

---

# HandsOff: Labeled Dataset Generation with No Additional Human Annotations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Because of their success in producing realistic images, generative adversarial net-  
2        works (GANs) have recently been leveraged to generate labeled synthetic datasets.  
3        However, existing dataset generation methods do not sufficiently leverage existing  
4        images with high quality labels, which often limits either the practicality of the  
5        system or the complexity of generated labels. We propose the HandsOff framework,  
6        which is capable of producing an unlimited number of synthetic images and corre-  
7        sponding labels after being trained on a small of number of *pre-existing labeled*  
8        *images*. Our framework avoids the practical drawbacks of similar frameworks  
9        while retaining the ability to generate rich pixel-wise labels, such as segmentation  
10        masks. This capability is achieved by unifying the field of GAN inversion with  
11        synthetic dataset generation, providing a new application for GAN inversion tech-  
12        niques. We demonstrate the efficacy of our framework on semantic segmentation  
13        tasks by generating labeled image datasets, and training and evaluating the perfor-  
14        mance of a downstream task. Our method achieves state-of-the-art performance in  
15        synthetic data trained semantic segmentation on both the CelebAMask-HQ dataset  
16        and Car-Parts-Segmentation dataset, and produces high quality segmentations  
17        in both domains. In addition, our framework uses significantly fewer computa-  
18        tional resources than prior work, demonstrating the supremacy of our approach in  
19        performance, annotation, and computation.

## 20    1 Introduction

21    The strong empirical performance of modern machine learning models has been enabled, in large part,  
22    by vast quantities of hand labeled data. Labeling massive datasets, such as ImageNet [1], requires a  
23    large time and cost investment. In contrast, collecting large quantities of *unlabeled* data is relatively  
24    easy. As a result, large quantities of unlabeled data exist alongside a small number of existing labeled  
25    images in many domains [1–3]

26    Recently, generative adversarial networks (GANs) [4–6], such as StyleGAN [7] and its variants  
27    [8–10], have demonstrated an ability to generate highly realistic images in numerous domains.  
28    Remarkably, the latent space of these networks form rich representations of images in a disentangled  
29    manner [11–13], which can be utilized to edit or remove complex semantic attributes in generated  
30    images. The ability to identify semantically meaningful parts of generated images in the latent  
31    space suggests that such representations could be used to generate pixel-level labels. This capability,  
32    coupled with GANs’ ability to generate vast troves of high quality images, could serve as the basis  
33    for generating synthetic image *datasets*.

34    In this work, we propose the HandsOff dataset generating framework, which is capable of producing  
35    synthetic images with corresponding labels. HandsOff leverages the expressive power of the GAN  
36    latent space, the existence of high quality labeled images, and recent advances in the field of GAN

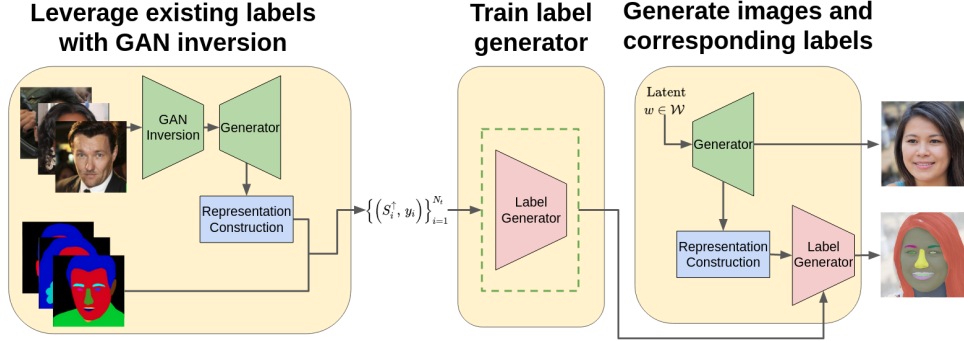


Figure 1: The HandsOff framework. GAN inversion is used to obtain training image latent codes  $w_i$ , which are then used to form hypercolumn representations  $S_i^\dagger$ . The label generator is then trained with the hypercolumn representations and original labels. To generate datasets, the trained label generator is used in conjunction with a StyleGAN2 generator to produce image-label pairs.

