

A Study of Artifacts on Melanoma Classification under Diffusion-Based Perturbations

Qixuan Jin

Massachusetts Institute of Technology, United States

QIXUANJ@MIT.EDU

Abstract

In melanoma classification, deep learning models have been shown to rely on non-medical artifacts (e.g., surgical markings) rather than clinically relevant features (e.g., lesion asymmetry), compromising their generalizability. In this work, we investigate the impact of artifacts on melanoma classification under two settings: (1) input disruptions, such as bounding boxes and frequency-based filtering, which isolate artifacts by region or frequency, and (2) a novel diffusion-based perturbation method that selectively introduces isolated artifacts into images, generating controlled pairs for direct comparison. We systematically analyze artifact biases in three benchmark datasets: ISIC 2018, HAM10000, and PH2. Our findings reveal that whole-image training outperforms lesion-only or background-only approaches, low-frequency features are essential for melanoma prediction, and classifiers are more sensitive to perturbations for the artifacts of ink markings, rulers, and patches. These results emphasize the need for systematic artifact assessment and provide insights for improving the robustness of melanoma classification models.

Data and Code Availability This paper uses the ISIC 2018 dataset (Codella et al., 2019) and the HAM10000 dataset (Tschandl et al., 2018), which are both available on the [ISIC repository](#). The paper uses the PH2 dataset (Mendonça et al., 2013), which is available at their [website](#). Our GitHub repository is at https://github.com/QixuanJin99/dermoscopic_artifacts/tree/main.

Institutional Review Board (IRB) Our research does not require IRB approval.

1. Introduction

Dermoscopy is a widely used imaging technique for melanoma diagnosis, enhancing the visualization of

subsurface skin structures (Sonthalia et al., 2019). Clinicians rely on diagnostic frameworks such as the ABCD rule (Nachbar et al., 1994) and the Seven-Point Checklist (Argenziano et al., 1998; Di Leo et al., 2010) to systematically assess lesion characteristics. While machine learning models have achieved comparable performance to dermatologists in melanoma classification (Haenssle et al., 2018), these models often rely on non-medical artifacts, such as surgical markings and rulers, rather than clinically relevant features (Winkler et al., 2019; Bevan and Atapour-Abarghouei, 2021).

Spurious correlations can degrade model generalization in real world settings where the spurious feature is absent in new data where artifacts are absent (Boland et al., 2024; Singla and Feizi, 2021; Minderer et al., 2020). In medical settings, ensuring generalizability is critical, as systematic biases across hospitals and populations can significantly impact real-world performance (Yellido, 2020; Salahuddin et al., 2022). Therefore, determining the impact of artifact reliance is essential for developing robust melanoma classification models. Dermatological datasets provide an ideal setting for studying artifact-driven biases, as diagnostic decisions rely on medical images and non-medical features can mislead deep learning classifiers.

In a prior work, Bissoto et al. (2020) manually annotated seven artifacts commonly present in dermoscopic images (Figure 1), finding that melanoma classifiers rely on both clinically relevant lesion features and non-medical artifacts for prediction. We make use of this setting to investigate the impact of artifacts in additional dermatological datasets, and new perturbation settings.

First, we assess whether background artifacts are predictive of melanoma, and how well they generalize under specific dataset disruptions, such as bounding box constraints and low- and high-pass filtering. We evaluate performance of ResNet-50 (He et al., 2016) models on melanoma classification using the ISIC

80 2018 ([Codella et al., 2019](#)), PH2 ([Mendonça et al., 2013](#)), and HAM10000 ([Tschandl et al., 2018](#)) benchmark datasets. Models trained on ISIC and PH2 are each evaluated on the remaining datasets to measure transfer performance. We find that non-medical artifacts do not inherently degrade generalization in this setting, as some artifacts transfer between datasets. For instance, artifacts near the outer edges of images retain their predictive value when transferred from ISIC to PH2. Our findings are in strong contrast to other work in the spurious correlation literature, and indicate that the effect of non-medical features on transfer performance should be assessed by model development and deployment teams prior to widespread use.

95 Second, we isolate the impact of specific artifacts on classification performance by finetuning Stable Diffusion 1.5 ([Rombach et al., 2022](#)) with Dreambooth ([Ruiz et al., 2023](#)) and LoRA ([Hu et al., 2021](#)) to inpaint isolated artifacts onto existing dermoscopic images. Our diffusion models are capable of generating the common non-medical artifacts of dark corners, gel bubbles, ink markings, colored patches, and rulers. To the best of our knowledge, no prior work has utilized diffusion models specifically for adding non-medical artifacts onto dermoscopic images. We evaluate the transfer performance of classifiers trained on PH2 and ISIC when applied to PH2 images perturbed by these artifacts. Our findings reveal that classifiers are especially sensitive to diffusion-based perturbations involving ink markings, rulers, and patches. Given that certain artifacts have a greater impact on transfer performance, systematic artifact assessment is essential for developing targeted mitigation strategies before real-world deployment.

116 We present the following contributions:

- 117 1. Comprehensively evaluated melanoma classification performance under different input disruptions and transfer across three public benchmark datasets.
- 122 2. Expanded upon manual annotations of common non-medical dermoscopic artifacts with the labeling of the PH2 dataset. We will release this dataset for public research.
- 125 3. Developed diffusion models capable of selectively augmenting dermoscopic images with specific artifacts (e.g. dark corners, gel bubbles, ink, patches, and rulers), providing a flexible frame-

work for future researchers to dynamically augment their datasets.

2. Related Work

2.1. Machine Learning for Melanoma Classification

Machine learning methods for melanoma classification on dermoscopic images generally fall into two categories: 1) Enhancing existing medical algorithms 2) image-based classification with vision models. [Kawahara et al. \(2018\)](#) uses a multitask deep convolutional neural network to predict the seven-point checklist criteria and thus perform the skin lesion diagnosis. [Almaraz-Damian et al. \(2020\)](#) fuses hand-crafted ABCD features with vision model embeddings for improved performance. Despite these methods' medical grounding, dermatologist-labeled features like border irregularity and asymmetry have inherent degrees of subjectivity and high-quality ground truth labels are costly to obtain at scale.

In parallel, a great variety of convolutional neural networks (CNN) and vision transformers (ViT) models have been adapted from natural image settings to predict melanoma in dermoscopic images with high accuracy ([Arshed et al., 2023; Budhiman et al., 2019; Gamage et al., 2024](#)). These deep image classifiers have demonstrated comparable or improved performance over the medical-augmented algorithms and clinicians ([Haenssle et al., 2018](#)). Specifically, the CNNs demonstrated higher specificity with comparable sensitivity to dermatologists, with slightly better performance overall.

In this work, we focus on melanoma classification with vision-based models, and explore how these models can utilize non-dermoscopic features outside the domain of medical algorithms to predict melanoma with better accuracy.

2.2. Artifacts in Dermoscopic Images

Non-medical artifacts have long been investigated for spurious association with disease prognosis in dermatology. Past work illustrates how a market-approved CNN classifier spuriously associates large dark corner artifacts with malignancy, leading to decreased performance for images with this artifact ([Sies et al., 2021](#)). Surgical ink markings have also been found to be positively associated with melanoma predictions ([Winkler et al., 2019](#)). We take inspiration

from the prior work [Bissoto et al. \(2020\)](#), which defined a set of non-medical artifacts present in the ISIC dataset (dark corners, hair, gel border, gel bubbles, ink markings, ruler, and patches) and demonstrated how model-based debiasing methods cannot fully remove dependency on artifacts.

2.3. Denoising Diffusion Probabilistic Models

Denoising diffusion probabilistic models (diffusion models) demonstrate state-of-the-art performance for text-to-image generation ([Dhariwal and Nichol, 2021](#)). In particular, Stable Diffusion ([Rombach et al., 2022](#)) has been widely used for art creation ([Liao et al., 2022](#)) and image editing ([Yang et al., 2023; Brooks et al., 2023](#)). Methods such as Textual Inversion ([Gal et al., 2022](#)) and Dreambooth ([Ruiz et al., 2023](#)) allow pretrained diffusion models to learn specific, new concepts given a few sample images. Low-Rank Adaptation (LoRA) is a parameter-efficient fine-tuning method that modifies a set of low-rank matrices with frozen pretrained model weights to reduce computational intensity ([Hu et al., 2021](#)). In this work, we use a Stable Diffusion backbone and fine-tune the model with Dreambooth and LoRA to learn a specific artifact’s visual features with a custom token.

