
Rethinking Robust Contrastive Learning from the Adversarial Perspective

Fatemeh Ghofrani^{*1} Mehdi Yaghouti^{*1} Pooyan Jamshidi¹

Abstract

To advance the understanding of robust deep learning, we delve into the effects of adversarial training on self-supervised and supervised contrastive learning alongside supervised learning. Our analysis uncovers significant disparities between adversarial and clean representations in standard-trained networks across various learning algorithms. Remarkably, adversarial training mitigates these disparities and fosters the convergence of representations toward a universal set, regardless of the learning scheme used. Additionally, increasing the similarity between adversarial and clean representations, particularly near the end of the network, enhances network robustness. These findings offer valuable insights for designing and training effective and robust deep learning networks. Our code is released at <https://github.com/softsys4ai/CL-Robustness>.

1. Introduction

Self-supervised learning has significantly improved in recent years, leading to state-of-the-art performance in various applications. While this paves the way to learn effective representations from massively available unlabeled data, the vulnerability to adversarial attack is still a fatal threat. Adversarial training is proven to be an effective defense method in supervised learning. This method can be interpreted as a min-max optimization problem (Madry et al., 2017), wherein the model parameters are updated iteratively by minimizing a training loss against the adversarial perturbations generated by maximizing an adversary loss function. While it is standard practice to use the same loss function

for both training and generating adversarial attacks, some works have explored the use of dissimilar loss functions to investigate robust training (Pal et al., 2021). Hendrycks et al. (Hendrycks et al., 2019) introduced a self-supervised term into the training loss to improve the robustness of a supervised model. Chen et al. (Chen et al., 2020b) were the first to apply adversarial training on a self-supervised model to achieve robust pre-trained encoders that can be used for downstream tasks through fine-tuning.

In recent years, the study of adversarial training on the robustness of various contrastive learning schemes has attracted great attention. The main idea of Contrastive Learning (CL) is to benefit from comparing semantically similar against dissimilar samples to learn the proper representations. Some recent works employed the generated adversarial examples as a similar match of the anchor data point to improve the robustness of the model. Kim et al. (Kim et al., 2020) were the first to utilize the contrastive loss to generate adversarial examples without any label for robustifying SimCLR (Chen et al., 2020a) framework. Moshavash et al. (Moshavash et al., 2021), Wahed et al. (Wahed et al., 2022), and Goyal et al. (Goyal et al., 2021) have applied the same technique to Momentum Contrast (MOCO) (He et al., 2020), Swapping Assignments between Views (SwAV) (Caron et al., 2020) and Bootstrap Your Own Latents (BYOL) (Grill et al., 2020), respectively. Fan et al. (Fan et al., 2021) introduced an additional regularization term in contrastive loss to enhance cross-task robustness transferability. They use a pseudo-label generation technique to avoid using labels in adversarial training of downstream tasks. Similarly, Jiang et al. (Jiang et al., 2020) considered a linear combination of two contrastive loss functions to study the robustness¹ under different pair selection scenarios.

One of the central steps in any contrastive learning scheme is the selection of positive and negative pairs. Without label information, a positive pair is often obtained by data augmentation, while the negative samples are randomly chosen from a mini-batch. However, this random selection strategy can lead to choosing the false-negative pairs when two samples are taken from the same class. In (Gupta et al.,

^{*}Equal contribution ¹Department of Computer Science and Engineering, University of South Carolina, Columbia, USA. Correspondence to: Fatemeh Ghofrani <ghofrani@email.sc.edu>, Mehdi Yaghouti <yaghouti@mailbox.sc.edu>.

2nd AdvML Frontiers workshop at 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

¹We use the word “robust” as shorthand for “adversarially robust” throughout the paper.

2022), Gupta et al. empirically demonstrate that the adversarial vulnerability of contrastive learning is related to employing false negative pairs during training. One remedy to this dilemma is to leverage the label information to extend the self-supervised contrastive loss into a supervised contrastive (SupCon) loss introduced by (Khosla et al., 2020). SupCon loss contrasts embeddings from the same class as positive samples against embeddings from different classes as negative samples. It has shown that SupCon outperforms cross-entropy loss in terms of accuracy and hyper-parameters stability (Khosla et al., 2020). Islam et al. (Islam et al., 2021) have conducted a broad study to compare the transferability of learned representations by cross-entropy, SupCon, and standard contrastive losses for several downstream tasks. Zhong et al. (Zhong et al., 2022) have designed a series of robustness tests, including data corruptions, ranging from pixel-level gamma distortion to patch-level shuffling and dataset-level distribution shift to quantify differences between contrastive learning and supervised learning frameworks.

Some very recent literature has started explaining and understanding robust networks. Jones et al. (Jones et al., 2022) have shown that irrespective of architecture or random initialization, adversarial robustness is a significant constraint on the learned function of a network. The research conducted by (Cianfarani et al., 2022) provides insights into the effects of adversarial training on representations. It emphasizes the lack of specialization in robust representations and the significant impact of overfitting on deeper layers during robust training. However, it is important to note that the primary focus of this study is on supervised learning algorithms. To address this research gap, our study investigates the effect of robust training on models trained using different learning schemes (contrastive, supervised-contrastive, and supervised) through a unified lens.

Contributions In this work, we conduct several comprehensive experiments to compare the robustness of contrastive and supervised contrastive with standard supervised learning under different training scenarios. Our research utilizes explanatory tools such as CKA (Centered Kernel Alignment) (Kornblith et al., 2019; Nguyen et al., 2020; Subramanian, 2021) and linear probing (Alain & Bengio, 2016) to investigate and contrast the layer-wise representations learned through various training methods. The design and implementation of the experiments are motivated by the following research questions:

- Q1: *Is there anything special about the learned representation with contrastive learning in terms of adversarial robustness?*
- Q2: *To what extent does employing the label information benefit or deteriorate the robustness of contrastive learning representations?*

Q3: *How does adversarial training affect (similarities and differences) the learned representations in supervised and contrastive learning?*

Our key findings can be summarized as follows:

- R1: Our results show that **contrastive learning without label information is less robust than other learning schemes in standard training**. However, combining the standard contrastive loss with either supervised cross-entropy or supervised contrastive loss can improve the robustness of the learned representations by leveraging the label information. (Section 3.2)
- R2: From our results, we can observe the **significant positive impact of full adversarial fine-tuning on the robustness of representations learned by contrastive learning**. However, full adversarial fine-tuning is ineffective in supervised contrastive or standard supervised learning schemes. (Section 3.3.1)
- R3: Our study reveals important insights regarding the impact of adversarial training on representations in different learning schemes. We observed **substantial differences between adversarial and clean representations in standard-trained networks across various learning schemes**. However, **after adversarial training, we observed a remarkable similarity between adversarial and clean representations**. This indicates that **regardless of the learning scheme utilized, adversarial training facilitates the convergence of representations towards a universal set, characterized by features² that consistently emerge across different models and tasks (Olah et al., 2020)**. Additionally, we found that **increasing the similarity between adversarial and clean representations, especially at the end of the network, improves the robustness of the network**. These findings offer valuable insights into designing and training more efficient and effective robust networks. (Section 3.3.2)

2. Methodology

In this section, we explain the methodology of our comparative study on the robustness of the three following learning schemes:

- *Contrastive Learning (CL)*: In the standard framework of SimCLR, contrastive learning trains a base encoder by minimizing a contrastive loss over the representations projected into a latent space (Figure 1a). The extracted features will train a linear classifier on a downstream task.

²A representation consists of all the features found in a layer.

