# It's All In the Teacher: Zero-Shot Quantization Brought Closer to the Teacher

Kanghyun Choi[1], Hye Yoon Lee[1], Deokki Hong[1], Joonsang Yu[2],
Noseong Park[1], Youngsok Kim[1], and Jinho Lee[1*]

[1]College of Computing, Yonsei University    [2]CLOVA ImageVision, CLOVA AI Lab, NAVER

[1]{kanghyun.choi,hylee817,dk.hong,noseong,youngsok,leejinho}@yonsei.ac.kr

[2]joonsang.yu@navercorp.com

## Abstract

*Model quantization is considered as a promising method to greatly reduce the resource requirements of deep neural networks. To deal with the performance drop induced by quantization errors, a popular method is to use training data to fine-tune quantized networks. In real-world environments, however, such a method is frequently infeasible because training data is unavailable due to security, privacy, or confidentiality concerns. Zero-shot quantization addresses such problems, usually by taking information from the weights of a full-precision teacher network to compensate the performance drop of the quantized networks. In this paper, we first analyze the loss surface of state-of-the-art zero-shot quantization techniques and provide several findings. In contrast to usual knowledge distillation problems, zero-shot quantization often suffers from 1) the difficulty of optimizing multiple loss terms together, and 2) the poor generalization capability due to the use of synthetic samples. Furthermore, we observe that many weights fail to cross the rounding threshold during training the quantized networks even when it is necessary to do so for better performance. Based on the observations, we propose AIT, a simple yet powerful technique for zero-shot quantization, which addresses the aforementioned two problems in the following way: AIT i) uses a KL distance loss only without a cross-entropy loss, and ii) manipulates gradients to guarantee that a certain portion of weights are properly updated after crossing the rounding thresholds. Experiments show that AIT outperforms the performance of many existing methods by a great margin, taking over the overall state-of-the-art position in the field.*

## 1. Introduction

Deep neural network quantization [13,22,36,69] is a powerful tool for improving the computational efficiency of deep neural networks (DNNs). When being accompanied with the low-bitwidth hardware design [28,49,58], the latency and energy consumption of DNNs can be greatly reduced.

One problem of quantized models is, however, that they often suffer from the significant drop in accuracy, mainly due to quantization errors [36]. A popular way to address the problem is to further train or calibrate the model with training data [8,23,52,65,68,69]. During the fine-tuning procedure, the forward pass is performed with quantized values whereas the backpropagation is done with floating-point values to recover the accuracy loss in the initial quantization.

Unfortunately, such fine-tuning methods, which assume the full availability of training data at the time of quantization, are often not feasible in reality. Many models are disclosed to public only with their trained weights, and the dataset may contain proprietary, confidential, or sensitive data that fundamentally prohibit sharing.

Zero-shot quantization (or data-free quantization) [3,4,9, 10,44,61,66,67,70] is therefore a necessary technique for quantization. It assumes that only the architecture and the pre-trained weights are available at the time of quantization. Current successful approaches are mainly led by generative approaches [4,9,38,61,66,70]. Using synthetic samples from generators, knowledge distillation [21] is applied against full-precision models. It is known that the state-of-the-art methodology achieves almost similar performance to that of the data-driven approaches (i.e., quantization with real samples) for 5-bit fixed-point quantization, and comparable performance on 4-bit fixed-point setting. [9].

However, the recipe of the fine-tuning in zero-shot quantization is mainly adopted from common knowledge distillation problems [7,19,21] that consider neither quantization nor synthetic samples. As in the knowledge distillation, the loss function of the zero-shot quantization is habitually built as a combination of the cross-entropy (CE) against the hard label and the Kullback–Leibler (KL) divergence against the full-precision network's output.[1] It works well

---

[1]In the remainder of this paper, we refer to CE as the cross-entropy against the hard label and KL as the KL divergence against the full-precision network unless otherwise stated.

---

*Corresponding author

in practice, but there are no detailed studies to introspect the appropriateness of the loss in the context of zero-shot quantization. Therefore, more analyses on those solutions are needed. Moreover, the distribution of synthetic samples can be different from that of the original data. In such a case, they can be considered a type of adversarial samples (also see Fig. 6 for examples) and thus, the quantized network produces a huge generalization gap.

To our knowledge, we for the first time perform in-depth analyses on the loss surface of the zero-shot quantization problem. Through the analyses, we find several key observations for better quantization. First, quantized models often have difficulty optimizing multiple loss terms, and the loss terms fail to cooperate — in other words, the angle between the gradients of CE and KL is quite large in many cases. Second, KL usually has a much flatter loss surface than that of CE, having a better potential for generalizability.

To this end, we propose a method to address such problems of the zero-shot quantization, called *AIT (All In the Teacher)*. While pursuing a flatter surface of the loss curve, AIT lets the quantized student model get closer to the full-precision teacher model. To be more specific, we exclude CE from the loss, and apply our proposed <u>gradient inundation</u> with KL only. In addition, gradient inundation is designed to grow the gradients of KL in such a manner that a certain portion of weights are guaranteed to be updated in each layer. As a result, the quantized model approaches closer to the full-precision teacher, and our method takes over the state-of-the-art position for various datasets. Our contributions can be summarized as follows:

- We analyze the first and second-order loss surfaces, i.e., gradient and Hessian, of the zero-shot quantization problem. To the best of our knowledge, we are the first to closely investigate the loss function in the zero-shot quantization problem.

- We identify that the gradients from CE and KL form a large angle from the beginning to the end of the fine-tuning. This implies that the quantized network is suffering from their trade-off instead of benefiting from them working in harmony.

- We analyze the local curvature of the loss surface and observe that the two losses of our interest exhibit a great amount of curvature difference.

- We observe that the quantized student suffers from infrequent updates, where only a few layers are changing their integer weights and the remaining layers are stuck below rounding thresholds.

- Based on these findings, we propose AIT which excludes the cross-entropy loss, and manipulates the gradients using our proposed gradient inundation method

such that the quantized student model can faithfully resemble the full-precision teacher model.

- We perform a thorough evaluation of AIT. The results show that AIT outperforms the existing algorithms by a great margin, showing the state-of-the-art performance on the zero-shot quantization problem.

## 2. Background and Related Work

### 2.1. Quantization

Quantization of neural networks have been studied for a while, and there are numerous methods [5, 11, 16, 18, 29, 48]. In this work, we consider symmetric, uniform quantization that is known to be much easier to build hardware architectures for. With $n$ bits, a weight parameter $\theta$ is represented by one of $2^n$ ranges. We use $\theta^q$ to denote the quantized weights as outputs of a quantization function $Quant()$. For $Quant()$, we use a simple yet efficient function as the following [23]:

$$\theta^q = Quant(\theta) = \lfloor \theta \times S - z \rceil, \tag{1}$$

$$S = \frac{2^n - 1}{\theta_{max} - \theta_{min}}, \tag{2}$$

$$z = S \times \theta_{min} + 2^{n-1}, \tag{3}$$

where $S$ is the scaling factor to convert the range of $\theta$ to $n$ bit, and $z$ decides which quantized value zero is mapped to. After quantization, the quantized integer value represents $\theta' \in \mathbb{R}$ obtained by dequantization:

$$\theta' = \frac{(\theta^q + z)}{S}. \tag{4}$$

The procedure is the same for activation values, except that the minimum and the maximum are obtained from observing activations from a few batches and taking a moving average.

### 2.2. Zero-shot Quantization

Even though quantization has been shown to be effective even for extremely low bits [13, 39, 52, 69], they usually require training data to fine-tuning or calibration. Zero-shot quantization is a method to relax the privacy or confidentiality problem of the training data.

Earlier methods for zero-shot quantization were focused on how to build a good quantization function $Quant()$, by using schemes such as weight equalization, bias correction, or range adjustments [3, 44, 67]. Among them, ZeroQ [4] was the first work to introduce the notion of distilled data that is designed to match the batch-norm stats of the original full-precision network. With this scheme, choosing the adequate mixed-precision quantization for each layer has been proposed together. On top of ZeroQ, DSG [66] added diverse sample generation to improve the performance.
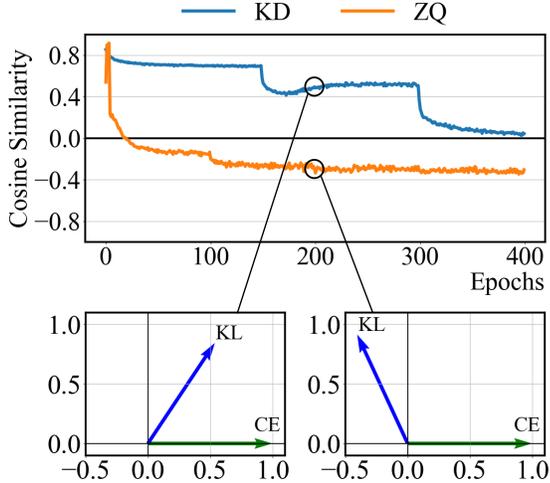
Figure 1. Plot of cosine similarity between cross-entropy on hard labels KL divergence on the full-precision model. At the bottom are snapshots of the relative angle between KL and CE of knowledge distillation (KD, bottom left) and zero-shot quantization (ZQ, bottom right) at epoch 200.

