

QUANTILE ACTIVATION: CORRECTING A FAILURE MODE OF ML MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

An established failure mode for machine learning models occurs when the same features are equally likely to belong to class 0 and class 1. In such cases, existing ML models cannot correctly classify the sample. However, a solvable case emerges when the probabilities of class 0 and 1 vary with the *context distribution*. To the best of our knowledge, standard neural network architectures like MLPs or CNNs are not equipped to handle this.

In this article, we propose a simple activation function, quantile activation (QACT), that addresses this problem without significantly increasing computational costs. The core idea is to “adapt” the outputs of each neuron to its *context distribution*. The proposed quantile activation, QACT, produces the *relative quantile* of the sample in its context distribution, rather than the actual values, as in traditional networks.

A practical example where the same sample can have different labels arises in cases of inherent distribution shift. We validate the proposed activation function under such shifts, using datasets designed to test robustness against distortions—CIFAR10C, CIFAR100C, MNISTC, TinyImagenetC. Our results demonstrate significantly better generalization across distortions compared to conventional classifiers, across various architectures. Although this paper presents a proof of concept, we find that this approach unexpectedly outperforms DINOv2 (small) under large distortions, despite DINOv2 being trained with a much larger network and dataset.

1 INTRODUCTION

Thanks to deep learning approaches, machine learning has been adopted across wide variety of domains. However, there is a significant failure mode within the standard framework of machine learning – A specific sample can potentially be labelled 0 or 1 depending on the context. Hence, standard ML models cannot handle these conflicts effectively.

A simple toy example to illustrate the failure mode of ML systems: We use a simple toy example to describe this failure mode. Consider the distribution generated as follows (figure 1a):

$$\begin{aligned}
 \mu_1 &\sim \mathcal{U}(S^1) && \text{random sample from uniform distribution on the circle} \\
 \mu_2 &= R(30^\circ)\mu_1 && \mu_2 \text{ is obtained by rotating } \mu_1 \text{ by } 30^\circ \\
 \text{Class 0} &\sim \mathcal{N}(\mu_1, 0.1\mathbf{I}) && \text{Class 0 generated using normal with mean } \mu_1 \text{ and stdev } 0.1 \\
 \text{Class 1} &\sim \mathcal{N}(\mu_2, 0.1\mathbf{I}) && \text{Class 1 generated using normal with mean } \mu_2 \text{ and stdev } 0.1
 \end{aligned} \tag{1}$$

Why is this classification problem hard? Note that any point x from the support of the above distribution is equally likely to be class 0 or class 1. So, one cannot construct any *fixed* function depending only on the input features x . Thus, most of existing ML frameworks, which insist on learning a fixed function, fail. Nevertheless, this is a valid distribution where classification with accuracy 1 is theoretically possible when one can reconstruct the latent μ_1, μ_2 .

Specifically, the current neural network architectures fail as well. Consider training a simple MLP on this dataset using gradient descent. That is, each batch (of size B) of samples is generated as - Sample μ_1, μ_2 , then sample $B/2$ points each from $\mathcal{N}(\mu_1, 0.1\mathbf{I})$ and $\mathcal{N}(\mu_2, 0.1\mathbf{I})$ respectively. Since, we

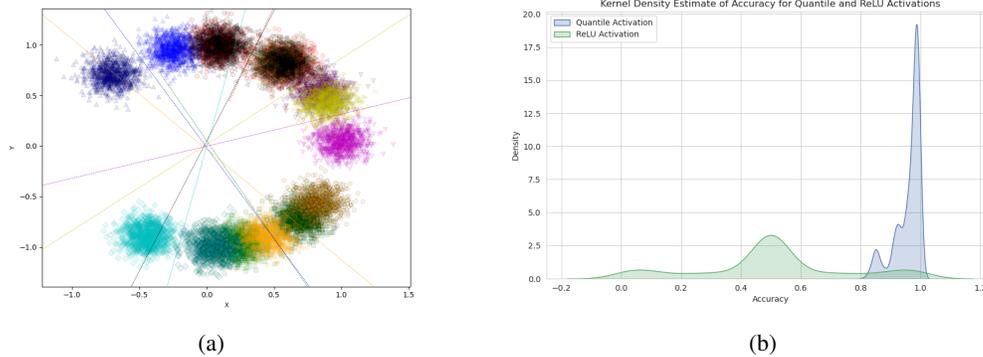


Figure 1: A simple toy example to illustrate where ML systems fail. (a) The distribution is a mixture of Gaussian distributions whose centers (μ_1, μ_2) are separated by 30° . The centres themselves can lie anywhere on the unit circle. (please refer to the text for exact description). (b) Histogram of accuracy over 1000 different combinations of μ_1, μ_2 for both ReLU activation and after incorporating QACT. Clearly, ReLU activation alone cannot perform better than random guess. Incorporating QACT on the other hand can easily infer the latent μ_1, μ_2 .

would have that a specific sample is equally likely to belong class 0 or 1, one would learn (probability) $p = 0.5$ for all the samples. This is verified in figure 1b where we use ReLU activation.

Remark: We remark that the failure mode described above is one of the possible failure modes, but there can exist other failure modes as well. In this article, we are interested in rectifying the above failure mode alone.

Fixing the failure mode and defining the “context”: The only way to fix the above failure mode is to infer the μ_1, μ_2 from the batch of samples. We define *context* as the batch of samples according to which each sample is processed. To our knowledge current neural network architectures such as MLP/CNN do not consider this. While transformers considers this to some extent via the self-attention module, it is still sample specific and is very expensive computationally to obtain self-attention for the entire batch. Moreover, vision transformers do not consider attention across different samples in the batch.

Quantile Activation can identify the context distribution from each batch: To rectify this shortcoming, we propose quantile activation, QACT, in this article. The key idea is that – at each neuron, the final activation value depends on the activations of the entire batch of samples. Specifically, we use the relative quantile of the pre-activation¹ with respect to other pre-activations in the batch. (Details in section 2). Figure 1b shows that this simple change can allow the neural network to learn in spite of the contradictory labels.

Distribution Shift as a failure mode: The inability of current network architectures to handle distribution shift is a manifestation of the failure mode discussed above. For instance consider the above distribution (figure 1a), where, if the pair of Gaussian distributions are rotated by 30° (a.k.a distribution shift), the class 0 and class 1 overlap. We validate the proposed quantile activation on the standard distribution shift dataset - CIFAR10C. Figure 2 illustrates the results obtained using ReLU activations and QACT. As severity increases (w.r.t Gaussian Noise), we observe that ReLU activation loses the class structure. On the other hand, the proposed QACT framework does not suffer from this and the class structure is preserved with QACT.

Remark 1: Although QACT deals with activation of neurons, it is fundamentally different compared to the usual activation functions such as ReLU/Gelu etc. QACT is designed for classification tasks

¹We use the following convention – “Pre-activations” denote the inputs to the activation functions and “Activations” denote the outputs of the activation function.

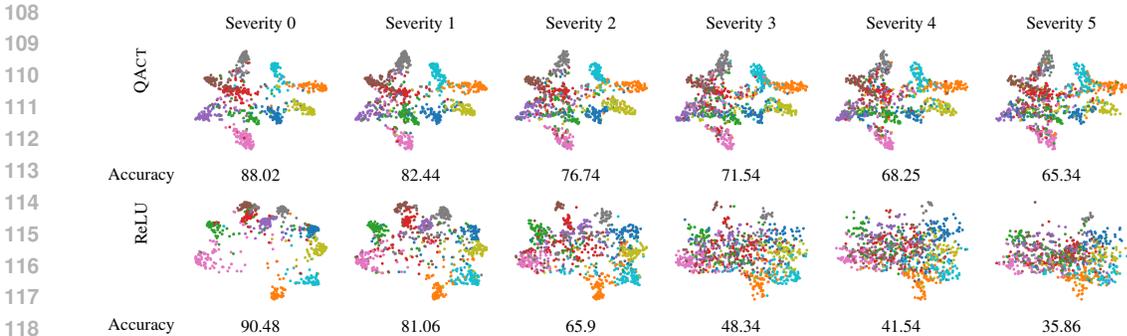


Figure 2: Comparing TSNE plots of QACT and ReLU activation on CIFAR10C with Gaussian distortions. Observe that QACT maintains the class structure extremely well across distortions, while the usual ReLU activations loses the class structure as severity increases.

where context plays an important role and hence considering “pure” classification tasks (i.e where the underlying distribution and ground-truth are unknown but fixed) is out of scope in this article.

Remark 2: Quantile activation is different from existing quantile neural network based approaches, such as regression (Prashanth et al., 2022), binary quantile classification (Tambwekar et al., 2022), Anomaly Detection (Li & van Leeuwen, 2023; Seo et al., 2022). Our approach is achieving best in-class performance by incorporating context distribution in the classification paradigm. Our approach is also markedly different from Machine unlearning which is based on selective forgetting of certain data points or retraining from scratch (Seetha et al., 2024).

Contributions: A decent amount of literature on neuronal activation is available. However, to the best of our knowledge, none matches the central idea proposed in this work.

