# THAT: <u>T</u>wo <u>H</u>ead <u>A</u>dversarial <u>T</u>raining for Improving Robustness at Scale

Zuxuan Wu[1,3]     Tom Goldstein[2]     Larry S. Davis[2]     Ser-Nam Lim[3]

[1] Fudan University     [2] University of Maryland     [3] Facebook AI

## Abstract

*Many variants of adversarial training have been proposed, with most research focusing on problems with relatively few classes. In this paper, we propose Two Head Adversarial Training (THAT), a two-stream adversarial learning network that is designed to handle the large-scale many-class ImageNet dataset. The proposed method trains a network with two heads and two loss functions; one to minimize feature-space domain shift between natural and adversarial images, and one to promote high classification accuracy. This combination delivers a hardened network that achieves state of the art robust accuracy while maintaining high natural accuracy on ImageNet. Through extensive experiments, we demonstrate that the proposed framework outperforms alternative methods under both standard and "free" adversarial training settings.*

## 1. Introduction

Convolutional neural networks have demonstrated remarkable performance in a multitude of computer vision tasks like image classification [16, 48], object detection [34, 3, 24], *etc*. Despite the power of these state-of-the-art models, they are found to be extremely unstable to input perturbations [17, 22, 12]. This fragility can be exploited by crafting adversarial examples, which are optimized to manipulate networks while appearing innocuous to humans.

One popular method to mitigate the brittleness of neural networks is adversarial training [25, 14, 20], in which network parameters are updated using adversarially perturbed images. This produces a hardened network that is robust to adversarial perturbations in the pixel space. While adversarial training is able to increase the robustness of classifiers, it often reduces accuracy on clean images at test time [51, 40, 2]. This accuracy loss is believed to be in part because of fundamental tradeoffs between accuracy and robustness [35, 40, 2], and in part because of domain shift between clean and adversarial image distributions [45, 27, 10].

In light of this, we propose to align features of clean images with their adversarially perturbed counterparts for improved clean accuracy and robustness. Our approach is
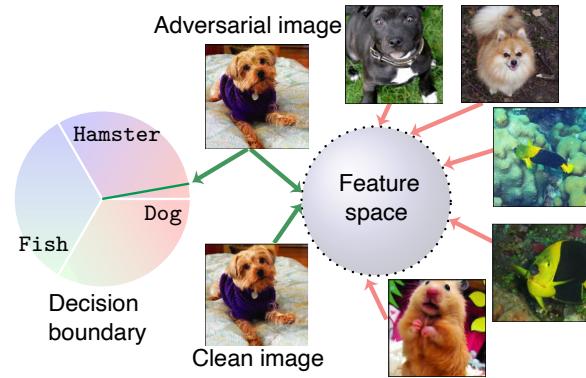


Figure 1: **A conceptual overview of our framework.** We train a neural network with two heads, each with its own loss function. The first loss term acts on the "feature extraction head," and promotes high feature similarity between an adversarial image and its clean counterpart. The second loss term acts on the "classification head," and enforces that adversarial images receive correct class labels.

inspired by recent self-supervised learning frameworks [5, 15], in which two random crops of the same image are considered as a positive pair and their distance in feature space is minimized to learn generic feature representations. Here, a clean image and its adversarial copy form a positive pair, and by contrasting it with other pairs, we modulate standard adversarial training to prevent the adversarial feature distribution from drifting away from the natural distribution.

With this in mind, we introduce Two Head Adversarial Training (THAT), an adversarial training framework that uses multiple training objectives to boost robustness. THAT adversarially trains a robust encoder using two parallel heads, a feature head and a classification head, both sitting on a shared ResNet backbone. For each training step, clean images are passed through a naturally trained clean network to produce features. Then, adversarial examples are made by attacking the classification head of the robust encoder. A loss function is then computed that contains two terms: a contrastive loss that promotes similarity between the natural features from the clean encoder and the adversarial features from the robust encoder, and a classi-

fication loss from the classification head. The former term enforces that feature representations of adversarial images are aligned with natural feature representations (*i.e.*, domain shift is minimized), while the second term ensures that the feature representations contain the information needed for accurate classification. The robust encoder is then updated to minimize the combined loss.

At test time, we consider two different classification modes for defense. In addition to using a standard classification head, we also consider a nearest-neighbor-based classification procedure that relies on the feature extraction head.

We conduct extensive experiments on ImageNet [8] using standard adversarial training updates, and also using accelerated (*a.k.a*, "free" [36]) adversarial updates. In both settings, we demonstrate that THAT outperforms other adversarial training frameworks in terms of both clean accuracy and robustness with different backbones.

## 2. Background and Related Work

**Adversarial robustness.** To mitigate threats posed by adversarial examples [29, 38], adversarial training [14, 25, 36, 43, 11, 50, 47] solves a min-max optimization problem in which adversarial examples are crafted to maximize the training loss, and these examples are then used to update network parameters during loss minimization. This process can be interpreted as approximately solving the saddle-point optimization problem:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(X,y)\sim\mathcal{D}} \left[ \max_{\|\delta\|_\infty < \epsilon} \ell(f(X + \delta), y; \boldsymbol{\theta}) \right], \quad (1)$$

where $X$ is a clean image drawn from the training set $\mathcal{D}$ with $y$ as its label, $f$ represents a model parameterized by weights $\boldsymbol{\theta}$, $\delta$ is the adversarial perturbation, and $\ell$ denotes the cross-entropy loss.