Diffusion models have also been adapted for dermatology applications. Derm-T2IM is a Stable Diffusion model fine-tuned with DreamBooth and LoRA to generate images of benign and malignant skin lesions ([Farooq et al., 2024](#)). Additionally, diffusion models have been used to augment dermatology images, including those not captured with a dermatoscope, to improve representation for underrepresented groups ([Sagers et al., 2022](#)). However, to the best of our knowledge, there are no publicly available Stable Diffusion models designed to generate non-medical artifacts in dermoscopic images. This work aims to bridge that gap by developing a tailored diffusion model for artifact generation in dermoscopy.

3. Method

3.1. Datasets and Preprocessing

International Skin Imaging Collaboration (ISIC) is a collection of dermoscopic images used for classification, segmentation, and detection tasks ([Codella et al., 2019](#)). We specifically use the training split ($n=2,594$) from the ISIC 2018 Task 1-2, as only this split is labeled with artifacts by [Bissoto et al.](#)

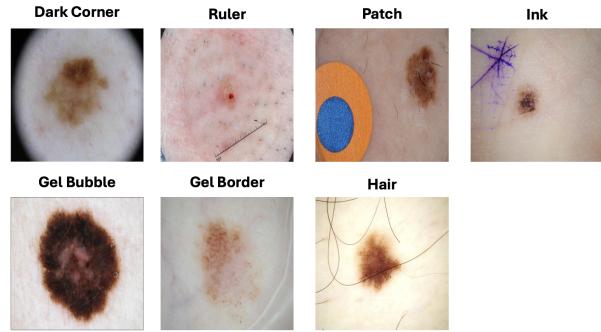


Figure 1: The types of dermoscopic artifacts.

([2020](#)). The HAM10000 dataset is a dataset of 10,015 dermoscopic images collected from diverse sources ([Tschanl et al., 2018](#)). The PH2 dataset is a collection of 200 dermoscopic images collected at Dermatology Service of Hospital Pedro Hispano, Portugal with human-labeled dermoscopic attributes ([Mendonça et al., 2013](#)).

Since we wish to investigate the robustness of image classifiers under different input image augmentations, we preprocess all datasets under 12 settings (“dataset modes”). The original image is mode “whole”. With lesion segmentations, we create a dataset with only the skin lesion (“lesion”) or only the background (“background”). We follow prior work ([Bissoto et al., 2020](#)) and compute the outer boundaries of the segmentation and convert the mask into a bounding box that removes the lesion boundary information. For “bbox70” and “bbox90”, the bounding box is expanded until approximately 70% or 90% of the image is masked out. We additionally evaluate high/low frequency features with a high/low-pass Gaussian filter with $\sigma = 1$. We use the filters on “whole”, “lesion”, and “background” images to obtain the remaining six dataset modes.

For melanoma classification, we use the benign and malignant lesion labels in the ISIC 2018 dataset. Based on the ISIC definition, we divide the HAM10000 into benign (‘Benign Keratosis-like Lesions’, ‘Dermatofibroma’, ‘Melanocytic Nevi’, ‘Vascular Lesions’) and malignant (‘Actinic Keratoses’, ‘Basal Cell Carcinoma’, ‘Melanoma’).

3.2. Classification Models

We frame the problem of melanoma classification as a binary classification task. Prior work ([Gazioğlu and Kamaşak, 2021](#)) shows the robustness of both

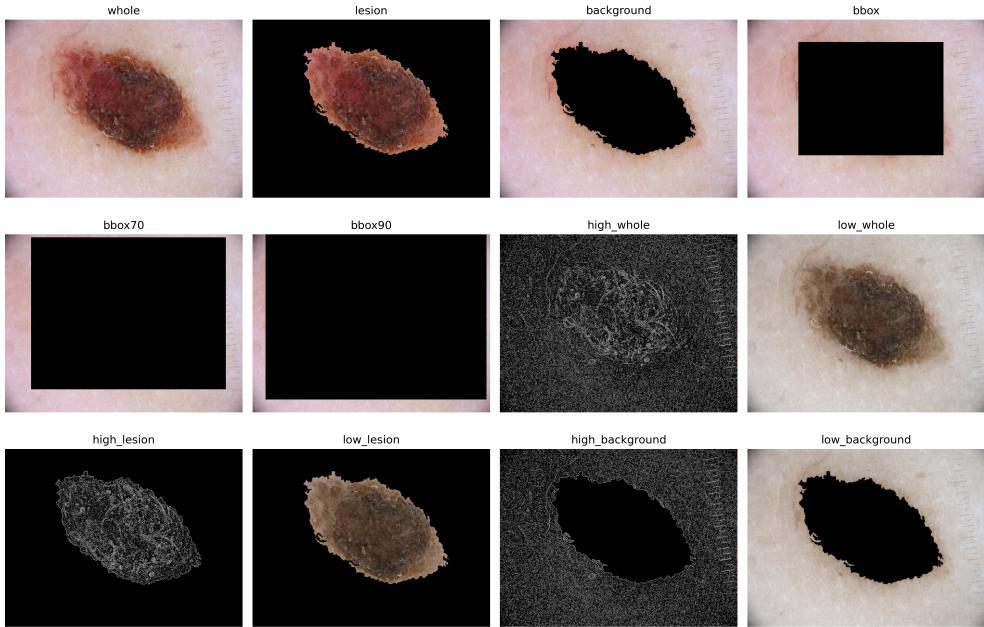


Figure 2: Visualization of the dataset modes we use to investigate the classifier’s robustness under different data preprocessing settings.

257 ResNet50 (He et al., 2016) and DenseNet121 (Huang
 258 et al., 2017) under certain noise settings. We focus on
 259 transferring ImageNet (Deng et al., 2009) pretrained
 260 ResNet50 classifier models to the dermoscopic do-
 261 main. We train with an Adam optimizer at a learning
 262 rate of $1e-4$ for 10 epochs and a batch size of 32. All
 263 images are resized to 224×224 and normalized to the
 264 ImageNet distribution. We finetune from pretrained
 265 ImageNet weights.

266 3.3. Diffusion Models for Artifact Generation 267 and Classifier Explainability

268 To investigate the trained classifier’s robustness to
 269 perturbations of specific artifacts, we propose an ex-
 270 plainability procedure that breaks the association be-
 271 tween artifact and prediction label through the in-
 272 discriminately addition of the artifact to all images.
 273 To automate the addition of artifacts in a realis-
 274 tic manner, we train diffusion models for each ar-
 275 tifact. Specifically, Dreambooth (Ruiz et al., 2023)
 276 is a method that can fine-tune text-to-image diffusion
 277 models to associate specific visual concepts with the
 278 provided text tokens given a set of example images.
 279 To prevent the association of existing visual concepts
 280 with our artifact concepts, we pick our Dreambooth

tokens from a list of rare single tokens (2kpr, 2024).
 281 We define the tokens of {patches: olis, dark corner:
 282 lun, ruler: dits, ink: httr, gel bubble: sown}. Since
 283 we are learning single visual concepts, we can speed
 284 up the training process with Low-Rank Adaptation
 285 (LoRA), in which we tune low-rank matrices instead
 286 of the full diffusion model (Hu et al., 2021).

287 For the base model, we use Stable Diffusion 1.5
 288 (Rombach et al., 2022) with the HuggingFace imple-
 289 mentation (“runwayml/stable-diffusion-v1-5”) and
 290 code from the diffusers library (von Platen et al.,
 291 2022). We finetune both the UNet (Ronneberger
 292 et al., 2015) and text encoder with LoRA matrices
 293 of rank 64 and $\alpha = 32$. We finetune for 3 to 5 epochs
 294 with a constant learning rate of $1e-4$, a batch size of
 295 2, and an image resolution of 512. For text condition-
 296 ing, we use the instance prompt template: “a dermo-
 297 scopic image of token class”. The token is the rare
 298 artifact token and the class is either “benign” or “ma-
 299 lignant” according to the melanoma label. To prevent
 300 large deviations from the melanoma class, we train
 301 with prior preservation with a prior loss weight of 0.3
 302 and 200 generated class images. The class prompt
 303 follows the same template but excludes the artifact
 304 token. For compute, we train using seed=0, mixed

306 precision (fp16) on a single NVIDIA A100 GPU, with
 307 training completing in approximately one hour.

