# Generalized Design Choices for Deepfake Detectors

Lorenzo Pellegrini*§, Serafino Pandolfini*, Davide Maltoni*,
Matteo Ferrara*, Marco Prati†, Marco Ramilli†

*Department of Computer Science and Engineering (DISI)
University of Bologna, Italy

†IdentifAI, Italy

*Abstract*—The effectiveness of deepfake detection methods often depends less on their core design and more on implementation details such as data preprocessing, augmentation strategies, and optimization techniques. These factors make it difficult to fairly compare detectors and to understand which factors truly contribute to their performance. To address this, we systematically investigate how different design choices influence the accuracy and generalization capabilities of deepfake detection models, focusing on aspects related to training, inference, and incremental updates. By isolating the impact of individual factors, we aim to establish robust, architecture-agnostic best practices for the design and development of future deepfake detection systems. Our experiments identify a set of design choices that consistently improve deepfake detection and enable state-of-the-art performance on the AI-GenBench benchmark.

*Index Terms*—Deepfake detection, AI-generated image, AI-GenBench benchmark, design choices.

Fig. 1: The different dimensions explored in this work, to optimize the training, inference and incremental update (grayed) of deepfake detectors.

## I. INTRODUCTION

The rapid advancement of generative models has led to an unprecedented ability to produce realistic synthetic images that are increasingly difficult to distinguish from human-generated (real) content. Diffusion-based architectures and large-scale text-to-image systems such as Stable Diffusion, Midjourney, and DALL-E have significantly lowered the barrier to generating high-quality content. In particular, this has enabled professionals as well as the general public to create new media content conditioned by a text prompt and other semantic inputs. These technologies have often been misused to disseminate disinformation, raising societal concerns and establishing the detection of AI-generated content as an important research area [1], [2], [3], which is essential for preserving trust, accountability, and authenticity in digital media.

Over the past few years, numerous detection approaches have been proposed, ranging from handcrafted forensic cues to deep neural networks specifically trained to distinguish real from synthetic content. Despite promising progress, detection performance often varies widely across studies, datasets, and model architectures. In many cases, the reported success of a particular method depends less on the core detection idea than on specific, and sometimes implicit, implementation details—such as the choice of data augmentations, preprocessing, or training strategy. This lack of systematic evaluation makes it difficult to identify which design factors truly contribute to
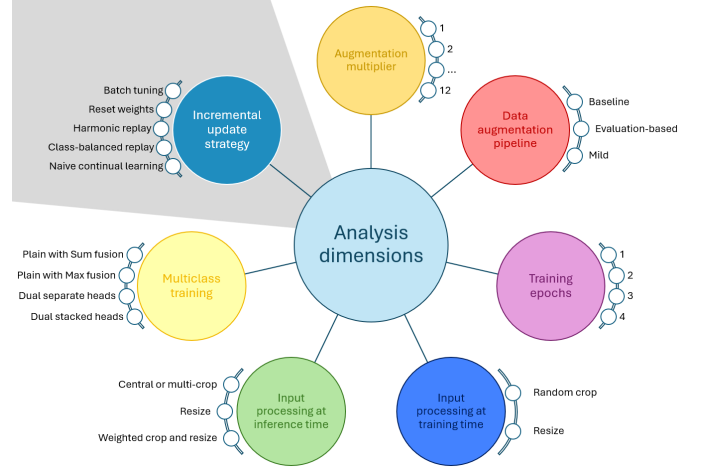
generalization across different types of generators and model architectures. Moreover, existing works often focus on training detection models on content produced by a limited (or even a single) set of hand-picked generators and then testing such models on images from other generators.

In this work, we present a comprehensive empirical study aimed at isolating general design principles from the influence of specific model architectures or fake data generators. To this purpose, we adopt the recent AI-GenBench benchmark [4], which temporally orders image generators to simulate the release of new generative models over time. With this setup, we systematically evaluate how various training and inference-time choices affect the generalization capability of detection models on both old (such as GAN-based) and recent generative techniques. Specifically, we analyze the impact of (i) the data augmentation pipeline, (ii) the training duration and augmentation multiplier, (iii) the preprocessing strategies such as cropping versus resizing, and (iv) the use of multiclass labels and related strategies. As a further dimension of analysis, we consider the incremental training of deepfake detectors. This is a relevant issue for the practical deployment of detectors that must be frequently updated to cope with novel and emerging generation techniques. Since retraining from scratch on a steadily increasing dataset can be very resource-consuming,

§Corresponding author: Lorenzo Pellegrini (l.pellegrini@unibo.it)
Project repository: https://github.com/MI-BioLab/AI-GenBench

we evaluate how detection models can be trained incrementally in a sample-efficient way while preserving detection capabilities on older generators, a scenario commonly referred to as *Continual Lifelong Learning*. An overview of the directions explored in this work is reported in Figure 1.

Our goal is not to introduce a new detection method, but rather to identify a set of robust, architecture-agnostic best practices that consistently enhance performance and generalization across diverse model families. To this end, we systematically evaluate multiple pre-trained vision backbones such as ResNet-50 CLIP, ViT-L CLIP, and DINOv2. The findings of this study offer actionable insights for the development of future detection systems, equipping both researchers and practitioners with a solid foundation for designing robust detectors of AI-generated images. To the best of our knowledge, this is the first model-agnostic study that systematically evaluates all those dimensions of deepfake detection, covering both training and evaluation mechanisms.

This work is organized as follows. Section II introduces the background and reviews relevant related work. Section III describes the experimental setup, detailing the adopted benchmark and the main dimensions of analysis. Section IV presents the results of these experiments. We discuss Incremental Update Strategies in Section V, as this topic is mostly orthogonal to the other dimensions. Section VI reports the detection performance achieved by the "best of" configuration, which combines the most effective approaches identified in our analysis. Finally, Section VII summarizes our findings and outlines future research directions.

## II. RELATED WORK

Deepfake detection aims to distinguish AI-generated content from authentic media including images, videos, and audio. In this section, we review relevant literature in two key areas of deepfake image detection: (i) detection methods, which focus on the backbones and algorithmic strategies used to identify manipulated content, and (ii) benchmarks, which provide datasets and evaluation protocols essential for developing, validating, and comparing detection techniques. A final subsection is devoted to a brief review of relevant continual learning techniques.

### A. Detection Methods

Early methods for differentiating synthetic from authentic images predominantly relied on Convolutional Neural Networks (CNNs) trained on large-scale datasets [3]. While these approaches achieve high accuracy under conditions closely aligned with the training distribution, their performance degrades significantly in real-world scenarios. In particular, they tend to be vulnerable to common image degradations such as compression, resizing, blurring, and cropping that frequently occur when images are shared via social media or instant messaging platforms. In such scenarios, detection systems often struggle to generalize to images generated by previously unseen models [5]. To mitigate these limitations, incorporating carefully designed data augmentation strategies during training has proven effective. These augmentations not only enhance

robustness against image-level perturbations but also improve cross-generator generalization [6]. Consequently, the design of the training augmentation pipeline is one of the main elements investigated in our study. While ImageNet-pretrained CNN backbones dominated early research in this field, recent work has explored large models based on different structures, such as Vision Transformers, and different pre-training strategies, such as vision-language models like CLIP. These demonstrate strong performance even when trained on data from a single generator, thanks to their rich feature representations and superior transferability [7]. Furthermore, recent studies suggest that foundational vision backbones such as DINOv2, trained using self-supervised learning, may be particularly effective for deepfake detection [4].