In (Challa et al., 2023), the authors propose an approach to calibrate a pre-trained classifier $f_{\theta}(\mathbf{x})$ by extending it to learn a *quantile function*, $Q(\mathbf{x}, \theta, \tau)$ (τ denotes the quantile), and then estimate the probabilities using $\int_{\tau} I[Q(\mathbf{x}, \theta, \tau) \geq 0.5] d\tau^2$. They show that this results in probabilities which are robust to distortions.

1. In this article, we extend this approach to the level of a neuron, by suitably deriving the forward and backward propagation equations required for learning (section 2).
2. We then show that a suitable incorporation of our extension produces context dependent outputs at the level of each neuron of the neural network.
3. Our approach contributes to achieving better generalization across distributions and is more robust to distortions, across architectures. We evaluate our method using different architectures and datasets, and compare with the current state-of-the-art – DINOv2. We show that QACT proposed here is more robust to distortions than DINOv2, even if we have considerably less number of parameters (22M for DINOv2 vs 11M for Resnet18). Additionally, DINOv2 is trained on 20 odd datasets, before being applied on CIFAR10C; in contrast, our framework is trained on CIFAR10, and produces more robust outcome (see figures 4,6).
4. We also adapt QACT to design a classifier which returns better calibrated probabilities. We show that, unlike the relevant, most recent baselines (RELU, DINOv2 (small)), QACT achieves constant calibration error across different severity of distortions.

Related Works: This work aims to address the failure mode described earlier. To the best of our knowledge, no existing literature specifically addresses this issue. However, we believe this issue to be widely persistent in several practical domains. We use distribution shift to validate the proposed approach, which is part of a bigger problem of domain generalization.

Related Works on Domain Generalization (DG): The problem of domain generalization tries to answer the question – Can we use a classifier trained on one domain across several other related

² $I[\cdot]$ denotes the indicator function

domains? The earliest known approach for this is *Transfer Learning* (Pan & Yang, 2010; Zhuang et al., 2021), where a classifier from a single domain is applied to a different domain with/without fine-tuning. Several approaches have been proposed to achieve DG, such as extracting domain-invariant features over single/multiple source domains (Ghifary et al., 2015; Akuzawa et al., 2019; Dou et al., 2019; Piratla et al., 2020; Hu et al., 2019), Meta Learning (Huang et al., 2020; Dou et al., 2019), Invariant Risk Minimization (Arjovsky et al., 2019). Self-supervised learning is another proposed approach which tries to extract features on large scale datasets in an unsupervised manner, the most recent among them being DINOv2 (Oquab et al., 2023) which is the current state-of-the-art³. Very large foundation models, such as GPT-4V, are also known to perform better with respect to distribution shifts (Han et al., 2023). Nevertheless, to the best of our knowledge, none of these models incorporates context distributions for classification.

2 QUANTILE ACTIVATION

Rethinking Outputs from a Neuron: To recall – if \mathbf{x} denotes the input, a typical neuron does the following – (i) Applies a linear transformation with parameters w, b , giving $w^t\mathbf{x} + b$ as the output, and (ii) applies a rectifier g , returning $g(w^t\mathbf{x} + b)$. Typically, g is taken to be the ReLU activation - $g_{relu}(x) = \max(0, x)$. Intuitively, we expect that each neuron captures an “abstract” feature, usually not understood by a human observer.

An alternate way to model a neuron is to consider it as predicting a latent variable y , where $y = 1$ if the feature is present and $y = 0$ if the feature is absent. Mathematically, we have the following model:

$$z = w^t\mathbf{x} + b + \epsilon \quad \text{and} \quad y = I[z \geq 0] \quad (2)$$

This is very similar to the standard latent variable model for logistic regression, with the main exception being, the *outputs* y are *not known* for each neuron beforehand. If y is known, it is rather easy to obtain the probabilities – $P(z \geq 0)$. Can we still predict the probabilities, even when y itself is a latent variable?

The authors in (Challa et al., 2023) propose the following algorithm to estimate the probabilities:

1. Let $\{\mathbf{x}_i\}$ denote the set of input samples from the input distribution \mathbf{x} and $\{z_i\}$ denote their corresponding latent outputs, which would be from the distribution z
2. Assign $y = 1$ whenever $z > (1 - \tau)^{th}$ quantile of z , and 0 otherwise. For a specific sample, we have $y_i = 1$ if $z_i > (1 - \tau)^{th}$ quantile of $\{z_i\}$
3. Fit the model $Q(x, \tau; \theta)$ to the dataset $\{((\mathbf{x}_i, \tau), y_i)\}$, and estimate the probability as,

$$P(y_i = 1) = \int_{\tau=0}^1 I[Q(x, \tau; \theta) \geq 0.5] d\tau \quad (3)$$

The key idea: Observe that in step 2., the labelling is done without resorting to actual ground-truth labels. This allows us to obtain the probabilities on the fly for any set of parameters, only by considering the quantiles of z .

Defining the Quantile Activation QACT Let z denote the pre-activation of the neuron, and let $\{z_i\}$ denote the samples from this distribution. Let F_z denote the cumulative distribution function (CDF), and let f_z denote the density of the distribution. Accordingly, we have that $F_z^{-1}(\tau)$ denotes the τ^{th} quantile of z . Using step (2) of the algorithm above, we define,

$$\text{QACT}(z) = \int_{\tau=0}^1 I[z > F_z^{-1}(1 - \tau)] d\tau \stackrel{\text{Substitute}}{\tau \rightarrow (1-\tau)} \int_{\tau=0}^1 I[z > F_z^{-1}(\tau)] d\tau \quad (4)$$

Computing the gradient of QACT: However, to use QACT in a neural network, we need to compute the gradient which is required for back-propagation. Let τ_z denote the quantile at which $F_z^{-1}(\tau_z) = z$. Then we have that $\text{QACT}(z) = \tau_z$ since $F_z^{-1}(\tau)$ is an increasing function. So, we

³as per <https://paperswithcode.com/sota/domain-generalization-on-imagenet-c> accessed on 26 September 2024

Algorithm 1 Forward Propagation for a single neuron

Input: $[z_i]$ a vector of pre-activations, $0 < \tau_1 < \tau_2 < \dots < \tau_{n_\tau} < 1$ - a list of quantile indices at which we compute the quantiles.

Append two large values, c and $-c$, to the vector $[z_i]$.

Count n_+ = number of positive values, n_- = number of negative values, and assign the weight $w_+ = 1/n_+$ to the positive values, and $w_- = 1/n_-$ to the negative values.

Compute *weighted* quantiles $\{q_i\}$ at each of $\{\tau_i\}$ over the set $\{z_i\} \cup \{c, -c\}$

Compute $\text{QACT}(z_i)$ using the function,

$$\text{QACT}(x) = \frac{1}{n_\tau} \sum_i I[x \geq q_i] \quad (6)$$

Remember $[z_i], w_+, w_-, [\text{QACT}(z_i)]$ for backward propagation.

return $[\text{QACT}(z_i)]$

Algorithm 2 Backward Propagation for a single neuron

Input: grad_output , $0 < \tau_1 < \tau_2 < \dots < \tau_{n_\tau} < 1$ - a list of quantile indices at which we compute the quantiles.

Context from Forward Propagation: $[z_i], w_+, w_-, [\text{QACT}(z_i)]$

Obtain a weighted sample from $[z_i]$ with weights w_+, w_- - (say) S .

Obtain a kernel density estimate, using points from S , at each of the points in z_i - (say) $\hat{f}_z(z_i)$

Set,

$$\text{grad_input} = \text{grad_output} \odot [\hat{f}_z(z_i)] \quad (7)$$

return grad_input

have that $\text{QACT}(F_z^{-1}(\tau)) = \tau$. In other words, we have that $\text{QACT}(z)$ is $F_z(z)$, which is nothing but the CDF of z . Hence, we have,

$$\frac{\partial \text{QACT}(z)}{\partial z} = f_z(z) \quad (5)$$

where $f_z(z)$ denotes the density of the distribution.

Grounding the Neurons: With the above formulation, observe that since QACT is identical to CDF, it follows that, $\text{QACT}(z)$ is always a uniform distribution between 0 and 1, irrespective of the distribution z . When training numerous neurons in a layer, this could cause all the neurons to learn the same behaviour. Specifically, if, half the time, a particular abstract feature is more prevalent than others, QACT (as presented above) would not be able to learn this feature. To correct this, we *enforce that positive values and negative values have equal weight*. Given the input distribution z , We perform the following transformation before applying QACT. Let

$$z^+ = \begin{cases} z & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad z^- = \begin{cases} z & \text{if } z < 0 \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

denote the truncated distributions. Then,

$$z^\ddagger = \begin{cases} z^+ & \text{with probability 0.5} \\ z^- & \text{with probability 0.5} \end{cases} \quad (9)$$

From definition of z^\ddagger , we get that the median of z^\ddagger is 0. This grounds the input distribution to have the same positive and negative weight.

Dealing with corner cases: It is possible that during training, some neurons either only get positive values or only get negative values. However, for smooth outputs, one should still only give the weight of 0.5 for positive values. To handle this, we include two values c (large positive) and $-c$ (large negative) for each neuron. Since, the quantiles are conventionally computed using linear interpolation, this allows the outputs to vary smoothly. We take $c = 100$ in this article.

