
Domain-Compatible Synthetic Data Generation for Infrequent Objects Detection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The recent advances in generative models have resulted in massive progress in
2 the quality of the generated images to the point that in many cases they cannot be
3 easily distinguished from real images. Despite this quality improvement, using AI
4 generated images for the purpose of training robust down-stream computer vision
5 models for real-world applications has proven to be very challenging. The AI
6 generated images usually lack the required diversity and scene complexity that is
7 crucial for many real-world applications, specifically the ones with safety concerns.
8 The difficulty of this challenge grows significantly when the underlying application
9 involves detection of some specific objects that appear with critically low frequency
10 in the available real datasets. This paper studies a new approach for generating
11 diverse, complex and domain-compatible synthetic images for detecting infrequent
12 objects by employing a diffusion-based generative model pretrained on a generic
13 dataset. More specifically, the impact of using the generated synthetic images with
14 the proposed approach in solving the real world problem of detecting emergency
15 vehicles in road scenes is investigated. Furthermore, the challenges of generating
16 synthetic datasets with the proposed approach will be thoroughly discussed.

17

1 Introduction

18 Detection of some domain specific and infrequent objects can be a crucial part of many computer
19 vision based systems. An example of such scenario is the detection of emergency vehicles for an
20 autonomous driving car application. Since the number of images containing the specific objects of
21 interest in the available datasets is critically limited, generating supplementary synthetic images is a
22 viable solution for training robust downstream object detection models.

23 Employing deep generative models to generate synthetic images for training downstream models in a
24 real-world application imposes some key challenges listed as follows:

25 **Insufficient samples to train the generative model** A deep generative model relies on a large
26 training dataset covering different varieties of the object of interest to be able to generate realistic
27 images. In the case of infrequent objects, the lack of sufficient training images is the reason synthetic
28 images are required in the first place.

29 **Insufficient diversity and scene complexity** The majority of of recent advancements in improving
30 the performance of generative models have been focused on enhancing the quality of the generated
31 images and making them more photo-realistic. The AI-generated images usually lack the required
32 scene complexity and diversity which is essential for training robust downstream models (Block et al.
33 2006).For the same reason there is normally a distribution shift between the generated images and the
34 real ones in terms of complexity and diversity (Joshi 2019).

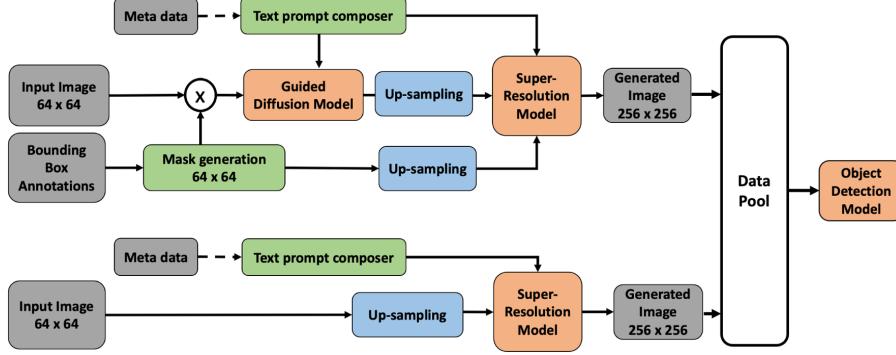


Figure 1: Block diagram of the architecture of the proposed approaches.

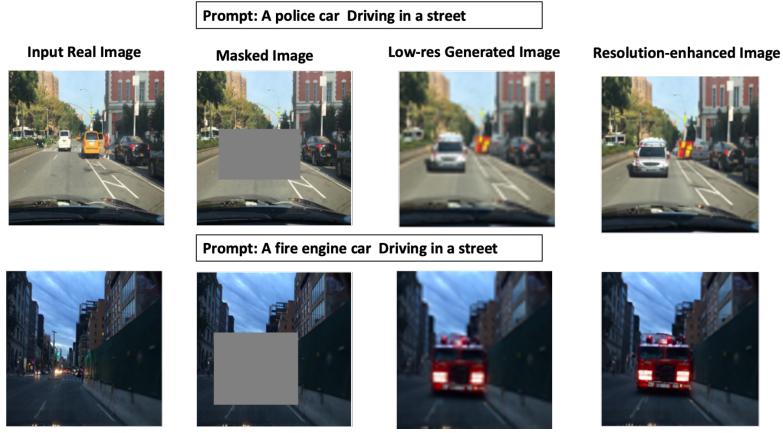


Figure 2: A few examples of input and output endpoints for Approach 1 .