Early adversarial training methods crafted adversarial examples with the Fast Gradient Sign method (FGSM) [14], but this approach results in networks that are easily broken by multi-step attacks. Madry *et al*. use project gradient descent (PGD) to generate adversarial examples to improve robustness for multi-step attacks [25]. PGD performs gradient ascent on input images in the signed gradient direction with respect to the classification loss, and then clips the perturbation to enforce an $\ell_\infty$-norm constraint. Kannan *et al*. further regularize adversarial training by penalizing the difference between logits from clean images and their adversarial variants [20]. Xie *et al*. introduce denoising blocks that use self-attention in feature maps to improve adversarial training [46]. Balaji *et al*. improve adversarial training for ImageNet by imposing different robustness criteria for each training sample [2]. The TRADES method achieves high levels of robustness by training with a loss that pro-

motes similarity between the predicted label scores for natural and adversarial images [51]. Due to the huge computational overhead of adversarial training methods, Shafahi *et al*. introduce a "free" adversarial training strategy that updates model weights and generates adversarial examples by replaying mini-batches [36].

In our work, we build upon standard and free adversarial training settings and demonstrate that explicitly contrasting features improves both clean accuracy and robustness. Dubey [10] *et al*. show that nearest neighbor search can be used for deflecting strong attacks, but this requires an extremely large gallery of reference images ($\sim$ 1 billion) to achieve good performance. Instead, we show that our network is able achieve strong defense via nearest neighbors using the much smaller ImageNet dataset. In addition, Mao *et al*. use a triplet loss for better adversarial training [27], which requires manual hard negative mining. Our contrastive loss automatically performs hard negative mining with a number of negative samples.

**Representation learning with contrastive losses.** Extensive studies have been conducted to learn visual representations in a self-supervised manner [30, 4, 23, 28, 49, 13, 9]. Among these approaches, contrastive learning methods [44, 15, 5, 39] currently achieve state-of-the-art results by maximizing the agreement of positive pairs (two random crops of the same image) relative to a large number of negative pairs. Wu *et al*. perform instance discriminative tasks based on the entire ImageNet training set [44]. MoCo introduces a memory bank to maintain consistent representations of negative samples with the help of a momentum encoder [15]. In this paper, we focus on improving the robustness of neural networks with contrastive losses. We consider a clean image and its adversarially perturbed variant as a positive pair rather than two stochastic data augmentations.

**Adversarial training and self-supervised learning.** There are some very recent studies exploring adversarial training with self-supervised learning [18, 19, 21, 6]. Chen *et al*. use several self-supervised losses like predicting rotations, permutations, to pretrian a ResNet and study its robustness [6]. [18, 19, 21, 6] focus on better pretraining with adversarial examples for downstream tasks on small datasets. In contrast, we use a contrastive branch to align features of clean and adversarial images to boost the performance of standard adversarial training. In addition, we also target at large datasets for large-scale adversarial training.

**Theoretical motivation and relation to TRADES.** A successful and well-known training strategy for adversarial defense is TRADES [51]. This method is related to THAT in that it can be interpreted as applying a regularization to compare logits (as opposed to features). However, while TRADES has been highly successful at robust optimization for CIFAR-10 and MNIST, which have relatively few
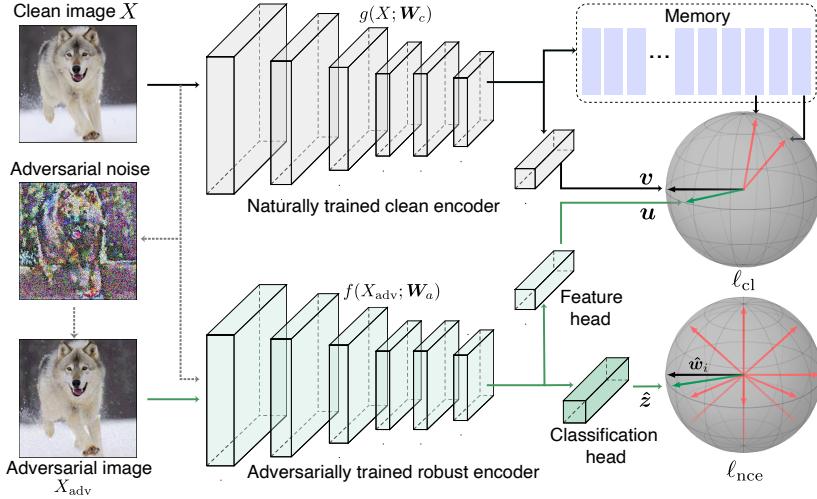
Figure 2: **An overview of the framework.** Given a pair containing a clean image and its adversarial variant, features from both images are mapped to locations on a hypersphere. After robust optimization, both feature vectors lie in close proximity. At the same time, the classification logits from the clean and adversarial image are projected onto a sphere, and forced to lie close to the ground truth label during training. See texts for more details.

classes, we have found the TRADES objective to lead to unstable and non-convergent behavior when used on ImageNet. This problem was also observed in [33]. Here, we explore theoretical reasons for the failure of TRADES and discuss how these weaknesses motivate the use of THAT.

