
Memorization Detection Benchmark for Generative Image models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generative models in medical imaging offer significant potential for data augmenta-
2 tion and privacy preservation, but they also pose risks of patient data memorization.
3 This study presents a comprehensive, data-driven approach to evaluate and charac-
4 terize the memorization behavior of generative models. We systematically compare
5 various network architectures, loss functions, pretraining datasets, and distance
6 metrics to identify optimal configurations for detecting potential privacy concerns
7 in synthetic images. Our analysis reveals that self-supervised contrastive networks
8 using Triplet Margin loss in models like DinoV2, DenseNet121, and ResNet50,
9 when paired with Bray-Curtis or Standardized Euclidean distance metrics, demon-
10 strate superior performance in detecting augmented copies of training images. We
11 further apply our methodology to characterize the memorization behavior of a con-
12 ditional diffusion image transformer model trained on mammography data. This
13 work contributes a robust framework for evaluating generative models in medical
14 imaging, offering a crucial tool for assessing the risk of patient data leakage in
15 synthetic datasets.¹

16 1 Introduction

17 The advent of generative models has a lot of potential in healthcare and medical imaging initiatives,
18 promising enhanced data sharing, expanded datasets, and improved training data diversity [1].
19 However, these advancements come with significant privacy implications, especially given the
20 sensitive nature of patient information. A key concern is the phenomenon of model memorization
21 [2, 3], where generative models inadvertently reproduce specific details from their training data,
22 potentially compromising patient confidentiality and undermining the core purpose of synthetic data
23 generation.

24 Recent research has demonstrated that a wide range of generative models, including GANs, VAEs,
25 and diffusion models, are vulnerable to memorization [4, 5, 6, 7]. Of particular note, diffusion models
26 [8], despite their impressive image quality, have shown a higher propensity for memorization [6].
27 This finding underscores the intricate interplay between model sophistication, output quality, and data
28 privacy. Furthermore, conventional evaluation metrics such as Inception Score (IS) [9] and Fréchet
29 Inception Distance (FID) [10] fall short in detecting these memorization issues, potentially masking
30 critical privacy vulnerabilities in emerging image generation techniques.

31 A common misconception is that memorization can be effectively addressed by simply monitoring
32 validation errors and preventing overfitting. However, this approach overlooks the fundamental
33 differences between these two phenomena [11]. While overfitting manifests as a global issue where

¹The code for this study is available at <https://github.com/molinamarcvdb/ImageFeatureExtractionBenchmark>

34 models excel on training data at the expense of generalization, memorization is a more nuanced
35 problem. It involves the model assigning disproportionately high probabilities to specific training
36 instances. Intriguingly, a model’s tendency to memorize can actually increase even as its validation
37 performance improves, particularly during the initial stages of training [11]. This paradoxical
38 relationship highlights the need for specialized strategies to identify and mitigate memorization,
39 distinct from traditional overfitting prevention techniques.

40 Our research builds upon recent advances in self-supervised contrastive learning for memorization
41 detection [5, 12], offering a comprehensive benchmark. We propose a novel approach to evaluate
42 the efficacy and resilience of self-supervised networks through systematic image augmentations.
43 Our study compares the performance of various state-of-the-art pretrained network architectures,
44 including ResNet50 [13] and DinoV2 [14]. We also investigate the influence of different loss
45 functions, including distance-based and entropy-based formulations, and examine the impact of
46 pretraining on natural versus medical image datasets. By comparing a range of similarity, distance,
47 and information-theoretic metrics, we aim to identify the most sensitive indicators for detecting and
48 characterizing training data memorization. To demonstrate the practical application of our findings,
49 we employ the best-performing method to analyze the memorization patterns in a diffusion model.

50 **2 Related work**

51 **2.1 Model Memorization**

52 The phenomenon of model memorization has been extensively studied in machine learning, particu-
53 larly in supervised learning contexts. Neural networks have demonstrated the capacity to memorize
54 entire datasets, including those with random labels [15]. This memorization is not uniform across all
55 data points; outliers and mislabeled samples are more likely to be memorized [16]. Memorization and
56 generalization might also depend on network architecture and optimization procedure, but also on the
57 data itself [17]. Moreover, some level of memorization in supervised learning has been shown to be
58 important for generalization in several standard benchmarks [18]. In generative models, memorization
59 presents unique challenges, as models that closely replicate training data may still achieve favorable
60 scores on standard quality and diversity metrics. Recent work has demonstrated that GANs, VAEs,
61 and diffusion models as well as vision language models are all susceptible to memorizing training
62 data [4, 5, 6, 7, 19]. Therefore, creating a memorization metric to be monitored during training would
63 enable a more comprehensive assessment of the generative model performance.

64 **2.2 Memorization Detection Methods**

65 Various approaches have been proposed to detect and quantify memorization in generative image
66 models. Correlation-based methods, such as the structural similarity index measure (SSIM) employed
67 by [20, 21, 22], offer a straightforward approach to assessing similarity between generated and training
68 images. However, these methods were initially developed to measure diversity not memorization
69 behaviour, and may be sub-optimal to detect generated samples which are mere augmented versions
70 of the training data (e.g., rotation or flipping).

71 More sophisticated approaches leverage self-supervised learning and contrastive methods. In [5] the
72 authors introduced a framework that uses contrastive learning to map images to a lower-dimensional
73 embedding space, allowing for the detection of copies that may include rotated or flipped variants of
74 training images. This method was further explored in [12], which investigated the effects of various
75 hyperparameters and training setups on memorization as well as mitigation strategies.

76 **2.3 Mitigation Strategies**

77 Various approaches have been proposed to mitigate memorization in generative models. These
78 include using exclusively augmented images during training [5], implementing Differentially Private
79 Stochastic Gradient Descent (DP-SGD) [23], and applying standard regularization techniques like
80 dropout and weight decay. Additionally, novel methods such as Privacy Distillation have been
81 introduced [24]. This two-step approach involves training an initial diffusion model on real data,
82 generating and refining synthetic samples to exclude identifiable information, and then using these
83 refined samples to train a second model. This method aims to reduce re-identification risk while
84 maintaining downstream performance.

85 However, these mitigation strategies often involve trade-offs. DP-SGD can compromise image quality
 86 or lead to model divergence [25], while data augmentation may complicate similarity assessments
 87 between synthetic and original images. The Privacy Distillation approach, while promising, may
 88 result in reduced quality of the final synthetic samples. Finally, factors such as over-training, dataset
 89 size, and augmentation techniques also significantly influence memorization and should be carefully
 90 addressed [5, 6, 12].

91 3 Methods

92 3.1 Problem Formulation

93 Let $\mathcal{X} = \{x_1, \dots, x_N\}$ represent a set of N training images, $\mathcal{X}_v = \{v_1, \dots, v_K\}$ denote a set of K
 94 validation images, and $\mathcal{G} = \{g_1, \dots, g_M\}$ be a set of M generated images. We train a Self-Supervised
 95 Contrastive Network (SSCN) to learn an embedding function $\phi : \mathcal{I} \rightarrow \mathbb{R}^d$, where \mathcal{I} is the image
 96 space and d is the embedding dimension, by minimizing a contrastive loss function $L(\phi; \mathcal{X})$.

97 Given a similarity metric $s : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, we compute the similarity between training and generated
 98 images as $S(x, g) = s(\phi(x), \phi(g))$ for $x \in \mathcal{X}, g \in \mathcal{G}$, and baseline similarities between training and
 99 validation images as $S_{base}(x, v) = s(\phi(x), \phi(v))$ for $x \in \mathcal{X}, v \in \mathcal{X}_v$. To prevent memorization of
 100 synthetic data, we set a threshold τ as the p -th percentile of the S_{base} distribution.

