

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 THE ATTACK MEANS NOTHING: TEST-TIME ADVER- SARIAL DEFENSE IMPROVES ZERO-SHOT ADVERSAR- IAL ROBUSTNESS FOR MEDICAL VISION-LANGUAGE MODELS

014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
**Anonymous authors**  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
Paper under double-blind review

## ABSTRACT

014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
Vision-language models (VLMs), exemplified by CLIP, have achieved remarkable zero-shot generalization but remain highly vulnerable to imperceptible adversarial perturbations, posing significant safety threats, particularly in medical scenarios. In this paper, we first prove that VLMs are much more robust than adversarial attacks when faced with weak transformations. Building upon this insight, we propose the **The Attack Means Nothing** (TAME), a simple yet effective test-time defense paradigm for improving the zero-shot adversarial robustness of medical VLMs. We conduct comprehensive experiments on 11 medical datasets across 9 imaging modalities against three representative white-box attacks (PGD, C&W, and AutoAttack). The BiomedCLIP with a backbone of ViT-B/16 is utilized as the victim model. Extensive experiment results demonstrate that our TAME consistently outperforms other defense methods across all attack types, boosting the vanilla BiomedCLIP by +47.47% under PGD, +46.73% under C&W, and +47.79% under AutoAttack, while maintaining competitive clean accuracy. These significant improvements also suggest a potential risk of label leakage during attacks. Furthermore, our TAME is plug-and-play and can be integrated with other adversarially fine-tuned VLMs to enhance their defense capabilities. These findings support a practical and generalizable approach to deploying medical VLMs in clinical scenarios with the presence of adversaries. Codes will be available on GitHub.

## 1 INTRODUCTION

037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
Recent advancements in vision-language models (VLMs) Zhao et al. (2025); Lai et al. (2024) have demonstrated significant success and potential for medical image analysis Koleilat et al. (2025); Stevens et al. (2024). Unlike traditional supervised learning focused on closed-set tasks, VLMs, such as Contrastive Language-Image Pre-training (CLIP) Radford et al. (2021a), enable the exploration of open-set visual concepts, yielding strong zero-shot generalization capabilities. Unfortunately, some studies Zhang et al. (2022); Zhao et al. (2023); Yin et al. (2023) reveal that adding even imperceptible adversarial perturbations to input images can severely degrade VLM's inference ability. This poses critical safety risks, especially in medical scenarios Dong et al. (2024), which may lead to serious misdiagnosis and hinder models from being deployed in real-world applications (see Figure 1).

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
Extensive research has explored adversarial training Chen et al. (2020) as an effective defense strategy, which can be broadly categorized into two categories: adversarial fine-tuning (AFT) Mao et al. (2023); Wang et al. (2024a); Schlarbmann et al. (2024); Wang et al. (2024c) and adversarial prompt tuning (APT) Li et al. (2024); Zhang et al. (2024); Zhou et al. (2024); Wang et al. (2024b); Zhou et al. (2024). AFT methods aim to establish a min-max game between the VLM and an adversary, fine-tuning the pre-trained VLM on generated adversarial examples to achieve transferable robustness across downstream tasks. However, most of these methods require substantial computational resources and inevitably degrade the model's generalization to testing data from unseen distributions. APT methods attempt to train learnable textual or visual prompts by aligning adversarial image embeddings with corresponding text prompts while keeping the model backbone frozen. Al-

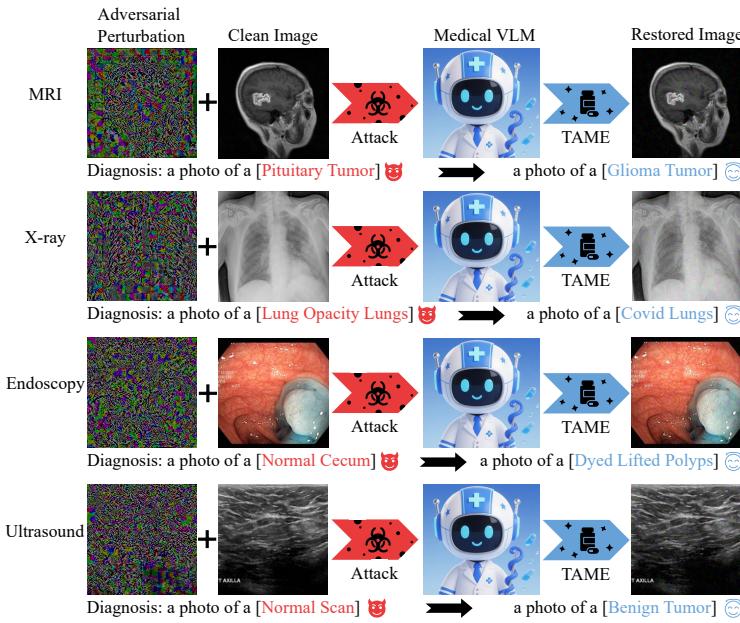


Figure 1: Adversarial attacks disturb model inference by adding imperceptible perturbations to the input image, leading to serious misdiagnosis. Our TAME enables the medical VLM to remain robust against adversarial attacks during inference without extra training on predefined adversarial data.

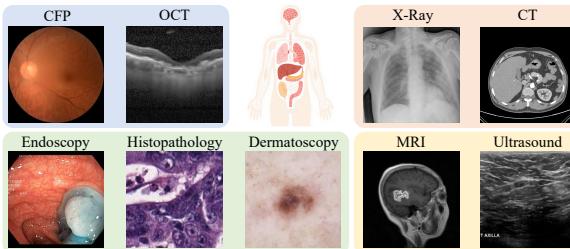


Figure 2: The medical imaging modalities used in this study. CFP: Color fundus photography. OCT: Optical coherence tomography. CT: Computed tomography. MRI: Magnetic resonance imaging.

though they reduce training costs, their effectiveness is constrained by predefined data distributions, limiting the adaptability to out-of-distribution environments. Consequently, achieving low-cost and effective adversarial robustness remains an open challenge.

Test-time adversarial defense (TAD) Alfarra et al. (2022); Pérez et al. (2021); Wu et al. (2021); Guo et al. (2018); Xing et al. (2025); Wang et al. (2025); Mao et al. (2021) has emerged as a promising paradigm to boost zero-shot adversarial robustness in a low-cost manner, as TAD requires only test data during the inference phase. Training-free TAD methods Pérez et al. (2021); Guo et al. (2018) assemble several image transformations to make it difficult for adversaries to circumvent the defense strategy. Training-based approaches Alfarra et al. (2022); Xing et al. (2025); Wu et al. (2021); Sheng et al. (2025); Mao et al. (2021) mainly focus on modifying the input image or training a prompt to counteract attacks. Despite their efforts, almost all existing TAD methods are designed for conventional networks like convolutional neural networks, with insufficient exploration of VLM. Furthermore, medical VLMs are typically utilized to process a wide range of modalities, as illustrated in Figure 2, posing a challenge to the defense method's generalizability across various modalities.

The key to addressing these issues is to identify commonalities of adversarial images to establish a general test-time defense paradigm for VLMs. In this paper, we first conduct a toy experiment on multiple datasets by applying several transformations to both clean and adversarial images. We observe that although transformations with large magnitude significantly disturb model predictions on

108 both images, this abnormal effect still appears on the adversarial images even under transformations  
 109 of minor magnitude. We term this phenomenon ‘**semantic fragility**’ of adversarial perturbations,  
 110 which can be interpreted as these perturbations being highly specific to the corresponding input  
 111 images. The VLM, trained on extensive and diverse data, exhibits inherent robustness to such mi-  
 112 nor transformations on clean images. In contrast, adversarial perturbations are over-fitted to both the  
 113 specific input and the current model parameters, rendering the semantic content within the perturbed  
 114 image embeddings highly susceptible to even slight alteration. Based on this observation, we pro-  
 115 pose **The Attack Means Nothing (TAME)**, a test-time defense paradigm for medical VLMs. TAME  
 116 counteracts adversarial attacks by training an adversarial restoration map for each adversarial im-  
 117 age in a single iteration. Specifically, we first introduce the adversarial restoration map to the input  
 118 image to produce the restored image and then minimize the KL divergence between the predicted  
 119 probability distributions of each restored image and its transformed version. Through this training  
 120 process, the trained adversarial restoration map learns to restore the model’s invariance to minor  
 121 transformations, thereby reinstating the inference capability of VLMs. Furthermore, the adversarial  
 122 restoration map should also minimize its effect on clean images, thereby avoiding significant per-  
 123 formance degradation induced by excessive image modification. To address this issue, we design a  
 124 dynamic weighting mechanism that adaptively allocates weights according to the degree of semantic  
 125 fragility exhibited by the input image. Comprehensive experiments are conducted across 11 medical  
 126 classification datasets, including 9 modalities (see Figure 2), to evaluate TAME and other state-of-  
 127 the-art methods against three typical adversaries (*i.e.*, Projected Gradient Descent PGD Madry et al.  
 128 (2018), C&W Carlini & Wagner (2017), and AutoAttack Croce & Hein (2020)) that aim to max-  
 129 imize the classification loss in a white-box setting. Extensive experiment results demonstrate the  
 130 effectiveness and superiority of our TAME across diverse scenarios.

131 The three key distributions of this paper are summarized as follows.

- 132 • We propose a simple yet effective method to enhance the zero-shot adversarial robustness  
 133 of medical Vision-Language Models (VLMs), which can be utilized as a plug-and-play  
 134 module without additional training.
- 135 • Based on observed commonalities in adversarial images, we propose TAME to protect  
 136 VLMs against multiple attack types alongside a dynamic weighting mechanism main-  
 137 taining performance on clean images.
- 138 • Extensive experiments on 11 medical classification datasets across 9 modalities demon-  
 139 strate the superiority of our TAME over other existing defense methods.

## 141 2 RELATED WORK

### 142 2.1 ADVERSARIAL TRAINING

143 Adversarial training enhances the adversarial robustness of the model by training on the predefined  
 144 adversarial samples, which can be broadly classified into adversarial fine-tuning Mao et al. (2023);  
 145 Wang et al. (2024a); Schlarmann et al. (2024); Wang et al. (2024c) and adversarial prompt tuning Li  
 146 et al. (2024); Zhang et al. (2024); Zhou et al. (2024); Wang et al. (2024b).

147 **Adversarial Fine-Tuning (AFT)** AFT improves the adversarial robustness by fine-tuning the VLM  
 148 on adversarial samples generated by an adversary. Mao *et al.* Mao et al. (2023) fine-tuned the  
 149 vision encoder of CLIP using adversarial contrastive learning with text-guided supervision on a  
 150 small set of adversarial samples. Wang *et al.* Wang et al. (2024a) proposed a pre-trained model  
 151 guided adversarial fine-tuning method, which distills the general knowledge from the original pre-  
 152 trained model to the target model to mitigate the over-fitting. Schlarmann *et al.* Schlarmann et al.  
 153 (2024) attempted to minimize the distance between the original and fine-tuned image embeddings  
 154 during adversarial training to preserve the performance of the fine-tuned model on clean data.

155 **Adversarial Prompt Tuning (APT)** APT learns trainable prompts to maintain alignment under  
 156 attack by exposing the model to adversarial samples while freezing the model parameters. Zhou *et*  
 157 *al.* Zhou et al. (2024) presented to learn adversarially correlated text supervision by enhancing the  
 158 consistency of multi-modal features and encouraging distinguishability between features of clean  
 159 and adversarial data. Zhang *et al.* Zhang et al. (2024) aligned learnable text prompts with adversarial  
 160 image embeddings to improve resistance against white-box and black-box adversarial attacks. Li *et*

162 *al.* Li et al. (2024) demonstrated the high sensitivity of both adversarial attacks and defenses to the  
 163 specific text prompts used in VLMs and proposed to improve adversarial robustness by learning  
 164 robust text prompts. Unlike these adversarial training methods, our TAME aims to achieve zero-  
 165 shot adversarial robustness using only test data in a low-cost and general manner, improving the  
 166 performance across various medical scenarios.

## 169 2.2 TEST-TIME ADVERSARIAL DEFENSE (TAD)

171 TAD aims to protect the pre-trained model from adversarial attacks in a low-cost manner during  
 172 inference, including two branches: training-free and training-based methods. Training-free meth-  
 173 ods Pérez et al. (2021); Guo et al. (2018) typically refer to designing the image transformation  
 174 strategy. Pérez et al. Pérez et al. (2021) proposed a transformation ensemble method achieving  
 175 consistent improvements in adversarial robustness across datasets and adversaries while preserving  
 176 clean data performance. Guo et al. Guo et al. (2018) found that total variance minimization and  
 177 image quilting are effective against several attacks, particularly on the model trained on such trans-  
 178 formation strategies. Training-based methods Alfarra et al. (2022); Xing et al. (2025); Wu et al.  
 179 (2021); Wang et al. (2025); Mao et al. (2021) primarily prevent the model from adversarial attacks  
 180 by modifying the input image or training a prompt. Alfarra et al. Alfarra et al. (2022) presented the  
 181 anti-adversary layer to generate a perturbed input image in the opposite direction of the adversarial  
 182 one to counter the attacks. Xing et al. Xing et al. (2025) maximized the classification loss on the  
 183 test image to counterattack adversaries and prevented further counterattacking on clean data using a  
 184 threshold. Sheng et al. Sheng et al. (2025) reformulated the marginal entropy objective to train the  
 185 textual prompts and proposed a reliability-weighted ensembling strategy that aggregates information  
 186 from trustworthy augmented views to enhance defense. Although TTC and R-TPT are designed for  
 187 VLMs, they still exhibit limitations: 1) TTC requires a hyperparameter to distinguish clean and  
 188 adversarial images during inference, leading to suboptimal performance in adversarial robustness  
 189 when misclassifications occur. 2) R-TPT relies on low-entropy predictions for pointwise entropy  
 190 minimization, however, adversarial images typically yield high-confidence but incorrect predictions,  
 191 potentially reinforcing wrong decisions. In contrast, our TAME explores the inherent defect of ad-  
 192 versarial attacks in VLMs and requires only one gradient backpropagation step. It employs a general  
 193 training objective for both clean and adversarial images to achieve zero-shot adversarial robustness.

## 194 3 METHODOLOGY

### 196 3.1 PRELIMINARIES

197 **Zero-shot inference of CLIP.** Let  $f_{\theta_v}$  and  $f_{\theta_t}$  denote the CLIP’s vision encoder and text encoder,  
 198 respectively, where  $\theta_v$  and  $\theta_t$  are their corresponding parameters. Given an input image  $I$  and a set  
 199 of  $k$  possible classes  $C = \{c_1, c_2, \dots, c_k\}$ ,  $I$  can be classified in a zero-shot manner by computing the  
 200 cosine similarity between the produced image embedding and the text embeddings of  $C$  wrapped in  
 201 a template(*e.g.*, “a photo of [CLASS]”). Specifically, the cosine similarity score between the image  
 202 embedding and the text embedding of  $i$ -th class  $c_i$  can be formulated as:

$$203 \quad 204 \quad 205 \quad 206 \quad 207 \quad 208 \quad 209 \quad 210 \quad 211 \quad 212 \quad 213 \quad 214 \quad 215 \quad 216 \quad 217 \quad 218 \quad 219 \quad 220 \quad 221 \quad 222 \quad 223 \quad 224 \quad 225 \quad 226 \quad 227 \quad 228 \quad 229 \quad 230 \quad 231 \quad 232 \quad 233 \quad 234 \quad 235 \quad 236 \quad 237 \quad 238 \quad 239 \quad 240 \quad 241 \quad 242 \quad 243 \quad 244 \quad 245 \quad 246 \quad 247 \quad 248 \quad 249 \quad 250 \quad 251 \quad 252 \quad 253 \quad 254 \quad 255 \quad 256 \quad 257 \quad 258 \quad 259 \quad 260 \quad 261 \quad 262 \quad 263 \quad 264 \quad 265 \quad 266 \quad 267 \quad 268 \quad 269 \quad 270 \quad 271 \quad 272 \quad 273 \quad 274 \quad 275 \quad 276 \quad 277 \quad 278 \quad 279 \quad 280 \quad 281 \quad 282 \quad 283 \quad 284 \quad 285 \quad 286 \quad 287 \quad 288 \quad 289 \quad 290 \quad 291 \quad 292 \quad 293 \quad 294 \quad 295 \quad 296 \quad 297 \quad 298 \quad 299 \quad 300 \quad 301 \quad 302 \quad 303 \quad 304 \quad 305 \quad 306 \quad 307 \quad 308 \quad 309 \quad 310 \quad 311 \quad 312 \quad 313 \quad 314 \quad 315 \quad 316 \quad 317 \quad 318 \quad 319 \quad 320 \quad 321 \quad 322 \quad 323 \quad 324 \quad 325 \quad 326 \quad 327 \quad 328 \quad 329 \quad 330 \quad 331 \quad 332 \quad 333 \quad 334 \quad 335 \quad 336 \quad 337 \quad 338 \quad 339 \quad 340 \quad 341 \quad 342 \quad 343 \quad 344 \quad 345 \quad 346 \quad 347 \quad 348 \quad 349 \quad 350 \quad 351 \quad 352 \quad 353 \quad 354 \quad 355 \quad 356 \quad 357 \quad 358 \quad 359 \quad 360 \quad 361 \quad 362 \quad 363 \quad 364 \quad 365 \quad 366 \quad 367 \quad 368 \quad 369 \quad 370 \quad 371 \quad 372 \quad 373 \quad 374 \quad 375 \quad 376 \quad 377 \quad 378 \quad 379 \quad 380 \quad 381 \quad 382 \quad 383 \quad 384 \quad 385 \quad 386 \quad 387 \quad 388 \quad 389 \quad 390 \quad 391 \quad 392 \quad 393 \quad 394 \quad 395 \quad 396 \quad 397 \quad 398 \quad 399 \quad 400 \quad 401 \quad 402 \quad 403 \quad 404 \quad 405 \quad 406 \quad 407 \quad 408 \quad 409 \quad 410 \quad 411 \quad 412 \quad 413 \quad 414 \quad 415 \quad 416 \quad 417 \quad 418 \quad 419 \quad 420 \quad 421 \quad 422 \quad 423 \quad 424 \quad 425 \quad 426 \quad 427 \quad 428 \quad 429 \quad 430 \quad 431 \quad 432 \quad 433 \quad 434 \quad 435 \quad 436 \quad 437 \quad 438 \quad 439 \quad 440 \quad 441 \quad 442 \quad 443 \quad 444 \quad 445 \quad 446 \quad 447 \quad 448 \quad 449 \quad 450 \quad 451 \quad 452 \quad 453 \quad 454 \quad 455 \quad 456 \quad 457 \quad 458 \quad 459 \quad 460 \quad 461 \quad 462 \quad 463 \quad 464 \quad 465 \quad 466 \quad 467 \quad 468 \quad 469 \quad 470 \quad 471 \quad 472 \quad 473 \quad 474 \quad 475 \quad 476 \quad 477 \quad 478 \quad 479 \quad 480 \quad 481 \quad 482 \quad 483 \quad 484 \quad 485 \quad 486 \quad 487 \quad 488 \quad 489 \quad 490 \quad 491 \quad 492 \quad 493 \quad 494 \quad 495 \quad 496 \quad 497 \quad 498 \quad 499 \quad 500 \quad 501 \quad 502 \quad 503 \quad 504 \quad 505 \quad 506 \quad 507 \quad 508 \quad 509 \quad 510 \quad 511 \quad 512 \quad 513 \quad 514 \quad 515 \quad 516 \quad 517 \quad 518 \quad 519 \quad 520 \quad 521 \quad 522 \quad 523 \quad 524 \quad 525 \quad 526 \quad 527 \quad 528 \quad 529 \quad 530 \quad 531 \quad 532 \quad 533 \quad 534 \quad 535 \quad 536 \quad 537 \quad 538 \quad 539 \quad 540 \quad 541 \quad 542 \quad 543 \quad 544 \quad 545 \quad 546 \quad 547 \quad 548 \quad 549 \quad 550 \quad 551 \quad 552 \quad 553 \quad 554 \quad 555 \quad 556 \quad 557 \quad 558 \quad 559 \quad 560 \quad 561 \quad 562 \quad 563 \quad 564 \quad 565 \quad 566 \quad 567 \quad 568 \quad 569 \quad 570 \quad 571 \quad 572 \quad 573 \quad 574 \quad 575 \quad 576 \quad 577 \quad 578 \quad 579 \quad 580 \quad 581 \quad 582 \quad 583 \quad 584 \quad 585 \quad 586 \quad 587 \quad 588 \quad 589 \quad 590 \quad 591 \quad 592 \quad 593 \quad 594 \quad 595 \quad 596 \quad 597 \quad 598 \quad 599 \quad 600 \quad 601 \quad 602 \quad 603 \quad 604 \quad 605 \quad 606 \quad 607 \quad 608 \quad 609 \quad 610 \quad 611 \quad 612 \quad 613 \quad 614 \quad 615 \quad 616 \quad 617 \quad 618 \quad 619 \quad 620 \quad 621 \quad 622 \quad 623 \quad 624 \quad 625 \quad 626 \quad 627 \quad 628 \quad 629 \quad 630 \quad 631 \quad 632 \quad 633 \quad 634 \quad 635 \quad 636 \quad 637 \quad 638 \quad 639 \quad 640 \quad 641 \quad 642 \quad 643 \quad 644 \quad 645 \quad 646 \quad 647 \quad 648 \quad 649 \quad 650 \quad 651 \quad 652 \quad 653 \quad 654 \quad 655 \quad 656 \quad 657 \quad 658 \quad 659 \quad 660 \quad 661 \quad 662 \quad 663 \quad 664 \quad 665 \quad 666 \quad 667 \quad 668 \quad 669 \quad 670 \quad 671 \quad 672 \quad 673 \quad 674 \quad 675 \quad 676 \quad 677 \quad 678 \quad 679 \quad 680 \quad 681 \quad 682 \quad 683 \quad 684 \quad 685 \quad 686 \quad 687 \quad 688 \quad 689 \quad 690 \quad 691 \quad 692 \quad 693 \quad 694 \quad 695 \quad 696 \quad 697 \quad 698 \quad 699 \quad 700 \quad 701 \quad 702 \quad 703 \quad 704 \quad 705 \quad 706 \quad 707 \quad 708 \quad 709 \quad 710 \quad 711 \quad 712 \quad 713 \quad 714 \quad 715 \quad 716 \quad 717 \quad 718 \quad 719 \quad 720 \quad 721 \quad 722 \quad 723 \quad 724 \quad 725 \quad 726 \quad 727 \quad 728 \quad 729 \quad 730 \quad 731 \quad 732 \quad 733 \quad 734 \quad 735 \quad 736 \quad 737 \quad 738 \quad 739 \quad 740 \quad 741 \quad 742 \quad 743 \quad 744 \quad 745 \quad 746 \quad 747 \quad 748 \quad 749 \quad 750 \quad 751 \quad 752 \quad 753 \quad 754 \quad 755 \quad 756 \quad 757 \quad 758 \quad 759 \quad 760 \quad 761 \quad 762 \quad 763 \quad 764 \quad 765 \quad 766 \quad 767 \quad 768 \quad 769 \quad 770 \quad 771 \quad 772 \quad 773 \quad 774 \quad 775 \quad 776 \quad 777 \quad 778 \quad 779 \quad 780 \quad 781 \quad 782 \quad 783 \quad 784 \quad 785 \quad 786 \quad 787 \quad 788 \quad 789 \quad 790 \quad 791 \quad 792 \quad 793 \quad 794 \quad 795 \quad 796 \quad 797 \quad 798 \quad 799 \quad 800 \quad 801 \quad 802 \quad 803 \quad 804 \quad 805 \quad 806 \quad 807 \quad 808 \quad 809 \quad 810 \quad 811 \quad 812 \quad 813 \quad 814 \quad 815 \quad 816 \quad 817 \quad 818 \quad 819 \quad 820 \quad 821 \quad 822 \quad 823 \quad 824 \quad 825 \quad 826 \quad 827 \quad 828 \quad 829 \quad 830 \quad 831 \quad 832 \quad 833 \quad 834 \quad 835 \quad 836 \quad 837 \quad 838 \quad 839 \quad 840 \quad 841 \quad 842 \quad 843 \quad 844 \quad 845 \quad 846 \quad 847 \quad 848 \quad 849 \quad 850 \quad 851 \quad 852 \quad 853 \quad 854 \quad 855 \quad 856 \quad 857 \quad 858 \quad 859 \quad 860 \quad 861 \quad 862 \quad 863 \quad 864 \quad 865 \quad 866 \quad 867 \quad 868 \quad 869 \quad 870 \quad 871 \quad 872 \quad 873 \quad 874 \quad 875 \quad 876 \quad 877 \quad 878 \quad 879 \quad 880 \quad 881 \quad 882 \quad 883 \quad 884 \quad 885 \quad 886 \quad 887 \quad 888 \quad 889 \quad 890 \quad 891 \quad 892 \quad 893 \quad 894 \quad 895 \quad 896 \quad 897 \quad 898 \quad 899 \quad 900 \quad 901 \quad 902 \quad 903 \quad 904 \quad 905 \quad 906 \quad 907 \quad 908 \quad 909 \quad 910 \quad 911 \quad 912 \quad 913 \quad 914 \quad 915 \quad 916 \quad 917 \quad 918 \quad 919 \quad 920 \quad 921 \quad 922 \quad 923 \quad 924 \quad 925 \quad 926 \quad 927 \quad 928 \quad 929 \quad 930 \quad 931 \quad 932 \quad 933 \quad 934 \quad 935 \quad 936 \quad 937 \quad 938 \quad 939 \quad 940 \quad 941 \quad 942 \quad 943 \quad 944 \quad 945 \quad 946 \quad 947 \quad 948 \quad 949 \quad 950 \quad 951 \quad 952 \quad 953 \quad 954 \quad 955 \quad 956 \quad 957 \quad 958 \quad 959 \quad 960 \quad 961 \quad 962 \quad 963 \quad 964 \quad 965 \quad 966 \quad 967 \quad 968 \quad 969 \quad 970 \quad 971 \quad 972 \quad 973 \quad 974 \quad 975 \quad 976 \quad 977 \quad 978 \quad 979 \quad 980 \quad 981 \quad 982 \quad 983 \quad 984 \quad 985 \quad 986 \quad 987 \quad 988 \quad 989 \quad 990 \quad 991 \quad 992 \quad 993 \quad 994 \quad 995 \quad 996 \quad 997 \quad 998 \quad 999 \quad 1000 \quad 1001 \quad 1002 \quad 1003 \quad 1004 \quad 1005 \quad 1006 \quad 1007 \quad 1008 \quad 1009 \quad 1010 \quad 1011 \quad 1012 \quad 1013 \quad 1014 \quad 1015 \quad 1016 \quad 1017 \quad 1018 \quad 1019 \quad 1020 \quad 1021 \quad 1022 \quad 1023 \quad 1024 \quad 1025 \quad 1026 \quad 1027 \quad 1028 \quad 1029 \quad 1030 \quad 1031 \quad 1032 \quad 1033 \quad 1034 \quad 1035 \quad 1036 \quad 1037 \quad 1038 \quad 1039 \quad 1040 \quad 1041 \quad 1042 \quad 1043 \quad 1044 \quad 1045 \quad 1046 \quad 1047 \quad 1048 \quad 1049 \quad 1050 \quad 1051 \quad 1052 \quad 1053 \quad 1054 \quad 1055 \quad 1056 \quad 1057 \quad 1058 \quad 1059 \quad 1060 \quad 1061 \quad 1062 \quad 1063 \quad 1064 \quad 1065 \quad 1066 \quad 1067 \quad 1068 \quad 1069 \quad 1070 \quad 1071 \quad 1072 \quad 1073 \quad 1074 \quad 1075 \quad 1076 \quad 1077 \quad 1078 \quad 1079 \quad 1080 \quad 1081 \quad 1082 \quad 1083 \quad 1084 \quad 1085 \quad 1086 \quad 1087 \quad 1088 \quad 1089 \quad 1090 \quad 1091 \quad 1092 \quad 1093 \quad 1094 \quad 1095 \quad 1096 \quad 1097 \quad 1098 \quad 1099 \quad 1100 \quad 1101 \quad 1102 \quad 1103 \quad 1104 \quad 1105 \quad 1106 \quad 1107 \quad 1108 \quad 1109 \quad 1110 \quad 1111 \quad 1112 \quad 1113 \quad 1114 \quad 1115 \quad 1116 \quad 1117 \quad 1118 \quad 1119 \quad 1120 \quad 1121 \quad 1122 \quad 1123 \quad 1124 \quad 1125 \quad 1126 \quad 1127 \quad 1128 \quad 1129 \quad 1130 \quad 1131 \quad 1132 \quad 1133 \quad 1134 \quad 1135 \quad 1136 \quad 1137 \quad 1138 \quad 1139 \quad 1140 \quad 1141 \quad 1142 \quad 1143 \quad 1144 \quad 1145 \quad 1146 \quad 1147 \quad 1148 \quad 1149 \quad 1150 \quad 1151 \quad 1152 \quad 1153 \quad 1154 \quad 1155 \quad 1156 \quad 1157 \quad 1158 \quad 1159 \quad 1160 \quad 1161 \quad 1162 \quad 1163 \quad 1164 \quad 1165 \quad 1166 \quad 1167 \quad 1168 \quad 1169 \quad 1170 \quad 1171 \quad 1172 \quad 1173 \quad 1174 \quad 1175 \quad 1176 \quad 1177 \quad 1178 \quad 1179 \quad 1180 \quad 1181 \quad 1182 \quad 1183 \quad 1184 \quad 1185 \quad 1186 \quad 1187 \quad 1188 \quad 1189 \quad 1190 \quad 1191 \quad 1192 \quad 1193 \quad 1194 \quad 1195 \quad 1196 \quad 1197 \quad 1198 \quad 1199 \quad 1200 \quad 1201 \quad 1202 \quad 1203 \quad 1204 \quad 1205 \quad 1206 \quad 1207 \quad 1208 \quad 1209 \quad 1210 \quad 1211 \quad 1212 \quad 1213 \quad 1214 \quad 1215 \quad 1216 \quad 1217 \quad 1218 \quad 1219 \quad 1220 \quad 1221 \quad 1222 \quad 1223 \quad 1224 \quad 1225 \quad 1226 \quad 1227 \quad 1228 \quad 1229 \quad 1230 \quad 1231 \quad 1232 \quad 1233 \quad 1234 \quad 1235 \quad 1236 \quad 1237 \quad 1238 \quad 1239 \quad 1240 \quad 1241 \quad 1242 \quad 1243 \quad 1244 \quad 1245 \quad 1246 \quad 1247 \quad 1248 \quad 1249 \quad 1250 \quad 1251 \quad 1252 \quad 1253 \quad 1254 \quad 1255 \quad 1256 \quad 1257 \quad 1258 \quad 1259 \quad 1260 \quad 1261 \quad 1262 \quad 1263 \quad 1264 \quad 1265 \quad 1266 \quad 1267 \quad 1268 \quad 1269 \quad 1270 \quad 1271 \quad 1272 \quad 1273 \quad 1274 \quad 1275 \quad 1276 \quad 1277 \quad 1278 \quad 1279 \quad 1280 \quad 1281 \quad 1282 \quad 1283 \quad 1284 \quad 1285 \quad 1286 \quad 1287 \quad 1288 \quad 1289 \quad 1290 \quad 1291 \quad 1292 \quad 1293 \quad 1294 \quad 1295 \quad 1296 \quad 1297 \quad 1298 \quad 1299 \quad 1300 \quad 1301 \quad 1302 \quad 1303 \quad 1304 \quad 1305 \quad 1306 \quad 1307 \quad 1308 \quad 1309 \quad 1310 \quad 1311 \quad 1312 \quad 1313 \quad 1314 \quad 1315 \quad 1316 \quad 1317 \quad 1318 \quad 1319 \quad 1320 \quad 1321 \quad 1322 \quad 1323 \quad 1324 \quad 1325 \quad 1326 \quad 1327 \quad 1328 \quad 1329 \quad 1330 \quad 1331 \quad 1332 \quad 1333 \quad 1334 \quad 1335 \quad 1336 \quad 1337 \quad 1338 \quad 1339 \quad 1340 \quad 1341 \quad 1342 \quad 1343 \quad 1344 \quad 1345 \quad 1346 \quad 1347 \quad 1348 \quad 1349 \quad 1350 \quad 1351 \quad 1352 \quad 1353 \quad 1354 \quad 1355 \quad 1356 \quad 1357 \quad 1358 \quad 1359 \quad 1360 \quad 1361 \quad 1362 \quad 1363 \quad 1364 \quad 1365 \quad 1366 \quad 1367 \quad 1368 \quad 1369 \quad 1370 \quad 1371 \quad 1372 \quad 1373 \quad 1374 \quad 1375 \quad 1376 \quad 1377 \quad 1378 \quad 1379 \quad 1380 \quad 1381 \quad 1382 \quad 1383 \quad 1384 \quad 1385 \quad 1386 \quad 1387 \quad 1388 \quad 1389 \quad 1390 \quad 1391 \quad 1392 \quad 1393 \quad 1394 \quad 1395 \quad 1396 \quad 1397 \quad 1398 \quad 1399 \quad 1400 \quad 1401 \quad 1402 \quad 1403 \quad 1404 \quad 1405 \quad 1406 \quad 1407 \quad 1408 \quad 1409 \quad 1410 \quad 1411 \quad 1412 \quad 1413 \quad 1414 \quad 1415 \quad 1416 \quad 1417 \quad 1418 \quad 1419 \quad 1420 \quad 1421 \quad 1422 \quad 1423 \quad 1424 \quad 1425 \quad 1426 \quad 1427 \quad 1428 \quad 1429 \quad 1430 \quad 1431 \quad 1432 \quad 1433 \quad 1434 \quad 1435 \quad 1436 \quad 1437 \quad 1438 \quad 1439 \quad 1440 \quad 1441 \quad 1442 \quad 1443 \quad 1444 \quad 1445 \quad 1446 \quad 1447 \quad 1448 \quad 1449 \quad 1450 \quad 1451 \quad 1452 \quad 1453 \quad 1454 \quad 1455 \quad 1456 \quad 1457 \quad 1458 \quad 1459 \quad 1460 \quad 1461 \quad 1462 \quad 1463 \quad 1464 \quad 1465 \quad 1466 \quad 1467 \quad 1468 \quad 1469 \quad 1470 \quad 1471 \quad 1472 \quad 1473 \quad 1474 \quad 1475 \quad 1476 \quad 1477 \quad 1478 \quad 1479 \quad 1480 \quad 1481 \quad 1482 \quad 1483 \quad 1484 \quad 1485 \quad 1486 \quad 1487 \quad 1488 \quad 1489 \quad 1490 \quad 1491 \quad 1492 \quad 1493 \quad 1494 \quad 1495 \quad 1496 \quad 1497 \quad 1498 \quad 1499 \quad 1500 \quad 1501 \quad 1502 \quad 1503 \quad 1504 \quad 1505 \quad 1506 \quad 1507 \quad 1508 \quad 1509 \quad 1510 \quad 1511 \quad 1512 \quad 1513 \quad 1514 \quad 1515 \quad 1516 \quad 1517 \quad 1518 \quad 1519 \quad 1520 \quad 1521 \quad 1522 \quad 1523 \quad 1524 \quad 1525 \quad 1526 \quad 1527 \quad 1528 \quad 1529 \quad 1530 \quad 1531 \quad 1532 \quad 1533 \quad 1534 \quad 1535 \quad 1536 \quad 1537 \quad 1538 \quad 1539 \quad 1540 \quad 1541 \quad 1542 \quad 1543 \quad 1544 \quad 1545 \quad 1546 \quad 1547 \quad 1548 \quad 1549 \quad 1550 \quad 1551 \quad 1552 \quad 1553 \quad 1554 \quad 1555 \quad 1556 \quad 1557 \quad 1558 \quad 1559 \quad 1560 \quad 1561 \quad 1562 \quad 1563 \quad 1564 \quad 1565 \quad 1566 \quad 1567 \quad 1568 \quad 1569 \quad 1570 \quad 1571 \quad 1572 \quad 1573 \quad 1574 \quad 1575 \quad 1576 \quad 1577 \quad 1578 \quad 1579 \quad 1580 \quad 1581 \quad 1582 \quad 1583 \quad 1584 \quad 1585 \quad 1586 \quad 1587 \quad 1588 \quad 1589 \quad 1590 \quad 1591 \quad 1592 \quad 1593 \quad 1594 \quad 1595 \quad 1596 \quad 1597 \quad 1598 \quad 1599 \quad 1600 \quad 1601 \quad 1602 \quad 1603 \quad 1604 \quad 1605 \quad 1606 \quad 1607 \quad 1608 \quad 1609 \quad 1610 \quad 1611 \quad 1612 \quad 1613 \quad 1614 \quad 1615 \quad 1616 \quad 1617 \quad 1618 \quad 1619 \quad 1620 \quad 1621 \quad 1622 \quad 1623 \quad 1624 \quad 1625 \quad 1626 \quad 1627 \quad 1628 \quad 1629 \quad 1630 \quad 1631 \quad 1632 \quad 1633 \quad 1634 \quad 1635 \quad 1636 \quad 1637 \quad 1638 \quad 1639 \quad 1640 \quad 1641 \quad 1642 \quad 1643 \quad 1644 \quad 1645 \quad 1646 \quad 1647 \quad 1648 \quad 1649 \quad 1650 \quad 1651 \quad 1652 \quad 1653 \quad 1654 \quad 1655 \quad 1656 \quad 1657 \quad 1658 \quad 1659 \quad 1660 \quad 1661 \quad 1662 \quad 1663 \quad 1664 \quad 1665 \quad 1666 \quad 1667 \quad 1668 \quad 1669 \quad 1670 \quad 1671 \quad 1672 \quad 1673 \quad 1674 \quad 1675 \quad 1676 \quad 1677 \quad 1678 \quad 1679 \quad 1680 \quad 1681 \quad 1682 \quad 1683 \quad 1684 \quad 1685 \quad 1686 \quad 1687 \quad 1688 \quad 1689 \quad 1690 \quad 1691 \quad 1692 \quad 1693 \quad 1694 \quad 1695 \quad 1696 \quad 1697 \quad 1698 \quad 1699 \quad 1700 \quad 1701 \quad 1702 \quad 1703 \quad 1704 \quad 1705 \quad 1706 \quad 1707 \quad 1708 \quad 1709 \quad 1710 \quad 1711 \quad 1712 \quad 1713 \quad 1714 \quad 1715 \quad 1716 \quad 1717 \quad 1718 \quad 1719 \quad 1720 \quad 1721 \quad 1722 \quad 1723 \quad 1724 \quad 1725 \quad 1726 \quad 1727 \quad 1728 \quad 1729 \quad 1730 \quad 1731 \quad 1732 \quad 1733 \quad 1734 \quad 1735 \quad 1736 \quad 1737 \quad 1738 \quad 1739 \quad 1740 \quad 1741 \quad 1742 \quad 1743 \quad 1744 \quad 1745 \quad 1746 \quad 1747 \quad 1748 \quad 1749 \quad$$

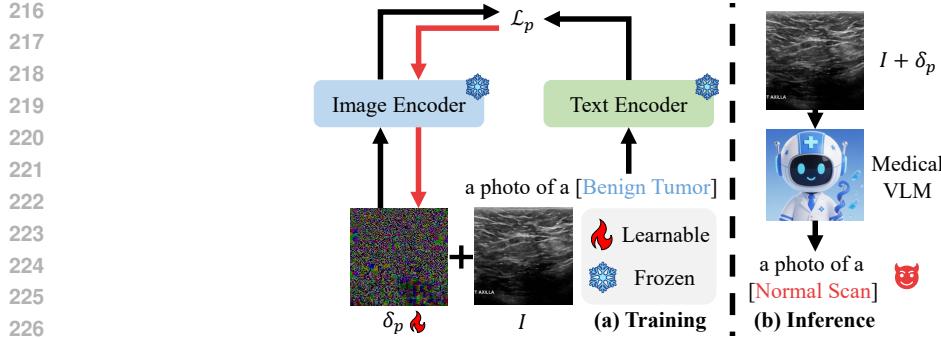


Figure 3: (a) Pipeline of training an adversarial perturbation map  $\delta_p$  for a specific image  $I$  with its corresponding label. (b) Pipeline of inference with an adversarial input image  $I + \delta_p$ . The black and red arrows indicate the data flow and gradient flow, respectively.