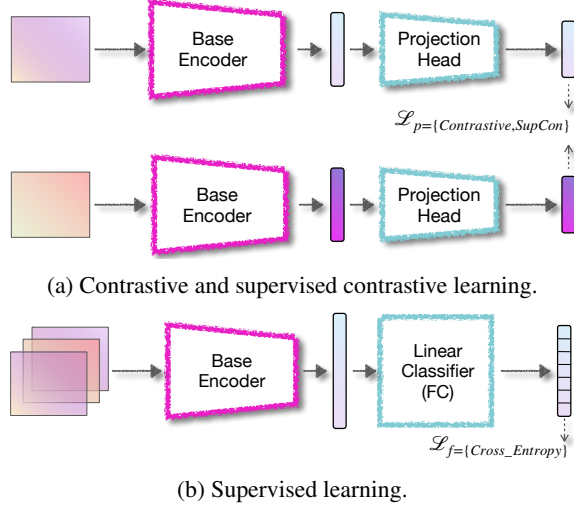


Figure 1. Training process of the studied learning schemes.

- **Supervised Contrastive Learning (SCL):** A supervised extension of contrastive learning introduced in (Khosla et al., 2020), to avoid false positive pairs selection by leveraging the label information.
- **Supervised Learning (SL):** The network consists of a base encoder followed by a fully connected layer as a linear classifier (see Figure 1b). In this case, cross-entropy between the true and predicted labels is utilized for training the network parameters.

The training process in contrastive and supervised contrastive learning includes the following two phases:

Pretraining Phase: The goal of this phase is to train the base encoder parameters θ_b by minimizing a self-supervised loss $\mathcal{L}_p(\theta_b, \theta_{ph})$ over a given dataset \mathcal{D}_p . Here θ_{ph} is the parameters vector of the projection head used to map the base encoder output into a low dimensional latent space where the \mathcal{L}_p is applied.

Supervised Fine-tuning Phase: The goal of this phase is to train the linear classifier parameters θ_c by minimizing the supervised loss $\mathcal{L}_f(\theta_c)$ over a labeled dataset \mathcal{D}_f . The linear classifier learns to map the representations extracted during the pretraining phase to the labeled space, where \mathcal{L}_f is the cross-entropy loss.

We examine the standard and robust training variations of the aforementioned training phases to compare the adversarial robustness across different learning schemes. Table 1 summarises all the studied training combinations for different possible scenarios of training phases in contrastive and supervised contrastive learning schemes.

In the *standard training of pretraining phase*, \mathcal{L}_p is given by $\mathcal{L}_{CL}(\theta_b, \theta_{ph}; \mathbf{x}', \mathbf{x}'')$ and $\mathcal{L}_{SCL}(\theta_b, \theta_{ph}; \mathbf{x}', \mathbf{x}'', \mathbf{y})$ in-

roduced in (Chen et al., 2020a) and (Khosla et al., 2020) for the contrastive and supervised contrastive schemes, respectively. Here the \mathbf{x}' and \mathbf{x}'' are two transformed views of the same minibatch and label vector \mathbf{y} is leveraged in \mathcal{L}_{SCL} to avoid false positive pair selection, as explained in (Khosla et al., 2020).

In the scenario of *robust pretraining or adversarial training of the representations*, we adopt a loss function \mathcal{L}_p inspired by Adversarial Contrastive Learning (Jiang et al., 2020). This loss function is formulated as a linear combination of two terms, as defined below:

$$\mathcal{L}_p(\theta_b, \theta_{ph}) = \alpha \mathcal{L}_{CL}(\theta_b, \theta_{ph}; \mathbf{x}', \mathbf{x}'') + \beta \mathcal{L}_{CL}(\theta_b, \theta_{ph}; \mathbf{x}, \mathbf{x}_{adv}) \quad (1)$$

and

$$\mathcal{L}_p(\theta_b, \theta_{ph}) = \alpha \mathcal{L}_{SCL}(\theta_b, \theta_{ph}; \mathbf{x}', \mathbf{x}'', \mathbf{y}) + \beta \mathcal{L}_{SCL}(\theta_b, \theta_{ph}; \mathbf{x}, \mathbf{x}_{adv}, \mathbf{y}) \quad (2)$$

for the contrastive and supervised contrastive schemes, respectively. In these equations, \mathbf{x}' and \mathbf{x}'' are the transformed views of \mathbf{x} , while \mathbf{x}_{adv} is the PGD attack generated by maximizing the associated loss function iteratively over each given minibatch \mathbf{x} as follows:

$$\mathbf{x}^{t+1} = \Pi_{\mathbf{x}+\mathcal{S}}(\mathbf{x}^t + \alpha \text{sgn}(\nabla_{\mathbf{x}} \mathcal{L}_{CL}(\theta_b, \theta_{ph}; \mathbf{x}, \mathbf{x}_{adv}))) \quad (3)$$

and

$$\mathbf{x}^{t+1} = \Pi_{\mathbf{x}+\mathcal{S}}(\mathbf{x}^t + \alpha \text{sgn}(\nabla_{\mathbf{x}} \mathcal{L}_{SCL}(\theta_b, \theta_{ph}; \mathbf{x}, \mathbf{x}_{adv}, \mathbf{y}))) \quad (4)$$

for the contrastive and supervised contrastive schemes, respectively. Here, $\Pi_{\mathbf{x}+\mathcal{S}}$ denotes projecting perturbations into the set of allowed perturbations \mathcal{S} and α is the step size. In this setup, the minimization of \mathcal{L}_p corresponds to the simultaneous minimization of each term weighted by coefficients of α and β . In this study, we take $\alpha = \beta$ to avoid unjustified prioritization between the transformed and adversarial terms.

We consider two alternatives for the *robust training of the fine-tuning phase*. In the *partial adversarial training*, we only update the linear classifier parameters θ_c by minimizing the loss function,

$$\mathcal{L}_f(\theta_c) = \alpha \mathcal{L}_{CE}(\theta_c; \mathbf{x}, \mathbf{y}) + \beta \mathcal{L}_{CE}(\theta_c; \mathbf{x}_{adv}, \mathbf{y}) \quad (5)$$

As the second alternative, *full adversarial fine-tuning* utilizes the following loss function,

$$\mathcal{L}_f(\theta_b, \theta_c) = \alpha \mathcal{L}_{CE}(\theta_b, \theta_c; \mathbf{x}, \mathbf{y}) + \beta \mathcal{L}_{CE}(\theta_b, \theta_c; \mathbf{x}_{adv}, \mathbf{y}) \quad (6)$$

Table 1. Summary of the training scenarios.

Scenarios	Pretraining Phase	Finetuning Phase
ST	Standard Training	Standard Training (with fixed θ_b)
AT	Adversarial Training	Standard Training (with fixed θ_b)
Partial-AT	Adversarial Training	Partial Adversarial Training (with fixed θ_b)
Full-AT	Adversarial Training	Full Adversarial Training

to readjust the base encoder parameters θ_b and train the linear classifier. In these equations, x_{adv} is the PGD attack generated by maximizing the cross-entropy loss iteratively over each given minibatch x . Here we take $\alpha = \beta$ in parallel with the pretraining phase.

We investigate both the standard and robust training approaches of the aforementioned training phases to compare the adversarial robustness among various learning schemes.

3. Experiments

Our goal is to understand whether there are differences in how contrastive learning learns the representation from data compared to supervised learning from the adversarial perspective. To this end, we conduct extensive experiments to evaluate the robustness of Contrastive Learning (CL), Supervised Contrastive Learning (SCL), and Supervised Learning (SL) under different training scenarios as shown in Table 1 on CIFAR-10 and CIFAR-100 image classification benchmarks. In all the subsequent experiments, we train the base encoder and the linear classifier on the same dataset. Our experimental setup is provided in Appendix A.

3.1. Threat Models

To evaluate the robustness of each scenario, we consider two different threat models:

- Threat Model-I (end-to-end attack generated by \mathcal{L}_{CE}): In this threat model, the attacker has complete knowledge of architecture and network parameters in the base encoder and linear classifier. This excludes any knowledge about the projection head utilized in the pretraining phase. The attacks are generated end-to-end through the utilization of the cross-entropy loss.
- Threat Model-II (attack generated against base encoder by \mathcal{L}_p): In Threat Model-II, the attacker possesses complete knowledge of all components in the pretraining phase, including the architecture, network parameters, loss function, and training dataset.