101 For evaluation, we define a set of severely augmented images $\mathcal{X}_a = \{a_1, \dots, a_L\}$, where each a_i
 102 is derived from \mathcal{X} using strong augmentations. We monitor the percentage of augmented images that
 103 match their corresponding original images in \mathcal{X} according to the similarity threshold τ .

104 3.2 Self-Supervised Contrastive Network

105 3.2.1 Architecture

106 The SSCN comprises a backbone network $f_\theta : \mathcal{I} \rightarrow \mathbb{R}^d$, followed by a projection head $g_\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$.
 107 The backbone extracts features from the input images, while the projection head maps these features
 108 to a lower-dimensional embedding space. The complete network is represented as:

$$h_{\theta, \phi}(x) = g_\phi(f_\theta(x)) \quad (1)$$

109 We experiment with several backbone architectures, including ResNet50 [13], DenseNet121 [26],
 110 Inception V3 [27], CLIP Image Encoder [28], and DinoV2 [14]. The projection head is a linear layer
 111 defined as $g_\phi(z) = Wz + b$, where $W \in \mathbb{R}^{k \times d}$ and $b \in \mathbb{R}^k$.

112 To explore the impact of domain-specific knowledge, we use backbones pretrained on both natural
 113 images (ImageNet [29]) and medical images (RadImageNet [30]). This comparison allows us to
 114 evaluate the transfer learning benefits of using medical-domain-specific pretraining.

115 3.2.2 Loss Functions

116 To structure the embedding space, we employ and compare two popular contrastive losses: the Triplet
 117 Margin Loss [31] and InfoNCE Loss[32]. Both losses aim to pull semantically similar data points
 118 closer while pushing dissimilar points farther apart.

119 **Triplet Margin Loss.** This loss function ensures that the distance between an anchor-positive pair
 120 is smaller than the distance between the anchor-negative pair, with a margin m . Specifically, for an
 121 anchor a , a positive example p , and a negative example n , the loss is defined as:

$$L_{\text{triplet}}(a, p, n) = \max(0, m + d(a, p) - d(a, n)) \quad (2)$$

122 where $d(\cdot, \cdot)$ is the Euclidean distance and m is the margin parameter. This encourages positive pairs
 123 to be closer together while keeping negatives farther apart in the embedding space.

124 **InfoNCE Loss.** InfoNCE (Information Noise-Contrastive Estimation) compares each anchor repre-
 125 sentation z_i with one positive sample z_j^+ and $N - 1$ negative samples $\{z_j^-\}$, using Cosine similarity
 126 between the embeddings. The objective is to maximize the probability that the positive pair is more
 127 similar than the negative ones. This probability is expressed as:

$$P(i|j) = \frac{\exp(s(z_i, z_j^+)/\tau)}{\exp(s(z_i, z_j^+)/\tau) + \sum_{z_j^-} \exp(s(z_i, z_j^-)/\tau)} \quad (3)$$

128 where τ is a temperature parameter that controls the smoothness of the distribution, and $s(z_i, z_j)$ is
 129 the Cosine similarity between anchor z_i and positive or negative samples.

130 The InfoNCE loss is computed as the negative log-likelihood of the positive pair:

$$L_{\text{InfoNCE}}(z_i, z_j^+, \{z_j^-\}) = -\log P(i|j) \quad (4)$$

131 3.2.3 Training Procedure

132 The training process is conducted over 100 epochs. For each epoch, mini-batches are sampled from
 133 the training set. Each batch undergoes a series of stochastic augmentations, including rotation, scaling,
 134 flipping, affine transformations, bias field distortion, gamma correction, noise addition, and blurring.
 135 These augmentations enhance the network’s ability to learn invariant features and generalize better.

136 The model computes embeddings for both the original and augmented batches, then calculates the
 137 loss (either Triplet or InfoNCE) based on these embeddings. Network parameters are updated using
 138 the AdamW optimizer with an initial learning rate of 10^{-4} , which is decayed exponentially with a
 139 factor of 0.99 after each epoch.

140 We implemented the model using PyTorch and distributed the training across two NVIDIA RTX 4090
 141 GPUs. A batch size of 128 was used for most experiments, except for CLIP and DinoV2 models,
 142 where it was reduced to 64 due to memory constraints. For the InfoNCE loss, we set the temperature
 143 $\tau = 0.5$, while for the triplet margin loss, we used a margin $m = 0.05$ with hard negative mining.
 144 The backbone was frozen during the first 5 epochs to ensure proper warm-up of the linear layer.

145 3.3 Embedding Similarity Analysis

146 To comprehensively evaluate the similarity between the learned embeddings, we employed and
 147 compared the following distance and similarity metrics: Bray-Curtis distance, Canberra distance,
 148 Chebyshev distance, City Block (Manhattan) distance, Correlation distance, Cosine similarity, Dice
 149 similarity coefficient, Euclidean distance, Jensen-Shannon divergence, Mahalanobis distance, Match-
 150 ing distance, Minkowski distance, Standardized Euclidean distance (SEuclidean), and Squared
 151 Euclidean distance.

152 3.3.1 Similarity Distributions

153 For each trained model, we compute the similarity metrics between the training set and its adversarial
 154 (augmented) counterpart, the validation set and its adversarial counterpart, and for baseline similarity
 155 level assessment between the training and validation sets. These result in a distribution of the highest
 156 similarity score for each image enabling to test whether the contrastive model is capable of detecting
 157 augmented image copies and assess quantitatively the memorization degree by comparing with the
 158 train-val distribution. When aggregating over networks, losses, pretrainign and/or metrics we report
 159 the mean validation (augmented) detection with error bars representing 95 % confidence intervals,
 160 and significance test are calculated using two-tailed t-test.

161 **Detection of Augmented Copies** To evaluate the effectiveness of our similarity metrics in identify-
 162 ing augmented copies, we implement a threshold-based detection method. Let $\mathcal{X}_{aug} = \{x'_1, \dots, x'_N\}$
 163 and $\mathcal{X}_{v,aug} = \{v'_1, \dots, v'_K\}$ represent the augmented versions of the training and validation sets,
 164 respectively. Given our similarity metric s and embedding function ϕ , we compute the similarity
 165 $S(x, x') = s(\phi(x), \phi(x'))$ between each original image $x \in \mathcal{X}$ and its augmented version $x' \in \mathcal{X}_{aug}$.

166 We flag x' as a potential copy if $S(x, x') > \tau$ for any $x \in \mathcal{X}$, where τ is set as the p -th percentile of
167 the baseline similarity distribution $S_{base}(x, v) = s(\phi(x), \phi(v))$ for $x \in \mathcal{X}, v \in \mathcal{X}_v$.

168 Our benchmark aims to detect all images in \mathcal{X}_{aug} and $\mathcal{X}_{v,aug}$ as copies of their original counterparts
169 when using a τ equal to the 5-th percentile of $S_{base}(x, v)$. By comparing the detection rates between
170 \mathcal{X}_{aug} and $\mathcal{X}_{v,aug}$, the model’s generalizability and robustness of our similarity metrics in identifying
171 augmented copies can be assessed.

172 3.4 Dataset

173 Our study utilized an anonymized X-ray mammography dataset comprising 7,184 scans from 1,718
174 unique patients. The images were obtained and stored in DICOM format with a median shape of
175 2800 x 2082 pixels and median spacing of 0.065 x 0.065 mm.

176 The dataset includes two primary classes of mammography scans: normal scans and scans with
177 calcification. To ensure the integrity of our evaluation, we performed a patient-aware train-validation
178 split, ensuring that scans from the same patient were not distributed across different sets.