Later, GDFQ [61] adopted generative models [43, 47] to create better samples. The generator $G$ and the quantized model $Q$ are jointly trained with the following loss functions:

$$\mathcal{L}_{GDFQ}(G) = (1 - \alpha)\mathcal{L}_{CE}^{P}(G) + \alpha\mathcal{L}_{BNS}^{P}(G), \quad (5)$$

$$\mathcal{L}_{GDFQ}(Q) = (1 - \delta)\mathcal{L}_{CE}^{Q}(Q) + \delta\mathcal{L}_{KL}^{Q}(Q), \quad (6)$$

where both utilizes cross-entropy ($L_{CE}$), while the generator matches the batch normalization stats from the full-precision model ($L_{BNS}$) and the quantized network optimizes KL divergence ($L_{KL}$). Variants of GDFQ are currently forming state-of-the-art family of the zero-shot quantization, by adopting better generator [70], adversarial training [38], or boundary-supporting sample generations [9]. In this work, we provide an in-depth analysis of the loss function $\mathcal{L}_{GDFQ}(Q)$, and a novel scheme to improve its performance.

## 3. Analyses on the Zero-shot Quantization

In this section, we provide in-depth analyses on the first/second-order loss surface of the state-of-the-art zero-shot quantization method. We studied the impact of CE and KL on the loss function of zero-shot quantization setting, and reveal that they are not cooperative but hinder each other (Sec. 3.1). Then, we investigate the difference of the local curvature with the lens of Hessian. Because the zero-shot quantization suffers from larger generalization gaps, finding the solution of the flatter minima is critical (Sec. 3.2).

### 3.1. Gradient Cosine Similarity

In this subsection, we attempt to find a partial answer to the question: <u>are CE and KL cooperative in quantization?</u> As discussed in Sec. 2.2, loss functions from current techniques for zero-shot quantization [9, 61, 70] are mainly composed of the CE against the hard labels, and KL divergence against the full-precision teacher.
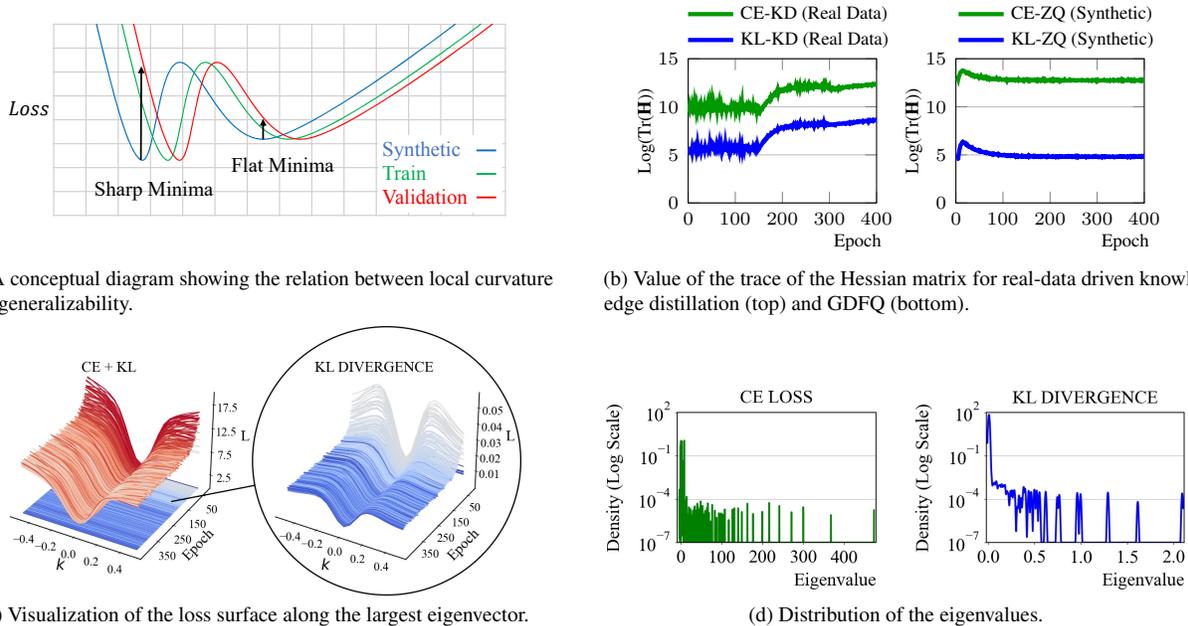
However, it has been discussed that better models do not necessarily make good teachers when the student has limited capacity. In such cases, the student often has to make trade-off between KL and CE [7]. As quantized models have much lower representational capability [23, 45], it could be difficult for them to optimize both terms. Furthermore, with synthetic samples not exactly matching the distribution of the real samples, the labels associated to each sample via hard labels and teacher outputs can be distinct, contributing to the difficulty of addressing both losses.

In such regard, [15] suggested using cosine similarity of the two gradients as a metric for determining whether an auxiliary can contribute toward a single main task. The authors suggest that the two losses should be used together only in steps where their cosine similarity is larger than zero (when they form an acute angle).

Thus, inspired by the proposals of [15], we analyzed the training of the zero-shot quantization as the following. Using GDFQ [61] as a representative for zero-shot quantization, we measured the cosine similarity of the $g_{KL}$ and $g_{CE}$ while training a quantized ResNet-20 model with synthetic data. The full-precision teacher is pre-trained, and the generator is jointly trained with the quantized student using Eq. (5) and Eq. (6) with $\alpha = \delta = 0.5$, respectively. For comparison, we have also measured the same metric from a common real data (CIFAR-10) based knowledge distillation with pre-trained ResNet-20 as a teacher and random initialized ResNet-20 as a student (self-distillation) using the same loss function.

The results are presented in Fig. 1. Using the real samples denoted as 'KD', the KL and CE terms induce gradients of the similar direction, supporting the common wisdom that the combination of them works well in practice. However, with the zero-shot quantization setting using synthetic samples denoted as 'ZQ', the cosine distance between them take negative values. The bottom two plots visualize the angle of the gradients. In both plots, the $g_{CE}$ is set to (1,0), and $g_{KL}$ is plotted to preserve the relative angle to the $g_{CE}$.

The trend persists throughout the training, as shown in the plot. Right after the training begins, the cosine distance of ZQ becomes negative and it is maintained until the end of training, while that of the KD is positive. This implies that combinations of the two losses do not cooperate well with each other, and using them together could potentially harm the model performance. Although we display only one case for clarity, the same trend was observed across many models and datasets. Refer to the Appendix for further results.

(a) A conceptual diagram showing the relation between local curvature and generalizability.

(b) Value of the trace of the Hessian matrix for real-data driven knowledge distillation (top) and GDFQ (bottom).

(c) Visualization of the loss surface along the largest eigenvector.

(d) Distribution of the eigenvalues.

Figure 2. Analysis of the loss surface.

## 3.2. Generalizability

The observation in Sec. 3.1 suggest that using only one of the two losses — KL divergence against the full-precision teacher or CE against the hard label — could be better for the problem. Some work [7] suggest modifying the teacher for distillation, but such method is unavailable in a zero-shot setting because we have no access to the training data

In such regard, we examine the generalizability of the loss terms. Including zero-shot quantization, diverse applications relying on synthetic samples [12, 35, 37, 50, 57, 60, 63] usually suffer from huge generalization gap, coming from the discrepancy in the data distributions. One can easily infer that quantization requires stronger generalization when performed under a zero-shot environment.

To evaluate the generalizability, we measure local curvature of the loss surface. Popularly measured with the Hessian matrix $\mathbf{H}$ ($\frac{\partial^2 L}{\partial \theta^2} \in \mathbb{R}^{n \times n}$, where $\theta$ is a vector of $n$ weight parameters), local curvature of the loss surface is a metric that is drawing much attention from the field, and it is thought to hold a key to better generalization [2, 6, 24–27, 30]. As illustrated in Fig. 2a, if the optimizer settles at a sharp minima, the performance at test time is likely to incur a larger degradation compared that of a flat minima. Such gap would be much larger with synthetic data incorrectly modeling the validation data distribution. In line with this finding, many literatures support the claim of smaller local curvature improving generalization [2, 6, 24–27, 30].

Fig. 2b plots $Tr(\mathbf{H})$, the trace of the Hessian matrix, approximated by PyHessian [62] implementing Lanczos algorithm [34]. We separate Hessian calculation for each

of CE and KL. The trace values are significantly different, where that of the KL is much smaller than the CE. The gap is notably larger in zero-shot quantization (right, ZQ) than on real-data knowledge distillation (left, KD). In addition, the distribution of the eigenvalues displayed in Fig. 2d also show a huge difference in the local curvature of the loss terms. While CE has longer tail for high eigenvalues, those of KL has more concentration to lower eigenvalues.