37 inversion. We empirically validate the power of our framework in two domains: faces and cars.  
 38 Furthermore, we explore directions for making the framework more computationally lightweight.

## 39 2 Related work

40 Our work is built on recent advances in GANs [4], which consist of a generator that synthesizes new  
 41 images, and a discriminator that discerns between real and generated images. Specifically, we utilize  
 42 the popular StyleGAN2 architecture [8], which synthesizes images by passing randomly sampled  
 43 inputs through a series of *style blocks*. StyleGAN2 is known for its numerous latent spaces, such as  
 44 the  $\mathcal{Z}$ ,  $\mathcal{W}$ ,  $\mathcal{W}+$ , and  $\mathcal{S}$  spaces. See [14] for a more detailed discussion.

45 GAN inversion is the process of mapping a real image onto the latent space of a GAN. The myriad of  
 46 inversion techniques range from encoder-based approaches [15–18], which utilize trained encoders  
 47 to map images directly to the latent space, to optimization-based approaches [19, 12, 13], which  
 48 directly optimize an image similarity based loss (e.g., LPIPS [20]) to obtain the latent code. Some  
 49 GAN inversion methods modify aspects of the generator, such as weights or noise injection values, to  
 50 increase image reconstruction quality. In our work, we exclusively use inversion methods that do not  
 51 modify the generator, since the generator must remain unperturbed to generate new images from the  
 52 original data distribution. We choose to invert images to the  $\mathcal{W}+$  space by learning a different latent  
 53 vector corresponding to each of the generator’s style blocks. As noted in [14], the  $\mathcal{W}+$  space is more  
 54 expressive and leads to higher quality reconstructed images. We primarily utilize encoder methods  
 55 to invert images, and use optimization methods to refine the encoder output in settings where finer  
 56 details are not preserved.

57 Numerous approaches utilize GANs to generate synthetic datasets [21–27], typically in the zero-shot  
 58 learning setting. We build upon DatasetGAN [28], which trains a label generator using representations  
 59 of an image formed from the GAN latent code. However, a considerable drawback of this framework  
 60 is that it requires *manual annotation of GAN generated images*. This is extremely burdensome, as  
 61 new annotations are required for every new domain in which a user wishes to synthesize datasets.  
 62 Furthermore, if the labeling paradigm changes, and the original labels cannot be directly mapped  
 63 to the new labels, then additional annotations are again required. Acquiring additional labels is  
 64 inconvenient, especially if labels already exist. EditGAN [29], a follow-up work to DatasetGAN,  
 65 hints that such a framework is possible. However, their focus is primarily on image editing, whereas  
 66 the HandsOff framework fully fleshes out the idea of unifying GAN inversion and dataset generation.

## 67 3 The HandsOff framework

68 The HandsOff framework, shown in Figure 1, builds upon the DatasetGAN framework by introducing  
 69 the ability to generate labeled synthetic data from existing labeled images. We use the term “label”  
 70 generically to emphasize that our approach applies broadly to label types such as segmentation  
 71 masks, keypoints, or any other pixel-level label. In our experiments, we focus on generating semantic  
 72 segmentation masks, a widely used and particularly resource intensive annotation to collect in bulk.