308 To generate artifacts for the perturbation method,
 309 we load the fine-tuned LoRA adapter into the in-
 310 painting pipeline. We generate images at a guidance
 311 scale of 10 and tune the strength parameter manu-
 312 ally through visualizations. We decide on the follow-
 313 ing strength parameters {patches : 0.85, dark corner :
 314 0.75, ruler : 0.65, ink : 0.7, gel bubble : 0.6}. In our
 315 perturbation analysis, we augment the PH2 dataset.
 316 We generate 5 augmented images per original image,
 317 and selected for the image that is most predictive
 318 of the artifact given artifact classifiers (ResNet50)
 319 trained on the ISIC ground truth artifact labels.

320 3.4. Evaluation

321 We use 5-fold cross-validation on each dataset to train
 322 and evaluate the classifier using non-overlapping train
 323 and test splits. We use area under the receiver oper-
 324 ating curve (AUROC) (Hanley and McNeil, 1982) as
 325 the primary performance metric, with error bounds
 326 computed as the standard deviation across the 5
 327 folds. We define evaluation results as “source” when
 328 the test split is drawn from the same dataset used for
 329 training, while “target” refers to evaluations on an
 330 out-of-domain test split.

331 4. Experiments

332 4.1. Melanoma Classification and Transfer 333 Performance under Dataset Disruptions

334 We evaluate the performance of a classifier for the
 335 classification of benign and malignant melanoma im-
 336 ages across two dimensions: (1) the effect of dataset
 337 perturbation by training the classifier on the train
 338 set and evaluating it on the test set from the same
 339 source dataset, and (2) the generalization ability of
 340 the perturbed classifier when evaluated on test sets
 341 from out-of-domain target datasets. Given that PH2-
 342 trained classifiers exhibit poor generalization across
 343 all settings, we primarily focus on ISIC-trained clas-
 344 sifiers when analyzing transfer performance.

345 First, we note that the classifier trained on whole
 346 images consistently outperforms classifiers trained on
 347 either lesion-only or background-only perturbed im-
 348 ages, both within the source dataset and in transfer
 349 to target datasets (Table 1 and Table 2). When eval-
 350 uated on the source dataset, the whole-image classi-
 351 fier achieves the highest AUROC (ISIC: 0.792, PH2:
 352 0.975), outperforming the lesion-only (ISIC: 0.746,

PH2: 0.952) and background-only classifiers (ISIC:
 353 0.752, PH2: 0.884). This trend extends to cross-
 354 dataset generalization. For the ISIC-trained models,
 355 the whole-image classifier exhibits superior transfer
 356 performance (HAM10000: 0.721, PH2: 0.858) com-
 357 pared to the lesion-only (HAM10000: 0.719, PH2:
 358 0.818) and background-only classifiers (HAM10000:
 359 0.714, PH2: 0.777). These results suggest that lesion
 360 and background features are both predictive, and the
 361 combination of the two improves model performance.
 362

363 Despite the clinical intuition that lesion features
 364 should be the most predictive (Williams et al.,
 365 2021; Longo et al., 2023; Tschandl et al., 2019),
 366 background-only classifiers remain surprisingly effec-
 367 tive. In the well-known ABCD rule (Nachbar et al.,
 368 1994), features such as asymmetry and irregular bor-
 369 ders are defined by the lesion *boundary*, a feature that
 370 is retained in our background-only perturbation. To
 371 evaluate the significance of this boundary, we use a
 372 bounding box to remove the lesion shape while pre-
 373 serving the surrounding background. In ISIC, per-
 374 formance drops from 0.752 (background) to 0.717
 375 (lesion) when the lesion boundary is removed. As
 376 the bounding box expands to cover 70% and 90% of
 377 the image, performance further declines to 0.681 and
 378 0.628, respectively, indicating that the model relies
 379 on more than just the outer edges of the image. A
 380 similar performance drop is observed in PH2, despite
 381 its small sample size ($n = 200$ samples) and larger
 382 error bounds (Table 2). These findings highlight the
 383 importance of lesion boundary features in melanoma
 384 classification and also suggests that the surrounding
 385 background contains predictive information.

386 Importantly, we observe that bbox classifiers
 387 trained on the ISIC dataset maintain strong per-
 388 formance when transferred to target datasets (ISIC
 389 source: 0.717, HAM10000 target: 0.687, PH2 target:
 390 0.821). The ability of bbox-based models to general-
 391 ize suggests that certain background features associ-
 392 ated with melanoma are consistently preserved across
 393 datasets. This finding suggests that non-medical arti-
 394 facts do not inherently degrade generalization, as
 395 some artifacts may be transferable between datasets.
 396 The impact of training models with non-medical arti-
 397 facts on transfer performance should therefore be
 398 assessed on a case-by-case basis.

399 Beyond localizing predictive regions, we examine
 400 whether the classifier relies on features at specific
 401 frequency ranges. Using high-pass and low-pass fil-
 402 ters, we extract high- and low-frequency signals to as-
 403 sess their contribution to classification performance.

Domain AUROC	Source ISIC	Target HAM10000	Target PH2
whole	0.792 ± 0.024	0.721 ± 0.061	0.858 ± 0.045
lesion	0.746 ± 0.040	0.719 ± 0.024	0.818 ± 0.099
background	0.752 ± 0.017	0.714 ± 0.038	0.777 ± 0.092
bbox	0.717 ± 0.024	0.687 ± 0.040	0.821 ± 0.051
bbox70	0.681 ± 0.026	0.686 ± 0.047	0.805 ± 0.013
bbox90	0.628 ± 0.016	0.616 ± 0.022	0.679 ± 0.064
whole (high freq)	0.530 ± 0.020	0.307 ± 0.027	0.442 ± 0.124
lesion (high freq)	0.590 ± 0.042	0.575 ± 0.077	0.502 ± 0.204
background (high freq)	0.523 ± 0.037	0.500 ± 0.127	0.682 ± 0.186
whole (low freq)	0.765 ± 0.017	0.741 ± 0.030	0.748 ± 0.052
lesion (low freq)	0.717 ± 0.028	0.685 ± 0.029	0.856 ± 0.086
background (low freq)	0.704 ± 0.035	0.624 ± 0.042	0.817 ± 0.078

Table 1: Performance in AUROC with standard deviations for the melanoma classifier trained on the ISIC dataset and evaluated on the ISIC, HAM10000, and PH2 datasets.

404 In ISIC, the classifier struggles to learn from high-
 405 frequency signals but maintains strong performance
 406 with only low-frequency information, showing an AU-
 407 ROC drop of less than 0.05 compared to the original
 408 image. Additionally, low-frequency models generalize
 409 well across datasets (ISIC source: 0.765, HAM10000
 410 target: 0.741, PH2 target: 0.748). In PH2, both
 411 high- and low-frequency models perform well, likely
 412 due to overfitting on the small training set. These
 413 results demonstrate that the model primarily relies
 414 on low-frequency signals for melanoma classification,
 415 suggesting that coarse patterns contribute more to
 416 prediction than fine-grained details.

417 **4.2. Artifact Prevalences and Correlations**

418 Given that the background contributes to melanoma
 419 classification, we examine whether the presence of
 420 specific non-medical artifacts may explain this pre-
 421 dictive signal.

422 Table 3 shows that artifact prevalence varies across
 423 datasets. Dark corners are significantly more com-
 424 mon in PH2 (94%) than in ISIC (37%), while rulers
 425 and patches are entirely absent from PH2. To
 426 assess potential associations between artifacts and
 427 melanoma, we compute the Pearson correlation be-
 428 tween individual artifacts and the melanoma label
 429 in ISIC and PH2 (Figure 3). We find that no sin-
 430 gle artifact exhibits a strong direct correlation with
 431 melanoma. However, several artifacts are correlated
 432 with each other. In ISIC, rulers frequently co-occur
 433 with dark corners ($\text{corr} = 0.31$) and gel borders (corr

= 0.34), while ink markings are often found alongside
 434 gel borders ($\text{corr} = 0.42$) and rulers ($\text{corr} = 0.40$).
 435 Since PH2 lacks rulers and contains very few ink-
 436 marked images, these inter-artifact correlations ob-
 437 served in ISIC do not transfer to PH2. These findings
 438 suggest that while no individual artifact strongly pre-
 439 dictes melanoma, the combined presence of artifacts
 440 may introduce patterns that influence model predic-
 441 tions.