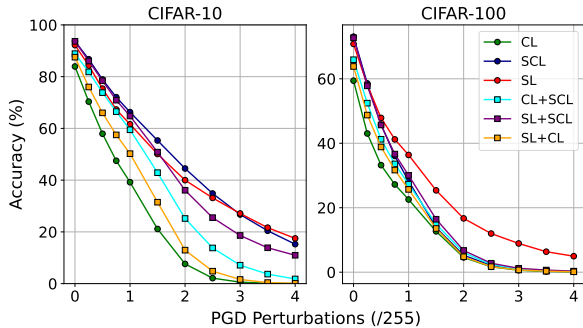


Figure 2. **Incorporating label information into contrastive learning enhances the robustness of the resulting representations.** We compare the test accuracy of different learning schemes on CIFAR10 and CIFAR100 datasets against adversarial examples through standard training settings. Contrastive learning without labels shows lower robustness compared to other learning schemes. We also observed that semi-supervised learning schemes SL+CL or SCL+CL achieve better robust performance than the CL scheme.

3.2. Contrastive Learning: Robustness and Label Impact Analysis through Standard Training

While it is widely known that neural networks trained through standard training are vulnerable to adversarial examples, the degree of vulnerability may differ among models trained using different learning methods. Hence, here we aim to investigate the vulnerability of various learning algorithms through standard training to evaluate their performance. The results as a function of the perturbation size under Threat Model-I are shown in Figure 2. We use 20-step l_∞ Projected Gradient Descent (PGD) (Madry et al., 2017) attacks with different perturbations to generate adversarial attacks during this experiment. In addition to the learning schemes mentioned before, the combination of them including the combination of supervised learning and supervised contrastive learning (denoted SL+SCL), and two other semi-supervised versions, including the combination of contrastive learning with supervised learning (denoted SL+CL) or with supervised contrastive learning (denoted CL+SCL), are investigated. These combinations are designed to answer this question: *Does employing the label information benefit the robustness of contrastive learning representations?* As we can see, contrastive learning without label information is less robust than other learning schemes. We also observed that semi-supervised learning schemes SL+CL or CL+SCL achieve better robust performance than the CL scheme. Appendix B visualizes the representations learned by all these learning schemes using t-SNE on the CIFAR-10 dataset.

We have excluded the semi-supervised versions from the subsequent experiments to prevent any potential confusion between the effects of label information and adversarial training. Appendix C provides more results under Threat

Model-II, where the attacks are generated against only the base encoder.

3.3. Adversarial Training: Comparing Representations

Here, we first compare the performance of different adversarial training scenarios. Subsequently, a set of explanatory tools (e.g., CKA and linear probing) is employed to inspect how adversarial training affects the learned representations in hidden layers.

3.3.1. DIRECT COMPARISON

This experiment aims to evaluate model robustness under Threat Model-I in the following scenarios: i) training a base encoder using adversarial training, then training the linear classifier separately after freezing the base encoder (AT); ii) training the base encoder by adversarial training, then training the linear classifier separately using adversarial training after freezing the base encoder (Partial AT); iii) training the base encoder by adversarial training, then robustifying the end-to-end model (Full AT). The latter case means that the base encoder parameters are first adversarially trained using \mathcal{L}_{CL} or \mathcal{L}_{SCL} , then those parameters are fine-tuned during adversarial training of the linear classifier where adversarial examples are generated using \mathcal{L}_{CE} . The results on the CIFAR-10 and CIFAR100 datasets against 20-step different PGD attacks are shown in Figure 3. Appendix D also provides the robustness of different scenarios on the CIFAR100 dataset against a state-of-the-art adversarial attack known as Auto-attack (Croce & Hein, 2020). Moreover, Appendix E provides more results under Threat Model-II, where the attacks are generated against only the base encoder. The results indicate two main observations: (i) CL under Full AT consistently outperforms other learning schemes in various evaluation scenarios, demonstrating a noticeable improvement in standard accuracy and robustness against adversarial attacks. The results from the previous section shed light on the positive impact of utilizing label information to enhance the robustness of the CL scheme during standard training. It is worth noting that Full AT also effectively incorporates label information into the robust network architecture, leading to a remarkable overall improvement in robustness. This finding is consistent with prior research (Zhai et al., 2019; Fan et al., 2021) which has investigated the utilization of pseudo labels for unlabeled data, demonstrating their effectiveness in enhancing the adversarial robustness of neural networks. By incorporating label information, the CL scheme benefits from an additional source of valuable supervision, strengthening its defense against adversarial attacks and enhancing overall robustness. (ii) There is a slight difference in the performance of SCL under AT and Full AT scenarios. This indicates that the representations learned by SCL from the AT scenario are already sufficient to achieve acceptable robustness.

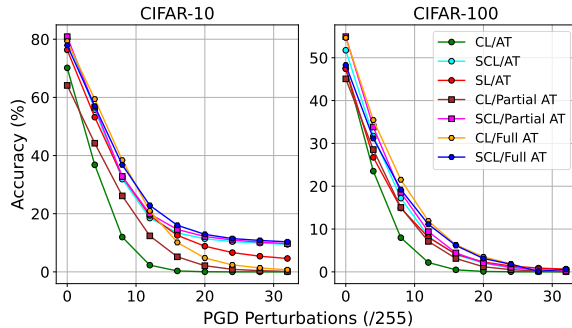


Figure 3. Full AT successfully integrates label information into the robust network architecture trained via contrastive learning, resulting in a remarkable overall improvement in robustness. We compare the test accuracy of different learning schemes against different PGD attacks on CIFAR10 and CIFAR100 datasets through different adversarial training scenarios. The results demonstrate that Full AT effectively incorporates label information into the adversarially-trained network obtained through the CL scheme. This incorporation leads to improved robustness. Moreover, there is a slight variation in the performance of the SCL under different adversarial training scenarios.

3.3.2. COMPARING CHARACTERISTICS

Previous results raise an important question from a representation learning perspective: what happens in the layers of neural networks when they are adversarially trained? We examine the internal layer representations learned by different robust learning schemes to shed light on this direction.

Adversarial and clean representations differ significantly in standard-trained networks. We begin our investigation by using CKA to study the internal representation structure of each model. CKA is a metric that measures the similarity between two sets of features. To answer how different learning schemes extract representations, we take every pair of layers X and Y within a model learned by different learning schemes and compute their CKA similarity on clean and adversarial examples. Figures 4 and 10 show the results as a heatmap for different learning algorithms under standard scenarios on clean and adversarial examples. Notably, we observe distinct differences in the internal representation structure between the three learning schemes: (i) the layers near the end of the network in SCL and SL schemes exhibit lower similarity with other layers compared to CL, and (ii) the dissimilarity between clean and adversarial representations in standard-trained networks highlights their vulnerability to adversarial examples. Previous research (Mitrovic et al., 2020) has assumed that data consists of content and style components, with only the content being relevant for unknown downstream tasks. Additionally, it is assumed that content and style are independent, meaning that style changes do not affect the underlying content. In contrast, adversarial perturbations

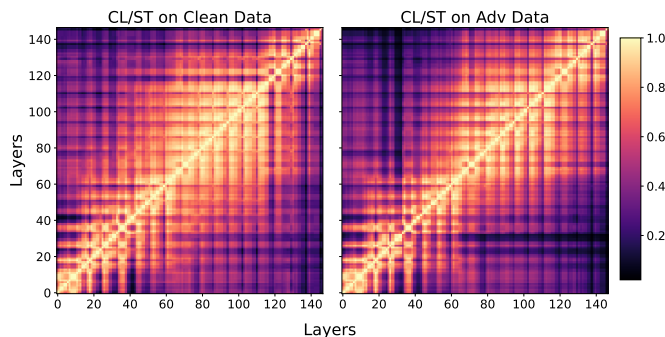


Figure 4. Standard-trained networks using different learning algorithms yield representations that display notable disparities between adversarial and clean examples. We compute the similarity of representations across all layer combinations in standard-trained networks that have been trained using the CL scheme, considering both clean and adversarial data.

are introduced as a specific distribution change in the natural data distribution (Zhang et al., 2020) where they alter the style while preserving the content, making them imperceptible to the human eye. Under these assumptions, the pronounced dissimilarity observed between adversarial and clean examples in standard-trained networks indicates an inability to extract content-related representations consistent across both examples.