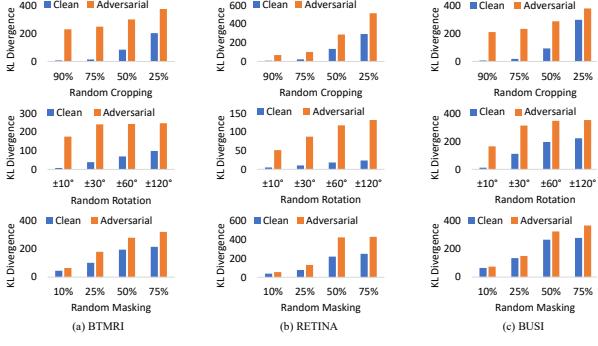


Figure 4: The KL divergence between BiomedCLIP's predictions before and after applying transformations across three datasets with various modalities.

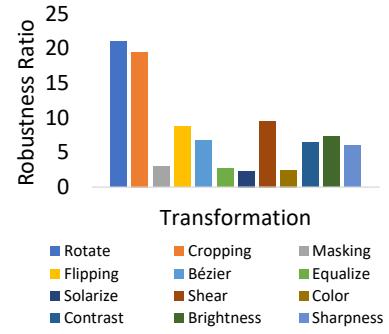


Figure 5: Averaged robustness ratio of 12 transformation strategies calculated on 11 datasets.

**Adversarial attack.** We focus on three typical adversarial attack methods under a white-box setting: PGD Madry et al. (2018), C&W Carlini & Wagner (2017), and AutoAttack Croce & Hein (2020). This setting assumes that adversaries have complete access to the architecture and parameters of the victim model, enabling direct gradient-based attacks. The pipeline is shown in Figure 3. The adversary learns an adversarial perturbation to increase the divergence between text and image embeddings by maximizing an adversarial perturbation loss  $\mathcal{L}_p$  (e.g., cross-entropy) as follows:

$$\delta_p = \arg \max_{\|\delta\|_\infty \leq \epsilon_p} \mathcal{L}_p(P(I + \delta), c_y) \quad (3)$$

where  $\epsilon_p$  and  $c_y$  denote the perturbation budget and the class label, respectively. The adversary then employs the victim VLM on the adversarial image  $\hat{I} = I + \delta_p$  to obtain an incorrect prediction.

### 3.2 SEMANTIC FRAGILITY OF ADVERSARIAL PERTURBATIONS

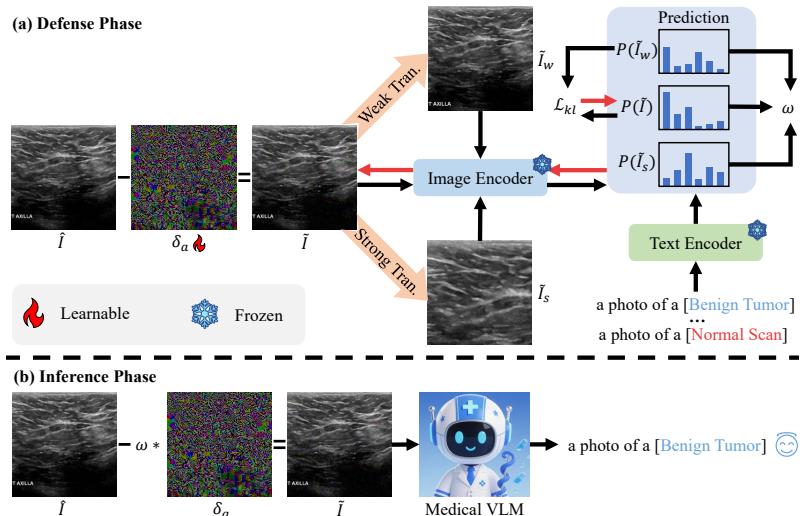
In this study, we found that VLM is more robust than adversarial perturbations, and the perturbed image embeddings are semantically fragile and highly sensitive to minor alterations. To demonstrate this, we first performed toy experiments on three representative datasets from three modalities (MRI, dermoscopy, and X-ray) using three types of transformations: random cropping, random rotation, and random masking. Specifically, we measured the robustness by calculating the symmetrical KL divergence between predictions before and after applying transformations as follows:

$$r(I, I_t) = \mathcal{L}_{kl}(P(I), P(I_t)) + \mathcal{L}_{kl}(P(I_t), P(I)), \quad (4)$$

where  $\mathcal{L}_{kl}$  indicates the KL divergence loss, and  $I_t$  denotes the transformed image. PGD and BiomedCLIP Zhang et al. (2023) are utilized as the adversarial attack method and the victim model, respectively. As illustrated in Figure 4, weak transformations (e.g.,  $\pm 10^\circ$  random rotation and 90% random cropping) produce low KL divergence for clean images but high divergence for adversarial

270 ones, with this gap decreasing under stronger transformations. This phenomenon can be attributed  
 271 to that adversarial perturbations are highly over-fitted to the specific image. However, this effect  
 272 is dependent on the transformation type. In contrast, random masking has a minimal effect at low  
 273 intensities, as it only alters a small portion of the adversarial perturbations.

274 To discuss the efficacy of various transformation strategies, we further defined a ratio  $R =$   
 275  $r(\hat{I}, \hat{I}_t)/r(I, I_t)$  to quantify the robustness discrepancy induced by a transformation  $t$  between clean  
 276 and adversarial images. A larger value of  $R$  indicates a stronger discriminative capacity of  $t$  in  
 277 distinguishing adversarial images from clean ones. We calculated  $R$  on 13 common transformation  
 278 strategies and displayed the results in Figure 5. It shows that an effective transformation strategy  
 279 yielding a high  $R$  should satisfy two criteria: (1) it should modify the values and/or spatial positions  
 280 of a majority of pixels, thereby amplifying its impact on adversarial images; and (2) it should  
 281 apply minimal distortion to preserve high robustness on clean images. Therefore, we selected ran-  
 282 dom cropping and random rotation, which exhibit high  $R$  values, to construct our defense strategy.  
 283 Comprehensive experiment results are provided in the Appendix.



301 Figure 6: (a) Pipeline of our TAME training the adversarial restoration map  $\delta_a$  for an adversarial  
 302 image  $\hat{I}$ . (b) Pipeline of inference with the weighted restored image  $\tilde{I}$ . The black and red arrows  
 303 indicate the data flow and gradient flow, respectively. ‘Tran.’: Abbreviation of ‘Transformation’.

### 306 3.3 THE ATTACK MEANS NOTHING (TAME)

#### 307 3.3.1 ADVERSARIAL RESTORATION

309 In this paper, we specifically target the preservation of VLM’s zero-shot inference robustness, where  
 310 the defender has neither access to task-specific training data nor annotations for test samples. Based  
 311 on the above observations (see Section 3.2), we proposed TAME, a simple yet effective method, as  
 312 illustrated in Figure 6. For each input adversarial image  $\hat{I}$ , we introduced a learnable adversarial  
 313 restoration map  $\delta_a$ , yielding the restored image  $\tilde{I} = \hat{I} - \delta_a$ , which is intended to approximate the  
 314 original clean image  $I$ . Due to the inability to directly obtain  $I$ , we adopted a compromise approach  
 315 that restores the model’s strong robustness on weak transformations by posing the consistency con-  
 316 straint to train  $\delta_a$  as follows:

$$318 \min_{\|\delta_a\|_\infty \leq \epsilon_a} \mathcal{L}_{kl}(P(\tilde{I}), P(\tilde{I}_w)), \quad (5)$$

320 where  $\epsilon_a$  and  $\tilde{I}_w$  indicate the defense budget and the weakly transformed  $\tilde{I}$ , respectively. Note that  
 321 the  $\mathcal{L}_{kl}(P(\tilde{I}_w), P(\tilde{I}))$  term is removed since  $P(\tilde{I})$  approximates  $P(\hat{I})$ , which may provide limited  
 322 supervision. This process to update  $\delta_a$  can be approximated by PGD Madry et al. (2018):

$$323 \delta_a^1 = \prod (\delta_a^0 - \alpha \text{sgn}(\nabla_{\delta_a} \mathcal{L}_{kl}(P(\tilde{I}), P(\tilde{I}_w)))), \quad (6)$$

324 where  $\alpha$  denotes the step-size, and the update step is fixed at 1 in this study. The initial perturbation  
 325 map  $\delta_a^0$  is randomly sampled from a uniform distribution  $U(-\epsilon_a, \epsilon_a)$ .  
 326

### 327 3.3.2 WEIGHTING MECHANISM

329 In this section, we attempt to address the risk that directly applying  $\delta_a$  to produce the restored  
 330 image may degrade the performance on clean images, thereby hindering model deployments. An  
 331 intuitive strategy is to leverage the divergent responses of clean and adversarial images to minor  
 332 transformations to distinguish them. However, it is hard to determine a universal threshold due to  
 333 the discrepancies among datasets. Recall the observation shown in Figure 4 that strong transformations  
 334 with large magnitude induce high KL divergence values on both clean and adversarial images.  
 335 Derived from this, such a high KL divergence value can be treated as an anchor to take the effect of  
 336 normalization across datasets. We then devised a dynamic weight coefficient  $\omega$  formulated by:  
 337

$$\omega = \frac{\mathcal{L}_{kl}(P(\tilde{I}), P(\tilde{I}_w))}{\mathcal{L}_{kl}(P(\tilde{I}), P(\tilde{I}_s))} \quad (7)$$

339 where  $\tilde{I}_s$  indicates the strongly transformed  $\tilde{I}$ . Due to the discrepancy between clean and adversarial  
 340 images, this ensures allocating larger weights to adversarial ones while avoiding excessive modification  
 341 of clean ones. To further amplify the effect of  $\omega$  on adversarial images and mitigate potential  
 342 instability from excessively large weights, we also truncated  $\omega$  using an empirically determined  
 343 threshold value of 0.5, where values exceeding this threshold were clipped to 1. Then the model can  
 344 perform inference on the weighted restored image as illustrated in Figure 6 (b). We summarize the  
 345 algorithm of our TAME in Algorithm 1. We summarize the algorithm of our TAME in Algorithm 1.  
 346

---

#### 347 **Algorithm 1:** TAME Algorithm.

348 **Input:** Current test image  $I$ , pre-trained VLM (including  $f_{\theta_v}$  and  $f_{\theta_t}$ ), defense budget  $\epsilon_a$ , and step-size  $\alpha$ .  
 349 1:  $\delta_a^0 \sim U(-\epsilon_a, \epsilon_a)$ .  
 350 2:  $\delta_a^1 = \prod(\delta_a^0 - \alpha \text{sgn}(\nabla_{\delta_a} \mathcal{L}_{kl}(P(I), P(I_w))))$ .  
 351 3:  $\delta_a = \text{clamp}(\delta_a^1, -\epsilon_a, \epsilon_a)$ .  
 352 4:  $\omega = \frac{\mathcal{L}_{kl}(P(I), P(I_w))}{\mathcal{L}_{kl}(P(I), P(I_s))}$ .  
 353 5:  $\omega = \omega \cdot \mathbb{1}_{\omega \leq 0.5} + 1 \cdot \mathbb{1}_{\omega > 0.5}$   
 354 **Output:**  $P(I - \omega \cdot \delta_a)$

---

## 356 4 EXPERIMENTS AND RESULTS

### 357 4.1 DATASETS AND METRIC

361 We conducted all experiments on the test sets of 11 diverse medical datasets spanning 9 imaging  
 362 modalities: Computerized Tomography (CTKidney Islam et al. (2022)), Dermatoscopy (DermaM-  
 363 NIST Codella et al. (2019); Tschandl et al. (2018)), Endoscopy (Kvasir Pogorelov et al. (2017)),  
 364 Color Fundus Photography (RETINA Köhler et al. (2013); Porwal et al. (2018)), Histopathology  
 365 (LC25000 Borkowski et al. (2019) and CHMNIST Kather et al. (2016)), Magnetic Resonance Imaging  
 366 (BTMRI Nickparvar (2021)), Optical Coherence Tomography (OCTMNIST Kermany et al.  
 367 (2018)), Ultrasound (BUSI Al-Dhabayani et al. (2020)), and X-Ray (COVID-QU-Ex Tahir et al.  
 368 (2021) and KneeXray Chen (2018)). The details are listed in Table 3. We utilized the classification  
 369 accuracy as the metric to evaluate our TAME and other competing methods.

### 370 4.2 IMPLEMENTATION DETAILS

372 We employed the pre-trained BiomedCLIP Zhang et al. (2023) with a ViT-B/16 backbone as the  
 373 victim model and reported the averaged results calculated across three trials (*i.e.*, setting the random  
 374 seed to 0, 1, and 2). The attack budget  $\epsilon_p$  was set to 1/255 to guarantee imperceptible perturbations,  
 375 and the update step for each attack method was set to 10. All experiments are conducted under the  
 376 white-box setting, where the adversary has full access to the victim model. In our TAME, the weak  
 377 transformation strategy is empirically defined as the combination of  $\pm 10^\circ$  random rotation and 90%  
 378 random cropping. The strong transformation strategy applies a more intensive combination of  $\pm 30^\circ$

378 random rotation and 50% random cropping. Since the defense is operated by the user at test time,  
 379 there is no need for the adversarial restoration map to be undetectable, allowing a large defense  
 380 budget. Therefore, we set both the defense budget  $\epsilon_a$  and the step-size  $\alpha$  to 8/255.  
 381

### 382 4.3 EXPERIMENTAL RESULTS

384 Table 1: Zero-shot adversarial robustness (%) of our TAME, the BiomedCLIP baseline, and other  
 385 competing TAD methods on 11 medical datasets. We report the mean and standard deviation cal-  
 386 culated across three trials. For each dataset, the highest performance under the Clean, PGD, C&W,  
 387 and AutoAttack (AA) settings is highlighted in red, blue, green, and purple, respectively.  
 388

Dataset	Attack	BiomedCLIP	Anti-Adv	HedgeDefense	TTC	R-TPT	TAME
<b>BTMRI</b>	Clean	56.79	41.62 $\pm$ 0.10	<b>58.77</b> $\pm$ 0.08	42.07 $\pm$ 0.43	54.30 $\pm$ 0.34	54.13 $\pm$ 1.09
	PGD	0.68 $\pm$ 0.07	9.88 $\pm$ 0.29	6.37 $\pm$ 0.10	<b>53.41</b> $\pm$ 0.08	48.67 $\pm$ 0.47	<b>61.21</b> $\pm$ 0.33
	C&W	0.68 $\pm$ 0.03	7.96 $\pm$ 0.35	7.55 $\pm$ 0.07	<b>53.37</b> $\pm$ 0.67	48.84 $\pm$ 0.22	<b>61.50</b> $\pm$ 0.43
	AA	0.06 $\pm$ 0.00	8.15 $\pm$ 0.22	7.65 $\pm$ 0.15	<b>56.92</b> $\pm$ 0.87	50.46 $\pm$ 0.26	<b>61.25</b> $\pm$ 0.19
<b>BUSI</b>	Clean	59.75	27.54 $\pm$ 0.00	49.72 $\pm$ 0.72	40.96 $\pm$ 1.21	45.48 $\pm$ 1.11	<b>62.71</b> $\pm$ 2.07
	PGD	0.00 $\pm$ 0.00	8.33 $\pm$ 0.20	2.83 $\pm$ 0.53	51.84 $\pm$ 1.60	34.88 $\pm$ 1.44	<b>68.08</b> $\pm$ 2.45
	C&W	0.00 $\pm$ 0.00	14.41 $\pm$ 0.35	5.37 $\pm$ 1.00	49.15 $\pm$ 1.83	34.32 $\pm$ 0.92	<b>70.90</b> $\pm$ 0.53
	AA	0.00 $\pm$ 0.00	8.47 $\pm$ 0.00	4.38 $\pm$ 0.53	55.09 $\pm$ 1.51	37.71 $\pm$ 1.51	<b>65.54</b> $\pm$ 0.20
<b>COVID-QU-Ex</b>	Clean	43.82	43.80 $\pm$ 0.01	<b>48.67</b> $\pm$ 0.21	31.50 $\pm$ 0.37	37.51 $\pm$ 0.15	36.38 $\pm$ 0.26
	PGD	0.00 $\pm$ 0.00	0.15 $\pm$ 0.05	0.62 $\pm$ 0.01	48.93 $\pm$ 0.28	25.99 $\pm$ 0.09	<b>54.41</b> $\pm$ 0.22
	C&W	0.00 $\pm$ 0.00	0.17 $\pm$ 0.05	0.66 $\pm$ 0.08	49.30 $\pm$ 0.31	26.23 $\pm$ 0.13	<b>53.70</b> $\pm$ 0.42
	AA	0.00 $\pm$ 0.00	0.20 $\pm$ 0.03	10.03 $\pm$ 0.06	40.51 $\pm$ 0.40	31.76 $\pm$ 0.29	<b>54.00</b> $\pm$ 0.48
<b>CTKIDNEY</b>	Clean	42.43	40.25 $\pm$ 0.01	42.56 $\pm$ 0.04	29.73 $\pm$ 0.14	<b>48.33</b> $\pm$ 0.03	40.36 $\pm$ 0.38
	PGD	0.87 $\pm$ 0.03	1.31 $\pm$ 0.04	2.36 $\pm$ 0.12	26.32 $\pm$ 0.29	40.98 $\pm$ 0.13	<b>53.01</b> $\pm$ 0.60
	C&W	0.88 $\pm$ 0.02	2.75 $\pm$ 0.09	2.89 $\pm$ 0.04	26.35 $\pm$ 0.39	41.10 $\pm$ 0.12	<b>52.02</b> $\pm$ 0.60
	AA	0.05 $\pm$ 0.00	0.68 $\pm$ 0.04	4.91 $\pm$ 0.12	32.58 $\pm$ 0.25	45.23 $\pm$ 0.31	<b>50.42</b> $\pm$ 0.61
<b>DermaMNIST</b>	Clean	38.80	38.65 $\pm$ 0.00	37.44 $\pm$ 0.14	15.69 $\pm$ 0.33	<b>43.09</b> $\pm$ 0.39	27.95 $\pm$ 0.63
	PGD	0.00 $\pm$ 0.00	0.07 $\pm$ 0.06	0.88 $\pm$ 0.02	<b>40.57</b> $\pm$ 0.24	21.00 $\pm$ 0.27	40.28 $\pm$ 0.59
	C&W	0.00 $\pm$ 0.00	0.13 $\pm$ 0.02	1.02 $\pm$ 0.18	39.98 $\pm$ 1.40	20.03 $\pm$ 0.12	<b>41.30</b> $\pm$ 0.60
	AA	0.00 $\pm$ 0.00	0.30 $\pm$ 0.04	6.23 $\pm$ 0.04	33.47 $\pm$ 0.48	33.44 $\pm$ 0.41	<b>41.99</b> $\pm$ 0.25
<b>Kvasir</b>	Clean	54.58	42.42 $\pm$ 0.00	<b>56.09</b> $\pm$ 0.12	26.64 $\pm$ 1.22	<b>56.28</b> $\pm$ 0.61	48.30 $\pm$ 1.00
	PGD	0.00 $\pm$ 0.00	2.19 $\pm$ 0.40	0.42 $\pm$ 0.07	46.16 $\pm$ 0.59	41.89 $\pm$ 0.45	<b>59.61</b> $\pm$ 0.08
	C&W	0.00 $\pm$ 0.00	2.75 $\pm$ 0.20	0.31 $\pm$ 0.14	43.11 $\pm$ 0.52	41.42 $\pm$ 0.65	<b>58.11</b> $\pm$ 0.22
	AA	0.00 $\pm$ 0.00	3.28 $\pm$ 0.04	4.06 $\pm$ 0.28	48.19 $\pm$ 0.31	47.86 $\pm$ 0.34	<b>63.72</b> $\pm$ 0.67
<b>CHMNIST</b>	Clean	<b>30.65</b>	20.39 $\pm$ 0.06	25.53 $\pm$ 0.14	25.42 $\pm$ 0.33	29.94 $\pm$ 0.55	21.77 $\pm$ 0.71
	PGD	0.00 $\pm$ 0.00	5.57 $\pm$ 0.30	0.18 $\pm$ 0.07	20.15 $\pm$ 0.09	16.51 $\pm$ 0.37	<b>25.97</b> $\pm$ 0.17
	C&W	0.02 $\pm$ 0.03	5.30 $\pm$ 0.36	0.40 $\pm$ 0.06	19.88 $\pm$ 0.24	16.58 $\pm$ 0.64	<b>24.98</b> $\pm$ 0.73
	AA	0.00 $\pm$ 0.00	3.32 $\pm$ 0.11	3.37 $\pm$ 0.03	24.96 $\pm$ 0.22	22.52 $\pm$ 0.38	<b>30.83</b> $\pm$ 0.32
<b>LC25000</b>	Clean	50.01	48.10 $\pm$ 0.04	<b>54.14</b> $\pm$ 0.07	32.77 $\pm$ 0.10	<b>50.12</b> $\pm$ 0.14	44.04 $\pm$ 0.23
	PGD	0.01 $\pm$ 0.00	1.21 $\pm$ 0.03	0.21 $\pm$ 0.02	32.21 $\pm$ 0.08	38.87 $\pm$ 0.05	<b>55.75</b> $\pm$ 0.56
	C&W	0.02 $\pm$ 0.01	1.48 $\pm$ 0.04	0.36 $\pm$ 0.03	30.89 $\pm$ 0.11	38.43 $\pm$ 0.14	<b>52.47</b> $\pm$ 0.41
	AA	0.01 $\pm$ 0.00	5.74 $\pm$ 0.01	8.76 $\pm$ 0.03	41.19 $\pm$ 0.33	43.14 $\pm$ 0.05	<b>54.62</b> $\pm$ 0.20
<b>RETINA</b>	Clean	26.26	26.37 $\pm$ 0.04	26.10 $\pm$ 0.51	29.15 $\pm$ 0.62	<b>32.89</b> $\pm$ 0.81	26.18 $\pm$ 0.30
	PGD	0.00 $\pm$ 0.00	9.75 $\pm$ 0.55	2.60 $\pm$ 0.45	<b>35.54</b> $\pm$ 0.82	20.27 $\pm$ 0.79	26.13 $\pm$ 0.53
	C&W	0.00 $\pm$ 0.00	8.89 $\pm$ 0.23	1.61 $\pm$ 0.10	<b>33.62</b> $\pm$ 0.70	21.69 $\pm$ 0.84	26.21 $\pm$ 0.48
	AA	0.00 $\pm$ 0.00	9.59 $\pm$ 0.41	8.63 $\pm$ 0.50	29.94 $\pm$ 1.11	26.68 $\pm$ 0.26	26.68 $\pm$ 0.21
<b>KneeXray</b>	Clean	29.47	8.86 $\pm$ 0.03	23.85 $\pm$ 0.05	24.88 $\pm$ 0.47	<b>40.84</b> $\pm$ 0.31	38.38 $\pm$ 0.44
	PGD	0.00 $\pm$ 0.00	0.24 $\pm$ 0.13	0.22 $\pm$ 0.11	<b>47.75</b> $\pm$ 0.12	27.70 $\pm$ 0.64	46.15 $\pm$ 0.33
	C&W	0.00 $\pm$ 0.00	1.07 $\pm$ 0.11	1.61 $\pm$ 0.27	<b>46.94</b> $\pm$ 0.39	28.92 $\pm$ 0.84	41.08 $\pm$ 0.06
	AA	0.00 $\pm$ 0.00	3.72 $\pm$ 0.06	18.74 $\pm$ 0.44	17.47 $\pm$ 1.07	35.65 $\pm$ 0.35	<b>39.61</b> $\pm$ 0.26
<b>OCTMNIST</b>	Clean	29.90	28.80 $\pm$ 0.00	26.23 $\pm$ 0.05	29.37 $\pm$ 0.79	25.40 $\pm$ 0.08	<b>34.10</b> $\pm$ 0.36
	PGD	6.27 $\pm$ 0.68	8.13 $\pm$ 0.33	24.83 $\pm$ 0.05	33.73 $\pm$ 0.38	25.17 $\pm$ 0.05	<b>39.40</b> $\pm$ 0.29
	C&W	6.37 $\pm$ 0.17	7.33 $\pm$ 0.33	25.17 $\pm$ 0.17	32.57 $\pm$ 0.41	25.17 $\pm$ 0.05	<b>39.63</b> $\pm$ 0.38
	AA	0.00 $\pm$ 0.00	0.33 $\pm$ 0.05	18.30 $\pm$ 0.14	20.40 $\pm$ 0.59	25.20 $\pm$ 0.14	<b>37.10</b> $\pm$ 0.57
<b>Average</b>	Clean	42.04	33.35 $\pm$ 0.01	40.83 $\pm$ 0.13	29.84 $\pm$ 0.14	<b>42.20</b> $\pm$ 0.10	39.49 $\pm$ 0.21
	PGD	0.71 $\pm$ 0.06	4.26 $\pm$ 0.03	3.77 $\pm$ 0.04	39.69 $\pm$ 0.29	31.08 $\pm$ 0.08	<b>48.18</b> $\pm$ 0.18
	C&W	0.72 $\pm$ 0.01	4.75 $\pm$ 0.07	4.27 $\pm$ 0.07	38.65 $\pm$ 0.24	31.16 $\pm$ 0.21	<b>47.45</b> $\pm$ 0.14
	AA	0.01 $\pm$ 0.00	3.98 $\pm$ 0.06	8.64 $\pm$ 0.05	36.68 $\pm$ 0.13	36.63 $\pm$ 0.11	<b>47.80</b> $\pm$ 0.11