35 **Generated images may require labeling** As opposed to synthetic images generated by rendering
 36 engines, AI-generated images may require an annotation process to be ready for a real application.

37 In order to tackle the above challenges and generate synthetic images that can be effectively used in
 38 real-world applications, in this paper we investigate three different approaches of using a generative
 39 model that has been only trained on a generic dataset. The proposed approaches can be used to
 40 generate a large, complex and widely diverse dataset from a small relevant real dataset. We use
 41 a diffusion-based model (Ho, Jain, and Abbeel 2020) (Kim and Ye 2021) (Dhariwal and Nichol
 42 2021) that can be conditioned on different information and be partially masked during the generative
 43 process to make carefully controlled changes in the real images in a systematic way. This allows
 44 the generation of a sufficiently large domain-compatible dataset that covers the required variety and
 45 complexity for training a robust downstream model. Since the proposed approach uses real images
 46 as the basis to create the synthetic images, there is no domain-shift between the generated images
 47 and the real dataset. Conditioning the generative process on a set of guiding text prompts as well as
 48 partially masking specific parts of the image during the process allows imposing a customized level of
 49 diversity while maintaining the domain characteristics and scene complexity of the real images. The
 50 proposed approach also allows either preserving the available annotations or automatically generating
 51 new annotations for the synthetically generated objects. We run several experiments to extensively
 52 assess the performance enhancement that the generated images provide to the final downstream object
 53 detection models.

54 2 Related Work

55 A wide variety of approaches have been investigated for synthetic data creation including simple
 56 rule-based algorithms, statistical models, computer simulations and data augmentation techniques.

57 One of the most commonly used approaches to generate synthetic image data is through the use of
58 photo-realistic 3D physics engines (Pollock et al. 2019) (Cisse et al. 2017). These engines can be
59 used to render images from 3D computer-aided design (CAD) models of the target objects. The
60 photo-realism achieved through these image rendering engines has reached a point where synthetic
61 images can be hardly distinguished from real ones (Hamesse et al. 2019). However, there are some
62 drawbacks to these synthetic data generation approaches that make them unsuitable for many practical
63 applications. These include, but are not limited to, requiring 3D asset development, challenges in
64 tuning design parameters (e.g. brightness) and lack of the required diversity and complexity in the
65 image background.

66 Deep generative models including generative adversarial networks (GANs) have been vastly studied
67 for synthetic image generation and synthetic augmentation (Vega-Marquez et al. 2019)(Frid-Adar et
68 al. 2018). In the field of medical imaging, GAN-based data augmentation has particularly been used
69 to improve sensitivity and specificity of models tried on small medical imaging datasets by 5-7%
70 (Bowles et al. 2018) (Frid-Adar et al. 2018). Class imbalance has been addressed by generating
71 additional examples of infrequent samples through adversarial autoencoders, a GAN variant (Lim et
72 al. 2018). More over, deep learning based style transfer has shown 2% improvements in classification
73 accuracy over traditional augmentation strategies (Zheng et al. 2019). Style transfer, in particular, is
74 capable of preserving image content while copying the style of a separate, unrelated image (Gatys,
75 Ecker, and Bethge 2015).

76 Denoising diffusion models were initially introduced by (Sohl-Dickstein et al. 2015). Recent work
77 has demonstrated the ability of diffusion models to compete and potentially outperform traditional
78 generative adversarial networks in realistic image generation and producing synthetic results indistin-
79 guishable from real images to human evaluators in some cases (Dhariwal and Nichol 2021) (Zhou et
80 al. 2019).

81 **3 Methodology**

82 In the proposed methodology for syntehetic image generation, first a pretrained diffusion model
83 (Dhariwal and Nichol 2021) (Nichol et al. 2021) is fine-tuned on a generic dataset which does not
84 necessarily include the infrequent target objects (we used a generic driving dataset (Yu et al. 2020)).
85 In order to condition the diffusion process on text, we use a CLIP model (Radford et al. 2021)
86 that perturbs the denoising process mean with the gradient of the dot product of the image and text
87 encoding with respect to the image. Next, we explore three different image manipulation approaches
88 with this model that allows generating synthetic images that contain a large variety of infrequent
89 objects of interest. These synthetic images are then used for training downstream object detection
90 models as shown in Figure 1. Finally, a text-conditioned super-resolution diffusion model is cascaded
91 with the generative model in the pipeline to increase the resolution of the generated images. The
92 proposed approaches are based on the assumption that a very small but domain-relevant real dataset
93 is available and synthetic images are generated by manipulating those real images. In fact, using this
94 small real data as the basis is essential in keeping the generated images in the target domain. In this
95 section, the three proposed image manipulation approaches will be explained in detail.