179 For preprocessing, all images were resized to square resolutions. During model training, images
180 were further resized to match the natural input resolution of the backbone networks, typically 224 x
181 224 pixels. This dataset provides a robust foundation for training and evaluating our self-supervised
182 contrastive network and conditional diffusion model for medical image synthesis.

183 3.5 Conditional Diffusion Model for Medical Image Synthesis

184 To enhance our dataset and evaluate the potential of generative models in medical imaging, we trained
185 a class-conditional diffusion model using our medical imaging data. This model was designed to
186 generate high-quality, synthetic medical images while preserving class-specific features.

187 **Training Process** We utilized a Diffusion Image Transformer (DiT) architecture [33], specifically
188 the DiT XL/2 variant (670M), comprising 28 Transformer layers with a hidden size dimension of
189 1152 and 16 attention heads. The model, was initially pretrained on ImageNet and then fine-tuned
190 on our medical imaging dataset for 100,000 steps with a learning rate of 1e-4, batch size of 2, with
191 horizontal flip as the only augmentation.

192 **Inference and Dataset Augmentation** At inference time, we used the trained model to upsample
193 our original dataset, effectively doubling its size. The resulting images were later processed via
194 the best performing SSCN to showcase the usability of such privacy detector methods and their
195 memorization characterization performance.

196 4 Results

197 In this study, we evaluated the performance of various deep learning models for a detection task,
198 comparing different network architectures, pretraining datasets (ImageNet and RadImageNet), and
199 loss functions (InfoNCE and Triplet). Our results reveal significant variations in performance across
200 these factors, with some clear trends emerging.

201 4.1 Network, Pretraining and Loss

202 The performance of self-supervised networks varied significantly across different architectures, loss
203 functions, and pretraining datasets (Figure 1). Consistently across all network architectures, the
204 Triplet loss outperformed InfoNCE, often by a substantial margin. This superiority of Triplet loss
205 over InfoNCE was found to be statistically significant ($p < 0.05$) for all tested network architectures
206 and pretraining datasets, with many comparisons showing highly significant differences ($p < 0.001$).

207 When comparing the best configurations of different network architectures, several significant
208 differences emerged. DinoV2 with ImageNet pretraining and Triplet loss achieved the highest overall
209 performance (0.722), closely followed by DenseNet121 (0.710) and ResNet50 with RadImageNet
210 pretraining (0.660). The differences between these top-performing models were not statistically
211 significant ($p > 0.05$), suggesting that they perform comparably well.

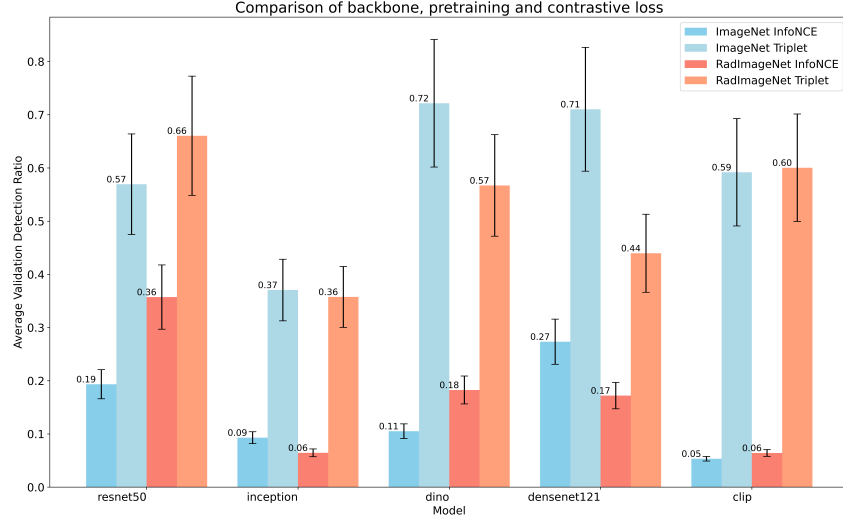


Figure 1: Comparison of network architectures performance with their best configurations.

212 However, significant differences were observed between the top-performing models and the Inception
 213 architecture. Inception, even in its best configuration (ImageNet, Triplet), performed significantly
 214 worse than ResNet50 ($p = 0.030$), DinoV2 ($p = 0.014$), and DenseNet121 ($p = 0.015$). The CLIP
 215 model, with its best configuration (RadImageNet, Triplet), showed intermediate performance (0.600)
 216 that was not significantly different from the top models but was marginally better than Inception
 217 ($p = 0.059$).

218 Interestingly, when focusing on the Triplet loss, the choice of pretraining dataset (ImageNet vs.
 219 RadImageNet) did not lead to statistically significant differences in performance for most architectures.
 220 This lack of significant difference in pretraining datasets for Triplet loss was consistent across all
 221 models, including ResNet50 ($p = 0.540$), CLIP ($p = 0.953$), Inception ($p = 0.875$), DinoV2
 222 ($p = 0.323$), and DenseNet121 ($p = 0.060$).

223 These findings indicate that while the choice of network architecture and loss function (Triplet
 224 vs. InfoNCE) has a significant impact on performance, the effect of pretraining dataset is more
 225 nuanced, particularly when using Triplet loss. The top-performing models (DinoV2, DenseNet121,
 226 and ResNet50) show comparable performance, significantly outperforming Inception, with CLIP
 227 falling in between. The robustness of Triplet loss to variations in pretraining data suggests it may
 228 offer more flexibility in the choice of pretraining dataset for self-supervised learning tasks across
 229 different network architectures.

230 4.2 Impact of Distance Metrics on Triplet Loss Performance

231 In addition to comparing network architectures and pretraining datasets, we also evaluated the
 232 performance of various distance metrics when using the Triplet loss function. The results, as
 233 illustrated in Figure 2, reveal substantial variations in performance across metrics, with the mean
 234 validation detection ratios and their respective confidence intervals showing clear differences.

235 The Bray-Curtis distance metric demonstrated the highest mean validation detection ratio of 0.8094
 236 (± 0.1036 CI), positioning it as the best performer. It was closely followed by the Jensen-Shannon
 237 divergence (0.7882 ± 0.1107 CI) and a group of Euclidean-based metrics, including Euclidean,
 238 Minkowski, and Squared Euclidean, which all achieved 0.7871 (± 0.1235 CI). These metrics consistently
 239 performed well across various configurations, highlighting their robustness when applied with
 240 models trained on Triplet Margin loss.

241 A slightly lower performance was observed with metrics such as the City Block (Manhattan) distance
 242 (0.7813 ± 0.1269 CI), the Canberra distance (0.7810 ± 0.1101 CI), and the Standardized Euclidean
 243 distance (0.7708 ± 0.1349 CI). Although these metrics exhibited detection ratios slightly below
 244 the top group, they still maintained strong performance, with detection ratios above 0.77. These

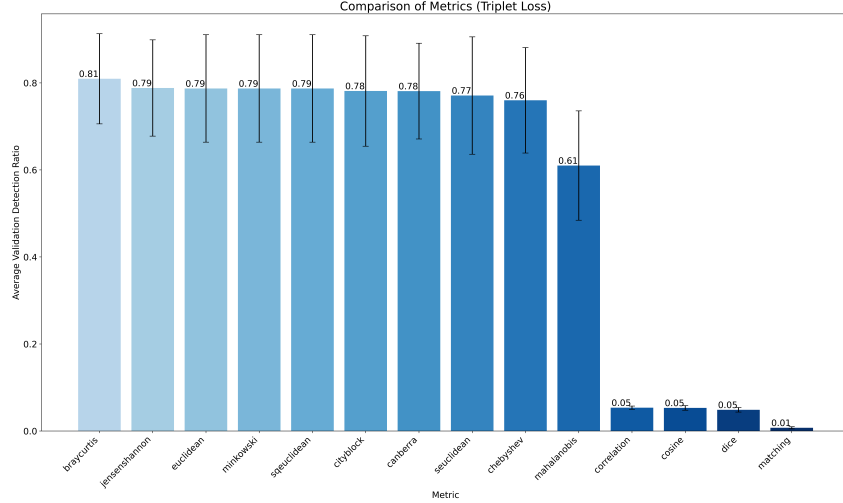


Figure 2: Comparison of distance metrics performance with Triplet loss in terms of mean validation detection ratio. Error bars represent confidence intervals.