Among the plethora of methods proposed in the literature, some of them focus on the modification of network architectures to better capture low- and high-level forensic traces [8], [9], while others improve training strategies [10], [11] or simulate generator-specific artifacts [12], [13]. Recent research has also explored formulating the detection task beyond traditional binary or multiclass classification. For instance, LASTED [14] employs a language-guided contrastive learning objective to align images with descriptive text prompts, thereby learning representations that generalize more effectively to unseen generators. Another strategy for improving generalization involves few-shot or incremental learning methods [15], [16], [17]. While promising, these approaches require access to images from generators, which may not always be available in the most challenging scenarios. An alternative line of research considers periodically retraining detectors while preserving the temporal order of generator releases [18]. This approach leverages forensic traces from known generators, which are often similar to those in newer models. Indeed, it is reasonable to believe that artificial fingerprints [19] from one generator can enable classifiers to generalize across entire families of models, not just individual ones [6].

### B. Benchmarks

Early benchmarks for synthetic image detection primarily focused on GAN-based generators and were often limited to specific domains, such as facial imagery—e.g., ForgeryNet [20], DiffusionFace [21], and DIFF [22]—or artistic content [23]. To overcome these domain-specific limitations, recent benchmarks emphasize generalization by introducing large-scale, diverse datasets [11], [24], [25], [26], or by providing open-source frameworks that facilitate the integration and evaluation of new generative models [27]. Additional efforts include comprehensive evaluations of existing datasets [28], as well as human perceptual studies on synthetic content detection [29]. Moreover, several influential datasets, although not explicitly designed as benchmarks, are widely adopted by the research community for training and evaluation purposes. These include both GAN-based [6] and diffusion-based image collections [7], [30], [31], [32].

A common evaluation protocol across many benchmarks is to assess the generalization capability of detection models by testing them on generators unseen during training.

However, this setup often overlooks the temporal evolution of generative techniques. In practice, new architectures are released continuously, making it essential to evaluate generalization under temporally realistic conditions: training on older generators and testing on newer ones according to their historical release timeline. This perspective was first introduced in [18], which demonstrated a significant drop in detection performance when models encountered major shifts in generative architectures. Building on this insight, AI-GenBench [4] introduced a temporal evaluation benchmark comprising 36 mainstream generative models released between 2017 and 2024. In our experiments, we adopt AI-GenBench as it provides a more realistic and forward-looking framework for assessing generalization over time. This benchmark is based on a protocol that defines the temporal order in which the generators are encountered. In addition, to allow for a fair comparison of different approaches, it also defines rules regarding the augmentation intensity to be used during training, the evaluation pipeline (and especially the augmentations used to introduce social media-alike distortions), and the metrics used to evaluate the ability of each model to both generalize to unseen generators and retain detection capabilities on older ones. More information on the experimental protocol will be given in Section III.

### C. Continual learning

Continual learning addresses the challenge of updating models with new data over time without losing performance on previously learned tasks, a problem commonly known as catastrophic forgetting [33]. Existing approaches can be broadly categorized into three groups: (i) regularization-based methods, which constrain weight updates to preserve prior knowledge; (ii) architectural methods, which expand the model to allocate capacity for new tasks; and (iii) replay-based (or rehearsal) strategies, which retain and interleave a subset of past samples during training. Among these, replay methods are particularly popular due to their simplicity and effectiveness. They typically rely on memory buffers with various sampling and replacement policies [34], [35]. In this work, we adopt replay-based strategies as a practical mechanism to mitigate forgetting in deepfake detection models, enabling adaptation to new generators without retraining on the entire data of past generators.

### III. EXPERIMENTAL SETUP

The evaluation is conducted using the AI-GenBench temporal framework. In this benchmark, the 36 image generators are ordered by release date and split into temporal windows, each containing four generators. The detection model is trained progressively: at each step $k$, the model is trained on all generators within the sliding windows $w_j, j \leq k$. This setup simulates a realistic scenario where detectors are periodically retrained to keep up with novel generative models. After each training step $k$, the model is evaluated to measure its ability to detect images from both past and future generators. The benchmark defines a set of three scenarios on which the relevant metrics are measured:

- *Next Period* - the detection performance is measured on the generators of the next sliding window ($w_{k+1}$).
- *Past Period* - performance is measured on the generators belonging to windows $w_j, j \leq k$.
- *Whole Period* - performance is measured on the generators belonging to both the past and next time windows ($w_j, j \leq k + 1$).

The benchmark proposes the *Area Under Receiver Operating Characteristic Curve (AUROC)* as the main metric, which is averaged across all steps to obtain a single compact value. The performance measured on the *Next Period* is particularly important as it measures the detector's ability to generalize to unseen generators, which will become available in the near future. For this reason, the authors of AI-GenBench consider the *average AUROC on the Next Period* as the main metric.

To identify which strategies generalize across different families of detectors, we consider the following well-known pre-trained vision (and language-vision) models: i) *ResNet-50 CLIP* by OpenAI [36], ii) *ViT-L/14 CLIP* from LAION models[1], and iii) *ViT-L/14 DINOv2* [37].

We focus on the following design dimensions:

- *Data augmentation pipeline* - evaluating the impact of transformations such as color jitter, Gaussian noise, blurring, geometric transformations, and especially JPEG compression.
- *Augmentation multiplier* - given an augmentation pipeline, systematic varying the number of diverse images presented to the model.
- *Training duration* - determining the optimal number of training epochs, for a given augmentation pipeline and multiplier.
- *Input processing at training time* - comparing training strategies based on image crops versus resized full images.
- *Input processing at inference time* - evaluating whether binary predictions are best obtained by (i) fusing scores from multiple image crops, (ii) using a resized version of the full image, or (iii) computing a weighted score from both multiple crops and (resized) full images.
- *Multiclass training* - investigating whether training the detection model on a multiclass problem using generator labels improves binary detection performance.
  - *Multiclass to binary* - strategies to fuse multiclass scores into a binary prediction.
  - *Multiclass and Binary training* - assessing the benefits of training the model using both the multiclass and binary losses with different multi-head approaches.
  - *MLP vs distance-based approach* - exploring whether replacing the MLP classification head with a distance-to-centroid scoring function improves evaluation-time robustness after multiclass training.

Following the AI-GenBench protocol, we adopt the AUROC on the *Next Period* (averaged across all steps) as the primary

---

[1]laion/CLIP-ViT-L-14-CommonPool.XL-s13B-b90K

evaluation metric because it captures the detector's ability to generalize to future, unseen generators.

## A. Data augmentation pipeline

Data augmentation plays a central role in enhancing the robustness of deepfake detection models, as it helps detectors generalize across different synthetic image generators and remain resilient to realistic image corruptions [6], [38], [39]. To systematically assess its impact, we evaluate three distinct training-time augmentation pipelines, while at evaluation time all models are tested using the mandatory AI-GenBench preprocessing pipeline.

*Baseline pipeline:* The baseline pipeline is identical to the default augmentation strategy used in the AI-GenBench paper and initially proposed by Corvi et al. in [30]. It applies relatively strong transformations in a probabilistic manner, including random resized cropping, color jitter, grayscale conversion, dropout, Gaussian noise, blurring, random rotations, and horizontal flipping. A single JPEG compression pass is also applied with quality uniformly sampled from the range $[30, 100]$. This configuration aims to improve generalization by exposing the detector to a wide spectrum of perturbations.

*Evaluation-based pipeline:* The second training pipeline is derived from the AI-GenBench evaluation pipeline, which was originally designed to simulate realistic degradations caused by upload, download, and re-encoding processes on social media or messaging platforms. This pipeline applies up to three successive JPEG compression passes with variable quality levels, combined with a softer set of augmentations compared to the baseline. For training purposes, we extend this pipeline by adding random horizontal flipping and random rotation. This design allows us to test whether training with more realistic and less aggressive augmentations can improve generalization while preserving robustness to real-world corruptions.

