

440 **A Additional Experiment: Adversarial Examples (in Adversarial Training)
Contain Adversarial Non-robust Features**

442 In Section 3, we collect misclassified adversarial examples on an early checkpoint A and evaluate them
 443 on a later checkpoint B with their correct labels to examine the memorization of those adversarial
 444 non-robust features. This is because adversarial examples used for AT consist of robust features from
 445 the original label y and non-robust features from the misclassified label \hat{y} , so B must have memorized
 446 those features (non-robust on A) to correctly classify them. In this section, we further provide a
 447 rigorous discussion on the non-robust features contained in the adversarial examples through an
 448 additional experiment.

449 **Experiment Design.** Recall that Ilyas et al. [15] leverage targeted PGD attack [22] to craft adversarial
 450 examples on a standard (non-robust) classifier and relabel them with their target labels to build a
 451 non-robust dataset. Finding standard training on it yields good accuracy on clean test data, they
 452 prove that those adversarial examples contain non-robust features corresponding to the target labels
 453 (Section 3.2 of their paper).

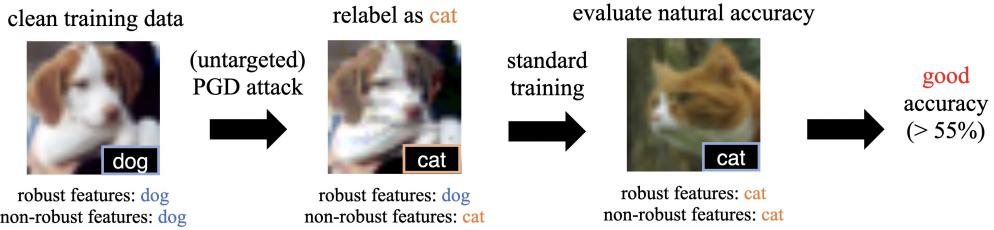


Figure 7: Mining non-robust features from adversarial examples in adversarial training. 1) Craft adversarial examples \hat{x}_i using untargeted PGD attack on a AT checkpoint. 2) Relabel them (only misclassified ones) with misclassified label \hat{y}_i to build a non-robust dataset $\{(\hat{x}_i, \hat{y}_i)\}$. 3) Perform standard training from the checkpoint. 4) Achieve good natural accuracy ($> 55\%$).

454 However, different from their settings where we know the target labels of the adversarial examples
 455 and they are generated on a non-robust classifier, we lay emphasis on mining non-robust features
 456 from adversarial examples generated on-the-fly in AT. To this end, we first craft adversarial examples
 457 on the AT checkpoint at some epoch using untargeted PGD attack [22] following the real setting of
 458 AT, then relabel the *misclassified* ones \hat{x}_i (e.g., a dog) with their misclassified labels $\hat{y}_i \neq y_i$ (e.g., cat)
 459 to build a non-robust dataset $\{(\hat{x}_i, \hat{y}_i)\}$. In order to capture non-robust features at exactly the training
 460 epoch these adversarial examples are used, we continue to perform standard training directly from
 461 the AT checkpoint on the non-robust dataset, and finally evaluate natural accuracy. See Figure 7 for
 462 an illustration of our experiment.

463 In details, we select the AT models at epoch 60 (before LR decay) and epoch 1,000 (after LR decay)
 464 to conduct the above experiment. We use untargeted PGD-20 with perturbation norm $\varepsilon_\infty = 16/255$
 465 to craft adversarial examples on those checkpoints from the *training data* of CIFAR-10 [17]. The
 466 attack we adopt is a little bit stronger than the attack of PGD-10 with $\varepsilon_\infty = 8/255$ that is commonly
 467 used to generate adversarial examples in baseline adversarial training, because the training robust
 468 accuracy rises to as high as 94.67% at epoch 1,000 (Figure 1a) and gives less than 2,700 misclassified
 469 examples, which is significantly insufficient for further training. Using the stronger attack, we obtain a
 470 success rate of 78.38% on the checkpoint at epoch 60 and a success rate of 34.86% on the checkpoint
 471 at epoch 1,000. Since the attack success rates are different, we randomly select a same number of
 472 misclassified adversarial examples for standard training for a fair comparison. The learning rate is
 473 initially set to 0.1 and decays to 0.01 after 20 epochs for another 10 epochs of fine-tuning. Given that
 474 the original checkpoints already have non-trivial natural accuracy, we also add two control groups
 475 that train with random labels instead of \hat{y} to exclude the influence that may brought by the original
 476 accuracy.