245 results indicate that they are viable alternatives, particularly in situations where domain-specific
 246 considerations or computational efficiency play a role in metric selection.

247 On the other hand, the Chebyshev distance (0.7599 ± 0.1213 CI) and the Mahalanobis distance
 248 (0.6099 ± 0.1257 CI) displayed notably lower performance. The lower mean detection ratios for
 249 these metrics suggest that they may not be as effective in this task when paired with the Triplet
 250 loss function. Furthermore, the correlation-based metrics, including Correlation, Cosine, Dice, and
 251 Matching, performed significantly worse, with detection ratios falling below 0.06. Notably, the
 252 Matching distance exhibited extremely poor performance (0.0073 ± 0.0025 CI), suggesting that
 253 correlation-based metrics are ill-suited for this particular detection task when using Triplet loss.

254 The statistical analysis of pairwise comparisons further reinforced these findings. The differences
 255 between the top-performing metrics—Bray-Curtis, Jensen-Shannon, and Euclidean-based—were not
 256 statistically significant ($p > 0.05$), indicating that their performances are comparable. However, these
 257 top-performing metrics were significantly superior to the lower-performing and poor-performing
 258 metrics, with highly significant differences observed when compared to Mahalanobis and correlation-
 259 based metrics ($p < 0.001$).

260 4.3 Best Combinations for Each Network Architecture

261 We present the best-performing combinations of network architecture, pretraining dataset, loss
 262 function and metric. Table 1 highlights the maximum validation detection achieved and the distance
 263 metric that produced this maximum value for each network configuration.

Table 1: Best combinations for each network architecture

Model	Pretraining	Loss	Val. Detection	Metric
DinoV2	ImageNet	Triplet	0.9971	Bray-Curtis
DenseNet121	ImageNet	Triplet	0.9842	SEuclidean
ResNet50	RadImageNet	Triplet	0.9568	Bray-Curtis
ResNet50	ImageNet	Triplet	0.8863	Bray-Curtis
CLIP	RadImageNet	Triplet	0.8806	City Block
CLIP	ImageNet	Triplet	0.8791	Euclidean
DinoV2	RadImageNet	Triplet	0.8604	Euclidean
DenseNet121	RadImageNet	Triplet	0.7281	Bray-Curtis
Inception	ImageNet	Triplet	0.5813	Canberra
ResNet50	RadImageNet	InfoNCE	0.5496	Euclidean

264 As shown in Table 1, the DinoV2 model pre-trained on ImageNet using the Triplet loss achieved
 265 the highest validation detection score (0.9971), with the Bray-Curtis distance metric. Similar trends
 266 are observed across other architectures, with DenseNet121 and ResNet50 also performing well with
 267 SEuclidean and Bray-Curtis metrics, respectively.

268 4.4 Memorization Characterization of Diffusion Models

269 Using the best-performing combinations identified for our dataset, the fine-tuned DinoV2 model was
 270 employed to analyze the memorization behavior of a DiT trained to generate synthetic mammography
 271 images (Figure 3). The augmented images are easily distinguishable from the training data, while
 272 the generated samples exhibit a slight shift towards the left of the training distribution. This shift
 273 suggests a degree of memorization, as the synthetic samples appear to be closer to the training data
 274 than the training data is to the validation images.

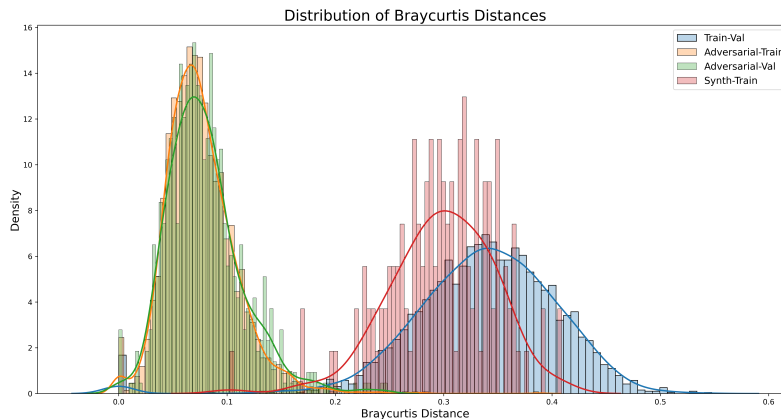


Figure 3: Memorization characterization by the two best-performing self-supervised contrastive networks, DinoV2 (left) and DenseNet121 (right), for generated samples by a DiT model.

275 5 Discussion

276 Our study presents a comprehensive, data-driven approach to evaluating and characterizing the
 277 memorization behavior of generative models in medical imaging. By systematically comparing
 278 various network architectures, loss functions, pretraining datasets, and distance metrics, we have
 279 identified optimal configurations for detecting potential privacy concerns in synthetic images. The
 280 results demonstrate that the developed method can identify all augmented images when using Triplet
 281 Margin loss with models like DinoV2, DenseNet121, and ResNet50, particularly when paired with
 282 the Bray-Curtis or Standardized Euclidean distance metrics. The ability to quantify the degree of
 283 memorization in generated images offers a method to assess the risk of patient data leakage in
 284 synthetic datasets. This approach can be integrated into the training pipeline of generative models,
 285 serving as an early warning system for memorization and potential privacy breaches.

286 **Limitations** As for limitations, our study is based on a private mammography dataset from various
 287 institutions. Although this dataset is substantial and diverse, the generalizability of our findings to
 288 other medical imaging modalities or natural image datasets remains to be validated. Future work
 289 should address these limitations by generating a foundational model that serves for both 2D and 3D
 290 data, multi-institutional and multi-modality datasets to avoid having to fine-tune the model for each
 291 dataset. A comparative analysis of various generative model architectures and stronger conditioning
 292 forms (text or segmentation) would provide a more comprehensive understanding of memorization
 293 behavior across generative models.