96 **3.1 Approach 1: Synthetic Infrequent Objects in a Real Background**

97 The idea behind this approach which is depicted in the upper part of Figure 1, is to generate instances
98 of the infrequent objects of interest inside a background sampled from the real dataset to maintain the
99 generated images in the same domain as the real dataset. The importance of this approach is that it can
100 be employed to generate a sufficiently large synthetic dataset even if the real dataset does not include
101 any images containing the infrequent target objects. The architecture of this approach consists of
102 four main components: A mask generator block, a text prompt composer unit, a text guided diffusion
103 generative model and a super-resolution model. The input image serving as background and the
104 corresponding annotations are first fed to a mask generator block which proposes a mask based on
105 the current bounding boxes in the image. The generated mask is then applied to the original image
106 and the resulted masked image is fed to the text conditioned diffusion model. The diffusion model
107 iteratively manipulates the masked part of the image following the input text prompt guidance until it
108 generates an instance of the target object inside the masked section which is well blended with the
109 background. The output of this model is then fed to a diffusion-based super-resolution model (Nichol



Figure 3: An examples of the steps of Approach 2 .



Figure 4: An example of changing the weather condition in the real image using Approach 3 .

110 and Dhariwal 2021) to enhance its resolution. The super-resolution model can also be conditioned
111 on the text prompt for improved enhancement. Figure 2 illustrates a few examples of the inputs and
112 output endpoints of the pipeline of this approach.

113 In the rest of this subsection, the mask generator and prompt composer blocks are described.

114 3.1.1 Mask generator block

115 This block proposes a region for masking the input image based on the available bounding boxes
116 in the annotations. In order to find a proper area for the placement of the target object, one or more
117 adjacent bounding boxes are randomly picked and merged together to make a target bounding box
118 while the following rules are met:

- 119 • The proposed bounding box should not cut any of the other bounding boxes to avoid
120 unrealistic coincidences between the generated objects and the ones in the background.
- 121 • If needed, the orientation of the bounding box should be compatible with the required
122 object alignment. Usually the orientation of the bounding box dictates the orientation of the
123 generated object and can be used as an additional factor for randomization.

124 Other customized rules can be easily integrated in this framework depending on the target application.

125 3.1.2 Text prompt composer unit

126 This block composes a text prompt to guide the diffusion process toward generating the desired target
127 image. Each composed prompt consists of five main components as follows:

128 **Subject** In approach 1, subject is randomly sampled from the list of infrequent target objects.

129 **Verb** Verb is randomly sampled from a list of possible actions relevant to the target object. For
130 example for a driving scene dataset, the possible verbs can be driving, crossing, parking, etc.

131 **Location** Represents the location of the target object in the image and it can be either extracted from
132 meta data (approach 1) or randomly sampled from possible options (approach 2).

133 **Condition** This field describes a global condition for the image. For example for a road scene dataset
134 this field can describe the weather condition, e.g. rainy, snowy, foggy, etc.

135 **Time** Optionally describes the time of day, e.g. morning, night, sunset, etc.



Figure 5: Examples demonstrating challenges with text condition image generation approaches.

136 3.2 Approach 2: Real Infrequent Objects in Synthetic Background

137 This approach can be also represented by the top part of 1. However instead of generating target
 138 objects in a real background, it generates a synthetic background for a real target object. The target
 139 object is first cropped from a real image and after random resizing is placed in a random position in a
 140 blank (all zeros) background. The resulted combinations is then fed to the diffusion model. There are
 141 two important differences between this approach and approach 1:

- 142 1. As opposed to approach 1, in this approach the mask only covers the real object and leaves
 143 everywhere else in the image available for the diffusion model’s generative manipulation.
 144 This results in the generation of a background that follows the text prompt guidance and
 145 blends well with the real object.
- 146 2. In this approach, the prompt composer unit randomly samples all of the background-related
 147 fields such as verb, location, condition and time from the the corresponding lists that are
 148 provided to the module based on the target application. The only field that will be extracted
 149 from the annotation is the type of target object that has been cropped from the real image.