477 **Results.** As shown in Figure 8, we find that standard training on the non-robust dataset $\{(\hat{x}_i, \hat{y}_i)\}$
 478 successfully converges to fairly good accuracy ($> 55\%$) on natural test images, *i.e.*, predicting cats as
 479 cats, no matter from which AT checkpoints (either epoch 60 or epoch 1,000) the adversarial examples
 480 are crafted. Also, we can see that the non-trivial natural accuracy has nothing to do with the original
 481 accuracy of the AT checkpoints, as the accuracy plummets to around 10% (random guessing) as
 482 soon as the standard training starts with random labels. This proves that *adversarial examples in*

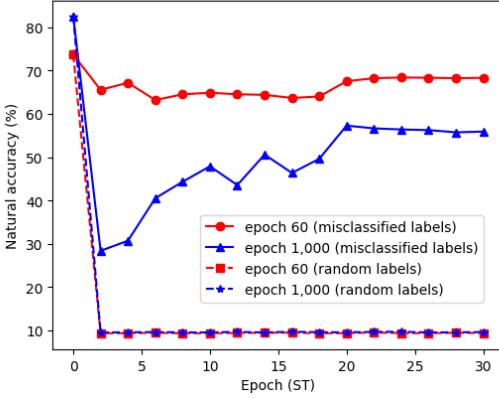


Figure 8: Natural accuracy during standard training. Standard training on the non-robust dataset built from the checkpoint at either epoch 60 or epoch 1,000 converges to fairly good natural accuracy ($> 55\%$). The failure of training with random labels proves that the good accuracy has nothing to do with the original natural accuracy of the AT checkpoint.

483 *adversarial training do contain non-robust features w.r.t. the classes to which they are misclassified,*
484 which strongly corroborates the validity of the experiments in Section 3.

485 **Discussion on the Influence Brought by Robust Features During Standard Training.** Since
486 untargeted PGD attack cannot be assigned with a target label, we cannot guarantee that the misclassi-
487 fied labels \hat{y} to be uniformly distributed regardless of the original labels (especially for real-world
488 datasets). This implies that the non-robust features are not completely decoupled from robust features,
489 *i.e.*, training on $\{(\hat{x}_i, \hat{y}_i)\}$ may take advantage of \hat{x}_i 's robust features from y_i through the correlation
490 between y_i and \hat{y}_i . However, we argue that robust features from y_i mingling with non-robust features
491 from \hat{y}_i only increases the difficulty of obtaining a good natural accuracy. This is because learning
492 through the shortcut will only wrongly map the robust features from y_i to \hat{y}_i that never equals to y_i ,
493 but during the evaluation, each clean test example x_i will always have robust and non-robust features
494 from y_i , and such wrong mapping will induce x_i to be misclassified to some \hat{y}_i due to the label of its
495 robust features. As a result, despite the negative influence brought by robust features, we still achieve
496 good natural accuracy at last, which further solidifies our conclusion.

497 B Experiment Details

498 In this section, we provide more experiment details that are omitted before due to the page limit.

499 B.1 Baseline Adversarial Training

500 In this paper, we mainly consider classification task on CIFAR-10 [17]. The dataset contains 60,000
501 32×32 RGB images from 10 classes. For each class, there are 5,000 images for training and 1,000
502 images for evaluation. Since we mainly aim to track the training dynamic to verify our understandings
503 in RO instead of formally evaluating an algorithm's performance, we do not hold a validation set in
504 most of our verification experiments and directly train on the full training set.

505 For baseline adversarial training, We use PreActResNet-18 [12] model as the classifier. We use
506 PGD-10 attack [22] with step size $\alpha = 2/255$ and perturbation norm $\varepsilon_\infty = 8/255$ to craft adversarial
507 examples on-the-fly. Following the settings in Madry et al. [22], we use SGD optimizer with
508 momentum 0.9, weight decay 5×10^{-4} and batch size 128 to train the model for as many as 1,000
509 epochs. The learning rate (LR) is initially set to be 0.1 and decays to 0.01 at epoch 100 and further
510 decays to 0.001 at epoch 150. For the version without LR decay used for comparison in our paper,
511 we simply keep the LR to be 0.1 during the whole training process.

512 Each model included in this paper is trained on a single NVIDIA GeForce RTX 3090 GPU. For
513 PGD-AT, it takes about 3d 14h to finish 1,000 epochs of training.

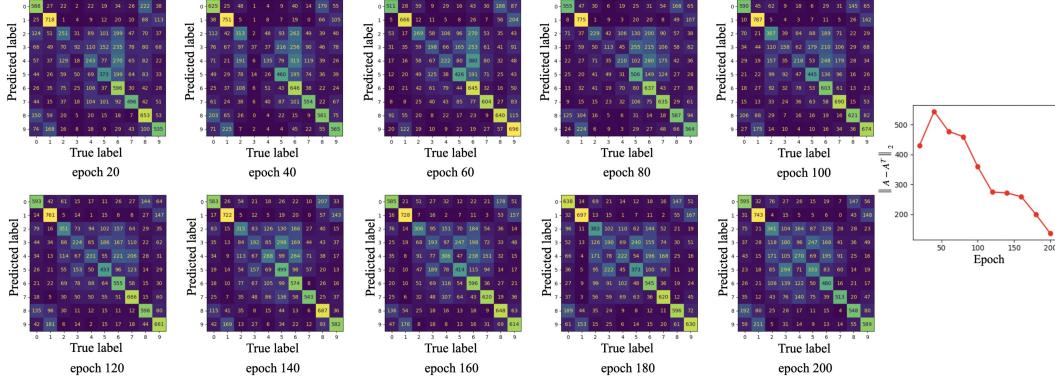


Figure 9: More test-time confusion matrices during the first 200 epochs of the training. After LR decays (the second row), the confusion matrix A immediately becomes symmetric, as the spectral norm $\|A - A^T\|_2$ w.r.t. the matrix decreases from ≥ 350 before epoch 100 to ≤ 150 after epoch 200.

