

---

# Medical Imaging Complexity and its Effects on GAN Performance

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The proliferation of machine learning models in diverse clinical applications has  
2       led to a growing need for high-fidelity, medical image training data. Such data  
3       is often scarce due to cost constraints and privacy concerns. Alleviating this  
4       burden, medical image synthesis via generative adversarial networks (GANs)  
5       emerged as a powerful method for synthetically generating photo-realistic images  
6       based on existing sets of real medical images. However, the exact image set size  
7       required to efficiently train such a GAN is unclear. In this work, we experimentally  
8       establish benchmarks that measure the relationship between a sample dataset  
9       size and the fidelity of the generated images, given the dataset’s distribution of  
10      image complexities. We analyze statistical metrics based on delentropy, an image  
11      complexity measure rooted in Shannon’s entropy in information theory. For our  
12      pipeline, we conduct experiments with two state-of-the-art GANs, StyleGAN 3 and  
13      SPADE-GAN, trained on multiple medical imaging datasets with variable sample  
14      sizes. Across both GANs, general performance improved with increasing training  
15      set size but suffered with increasing complexity.

## 16   1 Introduction

17   Machine learning in healthcare is a rapidly growing field with countless applications [28] including  
18   disease diagnosis [24], clinical treatment [26], drug development [19], and mental health [7]. The  
19   machine learning models driving these advances require the collection of high-quality, annotated  
20   medical training data, which persists as an arduous task due to privacy concerns surrounding sensitive  
21   patient data [23] and the time-intensive nature of labeling [5]. To address these issues, synthetic  
22   data—artificially generated information mimicking real-world data—has surfaced as a promising  
23   solution [8].

24   Currently, generative adversarial networks (GANs) remain one of the leading approaches to synthetic  
25   data generation [17]. Since its inception in 2014 [6], GANs have gained increasing attention in  
26   the medical research community due to their ability to synthesize medical images [29]. However,  
27   achieving results with high fidelity remains a difficult task factoring the lack of medical data and  
28   prevalence of smaller datasets in the medical domain. With limited data, a GAN’s efficacy is directly  
29   affected with consequences including mode collapse, where the generator produces a limited variety  
30   of outputs [20], and overfitting, where the GAN replicates training data rather than generalizing from  
31   it [33].

32   Various papers such as Wang et al. [32]’s transfer learning and Robb et al. [25]’s Few-Shot GAN  
33   (FSGAN) have addressed these issues as architecture-centric approaches, achieving increased training  
34   efficiency only as a result of the changes in the GAN’s structure. However, such approaches are  
35   ineffective when making alterations to a GAN’s internal structure are not feasible and when time  
36   constraints are present. As such, a data-centric approach by providing the GAN with the optimal

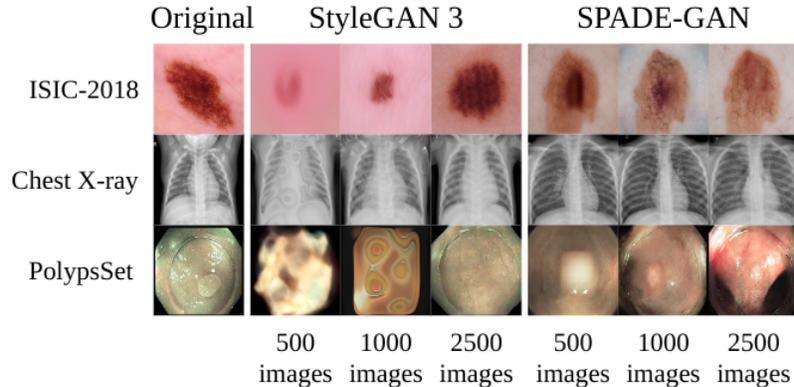


Figure 1: Comparison between original images and synthetic images from StyleGAN 3 and SPADE-GAN based on variable image set sizes.

37 amount of data to produce high-quality results is more appropriate. Nevertheless, the exact sample  
 38 set size required to train state-of-the-art GANs is obscure.

39 In this study, we introduce a data-centric optimization method to create efficient GAN training for  
 40 medical image synthesis. Our approach investigates how the image complexity distribution of a  
 41 medical image dataset can be utilized as a measure of training difficulty for a GAN. By doing so, we  
 42 can ascertain a correlation between the image complexities of the training images and the optimal  
 43 training set sizes by establishing benchmarks that evaluate the relationship between a sample training  
 44 set size and the fidelity of the generated images. We hypothesize that given a dataset of a specific  
 45 image complexity distribution, healthcare professionals can reference the closest image fidelity curve  
 46 to identify the optimal amount of experimental trials to produce superlative results. Ultimately, our  
 47 approach can avoid both undertraining and wasteful overtraining by constructing a data-efficient,  
 48 GAN training pipeline.

## 49 2 Background

50 **Generative Adversarial Networks (GANs)** Introduced by Goodfellow et al. [6], GANs are a class  
 51 of generative models that consist of two convolutional neural networks: a generator  $G$ , which aims  
 52 to transform its latent variable distribution  $p(z)$  to closely resemble the training data distribution  
 53  $p(x)$ , and a discriminator  $D$ , which differentiates between the ground truth and data generated by  $G$ .  
 54 Training is an adversarial process where  $G$  attempts to deceive  $D$  into classifying its outputs as real.  
 55 This two-player minimax game is represented by the following loss function:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\log D(\mathbf{x})] \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] . \quad (1)$$

56 Many papers have tried to address data scarcity and computational costs in GAN training architec-  
 57 turally. One approach proposed by Wang et al. [32] is transfer learning, which consists of fine-tuning  
 58 a pre-trained generator and discriminator to the desired domain. However, if the pre-trained models  
 59 do not align well with the target domain, this could result in even higher data and computational  
 60 demands [34]. Robb et al. [25] proposed another solution through their Few-Shot GAN (FSGAN),  
 61 achieving impressive adaptation even with extremely few training examples, albeit at the cost of  
 62 prolonged training times. This results in the reduced quality and diversity of the synthetic data when  
 63 time constraints are present.