443 **4.3. Training Diffusion Models for Artifact 444 Generation**

445 We first fine-tune Stable Diffusion 1.5 on the dermo-
 446 scopic image domain. Then, using DreamBooth and
 447 LoRA, we further fine-tune this base dermoscopis
 448 diffusion model to generate specific artifact types, lever-
 449 aging ground truth labels from the ISIC dataset.

450 Due to challenges in disentangling artifact concepts
 451 and achieving stable training convergence, we exclude
 452 results for the artifacts “hair” and “gel border”. For
 453 the remaining five artifacts, we use the trained dif-
 454 fusion models to augment the PH2 dataset by intro-
 455 ducing these artifacts through masked inpainting. We
 456 carefully tune the strength hyperparameter to achieve
 457 a balance between preserving the original image and
 458 effectively integrating the artifact (Figure 4).

Domain AUROC	Source PH2	Target ISIC	Target HAM10000
whole	0.975 ± 0.012	0.635 ± 0.050	0.638 ± 0.064
lesion	0.952 ± 0.047	0.623 ± 0.028	0.624 ± 0.034
background	0.884 ± 0.088	0.557 ± 0.080	0.505 ± 0.067
bbox	0.880 ± 0.080	0.565 ± 0.036	0.514 ± 0.041
bbox70	0.845 ± 0.046	0.566 ± 0.049	0.570 ± 0.016
bbox90	0.865 ± 0.081	0.544 ± 0.026	0.433 ± 0.030
whole (high freq)	0.903 ± 0.070	0.509 ± 0.023	0.609 ± 0.025
lesion (high freq)	0.877 ± 0.064	0.664 ± 0.055	0.625 ± 0.049
background (high freq)	0.908 ± 0.112	0.557 ± 0.049	0.493 ± 0.034
whole (low freq)	0.925 ± 0.041	0.638 ± 0.055	0.721 ± 0.047
lesion (low freq)	0.935 ± 0.087	0.582 ± 0.048	0.610 ± 0.079
background (low freq)	0.890 ± 0.071	0.590 ± 0.093	0.538 ± 0.056

Table 2: Performance in AUROC with standard deviations for the ResNet50 melanoma classifier trained on the PH2 dataset and evaluated on the ISIC, HAM10000, and PH2 datasets.

	PH2	ISIC
dark corner	0.940	0.370
hair	0.415	0.585
gel border	0.150	0.262
gel bubble	0.440	0.488
ruler	-	0.495
ink	0.010	0.170
patches	-	0.072

Table 3: Prevalence of artifacts in the PH2 and ISIC datasets. ISIC labels are from Bissoto et al. (2020) and PH2 labels are our manual labels.

4.4. Diffusion-Perturbed Explainability of Classifier Dependency of Artifacts

We evaluate the performance of classifiers trained on PH2 or ISIC data and tested on diffusion-perturbed PH2 datasets, where a specific artifact is added to each image to disrupt any spurious associations between melanoma and the artifact.

First, we find that classifiers trained on limited image regions are more sensitive to single-artifact perturbations after transfer. Specifically, the ISIC-trained classifier fails to generalize across all perturbations in the bbox90 setting, where 90% of the image is covered (Table 4). While the model can still learn weak associations in this heavily occluded setting, these associations do not hold when transferred to a different dataset.

Second, we observe that classifiers are more disrupted by the introduction of certain artifacts than others. Ink markings, rulers, and patches significantly degrade the performance of both PH2- and ISIC-trained classifiers, whereas the addition of dark corners and gel bubbles leads to a smaller decline.

For PH2-trained classifiers, this discrepancy may be attributed to artifact prevalence, where ink, rulers, and patches appear infrequently or are entirely absent. The lack of exposure to these artifacts during training leads to poor generalization when these artifacts are introduced in the perturbed datasets. However, classifiers trained on whole images or background-only images still exhibit reasonably strong transfer performance. While artifacts introduce noise, the whole-image and background settings retain enough diagnostic information for melanoma classification.

For ISIC-trained classifiers, the disruption caused by ink, rulers, and patches likely arises from spurious associations formed during training that do not transfer well to PH2. Despite this, whole-image and background-only models remain robust under perturbation. Whole-image classifiers likely generalize better because they preserve medically relevant features. Background-only models transfers well, suggesting that either the remaining unaltered artifact features continue to provide predictive value or that consistent background features contribute to melanoma classification across both datasets.

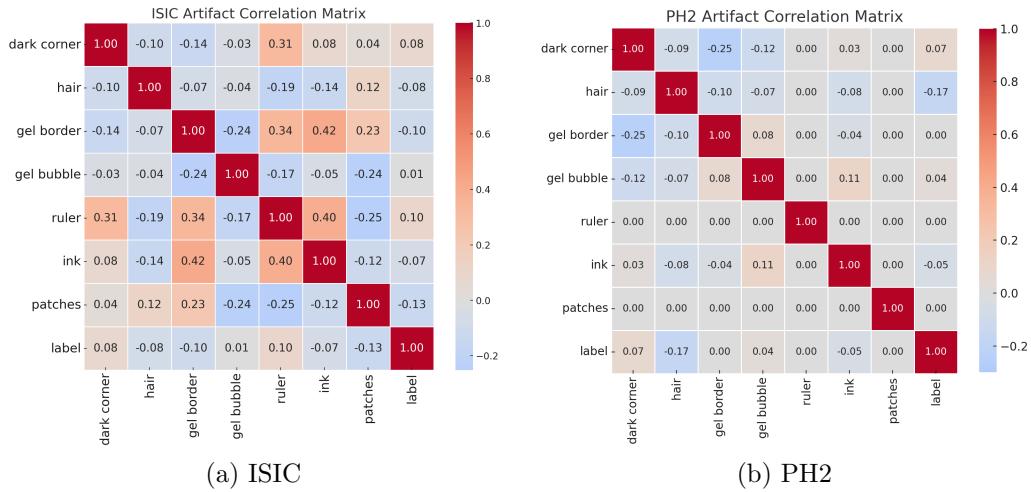


Figure 3: The correlation between the artifacts with themselves and with the melanoma label across the datasets. ISIC labels are from [Bissoto et al. \(2020\)](#) and PH2 labels are our manual labels.

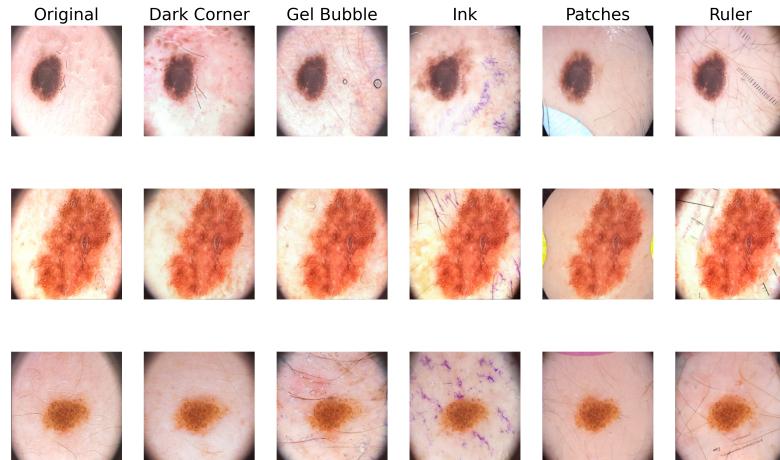


Figure 4: Visualization of three images perturbed by the artifact-specific diffusion model for the augmentation of the PH2 dataset.

These findings emphasize the importance of training classifiers on more comprehensive image regions to improve robustness to artifact perturbations. Models trained with limited image context struggle to generalize when artifacts are introduced in new datasets. Additionally, artifacts that are rare or absent during training can significantly degrade performance when encountered in out-of-domain datasets.

The specific spurious artifacts affecting transfer will depend on the datasets involved, underscoring the need for targeted mitigation strategies.

