti, M.; Coombes, T.; Katta, A.; Mullis, C.; Wortsman, M.; et al. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Shao, W.; Zhao, L.; He, Z.; Jiao, Z.; Chen, P.; and Ng, K.-T. 2023. OmniQuant: Omnidirectionally Calibrated Quantization for Large Language Models. In *The Eleventh International Conference on Learning Representations (ICLR)*.
- Sharma, P.; Ding, N.; Goodman, S.; and Soricut, R. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2556–2565.
- Tallic, C.; Blier, L.; and Ollivier, Y. 2023. Revisiting the Regularization Effect of Quantization. *arXiv preprint arXiv:2310.03113*.
- Teney, D.; Abbasi, E.; and van den Hengel, A. 2022. On the Pitfalls of Spurious Correlations for OOD Generalization. In *International Conference on Learning Representations (ICLR)*.
- Tishby, N.; Pereira, F. C.; and Bialek, W. 2000. The information bottleneck method. In *Proceedings of the 37th Annual Allerton Conference on Communication, Control, and Computing*, 368–377.
- Tu, W.; Deng, W.; and Gedeon, T. 2023. A closer look at the robustness of contrastive language-image pre-training (clip). *Advances in Neural Information Processing Systems*, 36: 13678–13691.
- Van Horn, G.; Mac Aodha, O.; Song, Y.; Cui, Y.; Sun, C.; Shepard, A.; Adam, H.; Perona, P.; and Belongie, S. 2018. The iNaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 8769–8778.
- Wang, H.; Ge, S.; Lipton, Z.; and Xing, E. P. 2019. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32.
- Wang, Q.; Lin, Y.; Chen, Y.; Schmidt, L.; Han, B.; and Zhang, T. 2024. A Sober Look at the Robustness of CLIPs to Spurious Features. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Xiao, G.; Lin, J.; Seznec, M.; Wu, H.; Demouth, J.; and Han, S. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning (ICML)*, 38087–38101. PMLR.
- Xiao, J.; Hays, J.; Ehinger, K. A.; Oliva, A.; and Torralba, A. 2010. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, 3485–3492. IEEE.
- Yang, J.; Wang, P.; Zou, D.; Zhou, Z.; Ding, K.; Peng, W.; Wang, H.; Chen, G.; Li, B.; Sun, Y.; Du, X.; Zhou,

- K.; Zhang, W.; Hendrycks, D.; Li, Y.; and Liu, Z. 2022. OpenOOD: Benchmarking Generalized Out-of-Distribution Detection. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, 30150–30164.
- Yang, J.; Zhou, K.; and Liu, Z. 2022. Full-Spectrum Out-of-Distribution Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16293–16302.
- Yin, D.; Gontareva, A.; Gontarev, I.; Kornblith, S.; Gu, S.; and Le, Q. V. 2019. A Fourier perspective on the generalization of deep neural networks. In *International Conference on Machine Learning (ICML)*, 7133–7142. PMLR.
- Zhang, J.; Yang, J.; Wang, P.; Wang, H.; Lin, Y.; Zhang, H.; Sun, Y.; Du, X.; Li, Y.; Liu, Z.; Chen, Y.; and Li, H. 2024. OpenOOD v1.5: Enhanced Benchmark for Out-of-Distribution Detection. *Journal of Data-centric Machine Learning Research*.
- Zhou, B.; Lapedriza, A.; Khosla, A.; Oliva, A.; and Torralba, A. 2017. Places: A 10 million image database for scene recognition. In *IEEE transactions on pattern analysis and machine intelligence*, volume 40, 1452–1464. IEEE.

Appendix

Abstract

This appendix provides a comprehensive extension to our main paper, focusing on explanations, reproducibility, and exhaustive results. We begin with a transparent discussion of our study’s limitations. The centerpiece is a novel frequency-domain analysis using the 2D Fourier spectrum, which reveals how Quantization-Aware Training (QAT) reshapes a model’s feature representation toward more robust, low-frequency concepts. To ensure full reproducibility, we meticulously document our experimental methodology, including a detailed justification for using simulated quantization and a full breakdown of model architectures, datasets, and hyperparameters. We also add the plots that did not fit in our main paper and that are very relevant. We then present additional analyses and visualizations, including a conceptual illustration of QAT’s preference for flat minima and detailed reliability diagrams. Finally, we provide comprehensive tables containing the complete numerical results for all experiments, covering in-distribution accuracy, calibration, out-of-distribution detection, and robustness for every model and method evaluated.

4.1 Limitations

Our study, while comprehensive, has certain limitations.

- **Visual Encoder Only:** Our primary limitation is that we only quantize the visual encoder. Quantizing the text encoder, or both encoders simultaneously, introduces additional challenges related to maintaining the alignment of the joint embedding space and represents an important avenue for future work.
- **Proxy Dataset:** All QAT methods were fine-tuned on CC3M as a proxy dataset. The choice of proxy data could influence the final performance, especially regarding catastrophic forgetting.
- **Dynamic Quantization:** In our work, we quantize all linear layers without prior analysis. Research has shown that for LLM, selective layers quantization is essential for quantization as some layers gain massive rounding errors when quantized (Li et al. 2022a). We did not include results using more specialized and complex quantization methods to avoid adding additional noise to our study. Including such complex methods would have made it difficult to separate the true impact of quantization from whatever additional methods were used.
- **Pure Zero-shot:** Some important baselines are missing from this work, some due to the complexity of implementation, and others due to relying on sample data e.g. Mahalanobis (Lee et al. 2018), LoCoOp (Kar et al. 2023), OmniQuant (Shao et al. 2023).

5 Frequency-Domain Analysis of Quantization Effects

To investigate the mechanistic underpinnings of how quantization affects feature representations, we analyze the 2D

Fourier spectrum of the Vision Transformer’s internal feature maps (Yin et al. 2019). Our analysis is grounded in the foundational principles of Fourier theory, which establishes a direct correspondence between a signal’s frequency content and its spatial rate of change. We extend this principle to the feature maps of a ViT, where low-frequency components at the spectrum’s origin represent coarse-grained visual concepts (e.g., broad shapes), while high-frequency components in the outer bands encode fine-grained details (e.g., sharp edges, textures). This spectral analysis thus allows for a quantitative measurement of the model’s capacity to represent features of varying granularity. The feature representations for our analysis were extracted from the output of the final residual block within the Vision Transformer encoder. Specifically, we capture the sequence of output tokens from the `vision_model.transformer.resblocks` module, immediately prior to the application of the terminal layer normalization (`ln_post`) and any subsequent projection layers. This layer provides the richest per-patch semantic representation while preserving the spatial topology necessary for our 2D spectral analysis, and we also argue the significance of our selection by only including the attention section of the ViT (which is known to be more sensitive to quantization) rather than the MLP section of the model (which is more quantize-friendly (Li et al. 2022a)).

Relative Spectral Error (RSE) To precisely quantify the deviation of quantized models from the full-precision baseline, we define the Relative Spectral Error (RSE). The RSE at each 2D frequency coordinate (u, v) is calculated as the normalized absolute difference between the quantized spectrum, S_{quant} , and the baseline spectrum, S_{FP32} :

$$\text{RSE}(u, v) = \frac{|S_{\text{quant}}(u, v) - S_{\text{FP32}}(u, v)|}{S_{\text{FP32}}(u, v) + \epsilon} \quad (1)$$

Here, S_{quant} represents either S_{PTQ} or S_{QAT} , and ϵ is a small constant (e.g., 10^{-9}) to ensure numerical stability where $S_{\text{FP32}}(u, v)$ is close to zero. The resulting 2D RSE map provides a spatial visualization of the error distribution across the frequency domain.

5.1 Initial Findings with Aggressive Quantization

We first investigate the effects of aggressive 6-bit quantization on a ViT-B/32 model. Figure 11 shows that while the full-precision (FP32) spectrum exhibits significant energy across all bands, Post-Training Quantization (PTQ) severely attenuates the high-frequency components, indicating a catastrophic loss of fine-grained representational capacity. Subsequent Quantization-Aware Training (QAT) only partially restores these components. The Relative Spectral Error (RSE) maps confirm that quantization error is maximized in these high-frequency bands. These findings suggest that aggressive quantization acts as a potent low-pass filter, forcing the model to rely on robust, low-frequency features, which may explain observed improvements in robustness to certain data shifts.

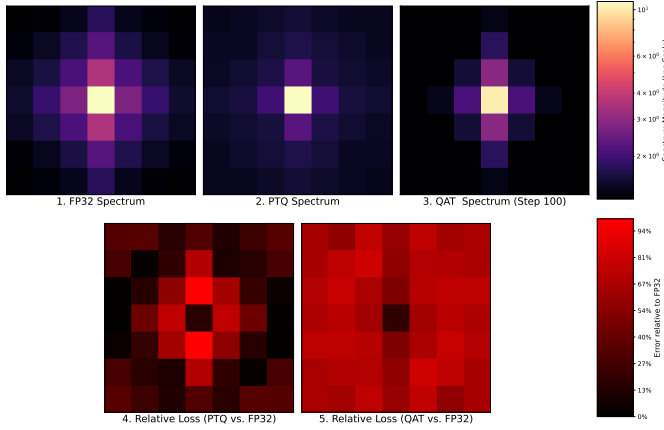


Figure 11: Spectral analysis of a ViT-B/32 model with 6-bit quantization. From left to right: the FP32 baseline spectrum, the PTQ spectrum showing severe high-frequency attenuation, and the partially restored QAT spectrum. Bottom row shows corresponding RSE maps.

5.2 Refined Analysis with Standard Quantization

While informative, these initial results arise from conditions—aggressive 6-bit precision and a low-resolution 7x7 feature grid—that may not reflect standard use cases. To disentangle these factors and study the adaptive nature of QAT with greater fidelity, we designed a refined experiment with standard w8a8 precision and a higher-resolution ViT-L/14 model, providing a 16x16 feature grid.

Our refined analysis, presented in Figure 12, reveals a more nuanced mechanism. The 8-bit PTQ induces a targeted, rather than catastrophic, degradation, manifesting as a distinct “error halo” in the mid-frequency bands of the RSE map. This confirms that even standard quantization preferentially degrades fine-grained features. The most striking result emerges from the QAT model. While it successfully eliminates the localized error halo, it does not converge to the original FP32 state. Instead, it finds a new equilibrium characterized by a low-magnitude, but globally distributed, error floor across the entire spectrum. This demonstrates that QAT is a sophisticated re-optimization process, as it reduces the coarse-grained features errors, compromising the fine-grained ones as an effective strategy to reduce global loss; the model finds a novel, stable representation within the 8-bit constraint, trading the targeted error for a minimal, systemic deviation. This suggests the robustness of quantized models arises not from a simple forgetting of fine-grained features, but from an active convergence to a new, more generalized internal representation that leverages coarse-grained, general features. Code for reproducibility of both plots is available.

5.3 Extended Experimental Setup

This section provides the full details of our experimental methodology to ensure clarity and reproducibility. Code will be provided for reproducibility, including all results (in CSV format), plotting notebooks, Quantization simulation, quan-

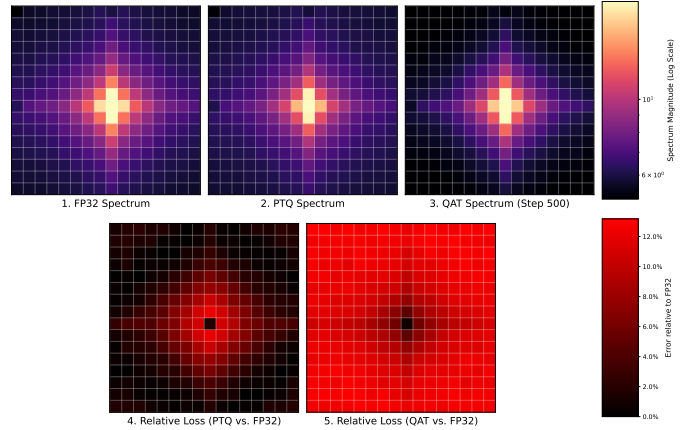


Figure 12: Spectral analysis of a ViT-L/14 model with 8-bit quantization. Top row: FP32, PTQ, and QAT spectra. Bottom row: RSE maps, where PTQ induces a mid-frequency “error halo” (left), and QAT eliminates this halo in favor of a low-magnitude, global error floor (right), indicating a complex adaptation.

tization simulation verification, and Fourier. We used an H100 cluster for the evaluations.