*Mild pipeline:* The third training pipeline is also derived from the AI-GenBench evaluation pipeline but, similar to the baseline pipeline, applies only a single final JPEG compression pass. The purpose of this intermediate configuration is to isolate the effect of repeated JPEG compression during training and determine whether multiple compression passes provide additional benefits compared to a simpler single-pass strategy.

At evaluation time, all models are tested exclusively using the mandatory AI-GenBench evaluation pipeline, independently of the training pipeline used.

## B. Augmentation multiplier and training duration

The augmentation multiplier ($am$) is a key hyperparameter introduced in the AI-GenBench framework to control the diversity of augmented images during training. Specifically, $am$ determines the number of unique augmented variants generated for each training image through deterministic augmentations. For instance, with $am = 4$ (the default setting in AI-GenBench), the effective size of the training dataset becomes $4 \times |D|$, where $|D|$ denotes the number of original training images.

In the original AI-GenBench setup, training is performed for a single epoch with $am = 4$, meaning the model sees exactly four distinct augmentations of each image. In our evaluation, we extend this analysis along two axes:

- *Varying augmentation multiplier* - we vary $am$ in the range $[1, 12]$ while keeping the number of epochs fixed at one. This isolates the effect of increasing augmentation diversity within a single pass over the dataset.
- *Varying number of epochs* - we vary the number of epochs in the range $[1, 4]$ while fixing $am = 4$. This isolates the effect of repeated passes over the augmented dataset while keeping augmentation diversity constant.

This setup enables us to disentangle the contribution of dataset expansion through augmentation from that of extended training duration, and to determine whether one or both factors are required to achieve optimal generalization performance.

## C. Input processing at training and inference time

An important design choice for deepfake detection models concerns how input images are processed before being fed to the backbone. At training time, we consider two main strategies:

- *Random crop* - a sub-region of the image is randomly cropped to match the model's input resolution. This strategy encourages the model to rely on fine-grained local artifacts and noise patterns that may reveal synthetic content.
- *Resize* - the entire image is resized to the model's input resolution, preserving global context. This allows the model to focus on semantic consistency and macroscopic distortions rather than local noise. However, severe downsizing may suppress subtle forensic cues.

At evaluation time, random crops are replaced by deterministic procedures:

- *Central crop or multi-crop* - either a single central crop or multiple crops followed by score fusion, approximating the training distribution of crop-based models.
- *Resize* - the full image is resized, mirroring the training setup of resize-based models.

In the original AI-GenBench paper, crop-trained models were evaluated using the multi-crop strategy (with single-crop also tested but found inferior), while resize-trained models were evaluated on resized images. Their findings suggest that the resize strategy generally yields superior performance.

We extend this analysis by evaluating both crop- and resize-trained models under both inference protocols. Specifically, for each trained model we generate predictions from multiple crops and from the resized image, then fuse the scores with equal weight. Importantly, this *Mixed* evaluation is applied separately to each model type: a crop-trained detector is never combined with a resize-trained detector. This setup allows us to test whether jointly leveraging local (crop-based) and global (resize-based) evidence at inference time improves robustness compared to relying on a single strategy.

## D. Multiclass training

Deepfake detection is typically formulated as a binary classification task: real versus synthetic. However, since training data often includes generator-specific labels, it is natural to ask whether reframing the problem as a multiclass task (real + $N$ generator classes) can improve binary detection performance. We investigate this question by evaluating both pure multiclass training and joint multiclass–binary approaches.

*Plain multiclass training:* In the first setting, we train the detector solely on the multiclass task, using the generator identity as the label. At evaluation time, the multiclass outputs must be converted into a binary prediction. We consider two fusion strategies:

- *Sum fusion* - the binary *fake* score is computed by summing the softmax scores of all fake classes (with class 0 representing the "real" class) encountered during training.
- *Max fusion* - the *fake* score is defined as the maximum softmax score among all fake classes encountered during training.

*Dual-head training:* In the second setting, we jointly train the model on both binary and multiclass objectives to assess whether generator-aware supervision can provide more discriminative features and whether enforcing both tasks jointly improves the final binary detection performance. To this end, we equip the backbone with two output heads and employ a combined loss retaining only the binary head at evaluation time. We experiment with two architectural variants:

- *Separate heads* - the backbone features two separate heads, one for binary classification and one for multiclass classification, trained simultaneously.
- *Stacked heads* - the binary head is placed on top of the multiclass head; specifically, multiclass logits are passed through a ReLU and then projected onto a binary prediction.

For both variants, we explore different loss weightings: equal weighting (0.5 each) and an asymmetric configuration where the binary loss dominates (0.75 binary, 0.25 multiclass), treating multiclass supervision as an auxiliary signal.

*MLP vs distance-based approach:* In addition to standard multiclass-to-binary fusion, we explore a distance-based alternative to the usual MLP classification head. The procedure consists of five steps:

1) *Training* - the model is trained using the plain multiclass setup described above, without any dual-head architecture or binary loss.
2) *Centroid extraction* - for each class (generator), we extract $c$ centroids from the training set, with $c \in [1, 3]$, to assess whether multiple centroids improve performance. Centroids are computed in the feature space immediately before the final classification layer. For $c > 1$, centroids are obtained by clustering training patterns using K-Means.
3) *Distance scoring* - for a test image, we compute its distance to each centroid of every class and convert this into an inverse-distance score, so that closer proximity corresponds to higher confidence.

TABLE I: Performance (average Next Period AUROC, %) of different augmentation pipelines across detector backbones. The best result per backbone is highlighted in bold.

| Pipeline | DINOv2 | ViT-L CLIP | ResNet-50 CLIP | Average |
|---|---|---|---|---|
| Baseline | 94.2 | 90.8 | 85.2 | 90.1 |
| Evaluation | **95.7** | **94.6** | **93.2** | **94.5** |
| Mild | 95.4 | 93.5 | 90.2 | 93.1 |

4) *Class-level aggregation* - for each class, we retain the maximum inverse-distance score among its centroids.
5) *Binary prediction* - to produce a real/fake score, we apply a fusion strategy across the $N$ fake classes, analogous to the multiclass-to-binary approaches:
   - *Sum fusion* - sum the scores of all fake classes.
   - *Max fusion* - take the maximum score among all fake classes.

## E. Baseline

All results are reported relative to a *Baseline* configuration. This setup follows the procedure proposed in the initial AI-GenBench experiments and consists of using the *Baseline pipeline* for data augmentation, training for a single epoch with an augmentation multiplier of $am = 4$, resizing the entire image to the model's input size for both training and evaluation, and directly optimizing the binary classification objective (i.e., using only the binary loss).

## IV. RESULTS

### A. Data augmentation pipeline

Figure 2 and Table I report the performance of the three training-time augmentation pipelines across all considered backbones. While the *baseline* pipeline, characterized by heavy perturbations and a single JPEG compression pass, achieves competitive results, we observe that the *evaluation-based* pipeline (which applies up to three JPEG compression passes combined with milder augmentations) consistently leads to higher AUROC scores on the Next Period metric (90.1% vs. 94.5% on average).

The *mild* pipeline, derived from the evaluation-based strategy but restricted to a single JPEG compression pass, achieves intermediate performance (93.1%), suggesting that repeated compression during training plays an important role in preparing detectors for real-world degradations.