#### 424 4.3.1 COMPARISON WITH OTHER TAD METHODS

426 We compared our TAME with the BiomedCLIP baseline, two TAD methods designed for the tradi-  
 427 tional models (Anti-Adv Alfarra et al. (2022) and HedgeDefense Wu et al. (2021)), and two TAD  
 428 methods tailored for VLMs (TTC Xing et al. (2025) and R-TPT Sheng et al. (2025)). Specifically,  
 429 we re-implemented all the competing methods using the same baseline and reproduced the results by  
 430 utilizing their open-source codes. As detailed in Table 1, the results reveal that (1) BiomedCLIP is  
 431 highly susceptible to adversarial attacks, which devastate its inference capabilities; (2) Anti-Adv and  
 HedgeDefense provide only marginal improvements, underscoring their limited defense ability for

432  
 433 Table 2: Zero-shot adversarial robustness (%) of our TAME, TTC, and R-TPT integrated with three  
 434 distinct AFT methods: FARE, PMG, and TeCoA. We report the mean and standard deviation cal-  
 435 culated across three trials. For each AFT method, the highest performance under the Clean, PGD,  
 436 C&W, and AutoAttack (AA) settings is highlighted in **red**, **blue**, **green**, and **purple**, respectively.

Method	Clean	PGD	C&W	AA
CLIP (ViT-B/32)	24.33	$0.07 \pm 0.02$	$0.13 \pm 0.01$	$0.13 \pm 0.00$
FARE	22.51	$6.09 \pm 0.01$	$6.02 \pm 0.01$	$5.71 \pm 0.00$
FARE + TTC	$23.09 \pm 0.22$	$16.97 \pm 0.08$	$16.40 \pm 0.03$	$22.82 \pm 0.03$
FARE + R-TPT	$22.70 \pm 0.02$	$16.77 \pm 0.03$	$17.05 \pm 0.05$	$19.46 \pm 0.07$
FARE + TAME (Ours)	$23.43 \pm 0.10$	$32.32 \pm 0.26$	$30.85 \pm 0.36$	$28.98 \pm 0.11$
PMG	<b>22.95</b>	$12.27 \pm 0.02$	$11.71 \pm 0.01$	$11.65 \pm 0.02$
PMG + TTC	$22.48 \pm 0.05$	$16.29 \pm 0.07$	$15.96 \pm 0.07$	$20.12 \pm 0.13$
PMG + R-TPT	$20.53 \pm 0.05$	$17.70 \pm 0.01$	$17.51 \pm 0.04$	$19.07 \pm 0.04$
PMG + TAME (Ours)	$21.02 \pm 0.17$	<b>21.40</b> $\pm 0.14$	<b>21.09</b> $\pm 0.22$	<b>20.80</b> $\pm 0.04$
TeCoA	22.56	$11.96 \pm 0.01$	$11.42 \pm 0.01$	$11.49 \pm 0.01$
TeCoA + TTC	$22.24 \pm 0.07$	$16.14 \pm 0.11$	$16.02 \pm 0.13$	$20.00 \pm 0.23$
TeCoA + R-TPT	<b>22.84</b> $\pm 0.05$	$19.16 \pm 0.07$	$19.07 \pm 0.08$	$21.34 \pm 0.10$
TeCoA + TAME (Ours)	$22.68 \pm 0.07$	<b>23.52</b> $\pm 0.05$	<b>23.43</b> $\pm 0.07$	<b>23.40</b> $\pm 0.10$

448  
 449 VLMs; and (3) our TAME consistently demonstrates strong adversarial robustness, delivering su-  
 450 perior performance in most scenarios and achieving the best overall accuracy across all attack types,  
 451 while maintaining accuracy on clean images with minor and acceptable degradation. If higher clean  
 452 accuracy is required, the defense budget can be reduced, as explored in Appendix C.1. Addition-  
 453 ally, an intriguing observation is that TAME’s overall accuracy under adversarial attacks (48.18%  
 454 for PGD, 47.45% for C&W, and 47.8% for AutoAttack) surpasses that of BiomedCLIP on clean  
 455 images (42.04%). This phenomenon indicates a potential risk of label leakage during the attack  
 456 process. We will discuss it in the Appendix B.

#### 457 4.3.2 EXTENSIBILITY ANALYSIS

458 To evaluate the extensibility of TAME, TTC, and R-TPT, we integrated each one with various adver-  
 459 sarially fine-tuned models. In this experimental setup, we utilized a pre-trained CLIP model with a  
 460 ViT-B/32 backbone as the base victim model. Due to the challenge of obtaining a fine-tuning dataset  
 461 that covers all downstream modalities, we implemented three AFT methods (*i.e.*, FARE Schlar mann  
 462 et al. (2024), PMG Wang et al. (2024a), and TeCoA Mao et al. (2023)) by fine-tuning the CLIP  
 463 vision encoder on adversarial images from the TinyImageNet dataset. The presence of a significant  
 464 discrepancy between the adversarial training data and testing data can serve to assess the adaptabil-  
 465 ity of AFT methods in generalizing to unseen testing adversarial images. For conciseness, we only  
 466 display the average accuracy across 11 datasets in Table 2, and the complete results can be found in  
 467 Table 8. The results indicate that (1) AFT methods provide only a partial defense against attacks,  
 468 which can be attributed to their limited adaptation capability when generalized to diverse test data;  
 469 (2) all TAD methods consistently boost the adversarial robustness of adversarially fine-tuned mod-  
 470 els, demonstrating the effectiveness of test-time adversarial defense; and (3) our TAME achieves  
 471 significantly superior robustness enhancements across all adversarial attack types compared to TTC  
 472 and R-TPT, regardless of the deployed victim model, underscoring its exceptional extensibility.

## 473 5 CONCLUSION

474 In this paper, we propose TAME, a novel test-time adversarial defense method designed to improve  
 475 the zero-shot adversarial robustness of medical vision-language models. By leveraging the semantic  
 476 fragility of adversarial perturbations, TAME effectively restores model predictions through an ad-  
 477 versarial restoration map trained specifically for each test image, requiring only a single update step.  
 478 To mitigate adverse effects on clean inputs, we further introduce an adaptive weighting mechanism  
 479 that balances the trade-off between adversarial robustness and clean accuracy, eliminating the need  
 480 for manual hyperparameter tuning. Extensive experiments across multiple adversarial attacks and  
 481 11 medical datasets spanning 9 imaging modalities demonstrate the superiority of our approach, in-  
 482 dicating that TAME not only outperforms existing defense strategies but also generalizes effectively  
 483 to adversarially fine-tuned models. Future work will investigate the extension of this paradigm to a  
 484 wider range of medical modalities and a more extensive suite of adversarial attacks.

486 REPRODUCIBILITY STATEMENT  
487488 As recommended, we state the reproducibility of this study here. All 11 medical datasets utilized in  
489 this paper are public, and the download links are shown as follows:  
490

- 491 • **BTMRI**: [https://drive.google.com/file/d/1\\_1JLZRUMczqZqoN-dNqkAzGzmi4ONoU5/view?usp=sharing](https://drive.google.com/file/d/1_1JLZRUMczqZqoN-dNqkAzGzmi4ONoU5/view?usp=sharing)
- 492 • **BUSI**: <https://drive.google.com/file/d/1hB5M7wcAUTV9EtjYrijACoQ36R6VmQaa/view?usp=sharing>
- 493 • **COVID-QU-Ex**: [https://drive.google.com/file/d/1zMLN5q5e\\_tmH-deSZQiY4Xq0M1EqCrML/view?usp=sharing](https://drive.google.com/file/d/1zMLN5q5e_tmH-deSZQiY4Xq0M1EqCrML/view?usp=sharing)
- 494 • **CTKidney**: <https://drive.google.com/file/d/1PBZ299k--mZL8JU7nhC1Wy8yEmlqmVDh/view?usp=sharing>
- 495 • **DermaMNIST**: <https://drive.google.com/file/d/1Jxd1-DWljunRDZ8fY80dl5zUMefriQxt/view?usp=sharing>
- 496 • **Kvasir**: [https://drive.google.com/file/d/1T\\_cqnNIjmGazNeg6gziarvCNWGsFEkRi/view?usp=sharing](https://drive.google.com/file/d/1T_cqnNIjmGazNeg6gziarvCNWGsFEkRi/view?usp=sharing)
- 497 • **CHMNIST**: [https://drive.google.com/file/d/1tyQiYQmqAGNaY4SCK\\_8U5vEbbaa1AD-g/view?usp=sharing](https://drive.google.com/file/d/1tyQiYQmqAGNaY4SCK_8U5vEbbaa1AD-g/view?usp=sharing)
- 498 • **LC25000**: <https://drive.google.com/file/d/1YIu5fqMXgyemisiL1L1HCvES2nVpCtun/view?usp=sharing>
- 499 • **RETINA**: <https://drive.google.com/file/d/18U-Gc22h5QryomNNzY4r4Qfrq52yf5EO/view?usp=sharing>
- 500 • **KneeXray**: [https://drive.google.com/file/d/1DBVraYJmxy2UcQ\\_nGLYvTB2reITOm453/view?usp=sharing](https://drive.google.com/file/d/1DBVraYJmxy2UcQ_nGLYvTB2reITOm453/view?usp=sharing)
- 501 • **OCTMNIST**: <https://drive.google.com/file/d/1mYZNWxbPxnnVvcwHQYybA8gdMzQAOem/view?usp=sharing>

516 We followed the data processing pipeline detailed in Koleilat et al. (2025), which is also open-  
517 source (<https://github.com/HealthX-Lab/BiomedCoOp/tree/main>). The results  
518 of the competing methods are reproduced by using their publicly available source codes, and the  
519 corresponding GitHub links are listed below:

- 520 • **Anti-Adv**: <https://github.com/MotasemAlfarra/Combating-Adversaries-with-Anti-Adversaries>
- 521 • **HedgeDefense**: [https://github.com/burcywu/hedge\\_defense](https://github.com/burcywu/hedge_defense)
- 522 • **TTC**: <https://github.com/Sxing2/CLIP-Test-time-Counterattack/tree/main>
- 523 • **R-TPT**: <https://github.com/TomSheng21/R-TPT/tree/main>
- 524 • **FARE**: <https://github.com/chs20/RobustVLM>
- 525 • **PMG**: <https://github.com/serendipity1122/Pre-trained-Model-Guided-Fine-Tuning-for-Zero-Shot-Adversarial-Robustness>
- 526 • **TeCoA**: <https://github.com/cvlab-columbia/ZSRobust4FoundationModel>

533 The hyper-parameter configurations of our TAME can be found in Section 4.2, and the code and the  
534 computational environment will be available on GitHub.  
535537 REFERENCES  
538539 Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultra-  
sound images. *Data in brief*, 28:104863, 2020.

- 540 Motasem Alfarra, Juan C Pérez, Ali Thabet, Adel Bibi, Philip HS Torr, and Bernard Ghanem. Com-  
 541 bating adversaries with anti-adversaries. In *Proceedings of the AAAI Conference on Artificial*  
 542 *Intelligence*, volume 36, pp. 5992–6000, 2022.
- 543
- 544 Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes Van Diest, Bram Van Ginneken, Nico  
 545 Karssemeijer, Geert Litjens, Jeroen AWM Van Der Laak, Meyke Hermsen, Quirine F Manson,  
 546 Maschenka Balkenhol, et al. Diagnostic assessment of deep learning algorithms for detection of  
 547 lymph node metastases in women with breast cancer. *Jama*, 318(22):2199–2210, 2017.
- 548 Andrew A Borkowski, Marilyn M Bui, L Brannon Thomas, Catherine P Wilson, Lauren A DeLand,  
 549 and Stephen M Mastorides. Lung and colon cancer histopathological image dataset (lc25000).  
 550 *arXiv preprint arXiv:1912.12142*, 2019.
- 551
- 552 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101–mining discriminative compo-  
 553 nents with random forests. In *European conference on computer vision*, pp. 446–461. Springer,  
 554 2014.
- 555 Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017*  
 556 *ieee symposium on security and privacy (sp)*, pp. 39–57. Ieee, 2017.
- 557
- 558 Pingjun Chen. Knee osteoarthritis severity grading dataset. *Mendeley Data*, 1(10.17632):30784984,  
 559 2018.
- 560
- 561 Tianlong Chen, Sijia Liu, Shiyu Chang, Yu Cheng, Lisa Amini, and Zhangyang Wang. Adversarial  
 562 robustness: From self-supervised pre-training to fine-tuning. In *Proceedings of the IEEE/CVF*  
 563 *Conference on Computer Vision and Pattern Recognition*, pp. 699–708, 2020.
- 564 Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. De-  
 565 scribing textures in the wild. In *Proceedings of the IEEE Conference on Computer Vision and*  
 566 *Pattern Recognition*, pp. 3606–3613, 2014.
- 567
- 568 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised  
 569 feature learning. In *Proceedings of the fourteenth international conference on artificial intelli-  
 570 gence and statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 571
- 572 Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gut-  
 573 man, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion  
 574 analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging  
 575 collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.
- 576
- 577 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble  
 578 of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–  
 2216. PMLR, 2020.
- 579
- 580 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
 581 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
 582 pp. 248–255. Ieee, 2009.
- 583
- 584 Junhao Dong, Junxi Chen, Xiaohua Xie, Jianhuang Lai, and Hao Chen. Survey on adversarial attack  
 585 and defense for medical image analysis: Methods and challenges. *ACM Computing Surveys*, 57  
 (3):1–38, 2024.
- 586
- 587 Sedigheh Eslami, Christoph Meinel, and Gerard De Melo. Pubmedclip: How much does clip benefit  
 588 visual question answering in the medical domain? In *Findings of the Association for Compu-  
 589 tational Linguistics: EACL 2023*, pp. 1181–1193, 2023.
- 590
- 591 Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE Trans-  
 592 actions on Pattern Analysis and Machine Intelligence*, 28(4):594–611, 2006.
- 593
- Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset (technical  
 report). *California Institute of Technology*, 2007.