Model Architectures We used publicly available pre-trained OpenCLIP models. WIT pretrained denotes OpenAI weights. For all experiments, quantization and fine-tuning were applied **only to the visual encoder**, while the text encoder was kept frozen at FP32 precision. The primary models evaluated are detailed in Table 2.

Architecture	Pre-training Source	OpenCLIP Model Tag
ViT-B/32	WIT (OpenAI)	openai
ViT-B/32	LAION	laion400m_e32
ViT-L/14	WIT (OpenAI)	openai
ViT-L/14	LAION	laion400m_e32
ViT-B/16	WIT (OpenAI)	openai
ViT-B/16	LAION	laion400m_e32

Table 2: Details of the primary CLIP models used in our evaluation.

Benchmark Dataset Taxonomy Our benchmark suite systematically covers different forms of distribution shift, and we use the same datasets that are recognised and used by the community. For more details on dataset preparation for OOD detection, please refer to citep . They are categorized along two dimensions: *Covariate Shift* and *Semantic Shift*, as shown in Figure 13. We also use a specialized dataset for spurious correlations (Hochlehnert et al. 2025).

Quantization and Training Configurations Our experiments involved two main families of QAT, plus PTQ methods referenced from the literature. All QAT was performed on a subset of the CC3M(Sharma et al. 2018) dataset, as its higher quality and better curation proved slightly

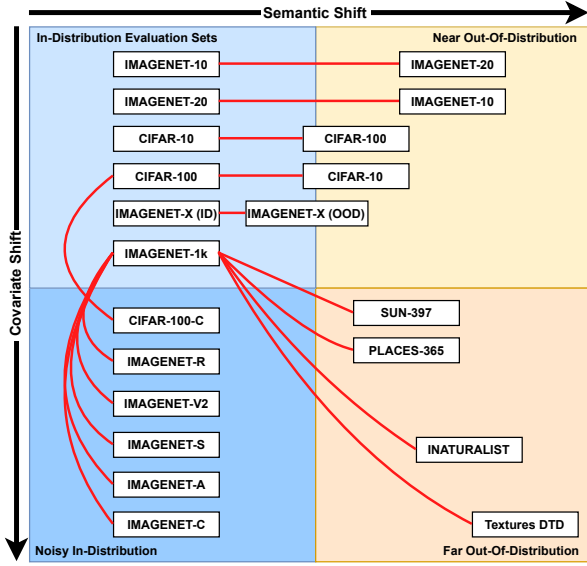


Figure 13: Systematic benchmark suite categorized by Co-variate and Semantic Shift. Arrows indicate the ID reference dataset for OOD detection.

more effective for short adaptation schedules compared to LAION400m (Schuhmann et al. 2022).

A crucial aspect of our methodology is the use of **simulated or “fake” quantization**. This approach was necessary because, at the time of our experiments, standard deep learning frameworks (e.g., PyTorch) and GPU hardware (e.g., NVIDIA H100 & RTX 3090) lack native support for performing computations with INT8, INT6 or INT4 *activations* or for models like Vision Transformers. While libraries like ‘bitsandbytes’ offer efficient kernels for weight-only quantization (e.g., for QLoRA), a full pipeline with quantized activations is not readily available. Therefore, simulation is the standard and necessary method in research to evaluate the potential performance of such models before dedicated hardware or software support exists.

The Role of Fake Quantization Fake quantization is a simulation technique that models the error introduced by quantization and de-quantization within a standard full-precision (FP32/FP16) training and inference loop. The process for a tensor x is as follows:

1. **Quantize:** The full-precision input tensor x is scaled, shifted, and rounded to the nearest integer value within the target bit-width’s range (e.g., $[-128, 127]$ for INT8).

$$x_{quant} = \text{round} \left(\frac{x}{\text{scale}} + \text{zero_point} \right)$$

2. **Clamp:** The integer values are clamped to the representable range of the target bit-width.
3. **De-quantize:** The clamped integer tensor is immediately converted back to a full-precision floating-point tensor.

$$x_{dequant} = (x_{quant} - \text{zero_point}) \times \text{scale}$$

The resulting tensor, $x_{dequant}$, has the same data type as the input but contains the precision loss that *would have occurred* in a true low-bit system. For the backward pass, a **Straight-Through Estimator (STE)** bypasses the non-differentiable ‘round’ function, enabling the model to learn weights that are robust to the simulated quantization noise.

Verification of Simulation Correctness. To ensure our simulation was valid, we implemented a verification utility. For any quantized model, this utility counts the number of unique values in both the weight and activation tensors (via forward hooks). For a successful INT8 simulation, this count must be $\leq 2^8 = 256$, which we confirmed for our implementations.

Motivation: Why Not Use Existing Kernels like QLoRA for Inference? While methods like QLoRA enable the *training* of models with 4-bit weights, they are primarily designed to reduce memory usage during the training phase, not to accelerate inference. As shown in our direct performance benchmark in Table 3, using these kernels for inference can lead to a significant *slowdown* compared to the FP32 baseline. The overhead of de-quantizing the weights on-the-fly for each computation outweighs the benefits of reduced memory bandwidth. True inference acceleration requires dedicated hardware and software support for low-precision matrix multiplications for multi-head attention, which is what our simulation-based study aims to evaluate in terms of potential accuracy and reliability before such support becomes widespread.

Post-Training Quantization (PTQ) PTQ methods (‘PTQ’, ‘SmoothQuant’, ‘Raana’) use a small calibration set to determine fixed quantization parameters. Our evaluation of these methods relies on the fake quantization simulation for inference to measure their impact on accuracy and reliability.

Full Quantization-Aware Training (QAT) This approach (‘QAT-LORA’ variants) involves fine-tuning using the fake quantization simulation. Although our QAT-LORA implementation is original, we suggest practitioners refer to more recent methods like LR-QAT (Bondarenko, Chiaro, and Nagel 2024), which are designed with inference efficiency in mind by re-fusing the LoRA matrices post-training.

Quantization Primitive: We implemented a custom Learnable Scale Quantization (LSQ) function with a gradient scaling factor of $\frac{1}{\sqrt{N \cdot Q_{max}}}$ to stabilize training.

Model Architecture: We apply LoRA adapters to the vision encoder’s linear layers. Our custom ‘EnhancedFakeQuantizer’ modules are then injected such that the base model’s path is fake-quantized while the parallel LoRA path operates on the full-precision input.

Training Hyperparameters: These models are trained with the parameters in Table 4. An optional knowledge distillation loss (MSE on normalized features) is used for ‘Distill’ variants.

Model	Quantization	Memory (MB)	Time (ms)	Speedup (vs. FP32)	Memory Gain (vs. FP32)
B/32	32-bit (FP32)	358.27	19.61	1.00x	0.0%
	16-bit (FP16)	297.79	17.67	1.11x	16.9%
	8-bit (INT8)	179.82	81.12	0.24x	49.8%
	4-bit (NF4)	123.80	47.70	0.41x	65.4%
B/16	32-bit (FP32)	579.51	12.33	1.00x	0.0%
	16-bit (FP16)	294.63	14.16	0.87x	49.2%
	8-bit (INT8)	177.66	90.90	0.14x	69.3%
	4-bit (NF4)	119.65	49.83	0.25x	79.3%
L/14	32-bit (FP32)	1312.47	27.28	1.00x	0.0%
	16-bit (FP16)	—	—	—	—
	8-bit (INT8)	—	—	—	—
	4-bit (NF4)	278.31	70.20	0.39x	78.8%

Table 3: Performance and memory benchmark of OpenAI CLIP models using non-simulated, kernel-based quantization. Inference time is the average latency per batch on an RTX 3090 GPU. Note the significant slowdown for 8-bit and 4-bit inference, which motivates our use of simulation to evaluate accuracy/reliability trade-offs, as current kernels are not optimized for inference speed. Failures for the L/14 model are attributed to a known bug in the tested library version.

Hyperparameter	Value
LoRA Rank (r)	8
LoRA Alpha (α)	16
Base Learning Rate	5×10^{-5}
LSQ Scale Learning Rate	1×10^{-7}
Optimizer	AdamW
Training Steps	100-8100
Batch size	100
Precision	W8A8 to W4A6
Distillation α	0.5 (if applicable)

Table 4: Hyperparameters for Full QAT methods.

Light Quantization-Aware Training (Light QAT) This approach (‘Basic-QAT’, ‘QAT-LSQ’) employs a unified ‘FakeQuantize’ module and a short training schedule (100 batches) to quickly adapt the model to quantization noise. ‘Distill’ variants use a KL-Divergence loss. Light Qat uses a smaller learning rate and a limited amount of adaptation samples, as QAT is an adaptation procedure and not a training one, it doesn’t follow the same general tools as for training, we wish to avoid catastrophic forgetting by reducing the amount of new semantic information introduced during the adaptation, we loop back the data and limit the number of unique samples, the experiment below shows that a single sample is enough to adapt the model to quantization, and increasing the number of unique samples doesn’t show a clear trend of improvement in this use-case.

5.4 Additional Analyses and Results

On QAT’s Preference for Flat Minima A significant body of work establishes that neural networks converging to “flat” minima in the loss landscape generalize better (Hochreiter and Schmidhuber 1997). As illustrated in Figure 15, QAT inherently penalizes sharp regions because the

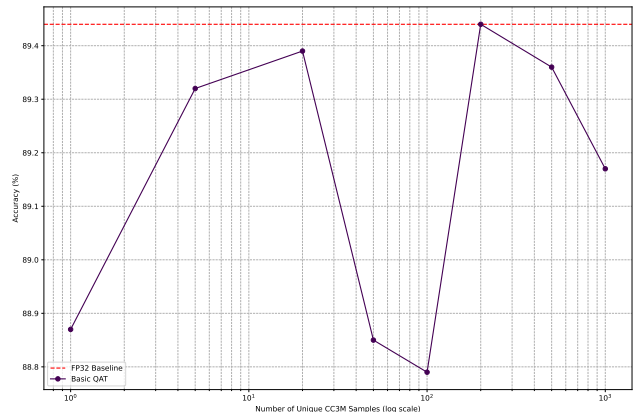


Figure 14: Evolution of ViT/B-32 Quantized model accuracy on Cifar10 dataset, relative to the amount of unique CC3M samples used during QAT adaptation.

Hyperparameter	Value
Learning Rate	1×10^{-6}
Optimizer	AdamW
Training Steps	100
Batch size	100
Unique samples used	100
Distillation α	0.5 (if applicable)

Table 5: Hyperparameters for Light QAT methods.

weight perturbation from quantization ($Q(w) - w$) causes a large increase in loss. This forces the optimizer to seek out flatter, more robust regions, which can sometimes lead to better generalization and accuracy than the original FP32 model.

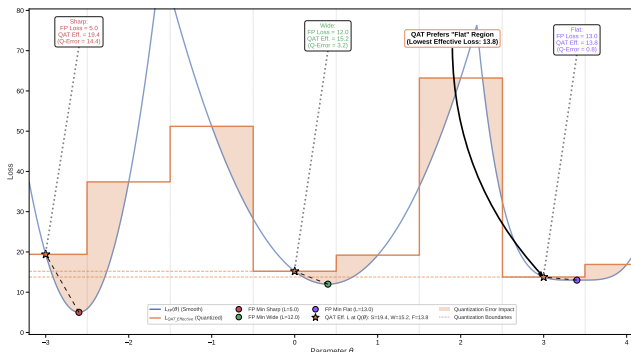


Figure 15: A conceptual illustration of how QAT forces the optimizer to abandon a sharp minimum in favor of a flatter, more robust solution.