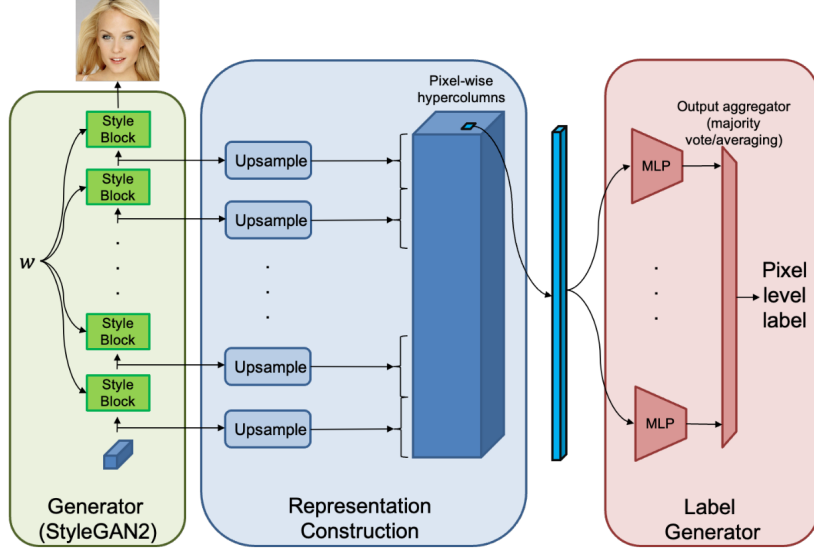


Figure 2: Construction of hypercolumns. Each intermediate style block output is upsampled to full image resolution and concatenated, resulting in a pixel-wise hypercolumn which can be fed into the label generator.

73 We utilize a frozen, pre-trained StyleGAN2 generator as the image generating backbone of the  
 74 framework. In order to generate labels, we pass the images’ latent codes through a label generator  
 75 that exploits the code’s unique semantic structure to efficiently generate high quality labels. Our  
 76 framework crucially departs from DatasetGAN’s during training. Rather than generating new latent  
 77 codes and manually annotating their corresponding images to train the label generator, we instead  
 78 use GAN inversion to obtain the latent codes corresponding to already labeled training images.  
 79 Specifically, assume we have  $N_t$  images  $X_1, \dots, X_{N_t}$  with corresponding labels  $y_1, \dots, y_{N_t}$ . After  
 80 passing these images through a GAN inverter, we obtain latent codes  $w_1, \dots, w_{N_t}$ , which are used to  
 81 form what we call the *hypercolumn representations*  $S_1^\dagger, \dots, S_{N_t}^\dagger$ . The label generator is then trained  
 82 on the  $\{(S_i^\dagger, y_i)\}_{i=1}^{N_t}$  pairs.

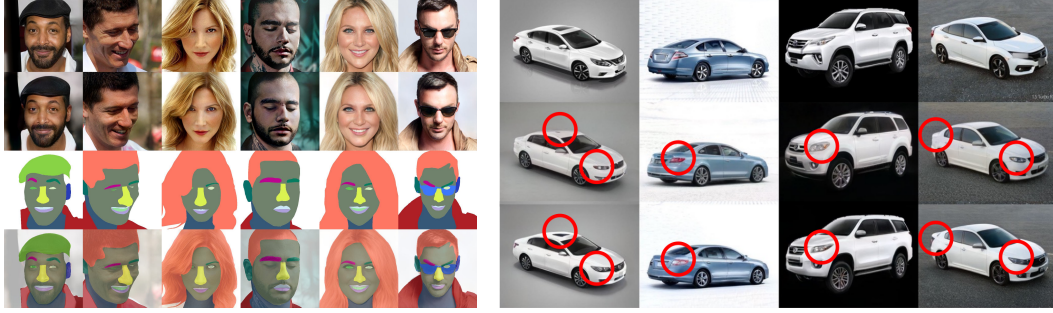
### 83 3.1 GAN inversion

84 The key step in the HandsOff framework is to employ GAN inversion in a new application area:  
 85 dataset generation. By utilizing GAN inversion in dataset generation, we *no longer require manual*  
 86 *annotation of GAN generated images*. Instead, a small number of pre-existing labels can be used to  
 87 generate massive synthetic datasets. By using pre-existing labels, practitioners not only avoid the  
 88 cost of acquiring labels, but also avoid the prerequisite of maintaining annotation workstreams in  
 89 their machine learning pipelines.