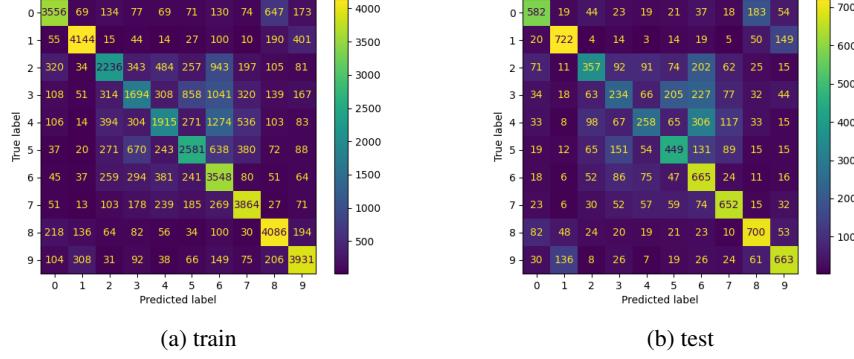


Figure 10: Confusion matrices of the training and the test data at an epoch before robust overfitting starts. They show nearly a same pattern of attacking preference among classes.

514 B.2 Verification Experiments for Our Minimax Game Perspective on Robust Overfitting

515 **Memorization of Adversarial Non-robust Features After LR Decay.** We craft adversarial examples
 516 on a checkpoint before LR decay (60th epoch) and evaluate the *misclassified* ones on a checkpoint
 517 after LR decay (≥ 150 th epoch) with their correct labels to evaluate the memorization of non-robust
 518 features in the training adversarial examples (see Appendix A for detailed discussions) in Section
 519 3.1 and 3.2. Following Appendix A, we adopt PGD-20 attack with perturbation norm $\varepsilon_\infty = 16/255$
 520 to craft adversarial examples which is stronger than the common attack setting we use in PGD-AT.
 521 We note that test adversarial examples crafted by a stronger attack indicates stronger extraction of
 522 the non-robust features, so they are more indicative of non-robust feature memorization when still
 523 correctly classified.

524 **Verification I: More Non-robust Features, Worse Robustness.** At the beginning of Section 3.2,
 525 we create synthetic datasets to demonstrate that memorizing the non-robust training features indeed
 526 harms test-time model robustness. To instill non-robust features into the training dataset, we minimize
 527 the adversarial loss w.r.t. the training data in a way that just like PGD attack, with the only difference
 528 that we minimize the adversarial loss instead of maximizing it. Since we only use a very small
 529 perturbation norm $\varepsilon_\infty \leq 4/255$, the added features are bound to be non-robust. For a fair comparison,
 530 we also perturb the training set with random uniform noise of the same perturbation norm to exclude
 531 the influence brought by (slight) data distribution shifts. We continue training from the 100-th baseline
 532 AT checkpoint (before LR decay) on each synthetic dataset for 10 epochs, and then evaluate model
 533 robustness with clean test data.

534 **Verification II: Vanishing Target-class Features in Test Adversarial Examples.** This is to say that
 535 when RO happens, we expect that test adversarial examples become less informative of the classes to

536 which they are misclassified according to our theory. To verify this, we first craft adversarial examples
 537 on a checkpoint T after RO begins, then evaluate the misclassified ones with their misclassified labels
 538 \hat{y} on the checkpoint saved at epoch 60. As a result, the accuracy reflects how much information
 539 (non-robust features) from \hat{y} the adversarial examples have to contain to be misclassified to \hat{y} on T .
 540 All adversarial examples evaluated in the experiments in Section 3.2.2 are crafted using PGD-10
 541 attack with perturbation norm $\epsilon_\infty = 8/255$.