- 594 Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens van der Maaten. Countering adversarial  
 595 images using input transformations. In *International Conference on Learning Representations*,  
 596 2018.
- 597 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recogni-  
 598 tion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
 599 770–778, 2016.
- 600 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset  
 601 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected  
 602 Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 603 Md Nazmul Islam, Mehedi Hasan, Md Kabir Hossain, Md Golam Rabiul Alam, Md Zia Uddin, and  
 604 Ahmet Soylu. Vision transformer and explainable transfer learning models for auto detection of  
 605 kidney cyst, stone and tumor from ct-radiography. *Scientific Reports*, 12(1):11440, 2022.
- 606 Jakob Nikolas Kather, Cleo-Aron Weis, Francesco Bianconi, Susanne M Melchers, Lothar R Schad,  
 607 Timo Gaiser, Alexander Marx, and Frank Gerrit Zöllner. Multi-class texture analysis in colorectal  
 608 cancer histology. *Scientific reports*, 6(1):1–11, 2016.
- 609 Daniel S Kermany, Michael Goldbaum, Wenjia Cai, Carolina CS Valentim, Huiying Liang, Sally L  
 610 Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, et al. Identifying medical diag-  
 611 noses and treatable diseases by image-based deep learning. *cell*, 172(5):1122–1131, 2018.
- 612 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International  
 613 Conference on Learning Representations*, 2014.
- 614 Thomas Köhler, Attila Budai, Martin F Kraus, Jan Odstrčilík, Georg Michelson, and Joachim  
 615 Hornegger. Automatic no-reference quality assessment for retinal fundus images using vessel  
 616 segmentation. In *Proceedings of the 26th IEEE international symposium on computer-based  
 617 medical systems*, pp. 95–100. IEEE, 2013.
- 618 Taha Koleilat, Hojat Asgariandehkordi, Hassan Rivaz, and Yiming Xiao. Biomedcoop: Learning  
 619 to prompt for biomedical vision-language models. In *Proceedings of the Computer Vision and  
 620 Pattern Recognition Conference*, pp. 14766–14776, 2025.
- 621 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained  
 622 categorization. In *Proceedings of the IEEE International Conference on Computer Vision Work-  
 623 shops*, pp. 554–561, 2013.
- 624 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
 625 2009.
- 626 Zhengfeng Lai, Joohi Chauhan, Brittany N Dugger, and Chen-Nee Chuah. Bridging the pathology  
 627 domain gap: Efficiently adapting clip for pathology image analysis with limited labeled data. In  
 628 *European Conference on Computer Vision*, pp. 256–273. Springer, 2024.
- 629 Lin Li, Haoyan Guan, Jianing Qiu, and Michael Spratling. One prompt word is enough to boost  
 630 adversarial robustness for pre-trained vision-language models. In *Proceedings of the IEEE/CVF  
 631 Conference on Computer Vision and Pattern Recognition*, pp. 24408–24419, 2024.
- 632 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
 633 Towards deep learning models resistant to adversarial attacks. In *International Conference on  
 634 Learning Representations*, 2018.
- 635 Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained  
 636 visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 637 Chengzhi Mao, Mia Chiquier, Hao Wang, Junfeng Yang, and Carl Vondrick. Adversarial attacks are  
 638 reversible with natural supervision. In *Proceedings of the IEEE/CVF International Conference  
 639 on Computer Vision*, pp. 661–671, 2021.

- 648 Chengzhi Mao, Scott Geng, Junfeng Yang, Xin Wang, and Carl Vondrick. Understanding zero-  
 649 shot adversarial robustness for large-scale models. In *The Eleventh International Conference on*  
 650 *Learning Representations*, 2023.
- 651
- 652 Msoud Nickparvar. Brain tumor mri dataset. *Kaggle*, 2021.
- 653
- 654 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number  
 655 of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pp.  
 656 722–729. IEEE, 2008.
- 657
- 658 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012*  
 659 *IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- 660
- 661 Juan C Pérez, Motasem Alfarra, Guillaume Jeanneret, Laura Rueda, Ali Thabet, Bernard Ghanem,  
 662 and Pablo Arbeláez. Enhancing adversarial robustness via test-time transformation ensembling.  
 663 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 81–91, 2021.
- 664
- 665 Konstantin Pogorelov, Kristin Ranheim Randel, Carsten Griwodz, Sigrun Losada Eskeland, Thomas  
 666 de Lange, Dag Johansen, Concetto Spampinato, Duc-Tien Dang-Nguyen, Mathias Lux, Peter  
 667 Thelin Schmidt, et al. Kvasir: A multi-class image dataset for computer aided gastrointestinal  
 668 disease detection. In *Proceedings of the 8th ACM on Multimedia Systems Conference*, pp. 164–  
 669 169, 2017.
- 670
- 671 Prasanna Porwal, Samiksha Pachade, Ravi Kamble, Manesh Kokare, Girish Deshmukh, Vivek  
 672 Sahasrabuddhe, and Fabrice Meriaudeau. Indian diabetic retinopathy image dataset (idrid): a  
 673 database for diabetic retinopathy screening research. *Data*, 3(3):25, 2018.
- 674
- 675 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 676 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 677 models from natural language supervision. In *International conference on machine learning*, pp.  
 678 8748–8763. PMLR, 2021a.
- 679
- 680 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 681 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 682 models from natural language supervision. In *International conference on machine learning*, pp.  
 683 8748–8763. PMLR, 2021b.
- 684
- 685 Christian Schlarbmann, Naman Deep Singh, Francesco Croce, and Matthias Hein. Robust clip: Un-  
 686 supervised adversarial fine-tuning of vision embeddings for robust large vision-language models.  
 687 In *International Conference on Machine Learning*, pp. 43685–43704. PMLR, 2024.
- 688
- 689 Lijun Sheng, Jian Liang, Zilei Wang, and Ran He. R-tpt: Improving adversarial robustness of vision-  
 690 language models through test-time prompt tuning. In *Proceedings of the IEEE/CVF Conference*  
 691 *on Computer Vision and Pattern Recognition*, pp. 29958–29967, 2025.
- 692
- 693 Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song,  
 694 David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, et al.  
 695 Bioclip: A vision foundation model for the tree of life. In *Proceedings of the IEEE/CVF Confer-  
 696 ence on Computer Vision and Pattern Recognition*, pp. 19412–19424, 2024.
- 697
- 698 Anas M Tahir, Muhammad EH Chowdhury, Amith Khandakar, Tawsifur Rahman, Yazan Qiblawey,  
 699 Uzair Khurshid, Serkan Kiranyaz, Nabil Ibtehaz, M Sohel Rahman, Somaya Al-Maadeed, et al.  
 Covid-19 infection localization and severity grading from chest x-ray images. *Computers in  
 700 biology and medicine*, 139:105002, 2021.
- 701
- 702 Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of  
 703 multi-source dermatoscopic images of common pigmented skin lesions. *Scientific data*, 5(1):1–9,  
 704 2018.
- 705
- 706 Sibo Wang, Jie Zhang, Zheng Yuan, and Shiguang Shan. Pre-trained model guided fine-tuning  
 707 for zero-shot adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer  
 708 Vision and Pattern Recognition*, pp. 24502–24511, 2024a.

- 702 Xin Wang, Kai Chen, Xingjun Ma, Zheneng Chen, Jingjing Chen, and Yu-Gang Jiang. Advqdet: De-  
 703 tecting query-based adversarial attacks with adversarial contrastive prompt tuning. In *Proceedings*  
 704 *of the 32nd ACM International Conference on Multimedia*, pp. 6212–6221, 2024b.
- 705
- 706 Xin Wang, Kai Chen, Jiaming Zhang, Jingjing Chen, and Xingjun Ma. TAPT: Test-time adversarial  
 707 prompt tuning for robust inference in vision-language models. In *Proceedings of the Computer*  
 708 *Vision and Pattern Recognition Conference*, pp. 19910–19920, 2025.
- 709 Zeyu Wang, Xianhang Li, Hongru Zhu, and Cihang Xie. Revisiting adversarial training at scale.  
 710 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 711 24675–24685, 2024c.
- 712 Boxi Wu, Heng Pan, Li Shen, Jindong Gu, Shuai Zhao, Zhifeng Li, Deng Cai, Xiaofei He, and Wei  
 713 Liu. Attacking adversarial attacks as a defense. *arXiv preprint arXiv:2106.04938*, 2021.
- 714
- 715 Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:  
 716 Large-scale scene recognition from abbey to zoo. In *2010 IEEE Computer Society Conference on*  
 717 *Computer Vision and Pattern Recognition*, pp. 3485–3492. IEEE, 2010.
- 718 Songlong Xing, Zhengyu Zhao, and Nicu Sebe. Clip is strong enough to fight back: Test-time  
 719 counterattacks towards zero-shot adversarial robustness of clip. In *Proceedings of the Computer*  
 720 *Vision and Pattern Recognition Conference*, pp. 15172–15182, 2025.
- 721
- 722 Ziyi Yin, Muchao Ye, Tianrong Zhang, Tianyu Du, Jinguo Zhu, Han Liu, Jinghui Chen, Ting Wang,  
 723 and Fenglong Ma. Vlattack: Multimodal adversarial attacks on vision-language tasks via pre-  
 724 trained models. *Advances in Neural Information Processing Systems*, 36:52936–52956, 2023.
- 725 Jiaming Zhang, Qi Yi, and Jitao Sang. Towards adversarial attack on vision-language pre-training  
 726 models. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 5005–  
 727 5013, 2022.
- 728
- 729 Jiaming Zhang, Xingjun Ma, Xin Wang, Lingyu Qiu, Jiaqi Wang, Yu-Gang Jiang, and Jitao Sang.  
 730 Adversarial prompt tuning for vision-language models. In *European Conference on Computer*  
 731 *Vision*, pp. 56–72. Springer, 2024.
- 732 Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Pre-  
 733 ston, Rajesh Rao, Mu Wei, Naveen Valluri, et al. Biomedclip: a multimodal biomedical  
 734 foundation model pretrained from fifteen million scientific image-text pairs. *arXiv preprint*  
 735 *arXiv:2303.00915*, 2023.
- 736
- 737 Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min  
 738 Lin. On evaluating adversarial robustness of large vision-language models. *Advances in Neural*  
 739 *Information Processing Systems*, 36:54111–54138, 2023.
- 740 Zihao Zhao, Yuxiao Liu, Han Wu, Mei Wang, Yonghao Li, Sheng Wang, Lin Teng, Disheng Liu,  
 741 Zhiming Cui, Qian Wang, et al. CLIP in medical imaging: A survey. *Medical Image Analysis*,  
 742 pp. 103551, 2025.
- 743
- 744 Yiwei Zhou, Xiaobo Xia, Zhiwei Lin, Bo Han, and Tongliang Liu. Few-shot adversarial prompt  
 745 learning on vision-language models. *Advances in Neural Information Processing Systems*, 37:  
 746 3122–3156, 2024.
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

## APPENDIX

## A DATASET DETAILS

Table 3 presents a summary of the 11 medical datasets used for evaluation in this study, encompassing 9 typical biomedical imaging modalities: MRI, ultrasound, X-ray, CT, dermatoscopy, endoscopy, histopathology, CFP, and OCT.

Table 3: Details of 11 datasets across 9 biomedical imaging modalities used in this study.

Modality	Dataset	Case Number
Magnetic Resonance Imaging (MRI)	BTMRI	1717
Ultrasound	BUSI	236
X-Ray	COVID-QU-Ex	1656
	KneeXray	6351
Computerized Tomography (CT)	CTKidney	3738
Dermatoscopy	DermaMNIST	2005
Endoscopy	Kvasir	1200
Histopathology	CHMNIST	1504
	LC25000	7500
Color Fundus Photography (CFP)	RETINA	1268
Optical Coherence Tomography (OCT)	OCTMNIST	1000

## B LABEL LEAKAGE BY ATTACKS

To validate the phenomenon of “label leakage”, we attacked BiomedCLIP using the PGD method with an attack budget of  $1/255$  and a step size of 10. Specifically, we evaluated three settings: (1) “Chance-level”: a chance-level baseline with random guessing; (2) “Random Noise”: classification using perturbations initialized from random noise; and (3) “Label As Target”: classification using adversarial perturbations generated by the PGD method. For the latter two, We employed a mini-ResNet He et al. (2016) (about  $0.3M$  parameters) as a simple classifier, utilized to predict the class labels from the input perturbations. This classifier is trained by an Adam Kingma & Ba (2014) optimizer using a learning rate of 0.001 for 10 epochs, with an 8:2 train-validation data split. As shown in Figure 7, the “Chance-level” achieves an accuracy of approximately  $1/k$ , where  $k$  denotes the number of categories. The accuracy of “Random Noise” is comparable to this baseline across most datasets, while “Label As Target” exhibits a significantly higher overall accuracy. This finding highlights the potential risk of label leakage via adversarial attacks. We argue that this leakage occurs since adversarial perturbations are optimized along gradient directions that are inherently label-aligned, thereby embedding class-related information at the pixel or feature level. Consequently, a defense approach that learns to recognize and reverse such information could transform adversarial perturbations into signals that are beneficial to the model’s performance. This insight suggests that future research should reconsider the supervision strategy of attack methods to mitigate the risk of label leakage.

## C ADDITIONAL EXPERIMENTAL RESULTS

## C.1 ABLATION STUDY

In this section, we will discuss the effect of the proposed weighting mechanism and analyze the sensitivity of our TAME to the defense budget  $\delta_a$  and the step-size  $\alpha$ . We repeated the experiments on 11 medical datasets using BiomedCLIP as the victim model, and the results are summarized in Table 4. It reveals that (1) the dynamic weight  $\omega$  preserves performance on clean images by sacrificing robustness to adversarial images, where the extremely high adversarial robustness intensifies the suspicion of label leakage during adversarial attacks; (2) our TAME is robust to the variation of  $\delta_a$  and  $\alpha$ ; and (3) the performance on clean and adversarial images generally exhibits opposite

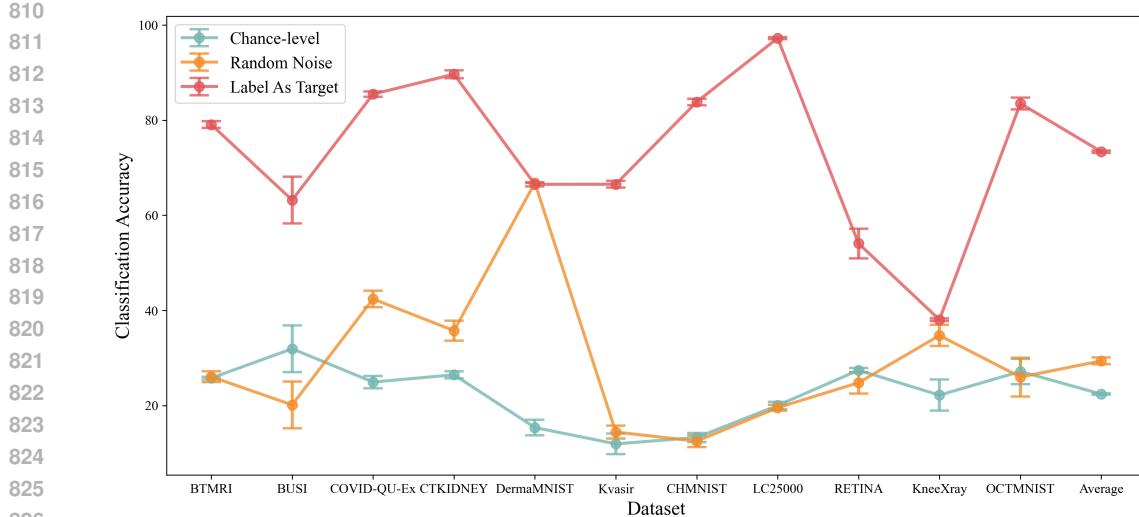


Figure 7: Classification accuracy across 11 datasets under various situations, where the error bars represent standard deviation calculated across 3 trials. Chance-level: the expected performance of making predictions by random guessing. Random Noise: using the perturbations initialized by random noise to train a simple classifier. Label As Target: using the adversarial perturbations produced by PGD to train a simple classifier.

trends as  $\delta_a$  and  $\alpha$  decrease. The trade-off issue between clean and adversarial robustness will be considered in our future work.

## C.2 LARGER ATTACK BUDGET

We enlarged the attack budget from 1/255 to 4/255 to evaluate the effectiveness of our TAME and other compared TAD methods in a more challenging environment. All the experiment configurations of all methods are frozen to avoid additional tuning. Since amplifying the attack budget will not affect the performance on clean images, the results under the clean setting are omitted. We reported the results in Table 6. The results demonstrate that our TAME still remains robust against all attack types even with such a larger attack budget and surpasses other compared methods across most datasets, achieving the best overall accuracy.

## C.3 DEPLOY TO OTHER MEDICAL VLMs

We further discussed the generalizability of the defense methods when deployed to other medical VLMs. The PubMedCLIP Eslami et al. (2023) with a ViT-B/32 backbone is introduced as the victim model, and other experimental configurations remain consistent with those in Table 1. The results displayed in Table 5 reveal that our TAME consistently achieves the strongest overall adversarial robustness against all attack types with an acceptable accuracy on clean images. The label leakage phenomenon can also be clearly observed that the overall accuracies of our TAME against PGD (41.41%), C&W (40.42%), and AutoAttack (36.19%) are much higher than the accuracy of the PubMedCLIP baseline on clean images (27.24%).

## C.4 ADVERSARIAL ROBUSTNESS ON NATURAL IMAGE TASKS

To evaluate the generalizability of our TAME in natural scenes, we followed previous adversarial defense works Xing et al. (2025); Wang et al. (2024a) and conducted experiments on 16 diverse natural image datasets, including four distinct tasks:

- **General Object Recognition:** CIFAR10 Krizhevsky et al. (2009), CIFAR100 Krizhevsky et al. (2009), STL10 Coates et al. (2011), ImageNet Deng et al. (2009), Caltech101 Fei-Fei et al. (2006), and Caltech256 Griffin et al. (2007).

864  
 865 Table 4: Zero-shot adversarial robustness (%) of various variants of our TAME on 11 medical  
 866 datasets. We report the mean and standard deviation calculated across three trials.