This could potentially lead to two conclusions: the loss surface of the KL divergence is much flatter, or the model has converged to the minima in the loss surface of the KL divergence. However, in our case, we believe it advocates for the former, based on an auxiliary experiment. It is commonly observed that near the minima, the disparity between gradients start to arise [17, 42]. Following the same regard, we measure the cosine distance of the gradient averaged within an epoch, compared to that of the previous epoch (inter-epoch cosine similarity). As displayed in Fig. 3a, the gradients of KL in the zero-shot quantization setting (ZQ) points to a consistent direction (large cosine similarity) compared to that of the real data distillation (KD), indicating that it has not reached the minima yet.

Fig. 2c illustrates a more direct visualization of the loss surface. From the Hessian matrix, we took the largest eigenvector $e$ at each epoch, we plotted the value of CE and KL on the left plot by calculating $L(\theta(t) + k \cdot e \cdot \hat{g}(t))$ for $k \in [-0.5, 0.5]$, where $\hat{g}$ is the average gradient along $e$. The left plot presents the CE in red color scheme and KL in blue color scheme. It is clearly shown that the surface is flatter in the KL surface especially near the end of training.
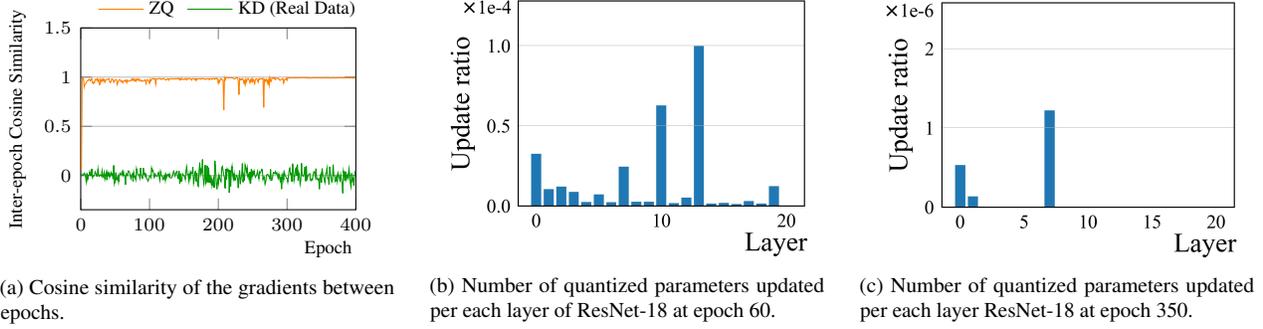
(a) Cosine similarity of the gradients between epochs.

(b) Number of quantized parameters updated per each layer of ResNet-18 at epoch 60.

(c) Number of quantized parameters updated per each layer ResNet-18 at epoch 350.

Figure 3. Analysis on the naive KL-only zero-shot quantization.

## 3.3. Summary

Summarizing the studies in this section, we first observed that the CE and KL form a large angle in the gradient space, and the quantized model has difficulty optimizing both directions. Furthermore, by measuring the stats from Hessian matrices, we conclude that KL has a much flatter loss surface for potentially better generalization, which is an important issue for generative zero-shot quantization methods.

## 4. AIT Method

In this section, we describe our AIT (All In the Teacher) method in detail. From the observations in Sec. 3, we first drop CE term from the loss, and apply a novel gradient inundation to bring the quantized model as close as possible to the full-precision teacher.

### 4.1. KL-only Zero-shot Quantization

Motivated by the experiments from Sec. 3.2, we ran GDFQ [61] with KL-only loss (i.e., $\delta = 1$). However, as will be shown later in Tab. 2, the performance severely degrades in all settings. We find an explanation from the experiments of Fig. 3a. Even toward the end of the training, the direction of the $g_{KL}$ remains consistent, and training for more epochs did not solve the problem. This indicates that the model did not converge at the minima of the KL surface.

Another set of experiments shown in Figs. 3b and 3c gives a closer look at the phenomenon. We count the average number of weight parameters that cross the rounding threshold (parameters whose quantized values have changed from the previous step). We make two observations: First, the portion of quantized values crossing the rounding threshold is extremely small. Even when training has not stabilized (epoch 60), only 0.0011% of weights are being updated each step. At a later epoch (350), the phenomenon becomes worse that only four values are updated in a whole epoch, which is only 1.8e-7% of weight updates per step during the epoch. In addition, the changes are extremely unbalanced, where all the updates are only occurring in just three layers.

We posit that this is from the quantized training process

that constrains integer value updates. During the training, the quantized network internally stores its full-precision values. The parameters are quantized for the forward pass of the backpropagation, and the gradients are applied to the internal full precision values. As the gradient values become smaller after a few epochs of training, the change in the parameters are usually not large enough to cross the threshold, and only a few layers are continuously making changes, stopping the model from moving towards a lower point in the loss surface.

### 4.2. Gradient Inundation

To address the problem of KL-only method, we propose gradient inundation (GI). Overall, we attempt to dynamically manipulate the gradients $g_l$ of each layer $l$, such that certain number of parameters are guaranteed to update in its integer value. With stochastic gradient descent, consider the update rule of parameter $\theta_{l,k}$ at step $k$ with learning rate $\eta$:

$$\theta_{l,k+1} = \theta_{l,k} - \eta \cdot g_{l,k}. \tag{7}$$

with gradient inundation, the modified rule is as the following: For the parameters $\theta_{l,k}$, quantized parameters $\theta_{l,k}^q$ and the corresponding gradients $g_{l,k}$ from layer $l$,

$$\theta_{l,k+1} = \theta_{l,k} - \eta \cdot g'_{l,k}, \tag{8}$$

$$g'_{l,k} = \kappa_l \cdot g_{l,k}, \tag{9}$$

$$\kappa_l = \arg\min_{\kappa_l} \|\Delta\theta_{l,k}^q - T\|, \tag{10}$$

$$\Delta\theta_{l,k}^q = \sum \mathbb{I}(\theta_{l,k}^q \neq \theta_{l,k+1}^q), \tag{11}$$

$$T = \rho \cdot dim(\theta_l), \tag{12}$$

where $\rho \in [0, 1]$ is a predetermined proportion that exceeds the quantization threshold, $\mathbb{I}()$ is the indicator function, and $dim(\theta_l)$ is the number of elements in $\theta_l$. Our goal is to find $\kappa_l$ that guarantees the number of parameter updates on a quantized layer $\Delta\theta_{l,k}^q$ exceeds a certain ratio $T$.

To quickly find an approximate solution, we applied a simple two-step heuristic. Firstly, starting from 1.0, $\kappa_l$ is doubled until $\Delta\theta_{l,k}^q > T$. Then, to satisfy Eq. 10, $\kappa_l$ is adjusted between $\kappa/2$ and $\kappa$ by binary search. For the sake of

| Dataset | Model (FP32 Acc.) | Bits | ZeroQ | GDFQ | GDFQ +AIT | Qimera | Qimera +AIT | ARC | ARC +AIT |
|---|---|---|---|---|---|---|---|---|---|
| CIFAR-10 | ResNet-20 | 4w4a | 79.30 | 90.25 | 91.23 ( +0.98 ) | **91.26** | 91.23 ( −0.03 ) | 88.55 | 90.49 ( +1.94 ) |
| | 93.89 | 5w5a | 91.34 | 93.38 | 93.41 ( +0.03 ) | **93.46** | 93.43 ( −0.03 ) | 92.88 | 92.89 ( +0.01 ) |
| CIFAR-100 | ResNet-20 | 4w4a | 47.45 | 63.39 | **65.80** ( +2.41 ) | 65.10 | 65.40 ( +0.30 ) | 62.76 | 61.05 ( −1.71 ) |
| | 70.33 | 5w5a | 65.61 | 66.12 | **69.26** ( +3.14 ) | 69.02 | **69.26** ( +0.24 ) | 68.40 | 68.40 ( +0.00 ) |
| ImageNet | ResNet-18 | 4w4a | 22.58 | 60.60 | 65.51 ( +4.91 ) | 63.84 | **66.83** ( +2.99 ) | 61.32 | 65.73 ( +4.41 ) |
| | 71.47 | 5w5a | 59.26 | 68.40 | 70.01 ( +1.61 ) | 69.29 | 69.22 ( −0.07 ) | 68.88 | **70.28** ( +1.40 ) |
| | ResNet-50 | 4w4a | 8.38 | 52.12 | 64.24 ( +12.12 ) | 66.25 | 67.63 ( +1.38 ) | 64.37 | **68.27** ( +3.90 ) |
| | 77.73 | 5w5a | 48.12 | 71.89 | 74.23 ( +2.34 ) | 75.32 | 75.54 ( +0.22 ) | 74.13 | **76.00** ( +1.87 ) |
| | MobileNetV2 | 4w4a | 10.96 | 59.43 | 65.39 ( +5.96 ) | 61.62 | 66.81 ( +5.19 ) | 60.13 | **66.47** ( +6.34 ) |
| | 73.03 | 5w5a | 59.88 | 68.11 | 71.70 ( +3.59 ) | 70.45 | 71.68 ( +1.23 ) | 68.40 | **71.96** ( +3.56 ) |

Table 1. Comparison on AIT with data-free quantization schemes.

computation efficiency, the total number of search steps is limited to five. In addition, to assure early phase stability of the training, we added a warm-up phase for the GI method. In the warm-up phase, the maximum of $\kappa$ is limited to 128 for more accurate solutions. When the generator requires a separate warm-up, the GI warm-up phase starts after the generator warm-up ends. Similar to learning rate exponential decay scheduling, we apply the exponential decay to $\rho$. We discuss the sensitivity to this in Sec. 5.4.