TRADES is based on the decomposition $\mathcal{R}_{rob} = \mathcal{R}_{nat} + \mathcal{R}_{bnd}$, which represents the robust accuracy of a model as the sum of the natural risk (*i.e.*, the 0-1 loss), and the "boundary" risk. The boundary risk is the probability that, for an input $X$, there exists an input point $X'$ within a ball around $X$ such that $X$ and $X'$ are assigned different labels. These terms are then upper bounded by a smooth loss and minimized. The multi-class TRADES objective is:

$$\min_{\boldsymbol{\theta}} \quad \mathbb{E}_{(X,y)\sim\mathcal{D}} \; \ell(f(X), \, y; \, \boldsymbol{\theta}) \qquad (2)$$

$$+ \, \lambda \max_{\|\delta\|_\infty < \epsilon} \mathcal{L}(f(X + \delta), f(X); \, \boldsymbol{\theta}), \qquad (3)$$

where $\mathcal{L}$ is a loss function that measures the discrepancy between class label distributions predicted for $X$ and its adversarial example $X + \delta$, and $\lambda$ is a scalar parameter.

For 2-class problems, the loss (2) is a tight convex relaxation of the non-convex 0-1 loss whenever $\mathcal{L}$ is a calibrated loss function that satisfies several weak assumptions [51]. However, the tightness result for TRADES is specific to binary classification; the relaxation is no longer tight for multi-class problems, and becomes extremely pathological for many-class problems like ImageNet.

To understand why this is, observe that TRADES uses the KL divergence for $\mathcal{L}(\cdot, \cdot)$ when solving multi-class problems. If $p$ and $q$ are the class distributions predicted by the network for $X$ and its adversarial example, respectively,

then the TRADES loss contains the KL-divergence term:

$$\mathcal{L}(f(X + \delta), f(X); \, \boldsymbol{\theta}) = \sum_i p_i \log(p_i/q_i), \qquad (4)$$

where the sum in Eqn. (4) is over class indices. For typical ImageNet images, the distributions $p$ and $q$ contain just a few relevant/large class probabilities, and then a "long tail" with nearly 1000 small probabilities that are generally discarded when making classification decisions. The KL loss used by TRADES sums over all of these class labels, and its value becomes highly dominated by the many contributions from classes with small probabilities. Furthermore, the smallest values of $q_i$, which appear in the denominator of Eqn. (4), make the loss most sensitive to unlikely classes. In practice, we often ignore such class labels with near-zero probability, and yet they account for a dominant share of the TRADES objective.

Our proposed method is similar in spirit to TRADES — it enforces similarity between network outputs for natural and adversarial images. But unlike TRADES, our proposed method uses a contrastive loss that is computed using the inner product between feature vectors, and is dominated by large entries in feature representations. This enables THAT to do contrastive learning while avoiding the pathological behaviors that TRADES suffers in the many-class regime.

## 3. THAT: Two Head Adversarial Training

In this section, we introduce THAT, a contrastive framework with a naturally trained clean encoder and an adversarially trained robust encoder. The method learns feature

representations that are robust to attacks while simultaneously achieving high classification accuracy. Each of these objectives is achieved with a loss function on a separate head of the network. We describe each component of our framework here.

**Naturally trained clean encoder.** The natural/clean encoder $g(X; \boldsymbol{W}_c)$, parameterized by weights $\boldsymbol{W}_c$, takes a clean image $X$ as inputs. Following [7], it uses a two-layer projection head on top of feature maps from the Res_5 stage of a ResNet to produce a 128D feature representation $\boldsymbol{v}$. The feature vector is further normalized with an $\ell_2$ norm. We set the weights of the clean encoder from a pre-trained self-supervised model [7]. The weights of this model remain frozen while the robust model is trained.

**Adversarially trained robust encoder.** The robust encoder $f(X_{\text{adv}}; \boldsymbol{W}_a)$, parameterized by weights $\boldsymbol{W}_a$, is a two-head architecture with a feature head and a classifier head. In particular, $\boldsymbol{W}_a = \{\boldsymbol{W}_{ab}, \boldsymbol{W}_{af}, \boldsymbol{W}_{ac}\}$, where $\boldsymbol{W}_{ab}$, $\boldsymbol{W}_{af}$ and $\boldsymbol{W}_{ac}$ represent the weight matrices for the base ResNet model, a feature head with two fully-connected layers, and the classifier head with one fully-connected layer, respectively. Given a clean image $X$, the robust encoder first generates its adversarially perturbed version $X_{adv} = X + \delta$ on-the-fly by attacking the classification head using multiple steps of PGD. The feature head then computes a 128 dimensional feature vector $\boldsymbol{u}$ for the perturbed image, which is normalized to have unit $\ell_2$ norm. At the same time, the separate classifier head produces a logit vector $\boldsymbol{z}$ for use in the classification loss.

**The contrastive loss.** The feature extraction head provides adversarial features $\boldsymbol{u}$ for the attack image, and the clean network provides a clean image representation $\boldsymbol{v}$. These two representations form a "positive pair" and should be aligned by the contrastive loss. We also form "negative pairs" by comparing $\boldsymbol{u}$ with representations from randomly selected clean images, which could be sampled from the same mini-batch or from an external memory bank. Here, we use a memory bank as it is has been demonstrated that using a large number of negative pairs is beneficial [15, 44]. Finally, we form the contrastive loss function [37, 31]:

$$\ell_{\text{cl}}(\boldsymbol{u}) = -\log \frac{\exp(\boldsymbol{u}^T \boldsymbol{v}/\tau)}{\exp(\boldsymbol{u}^T \boldsymbol{v}/\tau) + \sum_{\boldsymbol{v} \in \mathcal{V}_-} \exp(\boldsymbol{u}^T \boldsymbol{v}/\tau)}, \tag{5}$$

where $\mathcal{V}_-$ is the set of features from random negative samples. The parameter $\tau$ is the "temperature," controlling the sharpness of the distribution. The contrastive loss forces the adversarial representation $\boldsymbol{u}$ to be closer to its own clean base image than other images. This suppresses domain shift between the clean and adversarial images in feature space.