64 **Image Complexity** Objectively, image complexity can be defined as the variety of features and  
 65 details within an image. It has been shown that information entropy is a traditional, heuristic-based  
 66 method of calculating the complexities of images in small-scale datasets [16].

67 Traditional information entropy, or Shannon entropy, is a foundational abstraction in information  
 68 theory introduced by Shannon [27]. Used as a measure of uncertainty or “surprise” in data, it is the

69 variation in the distribution of pixel intensities of an image in grayscale format. The equation is  
70 defined as

$$H = - \sum_{i=0}^{n-1} p_i \log_b p_i, \quad (2)$$

71 where  $n$  denotes the number of gray levels (256 for 8-bit images),  $b$  stands for the logarithmic base  
72 (returning bits when  $b = 2$ ), and  $p_i$  is the probability of a pixel having gray level  $i$ . However,  
73 although Shannon entropy considers compositional image information, it fails to account for spatial  
74 information, specifically the relationship between neighbouring pixels [4].

75 Another entropy-based metric, the Gray Level Co-Occurrence Matrix (GLCM), unlike Shannon’s  
76 Entropy, is a measure of how often pairs of pixel values occur in a grayscale image distribution [9].  
77 Taking into account local spatial information, the GLCM is useful for textural analysis tasks such as  
78 feature extraction for medical image segmentation [14]. The GLCM entropy can be represented as

$$H_g = - \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} p_{(i,j)} \log_b p_{(i,j)}, \quad (3)$$

79 where  $p_{(i,j)}$  is the probability of a two pixels having gray levels  $i$  and  $j$  at a certain angle  $\theta$  and  
80 distance  $d$  away from each other.

### 81 3 Methodology

#### 82 3.1 Image Complexity Metric

83 We utilize Larkin’s delentropy, a metric identical to the Shannon entropy and the GLCM, but  
84 incorporating a new density function known as the *deledensity* [15]. By analyzing the relationship  
85 between the local and global features of an image, delentropy accounts for an image’s gradient vector  
86 field and pixel co-occurrence, encapsulating its spatial information as a whole. The deledensity, as a  
87 joint probability function is formulated as

$$p_{(i,j)} = \frac{1}{4WH} \sum_{w=0}^{W-1} \sum_{h=0}^{H-1} \delta_{i,d_x(w,h)} \delta_{j,d_y(w,h)}, \quad (4)$$

88 where  $d_x$  and  $d_y$  denote the derivative kernels in the x and y direction,  $\delta$  is the Kronecker delta  
89 to describe the binning operation required to generate a histogram, and  $H$  and  $W$  is the image’s  
90 dimensions (height and width) [13]. By obtaining this, we can then calculate delentropy as

$$DE = -\frac{1}{2} \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} p_{(i,j)} \log_b p_{(i,j)}, \quad (5)$$

91 such that  $I$  and  $J$  represent the number of bins (discrete cells) in the 2D distribution, and the  $\frac{1}{2}$  is  
92 derived from Papoulis’ generalized sampling expansion [21].

93 To interpret this measure, yielding a high delentropy suggests an image has a high range of variation  
94 in pixel intensities and more sophisticated details. A low delentropy can be interpreted as a result  
95 of having a uniform distribution of pixel intensities, indicating simple structure and a less-detailed  
96 image.

97 Prior to any calculations, each image is preprocessed into an 8-bit, grayscale image. This ensures  
98 delentropy can be calculated in a consistent, single-channel format throughout each dataset.

### 99 3.2 GAN Selection

100 Core to the experimental approach was the selection of two state-of-the-art GANs, SPADE-GAN  
101 [22] and StyleGAN 3 [11] on which to run the experimental pipeline. These networks have been  
102 widely adopted by the medical image synthesis community and empirically observed to produce  
103 superior-quality medical images when compared to predecessor GANs [30]. StyleGAN 3’s large  
104 community support and wide availability of its code repository along with its numerous configurations  
105 for different training settings was taken into account as well.

### 106 3.3 GAN Pipeline

107 Throughout the experiment, as our approach was data-centric, StyleGAN 3 and SPADE-GAN  
108 were run with the official, publicly available implementations with default hyperparameters and no  
109 augmentations to each network’s architecture.

110 **Preprocessing** We first set all images to a consistent 512x512 resolution. As such, training param-  
111 eters were based on the size of the preprocessed images, as documented in the official implementations.  
112 SPADE-GAN additionally relies on segmentation masks to produce synthetic data. We used pre-  
113 existing annotations for ISIC-2018 and the Polyps Set. Because the Chest X-ray dataset did not  
114 have such annotations, masks were generated using TorchXRyVision [3]. All experiments were  
115 performed on one NVIDIA A100 and three NVIDIA A40 GPUs.

116 **Training and Generation** The experimental pipeline is designed to identify the role of image  
117 dataset size in the image generation fidelity of selected GANs. To that end, for each GAN training  
118 run, all parameters are held constant with the exception of the image set size, which were set to  
119 500, 1000, and 2500 images, randomly sampled from the same dataset for each experimental run,  
120 respectively. For StyleGAN 3, all experimental runs were trained for 100 epochs; for SPADE-GAN,  
121 50 epochs. The trained adversarial network is then used to generate synthetic images, the fidelity of  
122 which is then evaluated for each training set size.