Overall, these results indicate that:

- while data augmentation is critical for robust detection, excessively strong augmentations (as in the baseline pipeline) may be counterproductive;
- augmentations that closely mimic realistic post-processing operations encountered in-the-wild provide more consistent improvements;
- introducing repeated JPEG compression passes during training effectively improves the generalization capabilities. These trends are observed across all three detector architectures, underscoring the generality of our findings.
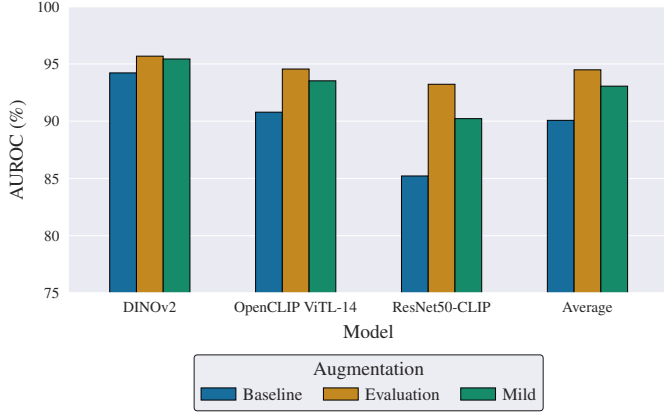
Fig. 2: Impact of different augmentation pipelines on Next Period AUROC. Results are shown for all three detector backbones.

### B. Augmentation Multiplier and Training Duration

Figure 3 and Tables II and III summarize the results obtained by varying the augmentation multiplier ($am$) in the range $[1, 16]$ while fixing the number of epochs to 1, and by varying the number of epochs in $[1, 4]$ while fixing $am = 4$, as in the AI-GenBench paper.

We observe that larger models, such as ViT-L CLIP and DINOv2, reach a performance plateau more quickly than smaller backbones like ResNet-50 CLIP. In particular, ResNet-50 CLIP continues to benefit from longer training schedules, whereas transformer-based models converge after only one or two epochs. When comparing the effect of increasing $am$ versus increasing the number of epochs, the two strategies appear equivalent. For example, configurations $am = 8$, *epochs=1* and $am = 4$, *epochs=2* yield similar AUROC values (96.36% vs 96.38% for DINOv2). This suggests that dataset expansion through augmentation and repeated exposure to the same augmented samples both provide sufficient "fuel" for training, with no clear advantage of one approach over the other. Overall, these results indicate that, while the augmentation multiplier can effectively replace longer training schedules, small-capacity models may still benefit from additional epochs before reaching their performance ceiling. Finally, it is worth noting that, while increasing the augmentation multiplier $am$ could be a reasonable choice for practical deployments, the AI-GenBench fairness rules prohibit values above $am = 4$ to constrain training data diversity and ensure comparability across approaches. In contrast, increasing the number of epochs is allowed.

### C. Input processing at training and inference time

The AI-GenBench paper established a strong baseline where models are trained on resized images (downsized to the model's input resolution) and evaluated in the same way. In their study, this resize-based setup outperformed an alternative configuration where models were trained on random crops and evaluated via multi-crop inference (with scores averaged across crops). Multi-crop inference, in turn, proved to be better than single (center)-crop inference.
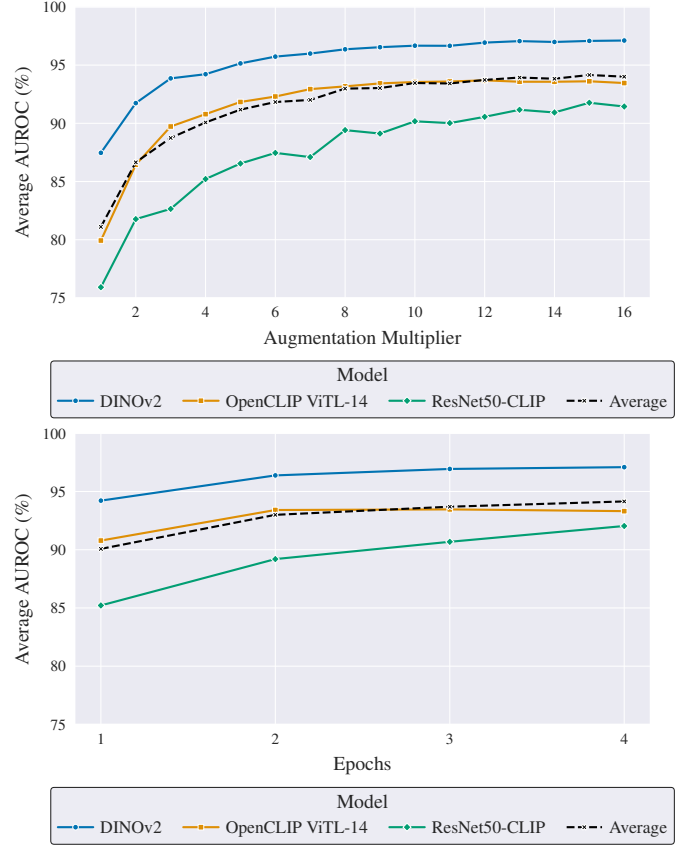


Fig. 3: Effect of varying the augmentation multiplier (top) and number of epochs (bottom) on generalization performance (average Next Period AUROC, %) across models.

TABLE II: Effect of augmentation multiplier ($am$) on the average Next Period AUROC (%). The best result per backbone is highlighted in bold.

| $am$ | DINOv2 | ViT-L CLIP | ResNet-50 CLIP | Average |
|------|--------|------------|----------------|---------|
| 1  | 87.47 | 79.93 | 75.92 | 81.11 |
| 2  | 91.73 | 86.44 | 81.77 | 86.65 |
| 3  | 93.86 | 89.73 | 82.65 | 88.74 |
| 4  | 94.22 | 90.79 | 85.21 | 90.07 |
| 5  | 95.14 | 91.83 | 86.54 | 91.17 |
| 6  | 95.73 | 92.30 | 87.46 | 91.83 |
| 7  | 95.99 | 92.93 | 87.10 | 92.01 |
| 8  | 96.36 | 93.17 | 89.41 | 92.98 |
| 9  | 96.54 | 93.43 | 89.12 | 93.03 |
| 10 | 96.67 | 93.54 | 90.17 | 93.46 |
| 11 | 96.66 | 93.59 | 90.02 | 93.42 |
| 12 | 96.93 | **93.69** | 90.56 | 93.73 |
| 13 | 97.06 | 93.57 | 91.16 | 93.93 |
| 14 | 96.98 | 93.57 | 90.93 | 93.83 |
| 15 | 97.08 | 93.61 | **91.77** | **94.15** |
| 16 | **97.11** | 93.46 | 91.44 | 94.00 |

In our work, we extend this comparison by introducing a hybrid evaluation strategy that combines both resized and cropped inputs. Specifically, at inference time, we generate five crops per image, average their prediction scores, and then combine this crop-based score with the resized-image score using equal weighting ($w = 0.5$ for both). This *Mixed* evaluation strategy is applied separately to both resize-trained

TABLE III: Effect of training duration (epochs) on the average Next Period AUROC (%). The best result per backbone is highlighted in bold.

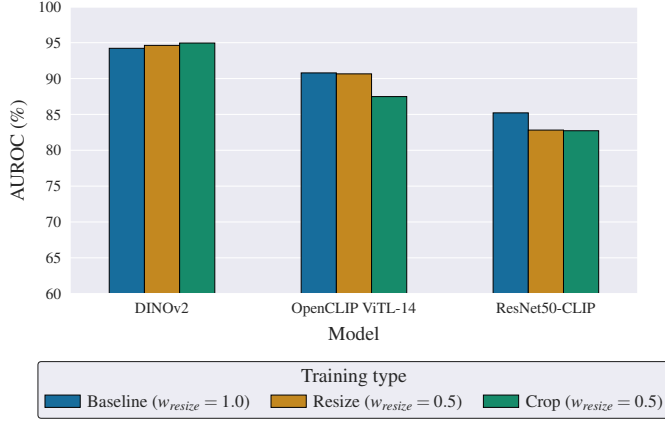| Epochs | DINOv2 | ViT-L CLIP | ResNet-50 CLIP | Average |
|---|---|---|---|---|
| 1 | 94.22 | 90.79 | 85.21 | 90.07 |
| 2 | 96.38 | 93.41 | 89.20 | 92.99 |
| 3 | 96.94 | **93.46** | 90.68 | 93.69 |
| 4 | **97.09** | 93.31 | **92.04** | **94.15** |



Fig. 4: Comparison of training and inference input processing strategies across backbones on the average Next Period AUROC (%). In the *Baseline* approach, the model is trained and evaluated on the resized version of images. In *Resize*, the model is trained on resized images but, at evaluation time, both the resized images and five of their crops are considered (fusing scores with equal weight between the resized image and the multi-crops). The *Crop* approach follows the same evaluation mechanism, but the model is trained on crops only.

and crop-trained models.