**Classifier loss.** The contrastive learning loss forces feature representations to be invariant to attacks, but it does

---

**Algorithm 1** Pseudocode of our approach in PyTorch style.

```
# g: naturally trained clean encoder
# f: adversarially trained robust encoder
# eps: adversarial perturbation epsilon
# K: number of steps for PGD
# mem: memory bank with clean image features

for x, y in loader: # x: data, y: labels
    # generate adversarial examples with K-step PGD
    x_adv = PGD_attack(f, x, y, K, eps)

    # compute features for clean and adversarial images
    feat_clean = g.forward(x)
    feat_adv, logits_adv = f.forward(x_adv)

    # compute contrastive losses
    loss_cl = cl_loss(feat_adv, feat_clean, mem)

    # compute classification losses
    loss_cls = nce_loss(logits_adv, y)

    loss = loss_cls + loss_cl
    loss.backward()
    optimizer.step()
```

---

not measure classification performance of these representations. To get good classification performance, the classifier head of the robust network produces its own training loss. For the classifier head, we use a "normalized" cross entropy [41, 32, 44], which measures the disparity between output logits and ones-hot vectors using a contrastive loss. This keeps gradient scaling and training dynamics of the classifier loss similar to the contrastive loss on the features. The normalized cross entropy loss is:

$$\ell_{\text{nce}}(\boldsymbol{z}) = -y_i \log \frac{\exp(\hat{\boldsymbol{z}}^T \hat{\boldsymbol{w}}_i/\eta)}{\sum_{i=1}^{C} \exp(\hat{\boldsymbol{z}}^T \hat{\boldsymbol{w}}_i/\eta)}, \tag{6}$$

where $\hat{\boldsymbol{z}}$ is the normalized logits based on $\boldsymbol{z}$, $\hat{\boldsymbol{w}}_i$ is the $i$ the column in $\boldsymbol{W}_{ac}$ representing the normalized weights for the $i$-th class, $y_i \in \mathbb{R}^{\{0,1\}}$ is the label for the $i$-th class, and $C$ is the total number of classes in the dataset. $\eta$ is a learnable parameter to control the sharpness of the distribution. Here, $\hat{\boldsymbol{w}}_i$ is considered as the class prototype, and we are forcing the logit vectors to lie close to it. This is similar in spirit to Eqn (5), in which features from an adversarial example are mapped to be close to features of its clean twin.

Finally, the combined training objective of THAT can be written as:

$$\min_{\boldsymbol{W}_a} \mathbb{E}_{(X,y) \sim \mathcal{D}} \, \ell_{\text{nce}}(f(X + \delta), y; \boldsymbol{W}_a) +$$
$$\ell_{\text{cl}}(g(X), f(X + \delta); \boldsymbol{W}_a) \tag{7}$$
$$\text{where} \quad \delta = \max_{\|\delta\|_\infty < \epsilon} \ell_{\text{nce}}(f(X + \delta), y; \boldsymbol{W}_a).$$

**Defense strategies.** Once THAT is trained, it is able to defend against strong PGD attacks during testing. We experiment with two different classification modes for defense during testing: (i) standard softmax-based defense using outputs from the classification head of the robust encoder; (ii) nearest-neighbor based defense using features from the

|          | Clean           | PGD-10          | PGD-30          | PGD-200         | PGD-1000        |
|----------|-----------------|-----------------|-----------------|-----------------|-----------------|
| **R50**  |                 |                 |                 |                 |                 |
| Standard AT | 50.81        | 47.78           | 39.31           | 38.09           | 37.74           |
| Ours     | 53.29 (+2.48)   | 49.62 (+1.84)   | 41.01 (+1.70)   | 40.22 (+2.13)   | 39.59 (+1.85)   |
| **R101** |                 |                 |                 |                 |                 |
| Standard AT | 56.21        | 51.40           | 42.68           | 41.08           | 40.86           |
| Ours     | 58.17 (+1.96)   | 52.70 (+1.30)   | 44.02 (+1.34)   | 42.79 (+1.71)   | 42.34 (+1.48)   |
| **R152** |                 |                 |                 |                 |                 |
| Standard AT | 57.61        | 52.13           | 43.86           | 42.73           | 42.14           |
| Ours     | 60.38 (+2.77)   | 54.69 (+2.56)   | 45.67 (+1.81)   | 44.51 (+1.78)   | 44.22 (+2.08)   |

Table 1: **Results and comparisons of our method with standard adversarial training** using different backbone networks.

feature head of the robust encoder. Below we introduce the nearest-neighbor classifier in detail.

The nearest-neighbor classifier computes feature representations for all training samples with the clean encoder and stores them in a memory bank $M_{\text{train}}$. Given a test image $\hat{X}$, we first compute its feature $u_t$ through the robust feature head. Then, the embedding is compared with those of all training samples to retrieve the top-$k$ nearest neighbors. $N_k$ similarity scores are computed based on dot products, and a weighted average of neighbor labels is computed. More precisely, the confidence of $\hat{X}$ belonging to the $c$-th class is defined as:

$$P(c|\hat{X}) = \sum_{i=1}^{N_k} u_t^T u^i \cdot y^i(c), \qquad (8)$$

where $u^i$ is the feature for $i$-th neighbor sample, and $y^i$ is the corresponding one-hot label.