Detailed Calibration Analysis Here we provide the detailed visualizations of the calibration process, fulfilling the promises from the main paper. Figure 16 shows the direct impact of QAT (Phase 1), illustrating how quantization alone squeezes confidences, often pulling them away from the ideal diagonal. Figure 17 shows the final state after our full two-phase re-calibration process (QAT followed by logit scale tuning). The points move much closer to the diagonal, demonstrating the necessity of the second tuning phase to correct for the distortions introduced by QAT and achieve a well-calibrated model.

Accuracy vs. NLL Trade-off The main paper’s teaser plot shows the trade-off between accuracy and ECE. Figure 18 provides another view, plotting accuracy change against Negative Log-Likelihood (NLL) change. NLL is a stricter metric that penalizes overconfident errors more heavily than ECE. The ideal quadrant (green) represents a simultaneous improvement in accuracy and a reduction in NLL. This plot further highlights the complex, multi-objective nature of deploying quantized models, where gains in one metric can often come at the cost of another.

OOD Detection Performance This section provides more detailed results on OOD detection to complement the main paper. Figure 21 shows the False Positive Rate at 95% True Positive Rate (FPR@95), a key metric for evaluating performance at a practical operating point. Lower is better. This plot reinforces the findings from the AUROC analysis, showing that VLM-specific methods are generally more robust to quantization and that the impact of quantization is highly dependent on the pre-training source.

Robustness to Covariate Shift Here we expand on the covariate shift results from the main paper. Figure 22 shows how Light QAT methods affect the model’s ability to *detect* corrupted data as OOD, measured by the change in AUROC. For low-level corruptions like noise and blur, QAT significantly improves detection ability (a large positive change). However, for more complex, semantic-style shifts like ImageNet-R and -Sketch, the effect is neutral or slightly negative. This reveals a nuanced trade-off: QAT makes the

model robust to simple noise at the cost of sensitivity to fine-grained stylistic changes.

Covariate Shift Data Generation Pipeline For experiments on datasets like ImageNet-C, it is crucial to apply corruptions in a way that mimics real-world sensor noise. Our ‘ResizeThenCorruptDataset’ class ensures this by applying corruptions *after* the standard resizing and cropping but *before* the final ‘ToTensor’ and normalization steps. This ‘‘Resize-then-Corrupt’’ order prevents artifacts from the corruption process being altered by subsequent resizing operations.

5.5 Comprehensive Tables

Below we present the detailed numerical results for some experiments (PTQ ViT/B32,B16,L14 and Light QAT ViT/B32). The data is grouped by experimental configuration, predictive quality, and OOD detection (Model, Pre-training, and ID/OOD datasets) to facilitate direct comparison between the different quantization methods within each scenario. Please note that we cannot add all tables due to time constraints, please refer to the CSVs in the code.

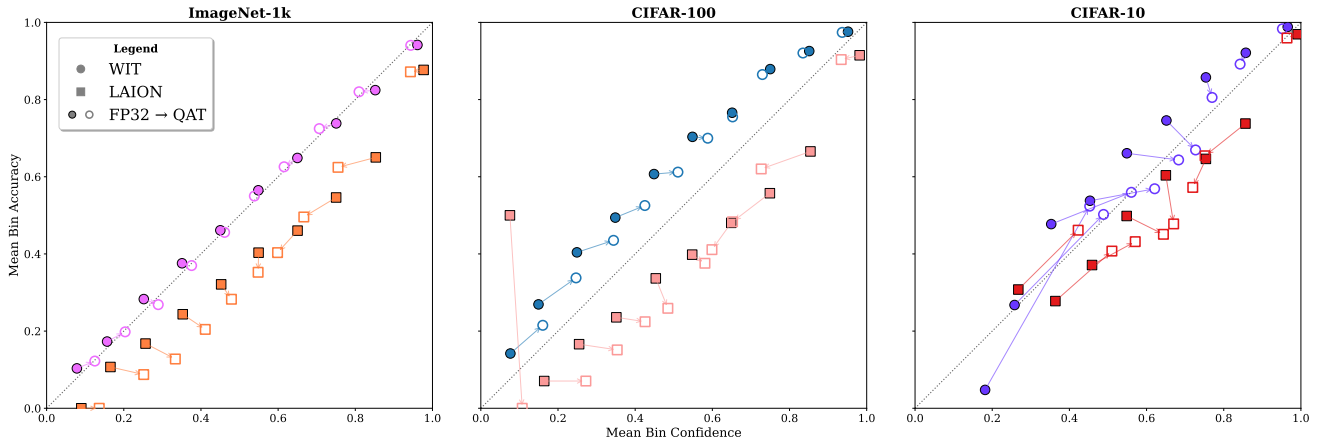


Figure 16: Dataset-specific reliability diagrams showing the direct impact of Quantization-Aware Training (Phase 1). This plot shows the calibration state after QAT has been applied but before logit scale tuning. Note how the confidences for both WIT (blue/purple) and LAION (red/orange) models are shifted away from their initial FP32 state, often towards the center, but are not yet aligned with the ideal diagonal.

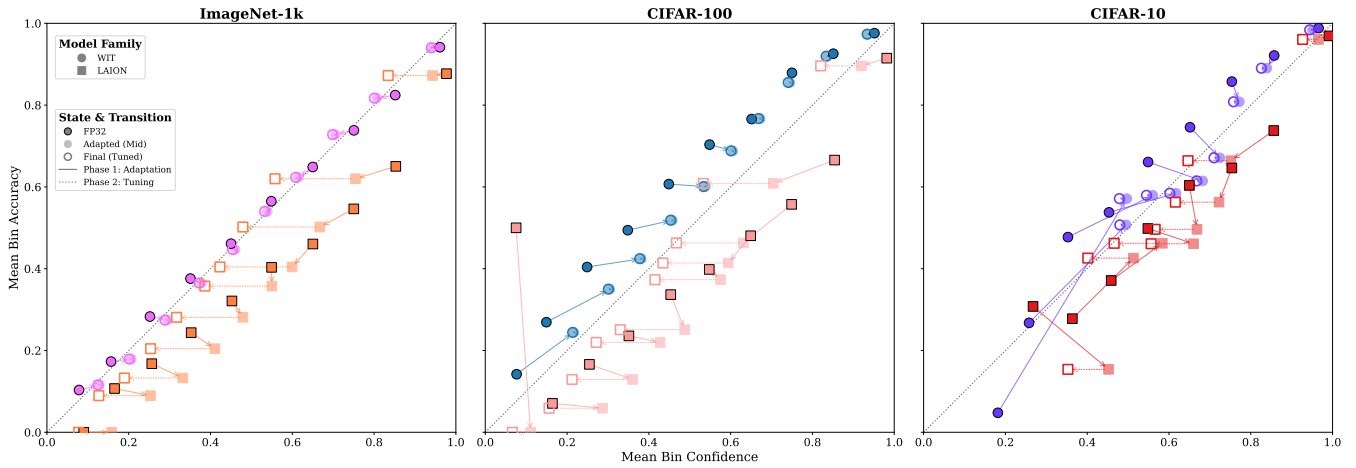


Figure 17: Dataset-specific reliability diagrams showing the final calibration state after the full two-phase process (QAT + Logit Scale Tuning). The "Final (Tuned)" points (hollow markers) show a significant improvement over the intermediate "Adapted (Mid)" state, moving much closer to the perfect calibration line (dashed diagonal). This demonstrates the effectiveness and necessity of the second logit-tuning phase.

Trade-offs in CLIP Quantization: Accuracy and Uncertainty

Impact of Quantization on Zero-Shot Accuracy vs. NLL

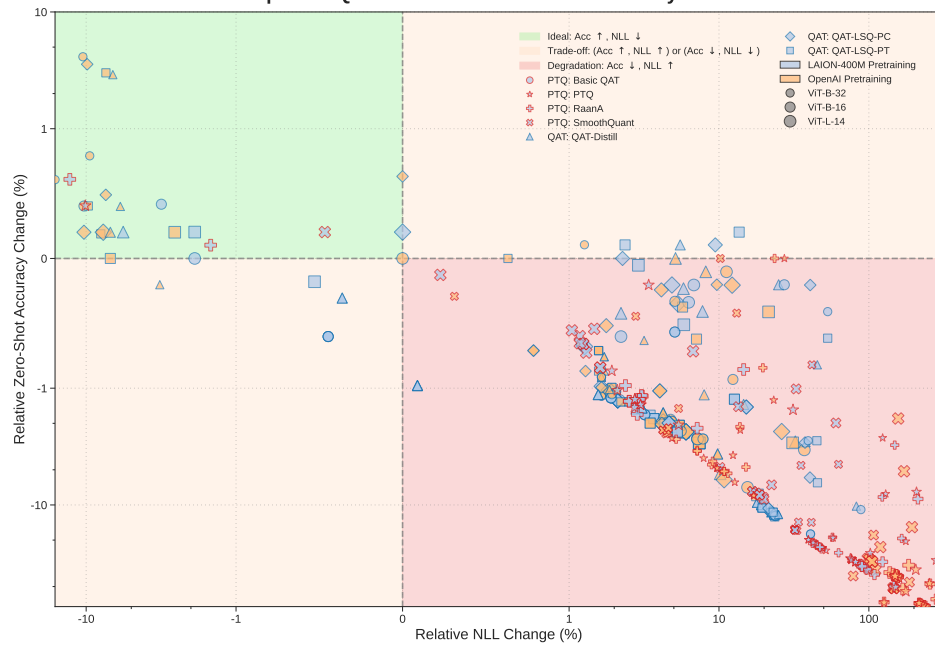


Figure 18: The trade-off between zero-shot accuracy and Negative Log-Likelihood (NLL). The ideal outcome (green, top-left) is a gain in accuracy and a reduction in NLL. Most methods fall into the trade-off or degradation quadrants, highlighting the challenge of improving both metrics simultaneously. Note the logarithmic scale for the NLL axis.