Dataset Mode	Artifact	PH2	ISIC	Dataset Mode	Artifact	PH2	ISIC
whole	original	0.975	0.858	whole (high freq)	original	0.903	0.442
	dark corner	0.978	0.816		dark corner	0.487	0.430
	gel bubble	0.973	0.841		gel bubble	0.414	0.522
	ink	0.959	0.788		ink	0.484	0.490
	patches	0.976	0.848		patches	0.427	0.417
	ruler	0.966	0.752		ruler	0.335	0.466
background	original	0.884	0.777	whole (low freq)	original	0.935	0.856
	dark corner	0.888	0.759		dark corner	0.836	0.873
	gel bubble	0.893	0.759		gel bubble	0.860	0.866
	ink	0.877	0.729		ink	0.781	0.803
	patches	0.860	0.757		patches	0.818	0.874
	ruler	0.890	0.765		ruler	0.855	0.837
bbox	original	0.880	0.821	background (high freq)	original	0.877	0.502
	dark corner	0.881	0.703		dark corner	0.842	0.376
	gel bubble	0.868	0.693		gel bubble	0.859	0.555
	ink	0.873	0.584		ink	0.840	0.474
	patches	0.896	0.756		patches	0.850	0.479
	ruler	0.880	0.642		ruler	0.847	0.595
bbox70	original	0.845	0.805	background (low freq)	original	0.890	0.817
	dark corner	0.832	0.773		dark corner	0.874	0.794
	gel bubble	0.852	0.787		gel bubble	0.850	0.778
	ink	0.826	0.667		ink	0.893	0.768
	patches	0.810	0.711		patches	0.803	0.790
	ruler	0.856	0.783		ruler	0.889	0.817
bbox90	original	0.865	0.679				
	dark corner	0.875	0.573				
	gel bubble	0.824	0.595				
	ink	0.806	0.546				
	patches	0.809	0.531				
	ruler	0.817	0.498				

Table 4: Performance in AUROC for ResNet50 classifiers evaluated on the PH2 dataset augmented with Dreambooth artifacts. PH2 and ISIC columns indicate the *training* dataset for the melanoma classifier. The models with low performance (below 0.5 AUROC) and models with the worst performance in each artifact subgroup are bolded.

5. Discussion

5.1. Takeaways

In this work, we present a comprehensive evaluation of the robustness of melanoma classification models under different dataset disruptions and diffusion-based perturbations of isolated non-medical artifacts. For researchers or clinicians hoping to implement

deep learning-based models for melanoma classification in practice, we offer the following suggestions: (1) completely removing background information is not an effective solution for mitigating artifacts, as classifiers trained on whole images are more robust under both dataset shift and artifact perturbations; (2) the identification of spurious artifacts requiring correction is dataset-specific and should be systemat-

523
524
525
526
527
528
529
530

531 ically assessed; and (3) mild image compression during
 532 data preprocessing is unlikely to degrade classifier performance, as models primarily rely on low-
 533 frequency features.
 534

535 5.2. Limitations

536 Our proposed diffusion-based perturbation explain-
 537 ability method has several inherent limitations. First,
 538 it requires prior knowledge of the artifact to be anal-
 539 yzed, as well as example images containing the ar-
 540 tifact to train the diffusion model. However, this re-
 541 quirement aligns with the broader challenges of devel-
 542 oping medical deep learning models, in which defining
 543 and verifying what counts as a “non-medical” feature
 544 necessitates clinical expertise. For example, irregular
 545 lesion borders are a clinically significant feature in
 546 melanoma classification, as melanoma often presents
 547 with jagged or poorly defined edges. However, a non-
 548 clinical expert might misinterpret this as a spurious
 549 feature, assuming the model is relying on segmen-
 550 tation artifacts rather than a meaningful diagnostic
 551 signal.

552 Second, we experienced difficulties with finetuning
 553 the diffusion model for the artifacts of “hair” and “gel
 554 border”, suggesting that some artifacts are inherently
 555 more challenging to learn. While we used standard
 556 Dreambooth and LoRA finetuning techniques, more
 557 specialized diffusion-based methods may help address
 558 these limitations. Lastly, the diffusion model may
 559 unintentionally introduce additional diffusion-specific
 560 artifacts. To improve the quality of artifact perturba-
 561 tions, highly targeted artifact classification or detec-
 562 tion models could be implemented to filter the gener-
 563 ated images. We leave this direction for future work.

564 6. Conclusion

565 Past work have shown that non-medical artifacts,
 566 such as surgical markings and dark corners, can in-
 567 troduce biases that hinder model generalization in
 568 melanoma classification. In this study, we systemat-
 569 ically evaluated the impact of common artifacts and
 570 introduced a diffusion-based perturbation method to
 571 isolate and analyze their effects. Our findings show
 572 that whole-image classifiers outperform lesion-only
 573 or background-only models, emphasizing the impor-
 574 tance of preserving contextual information. Addi-
 575 tionally, melanoma classifiers rely primarily on low-
 576 frequency features, making them less sensitive to
 577 mild image compression. Through controlled artifact

578 perturbations, we identify ink markings, rulers, and
 579 patches as key sources of classifier sensitivity, while
 580 other artifacts, such as dark corners and gel bubbles,
 581 have a lesser impact. These results underscore the
 582 need for dataset-specific artifact analysis and miti-
 583 gation strategies to improve the robustness of deep
 584 learning models in melanoma classification.

585 7. Acknowledgements

586 This material is based upon work supported by the
 587 National Science Foundation Graduate Research Fel-
 588 lowship under Grant No. 2389810.

589 **References**

- 590 2kpr. Dreambooth Tokens, 2024. URL <https://github.com/2kpr/dreambooth-tokens>. Accessed: 2025-02-11.
- 593 Jose-Agustin Almaraz-Damian, Volodymyr Ponomaryov, Sergiy Sadovnichiy, and Heydy Castillejos-Fernandez. Melanoma and nevus skin lesion classification using handcraft and deep learning feature fusion via mutual information measures. *Entropy*, 22(4):484, 2020.
- 599 Giuseppe Argenziano, Gabriella Fabbrocini, Paolo Carli, Vincenzo De Giorgi, Elena Sammarco, and Mario Delfino. Epiluminescence microscopy for the diagnosis of doubtful melanocytic skin lesions: comparison of the abcd rule of dermatoscopy and a new 7-point checklist based on pattern analysis. *Archives of dermatology*, 134(12):1563–1570, 1998.
- 606 Muhammad Asad Arshed, Shahzad Mumtaz, Muhammad Ibrahim, Saeed Ahmed, Muhammad Tahir, and Muhammad Shafi. Multi-class skin cancer classification using vision transformer networks and convolutional neural network-based pre-trained models. *Information*, 14(7):415, 2023.
- 612 Peter J Bevan and Amir Atapour-Abarghouei. Skin deep unlearning: Artefact and instrument debiasing in the context of melanoma classification. *arXiv preprint arXiv:2109.09818*, 2021.
- 616 Alceu Bissoto, Eduardo Valle, and Sandra Avila. Debiasing skin lesion datasets and models? not so fast. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 740–741, 2020.
- 621 Christopher Boland, Keith A Goatman, Sotirios A Tsaftaris, and Sonia Dahdouh. There are no shortcuts to anywhere worth going: Identifying shortcuts in deep learning models for medical image analysis. In *Medical Imaging with Deep Learning*, 2024.
- 627 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.
- 632 Arief Budhiman, Suyanto Suyanto, and Anditya Ari-fianto. Melanoma cancer classification using resnet with data augmentation. In *2019 international seminar on research of information technology and intelligent systems (ISRITI)*, pages 17–20. IEEE, 2019.
- 638 Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.
- 646 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- 651 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- 655 Giuseppe Di Leo, Alfredo Paolillo, Paolo Sommella, and Gabriella Fabbrocini. Automatic diagnosis of melanoma: a software system based on the 7-point check-list. In *2010 43rd Hawaii international conference on system sciences*, pages 1–10. IEEE, 2010.
- 661 Muhammad Ali Farooq, Wang Yao, Michael Schukat, Mark A Little, and Peter Corcoran. Derm-t2im: Harnessing synthetic skin lesion data via stable diffusion models for enhanced skin disease classification using vit and cnn. *arXiv preprint arXiv:2401.05159*, 2024.
- 667 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- 672 Lahiru Gamage, Udhitha Isuranga, Dulani Meedeniya, Senuri De Silva, and Pratheepan Yogarajah. Melanoma skin cancer identification with explainability utilizing mask guided technique. *Electronics*, 13(4):680, 2024.
- 677 Bilge S Akkoca Gazioglu and Mustafa E Kamaşak. Effects of objects and image quality on melanoma classification using deep neural networks. *Biomedical Signal Processing and Control*, 67:102530, 2021.