Figure 4 and Table IV summarize the results. The relative effectiveness of each strategy depends on the backbone architecture:

- ResNet-50 CLIP - the resize-only baseline remains superior, and mixed evaluation does not provide improvements.
- ViT-L CLIP - crop-trained models underperform, while resize-trained models with mixed evaluation achieve results comparable to the resize-only baseline.
- DINOv2 - the ranking changes, crop-trained models with mixed evaluation achieve the best performance, followed by resize-trained mixed evaluation, with the resize-only baseline performing worst.

These findings suggest that the optimal input processing strategy is architecture-dependent. Larger backbones such as DINOv2 appear to benefit from incorporating multi-crop information, whereas smaller backbones like ResNet-50 are more stable when trained and evaluated solely on resized images.

### D. Multiclass training

An open question in deepfake detection is whether exploiting generator labels during training can improve binary

TABLE IV: Comparison of training and evaluation input strategies (average Next Period AUROC %). Here *Mixed* refers to the fusion of scores obtained from both crops and resized version of the full image. The first row corresponds to the *Baseline* model. The best result per backbone is highlighted in bold.

| Training input | Evaluation input | DINOv2 | ViT-L CLIP | ResNet-50 CLIP |
|---|---|---|---|---|
| Resize | Resize | 94.22 | **90.79** | **85.21** |
| Resize | Mixed | 94.63 | 90.65 | 82.81 |
| Crop | Mixed | **94.95** | 87.48 | 82.72 |

TABLE V: Plain multiclass training: comparison of fusion strategies (average Next Period AUROC %). Baseline refers to direct binary training. The best result per backbone is highlighted in bold.

| Fusion strategy | DINOv2 | ViT-L CLIP | ResNet-50 CLIP |
|---|---|---|---|
| Baseline | **94.22** | **90.79** | **85.21** |
| Sum fusion | 92.87 | 84.55 | 82.90 |
| Max fusion | 92.31 | 84.55 | 82.64 |

classification performance. While the task is typically formulated as a binary problem (real vs. synthetic), it can also be reframed as a multiclass problem (one real class plus one class per generator). The key challenge then becomes how to map multiclass outputs to a single binary prediction at evaluation time. To investigate this, we evaluate three strategies: *i)* plain multiclass training with fusion to binary, *ii)* dual-head training combining binary and multiclass losses, and *iii)* a distance-based approach using class centroids.

*1) Plain multiclass training:* In the first setup, models are trained to predict generator identities directly. At inference time, the multiclass outputs are mapped to a binary decision using either *Sum fusion* or *Max fusion*.

Figure 5 and Table V report the results. Direct binary training (the *baseline*) outperforms plain multiclass training followed by fusion across all models. The gap is especially pronounced for the ViT-L CLIP backbone, where binary training yields a significantly higher AUROC score ($90.79\%$ vs. $84.55\%$). For DINOv2 and ResNet-50 CLIP, the difference is smaller.

When comparing fusion strategies, *sum fusion* performs slightly better than *max fusion*, though neither closes the gap with binary training. These results suggest that while generator-specific supervision encourages richer representations, simply collapsing them into a binary decision at inference time is less effective than directly optimizing for the binary objective.

*2) Dual-head training:* Starting from the previous observation, we evaluated whether combining binary and multiclass supervision via a dual-head architecture can improve detection performance.

For each configuration (*Separate heads* and *Stacked heads*), we explored two loss-weighting schemes: *i)* equal weighting ($w_{\text{bin}} = w_{\text{multi}} = 0.5$), and *ii)* auxiliary weighting, where the

Fig. 5: Plain multiclass training: comparison of fusion strategies (sum vs. max) against binary baseline on the average Next Period AUROC (%).

TABLE VI: Dual-head training results (average Next Period AUROC %). *Aux* denotes down-weighting the multiclass loss ($w_{bin} = 0.75, w_{multi} = 0.25$). Baseline refers to training using a binary loss only. The best result per backbone is highlighted in bold.

| Configuration | DINOv2 | ViT-L CLIP | ResNet-50 CLIP |
|---|---|---|---|
| Baseline | **94.22** | 90.79 | **85.21** |
| Separate heads | 93.15 | 89.03 | 82.99 |
| Stacked heads | 91.17 | 84.95 | 77.70 |
| Separate heads (aux) | 94.21 | **92.35** | 84.09 |
| Stacked heads (aux) | 93.66 | 86.74 | 81.65 |

binary loss dominates ($w_{\text{bin}} = 0.75$, $w_{\text{multi}} = 0.25$) treating the multiclass signal as an auxiliary objective.

Figure 6 and Table VI report the results of the four configurations and the baseline. Several consistent patterns emerge:

- The *Separate heads* approach consistently outperforms the *Stacked heads* approach across all backbones.
- Using the multiclass loss as an auxiliary signal (*aux*) is superior to equal weighting, consistently across both dual-head approaches and all backbones.
- Among the four dual-head combinations, separate heads with auxiliary weighting achieves the best performance.

When compared to the pure binary baseline, this best dual-head strategy is competitive for DINOv2 (94.21% vs. 94.22%), significantly superior for ViT-L CLIP (92.35% vs. 90.79%), and slightly inferior for ResNet-50 CLIP (84.09% vs. 85.21%). These results suggest that employing an auxiliary supervision based on the generator label can benefit larger transformer-based detectors. Fine-tuning the relative loss weights may yield further improvements, but this is beyond the scope of the present study.

*3) MLP vs distance-based approach:* Figure 7 and Table VII summarize the results for the distance-based approach. Across all backbones, the *baseline* MLP trained directly on the binary task remains superior. Compared to the plain multiclass setup with fusion (see Table V), even the best centroid configuration underperforms on DINOv2 and ResNet-50 CLIP,
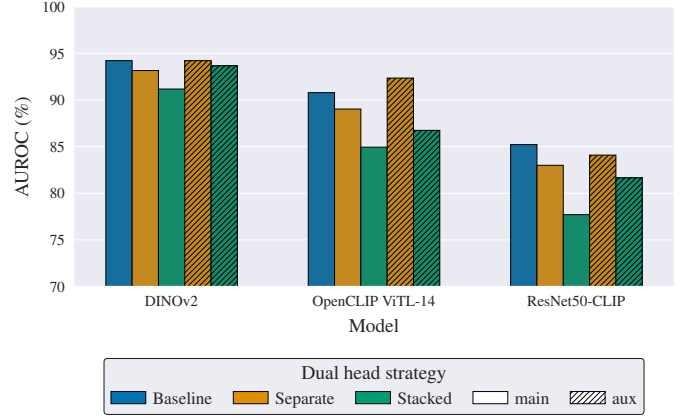