Adversarial and clean representations exhibit substantial similarity in adversarially trained networks, regardless of the learning schemes used. We also perform previous comparisons for robust models, taking every pair of layers X and Y within an adversarially trained model robustified by different learning schemes and computing their CKA similarity on clean and adversarial examples. Moreover, we utilize linear probing as a conceptual tool to better understand the dynamics within the neural network and the specific roles played by individual intermediate layers. Figures 5 and 11 highlight several observations from the results of the adversarial training experiments: (i) Cross-layer similarities are amplified compared to standard training regardless of learning schemes used. This is evident by the higher degree of brightness in the plots. (ii) In networks trained through adversarial training, the adversarial representations are significantly similar to clean representations. Previous research (Jones et al., 2022) has demonstrated that robust training effectively mitigates the impact of adversarial perturbations, resulting in similar representations for clean and adversarial examples in robust networks. However, their studies have only focused on the supervised learning scheme. Our results support and extend these findings by demonstrating that the similarity between clean and adversarial representations holds irrespective of the learning scheme used. (iii) When comparing the representations obtained from AT and its counterpart, Full AT (see Figure

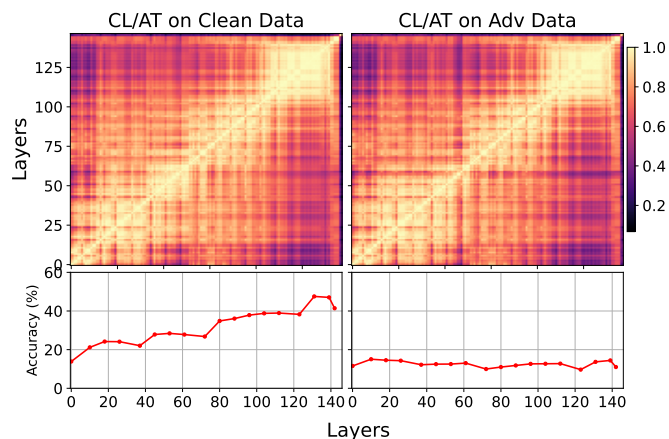


Figure 5. Regardless of the learning scheme employed, adversarially trained networks exhibit significant similarity between adversarial and clean representations. We compute the similarity of representations across all layer combinations in adversarially trained networks that have been trained using the CL scheme, considering both clean and adversarial data.

11), we observe a notable increase in long-range similarities within the CL framework. This enhancement in similarity translates to significant improvements in both standard and adversarial accuracy. Remarkably, Full AT significantly improves overall robustness by incorporating label information into the network. In contrast, the representations learned by SCL and SL under the AT and Full AT scenarios exhibit minor differences. These slight variations in representations result in marginal differences in performance, indicating that the label information utilized in AT already provides sufficient robustness for SCL and SL. Therefore, Full AT does not introduce additional information to enhance SCL and SL’s robustness further.

Increasing the similarity between adversarial and clean representations, especially near the end of the network, improves robustness. To gain a deeper understanding of the divergence between adversarial and clean representations, we compare each layer X in a model applied to clean data with its identical counterpart Y in the same model applied to adversarial examples. The results in the left-hand side of Figure 6 and Figure 12 illustrate that the adversarial representations in the network trained using standard training, exhibit significant dissimilarity from their clean counterparts, particularly towards the end of the network, regardless of the learning scheme used. In contrast, adversarial training significantly reduces the impact of adversarial perturbations, leading to similar representations for both clean and adversarial examples in robust networks. Notably, we observe a similarity drop across intermediate layers of CL/AT, which may explain its lower performance compared to other robust learning schemes. To confirm our explanation, we conducted a similar experiment on CL after Full AT

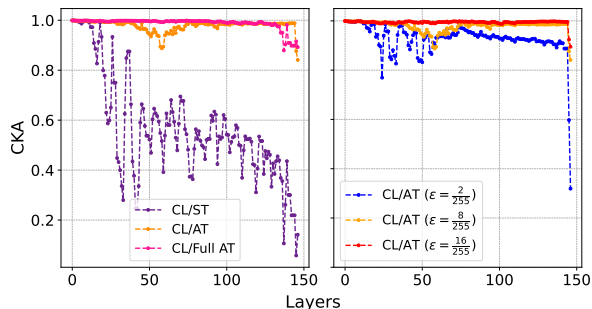


Figure 6. Increasing the similarity between adversarial and clean representations improves robustness, especially near the end of networks. Left) Comparing clean and adversarial representations in CL reveals significant dissimilarity in standard-trained networks. Adversarial training reduces this divergence significantly, but a drop in similarity is observed across intermediate layers, affecting its performance (CL/AT). Full AT effectively mitigates the similarity drop, enhancing overall model robustness. Right) During adversarial training, we increased perturbation budgets to vary the strength of adversarial attacks. This led to greater similarity between adversarial and clean representations, especially towards the end of the network.

(which significantly improves the robustness) and compared the resulting representations. As shown in the left-hand side of Figure 6, Full AT reduces the similarity drop, enhancing the model’s robustness. To gain further insights into this phenomenon, we conducted an ablation study by varying the strength of adversarial attacks during training through increased perturbation budgets (ϵ). As demonstrated in the right-hand side of Figure 6 and Figure 13, we observed that stronger adversarial perturbations led to an enhanced similarity between adversarial representations and their counterpart clean representations, particularly in the later layers of the network. This significant finding confirms the previous results reported in (Cianfarani et al., 2022) and highlights our novel contribution in extending this observation to contrastive learning schemes. Furthermore, from heatmaps in Figure 14, we can observe that increasing the strength of adversarial perturbations leads to the long-range similarity between different layers.

Unlike standard training, adversarial training converges toward a universal set of representations, regardless of the learning schemes utilized. Here, we perform cross-model comparisons to measure the similarities between all layers X of one model trained using a specific learning scheme and all layers Y of another model trained using a different learning scheme. The left-hand side of Figures 7 and 15 present the results for models trained using different standard learning schemes. We observed that, except for lower layers, the representations extracted by other layers were highly dissimilar. However, after applying adversarial training (shown on the right-hand side), the similarity be-

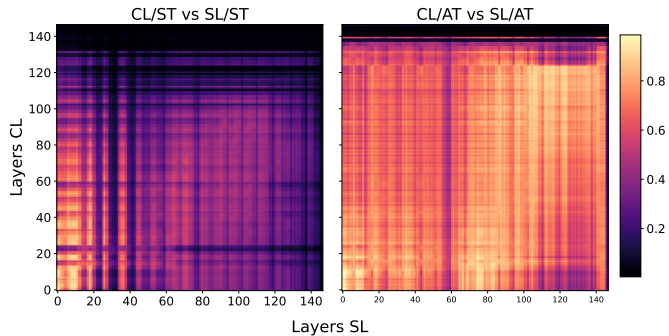


Figure 7. Adversarial training promotes convergence towards a universal set of representations. Standard-trained networks exhibit significant dissimilarity in adversarial representations across CL schemes, particularly in higher layers. However, after applying adversarial training, the similarity between layers increases, indicating a shift towards extracting a universal set of representations.

tween layers from networks trained with different learning schemes notably increased, indicating a tendency towards extracting a universal set of representations. This finding aligns with a previous study (Jones et al., 2022), which highlighted that robust networks converge towards a universal set of representations regardless of the architecture. Our results extend these observations to contrastive learning schemes, providing empirical evidence for a more general claim. Specifically, our study demonstrates that adversarial training promotes convergence toward a universal set of representations regardless of the learning schemes employed.