90 One requirement of the HandsOff framework places on GAN inversion techniques is preserving the  
 91 generators’ original weights. This requirement ensures that the generator produces *new* images from  
 92 the data distribution of the domain of interest. In our experiments, we utilize ReStyle [16], but we  
 93 emphasize that our framework is amenable to *any* GAN inverter that does not modify the generator  
 94 weights.

95 In settings where finer details are not preserved with ReStyle, we further refine the latent code  
 96 obtained from ReStyle by solving a regularized optimization problem, similar to approaches in  
 97 [29, 30]. Our use of optimization-based approaches is justified despite their slower inference times  
 98 because GAN inversion is only performed once at the beginning on a small number of training images.  
 99 In particular, let  $G$  be a frozen, pre-trained StyleGAN2 generator,  $X$  be the image to be inverted, and  
 100  $w^{(r)}$  be the latent code obtained by running ReStyle on  $X$ . We refine  $w^{(r)}$  by solving the following  
 101 optimization problem:

$$\min_w \mathcal{L}_{LPIPS}(X, G(w)) + \lambda_{\ell_2} \cdot \|X - G(w)\|_2^2 + \lambda_{reg} \cdot \|w - w^{(r)}\|_2^2. \quad (1)$$



(a) Visualization of inverting real face images (row 1) with ReStyle. The reconstructed images (row 2) align well semantically with the original segmentation masks (row 3), as shown in row 4. (b) Visualization of reconstruction quality before (row 2) and after (row 3) optimization refinement. Red circles indicate finer details that were improved to align better with the original image (row 1).

Figure 3: Visualization of image reconstruction quality for faces and car domains.

102 Above,  $\mathcal{L}_{LPIPS}$  is the LPIPS loss [20] and  $\|\cdot\|_2$  is the  $\ell_2$  norm. In practice, (1) is a highly non-convex  
 103 problem and using a set number of gradient descent iterations significantly refines the latent code.  
 104 This refinement approach generalizes beyond ReStyle – any encoder output can be refined by (1).

### 105 3.2 Hypercolumn representation

106 Within the StyleGAN2 generator, the latent code  $w$  is used to modulate convolution weights in  
 107 intermediate style blocks, which progressively grow an input to the final output image. For a  $1024$   
 108  $\times 1024$  resolution image, there are  $L = 18$  style blocks. We take the intermediate output of these  
 109 style blocks, upsample them channel-wise to the resolution of the full image, then concatenate  
 110 each upsampled intermediate output channel-wise to obtain pixel-wise hypercolumns, similar to the  
 111 approach in [28]. We denote the hypercolumn representation of the image as  $S^\dagger$ , with each pixel  $j$   
 112 now having a hypercolumn  $S^\dagger[j]$  of dimension  $C$ . This process is shown in Figure 2.

113 In practice, we cap the generated image resolution to  $512 \times 512$ , and downsample intermediate  
 114 outputs from the  $1024 \times 1024$  layers, as well as the original image. This is done for memory  
 115 considerations when storing and training with these hypercolumn representations, as the dimension  
 116 of each of the hypercolumns is relatively high ( $C = 6080$  for  $1024 \times 1024$  images). In Section 4.3,  
 117 we explore using only a subset of the channels from the intermediate output in order to form the  
 118 hypercolumn representation to alleviate these memory burdens.

### 119 3.3 Label generator

120 The label generator exploits the semantically rich latent space of the generator to efficiently produce  
 121 high quality labels for generated images. Because the latent codes already map to semantically  
 122 meaningful parts of generated images, complex vision models are not necessary to generate labels.  
 123 Specifically, we utilize an ensemble of  $M$  MLPs. The MLPs operate on a pixel-level, mapping a  
 124 pixel’s hypercolumn to a label. To generate a label for a synthetic image, we pass the hypercolumn  
 125 formed by latent code  $w$  through the  $M$  MLPs, and aggregate the outputs to produce a label.  
 126 Specifically, in the semantic segmentation setting with  $K$  parts, the MLP performs pixel-wise  
 127 classification, mapping pixel  $j$ ’s hypercolumn  $S^\dagger[j]$  to a label  $k \in \{1, \dots, K\}$ .