123 **Evaluation** The Fréchet Inception Distance (FID) [10] is a common metric used to evaluate the  
124 fidelity of the synthetically generated images for GANs [1]. Defined as the distance between the  
125 distributions of ground truth and the generated images respectively, in our paper we use the FID to  
126 assess the performance of the GAN (i.e. image fidelity) for each experimental run across both GANs.  
127 **A lower FID score signifies that a GAN is more proficient at generating synthetic data close to  
128 its target distribution.** From these data, we obtain fidelity curves for each dataset that describe how  
129 FID scores trend with increasing training set size.

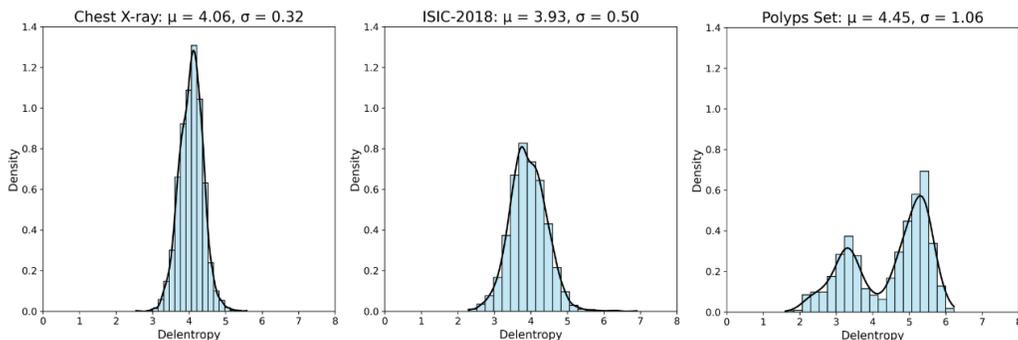


Figure 2: Delentropy distributions across each medical image dataset. A higher mean delentropy  $\mu$  indicates a dataset with more complex images.

## 130 4 Experimental Results

131 **Datasets** We employ three medical imaging datasets: International Skin Imaging Collaboration  
132 2018 Challenge (ISIC-2018) [2], Chest X-Ray Images (Chest X-ray) [12], and Colonoscopy Polyp

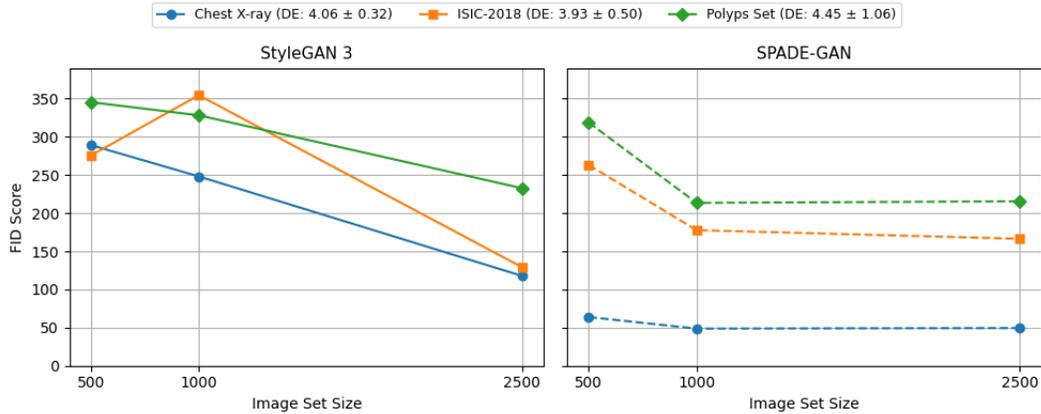


Figure 3: Fréchet Inception Distance (FID) curves comparing StyleGAN 3 and SPADE-GAN across each medical image dataset with varying sample sizes. Lower FID scores correspond to higher fidelity synthetic images.

133 Detection and Classification (Polyps Set) [31]. These datasets were chosen for their diversity in both  
 134 perceptual complexity, ranging from relatively skin lesions to complex colon polyps, and imaging  
 135 modality (dermoscopy vs. x-ray vs. colonoscopy).

136 We carry out delentropy calculations as described in Section 3 by using a publicly available imple-  
 137 mentation from Marchesoni [18]. To effectively capture the overall complexity of each image dataset,  
 138 we capture each dataset’s delentropy distribution as displayed in Fig. 2.

139 Across the experimental runs, FID scores consistently decreased with increasing dataset size. On  
 140 StyleGAN 3, synthesized images that had been generated by a GAN trained on 2500 images exhibited  
 141 an average FID score reduction of 48% when compared to those generated by a StyleGAN 3 that had  
 142 been trained on a mere 500 images (Fig. 3). SPADE-GAN experienced an analogous 31% FID score  
 143 reduction on average, though it is worth noting that FID reduction plateaued after only 1000 training  
 144 images.

145 Comparing both Fig. 2 and Fig. 3, one can see a general relationship between the delentropy distri-  
 146 bution and the training performance of both GANs. As the spread of image complexities increases  
 147 from a slender, peaked distribution to a broader, bimodal one, we see a corresponding increase in  
 148 FID scores for each dataset sample size. The Chest X-ray dataset with the most homogeneous image  
 149 complexities shown by a tall and narrow distribution, yields the lowest FID score after being trained  
 150 for 2500 images, indicating that both GANs had easier training runs with this dataset. On the contrast,  
 151 the Polyps Set—the dataset with the widest distribution and multiple complexity peaks—correlates  
 152 with the highest FID scores for each dataset sample size, which suggests that the GAN was faced  
 153 with a more challenging and unstable training run. Ultimately, this pattern shows a general inverse  
 154 relationship—GAN performance decreases with an increasing spread of image complexities within a  
 155 dataset.