4. Conclusion

This study compared the robustness of contrastive and supervised contrastive learning with standard supervised learning. Our results demonstrated the benefits of incorporating label information in contrastive learning for enhanced robustness. Adversarial training reduced disparities between adversarial and clean representations, leading to convergence toward a universal set of representations. The increased similarity between adversarial and clean representations improved robustness, especially in deeper layers. These findings offer valuable insights for optimizing robust learning schemes.

Acknowledgements

This work has been supported in part by NSF (Awards 2233873, 2007202, and 2107463), NASA (Award 80NSSC20K1720), and Chameleon Cloud.

References

Alain, G. and Bengio, Y. Understanding intermediate layers using linear classifier probes. *arXiv preprint*

- arXiv:1610.01644*, 2016.
- Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., and Joulin, A. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020a.
- Chen, T., Liu, S., Chang, S., Cheng, Y., Amini, L., and Wang, Z. Adversarial robustness: From self-supervised pre-training to fine-tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 699–708, 2020b.
- Cianfarani, C., Bhagoji, A. N., Sehwag, V., Zhao, B., Zheng, H., and Mittal, P. Understanding robust learning through the lens of representation similarities. *Advances in Neural Information Processing Systems*, 35:34912–34925, 2022.
- Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–2216. PMLR, 2020.
- Fan, L., Liu, S., Chen, P.-Y., Zhang, G., and Gan, C. When does contrastive learning preserve adversarial robustness from pretraining to finetuning? *Advances in Neural Information Processing Systems*, 34:21480–21492, 2021.
- Gowal, S., Huang, P.-S., van den Oord, A., Mann, T., and Kohli, P. Self-supervised adversarial robustness for the low-label, high-data regime. In *International conference on learning representations*, 2021.
- Grill, J.-B., Strub, F., Althé, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Gupta, R., Akhtar, N., Mian, A., and Shah, M. On higher adversarial susceptibility of contrastive self-supervised learning. *arXiv preprint arXiv:2207.10862*, 2022.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Hendrycks, D., Mazeika, M., Kadavath, S., and Song, D. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32, 2019.
- Islam, A., Chen, C.-F. R., Panda, R., Karlinsky, L., Radke, R., and Feris, R. A broad study on the transferability of visual representations with contrastive learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8845–8855, 2021.
- Jiang, Z., Chen, T., Chen, T., and Wang, Z. Robust pre-training by adversarial contrastive learning. *Advances in Neural Information Processing Systems*, 33:16199–16210, 2020.
- Jones, H. T., Springer, J. M., Kenyon, G. T., and Moore, J. S. If you’ve trained one you’ve trained them all: inter-architecture similarity increases with robustness. In *Uncertainty in Artificial Intelligence*, pp. 928–937. PMLR, 2022.
- Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., and Krishnan, D. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020.
- Kim, M., Tack, J., and Hwang, S. J. Adversarial self-supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:2983–2994, 2020.
- Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. Similarity of neural network representations revisited. In *International Conference on Machine Learning*, pp. 3519–3529. PMLR, 2019.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Mitrovic, J., McWilliams, B., Walker, J., Buesing, L., and Blundell, C. Representation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020.
- Moshavash, M., Eftekhari, M., and Bahraman, K. Momentum contrast self-supervised based training for adversarial robustness. *Journal of Computing and Security*, 8(1):33–43, 2021.
- Nguyen, T., Raghu, M., and Kornblith, S. Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. *arXiv preprint arXiv:2010.15327*, 2020.
- Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. <https://distill.pub/2020/circuits/zoom-in>.
- Pal, M., Jati, A., Peri, R., Hsu, C.-C., AbdAlmageed, W., and Narayanan, S. Adversarial defense for deep speaker recognition using hybrid adversarial training. In *ICASSP*

2021-2021 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6164–6168. IEEE, 2021.

Subramanian, A. torch_cka. <https://github.com/AntixK/PyTorch-Model-Compare>, 2021.

Wahed, M., Tabassum, A., and Lourentzou, I. Adversarial contrastive learning by permuting cluster assignments. *arXiv preprint arXiv:2204.10314*, 2022.

Zhai, R., Cai, T., He, D., Dan, C., He, K., Hopcroft, J., and Wang, L. Adversarially robust generalization just requires more unlabeled data. *arXiv preprint arXiv:1906.00555*, 2019.

Zhang, C., Zhang, K., and Li, Y. A causal view on robustness of neural networks. *Advances in Neural Information Processing Systems*, 33:289–301, 2020.

Zhong, Y., Tang, H., Chen, J., Peng, J., and Wang, Y.-X. Is self-supervised learning more robust than supervised learning? *arXiv preprint arXiv:2206.05259*, 2022.

Below, we provide supplementary information and additional results from the various sections.

A. Experiment Setup

Our experiment setup is similar to that used in (Chen et al., 2020a) for SimCLR and (Khosla et al., 2020) for SupCon, which are prominent works in contrastive learning. We use ResNet-50 as the base encoder for all scenarios and a two-layers MLP network as the projection head. The loss is optimized using the Adam optimizer with a learning rate of 0.0003. We train each model for 200 epochs using a mini-batch size of 128 for standard and 256 for adversarial scenarios. In all adversarial training scenarios, the adversarial perturbations are generated using a 5-step Projected Gradient Descent (PGD) attack under the l_∞ norm with a maximum perturbation limit of $\epsilon = 8/255$, unless a specific value of ϵ is specified. The models are evaluated against PGD attacks and state-of-the-art Auto-attacks (Croce & Hein, 2020) at the test time. We report the top-1 test accuracy for all scenarios to evaluate the mentioned scenarios.

B. t-SNE Visualization of Learning Schemes under Standard Training Scenario

Figure 8 visualizes the representations learned by CL, SCL, SL, SL+SCL, SL+CL, and CL+SCL learning schemes using t-SNE on the CIFAR-10 dataset. Here, we used labels to color the markers corresponding to each data point. The results depicted in the ST scenario clearly show much clearer class boundaries in the SCL and SL compared to the CL scheme. Furthermore, we can observe that employing some label information in both semi-supervised learning schemes, SL+CL and CL+SCL, can lead to separate classes more clearly. This suggests that semi-supervised learning versions make it difficult for an adversary to successfully perturb an image, resulting in a more robust prediction.

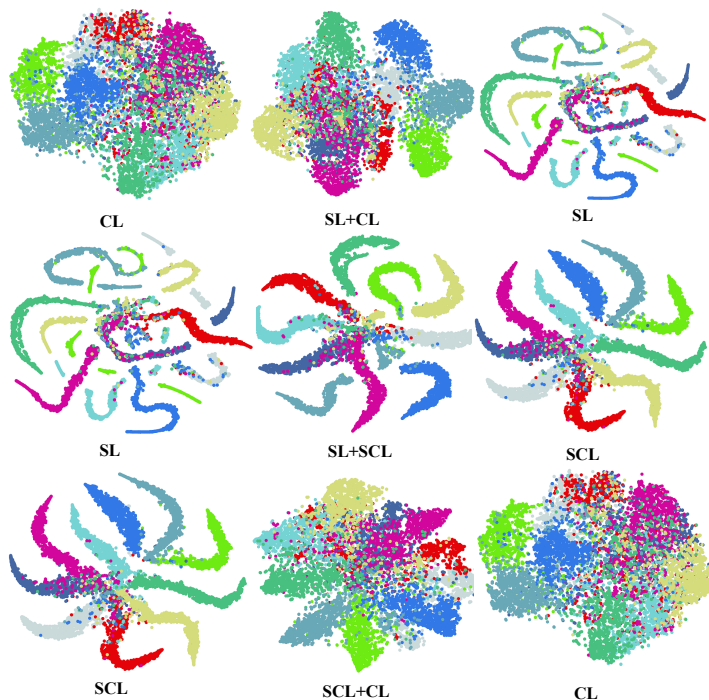


Figure 8. Semi-supervised learning schemes (SL+CL and SCL+CL) separate classes more clearly than contrastive learning (CL) schemes. We visualize the representations learned by different learning schemes (CL, SCL, SL, SL+SCL, SL+CL, and SCL+CL) using t-SNE on the CIFAR-10 dataset. The results show that in the ST scenario, both SCL and SL exhibit clearer class boundaries compared to CL. Furthermore, incorporating label information in the semi-supervised learning schemes (SL+CL and SCL+CL) enhances the separation of classes, indicating increased robustness against adversarial perturbations.