# 5. Experimental Results

## 5.1. Experimental Environments

We evaluate AIT on three datasets, CIFAR-10, CIFAR-100 [32], and ImageNet (ILSVRC2012 [33]). CIFAR-10 and CIFAR-100 contain 10 and 100 classes of images, respectively, and represent small-scale datasets. ImageNet has 1000 classes of images with 1.2M training samples and 50K validation samples, which represent large-scale dataset.

For CIFAR-10/100, we use commonly used ResNet-20 [20] model. For ImageNet, we use ResNet-18 and ResNet-50 to represent popular medium- and large-sized models, and MobileNetV2 [56] to represent a lightweight model. All pretrained models are from pytorchcv library [1]. For more results on various models, please refer to the Appendix.

For baselines, we use the official code provided by the authors of ZeroQ [4], GDFQ [61], ARC [70] and Qimera [9] with the identical settings. AIT is implemented using PyTorch [51] version 1.10.0. All experiments are conducted using NVIDIA RTX3090 and A6000 GPUs.

The generator is trained with the loss function Eq. (5) with $\alpha = 0.5$ using Adam optimizer [31] with learning rate of 0.001. For training the quantized student model, SGD with Nesterov [46] was used with momentum 0.9. For AIT, the hyperparameter $\rho$ was set to 0.001 and 0.0001 for CIFAR and ImageNet respectively, which both were decayed by 0.1 every 100 epochs. Experiments on CIFAR and ImageNet were run for 400 epochs on learning rate $\eta =$1e-4, with batchsize of 200 and 16, respectively.

## 5.2. Performance Comparison

AIT can be applied to most generative zero-shot quantization methods. In this section, we apply our method to three: GDFQ [61], the first method to suggest such approach, ARC [70], which improves the generator, and Qimera [9], the SOTA technique in the same family. We also include ZeroQ [4] for comparison. We report top-1 accuracies.

Overall, AIT achieves significant performance improvements in most settings tested, whether implemented on top of GDFQ, ARC or Qimera. Notably large improvements have been observed on ImageNet datasets, especially in 4w4a settings, because there still exists a large gap towards the full-precision (32bit) model. The largest gain was found for 4w4a ResNet-50 on top of GDFQ, with the gain of 12.12%p that seems to mainly come from the huge gap (25.61%p) GDFQ originally had between the full-precision model. For the results of lower-bit settings, refer to Appendix.

An interesting trend is that for the other two methods with better generators (Qimera, ARC), the performance gain on 4w4a setting is larger for smaller models (ResNet-50→ResNet-18→MobileNetV2). The improvements are (+1.38%p, +2.99%p, +5.19%p) for Qimera and (+3.90%p, +4.41%p, +6.34%p) for ARC in a descending order of model size. This indirectly supports our claim that quantized networks with smaller capacity have difficulties optimizing for multiple loss terms, and AIT can alleviate such effect.

In addition, the performance of ARC+AIT is better than Qimera+AIT for all ImageNet settings except one, even though Qimera outperforms ARC in their default settings. We find the reason from generator model size of ARC. While Qimera uses the exact same generator from GDFQ, ARC uses a larger generator model found by neural architecture search. The result demonstrate that AIT is making better use of the potential of the generator network. Small performance degradations were observed for CIFAR-10 on Qimera by 0.03%p. Since CIFAR-10 is a small dataset and the performance is already close to the fp32 model, we believe this is because there is not much room left to improve.

## 5.3. Ablation Study

Tab. 2 shows an ablation study performed over GDFQ. The ResNet family and MobileNet are denoted as 'RN' and 'MB' respectively. 'KL-only' drops CE from the original loss function of GDFQ, and let the quantized model optimize only on the KL divergence against the full-precision teacher. However, this results in a huge degradation in all settings. As analyzed in Sec. 4.1, this is due to the scarce, unbalanced quantized weight updates. By applying gradient inundation, the lost performance is more than recovered and the superior gain over the baseline is obtained (KL-only+GI).

'Baseline+GI' represents the gradient inundation applied on top of GDFQ without dropping the CE term, and 'CE only+GI' represents the same with KL term dropped from the loss. Unfortunately, they only result in performance degradation, because the baseline GDFQ with CE+KL loss or CE loss does not suffer from the aforementioned scarce weight update problems. Therefore, gradient inundation only makes detrimental changes to the quantized weights.

Seeing the effect of 'KL-only + GI', one might wonder if the weight update problem can be solved by simply increasing the learning rate. 'KL-only (high lr)' row shows the results of such experiments conducted with ×100 learning rate. In addition, Fig. 4 shows the distribution of the updates in each layer, in comparison to Fig. 3c. 'KL-only (high lr)' achieves a small gain but does not entirely solve the problem. First, increasing the learning rate incurs too frequent updates in a few layer which was already getting enough updates as shown in Figs. 4a and 4b, and many layers still not being updated. Moreover, further increasing the learning rate results in divergence of the model. Figs. 4c and 4d presents the number of updates with gradient inundation, where $\rho$ is depicted in dotted red lines. Gradient inundation tunes the gradients to the right level, leading to a better performance.

For a comprehensive comparison, we also tested adaptive optimizers such as Adam [31], RMSProp [59], and especially LARS [64], which adjusts learning rates per layer. The results of these optimizers for the quantized model are shown in Tab. 3. 'Baseline-' denotes the existing methods, and 'GI-' denotes modified optimizer with GI. The results show that 'Baseline-Adam' and 'Baseline-RMSProp' suffer from noticeable accuracy degradation, especially on the large models. 'Baseline-LARS' survived from such trend, it does not

| Dataset | CIFAR-10 | CIFAR-100 | ImageNet | | |
|---|---|---|---|---|---|
| Model | RN-20 | RN-20 | RN-18 | RN-50 | MB-V2 |
| Baseline-LARS | 90.01 | 63.84 | 58.94 | 52.98 | 59.58 |
| Baseline-SGD | 90.25 | 63.39 | 60.60 | 52.12 | 59.43 |
| Baseline-Adam | 91.12 | 57.39 | 40.97 | 30.16 | 26.35 |
| Baseline-RMSProp | 89.88 | 63.18 | 51.12 | 40.65 | 31.52 |
| GI-SGD (ours) | 91.23 | **65.80** | 65.51 | 64.24 | 65.39 |
| GI-Adam (ours) | **91.33** | 64.38 | 65.47 | **65.67** | 61.33 |
| GI-RMSProp (ours) | 90.82 | 65.78 | **65.73** | 62.42 | **65.41** |

Table 3. GI on Various Optimizers.

| Dataset | Model | $\rho$ | | | | |
|---|---|---|---|---|---|---|
| | | 0.005 | 0.001 | 0.0005 | 0.0001 | 0.00005 |
| CIFAR-100 | ResNet-20 | 63.01 | 65.80 | 65.41 | 65.04 | 65.30 |
| ImageNet | ResNet-18 | 57.95 | 60.11 | 64.48 | 65.51 | 65.92 |

Table 4. Sensitivity Analysis.

make significant differences compared to 'Baseline-SGD'. We further expand our research by applying GI to Adam and RMSProp, written as 'GI-Adam' and 'GI-RMSProp'. The results verify that our method has solid performance on various optimizers and still outperforms existing methods.

## 5.4. Sensitivity Analysis

The $\rho$ value controls the portion of quantized weights guaranteed to get updates in each layer. In Tab. 4, a sensitivity study of $\rho$ on top of GDFQ+AIT has been performed. The results show that AIT is not very sensitive to $\rho$, although there is some effect. Refer to the Appendix for more results.

Tab 5. shows the learning rate sensitivity of our method with comparison to the GDFQ. The results show that the GI method is robust to the learning rate changes on both datasets while steadily outperforms the baseline method.

## 5.5. Further Analysis

In this section, we present more details of AIT. Fig. 5a shows the KL divergence over the training in the baseline GDFQ, 'KL-only' and AIT. AIT is able to reach a lower KL distance. This supports our analysis from Sec. 4.1 that there is still room for KL to be optimized further.