128 The  $M$  MLPs are trained using a small number ( $\sim 50$ ) of pre-existing labeled images with a cross-  
 129 entropy loss. In Section 4.3, we explore modifications to the label generating architecture that result  
 130 in significant performance gains. We further experiment with reducing the number of MLPs in the  
 131 ensemble in Appendix C.

### 132 3.4 Downstream task

133 In order to benchmark the quality of our generated datasets, we quantify the downstream performance  
 134 of models trained exclusively on our datasets. After training a network on a generated dataset, we  
 135 evaluate its performance on a hold out test set of real images with human annotations. We refer to this  
 136 process as the *downstream task* and the trained model as the *downstream model*. Prior to training the



Table 1: Experimental results in face and car domains for semantic segmentation, reported in mIOU. Moving from 16 to 50 training images in HandsOff avoids the need to manually annotate 34 GAN generated images. We are unable to compare against DatasetGAN in the car domain because the labeling paradigm used to train DatasetGAN cannot be converted to the labeling paradigm in the Car-Part-Segmentation dataset.

	DatasetGAN 16 train	EditGAN 16 train	EditGAN (Encoder only) 16 train	HandsOff (Ours)
Faces	0.7013	0.7244	×	<b>0.7696</b> (16 train) <b>0.7748</b> (50 train)
Cars	N/A	0.6023	0.5368	<b>0.6222</b> (16 train) <b>0.6591</b> (50 train)

137 downstream network, we follow the approach of [28, 31], and filter out the top 10% most uncertain  
 138 images according to Shannon-Jenson divergence [32, 33].

## 139 4 Experimental results

140 In this section, we present qualitative results in GAN inversion, and downstream task performance in  
 141 the car and face domains. For faces, we utilize segmentation masks from CelebAMask-HQ [3], and  
 142 collapse the original 19 classes into 8 classes. For cars, we use the Car-Parts-Segmentation dataset  
 143 [34], and collapse the original 19 classes into 12 classes. Details about class collapse are described in  
 144 Appendix A.

### 145 4.1 GAN inversion images and original segmentation mask alignment

146 The underlying assumption of the HandsOff framework is that the semantic features in the recon-  
 147 structed images align well with the segmentation masks. Should the reconstructed semantic features  
 148 not align well, we would essentially be training the label generator with corrupted representations.  
 149 Therefore, a crucial first step is verifying the fidelity of image reconstructions.

150 In the face domain, we utilize ReStyle for GAN inversion. As seen in Figure 3a, the reconstructed  
 151 images from the latent codes obtained by ReStyle align very well with the semantic segmentation  
 152 masks from CelebAMask-HQ. In the car domain, we highlight the power of our latent code refinement  
 153 scheme, presented in Equation 1. ReStyle is unable to exactly preserve smaller details, such as the  
 154 shape of headlights or the presence of a sunroof, as shown in the first row of Figure 3b. As a result,  
 155 we refine the ReStyle output with 500 iterations of (1). As seen in the second row of Figure 3b, the  
 156 optimization program is able to fine-tune the mis-aligned details, highlighted in the red circles.

### 157 4.2 Downstream task performance

158 We now discuss the performance of the HandsOff framework on the downstream task. In particular,  
 159 we generate 10000 synthetic images and segmentation masks, filter out the top 10% most uncertain  
 160 images, and train DeepLabV3 for 20 epochs with the 9000 remaining images. Examples of generated  
 161 images and segmentation masks for the face domain can be found in Figure 4. For HandsOff, we  
 162 utilize  $M = 10$  MLPs in the label generating branch and report the best performing result among the  
 163 combinations of network layer widths experimented with later on in Section 4.3.