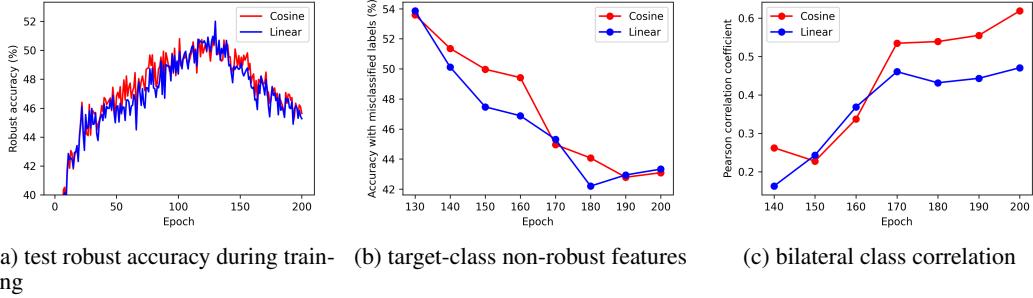
542 **Verification III: Bilateral Class Correlation.** To quantitatively analyze the correlation strength of
 543 bilateral misclassification described in Section 3.2.2, we first summarize all $y_i \rightarrow y_j$ misclassification
 544 rates into two confusion matrices $P^{\text{train}}, P^{\text{test}} \in \mathbb{R}^{C \times C}$ for the training and test data, respectively.
 545 Because we are mainly interested in the effect of the LR decay, we focus on the relative change on
 546 the test confusion matrix before and after LR decay, i.e., $\Delta P^{\text{test}} = P_{\text{after}}^{\text{test}} - P_{\text{before}}^{\text{test}}$. According
 547 to our theory, for each class pair (i, j) , there should be a strong correlation between the training
 548 misclassification of $i \rightarrow j$ before LR decay, i.e., $(P_{\text{before}}^{\text{train}})_{ij}$, and the increase in test misclassification
 549 of $j \rightarrow i$, i.e., $\Delta P_{ji}^{\text{test}}$, as $i \rightarrow j$ training misclassification (may due to intrinsic class bias, as will be
 550 further discussed below in details) induces $i \rightarrow j$ false mappings and creates $j \rightarrow i$ shortcuts. To
 551 examine their relationship, we plot the two variables $((P_{\text{before}}^{\text{train}})_{ij}, \Delta P_{ji}^{\text{test}})$ and compute their Pearson
 552 correlation coefficient ρ .

553 **Verification IV: Symmetrized Confusion Matrix.** In Section 3.2.2, we mention the growing
 554 symmetry of the test-time confusion matrix after LR decay as an evidence of the strengthening
 555 $y \rightarrow y'$ and $y' \rightarrow y$ correlation. Here we present more confusion matrices during the 200 epochs
 556 of the training in Figure 9, and it is very clear that the confusion matrices soon become symmetric
 557 after LR decay and RO starts. For a deeper comprehension of this phenomenon, we first visualize the
 558 confusion matrices of the training and the test data at an epoch before RO starts in Figure 10. They
 559 exhibit nearly a same pattern of attacking preference among classes (e.g., $y \rightarrow y'$) due to the bias
 560 rooted in the dataset, e.g., class 6 is intrinsically vulnerable in this case. For the test data, this intrinsic
 561 bias wouldn't be wiped out through learning due to the non-generalizability of the memorization of
 562 non-robust features in the training data, as discussed at the beginning of Section 3.2 (i.e., $y \rightarrow y'$
 563 bias still holds); and for the training data, this biased feature memorization will open shortcuts for
 564 test-time adversarial attack as discussed in Section 3.2 (i.e., $y' \rightarrow y$ begins). Combining both $y \rightarrow y'$
 565 and $y' \rightarrow y$, we arrive at the symmetry of test-time confusion matrix.

566 **Additional Results: Changing LR Decay Schedule.** Our discussion above is based on the piecewise
 567 LR decay schedule, in which the sudden decay of LR most obviously reflects our understandings.
 568 Besides, we also explore other LR decay schedules, including Cosine/Linear LR decay, to check
 569 whether different LR decay schedules will affect the observations and claims we made in this paper.
 570 For each schedule, we train the model for 200 epochs following the settings in Rice et al. [26]. As
 571 demonstrated in Figure 11a, we arrive at the same finding as Rice et al. [26] that with Cosine/Linear
 572 LR decay schedule, the training still suffers from severe RO after epoch 130. Then, we rerun
 573 the empirical verification experiments in Section 3.2.2 and find that under both the two LR decay
 574 schedules 1) the test adversarial examples indeed contain less and less target-class non-robust features
 575 as the training goes and RO becomes severer and severer (Figure 11b), 2) the bilateral class correlation
 576 becomes increasingly strong (Figure 11c) and 3) the confusion matrix indeed becomes symmetric
 577 (Figure 12). The results are almost the same as the results achieved when we adopt piecewise LR
 578 decay schedule because even though these LR decay schedules are mild, the LR eventually becomes
 579 small and makes the trainer \mathcal{T} overly strong to memorize the harmful non-robust features, indicating
 580 that our understandings in the cause of RO is fundamental and regardless of whatever LR decay
 581 schedule is used.

Table 3: Training with stronger attack and evaluating model robustness on CIFAR-10 under the perturbation norm $\epsilon_\infty = 8/255$ based on the PreActResNet-18 architecture.