867 Dataset	868 Attack	869 BiomedCLIP	870 Ablation study for TAME			
			871 w/o $\omega$	872 $\epsilon_a = \alpha = 4/255$	873 $\epsilon_a = \alpha = 2/255$	874 Ours
875 BTMRI	Clean	56.79	38.36 $\pm$ 1.62	54.07 $\pm$ 1.09	55.44 $\pm$ 0.98	54.13 $\pm$ 1.09
	PGD	0.68 $\pm$ 0.07	76.84 $\pm$ 0.22	58.71 $\pm$ 0.09	54.38 $\pm$ 0.06	61.21 $\pm$ 0.33
	C&W	0.68 $\pm$ 0.03	76.32 $\pm$ 0.29	58.40 $\pm$ 0.40	54.49 $\pm$ 0.48	61.50 $\pm$ 0.43
	AA	0.06 $\pm$ 0.00	71.27 $\pm$ 0.62	60.84 $\pm$ 0.32	57.39 $\pm$ 0.24	61.25 $\pm$ 0.19
876 BUSI	Clean	59.75	35.17 $\pm$ 3.51	61.72 $\pm$ 2.77	62.57 $\pm$ 2.69	62.71 $\pm$ 2.07
	PGD	0.00 $\pm$ 0.00	74.44 $\pm$ 0.87	62.99 $\pm$ 2.09	49.01 $\pm$ 2.30	68.08 $\pm$ 2.45
	C&W	0.00 $\pm$ 0.00	75.42 $\pm$ 1.51	65.68 $\pm$ 1.51	49.01 $\pm$ 2.23	70.90 $\pm$ 0.53
	AA	0.00 $\pm$ 0.00	70.48 $\pm$ 1.00	66.67 $\pm$ 1.78	58.61 $\pm$ 2.94	65.54 $\pm$ 0.20
877 COVID-QU-Ex	Clean	43.82	31.72 $\pm$ 0.30	43.39 $\pm$ 0.30	44.38 $\pm$ 0.28	36.38 $\pm$ 0.26
	PGD	0.00 $\pm$ 0.00	66.47 $\pm$ 0.18	47.54 $\pm$ 0.09	38.67 $\pm$ 0.38	54.41 $\pm$ 0.22
	C&W	0.00 $\pm$ 0.00	64.99 $\pm$ 0.04	47.14 $\pm$ 0.11	38.46 $\pm$ 0.36	53.70 $\pm$ 0.42
	AA	0.00 $\pm$ 0.00	60.33 $\pm$ 0.36	54.66 $\pm$ 0.69	46.12 $\pm$ 0.45	54.00 $\pm$ 0.48
878 CTKIDNEY	Clean	42.43	30.48 $\pm$ 0.47	43.39 $\pm$ 0.51	45.15 $\pm$ 1.15	40.36 $\pm$ 0.38
	PGD	0.87 $\pm$ 0.03	56.08 $\pm$ 0.09	48.66 $\pm$ 0.32	43.03 $\pm$ 0.36	53.01 $\pm$ 0.60
	C&W	0.88 $\pm$ 0.02	55.80 $\pm$ 0.06	48.21 $\pm$ 0.73	42.86 $\pm$ 0.36	52.02 $\pm$ 0.60
	AA	0.05 $\pm$ 0.00	53.14 $\pm$ 0.23	50.16 $\pm$ 0.36	47.68 $\pm$ 0.28	50.42 $\pm$ 0.61
879 DermaMNIST	Clean	38.80	25.06 $\pm$ 0.20	31.60 $\pm$ 0.42	34.20 $\pm$ 0.83	27.95 $\pm$ 0.63
	PGD	0.00 $\pm$ 0.00	43.67 $\pm$ 0.06	38.63 $\pm$ 0.44	28.91 $\pm$ 0.14	40.28 $\pm$ 0.59
	C&W	0.00 $\pm$ 0.00	45.50 $\pm$ 0.72	38.14 $\pm$ 0.28	28.50 $\pm$ 0.47	41.30 $\pm$ 0.60
	AA	0.00 $\pm$ 0.00	42.37 $\pm$ 0.21	45.22 $\pm$ 0.30	38.57 $\pm$ 0.31	41.99 $\pm$ 0.25
880 Kvasir	Clean	54.58	27.64 $\pm$ 1.46	48.94 $\pm$ 1.29	50.61 $\pm$ 1.09	48.36 $\pm$ 1.00
	PGD	0.00 $\pm$ 0.00	65.72 $\pm$ 0.45	53.72 $\pm$ 0.14	46.47 $\pm$ 0.44	59.61 $\pm$ 0.08
	C&W	0.00 $\pm$ 0.00	63.83 $\pm$ 0.54	52.30 $\pm$ 0.34	46.45 $\pm$ 0.17	58.11 $\pm$ 0.22
	AA	0.00 $\pm$ 0.00	66.92 $\pm$ 0.71	61.80 $\pm$ 0.32	56.03 $\pm$ 0.49	63.72 $\pm$ 0.67
881 CHMNIST	Clean	30.65	18.15 $\pm$ 0.79	29.77 $\pm$ 0.42	30.65 $\pm$ 0.30	21.77 $\pm$ 0.71
	PGD	0.00 $\pm$ 0.00	28.81 $\pm$ 0.11	30.52 $\pm$ 0.80	26.02 $\pm$ 0.58	25.97 $\pm$ 0.17
	C&W	0.02 $\pm$ 0.03	28.57 $\pm$ 0.32	30.96 $\pm$ 0.36	25.73 $\pm$ 0.74	24.98 $\pm$ 0.73
	AA	0.00 $\pm$ 0.00	31.61 $\pm$ 0.26	37.66 $\pm$ 1.00	34.35 $\pm$ 0.28	30.83 $\pm$ 0.32
882 LC25000	Clean	50.01	36.87 $\pm$ 0.38	46.71 $\pm$ 0.27	48.94 $\pm$ 0.02	44.04 $\pm$ 0.23
	PGD	0.01 $\pm$ 0.00	59.96 $\pm$ 0.33	51.22 $\pm$ 0.45	42.66 $\pm$ 0.18	55.75 $\pm$ 0.56
	C&W	0.02 $\pm$ 0.01	56.36 $\pm$ 0.06	48.99 $\pm$ 0.96	41.11 $\pm$ 0.14	52.47 $\pm$ 0.41
	AA	0.01 $\pm$ 0.00	54.81 $\pm$ 0.23	53.48 $\pm$ 0.22	48.17 $\pm$ 0.22	54.62 $\pm$ 0.20
883 RETINA	Clean	26.26	27.55 $\pm$ 0.58	26.10 $\pm$ 0.17	25.94 $\pm$ 0.82	26.18 $\pm$ 0.30
	PGD	0.00 $\pm$ 0.00	35.04 $\pm$ 0.15	20.22 $\pm$ 0.73	15.12 $\pm$ 0.38	26.13 $\pm$ 0.53
	C&W	0.00 $\pm$ 0.00	35.89 $\pm$ 0.78	20.24 $\pm$ 0.60	14.85 $\pm$ 0.30	26.21 $\pm$ 0.48
	AA	0.00 $\pm$ 0.00	34.33 $\pm$ 0.38	27.50 $\pm$ 0.67	21.85 $\pm$ 0.51	26.68 $\pm$ 0.21
884 KneeXray	Clean	29.47	37.80 $\pm$ 0.00	38.33 $\pm$ 0.30	37.92 $\pm$ 0.54	38.38 $\pm$ 0.44
	PGD	0.00 $\pm$ 0.00	50.14 $\pm$ 0.22	34.60 $\pm$ 0.54	25.08 $\pm$ 0.34	46.15 $\pm$ 0.33
	C&W	0.00 $\pm$ 0.00	44.67 $\pm$ 0.34	29.09 $\pm$ 0.45	20.27 $\pm$ 0.95	41.08 $\pm$ 0.06
	AA	0.00 $\pm$ 0.00	40.18 $\pm$ 0.16	39.77 $\pm$ 0.32	36.49 $\pm$ 0.25	39.61 $\pm$ 0.26
885 OCTMNIST	Clean	29.90	34.93 $\pm$ 0.25	33.90 $\pm$ 0.24	33.47 $\pm$ 0.95	34.10 $\pm$ 0.36
	PGD	6.27 $\pm$ 0.68	41.67 $\pm$ 0.68	39.23 $\pm$ 0.57	33.30 $\pm$ 0.08	39.40 $\pm$ 0.29
	C&W	6.37 $\pm$ 0.17	41.47 $\pm$ 0.42	39.53 $\pm$ 0.82	33.90 $\pm$ 0.08	39.63 $\pm$ 0.38
	AA	0.00 $\pm$ 0.00	38.30 $\pm$ 0.57	36.30 $\pm$ 0.92	34.87 $\pm$ 0.87	37.10 $\pm$ 0.57
886 Average	Clean	42.04	31.25 $\pm$ 0.25	41.63 $\pm$ 0.34	42.66 $\pm$ 0.12	39.49 $\pm$ 0.21
	PGD	0.71 $\pm$ 0.06	54.44 $\pm$ 0.06	44.19 $\pm$ 0.20	36.60 $\pm$ 0.19	48.18 $\pm$ 0.18
	C&W	0.72 $\pm$ 0.01	53.53 $\pm$ 0.14	43.52 $\pm$ 0.25	35.97 $\pm$ 0.20	47.45 $\pm$ 0.14
	AA	0.01 $\pm$ 0.00	51.25 $\pm$ 0.12	48.55 $\pm$ 0.29	43.65 $\pm$ 0.32	47.80 $\pm$ 0.11

- 904  
 905 • **Domain-specific Classification:** FGVC Aircraft Maji et al. (2013), EuroSAT Helber et al.  
 906 (2019), DTD Cimpoi et al. (2014), and PCAM Bejnordi et al. (2017).  
 907 • **Fine-grained Recognition:** OxfordPets Parkhi et al. (2012), Flowers102 Nilsback & Zis-  
 908 serman (2008), Food101 Bossard et al. (2014), and StanfordCars Krause et al. (2013).  
 909 • **Scene Understanding:** SUN397 Xiao et al. (2010) and Country211 Radford et al. (2021b).

910 The pre-trained CLIP served as the victim model, and the PGD method was utilized as the adversary.  
 911 Following Xing et al. (2025), we set the attack budget  $\epsilon_p$  and the number of update steps for PGD  
 912 to 1/255 and 10, respectively. It should be noted that our TAME was deployed directly without  
 913 any manual tuning. We compared our TAME with four AFT methods (CLIP-FT Xing et al. (2025),  
 914 TeCoA Mao et al. (2023), PMG Wang et al. (2024a), and FARE Schlarmann et al. (2024)) and four  
 915 TAD methods (TTE Pérez et al. (2021), Anti-Adv Alfarra et al. (2022), HD Wu et al. (2021), and  
 916 TTC Xing et al. (2025)). The results shown in Table 7 demonstrate that our TAME achieves superior  
 917 performance on 9 downstream datasets and the best overall accuracy.

918  
 919 Table 5: Zero-shot adversarial robustness (%) of our TAME, the PubMedCLIP baseline, and other  
 920 competing TAD methods on 11 medical datasets. We report the mean and standard deviation cal-  
 921 culated across three trials. For each dataset, the highest performance under the Clean, PGD, C&W,  
 922 and AutoAttack (AA) settings is highlighted in **red**, **blue**, **green**, and **purple**, respectively.

Dataset	Attack	PubMedCLIP	Anti-Adv	HedgeDefense	TTC	R-TPT	TAME
BTMRI	Clean	40.59	40.48 $\pm$ 0.00	<b>48.13</b> $\pm$ 0.11	40.26 $\pm$ 0.34	37.47 $\pm$ 0.21	34.75 $\pm$ 0.46
	PGD	0.37 $\pm$ 0.03	3.44 $\pm$ 0.30	8.17 $\pm$ 0.12	51.66 $\pm$ 0.34	27.08 $\pm$ 0.21	<b>62.90</b> $\pm$ 0.46
	C&W	0.41 $\pm$ 0.00	4.39 $\pm$ 0.19	8.58 $\pm$ 0.15	50.73 $\pm$ 0.22	26.91 $\pm$ 0.17	<b>58.80</b> $\pm$ 0.74
	AA	0.06 $\pm$ 0.00	3.94 $\pm$ 0.06	18.37 $\pm$ 0.10	40.46 $\pm$ 0.56	30.77 $\pm$ 0.12	<b>48.67</b> $\pm$ 0.38
BUSI	Clean	54.66	54.66 $\pm$ 0.00	<b>55.79</b> $\pm$ 0.40	50.57 $\pm$ 0.20	54.80 $\pm$ 0.20	43.79 $\pm$ 3.32
	PGD	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	4.10 $\pm$ 0.20	53.25 $\pm$ 0.53	53.53 $\pm$ 0.20	<b>72.32</b> $\pm$ 0.72
	C&W	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	10.31 $\pm$ 1.11	54.38 $\pm$ 2.30	52.82 $\pm$ 0.40	<b>74.44</b> $\pm$ 1.11
	AA	0.00 $\pm$ 0.00	3.81 $\pm$ 1.04	3.25 $\pm$ 0.20	24.72 $\pm$ 1.06	54.94 $\pm$ 0.20	<b>75.28</b> $\pm$ 0.72
COVID-QU-Ex	Clean	6.61	7.43 $\pm$ 0.00	6.79 $\pm$ 0.02	<b>15.78</b> $\pm$ 0.25	6.63 $\pm$ 0.02	11.13 $\pm$ 0.37
	PGD	0.00 $\pm$ 0.00	5.83 $\pm$ 0.16	0.04 $\pm$ 0.02	11.41 $\pm$ 0.15	5.67 $\pm$ 0.01	<b>17.17</b> $\pm$ 0.30
	C&W	0.00 $\pm$ 0.00	6.38 $\pm$ 0.04	0.09 $\pm$ 0.01	11.21 $\pm$ 0.10	5.69 $\pm$ 0.05	<b>16.35</b> $\pm$ 0.20
	AA	0.02 $\pm$ 0.00	1.79 $\pm$ 0.01	0.32 $\pm$ 0.03	12.32 $\pm$ 0.14	6.22 $\pm$ 0.03	<b>14.25</b> $\pm$ 0.24
CTKIDNEY	Clean	22.82	22.82 $\pm$ 0.00	18.98 $\pm$ 0.07	21.62 $\pm$ 0.19	<b>23.46</b> $\pm$ 0.10	21.85 $\pm$ 0.20
	PGD	0.44 $\pm$ 0.03	0.44 $\pm$ 0.03	1.75 $\pm$ 0.05	25.62 $\pm$ 0.27	19.73 $\pm$ 0.16	<b>37.69</b> $\pm$ 0.27
	C&W	0.95 $\pm$ 0.02	0.94 $\pm$ 0.02	4.11 $\pm$ 0.15	28.83 $\pm$ 0.07	19.15 $\pm$ 0.11	<b>39.41</b> $\pm$ 0.39
	AA	0.05 $\pm$ 0.00	0.05 $\pm$ 0.00	4.00 $\pm$ 0.10	20.35 $\pm$ 0.51	21.41 $\pm$ 0.01	<b>29.05</b> $\pm$ 0.38
DermaMNIST	Clean	16.36	16.06 $\pm$ 0.00	<b>27.48</b> $\pm$ 0.12	20.42 $\pm$ 0.61	18.29 $\pm$ 0.08	16.23 $\pm$ 0.39
	PGD	0.00 $\pm$ 0.00	0.13 $\pm$ 0.05	3.19 $\pm$ 0.15	15.79 $\pm$ 0.22	14.51 $\pm$ 0.07	<b>29.31</b> $\pm$ 0.09
	C&W	0.00 $\pm$ 0.00	0.23 $\pm$ 0.06	4.66 $\pm$ 0.22	15.35 $\pm$ 0.27	14.00 $\pm$ 0.16	<b>26.63</b> $\pm$ 0.27
	AA	0.00 $\pm$ 0.00	0.37 $\pm$ 0.02	15.16 $\pm$ 0.12	17.92 $\pm$ 0.41	16.68 $\pm$ 0.17	<b>24.82</b> $\pm$ 0.27
Kvasir	Clean	13.00	12.83 $\pm$ 0.00	12.92 $\pm$ 0.12	<b>13.58</b> $\pm$ 0.54	13.03 $\pm$ 0.08	9.81 $\pm$ 0.32
	PGD	0.00 $\pm$ 0.00	0.03 $\pm$ 0.04	0.33 $\pm$ 0.07	14.97 $\pm$ 0.42	12.17 $\pm$ 0.07	<b>19.03</b> $\pm$ 0.08
	C&W	0.00 $\pm$ 0.00	0.19 $\pm$ 0.04	0.64 $\pm$ 0.04	15.03 $\pm$ 0.28	12.17 $\pm$ 0.07	<b>18.53</b> $\pm$ 0.21
	AA	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.36 $\pm$ 0.10	16.08 $\pm$ 0.30	12.86 $\pm$ 0.17	<b>16.75</b> $\pm$ 0.12
CHMNIST	Clean	20.48	19.61 $\pm$ 0.00	17.60 $\pm$ 0.40	14.23 $\pm$ 0.14	22.78 $\pm$ 0.13	<b>27.68</b> $\pm$ 0.37
	PGD	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.16 $\pm$ 0.03	18.57 $\pm$ 0.42	18.48 $\pm$ 0.42	<b>34.62</b> $\pm$ 0.27
	C&W	0.00 $\pm$ 0.00	0.07 $\pm$ 0.00	0.24 $\pm$ 0.14	18.53 $\pm$ 0.13	18.28 $\pm$ 0.30	<b>33.42</b> $\pm$ 0.27
	AA	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	1.64 $\pm$ 0.21	17.42 $\pm$ 0.00	21.23 $\pm$ 0.31	<b>33.56</b> $\pm$ 0.36
LC25000	Clean	20.71	20.71 $\pm$ 0.00	20.59 $\pm$ 0.05	20.26 $\pm$ 0.05	19.92 $\pm$ 0.01	<b>22.62</b> $\pm$ 0.37
	PGD	1.07 $\pm$ 0.03	1.09 $\pm$ 0.04	4.17 $\pm$ 0.09	16.37 $\pm$ 0.08	19.43 $\pm$ 0.02	<b>38.40</b> $\pm$ 0.23
	C&W	1.62 $\pm$ 0.03	1.65 $\pm$ 0.02	4.49 $\pm$ 0.04	16.09 $\pm$ 0.02	19.37 $\pm$ 0.03	<b>38.09</b> $\pm$ 0.20
	AA	0.20 $\pm$ 0.00	0.24 $\pm$ 0.01	2.39 $\pm$ 0.04	19.71 $\pm$ 0.05	19.62 $\pm$ 0.03	<b>39.85</b> $\pm$ 0.07
RETINA	Clean	28.31	28.39 $\pm$ 0.00	28.86 $\pm$ 0.17	24.53 $\pm$ 0.23	<b>29.15</b> $\pm$ 0.10	28.79 $\pm$ 0.74
	PGD	0.00 $\pm$ 0.00	0.39 $\pm$ 0.06	1.52 $\pm$ 0.10	39.04 $\pm$ 1.01	20.48 $\pm$ 0.42	<b>50.50</b> $\pm$ 0.16
	C&W	0.00 $\pm$ 0.00	1.16 $\pm$ 0.15	1.81 $\pm$ 0.17	38.83 $\pm$ 0.83	20.45 $\pm$ 0.13	<b>50.34</b> $\pm$ 0.26
	AA	0.00 $\pm$ 0.00	2.68 $\pm$ 0.17	5.23 $\pm$ 0.16	24.16 $\pm$ 0.48	25.47 $\pm$ 0.19	<b>33.41</b> $\pm$ 0.21
KneeXray	Clean	38.65	<b>38.89</b> $\pm$ 0.00	38.51 $\pm$ 0.03	34.56 $\pm$ 0.22	38.65 $\pm$ 0.00	35.41 $\pm$ 0.45
	PGD	0.00 $\pm$ 0.00	0.18 $\pm$ 0.09	0.56 $\pm$ 0.17	44.44 $\pm$ 0.35	28.84 $\pm$ 0.16	<b>52.17</b> $\pm$ 0.44
	C&W	0.00 $\pm$ 0.00	0.48 $\pm$ 0.05	0.72 $\pm$ 0.17	43.84 $\pm$ 0.75	28.95 $\pm$ 0.17	<b>51.17</b> $\pm$ 0.56
	AA	0.00 $\pm$ 0.00	0.30 $\pm$ 0.00	1.87 $\pm$ 0.18	11.35 $\pm$ 0.45	36.11 $\pm$ 0.23	<b>43.32</b> $\pm$ 0.38
OCTMNIST	Clean	37.50	30.00 $\pm$ 0.00	<b>39.10</b> $\pm$ 0.65	29.73 $\pm$ 0.57	34.47 $\pm$ 0.09	27.97 $\pm$ 1.43
	PGD	0.00 $\pm$ 0.00	20.60 $\pm$ 0.71	4.43 $\pm$ 0.33	<b>43.63</b> $\pm$ 0.09	24.80 $\pm$ 0.75	41.40 $\pm$ 0.54
	C&W	0.00 $\pm$ 0.00	21.90 $\pm$ 0.42	4.63 $\pm$ 0.31	<b>43.80</b> $\pm$ 0.14	24.47 $\pm$ 0.31	37.40 $\pm$ 0.64
	AA	0.10 $\pm$ 0.00	19.80 $\pm$ 0.00	10.50 $\pm$ 0.86	<b>48.30</b> $\pm$ 1.63	28.73 $\pm$ 0.39	39.13 $\pm$ 0.41
Average	Clean	27.24	26.54 $\pm$ 0.00	<b>28.61</b> $\pm$ 0.05	25.96 $\pm$ 0.06	27.15 $\pm$ 0.03	25.46 $\pm$ 0.43
	PGD	0.17 $\pm$ 0.00	2.92 $\pm$ 0.03	2.58 $\pm$ 0.03	30.43 $\pm$ 0.19	22.25 $\pm$ 0.12	<b>41.41</b> $\pm$ 0.11
	C&W	0.27 $\pm$ 0.00	3.40 $\pm$ 0.04	3.66 $\pm$ 0.12	30.60 $\pm$ 0.34	22.02 $\pm$ 0.03	<b>40.42</b> $\pm$ 0.14
	AA	0.04 $\pm$ 0.00	3.00 $\pm$ 0.10	5.74 $\pm$ 0.07	22.98 $\pm$ 0.26	24.91 $\pm$ 0.05	<b>36.19</b> $\pm$ 0.10

## D COMPLETE RESULTS

961 Here, we display the complete results for Figure 4, Figure 5, and Table 2 in Figure 8, Figure 9, and  
 962 Table 8, respectively. As demonstrated in Figure 8, the semantic fragility of adversarial perturbations  
 963 is observable universally across all 11 datasets, as evidenced by high KL divergence under weak  
 964 transformations, particularly random cropping and random rotation. Additionally, it can be found  
 965 that the KL divergence of clean images increases at a markedly higher rate with magnitude than that  
 966 of their adversarial counterparts on most datasets. This provides powerful evidence for the design  
 967 of our dynamic weighting mechanism. Figure 9 reveals that both random rotation and random  
 968 cropping yield a higher robustness ratio across all datasets. This can be attributed to that these two  
 969 transformations alter the values and/or positions of most pixels in the image and are common in  
 970 the model training process, thereby leading to low/high robustness on adversarial/clean images. The  
 971 results in Table 8 indicate that our TAME method boosts the performance of three adversarially fine-  
 972 tuned models obtained by distinct AFT methods, achieving superior results on most datasets and the  
 973 highest overall accuracy against all attack types. Additionally, an important finding is the absence

972  
973  
974  
975  
976  
977  
Table 6: Zero-shot adversarial robustness (%) of our TAME, the BiomedCLIP baseline, and other  
978 competing TAD methods on 11 medical datasets with a larger attack budget of 4/255. We report the  
979 mean and standard deviation calculated across three trials. For each dataset, the highest performance  
980 under the PGD, C&W, and AutoAttack (AA) settings is highlighted in **blue**, **green**, and **purple**,  
981 respectively.