294 **References**

- 295 [1] Lennart R. Koetzier, Jie Wu, Domenico Mastrodicasa, Aline Lutz, Matthew Chung, W. Adam
296 Koszek, Jayanth Pratap, Akshay S. Chaudhari, Pranav Rajpurkar, Matthew P. Lungren, and
297 Martin J. Willeminck. Generating synthetic data for medical imaging. *Radiology*, 312(3),
298 September 2024.
- 299 [2] Gerrit J. J. van den Burg and Christopher K. I. Williams. On memorization in probabilistic deep
300 generative models. *CoRR*, abs/2106.03216, 2021.
- 301 [3] Ching-Yuan Bai, Hsuan-Tien Lin, Colin Raffel, and Wendy Kan. A large-scale study on training
302 sample memorization in generative modeling, 2021.
- 303 [4] Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and
304 Chiyuan Zhang. Quantifying memorization across neural language models, 2023.
- 305 [5] Salman Ul Hassan Dar, Arman Ghanaat, Jannik Kahmann, Isabelle Ayx, Theano Papavassiliu,
306 Stefan O. Schoenberg, and Sandy Engelhardt. Investigating data memorization in 3d latent
307 diffusion models for medical image synthesis, 2023.
- 308 [6] Muhammad Usman Akbar, Wuhao Wang, and Anders Eklund. Beware of diffusion models for
309 synthesizing medical images – a comparison with gans in terms of memorizing brain mri and
310 chest x-ray images, 2023.
- 311 [7] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein.
312 Diffusion art or digital forgery? investigating data replication in diffusion models, 2022.
- 313 [8] Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. An overview of diffusion models:
314 Applications, guided generation, statistical rates and optimization, 2024.
- 315 [9] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.
316 Improved techniques for training gans. *CoRR*, abs/1606.03498, 2016.
- 317 [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer,
318 and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium.
319 *CoRR*, abs/1706.08500, 2017.
- 320 [11] Nicholas Carlini, Chang Liu, Jernej Kos, Úlfar Erlingsson, and Dawn Song. The secret
321 sharer: Measuring unintended neural network memorization & extracting secrets. *CoRR*,
322 abs/1802.08232, 2018.
- 323 [12] Salman Ul Hassan Dar, Marvin Seyfarth, Jannik Kahmann, Isabelle Ayx, Theano Papavassiliu,
324 Stefan O. Schoenberg, and Sandy Engelhardt. Unconditional latent diffusion models memorize
325 patient imaging data, 2024.
- 326 [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
327 recognition. *CoRR*, abs/1512.03385, 2015.
- 328 [14] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
329 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran,
330 Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra,
331 Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick
332 Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without
333 supervision, 2024.
- 334 [15] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
335 deep learning requires rethinking generalization. *CoRR*, abs/1611.03530, 2016.
- 336 [16] Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In
337 *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, volume 80
338 of *STOC '20*, page 954–959. ACM, June 2020.
- 339 [17] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio,
340 Maxinder S. Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, and
341 Simon Lacoste-Julien. A closer look at memorization in deep networks, 2017.

- 342 [18] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering
343 the long tail via influence estimation. *CoRR*, abs/2008.03703, 2020.
- 344 [19] Casey Meehan, Kamalika Chaudhuri, and Sanjoy Dasgupta. A non-parametric test to detect
345 data-copying in generative models. *CoRR*, abs/2004.05675, 2020.
- 346 [20] Shaoyan Pan, Tonghe Wang, Richard L J Qiu, Marian Axente, Chih-Wei Chang, Junbo Peng,
347 Ashish B Patel, Joseph Shelton, Sagar A Patel, Justin Roper, and Xiaofeng Yang. 2d medical
348 image synthesis using transformer-based denoising diffusion probabilistic model. *Physics in
349 Medicine amp; Biology*, 68(10):105004, May 2023.
- 350 [21] Walter H. L. Pinaya, Petru-Daniel Tudosiu, Jessica Dafflon, Pedro F da Costa, Virginia Fernan-
351 dez, Parashkev Nachev, Sebastien Ourselin, and M. Jorge Cardoso. Brain imaging generation
352 with latent diffusion models, 2022.
- 353 [22] Firas Khader, Gustav Müller-Franzes, Soroosh Tayebi Arasteh, Tianyu Han, Christoph Haar-
354 burger, Maximilian Schulze-Hagen, Philipp Schad, Sandy Engelhardt, Bettina Baeßler, Se-
355 bastian Foersch, Johannes Stegmaier, Christiane Kuhl, Sven Nebelung, Jakob Nikolas Kather,
356 and Daniel Truhn. Denoising diffusion probabilistic models for 3d medical image generation.
357 *Scientific Reports*, 13(1), May 2023.
- 358 [23] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar,
359 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC
360 Conference on Computer and Communications Security, CCS'16*. ACM, October 2016.
- 361 [24] Virginia Fernandez, Pedro Sanchez, Walter Hugo Lopez Pinaya, Grzegorz Jacenków, Sotirios A.
362 Tsaftaris, and Jorge Cardoso. Privacy distillation: Reducing re-identification risk of multimodal
363 diffusion models, 2023.
- 364 [25] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr,
365 Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models,
366 2023.
- 367 [26] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks.
368 *CoRR*, abs/1608.06993, 2016.
- 369 [27] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna.
370 Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015.
- 371 [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
372 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
373 Sutskever. Learning transferable visual models from natural language supervision. *CoRR*,
374 abs/2103.00020, 2021.
- 375 [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
376 Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-
377 Fei. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014.
- 378 [30] Xueyan Mei, Zelong Liu, Philip M. Robson, Brett Marinelli, Mingqian Huang, Amish Doshi,
379 Adam Jacobi, Chendi Cao, Katherine E. Link, Thomas Yang, Ying Wang, Hayit Greenspan,
380 Timothy Deyer, Zahi A. Fayad, and Yang Yang. Radimagenet: An open radiologic deep
381 learning research dataset for effective transfer learning. *Radiology: Artificial Intelligence*, 4(5),
382 September 2022.
- 383 [31] Vassileios Balntas, Edgar Riba, Daniel Ponsa, and Krystian Mikolajczyk. Learning local feature
384 descriptors with triplets and shallow convolutional neural networks. In *Proceedings of the British
385 Machine Vision Conference 2016, BMVC 2016*, pages 119.1–119.11. British Machine Vision
386 Association, 2016.
- 387 [32] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive
388 predictive coding. *CoRR*, abs/1807.03748, 2018.
- 389 [33] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint
390 arXiv:2212.09748*, 2022.

391 **A Appendix**

392 **A.1 Backbone Comparison p-values**

393 Table 2 shows the p-values for the comparison between ImageNet and RadImageNet pretraining
 394 across different model backbones and loss functions.

Table 2: P-values for Backbone Comparison (ImageNet vs RadImageNet)

Model	Loss	p-value
ResNet50	InfoNCE	0.0205
ResNet50	Triplet	0.5401
Inception	InfoNCE	0.0428
Inception	Triplet	0.8749
DINO	InfoNCE	0.0147
DINO	Triplet	0.3230
DenseNet121	InfoNCE	0.0498
DenseNet121	Triplet	0.0599
CLIP	InfoNCE	0.1731
CLIP	Triplet	0.9528

395 The results indicate varying levels of statistical significance in the performance difference between
 396 ImageNet and RadImageNet pretraining across different model architectures and loss functions.
 397 P-values below 0.05 suggest a statistically significant difference:

- 398 • ResNet50, Inception, DINO, and DenseNet121 show statistically significant differences ($p <$
 399 0.05) when using InfoNCE loss.
- 400 • The Triplet loss generally shows no significant difference between ImageNet and RadIma-
 401 geNet pretraining across all models.
- 402 • CLIP shows no significant difference for either loss function.

403 These results suggest that the choice of pretraining dataset (ImageNet vs RadImageNet) may have a
 404 more pronounced effect when using InfoNCE loss, particularly for certain model architectures.

405 **A.2 Loss Function Comparison Results**

406 Table 3 presents the comparison between Triplet and InfoNCE loss functions across different model
 407 backbones and pretraining datasets.