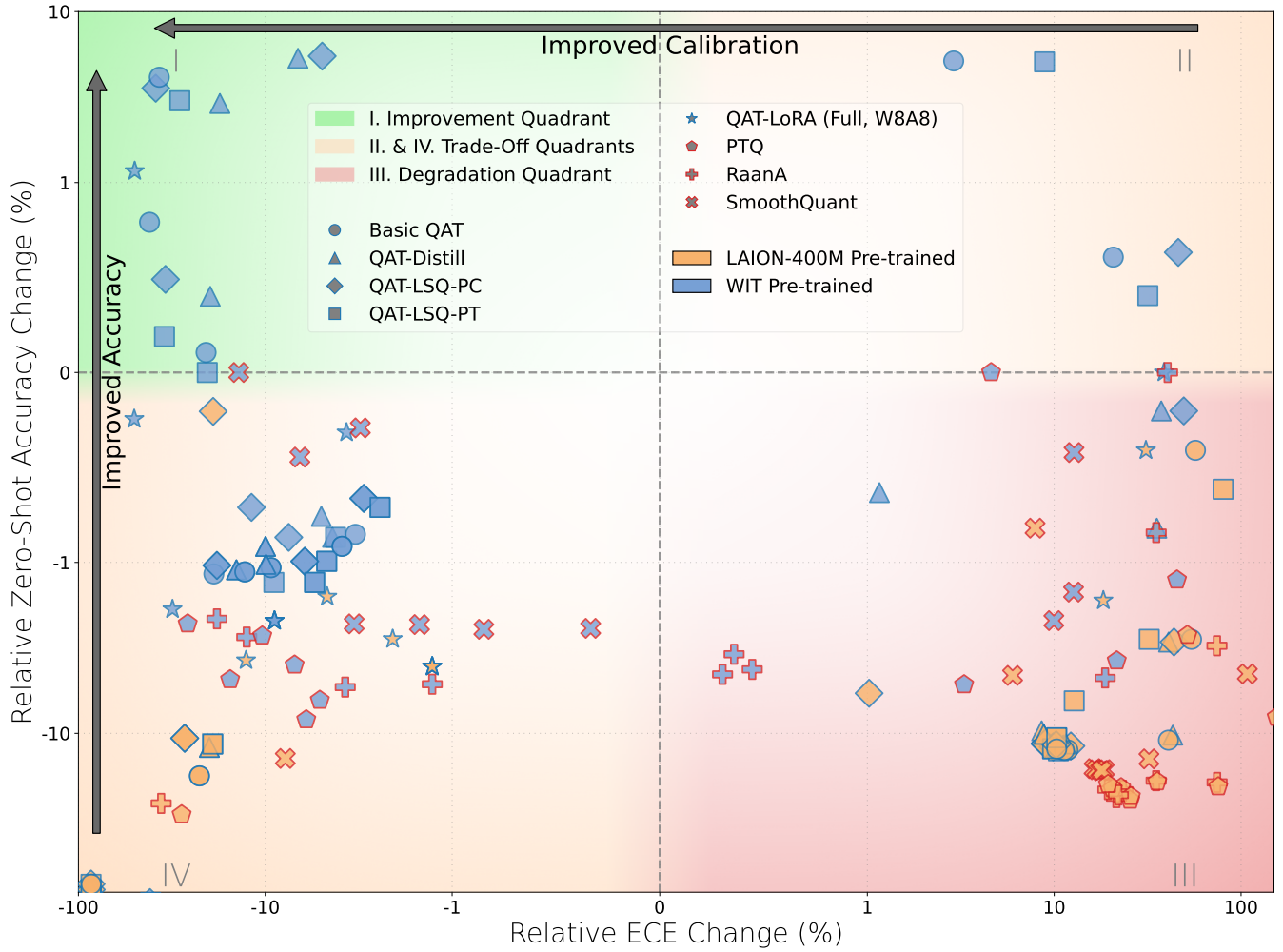


Figure 19: Full-size teaser figure: The dichotomous impact of quantization on zero-shot Performance. WIT models (blue) consistently improve in calibration (left), with several QAT methods achieving simultaneous accuracy gains. In contrast, LAION models (orange) show systematic degradation in calibration (right). The lack of pfalloints near the origin suggests that quantization is always impactful.

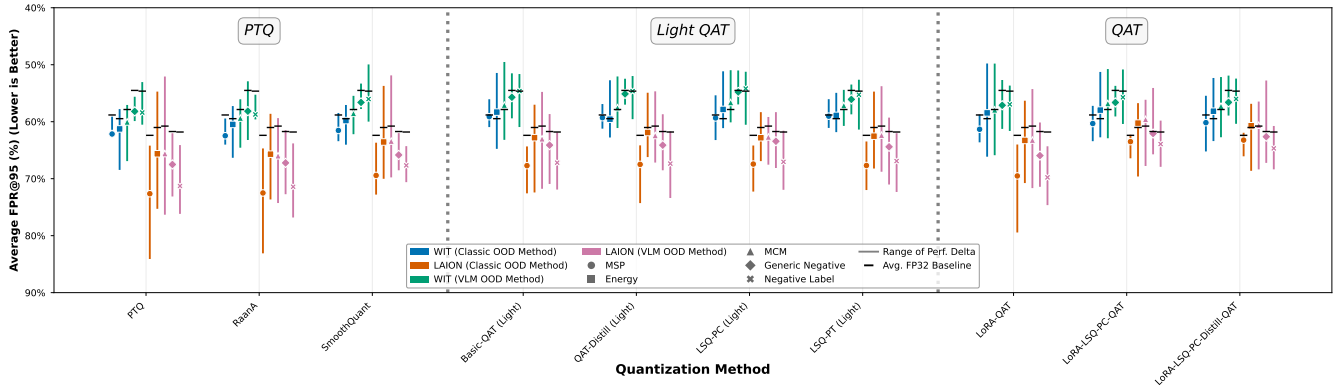


Figure 20: Average FPR@95 across all Far-OOD datasets. Lower values are better. This plot complements the AUROC results from the main paper, showing OOD detection performance at a practical operating point. The VLM-specific methods (teal and pink) consistently outperform classic methods (blue and orange), and their performance is more resilient to quantization.

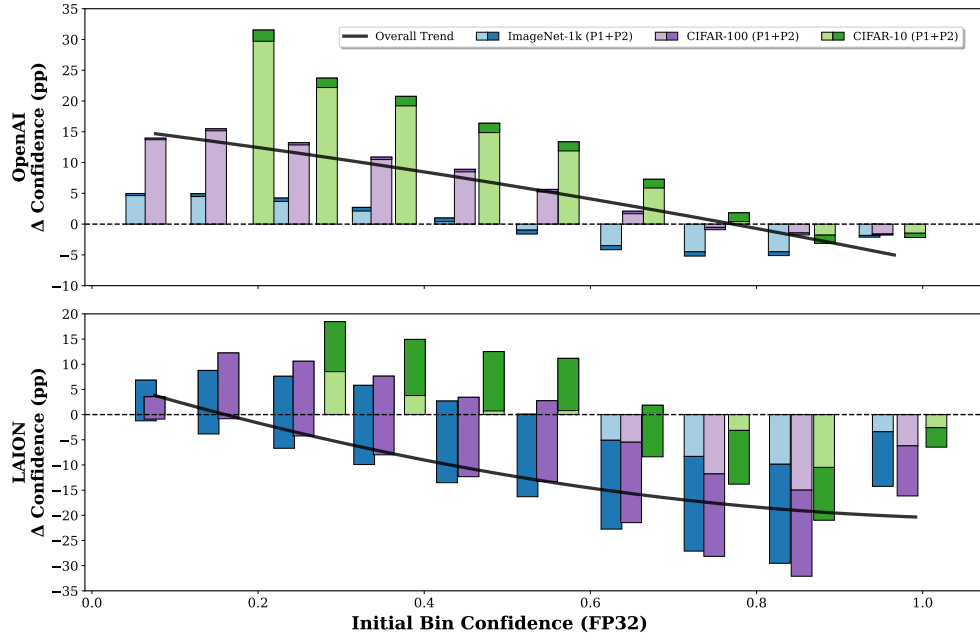


Figure 21: Confidence bins shift after QAT and Logit Adaptation, we can clearly see the trend where the overconfident LAION pretrained model gets its confidence corrected down, while the underconfident WIT pretrained gets its confidence increased, showing the uncertainty awareness introduced by the quantized-adapted model.

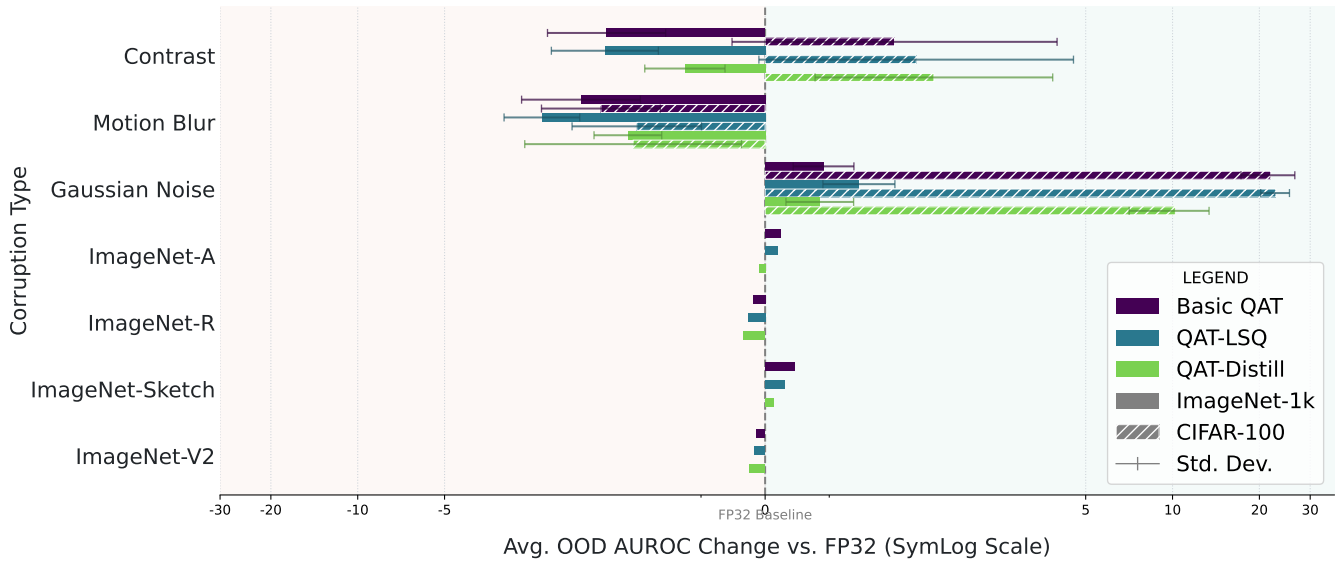


Figure 22: The change in OOD detection (AUROC) performance on various covariate shift datasets after applying Light QAT methods. Positive values indicate improved detection. QAT substantially improves the ability to detect low-level corruptions (top rows) but has a mixed-to-negative impact on detecting complex semantic shifts (bottom rows).

Scenario	Method	Out-of-Distribution Detection Performance									
		Traditional Methods					VLM-based Methods				
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	94.78	19.61	88.59	49.47	91.70	40.78	95.65	20.36	95.02	25.33
	PTQ	78.15	53.23	78.44	59.54	79.72	56.58	79.72	56.88	80.10	59.19
	SmoothQuant	84.66	41.02	81.06	54.51	83.67	48.36	86.93	41.79	86.65	44.49
	RaanA	79.12	51.58	80.08	57.33	81.18	54.48	80.83	54.45	81.21	57.18
ID: CIFAR100 OOD: CIFAR10	FP32	80.12	55.05	80.45	59.15	83.76	49.14	82.17	48.19	81.86	51.52
	PTQ	55.25	82.64	54.26	83.67	54.00	81.30	52.49	82.18	52.39	82.36
	SmoothQuant	60.41	77.82	60.62	79.97	61.42	74.85	58.26	76.53	59.50	77.24
	RaanA	55.29	82.57	54.83	83.16	54.49	81.12	52.49	81.97	52.59	82.60
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	79.00	63.18	74.20	72.46	76.15	70.30	77.35	70.65	78.44	68.74
	PTQ	61.87	82.56	59.56	86.53	60.95	84.90	60.62	84.92	60.51	84.15
	SmoothQuant	67.16	76.88	64.16	82.53	66.10	80.64	65.91	81.16	66.45	79.59
	RaanA	62.91	81.54	60.45	85.87	61.87	84.34	61.66	84.17	61.79	82.94
ID: ImageNet10 OOD: ImageNet20	FP32	96.95	10.60	96.09	26.40	96.18	25.80	95.36	28.80	96.60	21.40
	PTQ	86.97	44.40	87.55	52.60	87.97	52.80	85.42	57.60	87.31	55.60
	SmoothQuant	92.45	31.20	92.49	42.40	92.77	42.40	90.32	48.40	91.88	46.20
	RaanA	88.53	43.40	88.86	52.40	89.27	52.40	86.53	56.80	88.21	55.20
ID: ImageNet20 OOD: ImageNet10	FP32	95.23	23.80	92.97	53.00	93.32	52.80	92.02	64.30	94.43	42.80
	PTQ	87.27	38.40	83.58	69.70	84.85	64.90	83.48	63.70	85.18	54.70
	SmoothQuant	90.79	29.80	88.96	56.00	89.97	49.90	88.57	54.60	90.72	47.60
	RaanA	87.67	40.50	84.15	66.40	85.49	62.20	84.06	63.40	85.86	53.90
ID: IN1k OOD: DTD	FP32	80.69	77.70	64.74	79.50	70.93	77.44	85.18	54.30	84.80	63.40
	PTQ	73.94	78.27	63.10	76.21	69.55	73.52	75.52	63.93	77.36	64.39
	SmoothQuant	75.62	78.69	61.83	81.66	69.11	80.48	78.98	62.96	80.57	64.78
	RaanA	74.31	78.53	61.89	78.72	68.75	75.11	76.73	62.82	79.08	62.45
ID: IN1k OOD: iNaturalist	FP32	77.32	83.34	65.52	84.37	69.39	84.27	85.60	58.50	72.61	84.27
	PTQ	68.91	77.57	60.91	88.69	65.09	86.82	78.34	69.80	71.10	79.30
	SmoothQuant	72.66	76.23	61.38	90.12	66.43	87.71	81.51	66.43	73.69	81.04
	RaanA	69.05	77.18	61.03	88.30	65.22	86.49	78.45	69.24	71.63	79.14
ID: IN1k OOD: Places365	FP32	79.71	70.10	78.85	63.32	81.38	62.17	88.43	47.62	85.92	54.65
	PTQ	69.63	81.07	78.21	63.49	78.21	65.60	80.11	59.84	75.16	65.79
	SmoothQuant	71.63	79.89	76.26	68.22	77.96	67.32	84.11	54.38	76.88	65.89
	RaanA	68.17	82.54	76.17	66.03	75.97	67.23	81.11	58.45	72.83	68.11
ID: IN1k OOD: SUN397	FP32	80.34	70.28	82.60	54.64	84.54	55.22	86.49	53.34	87.34	51.63
	PTQ	67.94	81.72	78.58	60.31	77.50	63.07	77.28	62.57	73.76	66.50
	SmoothQuant	69.83	80.99	78.81	60.77	79.18	62.66	79.44	60.93	76.87	64.87
	RaanA	69.01	82.17	79.56	58.77	78.72	62.46	78.18	61.71	75.07	65.44