## 156 5 Discussion

157 SPADE-GAN outperformed StyleGAN 3 across *all datasets and training sizes*, with FID scores  
 158 averaging 33% lower, likely due to its architecture that incorporates segmentation masks for structural  
 159 information, whereas StyleGAN 3 trained on raw image data alone, making it more difficult to  
 160 generalize to high-delentropy datasets. Furthermore, ISIC-2018 being an outlier can be attributed  
 161 to its fluctuations in image complexity, reflected by the standard deviation in delentropy (Fig. 2).  
 162 Despite having a lower mean delentropy, its spread likely resulted in difficulties in GAN training and  
 163 learning the images’ distribution, contrasting with the Chest X-ray dataset.

164 While the experimental results generally reflected an intuitive understanding of how image complexity  
 165 and training data influence GAN training, the FID curves provide insightful details, offering a deeper  
 166 perspective on these effects. SPADE-GAN exhibits both better quality results than StyleGAN 3 in the

167 form of lower FID scores and more consistent training as evidenced by the smooth, non-overlapping  
168 FID curves (Fig. 3). As aforementioned, performance plateaued after 1000 training images, suggest-  
169 ing that additional training data past that point may not help increase GAN performance as measured  
170 by FID score. This is also apparent in the generated images themselves, which exhibit little perceptual  
171 difference between those generated after 1000 training images and those generated after 2500 (Fig.  
172 1). Contrast this with the StyleGAN 3 curves, which do not reach any noticeable plateau between 500  
173 and 2500 training images. In fact, the increasingly negative slope values of the StyleGAN 3 graphs  
174 imply that StyleGAN 3 begins to better capture the images' features at a point past 1000 images, the  
175 exact whereabouts of which would need to be determined by a separate study.

176 The FID curves generated by this set of experiments set up a useful benchmark to which other  
177 potential training image data sets can be compared. For training sets that are of similar delentropy  
178 distributions and used to train StyleGAN 3 or SPADE-GAN, it is not unreasonable to predict that  
179 their training curves will be similar to those represented in Fig. 3, though many more training set  
180 sizes and image sets are required before a truly comprehensive representation can be reached.

181 **Broader Impacts** Our research on GANs for medical image synthesis may have positive and  
182 negative societal implications. On the positive side, it can enhance healthcare outcomes by improving  
183 the training of machine learning models with realistic synthetic data, therefore protecting patient  
184 privacy. Contrarily, potential negative impacts include the risk of malicious use for generating  
185 fraudulent synthetic data and possibility of reinforcing biases due to a lack of diversity of representing  
186 patient populations. These considerations demonstrate the importance of addressing both the benefits  
187 and potential risks associated with the use of GANs in the medical domain.

## 188 6 Conclusion

189 In this work, we highlight the impact of image complexity on GAN performance in medical image  
190 synthesis. We empirically demonstrate the general inverse relationship of how higher image complex-  
191 ity leads to poorer image fidelity results and performance in GANs. By displaying FID curves, we  
192 show healthcare professionals the possibility for the use of our benchmarks to gauge an estimate of  
193 data training requirements to achieve desirable results based on image complexity.

## 194 7 Limitations

195 Due to limited resources, experiments were only run on 500, 1000, and 2500 training images, leading  
196 to coarse-grained results. An extended study with a larger range and finer-grained increments would  
197 better elucidate exactly how FID scores respond to changes in training image dataset size. The use of  
198 FID scores as an evaluation metric also has its limitations, not necessarily correlating with human  
199 perceptual interpretations, something that is extremely important in the medical field where human  
200 doctors are still largely the source of truth. Skandarani et al. [29] show that a lower FID score is not  
201 a good measure of how well synthetic images can perform on a downstream task. More research  
202 involving multiple evaluations of the same experimental setup is required.

## 203 References

- 204 [1] Ali Borji. Pros and cons of gan evaluation measures: New developments. *Computer Vision and Image*  
205 *Understanding*, 215:103329, 2022. ISSN 1077-3142. doi: <https://doi.org/10.1016/j.cviu.2021.103329>.  
206 URL <https://www.sciencedirect.com/science/article/pii/S1077314221001685>.
- 207 [2] Noel Codella, Veronica Rotemberg, Philipp Tschandl, M. Emre Celebi, Stephen Dusza, David Gutman,  
208 Brian Helba, Aadi Kallou, Konstantinos Liopyris, Michael Marchetti, Harald Kittler, and Allan Halpern.  
209 Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging  
210 collaboration (isic), 2018. URL <https://arxiv.org/abs/1902.03368>. Licensed under the Creative  
211 Commons Attribution-NonCommercial (CC-BY-NC) license.
- 212 [3] Joseph Paul Cohen, Joseph D. Viviano, Paul Bertin, Paul Morrison, Parsa Torabian, Matteo Guarrera,  
213 Matthew P Lungren, Akshay Chaudhari, Rupert Brooks, Mohammad Hashir, and Hadrien Bertrand.  
214 TorchXRyVision: A library of chest X-ray datasets and models. In *Medical Imaging with Deep Learning*,  
215 2022. URL <https://github.com/mlmed/torchxrayvision>.