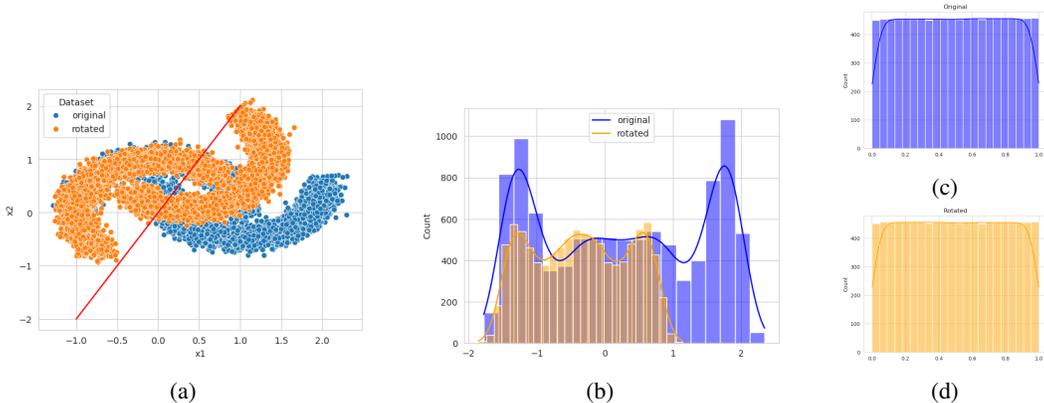


Figure 3: Intuition behind quantile activation. (a) shows a simple toy distribution of points (blue), its distortion (orange) and a simple line (red) on which the samples are projected to obtain activations. (b) shows the distribution of the pre-activations. (c) shows the distributions of the activations with QACT of the original distribution (blue). (d) shows the distributions of the activations with QACT under the distorted distribution (orange). Observe that the distributions match perfectly under small distortions. Note that even if the distribution matches perfectly, the quantile activation is actually a deterministic function.

Estimating the Density for Back-Propagation: Note that the gradient for the back propagation is given by the density of z^\ddagger (weighted distribution). We use the *Kernel Density Estimation* (KDE), to estimate the density. We, (i) First sample S points with weights w_+, w_- , and (ii) then estimate the density at all the input points $[z_i]$. This is point-wise multiplied with the backward gradient to get the gradient for the input. In this article we use $S = 1000$, which we observe gets reasonable estimates.

Computational Complexity: Computational Complexity (for a single neuron) is majorly decided by 2 functions – (i) Computing the quantiles has the complexity for a vector $[z_i]$ of size n can be performed in $\mathcal{O}(n \log(n))$. Since this is log-linear in n , it does not increase the complexity drastically compared to other operations in a deep neural network. (ii) Computational complexity of the KDE estimates is $\mathcal{O}(S n_\tau)$ where S is the size of sample (weighted sample from $[z_i]$) and n_τ is the number of quantiles, giving a total of $\mathcal{O}(n \log(n) + S n_\tau)$. In practice, we consider $S = 1000$ and $n_\tau = 100$ which works well, and hence does not increase with the batch size.

Remark: Algorithms 1, and 2 provide the pseudocode for the quantile activation. For stable training, in practice, we prepend and append the quantile activation with BatchNorm layers.

Why QACT is robust to distortions? To understand the idea behind quantile activation, consider a simple toy example in figure 3. For ease of visualization, assume that the input features (blue) are in 2 dimensions, and also assume that the line of the linear projection is given by the red line in figure 3a. Now, assume that the blue input features are rotated, leading to a different distribution (indicated here by orange). Since activations are essentially (unnormalized) signed distances from the line, we plot the histograms corresponding to the two distributions in figure 3b. As expected, these distributions are different. However, after performing the quantile activation in equation 4, we have that both are uniform distribution. This is illustrated in figures 3c and 3d. This behaviour has a normalizing effect across different distributions, and hence has better distribution generalization than other activations.

3 TRAINING WITH QACT

In the previous section, we described the procedure to adapt a single neuron to its context distribution. In this section we discuss how this extends to the Dense/Convolution layers, the loss functions to train the network and the inference aspect.

Extending to standard layers: The extension of equation 4 to dense outputs is straightforward. A typical output of the dense layer would be of the shape (B, N_c) - B denotes the batch size, N_c denotes the width of the network. The principle is - *The context distribution of a neuron is all the values which are obtained using the same parameters.* In this case, each of the values across the ‘ B ’ dimension are considered to be samples from the context distribution.

For a convolution layer, the typical outputs are of the form - (B, N_c, H, W) - B denotes the size of the batch, N_c denotes the number of channels, H, W denotes the sizes of the images. In this case we should consider all values across the 1st, 3rd and 4th dimension to be from the context distribution, since all these values are obtained using the same parameters. So, the number of samples would be $B \times H \times W$.

Loss Functions: One can use any differentiable loss function to train with quantile activation. We specifically experiment with the standard Cross-Entropy Loss, Triplet Loss, and the recently proposed Watershed Loss in (Challa et al., 2024) (see section 4). However, if one requires that the boundaries between classes adapt to the distribution, then learning similarities instead of boundaries can be beneficial. Both Triplet Loss and Watershed Loss fall into this category. We see that learning similarities does have slight benefits when considering the embedding quality.

Inference with QACT: As stated before, we want to assign a label for classification based on the context of the sample. There exist two approaches for this - (1) One way is to keep track of the quantiles and the estimated densities for all neurons and use it for inference. This allows inference for a single sample in the traditional sense. However, this also implies that one would not be able to assign classes based on the context at evaluation. (2) Another way is to make sure that, even for inference on a single sample, we include several samples from the context distribution, but only use the output for a specific sample. This allows one to assign classes based on the context. In this article, we follow the latter approach.

Quantile Classifier: Observe that the proposed QACT (without normalization) returns the values in $[0, 1]$ which can be interpreted as probabilities. Hence, one can also use this for the classification layer. Nonetheless, two changes are required - (i) Traditional softmax used in conjunction with negative-log-likelihood loss already considers “relative” activations of the classification in normalization. However, QACT does not. Hence, one should use Binary-Cross-Entropy loss with QACT, which amounts to one-vs-rest classification. (ii) Also, unlike a neuron in the middle layers, the bias of the neuron in the classification layer depends on the class imbalance. For instance, with 10 classes, one would have only 1/10 of the samples labelled 1 and 9/10 of the samples labelled 0. To address this, we require that the median of the outputs be at 0.9, and hence weight the positive class with 0.9 and the negative class with 0.1 respectively. In this article, whenever QACT is used, we use this approach for inference.

We observe that (figures 14 and 15) using quantile classifier on the learned features in general improves the consistency of the calibration error and also leads to the reducing the calibration error. In this article, for all networks trained with quantile activation, we use quantile classifier to compute the accuracy/calibration errors.

4 EVALUATION

To summarize, we make the following changes to the existing classification pipeline - (i) Replace the usual ReLU activation with QACT and (ii) Use triplet or watershed loss instead of standard cross-entropy loss. We expect this framework to learn context dependent features, and hence be robust to distortions. (iii) Also, use quantile classifier to train the classifier on the embedding for better calibrated probabilities.

Evaluation Protocol: To evaluate our approach, we consider the datasets developed for this purpose - CIFAR10C, CIFAR100C, TinyImagenetC (Hendrycks & Dietterich, 2019), MNISTC (Mu & Gilmer, 2019). These datasets have a set of 15 distortions at 5 severity levels. To ensure diversity we evaluate our method on 4 architectures - (overparametrized) LeNet, ResNet18 (He et al., 2016) (11M parameters), VGG (Simonyan & Zisserman, 2015)(15M parameters) and DenseNet

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

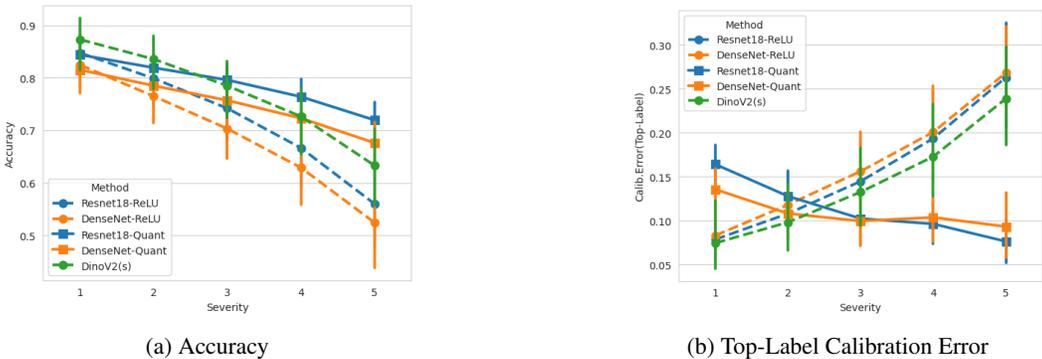


Figure 4: Comparing QACT with ReLU activation and DINOv2 (small) on CIFAR10C. We observe that, while at low severity of distortions QACT has a similar accuracy as existing pipelines, at higher levels the drop in accuracy is substantially smaller than existing approaches. With respect to calibration, we observe that the calibration error remains constant (up to standard deviations) across distortions.

(Huang et al., 2017) (1M parameters). The code to reproduce the results can be found at <https://anonymous.4open.science/r/QuantAct-534C>.