## 4. Experiments

**Datasets and metics.** We evaluate our framework on the ImageNet classification benchmark [8], which has $\sim$1.28 million images annotated into 1000 classes. We consider two adversarial training settings: (1) standard adversarial training (STANDARD AT) [25, 20] using $K$-step PGD to generate adversarial examples. This method increases model robustness but is computationally expensive (*i.e.*, $K$ times slower compared to natural image training). In this setting, for both training and testing, we consider *targeted attacks* by randomly selecting a targeted class uniformly following [46]; (2) The "free" adversarial training method [36] (FREE AT), which speeds up standard adversarial training by updating model weights and generating adversarial examples at the same time with mini-batch replay. Following [36], we consider *untargeted attacks* for both training and testing. We report top-1 classification accuracy on the 50k ImageNet validation data using both clean images and adversarially perturbed images and with

many-step PGD attacks as in [46, 20, 1].

**Implementation details.** We adopt Pytorch for implementation. Since adversarial training on ImageNet is expensive, we use distributed training with synchronized SGD. In particular, for standard adversarial training, we set the maximum perturbation for each pixel to $\epsilon = 16$, the step size to $\alpha = 4$, and the number of attack iterations to $K = 10$. We found that this achieves similar results (*i.e.*, robustness and clean accuracy) compared to using 30 attack iterations with a step size of 1, but can reduce training time by 3$\times$. We use a batch of size 4096 on 32 Tesla V100 32GB GPUs and train for 100 epochs as in [46]. The initial learning rate is set to 1.6, and is decayed by a factor of 10 at the 35, 60, and 90 epoch. For fast adversarial training, we set $\epsilon = 4$ and set the number of replays to 4, following [36]. We use a batch size of 2048 during training and train for 90 epochs (effectively 23 epochs with replay [36]). The clean encoder has the same backbone as the robust encoder, but is initialized from pre-trained self-supervised models [7] and fixed during training. We also show the clean encoder can be trained in Section 4.3.

### 4.1. Classifier-based Defense

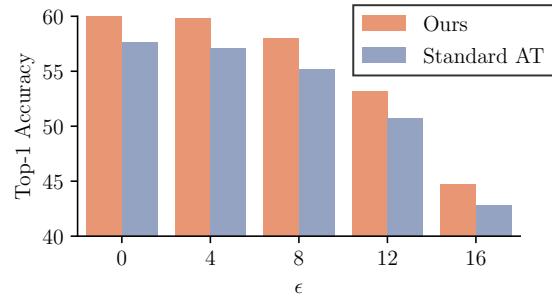**Standard adversarial training.** Table 1 summarizes the results of our approach and comparisons with standard



Figure 3: **Results of different** $\epsilon$ using a ResNet152 architecture, evaluated with PGD-200.

adversarial training (standard AT) using three different backbones, *i.e.* ResNet50 (R50), ResNet101 (R101), and ResNet152 (R152). We can see from the table that compared to standard AT, our method with a R152 backbone achieves a clean accuracy of 60.38% and an accuracy of 44.22% when evaluated with PGD-1000, offering 2% (absolute percentage points) compared to the standard AT baseline. We observe similar trends for both R50 and R101, confirming the generalization of our approach with different backbone networks. In addition, the performance of network models degrades when more attack iterations are used and tends to stabilize after 200 iterations. Comparing across different backbones, we can see that models with larger capacity perform better in terms of both clean accuracy and robustness, as observed in [47]. We also evaluate both our approach and standard AT with different maximum perturbation values (*i.e.* $\epsilon$). The results are shown in Figure 3. We see that our approach clearly outperforms the standard baseline model with different perturbation values.

Furthermore, we compare with the following state-of-the-art models on ImageNet: (1) ALP [20], which penalizes the outputs of clean images to be similar to those of adversarial images with a mean-squared error loss; (2) Feature Denoising [46], which adds non-local blocks [42] after each residual block in ResNet models ; (3) MBN-ALP [47], which reimplements the ALP algorithm by using different batch normalization statistics for clean and adversarial images. Results are summarized in Fig. 4.

We see from the figure that THAT achieves better results than alternative methods when strong attacks are presented at test time (*i.e.*, the number of attack iterations is greater than 200). Compared to Feature Denoising [46] which performs self-attention on feature maps, THAT offers 1.4% gain when evaluated against PGD-1000 and is slightly worse with PGD-10 attacks. We also experiment with non-local blocks and observe that they can slightly improve the performance for weak attacks and clean accuracy. However, adding non-local blocks makes training more computationally expensive as it flattens all the pixels into a huge vector to compute a dense graph. In addition, we also compare our R152 baseline with that in [46], and our method achieves better robustness against strong attacks while being 0.4% worse against PGD-10. Furthermore, both THAT and Feature Denoising outperform ALP and its variant by clear margins.