- 216 [4] Peichao Gao, Zhilin Li, and Hong Zhang. Thermodynamics-based evaluation of various improved shannon  
217 entropies for configurational information of gray-level images. *Entropy*, 20(1), 2018. ISSN 1099-4300.  
218 doi: 10.3390/e20010019. URL <https://www.mdpi.com/1099-4300/20/1/19>.
- 219 [5] Andrew Gilbert, Maciej Marciniak, Cristobal Rodero, Pablo Lamata, Eigil Samset, and Kristin Mcleod.  
220 Generating synthetic labeled data from existing anatomical models: An example with echocardiography  
221 segmentation. *IEEE Transactions on Medical Imaging*, 40(10):2783–2794, 2021. doi: 10.1109/TMI.2021.  
222 3051806.
- 223 [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron  
224 Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing*  
225 *systems*, 27, 2014.
- 226 [7] Sarah Graham, Colin Depp, Ellen E Lee, Camille Nebeker, Xin Tu, Ho-Cheol Kim, and Dilip V Jeste.  
227 Artificial intelligence for mental health and mental illnesses: an overview. *Current psychiatry reports*, 21:  
228 1–18, 2019.
- 229 [8] Xu Guo and Yiqiang Chen. Generative ai for synthetic data generation: Methods, challenges and the future.  
230 *arXiv preprint arXiv:2403.04190*, 2024.
- 231 [9] Robert M. Haralick, K. Shanmugam, and Its’Hak Dinstein. Textural features for image classification. *IEEE*  
232 *Transactions on Systems, Man, and Cybernetics*, SMC-3(6):610–621, 1973. doi: 10.1109/TSMC.1973.  
233 4309314.
- 234 [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans  
235 trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural*  
236 *Information Processing Systems (NeurIPS)*, volume 30, 2017. URL [https://arxiv.org/abs/1706.](https://arxiv.org/abs/1706.08500)  
237 08500.
- 238 [11] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and  
239 improving the image quality of stylegan, 2020. URL <https://arxiv.org/abs/1912.04958>. Licensed  
240 under the Nvidia Source Code License.
- 241 [12] Daniel Kermany, Kang Zhang, and Michael Goldbaum. Labeled optical coherence tomography (oct)  
242 and chest x-ray images for classification, 2018. URL <https://doi.org/10.17632/rsobjbr9sj.2>.  
243 Licensed under the Creative Commons Attribution 4.0 International (CC BY 4.0) license.
- 244 [13] T. M. Khan, S. S. Naqvi, and E. Meijering. Leveraging image complexity in macro-level neural network  
245 design for medical image segmentation. *Scientific Reports*, 12, 2022. doi: 10.1038/s41598-022-26482-7.  
246 URL <https://doi.org/10.1038/s41598-022-26482-7>.
- 247 [14] Z. Faizal khan and Sultan Refa Alotaibi. Computerised segmentation of medical images using neural net-  
248 works and glm. In *2019 International Conference on Advances in the Emerging Computing Technologies*  
249 *(AECT)*, pages 1–5, 2020. doi: 10.1109/AECT47998.2020.9194196.
- 250 [15] Kieran G. Larkin. Reflections on shannon information: In search of a natural information-entropy for  
251 images, 2016. URL <https://arxiv.org/abs/1609.01117>.
- 252 [16] Shipeng Liu, Liang Zhao, Dengfeng Chen, and Zhanping Song. Contrastive learning for image complexity  
253 representation, 2024. URL <https://arxiv.org/abs/2408.03230>.
- 254 [17] Yingzhou Lu, Minjie Shen, Huazheng Wang, Xiao Wang, Capucine van Rechem, Tianfan Fu, and Wenqi  
255 Wei. Machine learning for synthetic data generation: A review, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2302.04062)  
256 2302.04062.
- 257 [18] Franco Marchesoni. How do i know if an image after image enhancement is better  
258 than before?, March 2023. URL [https://ai.stackexchange.com/questions/39483/](https://ai.stackexchange.com/questions/39483/how-do-i-know-if-image-after-image-enhancement-is-better-than-before-image-pre)  
259 [how-do-i-know-if-image-after-image-enhancement-is-better-than-before-image-pre](https://ai.stackexchange.com/questions/39483/how-do-i-know-if-image-after-image-enhancement-is-better-than-before-image-pre).
- 260 [19] Galia Nordon, Gideon Koren, Varda Shalev, Eric Horvitz, and Kira Radinsky. Separating wheat from chaff:  
261 joining biomedical knowledge and patient data for repurposing medications. In *Proceedings of the AAAI*  
262 *Conference on Artificial Intelligence*, volume 33, pages 9565–9572, 2019.
- 263 [20] Zhaoqing Pan, Weijie Yu, Xiaokai Yi, Asifullah Khan, Feng Yuan, and Yuhui Zheng. Recent progress on  
264 generative adversarial networks (gans): A survey. *IEEE access*, 7:36322–36333, 2019.
- 265 [21] A. Papoulis. Generalized sampling expansion. *IEEE Transactions on Circuits and Systems*, 24(11):  
266 652–654, 1977. doi: 10.1109/TCS.1977.1084284.