Fig. 6: Dual-head training: performance comparison (average Next Period AUROC (%)) of head configurations (separate vs. stacked) and loss weighting schemes (equal vs. *aux*iliary) across backbones.
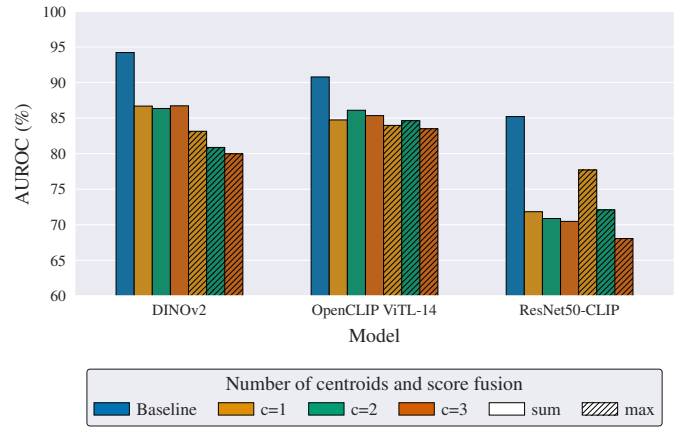


Fig. 7: Distance-based scoring: effect of the number of centroids per class ($c$) and fusion strategy (sum vs. max) across backbones. Metric: Next Period AUROC (%).

while for ViT-L CLIP shows only a slight improvement over the plain multiclass approach. However, it still falls notably short of the binary baseline.

Regarding fusion strategies within the distance-based approach, *sum fusion* tends to outperform *max fusion* on the larger backbones (DINOv2 and ViT-L CLIP), whereas *max fusion* performs relatively better for ResNet-50 CLIP (though still far below the MLP baseline). Overall, these mixed outcomes indicate that centroid-based scoring is generally less effective and less reliable than a standard MLP head trained on the binary task.

## V. INCREMENTAL UPDATE STRATEGY

An additional dimension in our study concerns how detectors can be efficiently and effectively updated as new generators are released over time. In the AI-GenBench evaluation framework, this corresponds to progressing through successive temporal windows, when (four) new generators are introduced at each step.

TABLE VII: Distance-based scoring vs. MLP baseline (average Next Period AUROC %). The best result per backbone is highlighted in bold.

| Centroids | Fusion | DINOv2 | ViT-L CLIP | ResNet-50 CLIP |
|-----------|--------|--------|------------|----------------|
| Baseline | – | **94.22** | **90.79** | **85.21** |
| 1 | Sum | 86.69 | 84.74 | 71.84 |
| 2 | Sum | 86.35 | 86.10 | 70.88 |
| 3 | Sum | 86.73 | 85.35 | 70.48 |
| 1 | Max | 83.15 | 83.98 | 77.71 |
| 2 | Max | 80.88 | 84.63 | 72.10 |
| 3 | Max | 79.98 | 83.50 | 68.06 |

*Baseline (batch tuning):* The standard setting, introduced in the AI-GenBench paper [4], consists of successive *batch tuning* steps on all data accumulated up to the current window, starting from the model weights obtained up to that moment. This strategy has two consequences: *i)* the model may benefit from the fact that past generators were already learned, and *ii)* past and new generators are equally represented in the cumulative training set. While effective, this approach is computationally expensive, as it requires retraining on the entire generator history at each step.

*Reset weights:* As a reference, we consider a variant in which, at each window, the detector is retrained on the cumulative data but initialized from the original pretrained weights rather than from the model obtained at the previous step. This design removes the bias in favor of earlier generators present in the baseline approach but discards previously consolidated knowledge that could be beneficial when adapting to new generators. In our experiments, we compare this strategy against both the baseline and continual learning approaches.

### A. Continual learning strategies

We explore several continual learning strategies aimed at balancing adaptability to new generators with retention of knowledge about older ones. These strategies are designed to be more efficient than the batch retraining baseline while mitigating catastrophic forgetting. In particular, we focus on replay-based strategies, which reduce forgetting by storing and replaying a subset of images from generators encountered in previous windows. We evaluate the performance of both a size-unbounded and a size-bounded replay strategy.

*Harmonic replay:* In this strategy the number of stored samples per generator decreases according to a harmonic schedule. Initially, all training samples for each generator are inserted into the replay buffer. Over time, the contribution of each generator is reduced by a factor of $1/i$, where $i$ is the number of windows elapsed since that generator was introduced. This allocation reserves more space for recent generators while gradually down-weighting older ones. We refer to this as a *Harmonic* schedule since the total buffer size grows without bound, following the harmonic series $N \cdot (1 + \frac{1}{2} + \frac{1}{3} + \dots)$.

*Class-balanced replay:* Here, we consider a bounded replay buffer of fixed size, maintained using a class-balanced strategy in which an equal number of examples is kept for all generators. This means that as new generators are introduced, the portion of the replay buffer allocated to earlier generators decreases. This strategy enforces a strict memory budget, making the computational cost directly proportional to the size of the buffer. To assess how performance scales with replay capacity, we evaluate different buffer sizes (e.g., 10000 and 20000 samples).

*Naive continual learning:* Finally, as a lower bound, we evaluate a naive continual tuning approach where the model is initialized with the weights from the previous step but fine-tuned only on data from generators in the current window, without any replay or regularization mechanism.

### B. Results

Results for Next Period and Past Period AUROC are summarized in Figures 8 and 9, with corresponding numerical values reported in Tables VIII and IX. As expected, the naive approach preserves adaptability to new generators but suffers from severe forgetting on older ones. Replay buffers substantially mitigate this effect, with both class-balanced and harmonic replay offering a favorable trade-off between buffer size and retention. Resetting the model weights generally performs worse than both the batch baseline and replay-based approaches, confirming that using previous weights is beneficial. Overall, continual learning with replay effectively balances adaptability and retention, while naive tuning alone is insufficient. Forgetting, as measured by the performance gap between the batch baseline and the continual learning approaches on the *Past Period*, is more pronounced for the smaller ResNet-50 CLIP model. This observation is consistent with recent findings in the literature showing that larger models, both in vision and language, exhibit greater resistance to forgetting [40].

Results in Table IX and Figure 9 confirm that both harmonic and class-balanced replay strategies greatly mitigate forgetting on generators introduced in past windows. Table X shows the computational impact of these approaches: replay strategies significantly reduce the computational cost of adapting models to new generators compared to full batch retraining. In addition, it should be noted that the harmonic strategy, which shows the best results, reduces the number of stored samples of generators over time, thus reducing their relevance at training time: generators become obsolete and content generated with those models is more easily recognizable. In addition, *Past Period* performance shows that detection capabilities for older generators can be retrained even with very few examples. Keeping these considerations in mind, decreasing the influence of older generators during training is both practical and natural.

### VI. INTEGRATING THE "BEST OF"

The main objective of this study was to identify design practices that consistently improve deepfake detection performance across different model architectures. By systematically evaluating the impact of individual training and inference-time choices, we aimed to isolate configuration elements that lead to the most transferable improvements.

Excluding continual learning experiments, which address a distinct adaptation scenario, our results highlight a configuration that achieves the most consistent performance across the
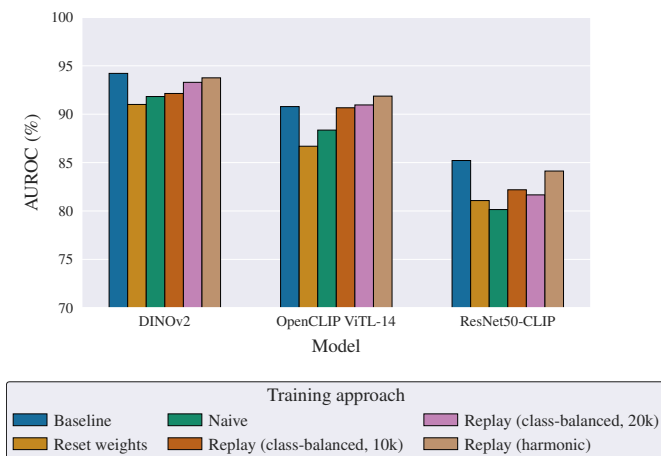
Fig. 8: Performance of different training regimes on the Next Period. Baseline refers to the standard approach of successive tuning on all the data of all sliding windows so far. *Reset weights* differs from baseline in which the model weights are reverted to their initial general pretraining. *Naive* and *Replay* refer to continual learning setups where the model is trained on the data of the current window only, either without protection techniques (Naive), or by counteracting forgetting by using replay data (Replay). Metric: Next Period AUROC (%).