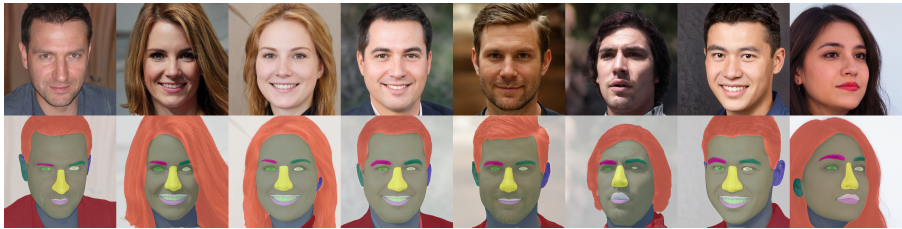
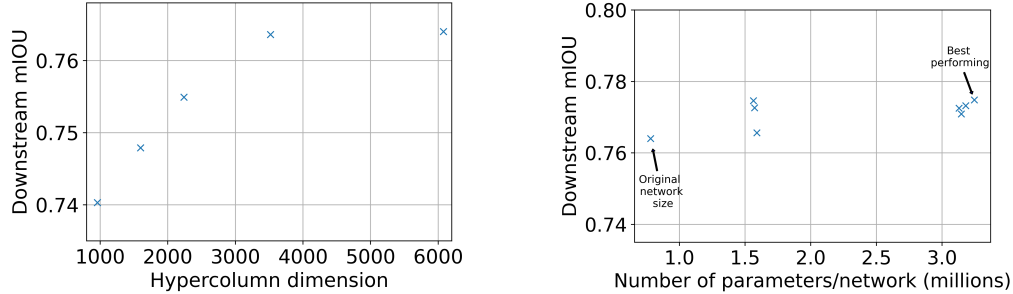


Figure 4: Examples of faces and their corresponding segmentation masks generated from the HandsOff framework trained on 50 images with *non-collapsed classes*. We do not perform class collapse here to highlight the level of detail that the label generator is able to achieve.



(a) Comparable performance is achieved with a 42% reduction in hypercolumn dimension. This reduction decreases the amount of memory used in training.

(b) Increasing layer widths results in relatively sizable performance gains, with the best performance occurring with intermediate layer widths of 512 and 256.

Figure 5: Framework ablations for hypercolumn dimension (left) and layer widths (right).

164 For faces, we split the 30000 images in CelebAMask-HQ into 3 sets: 50 images for training the label  
 165 generator, 450 images for validation, and 29500 images for testing. For cars, we retain the original  
 166 train (400 images) and test (100 images) splits from the Car-Parts-Segmentation dataset, using 20  
 167 images from the test set for validation. For both domains, we report mIOU of the downstream network  
 168 when trained on images generated from the HandsOff, DatasetGAN, and EditGAN frameworks.

169 As seen in Table 1, we achieve state-of-the-art performance in synthetic data trained semantic  
 170 segmentation in both the face and car domains. In particular, for faces, we highlight that increasing  
 171 the number of training examples from 16 to 50 results in a sizable performance increase. Within the  
 172 DatasetGAN framework, this would require manual annotation of 34 more images, a step that is not  
 173 required in the HandsOff framework. Note that we are unable to test the performance of EditGAN  
 174 (Encoder only) because the pre-trained EditGAN encoder weights were not publicly released.

175 In the car domain, the optimization based refinement of the ReStyle output results in a sizable  
 176 performance increase, highlighting the importance of strong alignment of reconstructed images with  
 177 the original segmentation masks. Furthermore, we are unable to run the DatasetGAN framework  
 178 with the Car-Part-Segmentation labels, because the original labeling paradigm used in training the  
 179 DatasetGAN framework in the car domain cannot be converted to the Car-Part-Segmentation labeling  
 180 paradigm. This highlights a key drawback of DatasetGAN, as discussed in Section 2.