- 267 [22] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-  
268 adaptive normalization, 2019. URL <https://arxiv.org/abs/1903.07291>. Licensed under the Cre-  
269 ative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) license.
- 270 [23] Vasileios C. Pezoulas, Dimitrios I. Zaridis, Eugenia Mylona, Christos Androutsos, Kosmas Apostolidis,  
271 Nikolaos S. Tachos, and Dimitrios I. Fotiadis. Synthetic data generation methods in healthcare: A  
272 review on open-source tools and methods. *Computational and Structural Biotechnology Journal*, 23:  
273 2892–2910, 2024. ISSN 2001-0370. doi: <https://doi.org/10.1016/j.csbj.2024.07.005>. URL <https://www.sciencedirect.com/science/article/pii/S2001037024002393>.  
274
- 275 [24] SM Rahman, Sifat Ibtisum, Ehsan Bazgir, and Tumpa Barai. The significance of machine learning in  
276 clinical disease diagnosis: A review. *arXiv preprint arXiv:2310.16978*, 2023.
- 277 [25] Esther Robb, Wen-Sheng Chu, Abhishek Kumar, and Jia-Bin Huang. Few-shot adaptation of generative  
278 adversarial networks. *arXiv preprint arXiv:2010.11943*, 2020.
- 279 [26] Junyuan Shang, Tengfei Ma, Cao Xiao, and Jimeng Sun. Pre-training of graph augmented transformers for  
280 medication recommendation. *arXiv preprint arXiv:1906.00346*, 2019.
- 281 [27] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):  
282 379–423, 1948. doi: 10.1002/j.1538-7305.1948.tb01338.x.
- 283 [28] Zheyuan Ryan Shi, Claire Wang, and Fei Fang. Artificial intelligence for social good: A survey. *arXiv*  
284 *preprint arXiv:2001.01818*, 2020.
- 285 [29] Youssef Skandarani, Pierre-Marc Jodoin, and Alain Lalande. GANs for Medical Image Synthesis: An  
286 Empirical Study. *Journal of Imaging*, 9(3):69, 2023. ISSN 2313-433X. doi: 10.3390/jimaging9030069.  
287 URL <https://www.mdpi.com/2313-433X/9/3/69>.
- 288 [30] Youssef Skandarani, Pierre-Marc Jodoin, and Alain Lalande. Gans for medical image synthesis: An  
289 empirical study. *Journal of Imaging*, 9, 2023. doi: 10.3390/jimaging9030069. URL <https://www.mdpi.com/2313-433X/9/3/69>.  
290
- 291 [31] Guanghui Wang. Replication Data for: Colonoscopy Polyp Detection and Classification: Dataset Creation  
292 and Comparative Evaluations, 2021. URL <https://doi.org/10.7910/DVN/FCBU0R>. Licensed under  
293 the CC0 1.0 Universal (CC0 1.0) Public Domain Dedication.
- 294 [32] Yaxing Wang, Chenshen Wu, Luis Herranz, Joost Van de Weijer, Abel Gonzalez-Garcia, and Bogdan  
295 Raducanu. Transferring gans: generating images from limited data. In *Proceedings of the European*  
296 *conference on computer vision (ECCV)*, pages 218–234, 2018.
- 297 [33] Ryan Webster, Julien Rabin, Loic Simon, and Frédéric Jurie. Detecting overfitting of deep generative  
298 networks via latent recovery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
299 *Recognition*, pages 11273–11282, 2019.
- 300 [34] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing  
301 He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.

302 **A Raw FID Scores**

<b>Dataset</b>	<b>Image Set Size</b>	<b>StyleGAN 3</b>	<b>SPADE-GAN</b>
Chest X-ray	500	289.20	63.90
	1000	248.20	48.74
	2500	117.85	49.49
ISIC-2018	500	275.67	263.09
	1000	354.61	177.82
	2500	129.30	166.45
Polyps Set	500	345.42	318.76
	1000	328.27	213.54
	2500	232.78	215.61

## 303 **NeurIPS Paper Checklist**

304 The checklist is designed to encourage best practices for responsible machine learning research,  
305 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
306 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
307 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
308 towards the page limit.

309 Please read the checklist guidelines carefully for information on how to answer these questions. For  
310 each question in the checklist:

- 311 • You should answer [Yes], [No], or [NA].
- 312 • [NA] means either that the question is Not Applicable for that particular paper or the  
313 relevant information is Not Available.
- 314 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

315 **The checklist answers are an integral part of your paper submission.** They are visible to the  
316 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
317 (after eventual revisions) with the final version of your paper, and its final version will be published  
318 with the paper.

319 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
320 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
321 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
322 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
323 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
324 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
325 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
326 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
327 please point to the section(s) where related material for the question can be found.

328 **IMPORTANT, please:**

- 329 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- 330 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 331 • **Do not modify the questions and only use the provided macros for your answers.**

### 332 **1. Claims**

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

335 Answer: [Yes]

336 Justification: They state our aim of establishing benchmarks for the relationship between  
337 dataset size and image fidelity in GAN-generated medical images, considering image  
338 complexity distribution in a variety of medical image datasets.

339 Guidelines:

- 340 • The answer NA means that the abstract and introduction do not include the claims  
341 made in the paper.
- 342 • The abstract and/or introduction should clearly state the claims made, including the  
343 contributions made in the paper and important assumptions and limitations. A No or  
344 NA answer to this question will not be perceived well by the reviewers.
- 345 • The claims made should match theoretical and experimental results, and reflect how  
346 much the results can be expected to generalize to other settings.
- 347 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
348 are not attained by the paper.

### 349 **2. Limitations**

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

351 Answer: [Yes]

352 Justification: The paper includes a dedicated "Limitations" section (Section 7) that discusses  
353 the constraints of the study. For instance, our limited range of training image set sizes and  
354 the limitations of using FID scores as the primary evaluation metric.