**Accelerated training results.** We now experiment with an accelerated adversarial training scheme to verify that THAT is compatible with different strategies. In particular, we train the proposed architecture on top of the "free" training framework introduced in [36]. The results are summarized in Table 2. Similarly to adversarial training, we observe THAT offers significant performance gains for both clean accuracy and robustness compared to the baseline method.
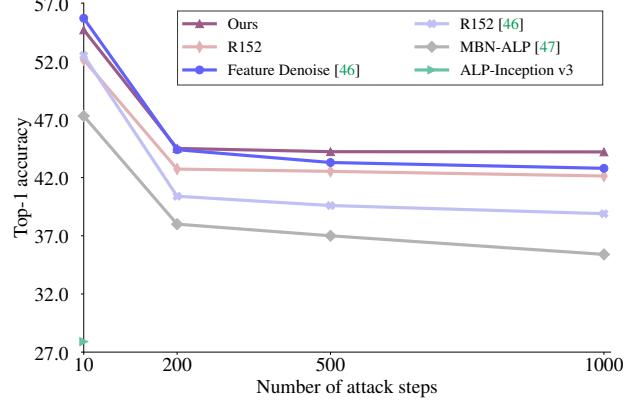


Figure 4: **Comparisons of THAT with state-of-the-art models** for standard adversarial training on ImageNet.

|  | Clean | PGD-10 | PGD-100 |
|---|---|---|---|
| **R50** | | | |
| Free [36] | 60.21 | 32.77 | 31.82 |
| Ours | 63.01 (+2.8) | 34.25 (+1.48) | 33.26 (+1.44) |
| **R101** | | | |
| Free [36] | 63.34 | 35.38 | 34.32 |
| Ours | 67.22 (+3.88) | 38.56 (+3.18) | 37.17 (+2.85) |
| **R152** | | | |
| Free [36] | 64.45 | 36.99 | 35.99 |
| Ours | 68.82 (+4.37) | 39.83 (+2.84) | 38.27 (+2.28) |

Table 2: **Results and comparisons with the "free" adversarial training [36]**. Here, the PGD attacks are untargeted.

In particular, R152 offers a 4% gain for clean accuracy and 2.28% improvement against PGD-100. Note that the results in Table 2 are not directly comparable to Table 1 since the attacks in the "free" setting are untargeted. In addition, we only evaluate 100 attack iterations as in [36] since the performance stabilizes.

### 4.2. Nearest neighbor classification

We demonstrate that features computed from the feature head of the robust encoder are able to facilitate nearest-neighbor based defense against strong PGD attacks. We compare THAT with a R152-CON model, which augments a R152 model with a feature head for contrastive learning. R152-Con is trained on clean images without adversarial training, offering an accuracy of 78.89% on the clean ImageNet validation set.

Table 3 summarizes the results of THAT and R152-Con using the top-50 nearest neighbors to classify test images. Interestingly, R152-Con offers a 15% top-1 accuracy against PGD-10 with nearest neighbor based defense, even though the model is not adversarially trained. This is much

| Method   | PGD-10 | PGD-30 | PGD-500 |
|----------|--------|--------|---------|
| R152-Con | 15.08  | 11.35  | 10.71   |
| Ours     | 35.38  | 29.24  | 27.67   |

Table 3: **Results and comparisons** of defending with nearest neighbors.

better than the 0.66% accuracy of its clean-trained softmax classifier, which indicates that nearest neighbor methods are indeed more resistant to PGD than softmax. THAT achieves 35.4% accuracy against PGD-10, which greatly exceeds naturally trained R152-Con. To compare to a KNN based robust model, we report results from [10] which similarly performs defense with nearest neighbors, but the features are used for classification rather then explicitly designed for contrasting learning. For fair comparisons we follow the setting in [10], which constrains perturbations within an $\ell_2$ ball instead of the (more standard) $\ell_\infty$ ball used in other sections of this paper. The results are summarized in Table 4. Using the same set of images for retrieval (*i.e.*, ImageNet training set with 1.28 million images), THAT outperforms [10] by 11.8% in top-1 accuracy, highlighting the effectiveness of features from the feature head for defense. In addition, with a R152 model, THAT achieves performance comparable to Dubey's KNN defense with one billion training images available for retrieval. Note that we did not retrain our network against $L2$ threat models.

| Method                    | Top-1 accuracy (PGD-10) |
|---------------------------|-------------------------|
| IG-1B-R50 [10]            | 46.2                    |
| ImageNet-1.3M (Ours-R152) | 45.6                    |
| ImageNet-1.3M-R50 [10]    | 23.5                    |
| IMageNet-1.3M (Ours-R50)  | 35.3                    |

Table 4: **Comparisons with state-of-the-art nearest neighbor defense.** Here, IG-1B denotes the dataset with 1 billion images from Instagram [26].

We also show qualitatively in Figure 5 both success (top two rows) and failure (bottom two rows) cases of randomly selected samples with the KNN defense. For both success and failure cases, we see that retrieved samples are indeed visually similar to the query image although the query image is perturbed with adversarial noise. The incorrect predictions for failure cases are largely due to the existence of fine-grained classes, rather than a robustness failure of the embedding. For example, in the third row of Fig. 5, different kinds of cats are retrieved but they belong to different species other than the "siamese cat".
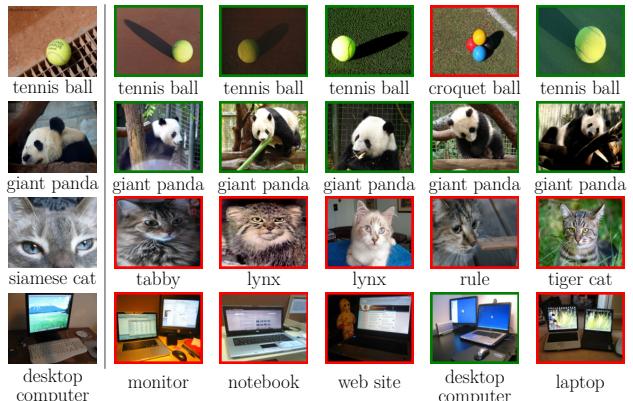


Figure 5: **Top-5 nearest neighbors retrieved given an adversarially perturbed query image.** Left: A query image which has been adversarially perturbed with PGD-10. Right: top-5 nearest neighbors of the query image from the training set of ImageNet.