Attack Strength	Natural			PGD-20			AutoAttack		
	best	final	diff	best	final	diff	best	final	diff
$\epsilon = 8/255$, PGD-10 (baseline)	83.50	84.94	-1.44	55.05	47.60	7.45	49.89	43.83	6.06
$\epsilon = 10/255$, PGD-12	80.66	82.48	-1.82	56.63	50.63	6.00	50.89	46.13	4.76
$\epsilon = 12/255$, PGD-15	78.17	80.25	-2.08	57.09	53.13	3.96	50.99	47.66	3.33
$\epsilon = 14/255$, PGD-17	73.92	76.70	-2.78	56.42	54.04	2.38	50.28	48.53	1.75
$\epsilon = 16/255$, PGD-20	69.51	73.11	-3.60	55.09	54.27	0.82	49.58	48.53	1.05



(a) test robust accuracy during training (b) target-class non-robust features (c) bilateral class correlation

Figure 11: Empirical verification of our explanation for robust overfitting when Cosine/Linear LR decay schedule is applied. (a) With Cosine/Linear LR decay schedule, the training still suffers from severe RO. (b) After RO begins, non-robust features in the test data become less and less informative of the classes to which they are misclassified. (c) Increasingly strong correlation between training-time $y \rightarrow y'$ misclassification and test-time $y' \rightarrow y$ misclassification increase.

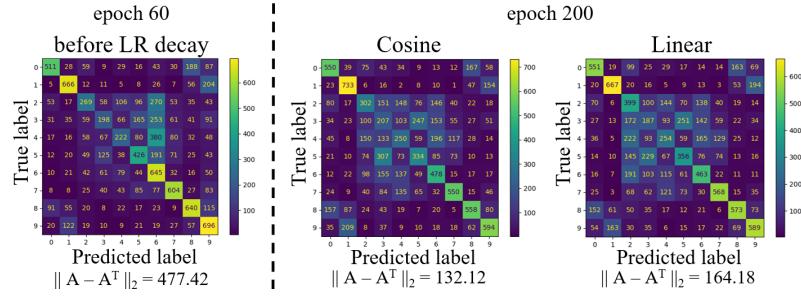


Figure 12: Test-time confusion matrix also becomes symmetric and implies that the bilateral correlation also exists when Cosine/Linear LR decay schedule is applied.

582 B.3 Experiments on the Effect of Stronger Training Attacker

583 In Section 4.2, we point out that using a stronger attacker in AT is able to mitigate RO to some extent
584 by neutralizing the trainer \mathcal{T} 's fitting power when it is overly strong. To achieve the results reported
585 in Figure 5d, we craft adversarial examples on-the-fly with more PGD iteration steps when ε is larger
586 (see Table 3), and further evaluate the best and last robustness of the WA models against PGD-20
587 and AA. Although RO is only partially mitigated and natural accuracy decreases when a stronger
588 attacker is applied as summarized in Addepalli et al. [1], it may be surprising to find from Table 3
589 that an attacker of appropriate strength may significantly boost the best WA robustness. This suggests
590 using a stronger attack could potentially be an interesting new path to stronger adversarial defense,
591 and we leave it for future work

592 B.4 Detailed Experiment Setup of ReBAT for Mitigating Robust Overfitting

593 **Datasets.** Beside CIFAR-10, we also include CIFAR-100 [17] and Tiny-ImageNet [6] for evaluation
594 of the effectiveness of ReBAT. CIFAR-100 shares the same training and test images with CIFAR-10,
595 but it classifies them into 100 categories, *i.e.*, 500 training images and 100 test images for each
596 class. Tiny-ImageNet is a subset of ImageNet [6] which contains labeled 64×64 RGB images from
597 200 classes. For each class, it includes 500 and 50 images for training and evaluation respectively.
598 Following Rice et al. [26], we hold out 1,000 images from the original CIFAR-10/100 training set,
599 and similarly 2,000 images from the original Tiny-ImageNet training set as validation sets.

600 **Training Strategy.** For CIFAR-10 and CIFAR-100, we follow exactly the same training strategy
601 as introduced in Appendix B.1, except that for ReBAT[strong] we adopt PGD-12 with perturbation
602 norm $\varepsilon_\infty = 10/255$ for training after LR decay. For Tiny-ImageNet, we follow the learning schedule
603 of Chen et al. [4], in which the model is trained for a total of 100 epochs and the LR decays twice (by
604 0.1) at epoch 50 and 80.

605 **Choices of Hyperparameters.** For KD+SWA [4], PGD-AT+TE [8], AWP [33] and WA+CutMix
606 [25], we strictly follow their original settings of hyperparameters. For MLCATWP [34], we simply

607 report the test results reported in their paper. Following Chen et al. [4], SWA/WA as well as the
 608 ReBAT regularization purposed in Section 4.1 start at epoch 105 (5 epochs later than the first LR
 609 decay where robust overfitting often begins), and following Rebuffi et al. [25] we choose the EMA
 610 decay rate of WA to be $\gamma = 0.999$. Please refer to Table 4 for our choices of the decay factor d and
 611 regularization strength λ . We notice that since CutMix improves the difficulty of learning, the model
 612 demands a relatively larger decay factor to better fit the augmented data. For Tiny-ImageNet, we also
 613 apply a larger λ after the second LR decay to better maintain the flatness of adversarial loss landscape
 614 and control robust overfitting. We provide more discussions on the choice of hyperparameters in
 615 Appendix C.1.