Table 6: PTQ Results: Out-of-Distribution (OOD) Detection for OpenAI pre-trained ViT-L-14. A/F denote AUROC/FPR95.

Scenario	Method	Out-of-Distribution Detection Performance									
		Traditional Methods					VLM-based Methods				
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	92.39	24.30	91.07	47.89	91.31	46.36	93.64	36.85	93.33	37.29
	PTQ	91.17	26.52	89.22	53.65	89.53	52.32	92.59	40.48	92.58	40.93
	SmoothQuant	92.28	24.82	90.73	48.24	91.00	47.42	93.51	38.48	93.32	37.45
	RaanA	91.78	25.17	89.99	51.21	90.28	50.13	93.05	39.07	92.97	39.08
ID: CIFAR100 OOD: CIFAR10	FP32	79.97	52.59	84.91	51.28	85.39	50.05	87.91	45.35	87.36	44.04
	PTQ	78.98	53.37	84.43	52.26	84.95	51.14	87.15	47.66	85.66	48.12
	SmoothQuant	79.69	53.03	84.95	51.28	85.47	50.06	87.18	47.00	86.61	45.82
	RaanA	78.67	55.66	83.95	52.43	84.48	51.09	86.84	49.10	85.26	49.40
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	75.49	65.81	75.12	73.71	75.59	72.86	76.68	71.67	78.52	69.09
	PTQ	74.95	66.59	74.72	73.38	75.21	72.69	76.20	71.48	78.09	68.92
	SmoothQuant	75.21	66.24	74.69	73.64	75.19	72.85	76.16	71.53	77.86	70.06
	RaanA	75.02	66.68	74.72	73.94	75.20	73.57	76.15	71.82	77.95	68.85
ID: ImageNet10 OOD: ImageNet20	FP32	96.30	17.40	96.49	19.80	96.52	19.80	97.21	16.80	97.43	12.40
	PTQ	96.13	18.00	96.43	21.00	96.46	20.60	97.20	16.00	97.43	11.60
	SmoothQuant	96.16	19.00	96.23	20.00	96.26	20.00	97.05	16.80	97.26	13.40
	RaanA	96.21	19.00	96.41	19.80	96.44	19.80	97.17	16.20	97.42	10.80
ID: ImageNet20 OOD: ImageNet10	FP32	93.59	23.90	93.13	44.50	93.27	44.40	94.05	37.80	94.88	34.10
	PTQ	93.19	24.40	92.52	43.40	92.66	42.50	93.65	37.20	94.49	33.90
	SmoothQuant	93.51	26.00	92.67	47.50	92.83	46.70	93.65	41.30	94.26	35.50
	RaanA	93.35	23.80	92.53	43.80	92.68	43.20	93.73	37.70	94.55	33.70
ID: IN1k OOD: DTD	FP32	71.78	81.36	84.45	61.26	84.63	61.58	80.21	68.01	68.03	87.42
	PTQ	71.50	79.62	84.47	59.93	84.62	60.27	81.79	62.76	70.91	82.22
	SmoothQuant	71.82	80.02	83.50	62.17	83.71	62.40	79.70	68.52	68.26	85.33
	RaanA	71.55	80.29	84.52	60.30	84.66	60.66	82.07	61.29	71.33	82.55
ID: IN1k OOD: iNaturalist	FP32	74.17	70.61	67.97	81.03	68.90	80.20	81.46	62.81	79.87	69.16
	PTQ	74.18	70.95	67.14	82.30	68.12	81.53	79.19	67.91	77.17	74.16
	SmoothQuant	74.18	70.54	67.30	81.66	68.27	81.19	80.39	65.27	78.02	73.06
	RaanA	74.07	71.01	67.40	81.65	68.35	80.78	79.81	66.95	77.79	73.47
ID: IN1k OOD: Places365	FP32	70.93	78.13	90.05	39.83	89.98	40.30	87.34	51.01	84.35	61.66
	PTQ	70.32	78.76	90.40	39.83	90.32	39.97	88.62	45.55	86.14	55.21
	SmoothQuant	70.13	78.61	89.98	39.99	89.90	40.38	87.21	52.01	84.99	57.90
	RaanA	70.15	78.72	90.01	40.64	89.93	41.26	87.34	50.21	85.61	56.22
ID: IN1k OOD: SUN397	FP32	70.67	79.87	88.36	41.08	88.34	41.60	87.17	47.93	84.56	62.05
	PTQ	70.28	80.30	88.22	42.07	88.18	42.99	86.61	49.42	84.43	60.64
	SmoothQuant	70.29	80.15	88.23	40.93	88.20	41.29	86.64	48.87	84.55	60.55
	RaanA	70.05	79.99	88.19	41.52	88.14	42.13	86.22	50.46	83.81	62.48

Table 7: PTQ Results: Out-of-Distribution (OOD) Detection for LAION-400M pre-trained ViT-L-14. A/F denote AU-ROC/FPR95.

Scenario		Method	Out-of-Distribution Detection Performance									
			Traditional Methods				VLM-based Methods					
			MSP		Energy		MCM		Gen.Neg.		NegLabel	
			A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	90.57	33.60	83.49	67.22	87.39	58.47	92.00	36.77	91.09	37.54	
	PTQ	63.01	75.92	64.65	82.22	66.21	78.02	70.76	72.62	70.25	72.22	
	SmoothQuant	78.78	52.99	75.25	71.90	78.91	63.48	83.43	52.09	82.42	54.56	
	RaanA	63.26	76.51	66.46	81.41	67.74	77.49	70.93	72.85	70.69	72.03	
ID: CIFAR100 OOD: CIFAR10	FP32	68.74	70.47	71.52	75.02	73.88	66.99	71.53	64.23	68.91	65.49	
	PTQ	55.84	89.90	52.94	90.71	54.33	86.98	51.59	90.50	53.46	90.27	
	SmoothQuant	53.45	86.77	58.33	85.99	57.33	81.73	53.23	83.66	52.07	85.50	
	RaanA	55.24	90.19	53.23	90.47	54.37	86.76	51.37	90.37	53.31	89.46	
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	74.23	70.75	69.24	77.94	71.72	75.67	74.12	72.71	73.78	72.51	
	PTQ	54.01	90.00	52.65	91.78	53.54	90.54	53.65	90.25	52.96	90.63	
	SmoothQuant	58.23	85.50	56.67	88.85	58.09	86.82	58.34	86.48	57.40	86.43	
	RaanA	54.65	88.68	53.18	91.40	54.13	89.83	54.34	89.82	53.65	89.94	
ID: ImageNet10 OOD: ImageNet20	FP32	97.64	10.20	96.82	18.20	97.03	17.40	96.63	25.40	97.62	12.20	
	PTQ	72.46	62.00	76.41	65.40	76.84	64.20	77.96	64.40	80.26	64.40	
	SmoothQuant	82.08	47.40	84.44	57.40	85.01	55.60	85.75	57.00	86.20	49.80	
	RaanA	75.95	60.00	77.61	60.40	78.46	61.40	79.99	62.40	81.75	61.40	
ID: ImageNet20 OOD: ImageNet10	FP32	95.42	15.90	94.16	34.80	94.82	30.00	93.80	40.20	95.44	26.40	
	PTQ	72.39	62.90	76.76	71.90	77.87	67.80	77.13	66.90	75.70	64.10	
	SmoothQuant	81.85	47.50	82.61	62.90	84.02	58.30	82.55	64.70	83.24	56.20	
	RaanA	73.05	64.80	77.36	68.70	78.63	65.00	77.54	68.20	76.47	65.90	
ID: IN1k OOD: DTD	FP32	78.95	80.22	72.63	69.47	77.94	69.21	85.92	57.17	87.94	58.39	
	PTQ	63.10	86.87	50.85	89.83	57.41	86.78	65.86	78.81	68.68	78.14	
	SmoothQuant	68.84	82.87	61.05	81.81	67.07	80.42	74.68	71.92	77.81	68.83	
	RaanA	63.60	86.31	52.49	88.84	58.77	85.35	67.11	77.52	69.48	76.98	
ID: IN1k OOD: iNaturalist	FP32	73.36	78.80	64.24	86.71	68.15	85.35	85.89	57.70	77.04	77.64	
	PTQ	57.72	86.75	49.04	93.98	52.62	90.44	72.47	78.18	69.43	82.49	
	SmoothQuant	63.13	85.42	57.89	89.79	61.67	88.17	75.59	75.40	71.45	84.23	
	RaanA	61.80	84.75	50.92	93.61	56.21	89.96	74.71	77.97	72.71	81.99	
ID: IN1k OOD: Places365	FP32	75.82	76.96	83.45	54.51	84.57	55.13	87.45	51.27	82.86	62.82	
	PTQ	58.98	88.70	67.44	78.35	65.98	78.69	73.03	72.98	60.85	81.72	
	SmoothQuant	61.05	86.16	73.06	71.59	71.92	72.03	75.30	69.02	67.68	75.97	
	RaanA	58.51	87.60	69.75	76.57	69.00	78.04	74.28	73.73	63.05	81.48	
ID: IN1k OOD: SUN397	FP32	75.21	78.80	82.12	54.36	83.42	57.12	84.87	59.73	83.83	63.00	
	PTQ	56.03	88.69	66.56	79.00	64.77	80.45	69.86	75.11	63.30	80.30	
	SmoothQuant	60.69	86.16	72.11	71.98	71.04	73.23	70.64	73.95	67.27	77.99	
	RaanA	57.10	88.48	65.72	80.84	64.83	80.54	70.25	74.37	64.56	79.55	

Table 8: PTQ Results: Out-of-Distribution (OOD) Detection for OpenAI pre-trained ViT-B-16. A/F denote AUROC/FPR95.