Table 3: Comparison of Triplet and InfoNCE Loss Functions

Model	Pretrain	p-value	Triplet Mean	InfoNCE Mean
ResNet50	ImageNet	0.0008	0.5694	0.1934
ResNet50	RadImageNet	0.0249	0.6604	0.3572
Inception	ImageNet	<0.0001	0.3705	0.0929
Inception	RadImageNet	<0.0001	0.3576	0.0645
DINO	ImageNet	<0.0001	0.7216	0.1051
DINO	RadImageNet	0.0006	0.5671	0.1826
DenseNet121	ImageNet	0.0016	0.7102	0.2733
DenseNet121	RadImageNet	0.0019	0.4396	0.1720
CLIP	ImageNet	<0.0001	0.5919	0.0533
CLIP	RadImageNet	<0.0001	0.6004	0.0642

408 The results show a consistent and statistically significant difference between the performance of
 409 Triplet and InfoNCE loss functions across all model architectures and pretraining datasets. Key
 410 observations include:

- 411 • All comparisons show p-values well below 0.05, indicating strong statistical significance in
 412 the difference between Triplet and InfoNCE loss performance.

- 413 • Triplet loss consistently outperforms InfoNCE loss across all models and pretraining datasets,
414 as evidenced by the higher mean values.
- 415 • The performance gap between Triplet and InfoNCE loss appears to be more pronounced for
416 some models (e.g., DINO, CLIP) compared to others.
- 417 • The choice of pretraining dataset (ImageNet vs RadImageNet) seems to influence the
418 magnitude of the difference between the two loss functions, though the trend of Triplet loss
419 outperforming InfoNCE remains consistent.

420 These findings suggest that the choice of loss function has a substantial impact on model performance,
421 with Triplet loss demonstrating superior results across various model architectures and pretraining
422 scenarios. This consistent pattern underscores the importance of loss function selection in the design
423 of contrastive learning frameworks for image analysis tasks.

424 A.3 Network Architecture Comparison Results

425 Table 4 presents the pairwise comparisons between different network architectures, considering their
426 performance with specific pretraining datasets and loss functions.

Table 4: Pairwise Comparison of Network Architectures

Model 1	Model 2	p-value	Model 1 Mean	Model 2 Mean
ResNet50 (RadImageNet, Triplet)	Inception (ImageNet, Triplet)	0.0298	0.6604	0.3705
ResNet50 (RadImageNet, Triplet)	DINO (ImageNet, Triplet)	0.7125	0.6604	0.7216
ResNet50 (RadImageNet, Triplet)	DenseNet121 (ImageNet, Triplet)	0.7606	0.6604	0.7102
ResNet50 (RadImageNet, Triplet)	CLIP (RadImageNet, Triplet)	0.6940	0.6604	0.6004
Inception (ImageNet, Triplet)	DINO (ImageNet, Triplet)	0.0139	0.3705	0.7216
Inception (ImageNet, Triplet)	DenseNet121 (ImageNet, Triplet)	0.0147	0.3705	0.7102
Inception (ImageNet, Triplet)	CLIP (RadImageNet, Triplet)	0.0588	0.3705	0.6004
DINO (ImageNet, Triplet)	DenseNet121 (ImageNet, Triplet)	0.9461	0.7216	0.7102
DINO (ImageNet, Triplet)	CLIP (RadImageNet, Triplet)	0.4465	0.7216	0.6004
DenseNet121 (ImageNet, Triplet)	CLIP (RadImageNet, Triplet)	0.4825	0.7102	0.6004

427 The results reveal interesting patterns in the performance of different network architectures:

- 428 • ResNet50 (RadImageNet, Triplet) shows significantly better performance than Inception
429 (ImageNet, Triplet) with a p-value of 0.0298.
- 430 • There is no statistically significant difference between ResNet50 (RadImageNet, Triplet)
431 and DINO, DenseNet121, or CLIP, as evidenced by high p-values (>0.05).
- 432 • Inception (ImageNet, Triplet) consistently underperforms compared to other architectures,
433 with statistically significant differences against DINO and DenseNet121 (p-values < 0.05).
- 434 • DinoV2, DenseNet121, and CLIP show comparable performance, with no statistically
435 significant differences among them (p-values > 0.05).
- 436 • The choice of pretraining dataset (ImageNet vs RadImageNet) appears to influence perfor-
437 mance, but the effect varies across architectures.

438 These findings suggest that:

- 439 1. ResNet50, DinoV2, DenseNet121, and CLIP demonstrate robust performance across differ-
440 ent pretraining scenarios when using Triplet loss.
- 441 2. Inception architecture may not be optimal for this particular task, consistently showing lower
442 performance.
- 443 3. The impact of pretraining dataset choice (ImageNet vs RadImageNet) may be architecture-
444 dependent and warrants further investigation.

445 Overall, these results underscore the importance of carefully selecting network architectures and
446 pretraining strategies in contrastive learning frameworks for image analysis tasks. The comparable
447 performance of several architectures (ResNet50, DINO, DenseNet121, CLIP) suggests that factors
448 beyond architecture, such as loss function and pretraining data, play crucial roles in determining
449 overall system performance.

450 **A.4 InfoNCE Loss Metric Results**

451 In this section we show the results on metric comparison for the models trained with InfoNCE loss
452 (Table 4). We observe that on average the results are much lower than when using Triplet Margin loss,
453 and distance metrics like Euclidean are preferred over correlation, Mahalanobis or Cosine similarity.

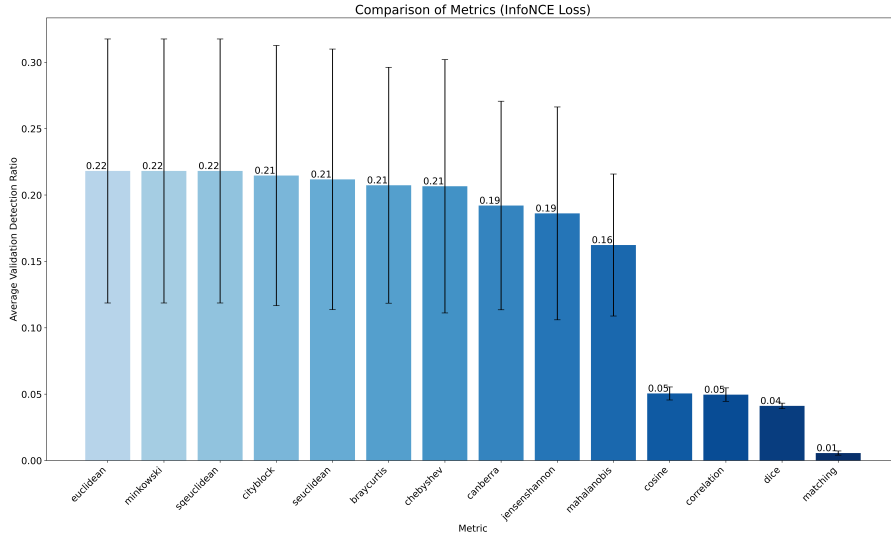


Figure 4: Comparison of distance metrics performance with InfoNCE loss in terms of mean validation detection ratio. Error bars represent confidence intervals.

454 **A.5 Class conditional Performance Metrics DiT vs StyleGAN2**

455 The generative imaging results, shown in Figure ??, indicate that the class-conditional DiT model
456 performs better or at least comparably across all relevant metrics to the unconditional StyleGAN2s.
457 DiT models learns more comprehensively the real image distribution and is less affected by mode-
458 collapse. Both models exhibit a tendency for memorization, as the generated data closely resembles
459 the training data more than the training data resembles the validation data. However, the degree of
460 memorization observed is not excessive after manually inspection.