### 4.3. Discussion

We conduct a set of experiments to analyze different components of THAT, and discuss the results.

**Ablation on losses of THAT.** We show results of THAT with different loss functions in Table 5. We observe that compared to our full framework, removing either the contrastive loss or the normalized softmax degrades the performance slightly, yet the resulting models still outperform the standard adversarial training baseline by at least 1% in clean accuracy and 0.6% against PGD-1000. This highlights the importance of both components for adversarial training. Furthermore, without the contrastive learning branch, our framework produces 58.92% and 42.78% accuracy when evaluated on clean images and against PGD-1000 attacks, respectively. The result is slightly worse compared to removing the normalized softmax (*i.e.*, a standard cross-entropy loss is used in the pipeline), suggesting that contrastive learning is relatively important. Note that we did not compare to TRADES [51] because we found it to be unstable and non-convergent, even after hyper-parameter searching. Similar attempts and failures to train TRADES on ImageNet are mentioned in [33]. See Section 2. As a workaround, we compare with standard adversarial training with KL divergence loss to force probability distributions from clean images to be close to those of adversarial images [1]. We see that it achieves high clean accuracy but is extremely vulnerable to strong PGD attacks. This is likely because of the strong emphasis that the KL loss puts on unlikely class labels (see Section 2).

---

[1]Note that this is different from TRADES as TRADES maximizes the KL divergence to generate adversarial examples, while we instead maximize classification loss as in our approach.

| Method | Clean | PGD-10 | PGD-1000 |
|---|---|---|---|
| Standard AT | 57.61 | 52.13 | 42.14 |
| Standard AT + KL | 71.05 | 14.13 | 0.00 |
| Ours w.o. CL | 58.92 | 53.32 | 42.78 |
| Ours w.o. NCE | 59.64 | 53.80 | 42.77 |
| Ours | 60.38 | 54.69 | 44.22 |

Table 5: **Ablating different components of THAT**.

**Number of negative samples for contrastive learning.**
Self-supervised learning methods suggest the number of
negative samples used for contrastive learning is important.
Therefore, we analyze the performance of THAT using dif-
ferent numbers of examples in the memory for contrastive
learning. The results are summarized in Table 6. We see
that increasing the number of samples in the memory is in-
deed beneficial for improving clean accuracy and robust-
ness. With 65536 samples in the memory, THAT offers a
clean accuracy of 60.68%, outperforming standard adver-
sarial training by 3%. A memory size of 32768 offers the
best trade-off between clean accuracy and robustness.

| # Samples | Clean | PGD-10 | PGD-1000 |
|---|---|---|---|
| 0 | 57.61 | 43.86 | 42.14 |
| 4096 | 59.77 | 45.40 | 43.15 |
| 8192 | 58.25 | **45.89** | 43.79 |
| 16384 | 60.26 | 45.71 | 43.59 |
| 32768 | 60.38 | 45.67 | **44.22** |
| 65536 | **60.68** | 45.69 | 43.43 |

Table 6: **Results of THAT** using different number of nega-
tive samples.

**Clean encoder.** Instead of freezing the weights of the clean
encoder during training, we also experiment with updat-
ing its weights to reflect the parameters of the robust en-
coder with momentum to ensure consistent representations
in memory [15]. The results are presented in Table 7.
This modified implementation that updates both networks
at once clearly beats standard adversarial training in terms
of robustness, and slightly in terms of clean accuracy. How-
ever our proposed framework beats both methods in terms
of clean accuracy and robustness.

**Loss surface visualization.** Figure 6 visualizes the loss sur-
face of a selected sample with both THAT (left side) and
standard adversarial training (right side). On the top row,
we show the cross-entropy loss projected on one random
(Rademacher) and one adversarial direction. On the bot-
tom, we project the loss along two random directions. We

| Method | Clean | PGD-10 | PGD-1000 |
|---|---|---|---|
| Standard AT | 57.61 | 52.13 | 42.14 |
| Ours-MoCo | 57.76 | 53.36 | 43.03 |
| Ours | 60.38 | 54.69 | 44.22 |

Table 7: **Results of updating the clean encoder** in a
MoCo [15] fashion.

see that the loss surface for THAT is more flat than standard
adversarial training. Since both models are adversarially
trained, the loss does not increase along the gradient direc-
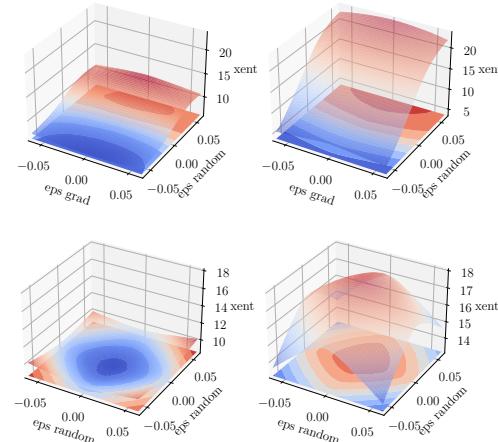tion. See more examples in the supplemental material.