Another observation can be found from Fig. 5b, where we have measured the angle between the CE and KL. For 'KL-only' and AIT, CE is calculated only for measurements, and did not affect the SGD updates. It is now shown that 'KL-only' and AIT, both have positive cosine similarities

| Dataset | CIFAR-10 | CIFAR-100 | ImageNet | | |
|---|---|---|---|---|---|
| Model | RN-20 | RN-20 | RN-18 | RN-50 | MB-V2 |
| Baseline (GDFQ) | 90.25 | 63.39 | 60.60 | 52.12 | 59.43 |
| KL-only | 90.06 | 58.93 | 58.49 | 42.64 | 47.03 |
| KL-only (high lr) | **92.20** | 62.20 | 65.34 | 61.68 | 64.70 |
| Baseline + GI | 89.32 | 59.05 | 55.01 | 44.09 | 43.57 |
| CE-only + GI | 90.89 | 51.57 | 52.72 | 27.86 | 33.88 |
| AIT (KL-only + GI) | 91.23 | **65.80** | **65.51** | **64.24** | **65.39** |

Table 2. Ablation Study.

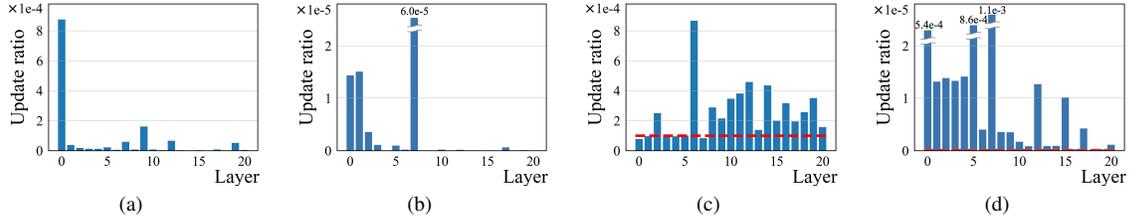| $\eta$ | Cifar-100 (RN-20) | | ImageNet (RN-18) | |
|---|---|---|---|---|
| | GDFQ | AIT | GDFQ | AIT |
| 1e-2 | 49.81 | 66.96 | 40.78 | 65.69 |
| 1e-3 | 58.10 | 66.21 | 40.90 | 65.57 |
| 1e-4 | 63.39 | 65.80 | 53.28 | 65.51 |
| 1e-5 | 61.08 | 65.92 | 59.32 | 65.70 |
| 1e-6 | 59.47 | 65.73 | 60.60 | 65.23 |

Table 5. Learning Rate Sensitivity Analysis.

Figure 4. Distribution of the updated quantized values with ×100 learning rate at an earlier epoch (a), later epoch (b) and with gradient inundation at an earlier epoch (c), later epoch (d).
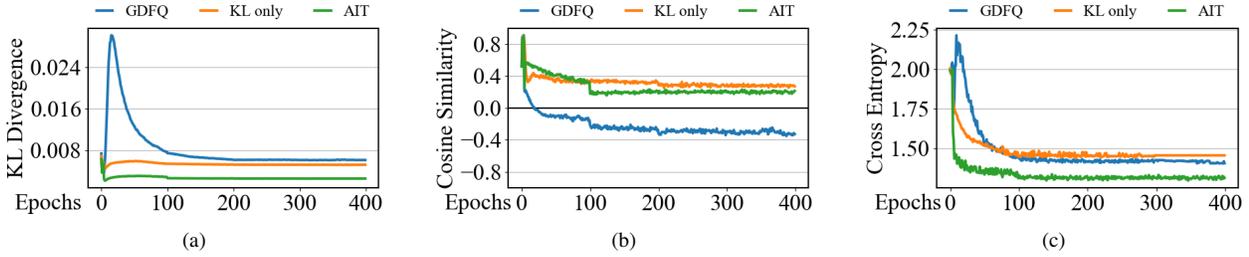


Figure 5. Further experiments. (a) KL divergence. (b) Cosine similarity of KL and CE for GDFQ, KL-only, and AIT. (c) Cross-entropy against the validation samples.
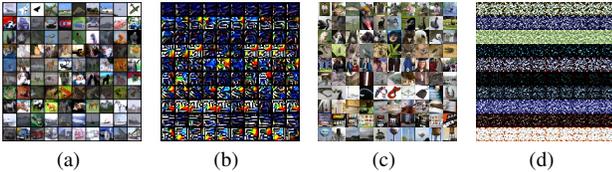


Figure 6. Comparison of samples. (a) Real CIFAR-10 samples (b) synthetic CIFAR-10 samples (c) Real ImageNet samples (d) synthetic ImageNet samples. Each row represents one of the 10 classes for CIFAR-10, and 10 randomly chosen classes of ImageNet.

between the losses. This represents that, although they are difficult to optimize concurrently in the beginning, as they get near the global minima, they become to share the same direction for optimization. AIT sacrifices a small amount of directional alignment, but outperforms 'KL-only' by guaranteeing quantized weight updates with gradient inundation. Lastly, Fig. 5c shows that AIT achieves lower CE in validation against the hard label, even though it has not seen the real data, and not optimized for CE.

## 6. Discussion

**Removing Cross-Entropy** against the hard label from the loss function as done in AIT could be a penalty, because some methods [9] rely on the sample labels. However, as we demonstrated in Sec. 5, AIT applied to Qimera was able to obtain significantly better performance despite the exclusion of the mixed labels. In addition, our method does not depend on per-image hard label, so it can be widely used for segmentation [40, 55] or object detection [53, 54].

**Privacy Leak** is one societal concern of zero-shot quantization because the generator creates synthetic samples that

follow the distribution of real data. As several input reconstruction techniques point out [14, 41], it could be that the synthetic samples can reconstruct the private training data. However, to the extent of our observation, there is no sign that AIT reconstructs the real data as shown in Fig. 6. Training method for the generator in AIT is no different from its baselines, since it does not alter the generator loss in Eq. (5), thus does not contribute any further to the privacy leak.

## 7. Conclusion

In this work, we analyzed the SOTA family of solutions for zero-shot quantization. Through a series of experiments and analyses, we found that current solutions can be improved through pursuing a flatter minima and guaranteeing weight updates during fine-tuning. We achieve the goal by bringing the quantized model closer to the full-precision model in terms of the KL divergence and designing AIT— before our method, people habitually have used a combination of CE and KL as the main loss of the zero-shot quantization. Experimental results show that AIT is effective and can be easily applied to existing algorithms.

## Acknowledgements

# References

[1] Computer vision models on PyTorch. 6

[2] Alessandro Achille, Matteo Rovere, and Stefano Soatto. Critical learning periods in deep neural networks. In International Conference for Learning Representations, 2019. 4

[3] Ron Banner, Yury Nahshan, Elad Hoffer, and Daniel Soudry. Post-training 4-bit quantization of convolution networks for rapid-deployment. arXiv preprint arXiv:1810.05723, 2018. 1, 2

[4] Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. ZeroQ: A novel zero shot quantization framework. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020. 1, 2, 6

[5] Zhaowei Cai, Xiaodong He, Jian Sun, and Nuno Vasconcelos. Deep learning with low precision by half-wave gaussian quantization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2017. 2

[6] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-SGD: Biasing gradient descent into wide valleys. In International Conference for Learning Representations, 2017. 4

[7] Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019. 1, 3, 4

[8] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. PACT: Parameterized clipping activation for quantized neural networks. arXiv preprint arXiv:1805.06085, 2018. 1

[9] Kanghyun Choi, Deokki Hong, Noseong Park, Youngsok Kim, and Jinho Lee. Qimera: Data-free quantization with synthetic boundary supporting samples. In Advances in Neural Information Processing Systems, 2021. 1, 3, 6, 8

[10] Yoojin Choi, Jihwan Choi, Mostafa El-Khamy, and Jungwon Lee. Data-free network quantization with adversarial knowledge distillation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020. 1

[11] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Towards the limit of network quantization. In International Conference for Learning Representations, 2017. 2

[12] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Dual-teacher class-incremental learning with data-free generative replay. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021. 4

[13] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. BinaryConnect: Training deep neural networks with binary weights during propagations. In Advances in Neural Information Processing Systems, 2015. 1, 2

[14] Alexey Dosovitskiy and Thomas Brox. Inverting visual representations with convolutional networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2016. 8

[15] Yunshu Du, Wojciech M Czarnecki, Siddhant M Jayakumar, Mehrdad Farajtabar, Razvan Pascanu, and Balaji Lakshminarayanan. Adapting auxiliary losses using gradient similarity. arXiv preprint arXiv:1812.02224, 2018. 3

[16] Julian Faraone, Nicholas Fraser, Michaela Blott, and Philip HW Leong. SYQ: Learning symmetric quantization for efficient deep neural networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018. 2

[17] Mahsa Forouzesh and Patrick Thiran. Early stopping by gradient disparity. 2020. 4

[18] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. arXiv preprint arXiv:1510.00149, 2015. 2

[19] Matan Haroush, Itay Hubara, Elad Hoffer, and Daniel Soudry. The knowledge within: Methods for data-free model compression. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020. 1

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2016. 6

[21] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In Advances in Neural Information Processing Systems Workshops, 2014. 1

[22] Kyuyeon Hwang and Wonyong Sung. Fixed-point feedforward deep neural network design using weights+ 1, 0, and-1. In International Workshop on Signal Processing Systems, 2014. 1