TABLE VIII: Next Period performance of different training regimes (AUROC %). Baseline is shown separately as a reference, while the best continual learning result per backbone is highlighted in bold. CB stands for *class-balanced*, followed by the number of samples in the replay buffer.

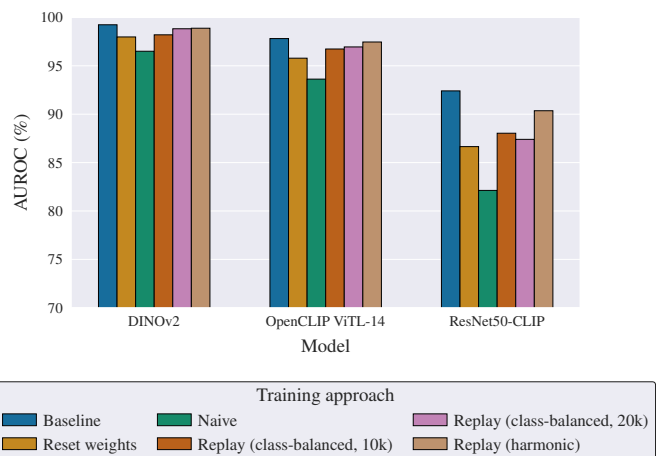| Configuration | DINOv2 | ViT-L CLIP | ResNet-50 CLIP |
|---|---|---|---|
| Baseline (reference) | 94.22 | 90.79 | 85.21 |
| Reset weights | 91.01 | 86.69 | 81.08 |
| Naive | 91.83 | 88.36 | 80.15 |
| Replay (CB, 10k) | 92.14 | 90.66 | 82.19 |
| Replay (CB, 20k) | 93.29 | 90.96 | 81.66 |
| Replay (harmonic) | **93.76** | **91.87** | **84.13** |



Fig. 9: Performance of different training regimes on the Past Period. Baseline refers to the standard approach of successive tuning on all the data of all sliding windows so far. *Reset weights* differs from baseline in which the model weights are reverted to their initial general pretraining. *Naive* and *Replay* refer to continual learning setups where the model is trained on the data of the current window only, either without protection techniques (Naive), or by counteracting forgetting by using replay data (Replay). Metric: Past Period AUROC (%).

TABLE IX: Past Period performance of different training regimes (AUROC %). Baseline is shown separately as a reference, while the best continual learning result per backbone is highlighted in bold. CB stands for *class-balanced*, followed by the number of samples in the replay buffer.

| Configuration | DINOv2 | ViT-L CLIP | ResNet-50 CLIP |
|---|---|---|---|
| Baseline (reference) | 99.24 | 97.81 | 92.41 |
| Reset weights | 97.98 | 95.79 | 86.65 |
| Naive | 96.50 | 93.62 | 82.13 |
| Replay (CB, 10k) | 98.20 | 96.74 | 88.03 |
| Replay (CB, 20k) | 98.83 | 96.95 | 87.40 |
| Replay (harmonic) | **98.88** | **97.46** | **90.36** |

three evaluated architectures. Specifically, while maintaining the augmentation multiplier at $am = 4$ in accordance with the AI-GenBench protocol, the best results were obtained with an extended training regimen of four epochs.

Among the tested augmentation pipelines, the *Evaluation-based* pipeline, featuring three probabilistic JPEG compression passes, proved to be the most effective. For input pre-processing, resizing the entire image to the model's input resolution emerged as the most reliable strategy across architectures, both during training and evaluation. Consequently, this approach was adopted for training the "best of" model. However, it should be noted that a hybrid strategy that combines full (resized) image training with prediction based on both the full image and a set of crops (five in our experiments) achieved comparable performance on larger models (DINOv2, VIT-L CLIP).

Direct optimization of the binary classification objective remains the most reliable approach. However, a dual-head configuration with an auxiliary head (and loss) jointly optimized on the multiclass generator labels achieves comparable results and offers the additional benefit of enabling *model attribution*.

When these optimal design choices were combined and applied to the *DINOv2* backbone, the highest-performing architecture among the evaluated models, the resulting configuration achieved an average AUROC of 97.36% on the *Next Period*, establishing the current state-of-the-art on AI-GenBench.

## VII. CONCLUSIONS

In this paper, we presented a comprehensive evaluation of the design factors that influence the performance and generalization of deepfake and AI-generated image detectors. Our analysis covered a wide range of training and inference choices, including augmentation strategies, preprocessing pipelines, training duration, multiclass supervision, and continual learning mechanisms.

Our findings reveal several principles that generalize across architectures. First, aligning the training distribution with

TABLE X: Relative computational cost per time window for different learning strategies, expressed as a percentage of the baseline cost (which increases linearly with each time window) and considering the AI-GenBench setup, where each sliding window carries 160000 new training examples. Lower values indicate lower computational burden. An estimation of the relative cost when training on the 25th window is also provided to simulate a longer 20-year benchmark. CB stands for class-balanced, followed by the number of samples in the replay buffer.

| Window | Baseline | Naive | CB (10k) | CB (20k) | Harmonic |
|--------|----------|-------|----------|----------|----------|
| 1 | 100% | 50.00% | 81.25% | 112.50% | 100.00% |
| 2 | 100% | 33.33% | 54.17% | 75.00% | 83.33% |
| 3 | 100% | 25.00% | 40.63% | 56.25% | 70.83% |
| 4 | 100% | 20.00% | 32.50% | 45.00% | 61.67% |
| 5 | 100% | 16.67% | 27.08% | 37.50% | 54.72% |
| 6 | 100% | 14.29% | 23.21% | 32.14% | 49.29% |
| 7 | 100% | 12.50% | 20.31% | 28.13% | 44.91% |
| 8 | 100% | 11.11% | 18.06% | 25.00% | 41.31% |
| **Average** | 100% | 22.86% | 37.15% | 51.44% | 63.26% |
| … | … | … | … | … | … |
| 24 (20 years) | 100% | 4.0% | 6.5% | 9.0% | 19.26% |

realistic degradation is beneficial: using the AI-GenBench *Evaluation* pipeline, which applies multiple JPEG compression passes, consistently outperforms more aggressive augmentation strategies. Second, training for four epochs while maintaining the standard augmentation multiplier ($am = 4$) offers an excellent performance–efficiency trade-off. Third, full-image resizing emerges as the most stable and reliable input processing strategy, outperforming crop-based or hybrid alternatives. Finally, direct optimization of the binary objective remains the most robust approach, although a dual-head configuration with an auxiliary multiclass loss can achieve comparable performance in larger models while enabling model attribution.

Beyond the standard batch setting, we examined how detectors can be periodically updated as new generative models become available. Our experiments demonstrate that the proposed *Harmonic replay* strategy achieves performance close to full retraining while significantly reducing computational costs and mitigating the influence of obsolete generators, making it a practical and scalable solution for maintaining detectors up to date with novel generative models.

By integrating the identified "best of" practices on the DINOv2 backbone, we obtained state-of-the-art performance on the AI-GenBench benchmark.