355 Guidelines:

- 356 • The answer NA means that the paper has no limitation while the answer No means that  
357 the paper has limitations, but those are not discussed in the paper.
- 358 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 359 • The paper should point out any strong assumptions and how robust the results are to  
360 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
361 model well-specification, asymptotic approximations only holding locally). The authors  
362 should reflect on how these assumptions might be violated in practice and what the  
363 implications would be.
- 364 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
365 only tested on a few datasets or with a few runs. In general, empirical results often  
366 depend on implicit assumptions, which should be articulated.
- 367 • The authors should reflect on the factors that influence the performance of the approach.  
368 For example, a facial recognition algorithm may perform poorly when image resolution  
369 is low or images are taken in low lighting. Or a speech-to-text system might not be  
370 used reliably to provide closed captions for online lectures because it fails to handle  
371 technical jargon.
- 372 • The authors should discuss the computational efficiency of the proposed algorithms  
373 and how they scale with dataset size.
- 374 • If applicable, the authors should discuss possible limitations of their approach to  
375 address problems of privacy and fairness.
- 376 • While the authors might fear that complete honesty about limitations might be used by  
377 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
378 limitations that aren't acknowledged in the paper. The authors should use their best  
379 judgment and recognize that individual actions in favor of transparency play an impor-  
380 tant role in developing norms that preserve the integrity of the community. Reviewers  
381 will be specifically instructed to not penalize honesty concerning limitations.

### 382 3. Theory Assumptions and Proofs

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

385 Answer: [NA]

386 Justification: This paper is primarily based on empirical results and does not present  
387 theoretical results requiring proofs.

388 Guidelines:

- 389 • The answer NA means that the paper does not include theoretical results.
- 390 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
391 referenced.
- 392 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 393 • The proofs can either appear in the main paper or the supplemental material, but if  
394 they appear in the supplemental material, the authors are encouraged to provide a short  
395 proof sketch to provide intuition.
- 396 • Inversely, any informal proof provided in the core of the paper should be complemented  
397 by formal proofs provided in appendix or supplemental material.
- 398 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 399 4. Experimental Result Reproducibility

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

403 Answer: [Yes]

404 Justification: The paper provides a detailed, step-by-step process of the GAN pipeline  
405 including the information on the datasets used, preprocessing steps, the number of training  
406 image sample size, training iterations, GAN architectures, and evaluation metrics.

407 Guidelines:

- 408 • The answer NA means that the paper does not include experiments.
- 409 • If the paper includes experiments, a No answer to this question will not be perceived  
410 well by the reviewers: Making the paper reproducible is important, regardless of  
411 whether the code and data are provided or not.
- 412 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
413 to make their results reproducible or verifiable.
- 414 • Depending on the contribution, reproducibility can be accomplished in various ways.  
415 For example, if the contribution is a novel architecture, describing the architecture fully  
416 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
417 be necessary to either make it possible for others to replicate the model with the same  
418 dataset, or provide access to the model. In general, releasing code and data is often  
419 one good way to accomplish this, but reproducibility can also be provided via detailed  
420 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
421 of a large language model), releasing of a model checkpoint, or other means that are  
422 appropriate to the research performed.
- 423 • While NeurIPS does not require releasing code, the conference does require all submis-  
424 sions to provide some reasonable avenue for reproducibility, which may depend on the  
425 nature of the contribution. For example
  - 426 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
427 to reproduce that algorithm.
  - 428 (b) If the contribution is primarily a new model architecture, the paper should describe  
429 the architecture clearly and fully.
  - 430 (c) If the contribution is a new model (e.g., a large language model), then there should  
431 either be a way to access this model for reproducing the results or a way to reproduce  
432 the model (e.g., with an open-source dataset or instructions for how to construct  
433 the dataset).
  - 434 (d) We recognize that reproducibility may be tricky in some cases, in which case  
435 authors are welcome to describe the particular way they provide for reproducibility.  
436 In the case of closed-source models, it may be that access to the model is limited in  
437 some way (e.g., to registered users), but it should be possible for other researchers  
438 to have some path to reproducing or verifying the results.

## 439 5. Open access to data and code

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

443 Answer: [Yes]

444 Justification: The paper provides open access to both the data and code in the supplementary  
445 material. Sufficient instructions are included to allow correct reproduction of the main  
446 experimental results.

447 Guidelines:

- 448 • The answer NA means that paper does not include experiments requiring code.
- 449 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
450 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 451 • While we encourage the release of code and data, we understand that this might not be  
452 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
453 including code, unless this is central to the contribution (e.g., for a new open-source  
454 benchmark).
- 455 • The instructions should contain the exact command and environment needed to run to  
456 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
457 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

- 458 • The authors should provide instructions on data access and preparation, including how  
459 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 460 • The authors should provide scripts to reproduce all experimental results for the new  
461 proposed method and baselines. If only a subset of experiments are reproducible, they  
462 should state which ones are omitted from the script and why.
- 463 • At submission time, to preserve anonymity, the authors should release anonymized  
464 versions (if applicable).
- 465 • Providing as much information as possible in supplemental material (appended to the  
466 paper) is recommended, but including URLs to data and code is permitted.

## 467 6. Experimental Setting/Details

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

471 Answer: [Yes]

472 Justification: The paper specifies the datasets used, preprocessing steps, GAN architectures  
473 (StyleGAN 3 and SPADE-GAN), number of training images (500, 1000, 2500), and number  
474 of epochs for each GAN, outlined in the methodology (Section. 3).

475 Guidelines:

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

## 481 7. Experiment Statistical Significance

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

484 Answer: [Yes]

485 Justification: This paper presents FID scores for different experimental conditions, and it  
486 also reports the variability in the image complexities of each medical imaging dataset as  
487 shown in Fig. 2.