[23] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018. 1, 2, 3

[24] Stanisław Jastrzębski, Devansh Arpit, Oliver Astrand, Giancarlo B Kerg, Huan Wang, Caiming Xiong, Richard Socher, Kyunghyun Cho, and Krzysztof J Geras. Catastrophic fisher explosion: Early phase fisher matrix impacts generalization. In International Conference on Machine Learning, 2021. 4

[25] Stanisław Jastrzębski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in SGD. In International Conference of Artificial Neural Networks, 2018. 4

[26] Stanisław Jastrzębski, Zachary Kenton, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. On the relation between the sharpest directions of dnn loss and the sgd step length. In International Conference for Learning Representations, 2018. 4

[27] Stanisław Jastrzębski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho, and Krzysztof Geras. The break-even point on optimization trajectories of deep neural networks. In International Conference for Learning Representations, 2020. 4

[28] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-Datacenter Performance Analysis of a Tensor Processing Unit. In International Symposium on Computer Architecture, 2017. 1

[29] Sangil Jung, Changyong Son, Seohyung Lee, Jinwoo Son, Jae-Joon Han, Youngjun Kwak, Sung Ju Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing quantization intervals with task loss. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019. 2

[30] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In International Conference for Learning Representations, 2017. 4

[31] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International Conference for Learning Representations, 2015. 6, 7

[32] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images, 2009. 6

[33] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In Advances in Neural Information Processing Systems, 2012. 6

[34] Cornelius Lanczos. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. 1950. 4

[35] Jaehoon Lee, Jihyeon Hyeong, Jinsung Jeon, Noseong Park, and Jihoon Cho. Invertible tabular GANs: Killing two birds with one stone for tabular data synthesis. In Advances in Neural Information Processing Systems, 2021. 4

[36] Darryl Lin, Sachin Talathi, and Sreekanth Annapureddy. Fixed point quantization of deep convolutional networks. In International Conference on Machine Learning, 2016. 1

[37] Xialei Liu, Chenshen Wu, Mikel Menta, Luis Herranz, Bogdan Raducanu, Andrew D Bagdanov, Shangling Jui, and Joost van de Weijer. Generative feature replay for class-incremental learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020. 4

[38] Yuang Liu, Wei Zhang, and Jun Wang. Zero-shot adversarial quantization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021. 1, 3

[39] Zechun Liu, Baoyuan Wu, Wenhan Luo, Xin Yang, Wei Liu, and Kwang-Ting Cheng. Bi-Real net: Enhancing the performance of 1-bit CNNs with improved representational capability and advanced training algorithm. In Proceedings of the European Conference on Computer Vision, 2018. 2

[40] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2015. 8

[41] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2015. 8

[42] Maren Mahsereci, Lukas Balles, Christoph Lassner, and Philipp Hennig. Early stopping without a validation set. arXiv preprint arXiv:1703.09580, 2017. 4

[43] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. Advances in Neural Information Processing Systems Workshops, 2014. 3

[44] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019. 1, 2

[45] Yury Nahshan, Brian Chmiel, Chaim Baskin, Evgenii Zheltonozhskii, Ron Banner, Alex M Bronstein, and Avi Mendelson. Loss aware post-training quantization. arXiv preprint arXiv:1911.07190, 2019. 3

[46] Yurii Evgen'evich Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. Dokl. Akad. Nauk SSSR, 1983. 6

[47] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier GANs. In International Conference on Machine Learning, 2017. 3

[48] Eunhyeok Park, Junwhan Ahn, and Sungjoo Yoo. Weighted-entropy-based quantization for deep neural networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2017. 2

[49] Eunhyeok Park, Dongyoung Kim, and Sungjoo Yoo. Energy-efficient neural network accelerator based on outlier-aware low-precision computation. In International Symposium on Computer Architecture, 2018. 1

[50] Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin Kim. Data synthesis based on generative adversarial networks. Proceedings of the VLDB Endowment, 2018. 4

[51] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In Advances in Neural Information Processing Systems Workshops, 2017. 6

[52] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. XNOR-Net: Imagenet classification using binary convolutional neural networks. In European conference on computer vision, 2016. 1, 2

[53] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2016. 8

[54] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In Advances in Neural Information Processing Systems, 2015. 8

[55] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In International Conference on Medical Image Computing and Computer-Assisted Intervention, 2015. 8

[56] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobilenetV2: Inverted residuals and linear bottlenecks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018. 6

[57] Abhilash Reddy Shankarampeta and Koichiro Yamauchi. Few-shot class incremental learning with generative feature replay. In International Conference on Pattern Recognition Applications and Methods, 2021. 4

[58] Hardik Sharma, Jongse Park, Naveen Suda, Liangzhen Lai, Benson Chau, Vikas Chandra, and Hadi Esmaeilzadeh. Bit

fusion: Bit-level dynamically composable architecture for accelerating deep neural network. In International Symposium on Computer Architecture, 2018. 1

[59] Tijmen Tieleman, Geoffrey Hinton, et al. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural networks for machine learning, 4(2):26–31, 2012. 7

[60] Gido M van de Ven, Zhe Li, and Andreas S Tolias. Class-incremental learning with generative classifiers. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021. 4

[61] Shoukai Xu, Haokun Li, Bohan Zhuang, Jing Liu, Jiezhang Cao, Chuangrun Liang, and Mingkui Tan. Generative low-bitwidth data free quantization. In European Conference on Computer Vision, 2020. 1, 3, 5, 6

[62] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. PyHessian: Neural networks through the lens of the hessian. In IEEE International Conference on Big Data, 2020. 4

[63] Hongxu Yin, Pavlo Molchanov, Jose M Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deep-inversion. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020. 4

[64] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. arXiv preprint arXiv:1708.03888, 2017. 7

[65] Dongqing Zhang, Jiaolong Yang, Dongqiangzi Ye, and Gang Hua. LQ-Nets: Learned quantization for highly accurate and compact deep neural networks. In Proceedings of the European conference on Computer Vision, 2018. 1

[66] Xiangguo Zhang, Haotong Qin, Yifu Ding, Ruihao Gong, Qinghua Yan, Renshuai Tao, Yuhang Li, Fengwei Yu, and Xianglong Liu. Diversifying sample generation for accurate data-free quantization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021. 1, 2

[67] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Chris De Sa, and Zhiru Zhang. Improving neural network quantization without retraining using outlier channel splitting. In International Conference on Machine Learning, 2019. 1, 2

[68] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen. Incremental network quantization: Towards lossless CNNs with low-precision weights. In International Conference on Learning Representations, 2017. 1

[69] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. DoReFa-Net: Training low bitwidth convolutional neural networks with low bitwidth gradients. arXiv preprint arXiv:1606.06160, 2016. 1, 2

[70] Baozhou Zhu, Peter Hofstee, Johan Peltenburg, Jinho Lee, and Zaid Alars. AutoReCon: Neural architecture search-based reconstruction for data-free compression. In International Joint Conference on Artificial Intelligence, 2021. 1, 3, 6

## A. Code

As a part of the supplementary materials, the code used to conduct experiments is attached as a separate zip archive. The zip archive contains the implementation of the AIT method on multiple zero-shot quantization backbones: GDFQ [4], ARC [6], and Qimera [1]. For reproducibility, the experiment environment setting and training scripts are included for all backbones. The code is under the terms of the GNU General Public License v3.0.

## B. Lower Bit-width Experiments

Further experiments on GDFQ and ARC were conducted in lower-bit settings. The experiment results are shown in Tab. 6 Following the main paper, the ResNet family and MobileNet are denoted as 'RN' and 'MB', respectively. We tested 3w3a and 3w4a quantization settings for the ImageNet experiments and further down to 2w2a and 2w3a for Cifar-10/100, which we found to be the lowest bits GDFQ and AIT converge.

## C. Experiments on Additional Network Models

We conducted a further evaluation of our method on various networks: InceptionV3 [3], SqueezeNext [2], and ShuffleNet [5]. The experimental results are shown in Table 7. Compared with the GDFQ baseline, our method still outperforms by a huge margin on all settings regardless of the quantization bitwidth. Furthermore, experimental results show that AIT is especially effective on smaller networks. This result again supports our observation in the main body that the limited capacity of a small network hinders the training phase from matching multiple loss terms simultaneously.

## D. Comparison with Label Smoothing

Label smoothing is a regularization technique that replaces one-hot label $y$ into a smooth label $y'$ by

$$y' = (1-c)y + c/K, \qquad (13)$$

where $K$ is the number of classes and $c$ is a label smoothing value. Label smoothing is known to help neural network training to avoid overfitting and increase generalization capability. Therefore, one might think that label smoothing can also help flatten the cross-entropy (CE) loss surface by its nature. To answer the question, we conducted comparative experiments with various label smoothing parameters. The experiments evaluate how the label smoothing affects the performance of GDFQ baseline and CE-only setting, which drops KL divergence from the training loss.