- 681 Holger A Haensle, Christine Fink, Roland Schnei-
682 derbauer, Ferdinand Toberer, Timo Buhl, Andreas
683 Blum, Aadi Kalloo, A Ben Hadj Hassen, Luc
684 Thomas, Alexander Enk, et al. Man against ma-
685 chine: diagnostic performance of a deep learn-
686 ing convolutional neural network for dermoscopic
687 melanoma recognition in comparison to 58 der-
688 matologists. *Annals of oncology*, 29(8):1836–1842,
689 2018.
- 690 James A. Hanley and Barbara J. McNeil. The mean-
691 ing and use of the area under a receiver operating
692 characteristic (roc) curve. *Radiology*, 143(1):29–36,
693 1982. doi: 10.1148/radiology.143.1.7063747.
- 694 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian
695 Sun. Deep residual learning for image recogni-
696 tion. In *Proceedings of the IEEE conference on com-
697 puter vision and pattern recognition*, pages 770–
698 778, 2016.
- 699 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
700 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
701 and Weizhu Chen. Lora: Low-rank adapta-
702 tion of large language models. *arXiv preprint
arXiv:2106.09685*, 2021.
- 704 Gao Huang, Zhuang Liu, Laurens Van Der Maaten,
705 and Kilian Q Weinberger. Densely connected con-
706 volutional networks. In *Proceedings of the IEEE
707 conference on computer vision and pattern recogni-
708 tion*, pages 4700–4708, 2017.
- 709 Jeremy Kawahara, Sara Daneshvar, Giuseppe Argen-
710 ziano, and Ghassan Hamarneh. Seven-point check-
711 list and skin lesion classification using multitask
712 multimodal neural nets. *IEEE journal of biomed-
713 ical and health informatics*, 23(2):538–546, 2018.
- 714 Peiyuan Liao, Xiuyu Li, Xihui Liu, and Kurt
715 Keutzer. The artbench dataset: Benchmarking
716 generative models with artworks. *arXiv preprint
arXiv:2206.11404*, 2022.
- 718 Caterina Longo, Riccardo Pampena, Elvira
719 Moscarella, Johanna Chester, Michela Starace,
720 Elisa Cinotti, Bianca Maria Piraccini, Giuseppe
721 Argenziano, Ketty Peris, and Giovanni Pellacani.
722 Dermoscopy of melanoma according to different
723 body sites: Head and neck, trunk, limbs, nail,
724 mucosal and acral. *Journal of the European
725 Academy of Dermatology and Venereology*, 37(9):
726 1718–1730, 2023.
- 727 Teresa Mendonça, Pedro M Ferreira, Jorge S Mar-
728 ques, André RS Marcal, and Jorge Rozeira. Ph
729 2-a dermoscopic image database for research and
730 benchmarking. In *2013 35th annual international
731 conference of the IEEE engineering in medicine
732 and biology society (EMBC)*, pages 5437–5440.
733 IEEE, 2013.
- 734 Matthias Minderer, Olivier Bachem, Neil Houlsby,
735 and Michael Tschannen. Automatic shortcut re-
736 moval for self-supervised representation learning.
737 In *International Conference on Machine Learning*,
738 pages 6927–6937. PMLR, 2020.
- 739 Franz Nachbar, Wilhelm Stolz, Tanja Merkle,
740 Armand B Cognetta, Thomas Vogt, Michael
741 Landthaler, Peter Bilek, Otto Braun-Falco, and
742 Gerd Plewig. The abcd rule of dermatoscopy:
743 high prospective value in the diagnosis of doubtful
744 melanocytic skin lesions. *Journal of the American
745 Academy of Dermatology*, 30(4):551–559, 1994.
- 746 Robin Rombach, Andreas Blattmann, Dominik
747 Lorenz, Patrick Esser, and Björn Ommer. High-
748 resolution image synthesis with latent diffusion
749 models. In *Proceedings of the IEEE/CVF confer-
750 ence on computer vision and pattern recogni-
751 tion*, pages 10684–10695, 2022.
- 752 Olaf Ronneberger, Philipp Fischer, and Thomas
753 Brox. U-net: Convolutional networks for biomedical
754 image segmentation. In *Medical image comput-
755 ing and computer-assisted intervention—MICCAI
756 2015: 18th international conference, Munich, Ger-
757 many, October 5–9, 2015, proceedings, part III 18*,
758 pages 234–241. Springer, 2015.
- 759 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael
760 Pritch, Michael Rubinstein, and Kfir Aberman.
761 Dreambooth: Fine tuning text-to-image diffusion
762 models for subject-driven generation. In *Proceed-
763 ings of the IEEE/CVF Conference on Computer
764 Vision and Pattern Recognition*, pages 22500–
765 22510, 2023.
- 766 Luke W Sagers, James A Diao, Matthew Groh,
767 Pranav Rajpurkar, Adewole S Adamson, and Ar-
768 jun K Manrai. Improving dermatology classifiers
769 across populations using images generated by large
770 diffusion models. *arXiv preprint arXiv:2211.13352*,
771 2022.
- 772 Zohaib Salahuddin, Henry C Woodruff, Avishek
773 Chatterjee, and Philippe Lambin. Transparency of

- 774 deep neural networks for medical image analysis:
 775 A review of interpretability methods. *Computers
 776 in biology and medicine*, 140:105111, 2022.
- 777 Katharina Sies, Julia K Winkler, Christine Fink,
 778 Felicitas Bardehle, Ferdinand Toberer, Felix KF
 779 Kommoß, Timo Buhl, Alexander Enk, Albert
 780 Rosenberger, and Holger A Haenssle. Dark corner
 781 artefact and diagnostic performance of a market-
 782 approved neural network for skin cancer classifica-
 783 tion. *JDDG: Journal der Deutschen Dermatologis-
 784 chen Gesellschaft*, 19(6):842–850, 2021.
- 785 Sahil Singla and Soheil Feizi. Salient imangenet:
 786 How to discover spurious features in deep learning?
 787 *arXiv preprint arXiv:2110.04301*, 2021.
- 788 Sidharth Sonthalia, Sara Yumeen, and Feroze
 789 Kaliyadan. Dermoscopy overview and extradiag-
 790 nostic applications. 2019.
- 791 Philipp Tschandl, Cliff Rosendahl, and Harald Kit-
 792 tler. The ham10000 dataset, a large collection of
 793 multi-source dermatoscopic images of common pig-
 794 mented skin lesions. *Scientific data*, 5(1):1–9, 2018.
- 795 Philipp Tschandl, Noel Codella, Bengü Nisa Akay,
 796 Giuseppe Argenziano, Ralph P Braun, Horacio
 797 Cabo, David Gutman, Allan Halpern, Brian Helba,
 798 Rainer Hofmann-Wellenhof, et al. Comparison of
 799 the accuracy of human readers versus machine-
 800 learning algorithms for pigmented skin lesion clas-
 801 sification: an open, web-based, international, diag-
 802 nostic study. *The lancet oncology*, 20(7):938–947,
 803 2019.
- 804 Alfredo Vellido. The importance of interpretability
 805 and visualization in machine learning for applica-
 806 tions in medicine and health care. *Neural comput-
 807 ing and applications*, 32(24):18069–18083, 2020.
- 808 Patrick von Platen, Suraj Patil, Anton Lozhkov,
 809 Kashif Rasul, William Falcon, Sayak Paul, and the
 810 Hugging Face Team. Diffusers: State-of-the-art
 811 pretrained diffusion models, 2022. URL <https://github.com/huggingface/diffusers>.
- 812
- 813 Natalie M Williams, Kristina D Rojas, John M
 814 Reynolds, Deukwoo Kwon, Jackie Shum-Tien, and
 815 Natalia Jaimes. Assessment of diagnostic accu-
 816 racy of dermoscopic structures and patterns used
 817 in melanoma detection: a systematic review and
 818 meta-analysis. *JAMA dermatology*, 157(9):1078–
 819 1088, 2021.
- Julia K Winkler, Christine Fink, Ferdinand Toberer,
 820 Alexander Enk, Teresa Deinlein, Rainer Hofmann-
 821 Wellenhof, Luc Thomas, Aimilios Lallas, Andreas
 822 Blum, Wilhelm Stolz, et al. Association between
 823 surgical skin markings in dermoscopic images and
 824 diagnostic performance of a deep learning convo-
 825 lutional neural network for melanoma recognitio
 826 *JAMA dermatology*, 155(10):1135–1141, 2019.
 827
- Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang,
 828 Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang
 829 Wen. Paint by example: Exemplar-based image
 830 editing with diffusion models. In *Proceedings of the
 831 IEEE/CVF Conference on Computer Vision and
 832 Pattern Recognition*, pages 18381–18391, 2023.
 833

834 **Appendix A. DenseNet121 Transfer
835 Performance under
836 Dataset and
837 Diffusion-based
838 Perturbations**

839 We repeat the dataset disruption and artifact pertur-
840 bation experiments using DenseNet121 (Huang et al.,
841 2017), another widely used model architecture for
842 melanoma classification.