150 Figure 3 illustrates the steps of this approach in an example.

151 3.3 Approach 3: Real Images Globally Altered

152 The third approach is represented by the bottom part of the block diagram in Figure 1. In this
 153 approach, certain aspects of the real images are altered as they are converted from low to high
 154 resolution by conditioning the super-resolution model to text prompts that guide the diffusion process
 155 toward those modifications. As suggested by the diagram, in this approach no masking is required as
 156 the entire input image is subject to the model’s subtle modifications. In order to propose suitable text
 157 prompts for randomized modifications to input images, the text composer unit randomly samples the
 158 condition field from a list of application-relevant conditions while rest of the fields are extracted form
 159 the annotations or meta-data if it is available. For example, multiple altered versions of an input real
 160 image can be generated synthetically by randomizing on weather condition or the time of the day.
 161 Figure 4 shows some examples of these modifications along with their corresponding text prompts.

162 4 Dataset

163 As explained in the previous section, the proposed approaches use a small real dataset as the basis to
 164 create the synthetic images. In this section, we introduce the real dataset that was used as a base for
 165 generating the synthetic images in all of our experiments.

166 4.1 LAVA and LAVA-emergency datasets

167 The LISA-Amazon Vehicle and Scene Attributes (LAVA) dataset (Ninad et al. 2021) has been
 168 collected as a part of a collaboration between the Amazon Machine Learning Solutions Lab with the

Table 1: The distribution of images and bounding boxes for the real and synthetic datasets.

| Dataset | Num. images | Medical | Fire | Police |
|------------|-------------|---------|------|--------|
| Real-Train | 215 | 47 | 42 | 126 |
| Real-Test | 539 | 270 | 68 | 215 |
| Type-1 | 1876 | 447 | 569 | 939 |
| Type-2 | 1875 | 620 | 306 | 949 |

Table 2: Downstream object detection performance for each dataset.

| Model and backbone | Dataset | Num. train images | mAP@0.50:0.95 | mAR@0.50:0.95 |
|---------------------------------|-----------|-------------------|---------------|---------------|
| SSD ResNet101 V1 FPN | R | 215 | 0 | 0.028 |
| SSD ResNet101 V1 FPN | R, S1 | 2091 | 0.147 | 0.441 |
| SSD ResNet101 V1 FPN | R, S2 | 2090 | 0.396 | 0.59 |
| SSD ResNet101 V1 FPN | R, S1, S2 | 3966 | 0.372 | 0.586 |
| SSD MobileNet V1 FPN | R | 215 | 0 | 0.095 |
| SSD MobileNet V1 FPN | R, S1 | 2091 | 0.129 | 0.331 |
| SSD MobileNet V1 FPN | R, S2 | 2090 | 0.475 | 0.637 |
| SSD MobileNet V1 FPN | R, S1, S2 | 3966 | 0.357 | 0.583 |
| EfficientDet D1 | R | 215 | 0.053 | 0.439 |
| EfficientDet D1 | R, S1 | 2091 | 0.136 | 0.523 |
| EfficientDet D1 | R, S2 | 2090 | 0.368 | 0.594 |
| EfficientDet D1 | R, S1, S2 | 3966 | 0.458 | 0.641 |
| Faster RCNN Inception ResNet V2 | R | 215 | 0.173 | 0.451 |
| Faster RCNN Inception ResNet V2 | R, S1 | 2091 | 0.454 | 0.723 |
| Faster RCNN Inception ResNet V2 | R, S2 | 2090 | 0.521 | 0.695 |
| Faster RCNN Inception ResNet V2 | R, S1, S2 | 3966 | 0.494 | 0.714 |

169 Laboratory of Intelligent and Safe Automobiles at the University of California, San Diego (UCSD) to
 170 build a large and richly annotated driving dataset with fine-grained vehicle, pedestrian, and scene
 171 attributes.

172 The LAVA dataset is annotated for all types of vehicles, traffic signs, traffic lights and pedestrians
 173 with 2D bounding boxes, class labels and some meta data. A subset of the LAVA dataset that covers
 174 all the images with emergency vehicles in them (in addition to other vehicles) was separated and used
 175 for generating synthetic images and training the downstream object detection models. We refer to
 176 this subset as LAVA-emergency dataset.

177 Table 1 shows the class distribution of the train and test splits of the LAVA-emergency dataset. It is
 178 essential to reserve a reasonable portion of the real dataset for testing to be able to reliably evaluate
 179 the impact of synthetic data generation approaches.