Table 2 shows the experimental results. For CIFAR-10 and CIFAR-100, label smoothing did not improve performance in any settings over the baseline GDFQ, whether with KL divergence or not. Some improvements were observed from ImageNet dataset, but the improvements were smaller than that of AIT. This shows that label smoothing helps flatten the loss surface to some degree, its effect was not enough to reach that of AIT.

| Dataset | Model | Bits | GDFQ | | ARC | |
|---|---|---|---|---|---|---|
| | | | Baseline | AIT | Baseline | AIT |
| ImageNet | RN-18 | 3w3a | 20.69 | 36.34 | 1.00 | 36.70 |
| | | 3w4a | 39.73 | 53.55 | 2.54 | 56.77 |
| | RN-50 | 3w3a | 0.21 | 1.31 | 0.20 | 3.98 |
| | | 3w4a | 26.85 | 37.50 | 1.37 | 49.34 |
| | MB-V2 | 3w3a | 5.50 | 13.83 | 0.20 | 30.35 |
| | | 3w4a | 26.87 | 37.77 | 0.22 | 47.41 |
| CIFAR-100 | RN-20 | 2w2a | 1.41 | 2.09 | 1.35 | 1.55 |
| | | 2w3a | 1.04 | 1.13 | 1.25 | 1.14 |
| | | 3w3a | 49.62 | 48.64 | 28.54 | 34.39 |
| | | 3w4a | 59.70 | 61.37 | 50.47 | 58.65 |
| CIFAR-10 | RN-20 | 2w2a | 16.48 | 15.57 | 16.18 | 13.47 |
| | | 2w3a | 37.64 | 40.98 | 20.87 | 20.42 |
| | | 3w3a | 80.70 | 80.49 | 52.99 | 51.78 |
| | | 3w4a | 90.02 | 90.20 | 82.10 | 82.98 |

Table 6. Low Bit-width Experiments Results

| Dataset | Model (FP32 Acc.) | Bits | GDFQ | GDFQ +AIT |
|---|---|---|---|---|
| ImageNet | InceptionV3 79.00 | 4w4a | 70.57 | 73.34 ( +2.77 ) |
| | | 5w5a | 77.25 | 77.67 ( +0.42 ) |
| | SqueezeNext 69.39 | 4w4a | 26.21 | 45.37 ( +19.16 ) |
| | | 5w5a | 56.07 | 62.76 ( +6.69 ) |
| | ShuffleNet 65.07 | 4w4a | 19.72 | 27.80 ( +8.08 ) |
| | | 5w5a | 45.92 | 48.97 ( +3.05 ) |

Table 7. Additional experiments on various network models.

## E. Further Analysis on $\rho$ Sensitivity

We deepen the sensitivity analysis with finer levels of $\rho$ values. The experiments are conducted five times per setting to demonstrate performance stability regarding $\rho$ values. The results in Tab. 8 show that our method can achieve a stable accuracy level without hand-crafted hyperparameter tuning.

## F. Gradient Cosine Similarity

Although the main body of the manuscript offers results for gradient cosine similarity measured on ResNet20 with CIFAR-10 dataset, we have done an extensive amount of experiments to study the distinct gradient directionality spotted in zero-shot quantization task. Here we share the results to further support our findings.

For CIFAR-10 and CIFAR-100 dataset, we used ResNet-20, ResNet-56, ResNeXt-29 32x4d, WRN28-10, and WRN40-8. On ImageNet, we evaluated on ResNet-18, ResNet-50, MobileNetV2, and InceptionV3. The experiment compares the directionality of loss functions in training these networks under two different settings: zero-shot quantization (ZQ) and knowledge distillation (KD). In the knowledge distillation setting, we used the same network for both the student and the teacher (self-distillation) for fair comparison against the Zero-shot quantization setting.

Fig. 7 shows the results for CIFAR-10, and Fig. 8 for CIFAR-100. Although the quantitative difference of cosine similarities and the details of its change throughout the training differs across different datasets and networks, one trend is consistent: KL divergence and cross-entropy disagrees with each other more under the zero-shot quantization setting. Such tendency is usually maintained throughout the training.

| $\rho$ | CIFAR-100 RN-20 | ImageNet RN-18 | $\rho$ | CIFAR-100 RN-20 | ImageNet RN-18 |
|---|---|---|---|---|---|
| 0.0005 | 65.41±0.20 | 64.48±0.28 | 0.00009 | 65.20±0.29 | 65.84±0.07 |
| 0.0004 | 65.55±0.15 | 65.23±0.10 | 0.00008 | 65.29±0.19 | 65.66±0.17 |
| 0.0003 | 65.44±0.34 | 65.41±0.53 | 0.00007 | 65.35±0.18 | 65.65±0.05 |
| 0.0002 | 65.21±0.27 | 65.85±0.07 | 0.00006 | 65.06±0.23 | 65.52±0.16 |
| 0.0001 | 65.04±0.13 | 65.51±0.09 | 0.00005 | 65.30±0.10 | 65.92±0.42 |

Table 8. Sensitivity Analysis on $\rho$.

| Dataset | Model | Method | $c^*$ | | | | AIT |
|---|---|---|---|---|---|---|---|
| | | | 0.00† | 0.10 | 0.30 | 0.50 | |
| CIFAR-10 | RN-20 | Baseline | 90.25 | 89.67 | 88.85 | 88.52 | 91.23 |
| | | CE only | 88.36 | 88.67 | 88.21 | 87.89 | |
| CIFAR-100 | RN-20 | Baseline | 63.39 | 60.50 | 59.11 | 58.53 | 65.80 |
| | | CE only | 56.76 | 60.10 | 59.13 | 57.81 | |
| ImageNet | RN-18 | Baseline | 60.60 | 62.41 | 62.57 | 62.25 | 65.51 |
| | | CE only | 60.33 | 62.48 | 62.27 | 62.18 | |

†No smoothing *Label smoothing parameter.

Table 9. Performance of GDFQ with label smoothing in 4w4a setting.

# G. Hessian Trace

In this paper, Hessian matrix was used to measure the local curvature of the loss surface and compare the generalizability of the two distinct loss terms. Since Hessian matrix itself is enormous in size and computations involving its entirety is considered almost infeasible, analyzing the trace value of the matrix is often the most preferred way to study its characteristics. Adding to our results on Section 3.2 of the main body, we share further analysis on the loss curvature using Hessian trace.

We conducted further analysis on CIFAR-10 and CIFAR-100 datasets, on four different network models: ResNet-20, ResNet-56, WRN-28, and WRN-40. For all cases, our findings are the same. KL divergence has much smaller local curvature than the cross-entropy, where the gap is larger in zero-shot quantization settings.

# References

[1] Kanghyun Choi, Deokki Hong, Noseong Park, Youngsok Kim, and Jinho Lee. Qimera: Data-free quantization with synthetic boundary supporting samples. In Advances in Neural Information Processing Systems, 2021. 12

[2] Amir Gholami, Kiseok Kwon, Bichen Wu, Zizheng Tai, Xiangyu Yue, Peter Jin, Sicheng Zhao, and Kurt Keutzer. Squeezenext: Hardware-aware neural network design. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, June 2018. 12

[3] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, June 2016. 12

[4] Shoukai Xu, Haokun Li, Bohan Zhuang, Jing Liu, Jiezhang Cao, Chuangrun Liang, and Mingkui Tan. Generative low-bitwidth data free quantization. In European Conference on Computer Vision, 2020. 12