Scenario	Method	Out-of-Distribution Detection Performance									
		Traditional Methods					VLM-based Methods				
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	89.98	32.09	89.36	52.87	89.85	51.21	91.93	43.94	91.42	35.67
	PTQ	66.55	76.38	73.97	76.40	74.44	75.49	75.86	70.91	69.87	72.94
	SmoothQuant	88.31	36.06	88.20	57.18	88.84	55.73	90.21	52.53	89.58	42.18
	RaanA	65.76	77.27	73.31	77.31	73.84	76.02	74.76	73.47	68.92	76.57
ID: CIFAR100 OOD: CIFAR10	FP32	73.87	66.03	81.29	63.22	81.95	61.82	78.60	64.67	79.79	58.63
	PTQ	60.58	82.57	59.17	85.41	60.47	83.97	57.59	88.01	61.18	80.53
	SmoothQuant	71.21	70.54	79.28	65.68	80.01	63.95	76.38	68.68	76.18	64.24
	RaanA	60.68	82.22	59.56	84.83	60.74	83.43	57.88	88.02	61.34	80.25
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	72.54	71.62	72.45	77.36	73.14	76.54	71.82	77.36	72.14	75.16
	PTQ	63.70	81.32	63.89	85.95	64.65	85.20	63.15	85.70	62.58	83.53
	SmoothQuant	70.39	73.27	70.05	80.37	70.85	79.24	69.55	79.17	69.46	77.92
	RaanA	62.58	81.64	62.86	86.76	63.52	85.93	61.62	86.64	61.21	83.81
ID: ImageNet10 OOD: ImageNet20	FP32	95.84	17.40	96.25	23.40	96.29	23.20	95.84	26.80	95.89	20.00
	PTQ	87.23	43.00	91.14	46.80	91.20	45.60	90.04	50.40	89.99	44.00
	SmoothQuant	94.38	20.80	95.60	26.00	95.64	25.60	95.01	30.80	95.11	25.60
	RaanA	87.33	44.00	91.32	47.20	91.37	47.60	90.12	49.00	89.98	47.00
ID: ImageNet20 OOD: ImageNet10	FP32	93.60	21.20	94.54	28.20	94.69	26.70	94.26	30.60	95.20	26.30
	PTQ	86.43	38.90	88.79	50.80	89.21	48.50	86.66	58.60	88.66	51.80
	SmoothQuant	91.55	26.90	93.35	36.30	93.53	35.70	93.54	32.60	94.37	27.00
	RaanA	86.32	39.80	89.06	52.20	89.43	50.10	86.98	57.10	88.54	49.30
ID: IN1k OOD: DTD	FP32	69.85	83.21	82.99	61.78	83.24	62.93	82.18	61.24	72.11	76.08
	PTQ	64.73	84.63	74.95	72.48	75.66	72.66	71.80	73.30	65.04	81.08
	SmoothQuant	68.87	82.58	79.41	67.82	80.03	68.59	80.27	61.31	70.53	77.26
	RaanA	62.64	85.79	74.20	74.44	74.69	75.60	70.81	73.57	61.22	85.12
ID: IN1k OOD: iNaturalist	FP32	71.78	73.63	65.81	83.41	67.11	82.66	73.62	76.13	74.16	76.32
	PTQ	64.24	81.93	59.36	88.24	60.69	87.83	64.07	81.61	60.95	84.51
	SmoothQuant	70.41	74.73	64.00	85.02	65.48	84.12	71.63	77.49	70.87	81.17
	RaanA	63.97	82.12	59.90	89.98	61.19	89.41	64.62	79.63	62.51	83.69
ID: IN1k OOD: Places365	FP32	68.23	80.60	88.11	44.49	88.01	45.42	81.50	62.97	74.52	76.56
	PTQ	61.32	86.05	84.30	52.29	83.98	54.21	71.36	77.68	64.82	83.67
	SmoothQuant	66.65	81.52	87.48	47.37	87.37	48.23	80.79	65.40	72.31	78.32
	RaanA	60.24	86.51	84.77	52.05	84.27	53.82	73.09	76.37	65.02	84.04
ID: IN1k OOD: SUN397	FP32	65.92	83.41	86.35	45.78	86.19	47.16	85.92	51.25	81.65	63.22
	PTQ	56.13	88.73	78.82	61.79	77.98	64.44	77.14	66.07	70.75	74.63
	SmoothQuant	63.28	84.94	84.52	52.73	84.30	54.36	83.62	57.47	78.43	66.85
	RaanA	55.80	89.12	76.80	64.46	76.11	67.06	77.05	67.17	70.81	73.91

Table 9: PTQ Results: Out-of-Distribution (OOD) Detection for LAION-400M pre-trained ViT-B-16. A/F denote AU-ROC/FPR95.

Out-of-Distribution Detection Performance											
Scenario	Method	Traditional Methods				VLM-based Methods					
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	90.56	34.04	86.00	61.30	88.36	55.10	91.88	39.47	91.00	37.72
	PTQ	89.46	36.47	85.50	59.62	87.91	54.74	91.06	42.84	89.90	40.26
	SmoothQuant	90.43	35.26	85.13	62.61	87.81	56.85	91.61	42.20	90.80	38.58
	RaanA	89.29	36.89	85.73	59.09	88.06	53.73	90.84	44.01	89.77	41.34
ID: CIFAR100 OOD: CIFAR10	FP32	72.39	70.08	73.32	70.56	76.21	66.38	75.59	65.61	71.32	69.38
	PTQ	71.51	70.98	74.55	70.90	76.89	66.64	75.64	66.72	69.70	71.23
	SmoothQuant	72.49	69.73	74.05	70.84	76.79	66.39	75.96	65.28	70.82	70.48
	RaanA	71.46	71.42	73.86	71.33	76.28	67.58	75.12	66.25	69.54	71.89
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	73.03	73.24	68.66	81.34	70.60	79.13	72.56	78.69	72.55	75.29
	PTQ	71.49	75.12	67.04	83.03	69.09	81.16	71.15	80.11	70.95	77.63
	SmoothQuant	72.42	74.44	68.13	81.75	70.13	79.69	72.00	78.66	71.76	76.13
	RaanA	71.17	75.10	67.03	82.37	69.02	80.66	71.10	80.04	70.55	77.32
ID: ImageNet10 OOD: ImageNet20	FP32	96.28	18.20	95.41	27.00	96.68	16.60	95.58	25.00	95.52	27.20
	PTQ	96.27	19.00	95.31	27.20	96.27	20.20	95.46	26.80	95.44	26.40
	SmoothQuant	96.17	19.80	95.46	24.60	96.68	17.20	95.61	23.80	95.59	24.80
	RaanA	96.21	20.20	95.37	25.60	96.42	19.20	95.57	23.40	95.52	25.40
ID: ImageNet20 OOD: ImageNet10	FP32	93.27	25.20	92.76	37.40	94.67	27.60	93.74	32.70	93.32	35.20
	PTQ	92.43	27.80	92.50	39.90	93.80	28.80	93.29	37.70	93.05	38.00
	SmoothQuant	92.67	29.80	92.19	39.70	94.17	26.90	93.36	35.30	94.17	38.10
	RaanA	92.20	30.40	92.11	38.40	93.61	31.50	93.07	41.40	93.61	36.90
ID: IN1k OOD: DTD	FP32	75.74	78.68	71.53	71.85	74.98	70.57	84.84	54.54	86.63	51.98
	PTQ	74.94	80.37	68.14	77.23	72.14	76.59	83.20	59.92	84.53	61.04
	SmoothQuant	76.14	78.63	70.34	73.56	74.27	72.44	84.51	55.66	86.53	53.55
	RaanA	74.54	80.07	68.58	76.63	72.40	76.11	83.30	59.36	84.58	59.62
ID: IN1k OOD: iNaturalist	FP32	71.34	76.14	59.91	86.56	63.39	85.54	77.41	72.99	73.56	78.67
	PTQ	70.96	76.44	59.72	87.61	63.43	86.32	76.64	75.01	72.30	81.64
	SmoothQuant	70.79	76.94	59.09	86.67	62.72	85.89	75.39	76.23	72.63	81.11
	RaanA	70.84	76.76	58.89	87.92	62.64	86.70	76.24	75.46	73.12	80.72
ID: IN1k OOD: Places365	FP32	75.00	76.47	86.43	47.10	86.90	48.23	82.50	59.85	78.99	66.33
	PTQ	74.18	77.88	84.02	53.87	84.80	54.66	80.83	62.88	77.44	69.38
	SmoothQuant	74.77	76.57	84.62	51.67	85.38	52.56	81.11	62.03	80.21	63.97
	RaanA	73.67	78.01	83.56	53.97	84.36	54.64	80.83	63.33	77.24	70.43
ID: IN1k OOD: SUN397	FP32	73.65	77.33	82.75	52.08	83.60	53.40	80.34	61.72	80.50	63.08
	PTQ	72.23	78.85	81.02	57.21	81.89	58.64	78.09	65.94	78.42	66.65
	SmoothQuant	73.35	77.49	80.76	56.54	81.91	57.47	78.79	64.99	79.49	65.34
	RaanA	72.94	78.51	81.43	56.09	82.40	56.98	79.02	65.39	79.23	65.51

Table 10: PTQ Results: Out-of-Distribution (OOD) Detection for OpenAI pre-trained ViT-B-32. A/F denote AUROC/FPR95.

Scenario	Method	Out-of-Distribution Detection Performance									
		Traditional Methods					VLM-based Methods				
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	88.65	37.58	88.38	59.87	88.90	58.14	90.40	48.46	89.25	46.25
	PTQ	76.01	59.28	83.30	61.75	83.55	60.72	84.39	56.14	76.34	60.80
	SmoothQuant	84.55	45.77	86.41	62.56	87.06	60.75	88.86	54.26	86.14	51.30
	RaanA	76.73	58.29	83.45	62.00	83.80	60.36	84.69	55.94	77.19	59.68
ID: CIFAR100 OOD: CIFAR10	FP32	73.91	69.47	79.34	67.00	80.18	64.95	78.60	71.44	76.42	64.60
	PTQ	63.04	80.43	64.64	81.27	65.64	80.50	62.96	82.02	60.69	79.43
	SmoothQuant	67.02	76.14	72.75	74.15	73.57	72.16	71.94	75.10	67.94	72.29
	RaanA	63.87	79.42	66.69	79.63	67.65	78.45	64.37	79.90	61.05	79.39
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	70.62	72.69	70.24	80.11	71.01	79.01	70.16	77.81	67.69	80.64
	PTQ	64.40	79.80	64.38	85.00	65.21	84.13	64.25	84.09	60.63	86.68
	SmoothQuant	66.25	77.82	65.99	84.33	66.81	83.44	65.99	82.50	62.68	85.78
	RaanA	64.49	79.95	64.55	85.27	65.39	84.37	64.54	83.82	60.81	86.52
ID: ImageNet10 OOD: ImageNet20	FP32	93.84	23.00	94.43	34.60	94.48	34.60	94.72	31.00	92.54	39.60
	PTQ	88.75	35.80	90.85	45.20	91.01	44.80	90.28	42.40	86.95	48.20
	SmoothQuant	91.10	33.40	91.80	43.60	91.93	43.60	92.24	37.80	89.37	45.20
	RaanA	88.19	38.40	91.09	46.20	91.22	44.80	90.20	42.00	86.66	46.60
ID: ImageNet20 OOD: ImageNet10	FP32	90.51	35.80	89.86	61.90	90.10	61.20	90.87	43.90	86.43	60.80
	PTQ	84.90	42.70	88.11	55.60	88.51	52.50	88.43	47.70	81.53	62.50
	SmoothQuant	86.28	41.20	87.79	54.60	88.17	52.30	87.85	48.60	81.68	63.50
	RaanA	85.23	40.00	88.05	59.50	88.42	59.80	87.82	50.20	81.44	62.20
ID: IN1k OOD: DTD	FP32	68.95	83.40	79.77	63.72	80.21	64.18	71.85	79.26	73.79	68.17
	PTQ	64.80	85.20	74.43	67.54	74.93	69.28	68.93	79.56	68.96	72.72
	SmoothQuant	66.36	84.69	76.86	66.04	77.41	67.11	69.58	79.54	70.43	71.44
	RaanA	64.25	85.70	75.91	66.15	76.19	68.04	69.35	79.13	66.84	74.49
ID: IN1k OOD: iNaturalist	FP32	70.75	74.67	60.90	85.31	62.60	84.32	69.25	79.83	67.49	80.24
	PTQ	65.26	81.34	58.66	88.29	60.30	87.48	68.51	82.09	65.26	83.20
	SmoothQuant	66.89	79.35	61.02	86.82	62.65	85.57	69.82	81.72	64.39	83.39
	RaanA	66.14	80.49	58.72	88.02	60.45	87.32	67.70	83.32	63.55	84.60
ID: IN1k OOD: Places365	FP32	66.73	81.39	85.48	48.52	85.42	50.42	81.56	63.67	80.81	59.51
	PTQ	59.78	86.64	84.51	53.30	83.99	54.36	79.46	66.31	69.65	74.91
	SmoothQuant	62.13	84.66	83.54	55.38	83.21	57.19	79.21	67.21	74.05	69.60
	RaanA	58.90	87.04	83.16	57.52	81.86	59.79	77.66	68.77	68.15	75.90
ID: IN1k OOD: SUN397	FP32	64.96	83.53	84.67	48.12	84.51	50.12	81.93	60.14	81.72	54.58
	PTQ	57.63	87.95	82.70	55.79	81.94	58.43	79.10	64.49	69.81	69.12
	SmoothQuant	60.08	86.78	82.49	53.51	81.96	56.16	80.53	63.21	74.55	65.13
	RaanA	56.18	88.86	82.44	54.16	81.49	57.11	78.63	65.76	67.93	70.22