Table 4: Choices of hyperparameters when training models on different datasets using different network architectures with ReBAT.

Network Architecture	Method	CIFAR-10/CIFAR-100	Tiny-ImageNet
PreActResNet-18	ReBAT	$d = 1.5, \lambda = 1.0$	$d = 4.0, \lambda_1 = 2.0, \lambda_2 = 10.0$
	ReBAT[strong]	$d = 1.7, \lambda = 1.0$	$d = 4.0, \lambda_1 = 2.0, \lambda_2 = 10.0$
	ReBAT+CutMix	$d = 4.0, \lambda = 2.0$	$d = 6.0, \lambda = 1.5$
WideResNet-34-10	ReBAT	$d = 1.3, \lambda = 0.5$	-
	ReBAT[strong]	$d = 1.3, \lambda = 0.5$	-
	ReBAT+CutMix	$d = 4.0, \lambda = 2.0$	-

616 B.5 Training Robust Accuracy

617 Figure 13 shows the robust accuracy change on the training data during the training. Compared
 618 with vanilla PGD-AT that yields training robust accuracy over 80% at epoch 200, ReBAT manages
 619 to suppress it to only 65%. It successfully prevents the trainer \mathcal{T} from learning the non-robust
 620 features *w.r.t.* the training data too fast and too well, and therefore significantly reduces the robust
 621 generalization gap (from $\sim 35\%$ to $\sim 9\%$) and mitigates RO.

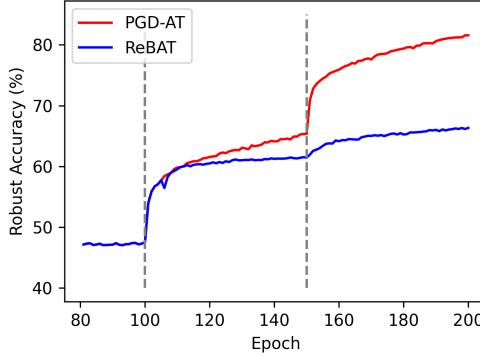


Figure 13: Training robust accuracy of PGD-AT and ReBAT on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$ based on the PreActResNet-18 architecture.

622 B.6 Training Efficiency

623 We also test and report the training time (per epoch) of several methods evaluated in this paper. For a
 624 fair comparison, all the compared methods are integrated into a universal training framework and
 625 each test runs on a single NVIDIA GeForce RTX 3090 GPU.

626 From Table 5, we can see that ReBAT requires nearly no extra computational cost compared with
 627 vanilla PGD-AT (136.2s v.s. 131.6s per epoch), implying that it is an efficient training method in
 628 practical. We also remark that KD+SWA, one of the most competitive methods that aims to address
 629 the RO issue, is not really computationally efficient as it requires to pretrain a robust classifier and a
 630 non-robust one as AT teacher and ST teacher respectively.

Table 5: Combining training time per epoch on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$ based on the PreActResNet-18 architecture.

Method	Training Time per Epoch (s)
PGD-AT	131.6
WA	132.1
KD+SWA	131.6+16.5+141.7
AWP	142.8
MLCAT _{wp}	353.3
ReBAT	136.2
WA+CutMix	168.6
ReBAT+CutMix	173.1

631 C More Experiments on ReBAT

632 In this section, we conduct extensional experiments on the proposed ReBAT method to further
633 demonstrate its effectiveness, efficiency and flexibility.

634 C.1 Additional Results on BoAT Loss

635 In Section B.4, we discuss the detailed configurations for the experiments in Figure 5a, where we
636 show that BoAT can largely mitigate robust overfitting. Here, we further summarize the performance
637 of best and final checkpoints of the original AT+WA method and our BoAT. As shown in Table 6,
638 BoAT not only boosts the best robustness by a large margin (0.67% higher against AA) but also
639 significantly suppresses RO (1.58% v.s. 6.06% against AA). To achieve the reported robustness, we
640 first use $\lambda_1 = 10.0$ after the first LR decay and then apply $\lambda_2 = 60.0$ after the second LR decay to
641 better maintain the flatness of adversarial loss landscape and control robust overfitting.

Table 6: Comparing model robustness w/ and w/o BoAT loss on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$ based on the PreActResNet-18 architecture.

Method	Natural			PGD-20			AutoAttack		
	best	final	diff	best	final	diff	best	final	diff
AT+WA	83.50	84.94	-1.44	55.05	47.60	7.45	49.89	43.83	6.06
ReBAT($d = 10$)	81.54	82.42	-0.88	55.29	53.43	1.86	50.56	48.98	1.58

642 C.2 Additional Results on the Effect of LR decay

643 Below we show that when a relatively large decay factor d is applied, i.e., the model has overly strong
644 fitting ability that results in robust overfitting, a large regularization coefficient λ should be chosen
645 for better performance. Table 7 reveals this relationship between d and λ . When $d = 1.3$, even a λ
646 as small as 1.0 will harm both the best and last robustness as well as natural accuracy, as $d = 1.3$
647 is already too small a decay factor that makes the model suffering from underfitting and naturally
648 requiring no more flatness regularization. On the other side, when d is relatively large, even a strong
649 regularization of $\lambda = 4.0$ is not adequate to fully suppress RO. Besides, comparing the situation of
650 $\lambda > 0$ and $\lambda = 0$ for a fixed $d(\geq 1.5)$, we emphasize that the purposed BoAT loss again exhibits its
651 apparent effectiveness in simultaneously boosting the best robustness and mitigating RO.