[5] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, June 2018. 12

[6] Baozhou Zhu, Peter Hofstee, Johan Peltenburg, Jinho Lee, and Zaid Alars. AutoReCon: Neural architecture search-based reconstruction for data-free compression. In International Joint Conference on Artificial Intelligence, 2021. 12
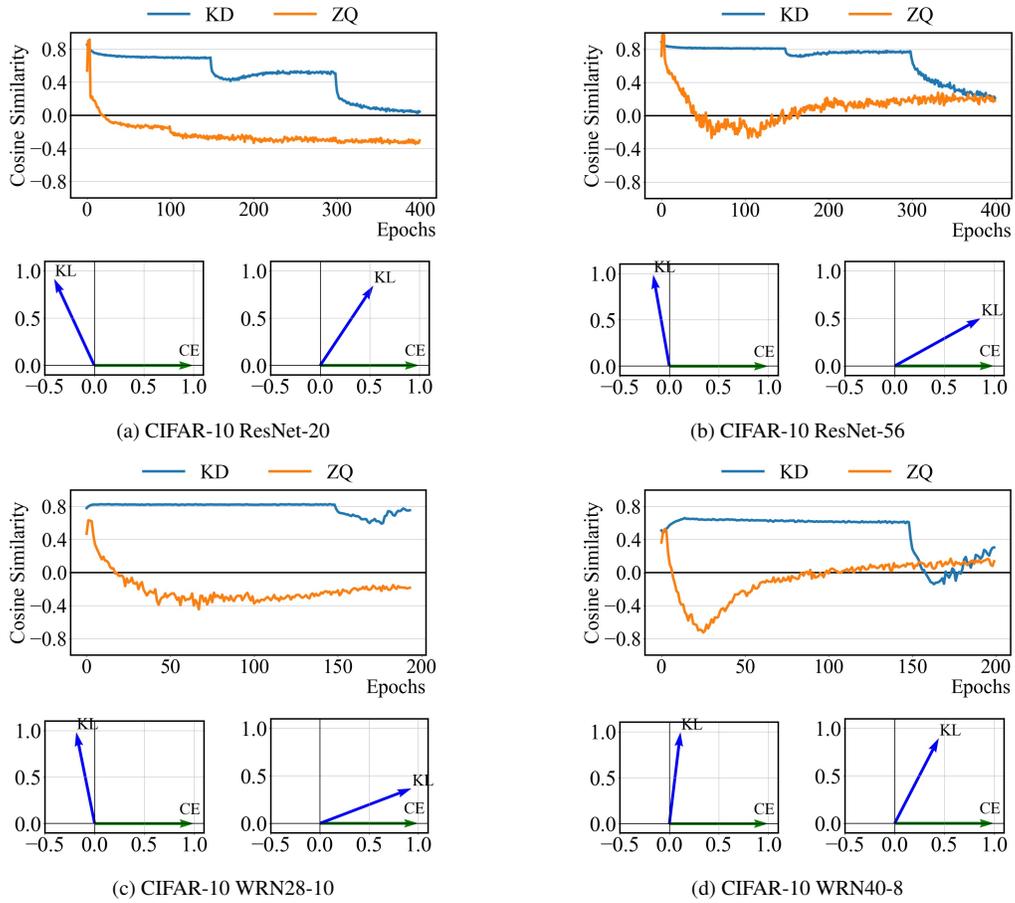
Figure 7. Gradient directionality of KL divergence and cross-entropy loss measured with CIFAR-10 dataset. In each setting, bottom left plots gradients under zero-shot quantization and bottom right plots gradients from knowledge distillation (self-distillation), captured from middle of the training.
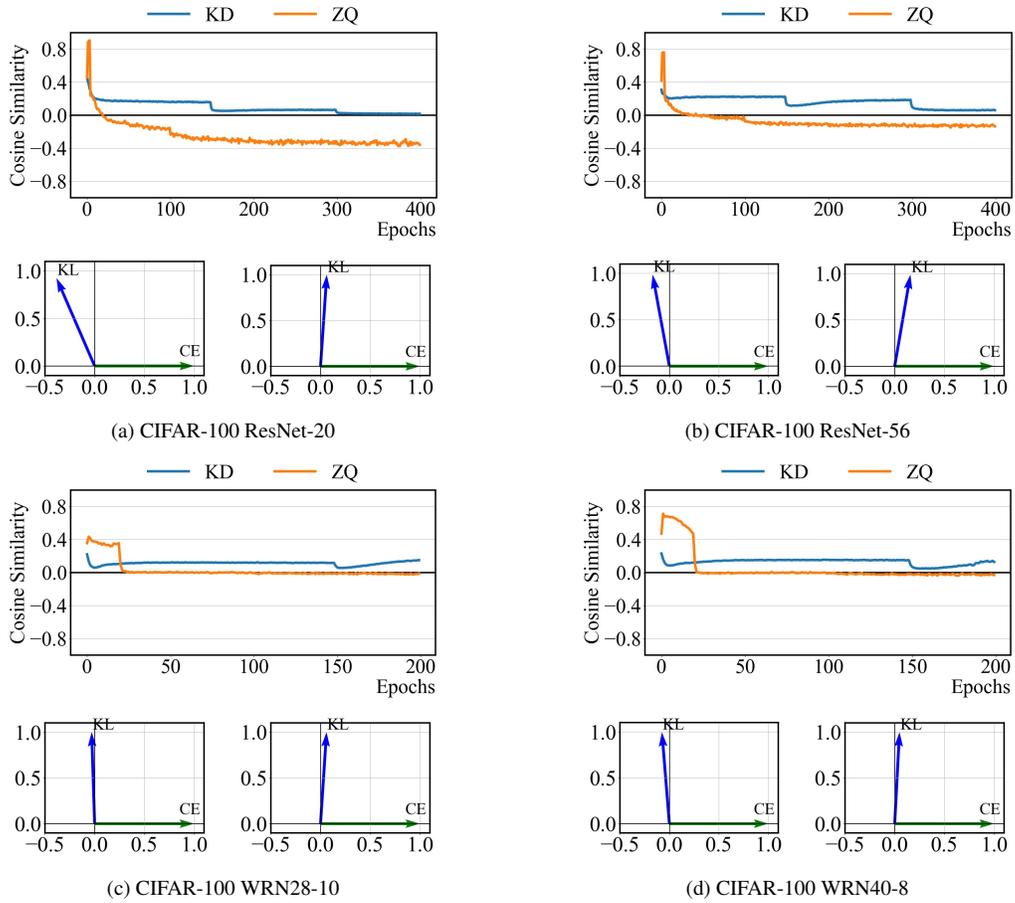
Figure 8. Gradient directionality of KL divergence and cross-entropy loss measured with CIFAR-100 dataset. In each setting, bottom left plots gradients under zero-shot quantization and bottom right plots gradients from knowledge distillation (self-distillation), captured from middle of the training.
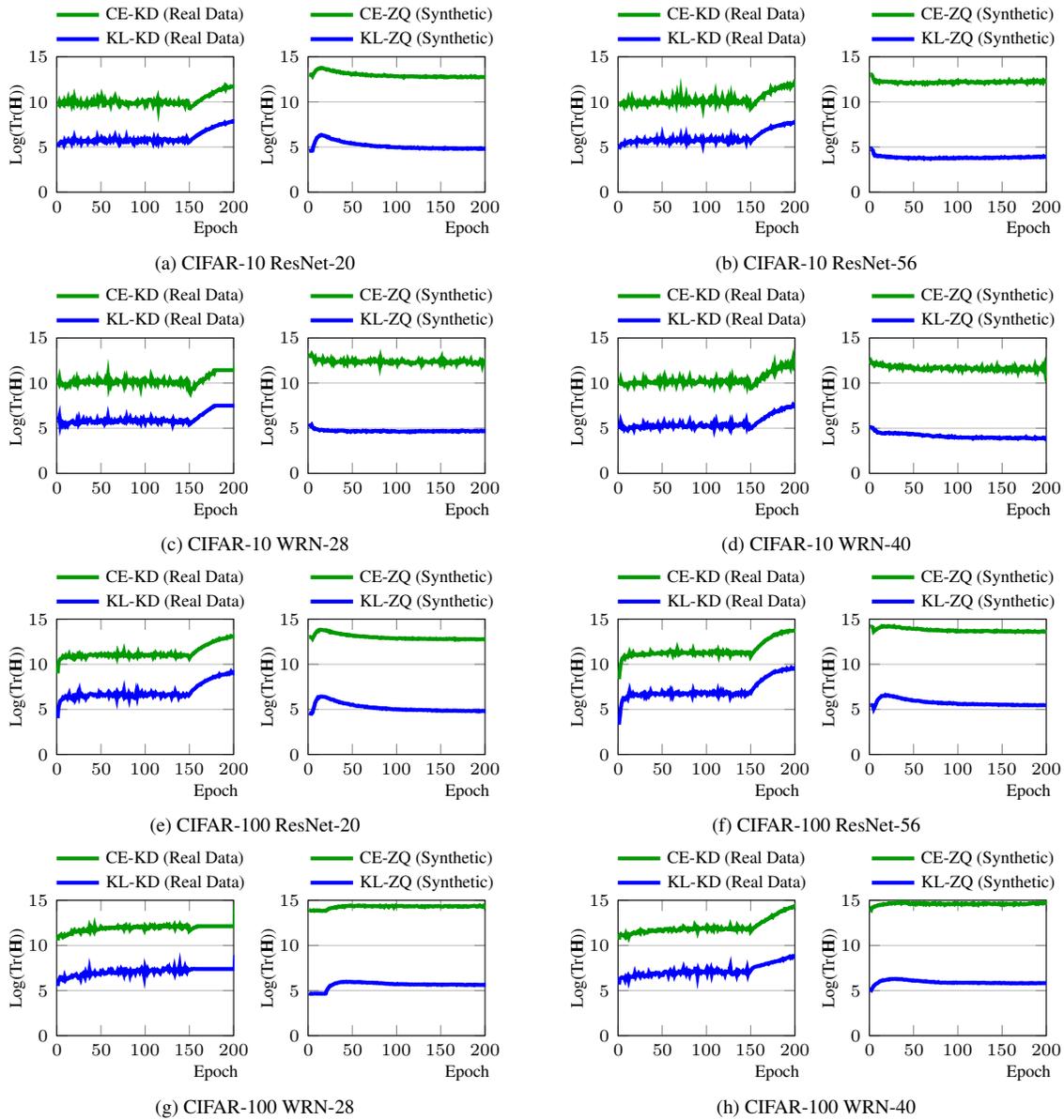
Figure 9. Hessian trace of KL divergence and cross-entropy, measured across diverse datasets and networks.