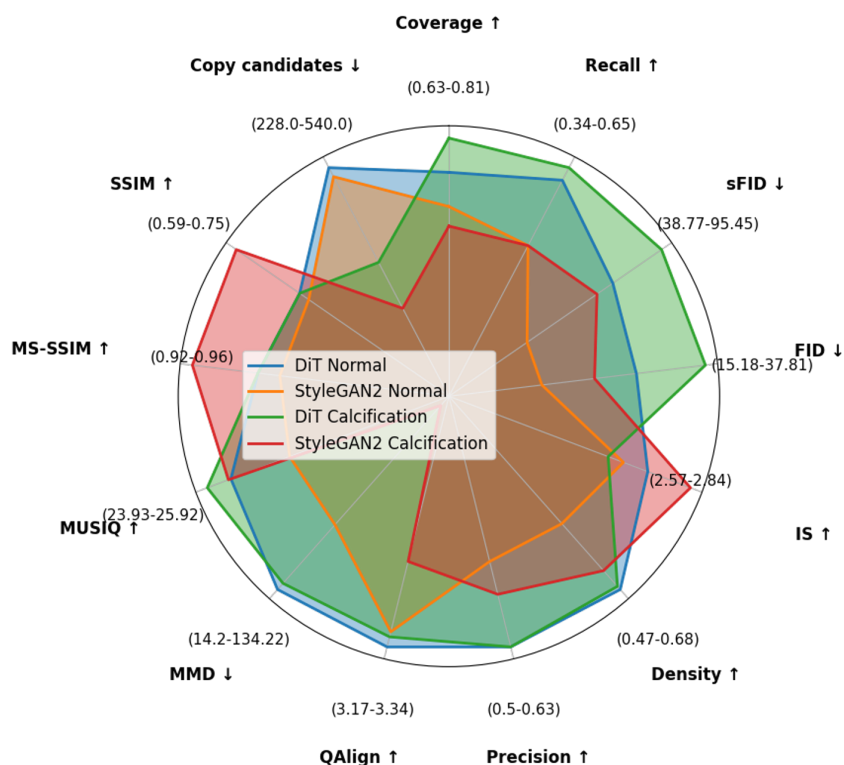


Figure 5: High-Resolution Image Synthesis Comparison Between StyleGAN and Diffusion Models (2048 x 2048 pixels).

461 NeurIPS Paper Checklist

462 1. Claims

463 Question: Do the main claims made in the abstract and introduction accurately reflect the
464 paper’s contributions and scope?

465 Answer: [Yes]

466 Justification: The main goals and scope are specified at the abstract and introduction by
467 acknowledging we build upon a previous work on the memorization topic and we focus on
468 the optimization of it via a thorough benchmark to ensure that non-trivial decisions about
469 the method configuration (Backbone, Loss, Pretraining, Similarity metric) are backed with
470 experiments.

471 Guidelines:

- 472 • The answer NA means that the abstract and introduction do not include the claims
473 made in the paper.
- 474 • The abstract and/or introduction should clearly state the claims made, including the
475 contributions made in the paper and important assumptions and limitations. A No or
476 NA answer to this question will not be perceived well by the reviewers.
- 477 • The claims made should match theoretical and experimental results, and reflect how
478 much the results can be expected to generalize to other settings.
- 479 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
480 are not attained by the paper.

481 2. Limitations

482 Question: Does the paper discuss the limitations of the work performed by the authors?

483 Answer: [Yes]

484 Justification: The paper comment on the limitations of our work in the last section of the
485 paper where we acknowledge the generalization aspect of our method remains unclear and
486 fine-tuning to each individual dataset might be required. And we propose as future work
487 to create a multi-modality medical imaging dataset to develop a contrastive foundational
488 model capable of generalize to unseen data and avoid the fine-tuning step.

489 Guidelines:

- 490 • The answer NA means that the paper has no limitation while the answer No means that
491 the paper has limitations, but those are not discussed in the paper.
- 492 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 493 • The paper should point out any strong assumptions and how robust the results are to
494 violations of these assumptions (e.g., independence assumptions, noiseless settings,
495 model well-specification, asymptotic approximations only holding locally). The authors
496 should reflect on how these assumptions might be violated in practice and what the
497 implications would be.
- 498 • The authors should reflect on the scope of the claims made, e.g., if the approach was
499 only tested on a few datasets or with a few runs. In general, empirical results often
500 depend on implicit assumptions, which should be articulated.
- 501 • The authors should reflect on the factors that influence the performance of the approach.
502 For example, a facial recognition algorithm may perform poorly when image resolution
503 is low or images are taken in low lighting. Or a speech-to-text system might not be
504 used reliably to provide closed captions for online lectures because it fails to handle
505 technical jargon.
- 506 • The authors should discuss the computational efficiency of the proposed algorithms
507 and how they scale with dataset size.
- 508 • If applicable, the authors should discuss possible limitations of their approach to
509 address problems of privacy and fairness.
- 510 • While the authors might fear that complete honesty about limitations might be used by
511 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
512 limitations that aren't acknowledged in the paper. The authors should use their best
513 judgment and recognize that individual actions in favor of transparency play an impor-
514 tant role in developing norms that preserve the integrity of the community. Reviewers
515 will be specifically instructed to not penalize honesty concerning limitations.

516 3. Theory Assumptions and Proofs

517 Question: For each theoretical result, does the paper provide the full set of assumptions and
518 a complete (and correct) proof?

519 Answer: [Yes]

520 Justification: We do define the problem in a closed format so that all assumptions and criteria
521 are clear, specifically for our detection benchmark via adversarial augmentations.

522 Guidelines:

- 523 • The answer NA means that the paper does not include theoretical results.
- 524 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
525 referenced.
- 526 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 527 • The proofs can either appear in the main paper or the supplemental material, but if
528 they appear in the supplemental material, the authors are encouraged to provide a short
529 proof sketch to provide intuition.
- 530 • Inversely, any informal proof provided in the core of the paper should be complemented
531 by formal proofs provided in appendix or supplemental material.
- 532 • Theorems and Lemmas that the proof relies upon should be properly referenced.

533 4. Experimental Result Reproducibility

534 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
535 perimental results of the paper to the extent that it affects the main claims and/or conclusions
536 of the paper (regardless of whether the code and data are provided or not)?

537 Answer: [Yes]

538 Justification: The paper provides a complete definition of the attained problem, a definition
539 of the networks, losses, pretraining, hyperparameters and metrics employed as well as the
540 resources needed to run all the experiments (Methods section).

541 Guidelines:

- 542 • The answer NA means that the paper does not include experiments.
- 543 • If the paper includes experiments, a No answer to this question will not be perceived
544 well by the reviewers: Making the paper reproducible is important, regardless of
545 whether the code and data are provided or not.
- 546 • If the contribution is a dataset and/or model, the authors should describe the steps taken
547 to make their results reproducible or verifiable.
- 548 • Depending on the contribution, reproducibility can be accomplished in various ways.
549 For example, if the contribution is a novel architecture, describing the architecture fully
550 might suffice, or if the contribution is a specific model and empirical evaluation, it may
551 be necessary to either make it possible for others to replicate the model with the same
552 dataset, or provide access to the model. In general, releasing code and data is often
553 one good way to accomplish this, but reproducibility can also be provided via detailed
554 instructions for how to replicate the results, access to a hosted model (e.g., in the case
555 of a large language model), releasing of a model checkpoint, or other means that are
556 appropriate to the research performed.
- 557 • While NeurIPS does not require releasing code, the conference does require all submis-
558 sions to provide some reasonable avenue for reproducibility, which may depend on the
559 nature of the contribution. For example
 - 560 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
561 to reproduce that algorithm.
 - 562 (b) If the contribution is primarily a new model architecture, the paper should describe
563 the architecture clearly and fully.
 - 564 (c) If the contribution is a new model (e.g., a large language model), then there should
565 either be a way to access this model for reproducing the results or a way to reproduce
566 the model (e.g., with an open-source dataset or instructions for how to construct
567 the dataset).
 - 568 (d) We recognize that reproducibility may be tricky in some cases, in which case
569 authors are welcome to describe the particular way they provide for reproducibility.
570 In the case of closed-source models, it may be that access to the model is limited in
571 some way (e.g., to registered users), but it should be possible for other researchers
572 to have some path to reproducing or verifying the results.