Figure 6: **Loss surface visualization** of THAT (left) and
standard AT (right).

## 5. Conclusion

We presented THAT, a two-stream contrastive learning
framework for improved robustness and clean accuracy.
THAT is trained using two loss functions; one to align the
feature distributions between natural and adversarial im-
ages, and one to promote good classification accuracy. The
resulting model is able to defend against strong adversarial
attacks at test time not only using the hardened classifier but
also using a KNN search. Through extensive experiments,
we demonstrate THAT achieves better results than alterna-
tive methods for ImageNet under a wide range of settings.

## References

[1] Anish Athalye, Nicholas Carlini, and David Wagner. Obfus-
cated gradients give a false sense of security: Circumventing
defenses to adversarial examples. In *ICML*, 2018. 5

[2] Yogesh Balaji, Tom Goldstein, and Judy Hoffman. Instance
adaptive adversarial training: Improved accuracy tradeoffs in
neural nets. *arXiv preprint arXiv:1910.08051*, 2019. 1, 2

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1

[4] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, 2018. 2

[5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 1, 2

[6] Tianlong Chen, Sijia Liu, Shiyu Chang, Yu Cheng, Lisa Amini, and Zhangyang Wang. Adversarial robustness: From self-supervised pre-training to fine-tuning. In *CVPR*, 2020. 2

[7] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 4, 5

[8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 2, 5

[9] Alexey Dosovitskiy, Philipp Fischer, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE TPAMI*, 2016. 2

[10] Abhimanyu Dubey, Laurens van der Maaten, Zeki Yalniz, Yixuan Li, and Dhruv Mahajan. Defense against adversarial images using web-scale nearest-neighbor search. In *CVPR*, 2019. 1, 2, 7

[11] Logan Engstrom, Andrew Ilyas, and Anish Athalye. Evaluating and understanding the robustness of adversarial logit pairing. *arXiv preprint arXiv:1807.10272*, 2018. 2

[12] Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A Wichmann. Generalisation in humans and deep neural networks. In *NeurIPS*, 2018. 1

[13] Spyros Gidaris, Andrei Bursuc, Nikos Komodakis, Patrick Pérez, and Matthieu Cord. Learning representations by predicting bags of visual words. In *CVPR*, 2020. 2

[14] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015. 1, 2

[15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 1, 2, 4, 8

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1

[17] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *ICLR*, 2019. 1

[18] Chih-Hui Ho and Nuno Vasconcelos. Contrastive learning with adversarial examples. *arXiv preprint arXiv:2010.12050*, 2020. 2

[19] Ziyu Jiang, Tianlong Chen, Ting Chen, and Zhangyang Wang. Robust pre-training by adversarial contrastive learning. *arXiv preprint arXiv:2010.13337*, 2020. 2

[20] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint arXiv:1803.06373*, 2018. 1, 2, 5, 6

[21] Minseon Kim, Jihoon Tack, and Sung Ju Hwang. Adversarial self-supervised contrastive learning. *arXiv preprint arXiv:2006.07589*, 2020. 2

[22] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *ICLR Workshops*, 2017. 1

[23] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *ECCV*, 2016. 2

[24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 1

[25] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018. 1, 2, 5

[26] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018. 7

[27] Chengzhi Mao, Ziyuan Zhong, Junfeng Yang, Carl Vondrick, and Baishakhi Ray. Metric learning for adversarial robustness. In *NeurIPS*, 2019. 1, 2

[28] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *CVPR*, 2020. 2

[29] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *CVPR*, 2016. 2

[30] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 2

[31] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 4

[32] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *CVPR*, 2018. 4

[33] Chongli Qin, James Martens, Sven Gowal, Dilip Krishnan, Krishnamurthy Dvijotham, Alhussein Fawzi, Soham De, Robert Stanforth, and Pushmeet Kohli. Adversarial robustness through local linearization. In *NeurIPS*, 2019. 3, 7

[34] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1

[35] Ali Shafahi, W Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. Are adversarial examples inevitable? In *ICLR*, 2019. 1

[36] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *NeurIPS*, 2019. 2, 5, 6

[37] Kihyuk Sohn. Improved deep metric learning with multiclass n-pair loss objective. In *NeurIPS*, 2016. 4

[38] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014. 2

[39] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *eccv*, 2020. 2

[40] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *ICLR*, 2019. 1

[41] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In *CVPR*, 2018. 4

[42] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 6

[43] Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. In *ICLR*, 2020. 2

[44] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018. 2, 4

[45] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan L Yuille, and Quoc V Le. Adversarial examples improve image recognition. In *CVPR*, 2020. 1

[46] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan Yuille, and Kaiming He. Feature denoising for improving adversarial robustness. In *CVPR*, 2019. 2, 5, 6

[47] Cihang Xie and Alan Yuille. Intriguing properties of adversarial training at scale. In *ICLR*, 2020. 2, 6

[48] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 1

[49] Xueting Yan, Ishan Misra, Ishan, Abhniav Gupta, Deepti Ghadiyaram, and Dhruv Mahajan. ClusterFit: Improving generalization of visual representations. In *CVPR*, 2020. 2

[50] Haichao Zhang and Jianyu Wang. Defense against adversarial attacks using feature scattering-based adversarial training. In *NeurIPS*, 2019. 2

[51] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. In *ICML*, 2019. 1, 2, 3, 7