C. Robustness through Standard Training under Threat Model-II

Threat Model-II is not applicable for the supervised learning scheme, as the base encoder and linear classifier are trained together end-to-end. Table 2 reports the results on CIFAR-10 and CIFAR-100 datasets against different 40-step PGD attacks. From the results, we can observe that the models trained using \mathcal{L}_{SCL} are more robust compared to \mathcal{L}_{CL} where the attacks are generated against the base encoder. This suggests that false negative pair selection in self-supervised contrastive learning leads to making the model less robust which is aligned with the results reported in (Gupta et al., 2022).

Table 2. **SCL is more robust than CL scheme against adversarial attacks generated against base encoder.** The performance of contrastive learning schemes, including CL and SCL, in the ST scenario evaluated under Threat Model-II (attack generated against base encoder). The best performance is highlighted in bold. Threat Model-II is not applicable to the SL scheme, as the base encoder and linear classifier are trained together end-to-end.

Models	Dataset	Standard Training	PGD (4/255)	PGD (4/255)	PGD (16/255)
SCL	CIFAR10	93.56	30.3	12.9	10.06
CL		84.27	19.58	9.61	7.19
SCL	CIFAR-100	73.38	7.87	3.28	2.16
CL		60.28	5.52	0.94	0.23

D. Performance of Different Adversarial Training Scenarios against a Range of Auto-Attacks

Here, we evaluate the robustness of different scenarios on the CIFAR100 dataset against different Auto-attacks. The results are shown in Figure 9. The comparative analysis provides evidence of Full AT being effective in improving the robustness of the CL-based network. This improvement is achieved through the integration of label information, resulting in enhanced robustness against adversarial attacks. Furthermore, slight variations in the performance of SCL are observed across different scenarios of adversarial training.

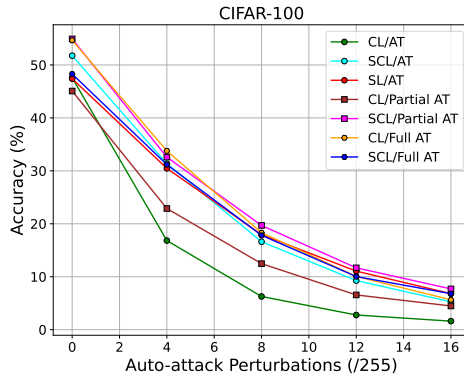


Figure 9. **Full AT effectively incorporates label information into the robust network trained via contrastive learning, significantly improving overall robustness.** The comparative analysis of different learning schemes’ test accuracy against various Auto-attacks on the CIFAR-100 dataset demonstrates the efficacy of Full AT in enhancing the robustness of the CL-based network by integrating label information. This integration leads to improved robustness against adversarial attacks. Additionally, slight variations in the performance of SCL are observed under different adversarial training scenarios.

E. Robustness through Adversarial Training under Threat Model-II (Non-Transferability of Cross-Task Robustness)

In contrastive learning, any objective that utilizes the learned representations of the base encoder is referred to as a downstream task. In this experiment, we compare model robustness where adversarial training is applied only to the base encoder in contrastive and supervised contrastive learning schemes. Table 3 shows the performance of the models on CIFAR-10 and CIFAR-100 datasets against different 40-step PGD attacks. Compared to Threat Model-I, our findings

in Threat Model-II indicate that the robust model achieved by applying adversarial training solely to the base encoder is not transferable to the downstream task. This failure in transferring robustness across the tasks, known as the cross-task robustness transferability challenge, has also been reported in (Fan et al., 2021).

Table 3. The robustness achieved through adversarial training solely applied to the base encoder does not transfer effectively to the downstream task. The performance of contrastive learning schemes, including CL and SCL in AT scenario, is compared to the baseline adversarially trained SL scheme in terms of top-1 accuracy on CIFAR-10 and CIFAR-100 datasets. The models are evaluated under Threat Model-I (end-to-end attack generated by cross-entropy loss) and threat model-II (attack generated against base encoder). The effectiveness of AT in enhancing the robustness of CL and SCL under Threat Model-II is evident. However, the results reveal that these robust models lose their robustness when subjected to end-to-end attacks or under Threat Model-I.

Models	Datasets	Standard Training	Adversarial Training	End-to-End Attack Generated by Cross-Entropy Loss			Attack Generated Against Base Encoder		
		Clean	Clean	PGD (4/255)	PGD (8/255)	PGD (16/255)	PGD (4/255)	PGD (8/255)	PGD (16/255)
SL		91.73	76.33	53.2	32.52	10.5	NA	NA	NA
SCL	CIFAR-10	93.56	80.22	11.99	55.83	31.45	59.36	62.25	59.38
CL		84.27	70.13	36.86	11.91	0.26	69.1	65.88	49.8
SL		69.56	47.4	26.68	14.8	3.8	NA	NA	NA
SCL	CIFAR-100	73.38	51.8	31.7	17.11	3.65	37.09	30.02	24.5
CL		60.28	47.58	23.58	7.83	0.43	46.18	42.44	26.9

F. Representation Structure of Different Learning Schemes

In this section, we provide all the figures related to the section 3.3.2 in the main body.