180 5 Experiments

181 5.1 Experimental Setup

182 For all of the experiments in this section, the LAVA-emergency dataset is used as a base for generating
 183 synthetic images using the approaches in section 3. The downstream task in our experiments is the
 184 detection of emergency vehicles including medical vehicles (ambulances), fire engines and police
 185 cars. These emergency vehicles appear with a critically low frequency in the road-scene datasets.
 186 For better understanding of the evaluation results, we group the synthetic data generation techniques
 187 into three general types. Type-1 (S1), represents the approaches wherein the emergency vehicles
 188 themselves are synthetically generated (only Approach 1). Type-2 (S2) represents all the approaches
 189 wherein the emergency vehicles are real but they have been placed in a synthetically generated or
 190 modified background (approach 2 and approach 3). Table 1 shows the distribution of generated data
 191 over different emergency vehicles categories.

192 In these experiments, for composing the text prompts, the weather condition is randomly and
 193 uniformly sampled from a list of 5 weather conditions namely, sunny, rainy, snowy , foggy and cloudy.

194 The location of the vehicle is randomly sampled from one of four options: street, road (each with
195 a probability of 0.35), parking (with a probability of 0.25) and bridge (with a probability of 0.05).
196 Each synthetic image is generated by applying 100 diffusion steps to the masked real input image (in
197 Approach 1 and 2). The resolution of the generated images is then enhanced by applying 30 addition
198 diffusion steps through the super-resolution model. Each experiment uses either only real data (R) or
199 a combination of it with one or more types of synthetic images. The objective of these experiments is
200 to evaluate how each of the synthetic data generation approaches improves the performance of the
201 downstream object detection models when combined with the real data.

202 5.2 Results

203 Table 2 shows the performance of various object detection algorithms trained on difference combinations
204 of real and synthetic images on the emergency-LAVA test set. As shown in this table
205 the single-stage detectors such as different flavors of SSD and EfficientDet are barely able to learn
206 anything from the small real training set. However, incrementally adding synthetic images to augment
207 the real training images remarkably improves the detector’s performance on the real test set.

208 The EfficientDet D1 model has monotonically increasing mAP and mAR as more synthetic data
209 is added. For SSD ResNet101, SSD MobileNet and Faster R-CNN models there is a considerable
210 performance improvement when trained on R, S1 or R, S2 compared to when they are only trained
211 on R. However, for these models, there is a slight drop in performance when they are trained on R,
212 S1, S2 compared to when they are trained on R, S2. As mentioned in section 3.1, in synthetic Type-1
213 images the emergency vehicles themselves are generated by the model and the generator model has
214 been trained on a generic dataset which contains vehicles from a variety of different countries in
215 the world. The LAVA-emergency test set however contains only emergency vehicles from Southern
216 California, and thus the discrepancy in performance when involving S1 in training along with S2
217 can be explained by the change in emergency vehicles characteristics from different geo-locations.
218 However, in S2 images, the emergency vehicles have been directly adopted from emergency-LAVA
219 training set and they are compatible with the emergency vehicles in the test set. Therefore, increasing
220 the number of Type-2 images always improves the performance of all of the object detection models.
221 This is seen experimentally when comparing the results from models trained with R, S1 and R, S2.
222 The models trained with R, S2 have consistently higher performance than those trained on R, S1 for
223 all models and backbones.

224 5.3 Practical Challenges

225 Although the synthetically generated images by the proposed approaches are realistic and diverse,
226 there are a few challenges that need to be considered depending on the target application as follows:

227 **Relative size of the objects** When an image generation process is conditioned on text, sometimes
228 the relative sizes of the generated objects can be slightly out of proportionate with respect to the
229 background objects, regardless of the type of the generative model. While some downstream vision
230 tasks such as object detection are not negatively impacted by this, some others may be impacted. The
231 top row of Figure 1 shows a few examples with slightly disproportionate objects.

232 **The number of the objects** One of the concepts that normally do not transfer properly between
233 language and vision spaces is the exact quantity of objects. Similar to the previous case, the exact
234 number of objects does not impact many of the vision tasks (e.g. object detection).

235 **The relative position of the objects** Similar to relative sizes of objects, their relative positions with
236 respect to each other can sometimes be unrealistic when the generative process is conditioned on text.
237 The bottom row of Figure 5 shows a few examples impacted by this effect.

238 6 Conclusions

239 In this work, a new approach for generating synthetic data for training downstream models in a
240 critically low data regime was studied. The experimental results showed that employing the synthetic
241 images generated by the proposed approach significantly improved the performance of all of the
242 investigated object detection models. Employing approaches similar to the proposed approach to

243 augment insufficiently small real datasets used in training the downstream computer vision models is
244 specifically crucial for applications with safety concerns.