Baselines for Comparison: To our knowledge, there exists no other framework which proposes classification based on context distribution. So, for comparison, we consider standard ReLU activation (Fukushima, 1970), pReLU (He et al., 2015), and SELU (Klambauer et al., 2017) for all the architectures stated above. Also, we compare our results with DINOv2 (small) (Oquab et al., 2023) (22M parameters) which is current state-of-the-art for domain generalization. Note that for DINOv2, architecture and datasets used for training are substantially different (and substantially larger) from what we consider in this article. Nevertheless, we include the results for understanding where our proposed approach lies on the spectrum. We consider the small version of DINOv2 to match the number of parameters with the compared models.

Metrics: We consider four metrics – Accuracy (ACC), calibration error (ECE) (Kumar et al., 2019) (both marginal and Top-Label) and mean average precision at K (MAP@K) to evaluate the embedding. For the case of ReLU/pReLU/SELU activation with Cross-Entropy, we use the logistic regression trained on the train set embeddings, and for QACT we use the calibrated linear classifier, as proposed above. We do not perform any additional calibration and use the probabilities. We discuss a selected set of results in the main article. Please see appendix C for more comprehensive results.

Calibration error measures the reliability of predicted probabilities. In simple words, if one predicts 100 samples with (say) probability 0.7, then we expect 70 of the samples to belong to class 1 and the rest to class 0. This is measured using either the marginal or top-label calibration error. We refer the reader to (Kumar et al., 2019) for details, which also provides an implementation to estimate the calibration error.

Remark: For all the baselines we use the standard Cross-Entropy loss for training. For inference on corrupted datasets, we retrain the last layer with logistic regression on the train embedding and evaluate it on test/corrupted embedding. For QACT, we as a convention use watershed loss unless otherwise stated, for training. For inference, we train the Quantile Classifier on the train embedding and evaluate it on test/corrupted embedding.

The proposed QACT approach is robust to distortions: In fig. 4 we compare the proposed QACT approach with predominant existing pipeline – ReLU+Cross-Entropy and DINOv2(small) on CIFAR10C. In figure 4a we see that as the severity of the distortion increases, the accuracy of ReLU and DINOv2 drops significantly. On the other hand, while at small distortions the results are comparable, as severity increases QACT performs substantially better than conventional approaches. At severity 5, QACT outperforms DINOv2. On the other hand, we observe that in figure 4b, the calibration error stays consistent across distortions.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

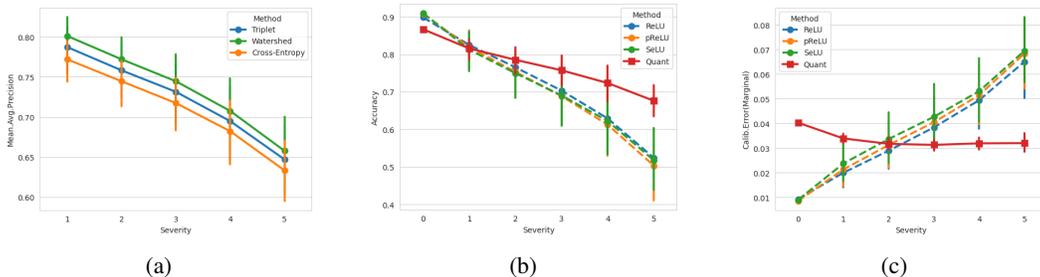


Figure 5: (a) Dependence on Loss functions. Here we compare watershed with other popular loss functions – Triplet and Cross-Entropy when used with QACT. We see that watershed performs slightly better with respect to MAP. (b) Comparing QACT with other popular activations – ReLU/pReLU/SELU with respect to accuracy. (c) Comparing QACT with other popular activations – ReLU/pReLU/SELU with respect to Calibration Error (Marginal). From both (b) and (c) we can conclude that QACT is notably more robust across distortions than several of the existing activation. All the plots use ResNet18 with CIFAR10C dataset.

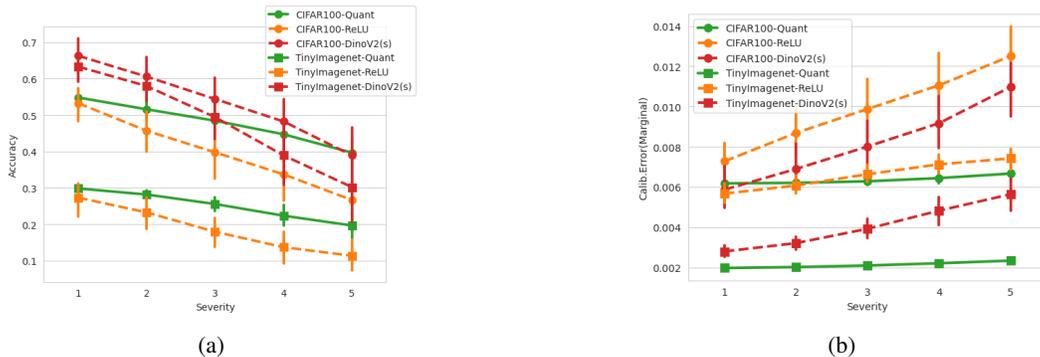


Figure 6: Results on CIFAR100C/TinyImagenetC. We compare QACT+watershed to ReLU and DinoV2 small on CIFAR100C/TinyImagenetC dataset with ResNet18. Note that the observations are consistent with CIFAR10C. (a) shows how accuracy changes across distortions. Observe that QACT is similar to DINOv2(s) with respect to embedding quality across all distortions, even if DINOv2 has 22M parameters as compared to Resnet18 11M parameters and is trained on larger datasets. (b) shows how calibration error (marginal) changes across severities. While other approaches lead to an increase in calibration error, QACT has similar calibration error across distortions.

How much does QACT depend on the loss function? Figure 5a compares the watershed classifier with other popular losses – Triplet and Cross-Entropy. We see that all the loss functions perform comparably when used in conjunction with QACT. We observe that watershed has a slight improvement when considering MAP and hence, we consider that as the default setting. However, we point out that QACT is compatible with several loss functions as well.

QACT vs ReLU/pReLU/SELU activations: To verify that most existing activations do not share the robustness property of QACT, we compare QACT with other activations in figures 5b and 5c. We observe that QACT is greatly more robust with respect to distortions in both accuracy and calibration error than other activation functions.

Results on Larger Datasets: To verify that our observations hold for larger datasets, we use CIFAR100C/TinyImagenetC to compare the proposed QACT+watershed with existing approaches. We observe on figure 6 that QACT performs comparably well as DINOv2, although DINOv2(s) has 22M parameters and is trained on significantly larger datasets. Moreover, we also observe that QACT has approximately constant calibration error across distortions, as opposed to a significantly increasing calibration error for ReLU or DINOv2.

5 CONCLUSION AND FUTURE WORK

To summarize, traditional classification systems do not consider the “context distributions” when assigning labels. In this article, we propose a framework to achieve this by – (i) Making the activation adaptive by using quantiles and (ii) Learning a kernel instead of the boundary for the last layer. We show that our method is more robust to distortions by considering MNISTC, CIFAR10C, CIFAR100C, TinyImagenetC datasets across varying architectures.

The scope of this article is to provide a proof of concept and a framework for performing inference in a context-dependent manner. We outline several potential directions for future research:

- I. The key idea in our proposed approach is that the quantiles capture the distribution of each neuron from the batch of samples, providing outputs accordingly. This poses a challenge for large datasets, and we have discussed two potential solutions: (i) remember the quantiles and density estimates for single sample evaluation, or (ii) ensure that a batch of samples from the same distribution is processed together. We adopt the latter method in this article. An alternative approach would be to *learn the distribution of each neuron* using auxiliary loss functions, adjusting these distributions to fit the domain at test time. This gives us more control over the network at test time compared to current workflows. If the networks are very large, where batch sizes cannot be big – there exists several strategies such as checkpointing to implicitly increase the batch size.
- II. Since the aim of the article was to establish a proof-of-concept, we did not focus on scaling, and use only a single GPU for all the experiments. To extend it to multi-GPU training, one needs to synchronize the quantiles across GPU, in a similar manner as that for Batch-Normalization. We expect this to improve the statistics, and to allow considerably larger batches of training.
- III. On the theoretical side, there is an interesting analogy between our quantile activation and how a biological neuron behaves. It is known that when the inputs to a biological neuron change, the neuron adapts to these changes (Clifford et al., 2007). Quantile activation does something very similar, which leads to an open question – can we establish a formal link between the adaptability of a biological neuron and the accuracy of classification systems?
- IV. Another theoretical direction to explore involves considering distributions not just at the neuron level, but at the layer level, introducing a high-dimensional aspect to the problem. The main challenge here is defining and utilizing *high dimensional quantiles*, which remains an open question (Koenker, 2005).