Table 11: PTQ Results: Out-of-Distribution (OOD) Detection for LAION-400M pre-trained ViT-B-32. A/F denote AU-ROC/FPR95.

Scenario	Method	Acc (%)	ECE	NLL
ID: CIFAR10	<i>FP32</i>	<i>93.59</i>	<i>4.260</i>	<i>0.229</i>
	PTQ	66.20	8.663	1.054
	SmoothQuant	79.09	8.274	0.674
	RaanA	66.81	8.059	1.033
ID: CIFAR100	<i>FP32</i>	<i>74.45</i>	<i>13.168</i>	<i>1.090</i>
	PTQ	33.56	8.895	3.013
	SmoothQuant	47.42	12.591	2.292
	RaanA	35.37	9.071	2.926
ID: ImageNet-X	<i>FP32</i>	<i>81.48</i>	<i>2.185</i>	<i>0.682</i>
	PTQ	52.23	5.503	2.084
	SmoothQuant	62.78	4.066	1.495
	RaanA	54.74	4.745	1.937

(a) OpenAI pre-trained ViT-L-14.

Scenario	Method	Acc (%)	ECE	NLL
ID: CIFAR10	<i>FP32</i>	<i>94.62</i>	<i>2.986</i>	<i>0.207</i>
	PTQ	93.17	3.709	0.272
	SmoothQuant	94.50	2.952	0.207
	RaanA	93.81	3.322	0.237
ID: CIFAR100	<i>FP32</i>	<i>78.03</i>	<i>11.065</i>	<i>0.946</i>
	PTQ	76.42	11.357	0.997
	SmoothQuant	77.52	10.929	0.957
	RaanA	76.31	11.340	1.014
ID: ImageNet-X	<i>FP32</i>	<i>80.38</i>	<i>9.394</i>	<i>0.902</i>
	PTQ	79.36	9.937	0.930
	SmoothQuant	79.95	9.477	0.916
	RaanA	79.44	9.854	0.930

(b) LAION-400M pre-trained ViT-L-14.

Table 12: In-Distribution PTQ Results for ViT-L-14 models (multiple scenarios shown).

Scenario	Method	Acc (%)	ECE	NLL
ID: CIFAR10	<i>FP32</i>	<i>88.13</i>	<i>6.917</i>	<i>0.410</i>
	PTQ	46.80	5.957	1.449
	SmoothQuant	72.17	10.087	0.850
	RaanA	46.19	5.025	1.457
ID: CIFAR100	<i>FP32</i>	<i>62.79</i>	<i>16.012</i>	<i>1.590</i>
	PTQ	19.72	6.513	3.600
	SmoothQuant	37.31	15.698	2.841
	RaanA	20.07	6.514	3.600
ID: ImageNet-X	<i>FP32</i>	<i>74.23</i>	<i>5.427</i>	<i>0.954</i>
	PTQ	24.67	1.087	3.751
	SmoothQuant	40.20	5.604	2.735
	RaanA	26.62	1.809	3.629

(a) OpenAI pre-trained ViT-B-16.

Scenario	Method	Acc (%)	ECE	NLL
ID: CIFAR10	<i>FP32</i>	<i>91.70</i>	<i>3.658</i>	<i>0.276</i>
	PTQ	64.41	9.813	1.052
	SmoothQuant	90.38	3.443	0.313
	RaanA	62.13	9.966	1.116
ID: CIFAR100	<i>FP32</i>	<i>71.43</i>	<i>10.843</i>	<i>1.119</i>
	PTQ	35.24	15.717	2.810
	SmoothQuant	68.02	10.652	1.233
	RaanA	35.83	16.640	2.764
ID: ImageNet-X	<i>FP32</i>	<i>75.68</i>	<i>9.526</i>	<i>1.030</i>
	PTQ	56.09	16.584	2.083
	SmoothQuant	70.59	11.676	1.261
	RaanA	52.44	18.272	2.302

(b) LAION-400M pre-trained ViT-B-16.

Table 13: In-Distribution PTQ Results for ViT-B-16 models (multiple scenarios shown).

Scenario	Method	Acc (%)	ECE	NLL
ID: CIFAR10	<i>FP32</i>	<i>89.76</i>	<i>5.031</i>	<i>0.338</i>
	PTQ	87.71	3.722	0.385
	SmoothQuant	89.36	4.703	0.347
	RaanA	87.84	4.118	0.384
ID: CIFAR100	<i>FP32</i>	<i>65.10</i>	<i>10.632</i>	<i>1.341</i>
	PTQ	63.35	9.529	1.406
	SmoothQuant	64.91	10.303	1.345
	RaanA	63.32	9.296	1.411
ID: ImageNet-X	<i>FP32</i>	<i>73.23</i>	<i>1.687</i>	<i>0.983</i>
	PTQ	70.48	2.052	1.105
	SmoothQuant	72.14	1.901	1.036
	RaanA	69.76	2.003	1.132

(a) OpenAI pre-trained ViT-B-32.

Scenario	Method	Acc (%)	ECE	NLL
ID: CIFAR10	<i>FP32</i>	<i>90.75</i>	<i>3.644</i>	<i>0.304</i>
	PTQ	72.15	6.395	0.840
	SmoothQuant	86.59	3.862	0.411
	RaanA	73.18	6.356	0.807
ID: CIFAR100	<i>FP32</i>	<i>70.59</i>	<i>11.284</i>	<i>1.122</i>
	PTQ	49.51	8.121	1.977
	SmoothQuant	60.69	10.403	1.503
	RaanA	52.45	7.214	1.827
ID: ImageNet-X	<i>FP32</i>	<i>71.91</i>	<i>10.641</i>	<i>1.207</i>
	PTQ	58.12	14.431	1.902
	SmoothQuant	61.76	14.062	1.705
	RaanA	58.30	14.379	1.897

(b) LAION-400M pre-trained ViT-B-32.

Table 14: In-Distribution PTQ Results for ViT-B-32 models (multiple scenarios shown).

Scenario	Method	Out-of-Distribution Detection Performance									
		Traditional Methods					VLM-based Methods				
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	FP32	90.56	34.04	85.99	61.30	88.36	55.10	91.88	39.47	90.97	37.57
	Basic QAT	90.92	34.39	88.07	53.28	90.44	46.76	92.38	39.18	91.52	36.39
	QAT-LSQ	90.83	33.92	87.05	57.61	89.22	51.65	92.25	41.33	91.18	38.79
	QAT-LSQ (PC)	90.56	34.18	88.29	52.99	90.12	48.20	92.20	39.19	91.05	37.45
	QAT-Distill	90.86	33.70	87.63	54.56	89.74	49.30	92.51	38.62	91.38	36.75
ID: CIFAR100 OOD: CIFAR10	FP32	72.39	70.08	73.32	70.56	76.21	66.38	75.59	65.61	70.98	68.86
	Basic QAT	75.26	67.32	78.53	64.38	80.46	61.27	77.84	64.76	72.65	69.25
	QAT-LSQ	75.08	67.31	78.26	66.03	80.16	63.19	77.55	65.12	72.09	69.12
	QAT-LSQ (PC)	75.16	66.63	77.63	65.85	79.81	62.53	78.27	62.11	72.96	66.72
	QAT-Distill	74.43	68.16	76.65	66.92	79.33	64.20	78.01	63.56	73.08	67.14
ID: ImageNet-X OOD: ImageNet-X OOD	FP32	73.03	73.24	68.66	81.34	70.60	79.13	72.56	78.69	72.82	74.97
	Basic QAT	72.67	73.35	69.25	80.50	71.01	78.42	72.45	78.39	72.63	75.14
	QAT-LSQ	72.63	73.68	68.94	80.44	70.79	78.72	72.66	77.66	72.76	74.73
	QAT-LSQ (PC)	72.69	72.98	69.21	80.26	71.03	78.19	72.59	77.94	72.83	74.66
	QAT-Distill	72.74	73.65	68.96	81.00	70.87	78.97	72.27	78.96	72.47	75.40
ID: ImageNet10 OOD: ImageNet20	FP32	96.28	18.20	95.41	27.00	95.52	27.20	95.58	25.00	96.72	17.40
	Basic QAT	96.38	18.40	95.94	22.00	96.07	20.60	95.97	22.00	96.95	14.40
	QAT-LSQ	96.55	16.20	95.87	25.20	95.98	24.00	95.86	24.40	96.85	15.40
	QAT-LSQ (PC)	96.25	18.60	95.83	24.80	95.93	24.00	95.80	25.40	96.83	14.00
	QAT-Distill	95.90	20.60	95.38	26.80	95.56	26.80	95.47	25.80	96.61	17.00
ID: ImageNet20 OOD: ImageNet10	FP32	93.27	25.20	92.76	37.40	93.32	35.20	93.74	32.70	94.87	27.00
	Basic QAT	93.08	26.10	93.17	34.60	93.54	36.90	93.73	37.60	94.73	25.80
	QAT-LSQ	93.19	26.70	93.04	37.80	93.61	33.40	93.83	36.90	94.66	27.80
	QAT-LSQ (PC)	92.69	29.60	93.08	35.80	93.58	33.80	93.98	31.60	94.84	24.90
	QAT-Distill	93.03	26.10	92.93	40.10	93.71	33.00	93.76	35.20	94.73	26.80
ID: IN1k OOD: DTD	FP32	75.74	78.68	71.53	71.85	74.98	70.57	84.84	54.54	86.16	53.69
	Basic QAT	74.59	80.81	70.04	75.27	73.72	73.41	84.12	56.10	85.33	55.93
	QAT-LSQ	74.46	80.96	70.40	74.13	73.69	73.45	84.24	57.03	85.14	57.73
	QAT-LSQ (PC)	75.10	80.32	71.05	72.89	74.37	72.21	84.52	55.52	85.60	55.48
	QAT-Distill	74.92	80.06	70.20	74.16	73.98	72.77	84.21	54.97	85.58	55.30
ID: IN1k OOD: iNaturalist	FP32	71.34	76.14	59.91	86.56	63.39	85.54	77.41	72.99	74.70	77.01
	Basic QAT	71.55	75.46	60.82	86.81	64.38	85.73	76.27	75.45	74.55	78.06
	QAT-LSQ	71.67	76.25	62.53	85.48	65.69	84.53	77.50	73.73	74.85	77.72
	QAT-LSQ (PC)	71.32	76.31	61.89	85.59	65.04	84.74	76.49	74.83	73.79	79.20
	QAT-Distill	71.57	76.08	60.09	87.52	63.68	86.08	76.84	74.42	74.95	77.78
ID: IN1k OOD: Places365	FP32	75.00	76.47	86.43	47.10	86.90	48.23	82.50	59.85	78.46	67.33
	Basic QAT	74.15	77.18	84.50	52.40	85.21	53.55	80.78	63.49	78.89	67.91
	QAT-LSQ	73.97	77.98	85.63	49.47	85.99	50.45	81.00	63.48	79.08	66.95
	QAT-LSQ (PC)	74.26	76.95	85.98	48.77	86.34	50.35	81.03	62.30	78.96	66.92
	QAT-Distill	74.35	77.15	85.51	50.41	86.20	51.47	81.42	61.55	77.84	68.21
ID: IN1k OOD: SUN397	FP32	73.65	77.33	82.75	52.08	83.60	53.40	80.34	61.72	79.58	66.08
	Basic QAT	72.36	79.03	81.30	55.59	82.21	57.13	78.74	64.39	78.05	68.67
	QAT-LSQ	72.61	78.28	82.23	53.87	82.89	54.96	79.16	65.10	78.00	69.11
	QAT-LSQ (PC)	72.83	78.49	81.88	53.66	82.62	55.65	79.11	64.15	78.65	68.26
	QAT-Distill	73.09	78.12	82.15	54.15	83.06	55.25	79.52	63.00	78.86	67.28