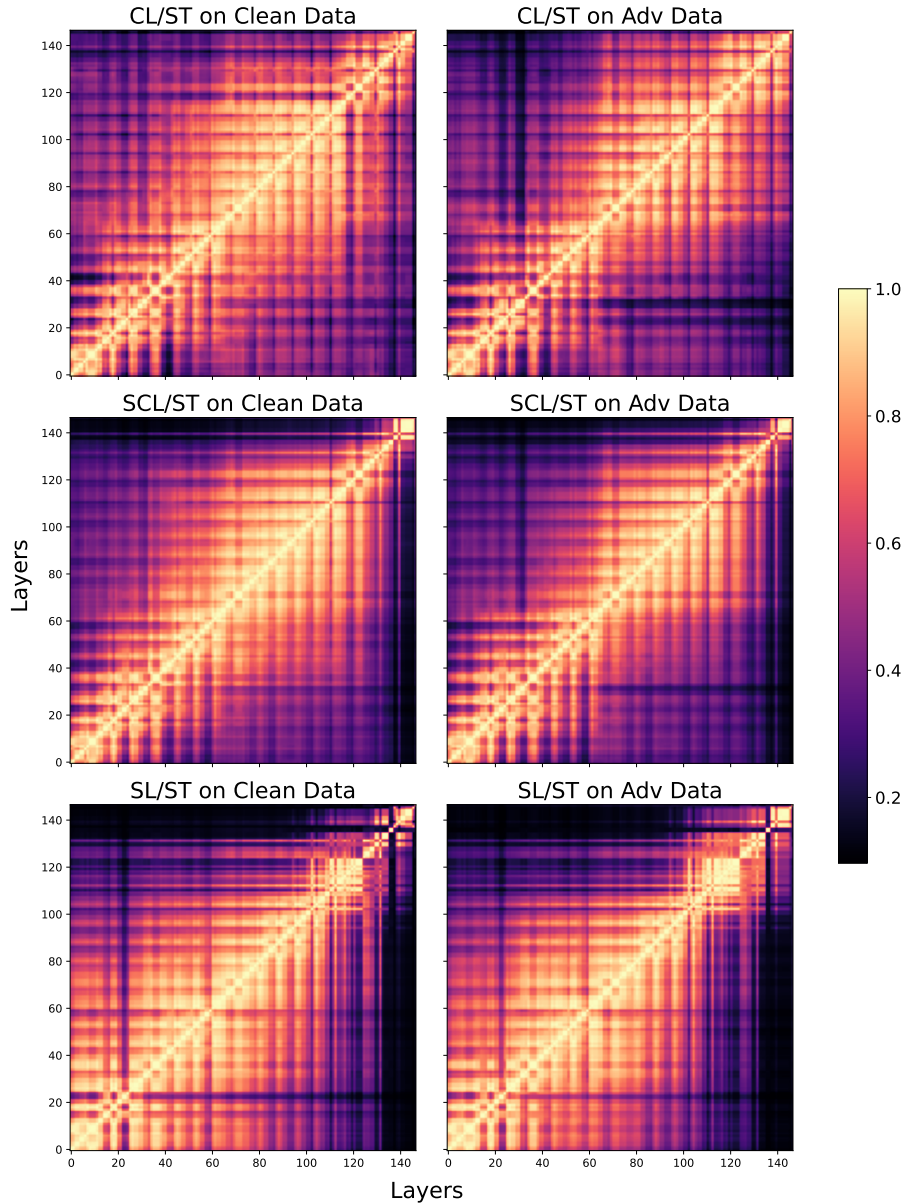


Figure 10. The representations obtained from standard-trained networks exhibit significant differences between adversarial and clean examples, regardless of the learning algorithm utilized. We compute the similarity of representations across all layer combinations in standard-trained networks that have been trained using different learning schemes, considering both clean and adversarial data. The three learning schemes (SCL, SL, and CL) have noticeable differences in their internal representation structures (the first column). CL demonstrates more consistent representations throughout the network when compared to SCL and SL. Moreover, standard-trained networks exhibit substantial dissimilarity between clean and adversarial representations (the first column vs. the second one).

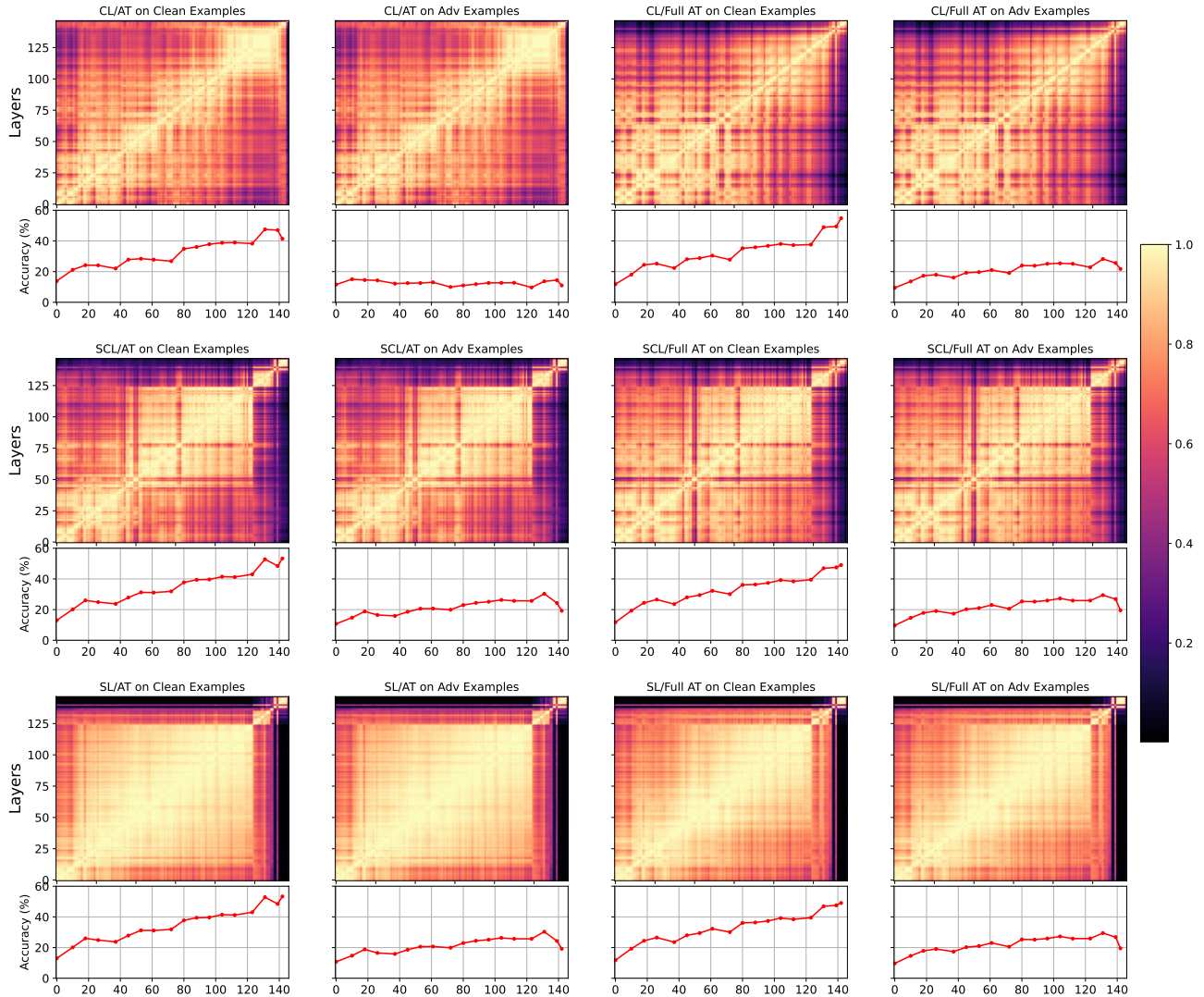


Figure 11. The similarity between adversarial and clean representations is substantial in adversarially trained networks, regardless of the learning scheme used. We analyze robust models by comparing layer pairs within different learning schemes and calculating their CKA similarity on clean and adversarial examples. Linear probing is employed to gain insights into the network dynamics and the roles of intermediate layers. The results demonstrate amplified cross-layer similarities compared to standard training, indicated by higher brightness levels in the plots. Additionally, networks trained through adversarial training exhibit significant similarities between adversarial and clean representations. Moreover, upon comparing the representations obtained from AT and its counterpart Full AT, we observe a significant enhancement in long-range similarities within CL. This improvement in similarity leads to substantial improvements in both standard and adversarial accuracy. In contrast, the representations learned by SCL and SL under AT and Full AT scenarios exhibit slight differences, resulting in minor variations in their performance.

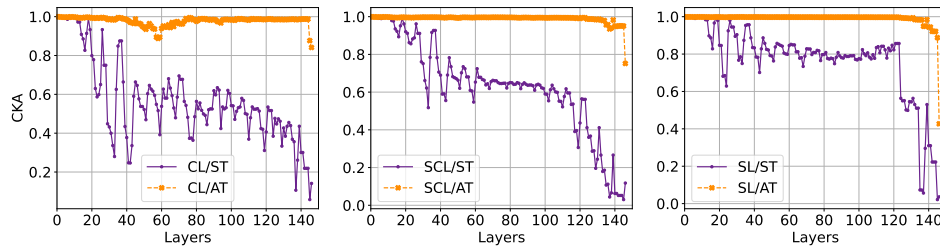


Figure 12. **Contrasting adversarial representations with their clean counterparts.** Comparing clean and adversarial representations in different layers of the model reveals significant dissimilarity in standard-trained networks. Adversarial training reduces this divergence, leading to similar representations for clean and adversarial examples in robust networks. However, there is a drop in similarity across intermediate layers for CL/AT, which may explain its lower performance compared to other robust learning schemes.