843 Overall, we observe similar trends to our ResNet50
844 results. Due to dataset size limitations, PH2 classi-
845 fiers generally do not transfer well to other datasets.
846 In ISIC, the whole-image classifier outperforms both
847 the lesion-only and background-only classifiers, not
848 only within the in-domain ISIC dataset but also
849 when transferred to out-of-domain datasets such as
850 HAM10000 and PH2 (Table 5). Additionally, ISIC
851 bounding box models exhibit reasonable transferabil-
852 ity across datasets, reinforcing the idea that certain
853 artifacts in the outer image regions contribute to a
854 degree of generalization. Finally, we find that low-
855 frequency models outperform high-frequency models,
856 consistent with our observations from ResNet50.

857 **Appendix B. Stacked Diffusion
858 Perturbations of Two
859 Artifact Effect on
860 Classifier Performance**

861 In this section, we study the effect of stacking
862 diffusion-based artifact perturbations. Given the al-
863 ready generated set of artifact-perturbed images from
864 Section 4.4, we stack a second artifact perturbation
865 through using our fine-tuned diffusion models to in-
866 paint a different, second artifact. For instance, a
867 heatmap entry with row “ink” and column “ruler”
868 means that we use the ink-specific diffusion model to
869 inpaint the artifact of “ink” on top of images already
870 perturbed with “ruler”.

871 In this section, we analyze the impact of stack-
872 ing multiple diffusion-based artifact perturbations.
873 Building on the artifact-perturbed images generated
874 in Section 4.4, we apply a second perturbation using
875 our fine-tuned diffusion models to inpaint a different
876 artifact on top of the existing perturbation. For ex-
877 ample, in the heatmap, an entry at the intersection of
878 row “ink” and column “ruler” indicates that the ink-
879 specific diffusion model was used to overlay an “ink”

880 artifact onto images that had already been perturbed
881 with “ruler.”

882 For the PH2-trained classifiers (Figure 5), we see
883 that the whole-image classifier remains robust even
884 under two artifact perturbations. When we exam-
885 ine the bounding box classifiers, we see that com-
886 binations of patches and dark corner (AUROC bbox:
887 0.832) and patches and gel bubbles (AUROC bbox70:
888 0.786, bbox90: 0.776) suffer from the greatest drop
889 in performance within each dataset mode. This in-
890 creased sensitivity to patches is consistent with our
891 findings from the single-artifact perturbation experi-
892 ments.

893 For the ISIC-trained classifiers (Figure 6), whole-
894 image and low-frequency classifiers are robust to arti-
895 fact perturbations, whereas high-frequency classifiers
896 fail to generalize. This finding is also consistent with
897 our single-artifact results. Interestingly, we observe
898 that both the bbox and bbox90 classifiers struggle to
899 transfer under two artifact perturbations. This dif-
900 fers from the single-artifact scenario, where the bbox
901 classifier showed better generalization. These find-
902 ings suggest that as artifact perturbations become
903 more complex, classifier generalization diminishes.

Domain AUROC	Source ISIC	Target HAM10000	Target PH2
whole	0.804 ± 0.029	0.781 ± 0.033	0.890 ± 0.082
lesion	0.758 ± 0.051	0.725 ± 0.027	0.810 ± 0.090
background	0.735 ± 0.033	0.733 ± 0.019	0.829 ± 0.076
bbox	0.720 ± 0.017	0.673 ± 0.038	0.810 ± 0.149
bbox70	0.661 ± 0.019	0.674 ± 0.057	0.762 ± 0.134
bbox90	0.634 ± 0.060	0.647 ± 0.047	0.680 ± 0.113
whole (high freq)	0.687 ± 0.034	0.567 ± 0.047	0.741 ± 0.033
lesion (high freq)	0.696 ± 0.014	0.624 ± 0.037	0.768 ± 0.152
background (high freq)	0.717 ± 0.009	0.623 ± 0.026	0.799 ± 0.068
whole (low freq)	0.807 ± 0.026	0.756 ± 0.018	0.927 ± 0.051
lesion (low freq)	0.790 ± 0.036	0.714 ± 0.047	0.900 ± 0.078
background (low freq)	0.723 ± 0.021	0.683 ± 0.013	0.849 ± 0.104

Table 5: Performance in AUROC with standard deviations for the DenseNet121 melanoma classifier trained on the ISIC dataset and evaluated on the ISIC, HAM10000, and PH2 datasets.

Domain AUROC	Source PH2	Target ISIC	Target HAM10000
whole	0.933 ± 0.042	0.573 ± 0.054	0.440 ± 0.045
lesion	0.956 ± 0.030	0.623 ± 0.038	0.700 ± 0.032
background	0.905 ± 0.054	0.624 ± 0.031	0.521 ± 0.052
bbox	0.888 ± 0.096	0.622 ± 0.015	0.554 ± 0.051
bbox70	0.871 ± 0.088	0.562 ± 0.056	0.456 ± 0.045
bbox90	0.832 ± 0.063	0.576 ± 0.023	0.552 ± 0.040
whole (high freq)	0.804 ± 0.100	0.511 ± 0.008	0.420 ± 0.062
lesion (high freq)	0.911 ± 0.104	0.595 ± 0.043	0.359 ± 0.022
background (high freq)	0.928 ± 0.078	0.564 ± 0.022	0.333 ± 0.014
whole (low freq)	0.930 ± 0.040	0.575 ± 0.054	0.551 ± 0.022
lesion (low freq)	0.936 ± 0.057	0.633 ± 0.038	0.723 ± 0.013
background (low freq)	0.869 ± 0.101	0.589 ± 0.031	0.536 ± 0.052

Table 6: Performance in AUROC with standard deviations for the DenseNet121 melanoma classifier trained on the PH2 dataset and evaluated on the ISIC, HAM10000, and PH2 datasets.

Dataset Mode	Artifact	PH2	ISIC	Dataset Mode	Artifact	PH2	ISIC
whole	original	0.933	0.890	whole (high freq)	original	0.804	0.741
	dark corner	0.947	0.846		dark corner	0.466	0.704
	gel bubble	0.948	0.888		gel bubble	0.505	0.738
	ink	0.946	0.837		ink	0.452	0.775
	patches	0.937	0.841		patches	0.351	0.798
	ruler	0.946	0.876		ruler	0.447	0.790
background	original	0.905	0.829	whole (low freq)	original	0.930	0.927
	dark corner	0.906	0.805		dark corner	0.928	0.893
	gel bubble	0.920	0.793		gel bubble	0.909	0.906
	ink	0.926	0.725		ink	0.903	0.848
	patches	0.904	0.767		patches	0.922	0.871
	ruler	0.920	0.812		ruler	0.915	0.888
bbox	original	0.888	0.810	background (high freq)	original	0.928	0.799
	dark corner	0.909	0.735		dark corner	0.865	0.744
	gel bubble	0.891	0.739		gel bubble	0.925	0.695
	ink	0.894	0.676		ink	0.827	0.712
	patches	0.919	0.689		patches	0.791	0.778
	ruler	0.910	0.734		ruler	0.864	0.799
bbox70	original	0.871	0.762	background (low freq)	original	0.869	0.849
	dark corner	0.861	0.695		dark corner	0.881	0.816
	gel bubble	0.878	0.753		gel bubble	0.869	0.802
	ink	0.854	0.613		ink	0.881	0.804
	patches	0.856	0.638		patches	0.877	0.841
	ruler	0.865	0.744		ruler	0.870	0.879
bbox90	original	0.832	0.680				
	dark corner	0.813	0.617				
	gel bubble	0.819	0.630				
	ink	0.819	0.592				
	patches	0.774	0.595				
	ruler	0.821	0.624				

Table 7: AUROC performance for DenseNet121 classifiers evaluated on the PH2 dataset augmented with Dreambooth artifacts. PH2 and ISIC columns indicate the *training* dataset for the melanoma classifier.