Dataset	Attack	BiomedCLIP	Anti-Adv	HedgeDefense	TTC	R-TPT	TAME (Ours)
BTMRI	PGD	0.02 $\pm$ 0.03	4.33 $\pm$ 0.24	0.04 $\pm$ 0.03	16.02 $\pm$ 0.50	<b>41.37</b> $\pm$ 0.27	38.38 $\pm$ 0.41
	C&W	0.02 $\pm$ 0.03	3.65 $\pm$ 0.10	0.04 $\pm$ 0.03	16.93 $\pm$ 0.47	<b>42.30</b> $\pm$ 0.34	38.73 $\pm$ 0.16
	AA	0.02 $\pm$ 0.03	7.20 $\pm$ 0.29	4.85 $\pm$ 0.29	19.92 $\pm$ 0.69	<b>48.71</b> $\pm$ 0.43	43.37 $\pm$ 0.34
BUSI	PGD	0.00 $\pm$ 0.00	0.14 $\pm$ 0.20	0.00 $\pm$ 0.00	12.15 $\pm$ 0.40	<b>26.98</b> $\pm$ 1.06	25.14 $\pm$ 1.63
	C&W	0.00 $\pm$ 0.00	0.85 $\pm$ 0.69	0.00 $\pm$ 0.00	14.12 $\pm$ 0.40	25.57 $\pm$ 1.11	<b>26.41</b> $\pm$ 1.31
	AA	0.00 $\pm$ 0.00	3.81 $\pm$ 0.00	2.97 $\pm$ 0.35	22.74 $\pm$ 1.56	39.83 $\pm$ 1.25	<b>41.38</b> $\pm$ 3.14
COVID-QU-Ex	PGD	0.00 $\pm$ 0.00	0.01 $\pm$ 0.01	0.00 $\pm$ 0.00	17.05 $\pm$ 0.38	12.03 $\pm$ 0.03	<b>27.66</b> $\pm$ 0.41
	C&W	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	17.17 $\pm$ 0.67	13.03 $\pm$ 0.13	<b>27.84</b> $\pm$ 0.17
	AA	0.00 $\pm$ 0.00	0.18 $\pm$ 0.02	9.77 $\pm$ 0.07	18.04 $\pm$ 0.07	19.94 $\pm$ 0.14	<b>34.06</b> $\pm$ 0.11
CTKIDNEY	PGD	0.00 $\pm$ 0.00	0.14 $\pm$ 0.07	0.00 $\pm$ 0.00	6.69 $\pm$ 0.23	29.23 $\pm$ 0.13	<b>39.42</b> $\pm$ 0.59
	C&W	0.00 $\pm$ 0.00	0.60 $\pm$ 0.07	0.00 $\pm$ 0.00	6.20 $\pm$ 0.09	29.28 $\pm$ 0.08	<b>39.18</b> $\pm$ 0.01
	AA	0.00 $\pm$ 0.00	0.35 $\pm$ 0.08	3.98 $\pm$ 0.09	7.31 $\pm$ 0.68	38.29 $\pm$ 0.56	<b>39.59</b> $\pm$ 0.25
DermaMNIST	PGD	0.00 $\pm$ 0.00	0.10 $\pm$ 0.07	0.00 $\pm$ 0.00	3.76 $\pm$ 0.15	5.60 $\pm$ 0.24	<b>12.67</b> $\pm$ 0.43
	C&W	0.00 $\pm$ 0.00	0.15 $\pm$ 0.07	0.02 $\pm$ 0.02	4.22 $\pm$ 0.20	5.50 $\pm$ 0.22	<b>13.40</b> $\pm$ 0.20
	AA	0.00 $\pm$ 0.00	0.25 $\pm$ 0.04	5.84 $\pm$ 0.23	8.21 $\pm$ 0.22	<b>25.09</b> $\pm$ 0.40	21.45 $\pm$ 0.99
Kvasir	PGD	0.00 $\pm$ 0.00	1.22 $\pm$ 0.26	0.00 $\pm$ 0.00	7.08 $\pm$ 0.72	25.22 $\pm$ 0.21	<b>28.31</b> $\pm$ 0.40
	C&W	0.00 $\pm$ 0.00	1.61 $\pm$ 0.17	0.00 $\pm$ 0.00	6.25 $\pm$ 0.20	25.42 $\pm$ 0.30	<b>27.00</b> $\pm$ 0.47
	AA	0.00 $\pm$ 0.00	2.69 $\pm$ 0.24	4.31 $\pm$ 0.08	15.17 $\pm$ 0.18	42.92 $\pm$ 0.36	<b>43.31</b> $\pm$ 0.55
CHMNIST	PGD	0.00 $\pm$ 0.00	4.14 $\pm$ 0.17	0.00 $\pm$ 0.00	1.91 $\pm$ 0.13	<b>8.53</b> $\pm$ 0.22	6.78 $\pm$ 0.73
	C&W	0.00 $\pm$ 0.00	4.19 $\pm$ 0.09	0.04 $\pm$ 0.03	2.08 $\pm$ 0.30	<b>8.22</b> $\pm$ 0.58	7.58 $\pm$ 0.29
	AA	0.00 $\pm$ 0.00	2.44 $\pm$ 0.19	3.75 $\pm$ 0.19	11.95 $\pm$ 0.49	<b>19.44</b> $\pm$ 0.58	17.95 $\pm$ 0.30
LC25000	PGD	0.00 $\pm$ 0.00	0.01 $\pm$ 0.00	0.00 $\pm$ 0.00	2.11 $\pm$ 0.13	27.40 $\pm$ 0.20	<b>33.71</b> $\pm$ 0.27
	C&W	0.00 $\pm$ 0.00	0.03 $\pm$ 0.02	0.00 $\pm$ 0.00	2.16 $\pm$ 0.14	26.97 $\pm$ 0.16	<b>33.75</b> $\pm$ 0.38
	AA	0.00 $\pm$ 0.00	4.84 $\pm$ 0.02	8.73 $\pm$ 0.04	12.45 $\pm$ 0.11	39.76 $\pm$ 0.17	<b>43.20</b> $\pm$ 0.21
RETINA	PGD	0.00 $\pm$ 0.00	4.60 $\pm$ 0.50	0.08 $\pm$ 0.06	11.86 $\pm$ 0.29	3.68 $\pm$ 0.39	<b>10.33</b> $\pm$ 0.29
	C&W	0.00 $\pm$ 0.00	4.57 $\pm$ 0.23	0.08 $\pm$ 0.06	11.04 $\pm$ 0.39	8.91 $\pm$ 0.06	<b>9.94</b> $\pm$ 0.17
	AA	0.00 $\pm$ 0.00	8.44 $\pm$ 0.06	6.97 $\pm$ 0.23	20.11 $\pm$ 0.87	<b>28.71</b> $\pm$ 0.46	17.64 $\pm$ 0.84
KneeXray	PGD	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	16.20 $\pm$ 0.76	6.08 $\pm$ 0.25	<b>34.88</b> $\pm$ 0.25
	C&W	0.00 $\pm$ 0.00	0.08 $\pm$ 0.06	0.00 $\pm$ 0.00	16.06 $\pm$ 0.64	32.45 $\pm$ 0.47	<b>33.39</b> $\pm$ 0.82
	AA	0.00 $\pm$ 0.00	3.70 $\pm$ 0.03	18.78 $\pm$ 0.65	19.97 $\pm$ 0.84	34.10 $\pm$ 0.30	<b>38.45</b> $\pm$ 0.16
OCTMNIST	PGD	0.00 $\pm$ 0.00	0.60 $\pm$ 0.28	0.43 $\pm$ 0.29	15.67 $\pm$ 1.01	25.20 $\pm$ 0.00	<b>30.20</b> $\pm$ 0.22
	C&W	0.07 $\pm$ 0.05	0.60 $\pm$ 0.08	0.43 $\pm$ 0.21	14.73 $\pm$ 0.45	25.20 $\pm$ 0.08	<b>30.50</b> $\pm$ 0.75
	AA	0.00 $\pm$ 0.00	0.40 $\pm$ 0.00	0.20 $\pm$ 0.22	23.13 $\pm$ 0.39	25.13 $\pm$ 0.05	<b>33.03</b> $\pm$ 0.37
Average	PGD	0.00 $\pm$ 0.00	1.39 $\pm$ 0.03	0.05 $\pm$ 0.02	10.04 $\pm$ 0.13	19.21 $\pm$ 0.18	<b>26.14</b> $\pm$ 0.16
	C&W	0.01 $\pm$ 0.01	1.48 $\pm$ 0.07	0.06 $\pm$ 0.02	10.09 $\pm$ 0.20	22.08 $\pm$ 0.00	<b>26.16</b> $\pm$ 0.32
	AA	0.00 $\pm$ 0.00	3.12 $\pm$ 0.02	6.38 $\pm$ 0.10	16.27 $\pm$ 0.28	32.90 $\pm$ 0.04	<b>33.95</b> $\pm$ 0.29

1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
Table 7: Zero-shot adversarial robustness (%) of our TAME, the CLIP baseline, and other competing  
AFT and TAD methods under the PGD attack on 16 datasets. We report the mean and standard  
deviation calculated across three trials. The results marked by  $\ddagger$  are inherited from Xing et al. (2025).  
The best and second-best results in each row are highlighted in **bold** and underline, respectively.

Dataset	CLIP $\ddagger$	Adversarial Fine-tuning (AFT)				Test-time Adversarial Defense (TAD)			TAME (Ours)
		CLIP-FT $\ddagger$	TeCoA $\ddagger$	PMG $\ddagger$	FARE $\ddagger$	TTE $\ddagger$	Anti-Adv $\ddagger$	HD $\ddagger$	
CIFAR10	0.74	3.34	33.61	40.66	19.65	<b>41.35</b> $\pm$ 6.14	12.39 $\pm$ 0.07	17.22 $\pm$ 0.45	28.75 $\pm$ 0.18
CIFAR100	0.26	0.90	18.95	22.52	11.40	<u>20.06</u> $\pm$ 4.03	<u>5.73</u> $\pm$ 0.04	3.86 $\pm$ 0.10	<b>27.64</b> $\pm$ 0.43
STL10	11.0	12.73	70.08	73.08	59.06	<b>78.48</b> $\pm$ 3.83	<u>37.42</u> $\pm$ 0.40	39.02 $\pm$ 0.30	<b>76.70</b> $\pm$ 0.23
ImageNet	1.15	0.93	18.89	21.43	14.00	<u>31.01</u> $\pm$ 4.40	8.67 $\pm$ 0.05	6.63 $\pm$ 0.05	<b>38.41</b> $\pm$ 0.07
Caltech101	14.67	14.21	55.51	61.08	50.74	<b>67.56</b> $\pm$ 3.88	34.81 $\pm$ 0.16	31.53 $\pm$ 0.22	65.78 $\pm$ 0.07
Caltech256	8.47	6.76	43.19	45.91	38.79	<u>60.09</u> $\pm$ 4.03	25.36 $\pm$ 0.17	23.48 $\pm$ 0.10	<b>60.11</b> $\pm$ 0.04
OxfordPets	1.04	2.10	38.35	41.18	31.07	<u>50.33</u> $\pm$ 7.30	20.42 $\pm$ 0.22	12.04 $\pm$ 0.16	57.87 $\pm$ 0.15
Flowers102	1.14	0.54	21.94	23.43	17.14	<u>35.88</u> $\pm$ 4.72	7.16 $\pm$ 0.41	7.29 $\pm$ 0.06	<b>39.14</b> $\pm$ 0.28
FGVCAircraft	0.00	0.00	2.49	2.22	1.35	<u>6.23</u> $\pm$ 1.37	1.27 $\pm$ 0.07	1.26 $\pm$ 0.07	<b>13.96</b> $\pm$ 0.23
StanfordCars	0.02	0.06	8.76	11.65	6.75	<u>22.36</u> $\pm$ 4.17	4.40 $\pm$ 0.30	2.71 $\pm$ 0.09	<b>33.01</b> $\pm$ 0.07
SUN397	1.14	0.94	19.39	22.58	14.91	<u>30.79</u> $\pm$ 4.43	8.05 $\pm$ 0.04	6.40 $\pm$ 0.06	<b>41.52</b> $\pm$ 0.04
Country211	0.04	0.03	1.78	2.12	0.85	<u>3.05</u> $\pm$ 0.89	0.67 $\pm$ 0.05	0.47 $\pm$ 0.02	<b>7.09</b> $\pm$ 0.04
Food101	0.70	0.42	13.90	18.57	11.65	<u>43.94</u> $\pm$ 6.97	13.12 $\pm$ 0.16	8.03 $\pm$ 0.11	57.84 $\pm$ 0.15
EuroSAT	0.03	0.04	11.96	12.60	10.67	<u>6.91</u> $\pm$ 2.13	2.15 $\pm$ 0.04	4.57 $\pm$ 0.09	12.19 $\pm$ 0.24
DTD	2.98	2.39	17.61	14.95	15.64	<u>23.90</u> $\pm$ 2.34	5.62 $\pm$ 0.07	11.63 $\pm$ 0.17	<b>27.32</b> $\pm$ 0.25
PCAM	0.08	1.11	48.24	46.18	16.23	<u>10.62</u> $\pm$ 3.22	4.97 $\pm$ 0.12	44.74 $\pm$ 0.17	<b>52.85</b> $\pm$ 0.20
Average	2.70	2.91	26.54	28.76	20.00	<u>33.28</u> $\pm$ 3.98	12.01 $\pm$ 0.04	13.81 $\pm$ 0.06	<b>39.17</b> $\pm$ 0.02

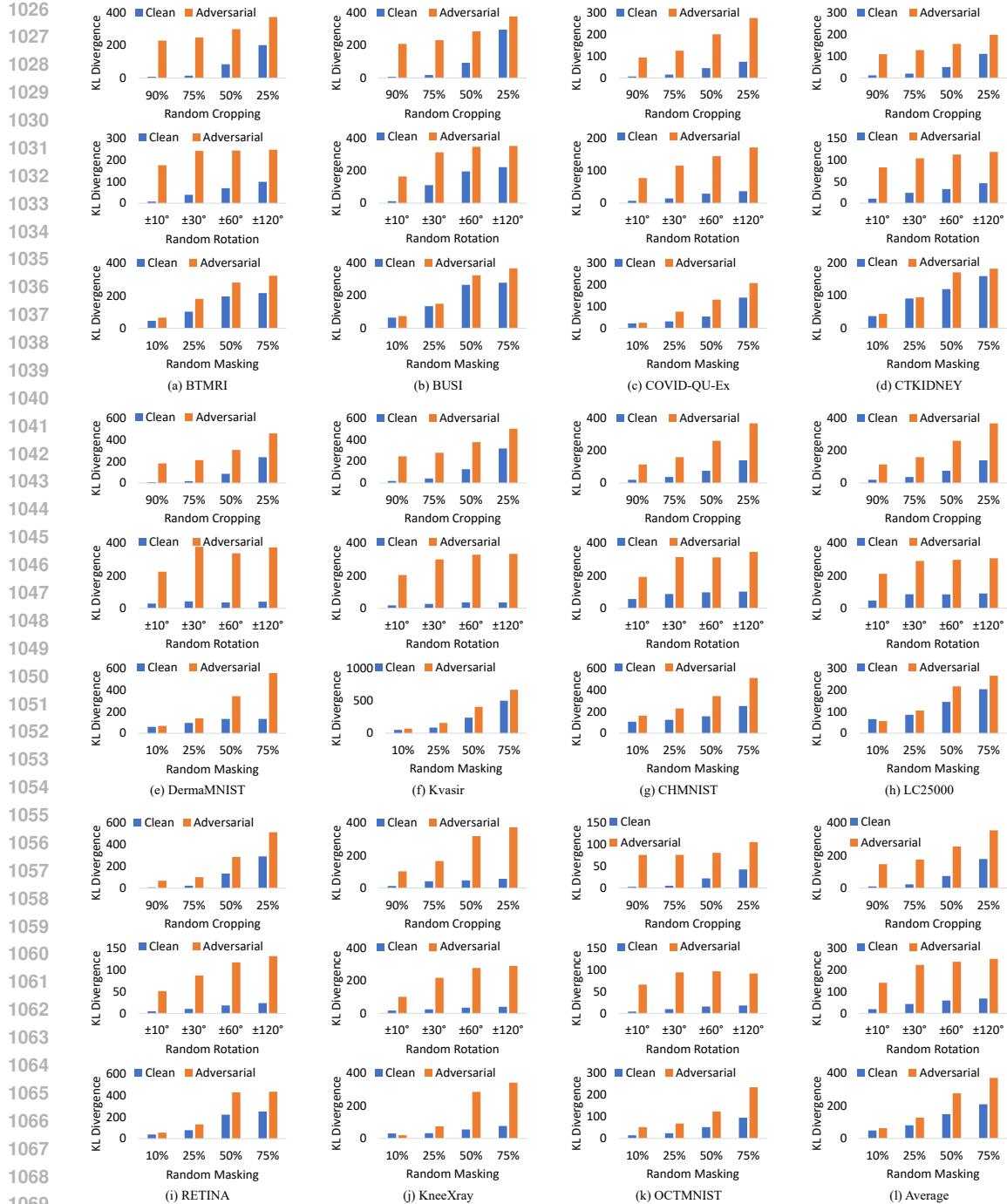


Figure 8: The KL divergence between BiomedCLIP’s predictions before and after applying transformations across all 11 datasets with various modalities.

of label leakage when PMG is employed as the AFT method, which suggests that the occurrence of this phenomenon depends on the specific victim model. Note that TeCoA fails on the BUSI dataset, classifying all samples into the same category with high confidence. Consequently, applying other attack or defense strategies yields identical results.

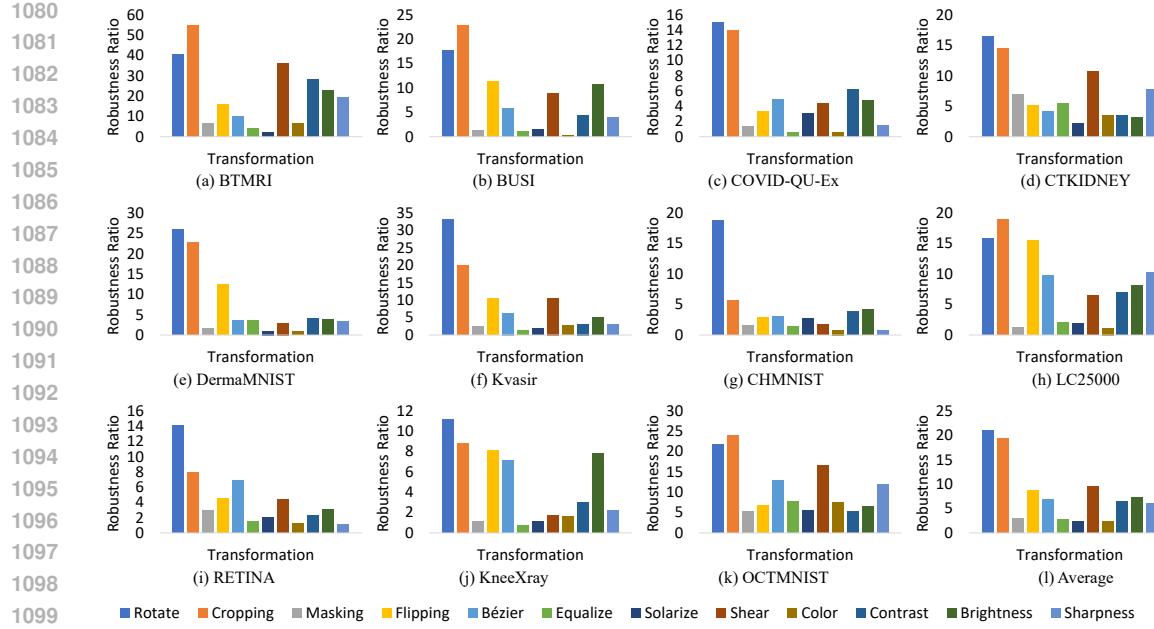


Figure 9: The robustness ratio of 12 transformation strategies across all 11 datasets.