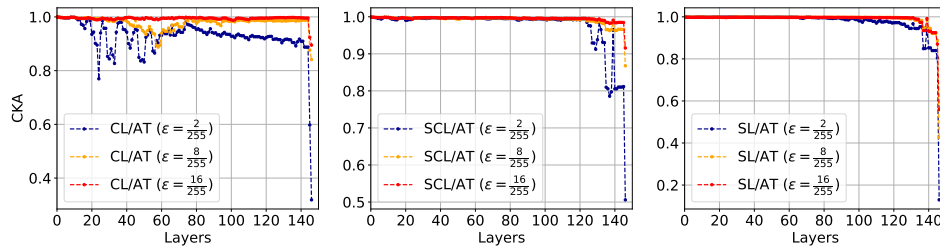


Figure 13. **Increasing the similarity between adversarial and clean representations improves robustness, especially near the end of networks.** Improving the similarity between adversarial and clean representations enhances robustness. During adversarial training, we increased perturbation budgets to vary the strength of adversarial attacks. This led to greater similarity between adversarial and clean representations, especially towards the end of the network.

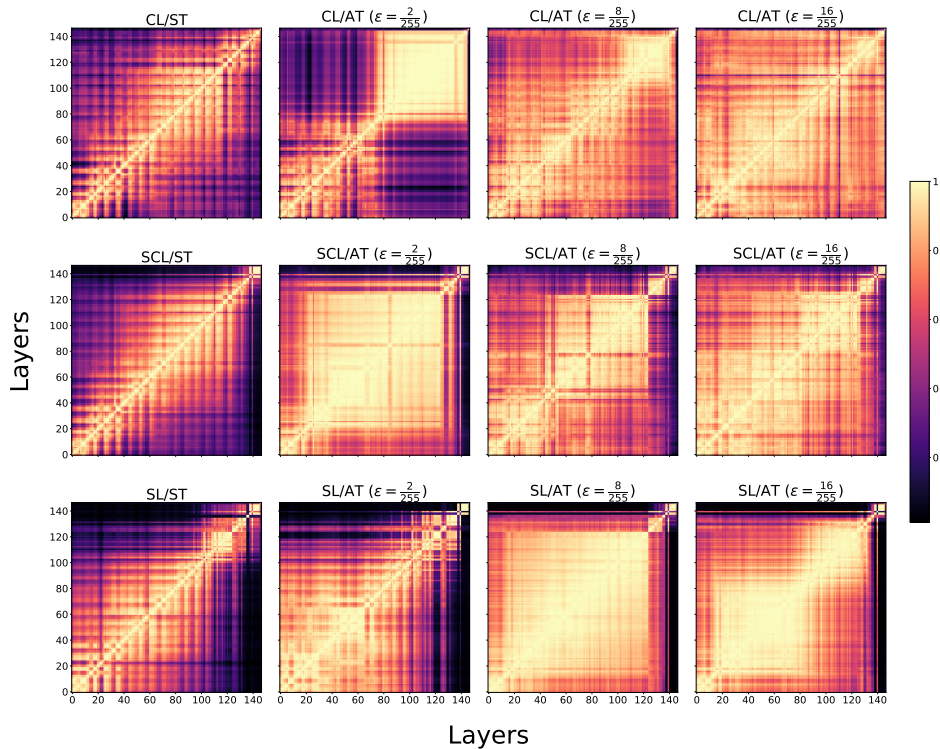


Figure 14. **Adversarial training promotes the emergence of long-range similarities between layers, regardless of the specific learning scheme employed.** We vary the strength of adversarial attacks during training. We observe that increasing the strength of adversarial perturbations leads to a consistent presence of long-range similarity between layers, independent of the learning scheme used.

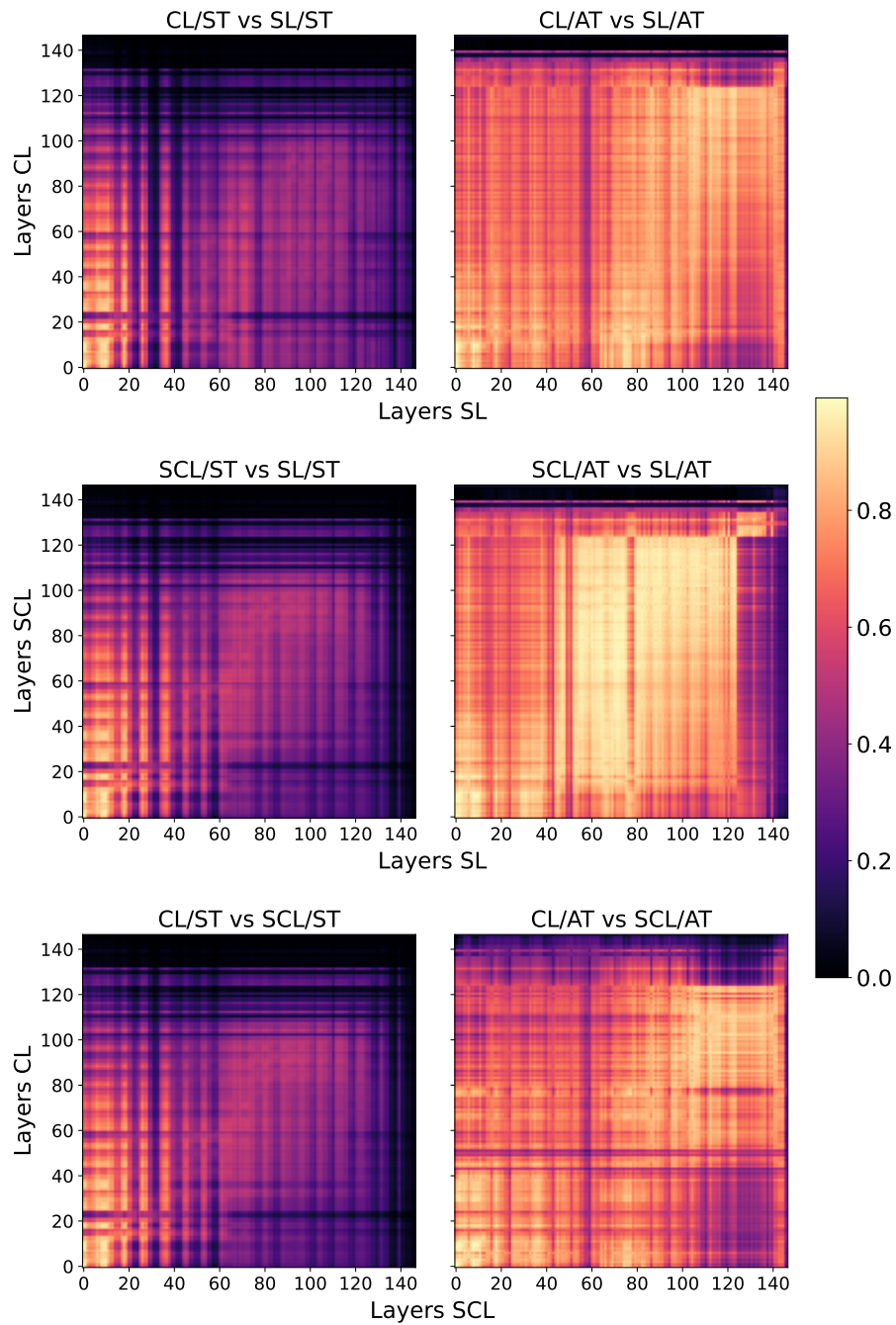


Figure 15. Unlike standard-trained networks, the ones trained through adversarial training show significant similarity in adversarial representations across different learning schemes. The cross-model CKA heatmap between standard-trained networks trained using different learning schemes highlights that these schemes extract distinct adversarial representations, particularly in a large number of higher layers within the network. Cross-model comparisons demonstrate that, after applying adversarial training, the similarity between layers from different learning schemes increases, suggesting a shift towards extracting a universal set of representations.