573 5. Open access to data and code

574 Question: Does the paper provide open access to the data and code, with sufficient instruc-
575 tions to faithfully reproduce the main experimental results, as described in supplemental
576 material?

577 Answer: [No]

578 Justification We do provide access to the code that enables to reproduce the results with any
579 dataset in a flexible manner. However, we do not disclose the dataset employed due to the
580 sensitivity of the medical images which although anonymized we are not allowed to share
581 openly.

582 Guidelines:

- 583 • The answer NA means that paper does not include experiments requiring code.
- 584 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
585 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 586 • While we encourage the release of code and data, we understand that this might not be
587 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not

- 588 including code, unless this is central to the contribution (e.g., for a new open-source
589 benchmark).
- 590 • The instructions should contain the exact command and environment needed to run to
591 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
592
 - 593 • The authors should provide instructions on data access and preparation, including how
594 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
 - 595 • The authors should provide scripts to reproduce all experimental results for the new
596 proposed method and baselines. If only a subset of experiments are reproducible, they
597 should state which ones are omitted from the script and why.
 - 598 • At submission time, to preserve anonymity, the authors should release anonymized
599 versions (if applicable).
 - 600 • Providing as much information as possible in supplemental material (appended to the
601 paper) is recommended, but including URLs to data and code is permitted.

602 6. Experimental Setting/Details

603 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
604 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
605 results?

606 Answer: [Yes]

607 Justification: We do describe all the hyper-parameters employed and how they were selected
608 as well as the optimizer and other training details selected via manual tuning in the methods
609 section.

610 Guidelines:

- 611 • The answer NA means that the paper does not include experiments.
- 612 • The experimental setting should be presented in the core of the paper to a level of detail
613 that is necessary to appreciate the results and make sense of them.
- 614 • The full details can be provided either with the code, in appendix, or as supplemental
615 material.

616 7. Experiment Statistical Significance

617 Question: Does the paper report error bars suitably and correctly defined or other appropriate
618 information about the statistical significance of the experiments?

619 Answer: [Yes]

620 Justification: We describe the statistical tests employed the meaning of the error bars which
621 support the main claims of our paper (Methods and results section, figure 1 2 and Tables in
622 supplementary)

623 Guidelines:

- 624 • The answer NA means that the paper does not include experiments.
- 625 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
626 dence intervals, or statistical significance tests, at least for the experiments that support
627 the main claims of the paper.
- 628 • The factors of variability that the error bars are capturing should be clearly stated (for
629 example, train/test split, initialization, random drawing of some parameter, or overall
630 run with given experimental conditions).
- 631 • The method for calculating the error bars should be explained (closed form formula,
632 call to a library function, bootstrap, etc.)
- 633 • The assumptions made should be given (e.g., Normally distributed errors).
- 634 • It should be clear whether the error bar is the standard deviation or the standard error
635 of the mean.
- 636 • It is OK to report 1-sigma error bars, but one should state it. The authors should
637 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
638 of Normality of errors is not verified.

- 639
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 640
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
- 641
- 642
- 643

644 8. Experiments Compute Resources

645 Question: For each experiment, does the paper provide sufficient information on the com-
646 puter resources (type of compute workers, memory, time of execution) needed to reproduce
647 the experiments?

648 Answer: [Yes]

649 Justification: We specify the hardware used to run all the experiments in the implementation
650 details of the methods section.

651 Guidelines:

- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 652
- 653
- 654
- 655
- 656
- 657
- 658
- 659

660 9. Code Of Ethics

661 Question: Does the research conducted in the paper conform, in every respect, with the
662 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

663 Answer: [Yes]

664 Justification: We have read and checked that we conform with the code of ethics of NeurIPS
665 and we preserve anonymity in the current submission for a fair and unbiased review process.

666 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 667
- 668
- 669
- 670
- 671

672 10. Broader Impacts

673 Question: Does the paper discuss both potential positive societal impacts and negative
674 societal impacts of the work performed?

675 Answer: [Yes]

676 Justification: We discuss the different aspects to consider about generative models memo-
677 rization in medical imaging and how by not implementing a safety measure patient data could
678 be leaked into the generated datasets which would violate privacy regulations (Introduction,
679 Related work and Discussion)

680 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- 681
- 682
- 683
- 684
- 685
- 686
- 687

- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

703 11. Safeguards

704 Question: Does the paper describe safeguards that have been put in place for responsible
705 release of data or models that have a high risk for misuse (e.g., pretrained language models,
706 image generators, or scraped datasets)?

707 Answer: [Yes]

708 Justification:

709 Guidelines:

- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

720 12. Licenses for existing assets

721 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
722 the paper, properly credited and are the license and terms of use explicitly mentioned and
723 properly respected?

724 Answer: [Yes]

725 Justification: We are the creators of all the models, code and data employed with the usage
726 of open source fully complying with the licensing of each pretrained backbone.

727 Guidelines:

- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

741 • If this information is not available online, the authors are encouraged to reach out to
742 the asset’s creators.

743 13. New Assets

744 Question: Are new assets introduced in the paper well documented and is the documentation
745 provided alongside the assets?

746 Answer: [Yes]

747 Justification: The new assets, the code, are properly documented in the GitHub repository
748 specified on the footnote of the first page of the paper.

749 Guidelines:

- 750 • The answer NA means that the paper does not release new assets.
- 751 • Researchers should communicate the details of the dataset/code/model as part of their
752 submissions via structured templates. This includes details about training, license,
753 limitations, etc.
- 754 • The paper should discuss whether and how consent was obtained from people whose
755 asset is used.
- 756 • At submission time, remember to anonymize your assets (if applicable). You can either
757 create an anonymized URL or include an anonymized zip file.

758 14. Crowdsourcing and Research with Human Subjects

759 Question: For crowdsourcing experiments and research with human subjects, does the paper
760 include the full text of instructions given to participants and screenshots, if applicable, as
761 well as details about compensation (if any)?

762 Answer: [NA]

763 Justification: Human subjects were not use in an instructive manner, no instructions to
764 specify.

765 Guidelines:

- 766 • The answer NA means that the paper does not involve crowdsourcing nor research with
767 human subjects.
- 768 • Including this information in the supplemental material is fine, but if the main contribu-
769 tion of the paper involves human subjects, then as much detail as possible should be
770 included in the main paper.
- 771 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
772 or other labor should be paid at least the minimum wage in the country of the data
773 collector.

774 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 775 Subjects

776 Question: Does the paper describe potential risks incurred by study participants, whether
777 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
778 approvals (or an equivalent approval/review based on the requirements of your country or
779 institution) were obtained?

780 Answer: [Yes]

781 Justification: The study contains retropective fully anonymized data based and all subjects
782 gave their consent to use their health data (images) for further research and was conducted
783 based on a IRB approval

784 Guidelines:

- 785 • The answer NA means that the paper does not involve crowdsourcing nor research with
786 human subjects.
- 787 • Depending on the country in which research is conducted, IRB approval (or equivalent)
788 may be required for any human subjects research. If you obtained IRB approval, you
789 should clearly state this in the paper.
- 790 • We recognize that the procedures for this may vary significantly between institutions
791 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
792 guidelines for their institution.

793
794

- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.