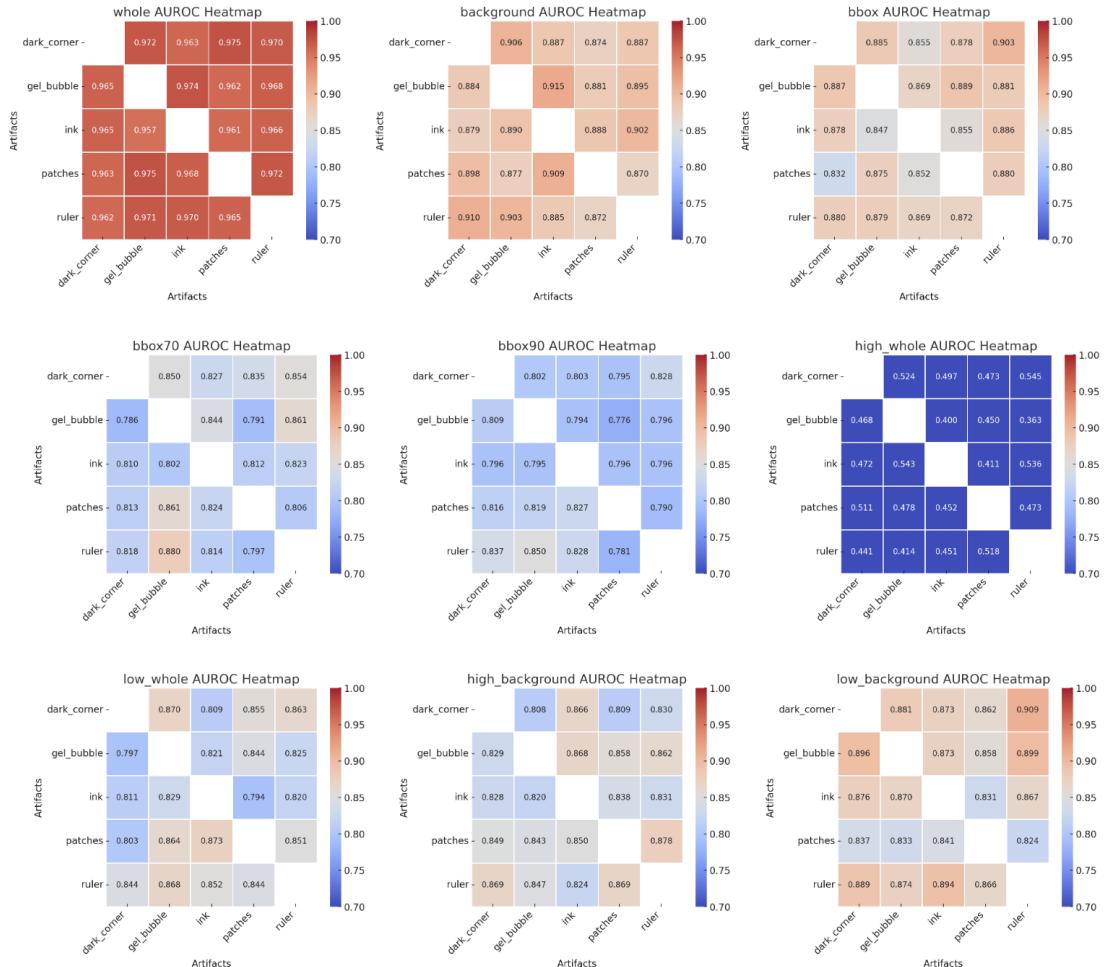


Figure 5: Heatmaps of test AUROC performances of the PH2 classifier under different dataset modes with two stacked artifact perturbations of the PH2 dataset. The x-axis is the first added artifact, and the y-axis is the second added artifact. Due to the stochasticity of diffusion generation, the matrix is not symmetric.

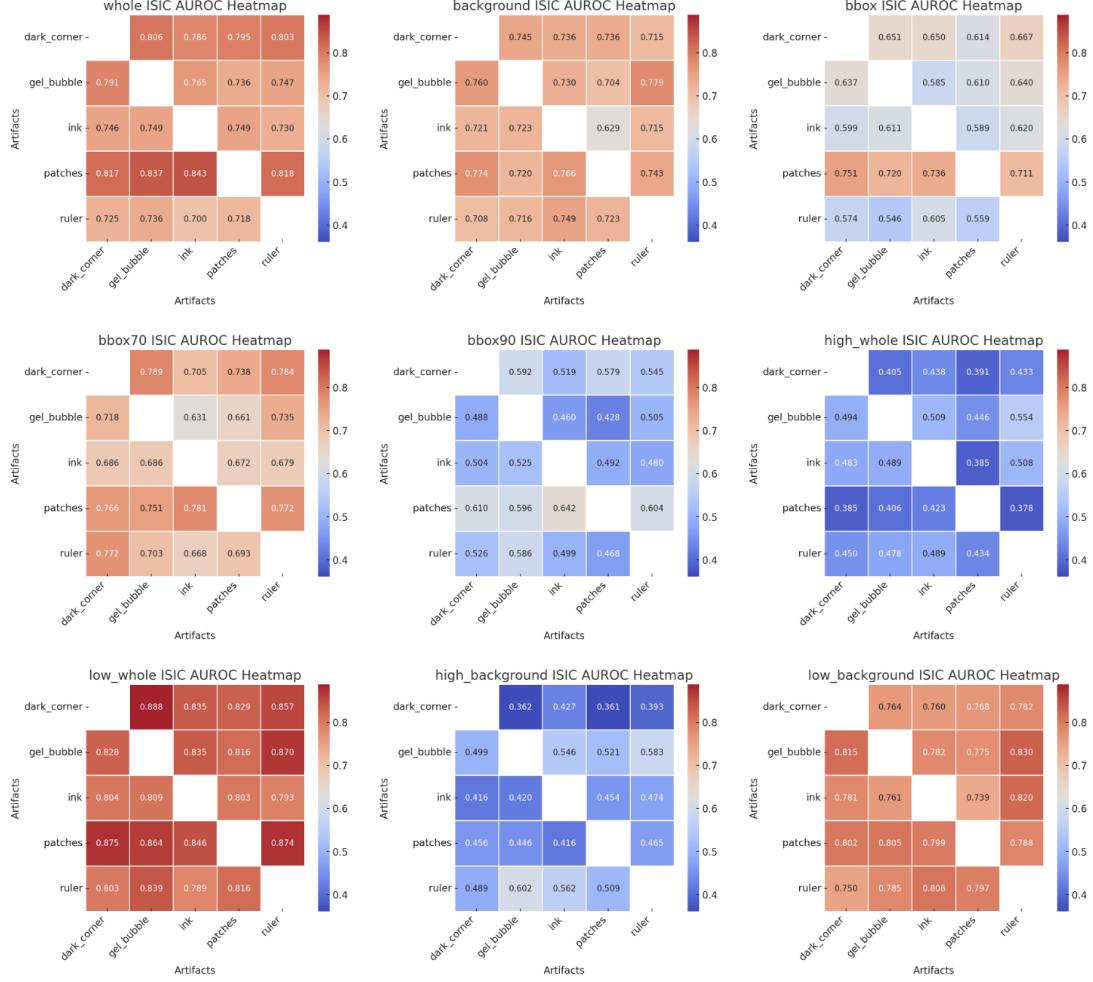


Figure 6: Heatmaps of test AUROC performances of the ISIC classifier under different dataset modes with two stacked artifact perturbations of the PH2 dataset. The x-axis is the first added artifact, and the y-axis is the second added artifact. Due to the stochasticity of diffusion generation, the matrix is not symmetric. Note that the scale of the colorbars for ISIC classifier is different from the PH2 classifiers, so colors cannot be directly compared across figures.