652 C.3 Additional Results on Adopting Stronger Training Time Attacker

653 In Section 4.3, we design two versions of ReBAT, and we name the stronger one that also uses a
654 stronger training attacker as ReBAT[strong]. Here we fix the hyperparameter used in ReBAT above
655 (note that we use a larger LR decay factor $d = 1.7$ in ReBAT[strong] for CIFAR-10, and now we still
656 use $d = 1.5$ as in ReBAT) and adjust the attacker strength to study its influence when combined with
657 ReBAT. According to Figure 14, though a slightly stronger attack (e.g., $\varepsilon = 9/255$) may marginally
658 improves the best and last robust accuracy, it heavily degrades natural accuracy, particularly when
659 much stronger attack is used. We deem that this is because it breaks the balance from the other side

Table 7: Changing decay factor d and regularization strength λ and evaluating model robustness on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$ based on the PreActResNet-18 architecture.

Method	Natural			PGD-20			AutoAttack		
	best	final	diff	best	final	diff	best	final	diff
ReBAT($d = 1.3, \lambda = 0.0$)	81.17	81.27	-0.10	56.43	56.23	0.20	50.80	50.75	0.05
ReBAT($d = 1.3, \lambda = 1.0$)	80.67	80.63	0.04	56.27	56.15	0.12	50.74	50.65	0.09
ReBAT($d = 1.3, \lambda = 4.0$)	78.50	78.36	0.14	55.10	55.04	0.06	50.31	50.30	0.01
ReBAT($d = 1.5, \lambda = 0.0$)	81.90	82.39	-0.49	56.21	55.95	0.26	50.81	50.58	0.23
ReBAT($d = 1.5, \lambda = 1.0$)	81.86	81.91	-0.05	56.36	56.12	0.24	51.13	51.22	-0.09
ReBAT($d = 1.5, \lambda = 4.0$)	79.68	79.88	-0.20	55.72	55.65	0.07	50.52	50.50	0.02
ReBAT($d = 4.0, \lambda = 0.0$)	83.05	85.38	-2.33	55.98	50.87	5.11	50.38	46.95	3.43
ReBAT($d = 4.0, \lambda = 1.0$)	82.46	84.84	-2.38	55.86	52.02	3.84	50.87	47.88	2.99
ReBAT($d = 4.0, \lambda = 4.0$)	80.99	84.07	-3.08	56.06	53.58	2.48	51.00	49.02	1.98

660 that the overly strong attacker \mathcal{A} dominates the adversarial game and results in an underfitting state
661 that harms both robust and natural accuracy.

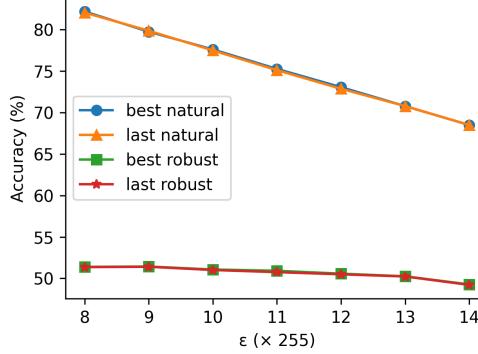


Figure 14: Using different adversarial attack strength in ReBAT.

662 C.4 Further Improving Natural Accuracy with Knowledge Distillation

663 Chen et al. [4] propose to adopt knowledge distillation (KD) [13] to mitigate RO and it is worth
664 mentioning that their method achieves relatively good natural accuracy according to Table 1 and 2.
665 Since our ReBAT method is orthogonal to KD, we propose to combine our techniques with KD to
666 further improve natural accuracy. Specifically, we simplify their method by using only a non-robust
667 standard classifier as a teacher (ST teacher) instead of using both a ST teacher and a AT teacher,
668 because i) a large sum of computational cost for training the AT teacher will be saved, ii) our main
669 goal is to improve natural accuracy so the ST teacher matters more and iii) ReBAT already use the
670 WA model as a very good teacher. This gives the final loss function as

$$\ell_{\text{ReBAT+KD}}(x, y; \theta) = (1 - \lambda_{\text{ST}}) \cdot \ell_{\text{BoAT}}(x, y; \theta) + \lambda_{\text{ST}} \cdot \text{KL}(f_\theta(x) \| f_{\text{ST}}(x)), \quad (4)$$

671 where f_{ST} indicates the ST teacher and λ_{ST} is a trade-off parameter.

Table 8: Combining our methods with knowledge distillation and evaluating model robustness on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$ based on the PreActResNet-18 architecture.