Table 15: Light QAT Results: Out-of-Distribution (OOD) Detection for OpenAI pre-trained ViT-B-32. A/F denote AU-ROC/FPR95.

Scenario	Method	Out-of-Distribution Detection Performance									
		Traditional Methods					VLM-based Methods				
		MSP		Energy		MCM		Gen.Neg.		NegLabel	
		A	F	A	F	A	F	A	F	A	F
ID: CIFAR10 OOD: CIFAR100	<i>FP32</i>	88.65	37.58	88.38	59.87	88.90	58.14	90.40	48.46	88.30	48.87
	Basic QAT	83.18	47.76	87.61	56.90	87.25	57.33	88.83	51.16	83.88	54.59
	QAT-LSQ	83.95	47.10	87.51	59.40	88.00	57.35	89.03	50.91	84.07	54.74
	QAT-LSQ (PC)	83.59	47.45	86.94	60.76	87.47	58.88	88.86	51.38	83.84	55.87
	QAT-Distill	82.70	49.47	87.09	59.18	86.11	60.09	88.72	52.31	83.62	55.67
ID: CIFAR100 OOD: CIFAR10	<i>FP32</i>	73.91	69.47	79.34	67.00	80.18	64.95	78.59	71.44	75.78	65.56
	Basic QAT	65.36	76.02	73.07	72.58	73.34	72.82	71.60	76.47	65.36	75.43
	QAT-LSQ	65.46	75.43	72.35	74.20	73.19	72.96	70.37	77.94	65.46	77.15
	QAT-LSQ (PC)	66.33	74.50	73.10	72.90	73.97	71.72	71.70	75.51	66.33	75.58
	QAT-Distill	66.28	75.04	73.41	72.19	73.36	71.35	72.16	75.73	66.28	75.13
ID: ImageNet-X OOD: ImageNet-X OOD	<i>FP32</i>	70.62	72.69	70.24	80.11	71.01	79.01	70.16	77.81	66.74	81.99
	Basic QAT	67.08	77.14	67.27	82.44	67.77	81.74	66.78	80.94	62.79	85.22
	QAT-LSQ	67.07	77.01	67.13	82.72	67.93	81.47	66.74	81.36	62.63	85.62
	QAT-LSQ (PC)	67.46	76.80	67.23	82.52	68.06	81.18	67.01	80.83	63.00	84.93
	QAT-Distill	67.51	76.75	67.33	82.26	68.30	81.15	67.04	81.06	62.91	85.45
ID: ImageNet10 OOD: ImageNet20	<i>FP32</i>	93.84	23.00	94.43	34.60	94.48	34.60	94.72	31.00	93.39	37.00
	Basic QAT	92.22	32.40	92.53	46.00	92.77	45.60	92.51	40.20	90.74	42.60
	QAT-LSQ	92.15	32.60	92.74	40.00	92.89	39.40	92.46	35.80	90.63	42.20
	QAT-LSQ (PC)	91.68	30.80	92.55	40.40	92.71	39.60	92.27	34.80	90.39	43.20
	QAT-Distill	91.41	30.80	92.70	37.80	92.73	39.40	92.50	37.80	90.51	41.40
ID: ImageNet20 OOD: ImageNet10	<i>FP32</i>	90.51	35.80	89.86	61.90	90.10	61.20	90.87	43.90	87.64	59.10
	Basic QAT	88.77	39.30	88.69	57.90	89.39	55.20	88.96	48.40	84.91	61.90
	QAT-LSQ	86.91	40.90	88.52	55.60	88.84	54.20	88.87	53.20	84.52	60.60
	QAT-LSQ (PC)	86.95	41.00	88.32	60.80	88.63	59.60	88.76	50.30	84.63	62.90
	QAT-Distill	86.92	39.10	88.74	55.80	88.95	55.10	88.95	49.70	84.65	64.70
ID: IN1k OOD: DTD	<i>FP32</i>	68.95	83.40	79.77	63.72	80.21	64.18	71.85	79.26	73.69	68.37
	Basic QAT	66.51	85.35	78.42	62.83	78.57	64.53	72.52	76.19	71.89	69.61
	QAT-LSQ	66.62	84.48	78.26	62.41	78.62	64.20	71.65	76.82	72.27	68.93
	QAT-LSQ (PC)	66.61	85.18	79.37	61.03	79.67	63.09	72.48	75.86	71.90	69.30
	QAT-Distill	66.80	85.16	78.21	63.07	78.29	65.10	72.78	76.23	72.62	68.49
ID: IN1k OOD: iNaturalist	<i>FP32</i>	70.75	74.67	60.90	85.31	62.60	84.32	69.25	79.83	66.28	80.94
	Basic QAT	67.42	79.30	61.87	85.41	63.17	84.72	70.44	80.21	65.58	80.98
	QAT-LSQ	68.28	77.59	60.63	86.29	62.33	85.30	69.74	80.10	65.77	81.12
	QAT-LSQ (PC)	68.02	78.29	61.26	85.99	62.88	85.05	70.06	79.41	65.84	80.68
	QAT-Distill	67.57	79.11	61.02	85.78	63.18	84.34	70.07	79.71	65.38	80.79
ID: IN1k OOD: Places365	<i>FP32</i>	66.73	81.39	85.48	48.52	85.42	50.42	81.56	63.67	79.58	61.00
	Basic QAT	62.00	84.94	85.01	51.33	84.79	52.89	81.52	61.78	74.72	66.96
	QAT-LSQ	62.12	85.08	84.83	52.28	84.43	53.95	81.76	61.54	75.40	65.61
	QAT-LSQ (PC)	62.59	84.73	85.14	50.70	84.76	52.29	81.94	61.10	75.76	64.53
	QAT-Distill	62.41	84.75	85.13	51.05	84.67	53.08	81.56	62.10	75.45	65.62
ID: IN1k OOD: SUN397	<i>FP32</i>	64.96	83.53	84.67	48.12	84.51	50.12	81.93	60.14	80.99	58.21
	Basic QAT	60.18	87.18	84.21	49.84	83.51	52.21	80.58	61.65	73.40	67.46
	QAT-LSQ	60.04	87.27	84.42	49.96	83.74	52.44	80.44	62.01	73.57	67.88
	QAT-LSQ (PC)	60.90	87.06	84.32	50.23	83.76	52.74	81.23	61.67	74.28	67.50
	QAT-Distill	59.96	87.27	84.14	49.93	83.64	52.17	80.43	62.40	72.89	68.94

Table 16: Light QAT Results: Out-of-Distribution (OOD) Detection for LAION-400M pre-trained ViT-B-32. A/F denote AU-ROC/FPR95.

Scenario	Method	OpenAI ViT-B-32			LAION ViT-B-32		
		Acc(%)	ECE	NLL	Acc(%)	ECE	NLL
CIFAR10	<i>FP32</i>	<i>89.76</i>	<i>5.03</i>	<i>0.34</i>	<i>90.75</i>	<i>3.64</i>	<i>0.30</i>
	Basic QAT	90.25	2.91	0.31	83.90	4.36	0.48
	QAT-LSQ	90.07	3.23	0.32	85.90	3.84	0.42
	QAT-LSQ (PC)	89.49	3.21	0.33	85.60	3.92	0.43
	QAT-Distill	90.09	3.49	0.31	81.78	4.76	0.54
CIFAR100	<i>FP32</i>	<i>65.10</i>	<i>10.63</i>	<i>1.34</i>	<i>70.59</i>	<i>11.28</i>	<i>1.12</i>
	Basic QAT	67.27	6.72	1.23	61.95	8.62	1.43
	QAT-LSQ	66.63	7.31	1.25	62.63	8.86	1.39
	QAT-LSQ (PC)	67.08	7.74	1.24	63.72	8.73	1.34
	QAT-Distill	66.83	8.40	1.25	63.45	8.86	1.37
ImageNet-X	<i>FP32</i>	<i>73.23</i>	<i>1.69</i>	<i>0.98</i>	<i>71.91</i>	<i>10.64</i>	<i>1.21</i>
	Basic QAT	72.70	1.28	1.00	64.21	12.91	1.57
	QAT-LSQ	72.93	1.77	0.99	64.31	12.73	1.56
	QAT-LSQ (PC)	73.15	1.73	0.99	64.66	12.88	1.55
	QAT-Distill	72.89	1.78	1.00	64.94	12.78	1.53
ImageNet10	<i>FP32</i>	<i>98.80</i>	<i>0.74</i>	<i>0.03</i>	<i>97.60</i>	<i>1.37</i>	<i>0.06</i>
	Basic QAT	99.20	0.75	0.03	98.20	1.83	0.07
	QAT-LSQ	99.20	0.84	0.03	97.20	1.79	0.09
	QAT-LSQ (PC)	99.00	0.68	0.03	97.60	1.22	0.09
	QAT-Distill	99.00	0.55	0.03	97.80	1.52	0.08
ImageNet20	<i>FP32</i>	<i>94.90</i>	<i>1.60</i>	<i>0.16</i>	<i>95.80</i>	<i>1.68</i>	<i>0.16</i>
	Basic QAT	94.70	1.36	0.16	93.80	1.96	0.20
	QAT-LSQ	94.60	1.47	0.17	92.20	2.84	0.24
	QAT-LSQ (PC)	94.70	1.17	0.16	92.90	2.44	0.23
	QAT-Distill	94.50	1.47	0.16	93.20	2.92	0.22
IN1k	<i>FP32</i>	<i>63.36</i>	<i>1.86</i>	<i>1.40</i>	<i>62.89</i>	<i>13.72</i>	<i>1.66</i>
	Basic QAT	62.79	1.72	1.42	54.97	15.25	2.06
	QAT-LSQ	62.74	1.67	1.43	54.82	14.98	2.07
	QAT-LSQ (PC)	62.87	1.72	1.42	55.51	15.20	2.03
	QAT-Distill	62.81	1.73	1.42	55.92	15.21	2.01

Table 17: Consolidated In-Distribution Light QAT Results for OpenAI and LAION-400M pre-trained ViT-B-32 models.