REFERENCES

- Kei Akuzawa, Yusuke Iwasawa, and Yutaka Matsuo. Adversarial invariant feature learning with accuracy constraint for domain generalization. In *European Conf. Mach. Learning*, 2019.
- Martín Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv:1907.02893*, 2019.
- Aditya Challa, Snehanu Saha, and Soma Dhavala. Quantprob: Generalizing probabilities along with predictions for a pre-trained classifier. *arXiv:2304.12766*, 2023.
- Aditya Challa, Sravan Danda, and Laurent Najman. A novel approach to regularising 1nn classifier for improved generalization. *arXiv:2402.08405*, 2024.
- Colin WG Clifford, Michael A Webster, Garrett B Stanley, Alan A Stocker, Adam Kohn, Tatyana O Sharpee, and Odelia Schwartz. Visual adaptation: Neural, psychological and computational aspects. *Vision research*, 2007.
- Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization via model-agnostic learning of semantic features. In *Neural Inform. Process. Syst.*, 2019.
- Kunihiko Fukushima. Correction to "visual feature extraction by a multilayered network of analog threshold elements". *IEEE Trans. Syst. Sci. Cybern.*, 1970.
- Muhammad Ghifary, W. Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *Proc. Int. Conf. Comput. Vision*, 2015.

- 540 Zhongyi Han, Guanglin Zhou, Rundong He, Jindong Wang, Tailin Wu, Yilong Yin, Salman H. Khan,
541 Lina Yao, Tongliang Liu, and Kun Zhang. How well does gpt-4v(ision) adapt to distribution shifts?
542 A preliminary investigation. *arXiv:2312.07424*, 2023.
- 543 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
544 human-level performance on imagenet classification. In *Proc. Int. Conf. Comput. Vision*, 2015.
- 545 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
546 recognition. In *Proc. Conf. Comput. Vision Pattern Recognition*, 2016.
- 547 Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common
548 corruptions and perturbations. In *Int. Conf. on Learning Representations*, 2019.
- 549 Shoubo Hu, Kun Zhang, Zhitang Chen, and Laiwan Chan. Domain generalization via multidomain
550 discriminant analysis. In *Uncertainty in Artificial Intelligence*, 2019.
- 551 Bincheng Huang, Si Chen, Fan Zhou, Cheng Zhang, and Feng Zhang. Episodic training for domain
552 generalization using latent domains. In *Int. Conf. on Cogni. Systems and Signal Process.*, 2020.
- 553 Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected
554 convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition,
555 CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 2261–2269. IEEE Computer Society, 2017.
556 doi: 10.1109/CVPR.2017.243. URL <https://doi.org/10.1109/CVPR.2017.243>.
- 557 Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing
558 neural networks. In *Neural Inform. Process. Syst.*, 2017.
- 559 Roger Koenker. *Quantile Regression*. Econometric Society Monographs. Cambridge University
560 Press, 2005. doi: 10.1017/CBO9780511754098.
- 561 Ananya Kumar, Percy Liang, and Tengyu Ma. Verified uncertainty calibration. In *Neural Inform.
562 Process. Syst.*, 2019.
- 563 Zhong Li and Matthijs van Leeuwen. Explainable contextual anomaly detection using quantile
564 regression forests. *Data Min. Knowl. Discov.*, 2023.
- 565 Norman Mu and Justin Gilmer. MNIST-C: A robustness benchmark for computer vision.
566 *arXiv:1906.02337*, 2019.
- 567 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
568 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas
569 Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael G.
570 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jégou, Julien Mairal, Patrick Labatut,
571 Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision.
572 *arXiv:2304.07193*, 2023.
- 573 Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.*,
574 2010.
- 575 Vihari Piratla, Praneeth Netrapalli, and Sunita Sarawagi. Efficient domain generalization via common-
576 specific low-rank decomposition. In *Int. Conf. Mach. Learning*, 2020.
- 577 Tejas Prashanth, Snehanshu Saha, Sumedh Basarkod, Suraj Aralihalli, Soma S. Dhavala, Sriparna
578 Saha, and Raviprasad Aduri. Lipgene: Lipschitz continuity guided adaptive learning rates for fast
579 convergence on microarray expression data sets. *IEEE ACM Trans. Comput. Biol. Bioinform.*,
580 2022.
- 581 Aditi Seetha, Satyendra Singh Chouhan, Emmanuel S Pilli, Vaskar Raychoudhury, and Snehanshu
582 Saha. Dievd-sf: Disruptive event detection using continual machine learning with selective
583 forgetting. *IEEE Transactions on Computational Social Systems*, 2024.
- 584 Hogeon Seo, Seunghyoung Ryu, Jiyeon Yim, Junghoon Seo, and Yonggyun Yu. Quantile autoencoder
585 for anomaly detection. In *AAAI Workshop on AI for Design and Manufacturing (ADAM)*, 2022.

594 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
595 recognition. In *Int. Conf. on Learning Representations*, 2015.
596

597 Anuj Tambwekar, Anirudh Maiya, Soma S. Dhavala, and Snehanshu Saha. Estimation and applica-
598 tions of quantiles in deep binary classification. *IEEE Trans. Artif. Intell.*, 2022.

599 Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and
600 Qing He. A comprehensive survey on transfer learning. *Proc. IEEE*, 2021.
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

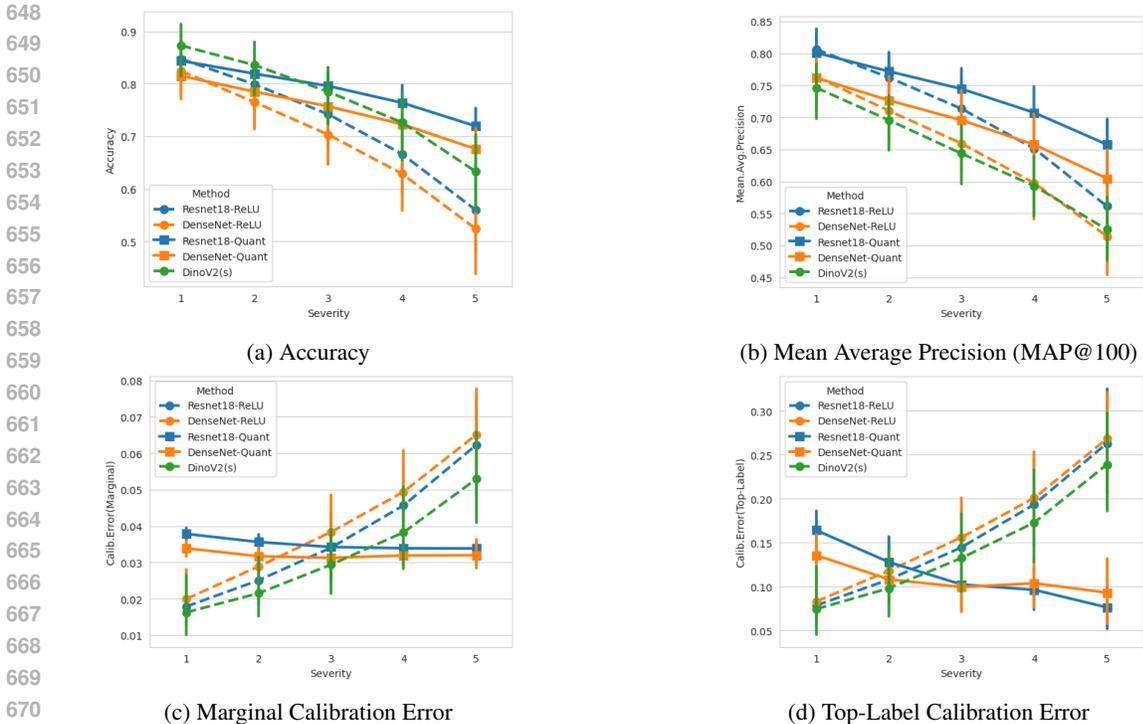


Figure 7: Comparing QACT with ReLU activation and DINOv2 (small).

A EXPERIMENT DETAILS FOR FIGURE 2

We consider the features obtained from ResNet18 with both QACT and ReLU activations for the datasets of CIFAR10C with gaussian_noise at all the severity levels. Hence, we have 6 datasets in total. To use TSNE for visualization, we consider 1000 samples from each dataset and obtain the combined TSNE visualizations. Each figure shows a scatter plot of the 2d visualization for the corresponding dataset.

B COMPUTE RESOURCES AND OTHER EXPERIMENTAL DETAILS

All experiments were performed on a single NVidia GPU with 32GB memory with Intel Xeon CPU (10 cores). For training, we perform an 80:20 split of the train dataset with seed 42 for reproducibility. All networks are initialized using default pytorch initialization technique.

We use Adam optimizer with initial learning rate $1e - 3$. We use ReduceLROnPlateau learning rate scheduler with parameters `factor=0.1`, `patience=50`, `cooldown=10`, `threshold=0.01`, `threshold_mode=abs`, `min_lr=1e-6`. We monitor the validation accuracy for learning rate scheduling. We also use `early_stopping` when the validation accuracy does not increase by 0.001.

C EXTENDED RESULTS SECTION

Comparing QACT + watershed and ReLU+Cross-Entropy: Figure 7 shows the corresponding results. The first experiment compares QACT + watershed with ReLU + Cross-Entropy on two standard networks – ResNet18 and DenseNet. With respect to accuracy, we observe that while at severity 0, ReLU + Cross-Entropy slightly outperforms QACT + watershed, as severity increases QACT + watershed is far more stable. We even outperform DINOv2(small) (22M parameters) at severity 5. Moreover, with respect to calibration error, we see a consistent trend across distortions. As (Challa et al., 2023) argues, this helps in building more robust systems compared to one where calibration error increases across distortions.