## E THE USE OF LARGE LANGUAGE MODELS

In this paper, we employed DeepSeek-V3 and GPT-5 as assist tools to polish writing and identify potential grammar or spelling errors.

Method	Attack	BTMRI	BUSI	COVID-QU-Es	CTKIDNEY	DermalMNIST	Kvasir	CHMNIST	LC25000	RETINA	KneeXray	OCTMNIST	Average
CLIP (ViT-B/32)	Clean	24.64	38.56	6.36	30.71	28.83	17.08	24.87	29.76	26.50	14.49	25.80	24.33
	PGD	0.02 $\pm$ 0.03	0.00 $\pm$ 0.00	0.19 $\pm$ 0.00	0.01 $\pm$ 0.01	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.07 $\pm$ 0.05	0.01 $\pm$ 0.01	0.45 $\pm$ 0.20	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.07 $\pm$ 0.02
	C&W	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.28 $\pm$ 0.02	0.03 $\pm$ 0.00	0.00 $\pm$ 0.00	0.00 $\pm$ 0.00	0.07 $\pm$ 0.05	0.01 $\pm$ 0.00	0.15 $\pm$ 0.21	0.02 $\pm$ 0.03	0.00 $\pm$ 0.00	0.13 $\pm$ 0.01
	AA	0.00 $\pm$ 0.00	0.42 $\pm$ 0.00	0.00 $\pm$ 0.00	0.21 $\pm$ 0.00	0.05 $\pm$ 0.00	0.00 $\pm$ 0.00	0.07 $\pm$ 0.00	0.32 $\pm$ 0.00	0.00 $\pm$ 0.00	0.12 $\pm$ 0.00	0.20 $\pm$ 0.00	0.13 $\pm$ 0.00
FARE	Clean	29.76	17.37	9.35	28.06	22.19	15.08	23.34	23.43	26.10	33.39	19.50	22.51
	PGD	19.74 $\pm$ 0.08	16.95 $\pm$ 0.00	0.84 $\pm$ 0.01	6.20 $\pm$ 0.02	3.16 $\pm$ 0.02	0.17 $\pm$ 0.00	12.59 $\pm$ 0.03	4.60 $\pm$ 0.01	2.44 $\pm$ 0.00	0.12 $\pm$ 0.00	0.20 $\pm$ 0.00	6.09 $\pm$ 0.01
	C&W	19.02 $\pm$ 0.11	16.95 $\pm$ 0.00	0.79 $\pm$ 0.01	5.69 $\pm$ 0.02	3.11 $\pm$ 0.02	0.29 $\pm$ 0.00	12.65 $\pm$ 0.03	4.84 $\pm$ 0.01	2.37 $\pm$ 0.00	0.18 $\pm$ 0.00	0.37 $\pm$ 0.05	6.02 $\pm$ 0.01
	AA	18.33 $\pm$ 0.03	16.95 $\pm$ 0.00	0.80 $\pm$ 0.01	5.37 $\pm$ 0.02	2.66 $\pm$ 0.02	0.17 $\pm$ 0.00	12.57 $\pm$ 0.00	3.82 $\pm$ 0.02	1.86 $\pm$ 0.04	0.12 $\pm$ 0.00	0.17 $\pm$ 0.05	5.71 $\pm$ 0.00
+ TTC	Clean	30.40 $\pm$ 0.46	17.37 $\pm$ 0.00	15.56 $\pm$ 0.52	25.28 $\pm$ 0.28	19.30 $\pm$ 0.31	15.14 $\pm$ 0.69	24.49 $\pm$ 0.54	28.95 $\pm$ 0.50	25.89 $\pm$ 0.21	32.27 $\pm$ 0.46	19.37 $\pm$ 0.42	23.09 $\pm$ 0.22
	PGD	27.84 $\pm$ 0.26	17.37 $\pm$ 0.00	9.06 $\pm$ 0.16	28.51 $\pm$ 0.05	10.47 $\pm$ 0.76	7.17 $\pm$ 0.53	15.38 $\pm$ 0.21	18.47 $\pm$ 0.40	23.73 $\pm$ 0.33	17.13 $\pm$ 0.20	11.50 $\pm$ 0.43	16.97 $\pm$ 0.08
	C&W	27.47 $\pm$ 0.10	17.37 $\pm$ 0.00	8.87 $\pm$ 0.21	28.80 $\pm$ 0.39	10.11 $\pm$ 0.42	7.31 $\pm$ 0.28	14.98 $\pm$ 0.11	15.55 $\pm$ 0.20	22.93 $\pm$ 0.27	16.47 $\pm$ 0.49	10.53 $\pm$ 0.54	16.40 $\pm$ 0.03
	AA	31.16 $\pm$ 0.29	17.23 $\pm$ 0.00	16.95 $\pm$ 0.00	28.46 $\pm$ 0.26	17.49 $\pm$ 0.80	12.78 $\pm$ 0.69	22.19 $\pm$ 0.35	28.29 $\pm$ 0.16	29.16 $\pm$ 0.32	28.46 $\pm$ 0.33	18.83 $\pm$ 0.20	22.82 $\pm$ 0.03
+ T-RPT	Clean	30.81 $\pm$ 0.10	17.37 $\pm$ 0.00	15.37 $\pm$ 0.15	28.03 $\pm$ 0.15	20.95 $\pm$ 0.11	16.39 $\pm$ 0.00	22.38 $\pm$ 0.08	28.18 $\pm$ 0.17	21.60 $\pm$ 0.06	32.91 $\pm$ 0.00	21.17 $\pm$ 0.41	22.70 $\pm$ 0.02
	PGD	27.02 $\pm$ 0.09	17.37 $\pm$ 0.00	10.17 $\pm$ 0.19	23.34 $\pm$ 0.12	12.60 $\pm$ 0.24	10.81 $\pm$ 0.40	17.89 $\pm$ 0.11	12.90 $\pm$ 0.22	22.53 $\pm$ 0.37	18.14 $\pm$ 0.07	11.67 $\pm$ 0.31	16.77 $\pm$ 0.02
	C&W	27.28 $\pm$ 0.12	17.37 $\pm$ 0.00	10.11 $\pm$ 0.12	23.15 $\pm$ 0.05	13.12 $\pm$ 0.07	11.69 $\pm$ 0.10	17.67 $\pm$ 0.03	12.63 $\pm$ 0.14	22.03 $\pm$ 0.21	21.01 $\pm$ 0.26	11.47 $\pm$ 0.49	17.05 $\pm$ 0.05
	AA	28.96 $\pm$ 0.22	17.37 $\pm$ 0.00	13.33 $\pm$ 0.05	25.42 $\pm$ 0.15	15.96 $\pm$ 0.29	13.53 $\pm$ 0.00	20.26 $\pm$ 0.28	14.38 $\pm$ 0.17	23.58 $\pm$ 0.23	25.22 $\pm$ 0.28	16.07 $\pm$ 0.25	19.46 $\pm$ 0.07
+ TAME (Ours)	Clean	36.48 $\pm$ 0.00	17.80 $\pm$ 0.35	12.55 $\pm$ 0.27	24.99 $\pm$ 0.19	17.09 $\pm$ 0.14	12.68 $\pm$ 0.17	25.26 $\pm$ 0.33	27.52 $\pm$ 0.04	21.35 $\pm$ 0.48	30.29 $\pm$ 0.47	27.27 $\pm$ 0.47	32.34 $\pm$ 0.10
	PGD	45.06 $\pm$ 1.09	21.19 $\pm$ 1.38	11.97 $\pm$ 0.18	33.16 $\pm$ 0.25	24.06 $\pm$ 0.12	20.39 $\pm$ 0.20	35.55 $\pm$ 0.08	43.44 $\pm$ 0.66	31.81 $\pm$ 0.42	57.05 $\pm$ 1.03	33.83 $\pm$ 0.68	32.52 $\pm$ 0.26
	C&W	44.73 $\pm$ 0.72	20.48 $\pm$ 0.40	11.90 $\pm$ 0.29	33.63 $\pm$ 1.03	22.46 $\pm$ 0.20	18.53 $\pm$ 0.31	31.87 $\pm$ 0.64	42.36 $\pm$ 0.45	31.04 $\pm$ 0.49	52.46 $\pm$ 0.12	29.90 $\pm$ 0.10	30.85 $\pm$ 0.33
	AA	41.04 $\pm$ 0.96	17.80 $\pm$ 0.35	13.28 $\pm$ 0.14	29.81 $\pm$ 0.05	22.46 $\pm$ 0.12	11.21 $\pm$ 0.69	29.17 $\pm$ 0.15	29.16 $\pm$ 0.64	34.49 $\pm$ 0.46	36.29 $\pm$ 0.10	29.70 $\pm$ 1.10	25.98 $\pm$ 0.11
PMG	Clean	27.84	17.80	26.97	24.88	16.91	15.08	22.54	19.72	21.06	36.35	23.30	22.95
	PGD	23.12 $\pm$ 0.00	17.37 $\pm$ 0.00	15.11 $\pm$ 0.01	14.54 $\pm$ 0.07	7.43 $\pm$ 0.04	15.18 $\pm$ 0.00	14.39 $\pm$ 0.06	8.65 $\pm$ 0.07	4.37 $\pm$ 0.06	8.30 $\pm$ 0.00	12.27 $\pm$ 0.02	12.71 $\pm$ 0.02
	C&W	22.97 $\pm$ 0.00	17.37 $\pm$ 0.00	14.12 $\pm$ 0.05	13.71 $\pm$ 0.07	5.17 $\pm$ 0.00	5.17 $\pm$ 0.00	15.16 $\pm$ 0.00	12.71 $\pm$ 0.00	7.49 $\pm$ 0.00	6.18 $\pm$ 0.00	8.60 $\pm$ 0.00	11.71 $\pm$ 0.00
	AA	23.04 $\pm$ 0.02	17.37 $\pm$ 0.00	14.37 $\pm$ 0.04	14.00 $\pm$ 0.03	5.67 $\pm$ 0.08	5.58 $\pm$ 0.07	14.58 $\pm$ 0.03	12.93 $\pm$ 0.05	7.78 $\pm$ 0.10	4.91 $\pm$ 0.03	7.90 $\pm$ 0.00	11.65 $\pm$ 0.02
+ TTC	Clean	29.14 $\pm$ 0.52	16.84 $\pm$ 0.60	24.45 $\pm$ 0.06	24.06 $\pm$ 0.14	16.40 $\pm$ 0.36	14.83 $\pm$ 0.18	22.03 $\pm$ 0.23	20.70 $\pm$ 0.15	21.58 $\pm$ 0.38	32.43 $\pm$ 1.44	23.00 $\pm$ 0.98	22.48 $\pm$ 0.50
	PGD	24.85 $\pm$ 0.54	17.37 $\pm$ 0.00	18.90 $\pm$ 0.04	23.74 $\pm$ 0.74	7.93 $\pm$ 0.16	9.00 $\pm$ 0.12	17.44 $\pm$ 0.33	15.97 $\pm$ 0.07	13.85 $\pm$ 0.32	17.71 $\pm$ 0.81	12.47 $\pm$ 0.12	16.29 $\pm$ 0.07
	C&W	24.71 $\pm$ 0.39	17.23 $\pm$ 0.00	18.89 $\pm$ 0.34	23.13 $\pm$ 0.17	7.08 $\pm$ 0.09	8.97 $\pm$ 0.31	16.93 $\pm$ 0.45	15.22 $\pm$ 0.20	12.72 $\pm$ 0.51	16.85 $\pm$ 0.30	12.87 $\pm$ 0.53	15.96 $\pm$ 0.07
	AA	25.17 $\pm$ 0.36	17.80 $\pm$ 0.00	22.81 $\pm$ 0.16	24.72 $\pm$ 0.11	11.74 $\pm$ 0.08	11.36 $\pm$ 0.48	20.37 $\pm$ 0.13	18.47 $\pm$ 0.13	17.59 $\pm$ 0.49	28.02 $\pm$ 0.22	13.93 $\pm$ 0.52	20.12 $\pm$ 0.23
+ T-RPT	Clean	27.53 $\pm$ 0.03	17.37 $\pm$ 0.00	27.38 $\pm$ 0.04	24.03 $\pm$ 0.04	13.02 $\pm$ 0.04	12.95 $\pm$ 0.17	22.34 $\pm$ 0.24	17.64 $\pm$ 0.06	13.85 $\pm$ 0.43	31.20 $\pm$ 0.20	18.50 $\pm$ 0.22	20.53 $\pm$ 0.03
	PGD	25.92 $\pm$ 0.17	17.37 $\pm$ 0.00	25.08 $\pm$ 0.12	22.31 $\pm$ 0.10	10.49 $\pm$ 0.06	12.03 $\pm$ 0.22	20.66 $\pm$ 0.03	16.13 $\pm$ 0.08	11.07 $\pm$ 0.04	21.14 $\pm$ 0.30	12.57 $\pm$ 0.05	17.70 $\pm$ 0.01
	C&W	26.15 $\pm$ 0.05	17.37 $\pm$ 0.00	24.83 $\pm$ 0.10	21.97 $\pm$ 0.16	9.66 $\pm$ 0.00	12.03 $\pm$ 0.10	19.88 $\pm$ 0.20	15.48 $\pm$ 0.03	10.49 $\pm$ 0.49	21.35 $\pm$ 0.33	13.37 $\pm$ 0.19	17.51 $\pm$ 0.04
	AA	26.08 $\pm$ 0.14	17.37 $\pm$ 0.00	25.81 $\pm$ 0.01	23.47 $\pm$ 0.08	11.14 $\pm$ 0.09	11.24 $\pm$ 0.09	21.03 $\pm$ 0.11	16.58 $\pm$ 0.04	12.09 $\pm$ 0.18	26.69 $\pm$ 0.13	16.00 $\pm$ 0.37	19.07 $\pm$ 0.04
+ TAME (Ours)	Clean	30.11 $\pm$ 0.41	23.84 $\pm$ 0.46	23.99 $\pm$ 0.41	12.24 $\pm$ 0.40	15.33 $\pm$ 0.20	21.16 $\pm$ 0.68	17.49 $\pm$ 0.33	21.06 $\pm$ 0.26	26.85 $\pm$ 0.91	21.80 $\pm$ 0.01	21.02 $\pm$ 0.17	21.02 $\pm$ 0.04
	PGD	35.74 $\pm$ 1.03	17.51 $\pm$ 0.20	23.26 $\pm$ 0.49	21.98 $\pm$ 0.12	9.90 $\pm$ 0.12	14.28 $\pm$ 0.69	20.81 $\pm$ 0.34	16.80 $\pm$ 0.15	17.61 $\pm$ 0.29	38.91 $\pm$ 0.17	18.63 $\pm$ 0.17	24.10 $\pm$ 0.14
	C&W	35.53 $\pm$ 1.03	17.66 $\pm$ 0.20	23.47 $\pm$ 0.20	21.75 $\pm$ 0.40	8.15 $\pm$ 0.12	14.31 $\pm$ 0.35	20.17 $\pm$ 0.49	17.44 $\pm$ 0.27	17.60 $\pm$ 0.46	39.29 $\pm$ 0.84	17.67 $\pm$ 0.87	21.09 $\pm$ 0.22
	AA	30.07 $\pm$ 0.34	17.37 $\pm$ 0.00	23.58 $\pm$ 0.27	22.37 $\pm$ 0.38	9.83 $\pm$ 0.14	14.67 $\pm$ 0.27	20.17 $\pm$ 0.17	16.09 $\pm$ 0.09	18.17 $\pm$ 0.15	36.41 $\pm$ 0.10	20.13 $\pm$ 0.02	20.80 $\pm$ 0.04
TeCoA	Clean	27.78	17.37	15.98	21.91	15.46	15.33	22.74	18.80	25.39	37.68	29.70	22.56
	PGD	27.26 $\pm$ 0.00	17.37 $\pm$ 0.00	4.25 $\pm$ 0.01	13.79 $\pm$ 0.03	7.50 $\pm$ 0.06	3.79 $\pm$ 0.07	13.83 $\pm$ 0.00	12.95 $\pm$ 0.03	10.88 $\pm$ 0.07	9.62 $\pm$ 0.03	10.40 $\pm$ 0.00	11.96 $\pm$ 0.00
	C&W	27.26 $\pm$ 0.00	17.37 $\pm$ 0.00	1.79 $\pm$ 0.00	13.33 $\pm$ 0.03	6.03 $\pm$ 0.00	3.30 $\pm$ 0.07	13.76 $\pm$ 0.00	12.69 $\pm$ 0.04	10.54 $\pm$ 0.04	9.64 $\pm$ 0.07	9.87 $\pm$ 0.05	11.42 $\pm$ 0.00
	AA	27.26 $\pm$ 0.00	17.37 $\pm$ 0.00	1.94 $\pm$ 0.00	13.48 $\pm$ 0.09	5.81 $\pm$ 0.10	3.29 $\pm$ 0.07	13.47 $\pm$ 0.06	12.79 $\pm$ 0.02	10.88 $\pm$ 0.07	9.70 $\pm$ 0.03	10.33 $\pm$ 0.12	11.49 $\pm$ 0.00
+ TTC	Clean	27.84 $\pm$ 0.05	17.37 $\pm$ 0.00	18.51 $\pm$ 0.29	22.00 $\pm$ 0.07	14.50 $\pm$ 0.39	16.61 $\pm$ 0.37	22.27 $\pm$ 0.05	18.90 $\pm$ 0.17	25.05 $\pm$ 0.19	34.20 $\pm$ 0.16	27.20 $\pm$ 0.45	22.44 $\pm$ 0.07
	PGD	27.70 $\pm$ 0.00	17.37 $\pm$ 0.00	11.13 $\pm$ 0.18	18.67 $\pm$ 0.28	7.46 $\pm$ 0.45	8.69 $\pm$ 0.14	16.13 $\pm$ 0.37	15.09 $\pm$ 0.04	17.14 $\pm$ 0.36	22.71 $\pm$ 0.30	15.43 $\pm$ 0.43	16.14 $\pm$ 0.11
	C&W	27.64 $\pm$ 0.07	17.37 $\pm$ 0.00	11.22 $\pm$ 0.19	18.58 $\pm$ 0.18	6.66 $\pm$ 0.27	8.55 $\pm$ 0.18	15.65 $\pm$ 0.31	15.07 $\pm$ 0.08	17.30 $\pm$ 0.23	23.01 $\pm$ 0.08	15.13 $\pm$ 0.29	16.02 $\pm$ 0.13
	AA	27.72 $\pm$ 0.08	17.37 $\pm$ 0.00	16.96 $\pm$ 0.21	20.44 $\pm$ 0.20	11.29 $\pm$ 0.11	13.78 $\pm$ 0.04	19.77 $\pm$ 0.39	15.67 $\pm$ 0.03	22.16 $\pm$ 0.07	32.31 $\pm$ 0.10	21.97 $\pm$ 0.45	20.00 $\pm$ 0.23
+ T-RPT	Clean	27.00 $\pm$ 0.03	22.63 $\pm$ 0.03	21.70 $\pm$ 0.06	12.74 $\pm$ 0.06	14.44 $\pm$ 0.14	12.21 $\pm$ 0.13	18.83 $\pm$ 0.09	28.16 $\pm$ 0.40	31.80 $\pm$ 0.03	32.67 $\pm$ 0.46	29.70 $\pm$ 0.00	22.86 $\pm$ 0.00
	PGD	27.45 $\pm$ 0.45	17.37 $\pm$ 0.00	18.40 $\pm$ 0.09	17.59 $\pm$ 0.09	10.32 $\pm$ 0.04	11.97 $\pm$ 0.14	13.83 $\pm$ 0.11	17.67 $\pm$ 0.02	22.16 $\pm$ 0.07	22.10 $\pm$ 0.22	27.37 $\pm$ 0.7	