Method	Natural			PGD-20			AutoAttack		
	best	final	diff	best	final	diff	best	final	diff
ReBAT ($\lambda_{\text{ST}} = 0.0$)	81.86	81.91	-0.05	56.36	56.12	0.24	51.13	51.22	-0.09
ReBAT+KD ($\lambda_{\text{ST}} = 0.4$)	83.59	83.64	-0.05	54.91	54.77	-0.14	50.70	50.99	-0.29
ReBAT+KD ($\lambda_{\text{ST}} = 0.5$)	84.12	84.20	-0.08	55.28	55.39	-0.11	50.47	50.72	-0.25
ReBAT+KD ($\lambda_{\text{ST}} = 0.6$)	84.34	84.72	-0.38	53.83	54.30	-0.47	50.15	50.37	-0.22

672 Table 8 compares the performance of ReBAT+KD when different λ_{ST} is applied, and clearly a large
 673 λ_{ST} results in improvement in natural accuracy and decreases robustness (may due to the theoretically
 674 principled trade-off between natural accuracy and robustness [37]). However, it is still noteworthy
 675 that when $\lambda_{ST} = 0.4$, a notable increase in natural accuracy ($\sim 1.7\%$) is achieved at the cost of only
 676 a small slide of $\sim 0.2\%$ in final robustness against AA. Also when $\lambda_{ST} = 0.5$, it achieves comparable
 677 natural accuracy with KD+SWA [4] but has much higher AA robustness. Moreover, we emphasize
 678 that RO is almost completely eliminated regardless of the trade-off, which is the main concern of
 679 this paper and demonstrates the superiority of our method against previous ones. An intriguing
 680 phenomenon is that nearly all the final results are better than the results on the best checkpoints
 681 selected by the validation set, which implies that in this training scheme AT enjoys the same property
 682 of “training longer, generalize better” as ST without any need of early stopping.

683 C.5 The Effect of Different Learning Rate Schedules

684 In the previous experiments we only investigate the piecewise LR decay schedule. However, a natural
 685 idea would be using mild LR decay schedules, e.g., Cosine and Linear decay schedule, instead of
 686 suddenly decaying it by a factor of d at some epoch in the piecewise decay schedule. As mentioned in
 687 Section 5, previous works have shown that changing LR decay schedule fails to effectively suppress
 688 RO whether with [31] or without WA [26] because the LR finally becomes small and endows the
 689 trainer with overly strong fitting ability. Therefore, here we continue to experiment with modified
 690 Cosine and Linear decay schedules that follow a similar LR scale of the piecewise LR decay schedule,
 691 and summarise the results in Table 9. To be specific, the LR still decays to 0.01 at epoch 150 and
 692 0.001 at epoch 200, though the two decay stages (from epoch 100 to 150 and 150 to 200) are designed
 693 to be gradual following the Cosine/Linear schedule. We also gradually increase the strength of
 694 ReBAT regularization from zero as the LR gradually decreases, following the “larger decay factor
 695 goes with stronger ReBAT regularization” principle that we introduced in Appendix C.2.

Table 9: Comparing our method with WA on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$
 based on PreActResNet-18 architecture, when Cosine/Linear LR decay schedule are applied.

Method	Cosine				Linear			
	Natural		PGD-20		AutoAttack		Natural	
	best	final	best	final	best	final	best	final
WA(d=10)	82.32	84.99	55.97	47.84	50.59	44.08	82.21	85.01
ReBAT(d=10)	82.06	82.22	56.12	54.44	50.98	49.81	81.98	82.13
							56.33	54.50
							51.08	49.63

696 It can be concluded from Table 9 that simply changing the LR decay schedule indeed improves the
 697 best robust accuracy from 49.89% to 50.59% and 50.67% against AA respectively, but it provides no
 698 help at all in mitigating RO as the final robust accuracy is still below 45% against AA. We also note
 699 that in this situation, the application of BoAT loss not only significantly mitigates RO but also further
 700 improves the best model robustness, which also proves its effectiveness.

701 C.6 Results on Different Network Architectures

702 In previous experiments we compare methods based on PreActResNet-18 and WideResNet-34-
 703 10 architecture, and here we also adopt VGG-16 [28] and MobileNetV2 [27] architecture. The
 704 significant improvement against baseline PGD-AT in both best and final robust accuracy and in
 705 mitigating RO reported in Table 10 further demonstrates that our method works on a wide range of
 706 network architectures.

Table 10: Comparing our method with PGD-AT on CIFAR-10 under the perturbation norm $\varepsilon_\infty = 8/255$
 based on VGG-16 and MobileNetV2 architecture.

Method	VGG-16				MobileNetV2			
	Natural		PGD-20		AutoAttack		Natural	
	best	final	best	final	best	final	best	final
PGD-AT	78.43	81.64	50.56	44.25	44.19	39.66	79.85	80.67
ReBAT	78.17	78.37	53.13	53.01	47.24	47.13	78.98	80.81
							53.18	52.57
							47.66	47.35