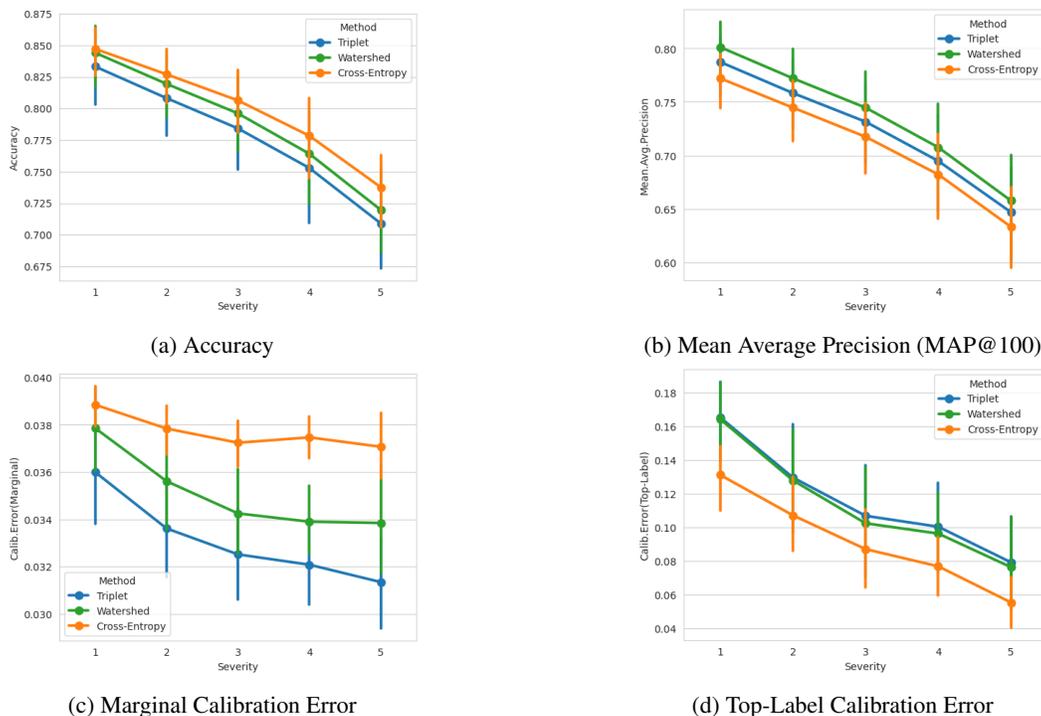


Figure 8: Triplet vs Watershed vs Cross-Entropy

Does loss function make a lot of difference? Figure 8 compares three different loss functions Watershed, Triplet and Cross-Entropy when used in conjunction with QACT. We observe similar trends across all loss functions. However, Watershed performs better with respect to Mean Average Precision (MAP) and hence we use this as a default strategy.

Why Mean-Average-Precision? – We argue that the key indicator of distortion invariance should be the quality of embedding. While, accuracy (as measured by a linear classifier) is a good metric, a better one would be to measure the Mean-Average-Precision. With respect to calibration error, due to the scale on the Y-axis, the figures suggest reducing calibration error. However, the standard deviations overlap, and hence, these are assumed to be constant across distortions.

How well does watershed perform when used with ReLU activation? Figure 9 shows the corresponding results. We observe that both the watershed loss and cross-entropy have large overlaps in the standard deviations at all severity levels. So, this shows that, when used in conjunction with ReLU watershed and cross-entropy loss are very similar. But in conjunction with QACT, we see that watershed has a slightly higher Mean-Average-Precision.

What if we consider an easy classification task? In figure 10, we perform the comparison of QACT+Watershed and ReLU and cross-entropy on MNISTC dataset. Across different architectures, we observe a lot less variation (standard deviation) of QACT+Watershed compared to ReLU and cross-entropy. This again suggests robustness against distortions of QACT+Watershed.

Comparing with other popular activations: Figures 11 and 12 shows the comparison of QACT with ReLU, pReLU and SeLU. We observe the same trend across ReLU, pReLU and SeLU, while QACT is far more stable across distortions.

Results on CIFAR100/TinyImagenetC: Figure 13 compares QACT+Watershed and ReLU+Cross-Entropy on CIFAR100C dataset. We also include the results of QACT+Cross-Entropy vs. ReLU+Cross-Entropy on TinyImagenetC. The results are consistent with what we observe on CIFAR10C, and hence, draw the same conclusions as before.