### 181 4.3 Practical modifications to HandsOff

182 We first experiment with keeping only a subset of the channels from the style block intermediate  
 183 outputs from the lower resolution layers. In the StyleGAN2 generator, the first 10 style block outputs  
 184 (which range from  $4 \times 4$  to  $128 \times 128$  resolutions) each contain 512 channels, comprising 5120 of the  
 185 6080 total channels. We quantify the effect of keeping zero or the first 64, 128, and 256 channels on  
 186 the downstream task performance in the face domain. As shown in Figure 5a, while utilizing only  
 187 higher resolution layers degrades performance considerably, we can remove 256 of the 512 channels  
 188 for the first 10 style blocks with very minimal loss in performance. This results in a hypercolumn  
 189 dimension 3520, which is a 42% reduction compared to the original dimension of 6080.

190 Finally, we investigate whether network layer widths impact downstream performance. The original  
 191 DatasetGAN framework utilizes 3-layer MLPs with intermediate dimensions of 128 and 32. We  
 192 explore 7 additional combinations of layer widths, as highlighted in Figure 5b. Generally, downstream  
 193 performance increases with increasing network widths.

## 194 5 Discussion

195 We present the HandsOff framework, which is capable of producing high quality labeled synthetic  
 196 datasets without requiring further annotation of images. We achieve state-of-the-art performance on  
 197 downstream tasks in this setting, and experiment with ways of making the framework more practical  
 198 from an implementation standpoint. While we focus primarily on generating semantic segmentation  
 199 datasets, nothing in this framework precludes generation of continuous pixel-wise labels, such as  
 200 depth maps which we plan to explore in future work.

201 **References**

- 202 [1] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
203 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual  
204 recognition challenge. *International journal of computer vision*, 2015.
- 205 [2] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul  
206 Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for  
207 autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on  
208 computer vision and pattern recognition*, 2020.
- 209 [3] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and  
210 interactive facial image manipulation. In *IEEE Conference on Computer Vision and Pattern  
211 Recognition (CVPR)*, 2020.
- 212 [4] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
213 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural  
214 information processing systems*, 2014.
- 215 [5] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity  
216 natural image synthesis. In *International Conference on Learning Representations*, 2018.
- 217 [6] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for im-  
218 proved quality, stability, and variation. In *International Conference on Learning Representations*,  
219 2018.
- 220 [7] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
221 adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and  
222 pattern recognition*, 2019.
- 223 [8] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila.  
224 Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF  
225 conference on computer vision and pattern recognition*, 2020.
- 226 [9] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila.  
227 Training generative adversarial networks with limited data. *Advances in Neural Information  
228 Processing Systems*, 2020.
- 229 [10] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen,  
230 and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information  
231 Processing Systems*, 2021.
- 232 [11] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for  
233 semantic face editing. In *Proceedings of the IEEE/CVF conference on computer vision and  
234 pattern recognition*, 2020.
- 235 [12] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the  
236 stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer  
237 Vision*, 2019.
- 238 [13] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embed-  
239 ded images? In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
240 recognition*, 2020.
- 241 [14] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. Gan  
242 inversion: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 243 [15] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an  
244 encoder for stylegan image manipulation. *ACM Transactions on Graphics (TOG)*, 2021.
- 245 [16] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder  
246 via iterative refinement. In *Proceedings of the IEEE/CVF International Conference on Computer  
247 Vision*, 2021.