488 Guidelines:

- 489 • The answer NA means that the paper does not include experiments.
- 490 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
491 dence intervals, or statistical significance tests, at least for the experiments that support  
492 the main claims of the paper.
- 493 • The factors of variability that the error bars are capturing should be clearly stated (for  
494 example, train/test split, initialization, random drawing of some parameter, or overall  
495 run with given experimental conditions).
- 496 • The method for calculating the error bars should be explained (closed form formula,  
497 call to a library function, bootstrap, etc.)
- 498 • The assumptions made should be given (e.g., Normally distributed errors).
- 499 • It should be clear whether the error bar is the standard deviation or the standard error  
500 of the mean.
- 501 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
502 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
503 of Normality of errors is not verified.
- 504 • For asymmetric distributions, the authors should be careful not to show in tables or  
505 figures symmetric error bars that would yield results that are out of range (e.g. negative  
506 error rates).
- 507 • If error bars are reported in tables or plots, The authors should explain in the text how  
508 they were calculated and reference the corresponding figures or tables in the text.

## 509 8. Experiments Compute Resources

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

513 Answer: [Yes]

514 Justification: The paper says that experiments were performed on "one NVIDIA A100 and  
515 three NVIDIA A40 GPUs.

516 Guidelines:

- 517 • The answer NA means that the paper does not include experiments.
- 518 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
519 or cloud provider, including relevant memory and storage.
- 520 • The paper should provide the amount of compute required for each of the individual  
521 experimental runs as well as estimate the total compute.
- 522 • The paper should disclose whether the full research project required more compute  
523 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
524 didn't make it into the paper).

## 525 9. Code Of Ethics

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

528 Answer: [Yes]

529 Justification: The paper conforms with the NeurIPS Code of Ethics such as providing our  
530 publicly used datasets.

531 Guidelines:

- 532 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 533 • If the authors answer No, they should explain the special circumstances that require a  
534 deviation from the Code of Ethics.
- 535 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
536 eration due to laws or regulations in their jurisdiction).

## 537 10. Broader Impacts

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

540 Answer: [Yes]

541 Justification: This paper explains both the positive and negative broader effects under the  
542 discussion (Section 5).

543 Guidelines:

- 544 • The answer NA means that there is no societal impact of the work performed.
- 545 • If the authors answer NA or No, they should explain why their work has no societal  
546 impact or why the paper does not address societal impact.
- 547 • Examples of negative societal impacts include potential malicious or unintended uses  
548 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
549 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
550 groups), privacy considerations, and security considerations.
- 551 • The conference expects that many papers will be foundational research and not tied  
552 to particular applications, let alone deployments. However, if there is a direct path to  
553 any negative applications, the authors should point it out. For example, it is legitimate  
554 to point out that an improvement in the quality of generative models could be used to  
555 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
556 that a generic algorithm for optimizing neural networks could enable people to train  
557 models that generate Deepfakes faster.
- 558 • The authors should consider possible harms that could arise when the technology is  
559 being used as intended and functioning correctly, harms that could arise when the  
560 technology is being used as intended but gives incorrect results, and harms following  
561 from (intentional or unintentional) misuse of the technology.

- 562 • If there are negative societal impacts, the authors could also discuss possible mitigation  
563 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
564 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
565 feedback over time, improving the efficiency and accessibility of ML).

## 566 11. Safeguards

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

570 Answer: [NA]

571 Justification: This paper does not release new models or datasets and poses no such risks.

572 Guidelines:

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

## 583 12. Licenses for existing assets

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

587 Answer: [Yes]

588 Justification: This paper properly credits used GAN models and datasets and can be found  
589 as a note in the reference of each asset.

590 Guidelines:

- 591 • The answer NA means that the paper does not use existing assets.
- 592 • The authors should cite the original paper that produced the code package or dataset.
- 593 • The authors should state which version of the asset is used and, if possible, include a  
594 URL.
- 595 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 596 • For scraped data from a particular source (e.g., website), the copyright and terms of  
597 service of that source should be provided.
- 598 • If assets are released, the license, copyright information, and terms of use in the  
599 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
600 has curated licenses for some datasets. Their licensing guide can help determine the  
601 license of a dataset.
- 602 • For existing datasets that are re-packaged, both the original license and the license of  
603 the derived asset (if it has changed) should be provided.
- 604 • If this information is not available online, the authors are encouraged to reach out to  
605 the asset's creators.

## 606 13. New Assets

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

609 Answer: [NA]

610 Justification: This paper does not introduce new assets such as datasets or models.

611 Guidelines:

- 612 • The answer NA means that the paper does not release new assets.

- 613 • Researchers should communicate the details of the dataset/code/model as part of their  
614 submissions via structured templates. This includes details about training, license,  
615 limitations, etc.
- 616 • The paper should discuss whether and how consent was obtained from people whose  
617 asset is used.
- 618 • At submission time, remember to anonymize your assets (if applicable). You can either  
619 create an anonymized URL or include an anonymized zip file.

#### 620 14. Crowdsourcing and Research with Human Subjects

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

624 Answer: [NA]

625 Justification: This paper does not involve any crowdsourcing or research with human subjects

626 Guidelines:

- 627 • The answer NA means that the paper does not involve crowdsourcing nor research with  
628 human subjects.
- 629 • Including this information in the supplemental material is fine, but if the main contribu-  
630 tion of the paper involves human subjects, then as much detail as possible should be  
631 included in the main paper.
- 632 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
633 or other labor should be paid at least the minimum wage in the country of the data  
634 collector.

#### 635 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 636 Subjects

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

641 Answer: [NA]

642 Justification: This paper does not involve any research with human subjects that would  
643 require IRB approval.

644 Guidelines:

- 645 • The answer NA means that the paper does not involve crowdsourcing nor research with  
646 human subjects.
- 647 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
648 may be required for any human subjects research. If you obtained IRB approval, you  
649 should clearly state this in the paper.
- 650 • We recognize that the procedures for this may vary significantly between institutions  
651 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
652 guidelines for their institution.
- 653 • For initial submissions, do not include any information that would break anonymity (if  
654 applicable), such as the institution conducting the review.