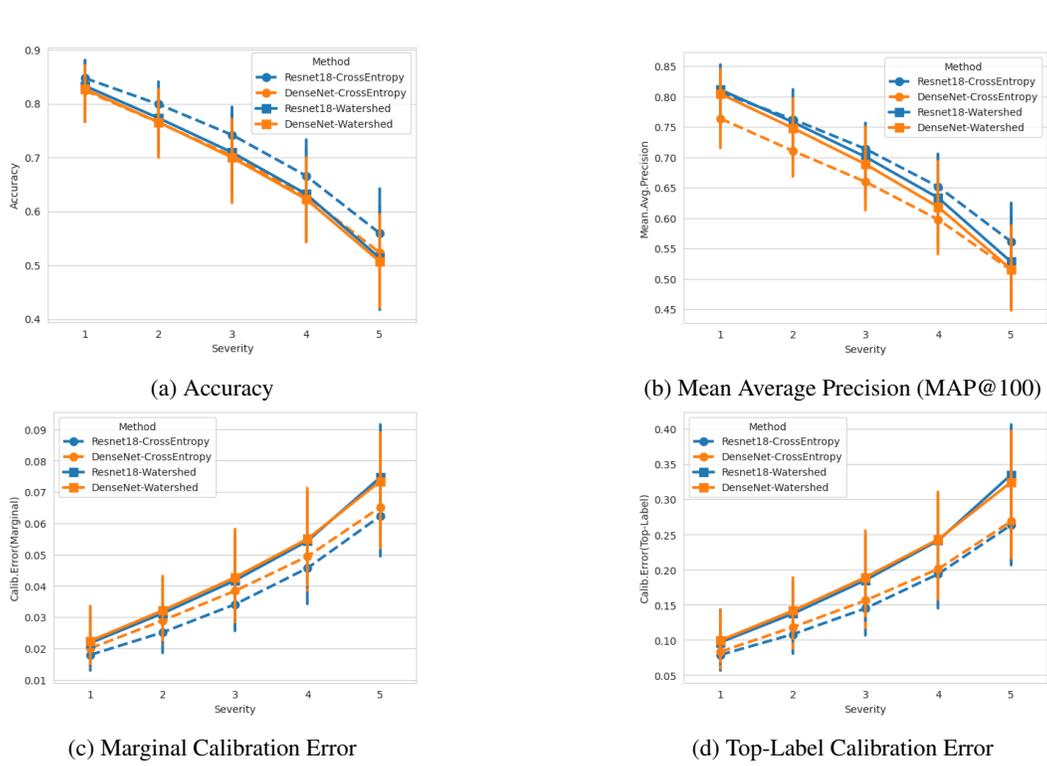


Figure 9: Watershed vs Cross-Entropy when using ReLU activation

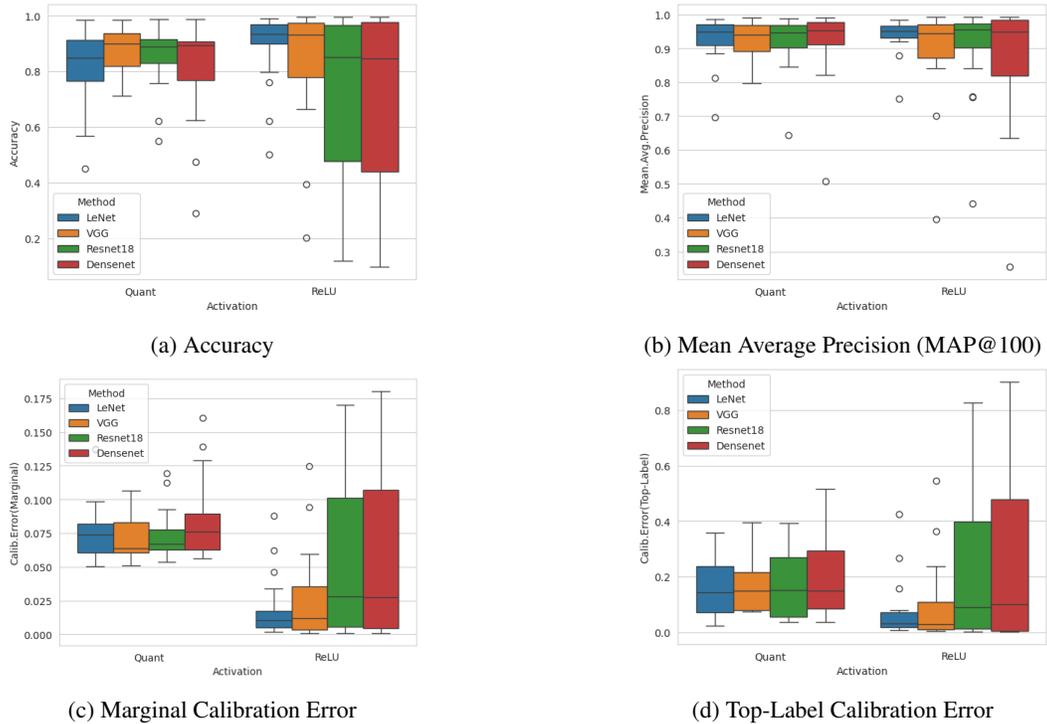


Figure 10: Results on MNIST

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

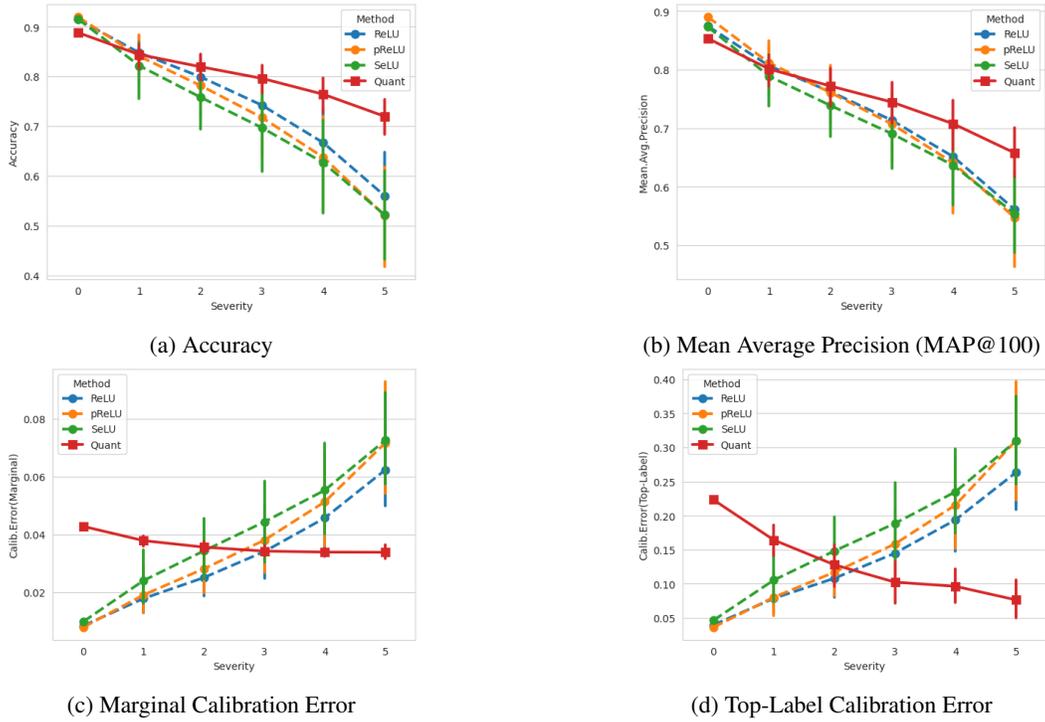


Figure 11: QACTvs ReLU vs pReLU vs Selu activations on ResNet18

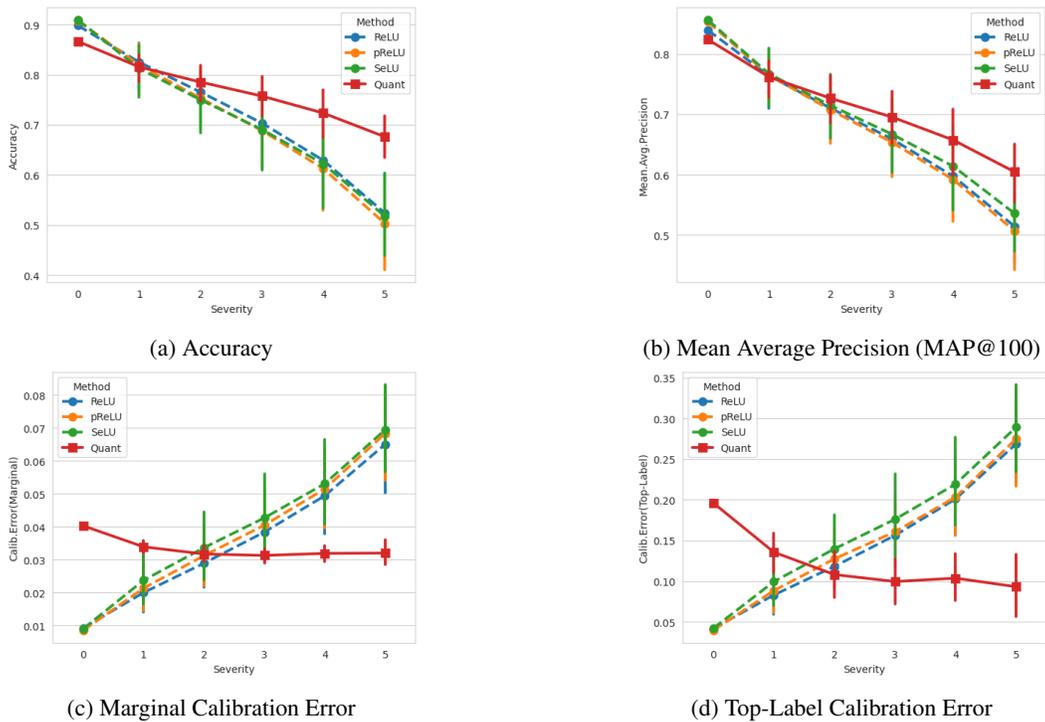


Figure 12: QACTvs ReLU vs pReLU vs Selu activations on Densenet

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

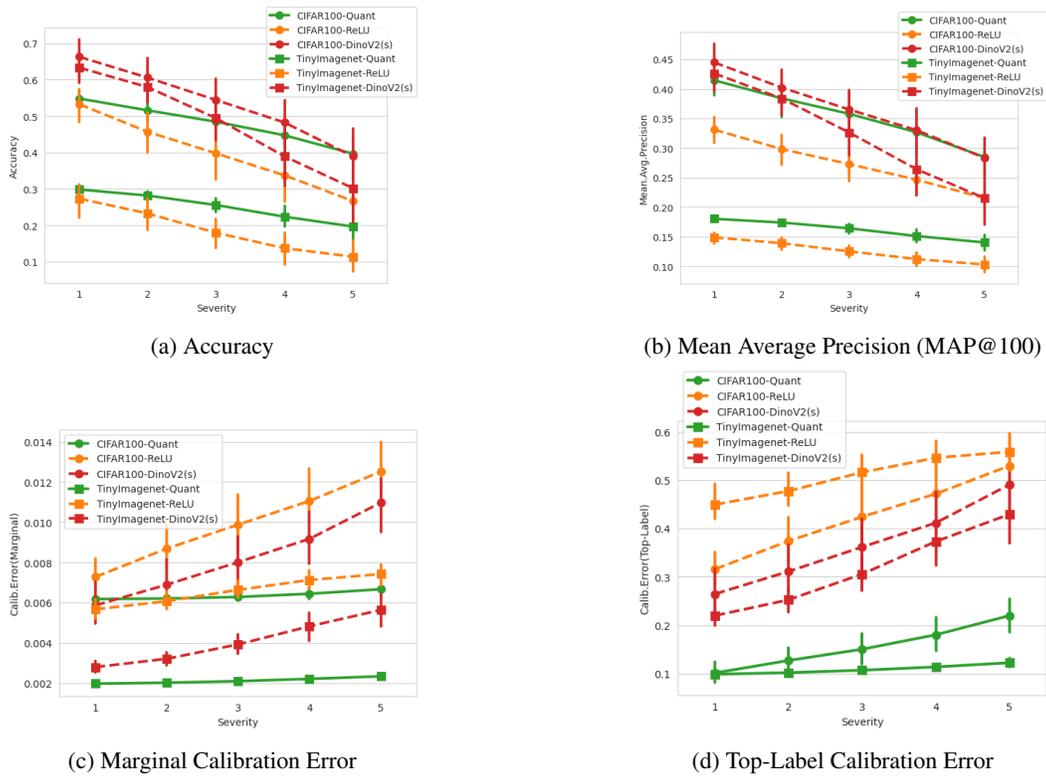


Figure 13: QACTvs ReLU on Resnet18+CIFAR100

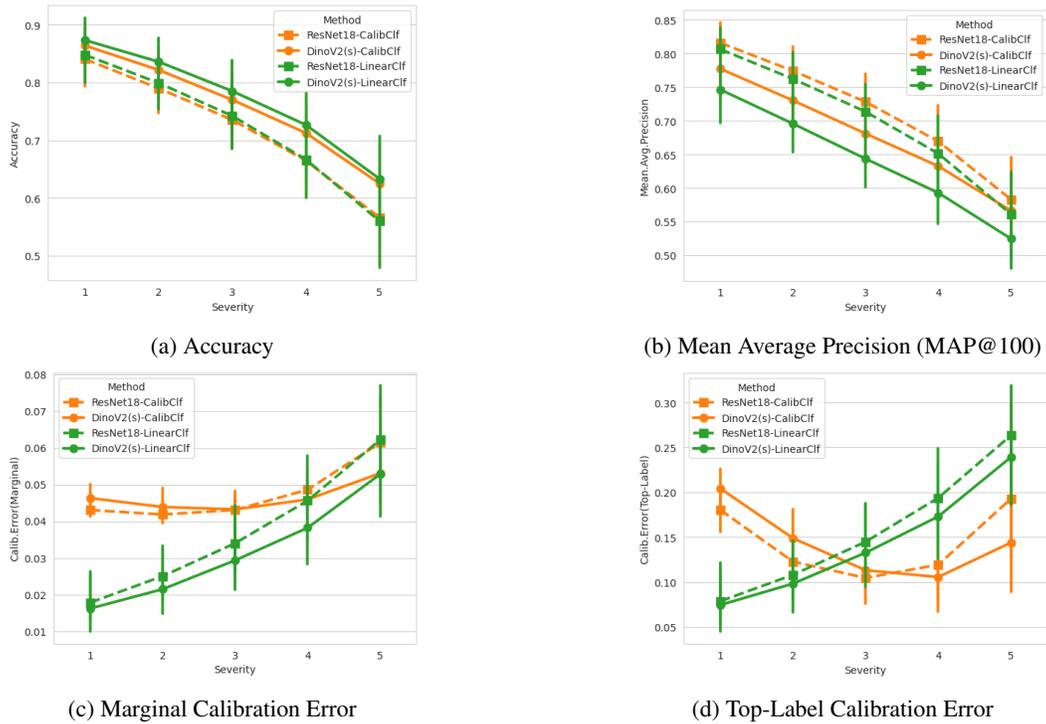


Figure 14: Effect of Quantile Classifier. We use ResNet18 and DinoV2 architectures on CIFAR10.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