Future work will extend this systematic analysis to more challenging settings, such as inpainting and localized manipulations, to assess whether the validated design choices remain optimal in these complex scenarios.

## Acknowledgment

## References

[1] Z. Epstein et al., "Art and the science of generative AI: A deeper dive," *Science*, vol. 380, 2023.

[2] C. Barrett et al., *Identifying and Mitigating the Security Risks of Generative AI*. Now Foundations and Trends, 2024.

[3] L. Lin et al., "Detecting multimedia generated by large AI models: A survey," *arXiv preprint arXiv:2204.06125*, 2024.

[4] L. Pellegrini et al., "Ai-genbench: A new ongoing benchmark for ai-generated image detection," in *IJCNN*, 2025.

[5] D. Tariang, R. Corvi, D. Cozzolino, G. Poggi, K. Nagano, and L. Verdoliva, "Synthetic Image Verification in the Era of Generative AI: What Works and What Isn't There Yet," *IEEE Security & Privacy*, vol. 22, pp. 37–49, 2024.

[6] S.-Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros, "CNN-generated images are surprisingly easy to spot... for now," in *CVPR*, 2020, pp. 8695–8704.

[7] U. Ojha, Y. Li, and Y. J. Lee, "Towards universal fake image detectors that generalize across generative models," in *CVPR*, 2023, pp. 24 480–24 489.

[8] C. Koutlis and S. Papadopoulos, "Leveraging Representations from Intermediate Encoder-blocks for Synthetic Image Detection," in *ECCV*, 2024, pp. 394–411.

[9] A. Sarkar, H. Mai, A. Mahapatra, S. Lazebnik, D. A. Forsyth, and A. Bhattad, "Shadows Don't Lie and Lines Can't Bend! Generative Models don't know Projective Geometry... for now," in *CVPR*, 2024.

[10] L. Baraldi, F. Cocchi, M. Cornia, L. Baraldi, A. Nicolosi, and R. Cucchiara, "Contrasting Deepfakes Diffusion via Contrastive Learning and Global-Local Similarities," in *ECCV*, 2024, pp. 199–216.

[11] D. Boychev and R. Cholakov, "ImagiNet: A Multi-Content Dataset for Generalizable Synthetic Image Detection via Contrastive Learning," *arXiv preprint arXiv:2407.20020*, 2024.

[12] A. S. Rajan, U. Ojha, J. Schloesser, and Y. J. Lee, "On the effectiveness of dataset alignment for fake image detection," *arXiv preprint arXiv:2410.11835*, 2024.

[13] F. Guillaro, G. Zingarini, B. Usman, A. Sud, D. Cozzolino, and L. Verdoliva, "A Bias-Free Training Paradigm for More General AI-generated Image Detection," *arXiv preprint arXiv:2412.17671*, 2024.

[14] H. Wu, J. Zhou, and S. Zhang, "Generalizable synthetic image detection via language-guided contrastive learning," *arXiv preprint arXiv:2305.13800*, 2025.

[15] F. Laiti, B. Liberatori, T. D. Min, and E. Ricci, "Conditioned prompt-optimization for continual deepfake detection," in *ICPR*, 2024.

[16] C. Li et al., "A continual deepfake detection benchmark: Dataset, methods, and essentials," in *WACV*, 2023, pp. 1339–1349.

[17] J. Tian et al., "Dynamic Mixed-Prototype Model for Incremental Deepfake Detection," in *ACM MM*, 2024, pp. 8129–8138.

[18] D. C. Epstein, I. Jain, O. Wang, and R. Zhang, "Online Detection of AI-Generated Images," in *ICCV Workshops*, Oct. 2023, pp. 382–392.

[19] R. Corvi, D. Cozzolino, G. Poggi, K. Nagano, and L. Verdoliva, "Intriguing properties of synthetic images: From generative adversarial networks to diffusion models," in *CVPR Workshops*, 2023, pp. 973–982.

[20] Y. He et al., "ForgeryNet: A Versatile Benchmark for Comprehensive Forgery Analysis," in *CVPR*, 2021, pp. 4360–4369.

[21] Z. Chen et al., "DiffusionFace: Towards a comprehensive dataset for diffusion-based face forgery analysis," *arXiv preprint arXiv:2403.18471*, 2024.

[22] H. Cheng, Y. Guo, T. Wang, L. Nie, and M. Kankanhalli, "Diffusion facial forgery detection," in *ACM Multimedia*, 2024, pp. 5939–5948.

[23] Y. Wang, Z. Huang, and X. Hong, "Benchmarking deepart detection," *arXiv preprint arXiv:2302.14475*, 2023.

[24] M. A. Rahman, B. Paul, N. H. Sarker, Z. I. A. Hakim, and S. A. Fattah, "ArtiFact: A Large-Scale Dataset with Artificial and Factual Images for Generalizable and Robust Synthetic Image Detection," in *ICIP*, 2023, pp. 2200–2204.

[25] M. Zhu et al., "GenImage: A Million-Scale Benchmark for Detecting AI-Generated Image," *NeurIPS*, vol. 36, pp. 77 771–77 782, 2023.

[26] Y. Hong and J. Zhang, "WildFake: A Large-scale Challenging Dataset for AI-Generated Images Detection," *arXiv preprint arXiv:2402.11843*, 2024.

[27] M. Schinas and S. Papadopoulos, "SIDBench: A Python framework for reliably assessing synthetic image detection methods," in *ACM International Workshop on Multimedia AI against Disinformation*, 2024, pp. 55–64.

[28] D. Park, H. Na, and D. Choi, "Performance Comparison and Visualization of AI-Generated-Image Detection Methods," *IEEE Access*, 2024.

[29] Z. Lu et al., "Seeing is not always believing: Benchmarking human and model perception of ai-generated images," *NeurIPS*, vol. 36, 2024.

[30] R. Corvi, D. Cozzolino, G. Zingarini, G. Poggi, K. Nagano, and L. Verdoliva, "On the detection of synthetic images generated by diffusion models," in *ICASSP*, 2023, pp. 1–5.

[31] Q. Bammey, "Synthbuster: Towards detection of diffusion model generated images," *IEEE Open Journal of Signal Processing*, vol. 5, pp. 1–9, 2024.

[32] G. Cazenavette, A. Sud, T. Leung, and B. Usman, "FakeInversion: Learning to Detect Images from Unseen Text-to-Image Models by Inverting Stable Diffusion," in *CVPR*, Jun. 2024, pp. 10 759–10 769.

[33] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in ser. Psychology of Learning and Motivation, G. H. Bower, Ed., vol. 24, Academic Press, 1989, pp. 109–165.

[34] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental Classifier and Representation Learning," in *CVPR*, 2017, pp. 5533–5542.

[35] A. Chaudhry et al., "On Tiny Episodic Memories in Continual Learning," *arXiv preprint arXiv:1902.10486*, 2019.

[36] A. Radford et al., "Learning Transferable Visual Models From Natural Language Supervision," in *ICML*, 2021, pp. 8748–8763.

[37] M. Oquab et al., "DINOv2: Learning Robust Visual Features without Supervision," *Transactions on Machine Learning Research*, 2024, ISSN: 2835-8856.

[38] S. Mandelli, N. Bonettini, P. Bestagini, and S. Tubaro, "Detecting gan-generated images by orthogonal training of multiple cnns," in *ICIP*, 2022, pp. 3091–3095.

[39] D. Gragnaniello, D. Cozzolino, F. Marra, G. Poggi, and L. Verdoliva, "Are gan generated images easy to detect? a critical analysis of the state-of-the-art," in *ICME*, 2021, pp. 1–6.

[40] V. V. Ramasesh, A. Lewkowycz, and E. Dyer, "Effect of scale on catastrophic forgetting in neural networks," in *ICLR*, 2022.