245 **References**

- 246 [1] Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B. & Chen, M. (2021) Glide: Towards
247 photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*.
- 248 [2] Joshi, C. (2019) Generative adversarial networks (GANs) for synthetic dataset generation with binary classes.
- 249 [3] Ho, J., Jain, A. & Abbeel, P. (2020) Denoising diffusion probabilistic models. Advances in *Neural Information
250 Processing Systems* 33 pp. 6840–6851.
- 251 [4] Kim, G. & Ye, J.C. (2021) Diffusionclip: Text-guided image manipulation using diffusion models.
- 252 [5] Dhariwal, P. & Nichol, A. (2021) Diffusion models beat gans on image synthesis. *Advances in Neural
253 Information Processing Systems* 15 pp. 8780–8794
- 254 [6] Nichol, A. Q. & Dhariwal, P. (2021) Improved denoising diffusion probabilistic models. In International
255 Conference on Machine Learning, pp. 8162–8171. PMLR.
- 256 [7] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021) Learning
257 transferable visual models from natural language supervision. In *International Conference on Machine Learning*
258 pp. 8748–8763. PMLR.
- 259 [8] Yu, F., Chen, H., Wang, X., Xian, W., Chen, Y., Liu, F., ... & Darrell, T. (2020) Bdd100k: A diverse driving
260 dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision
261 and pattern recognition* pp. 2636–2645.
- 262 [9] Ninad, K., Akshay, R., Jonathan, B., Jeremy, F., Mohan, T., Nachiket, D., Greer, R., Saman, S. & Suchitra,
263 S. (2021) Create a large-scale video driving dataset with detailed attributes using Amazon SageMaker Ground
264 Truth.
- 265 [10] Zheng, X., Chalasani, T., Ghosal, K., Lutz, S. & Smolic, A. (2019) STaDA: Style Transfer as Data
266 Augmentation. *CoRR*, *abs/1909.01056*.
- 267 [11] Frid-Adar, M., Klang, E., Amitai, M., Goldberger, J. & Greenspan, H. (2018) Synthetic data augmentation
268 using GAN for improved liver lesion classification. In *2018 IEEE 15th international symposium on biomedical
269 imaging (ISBI 2018)*, pp. 289–293.
- 270 [12] Pollok, T., Junglas, L., Ruf, B. & Schumann, A. (2019) UnrealGT: using unreal engine to generate ground
271 truth datasets. In *International Symposium on Visual Computing*, pp. 670–682. Springer, Cham.
- 272 [13] Cisse, M., Adi, Y., Neverova, N. & Keshet, J. (2017) Houdini: Fooling deep structured prediction models.
273 *arXiv preprint arXiv:1707.05373*.
- 274 [14] Hamesse, C., Lahouli, R., Fréville, T., Pairet, B. & Haelterman, R. (2019) Training Machine Learning
275 Algorithms for Computer Vision Tasks in Difficult Conditions: 3D Engines to the Rescue.
- 276 [15] Vega-Márquez, B., Rubio-Escudero, C., Riquelme, J. C. & Nepomuceno-Chamorro, I. (2019) Creation of
277 synthetic data with conditional generative adversarial networks. In *International Workshop on Soft Computing
278 Models in Industrial and Environmental Applications* pp. 231–240. Springer, Cham.
- 279 [16] Lim, S. K., Loo, Y., Tran, N. T., Cheung, N. M., Roig, G. & Elovici, Y. (2018) Doping: Generative data
280 augmentation for unsupervised anomaly detection with gan. In *2018 IEEE International Conference on Data
281 Mining (ICDM)* pp. 1122–1127. IEEE.
- 282 [17] Bowles, C., Chen, L., Guerrero, R., Bentley, P., Gunn, R., Hammers, A., ... & Rueckert, D. (2018) Gan aug-
283 mentation: Augmenting training data using generative adversarial networks. *arXiv preprint arXiv:1810.10863*.
- 284 [18] Gatys, L. A., Ecker, A. S. & Bethge, M. (2015) A neural algorithm of artistic style. *arXiv preprint
285 arXiv:1508.06576*.
- 286 [19] Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N. & Ganguli, S. (2015) Deep unsupervised learning using
287 nonequilibrium thermodynamics. In *International Conference on Machine Learning* pp. 2256–2265. PMLR.