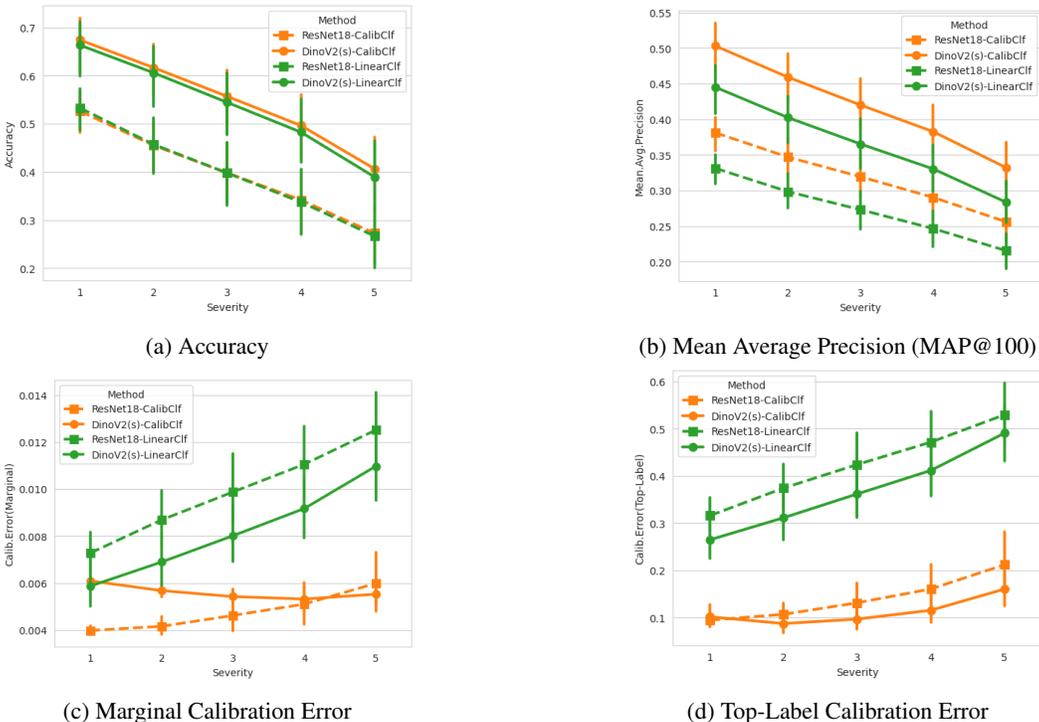


Figure 15: Effect of Quantile Classifier. We use ResNet18 and DinoV2 architectures on CIFAR100.

Dataset	Model/Method	0→1	0→2	0→3	0→4	0→5	1→2	1→3	1→4	1→5
CIFAR10	Resnet18-ReLU	6.70	11.57	17.27	24.82	35.49	4.87	10.57	18.12	28.79
	DinoV2(s)	6.92	10.66	15.74	21.58	30.93	3.74	8.82	14.67	24.02
	Resnet18-Quant	4.42	6.88	9.22	12.42	16.86	2.46	4.80	8.00	12.44
CIFAR100	Resnet18-ReLU	10.25	17.85	23.78	29.81	36.87	7.59	13.52	19.56	26.62
	DinoV2(s)	10.64	16.39	22.54	28.71	38.02	5.74	11.89	18.06	27.37
	Resnet18-Quant	5.83	9.07	12.20	15.93	21.00	3.24	6.38	10.10	15.17
TinyImagenet	Resnet18-ReLU	5.51	9.58	14.90	19.13	21.49	4.07	9.39	13.62	15.97
	DinoV2(s)	9.39	14.73	23.30	33.78	42.52	5.34	13.91	24.39	33.13
	Resnet18-Quant	3.27	4.86	7.26	10.43	13.18	1.59	3.99	7.16	9.90

Table 1: Measuring the drop in accuracy

Effect of Quantile Classifier: Figures 14 and 15 shows the effect of quantile classifier on standard ResNet10/DinoV2 outputs with CIFAR10C/CIFAR100C datasets. While the accuracy values are almost equivalent, we observe a “flatter” trend of the calibration errors, sometimes reducing the error as in the case of CIFAR100C.

Measuring Robustness of QACT: To measure the robustness of the proposed method we use the metric $Acc@Dist_i - Acc@Dist_j$ which measures the drop in accuracy when the distortion severity is increased from $i \rightarrow j$. A method is considered to be better if the values are *lower*, i.e if the drop in accuracy is smaller than comparative methods. Tables 1 and 2 shows the comparison between ReLU, QACT and DinoV2(s). We see that, in all the cases QACT outperforms both ReLU and DinoV2(s) at all possible $i \rightarrow j$.

Dataset	Model/Method	2→3	2→4	2→5	3→4	3→5	4→5
CIFAR10	Resnet18-ReLU	5.70	13.25	23.92	7.55	18.22	10.67
	DinoV2(s)	5.08	10.93	20.28	5.85	15.20	9.35
	Resnet18-Quant	2.34	5.54	9.98	3.20	7.64	4.45
CIFAR100	Resnet18-ReLU	5.93	11.97	19.03	6.04	13.10	7.06
	DinoV2(s)	6.15	12.32	21.63	6.17	15.48	9.31
	Resnet18-Quant	3.13	6.86	11.92	3.73	8.79	5.06
TinyImagenet	Resnet18-ReLU	5.32	9.55	11.90	4.23	6.58	2.35
	DinoV2(s)	8.57	19.05	27.79	10.48	19.22	8.74
	Resnet18-Quant	2.40	5.56	8.31	3.17	5.92	2.75

Table 2: Measuring the drop in accuracy (contd. from table 1)

D WATERSHED LOSS

The authors in (Challa et al., 2024) proposed a novel classifier – *watershed classifier*, which works by learning similarities instead of the boundaries. Below we give the brief idea of the loss function, and refer the reader to the original paper for further details.

- Let (\mathbf{x}_i, y_i) denote the samples in each batch, and let f_θ denote the embedding network. $f_\theta(\mathbf{x}_i)$ denotes the corresponding embedding.
- Starting from randomly selected seeds in the batch, propagate the labels to all the samples. Let \hat{y}_i denote the estimated samples. For each $f_\theta(\mathbf{x}_i)$ and for each label l , obtain the nearest neighbour in the samples in the set,

$$\mathcal{S}_l = \{f_\theta(\mathbf{x}_i) \mid \hat{y}_i = y_i = l\} \quad (10)$$

that is, all the samples of class l labelled correctly. Denote this nearest neighbour using $f_\theta(\mathbf{x}_{i,l,1nn})$.

- Then the loss is given by,

$$\text{Watershed Loss} = \frac{-1}{n_{\text{samples}}} \sum_{i=1}^{n_{\text{samples}}} \sum_{l=1}^L I[y_i = l] \log \left(\frac{\exp(-\|f_\theta(\mathbf{x}_i) - f_\theta(\mathbf{x}_{i,l,1nn})\|)}{\sum_{j=1}^L \exp(-\|f_\theta(\mathbf{x}_i) - f_\theta(\mathbf{x}_{i,j,1nn})\|)} \right) \quad (11)$$

Why Watershed Loss?: Observe that the loss in equation 11 implicitly learns representations consistent with the RBF kernel, which is known to be translation invariant. Minimizing this loss function, hence, will learn translation invariant kernels. This is important for obtaining networks robust to distortions.

If one uses (say) cross-entropy loss, then the features learned would be such that the classes are linearly separable. Contrast this with watershed, which instead learns a similarity between two points in a translation invariant manner.

Remark: Observe that the watershed loss is very similar to metric learning losses. The authors in (Challa et al., 2024) claim that this offers better generalization, and show that this is consistent with 1NN classifier. Moreover, they show that this classifier (without considering f_θ) has a VC dimension which is equal to the number of classes. While metric learning losses are similar, there is no such guarantee with respect to classification. This motivated our choice of using watershed loss over other metric learning losses.