- 248 [17] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and  
249 Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. In  
250 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021.
- 251 [18] Tengfei Wang, Yong Zhang, Yanbo Fan, Jue Wang, and Qifeng Chen. High-fidelity gan inversion  
252 for image attribute editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
253 and Pattern Recognition*, 2022.
- 254 [19] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the  
255 local semantics of gans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
256 Pattern Recognition*, 2020.
- 257 [20] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unrea-  
258 sonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE  
259 conference on computer vision and pattern recognition*, 2018.
- 260 [21] Yang Long, Li Liu, Fumin Shen, Ling Shao, and Xuelong Li. Zero-shot learning using  
261 synthesised unseen visual data with diffusion regularisation. *IEEE transactions on pattern  
262 analysis and machine intelligence*, 2017.
- 263 [22] Maxime Bucher, Stéphane Herbin, and Frédéric Jurie. Generating visual representations for  
264 zero-shot classification. In *Proceedings of the IEEE International Conference on Computer  
265 Vision Workshops*, 2017.
- 266 [23] Rafael Felix, Ian Reid, Gustavo Carneiro, et al. Multi-modal cycle-consistent generalized  
267 zero-shot learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*,  
268 2018.
- 269 [24] Mert Bulent Sariyildiz and Ramazan Gokberk Cinbis. Gradient matching generative networks  
270 for zero-shot learning. In *Proceedings of the IEEE/CVF conference on computer vision and  
271 pattern recognition*, 2019.
- 272 [25] Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba.  
273 Bigdatasetgan: Synthesizing imagenet with pixel-wise annotations. In *Proceedings of the  
274 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- 275 [26] Daiqing Li, Junlin Yang, Karsten Kreis, Antonio Torralba, and Sanja Fidler. Semantic segmenta-  
276 tion with generative models: Semi-supervised learning and strong out-of-domain generalization.  
277 In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- 278 [27] Rameen Abdal, Peihao Zhu, Niloy Mitra, and Peter Wonka. Labels4free: Unsupervised  
279 segmentation using stylegan, 2021.
- 280 [28] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso,  
281 Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal  
282 human effort. In *CVPR*, 2021.
- 283 [29] Huan Ling, Karsten Kreis, Daiqing Li, Seung Wook Kim, Antonio Torralba, and Sanja Fidler.  
284 Editgan: High-precision semantic image editing. In *Advances in Neural Information Processing  
285 Systems (NeurIPS)*, 2021.
- 286 [30] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A Efros. Generative visual  
287 manipulation on the natural image manifold. In *European conference on computer vision*, 2016.
- 288 [31] Weicheng Kuo, Christian Häne, Esther Yuh, Pratik Mukherjee, and Jitendra Malik. Cost-  
289 sensitive active learning for intracranial hemorrhage detection. In *International Conference on  
290 Medical Image Computing and Computer-Assisted Intervention*, 2018.
- 291 [32] Prem Melville, Stewart M Yang, Maytal Saar-Tsechansky, and Raymond Mooney. Active  
292 learning for probability estimation using jensen-shannon divergence. In *European conference  
293 on machine learning*, 2005.
- 294 [33] William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of  
295 ensembles for active learning in image classification. In *Proceedings of the IEEE conference on  
296 computer vision and pattern recognition*, 2018.

- 297 [34] Kitsuchart Pasupa, Phongsathorn Kittiworapanya, Napasin Hongngern, and Kuntpong Worarat-  
298 panya. Evaluation of deep learning algorithms for semantic segmentation of car parts. *Complex*  
299 *& Intelligent Systems*, 2021.



Figure 6: Examples of generated images with poor labels. These often contain components not seen in the training data (hats, multiple humans, children), suggesting that the distribution of parts in the training data has a noticeable impact on the generated labels.

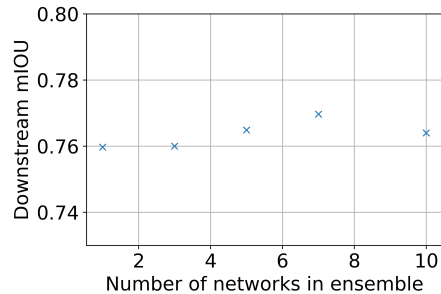


Figure 7: Downstream performance seems fairly robust to the number of MLPs in the label generator. Using fewer MLPs results in a decrease in time needed to train the framework.

## 300 A Dataset details

301 In our experiments, we collapse the original labels in each dataset in a smaller number of labeled  
 302 parts. For CelebAMask-HQ dataset, we remove any distinction between left/right in a number of  
 303 parts (e.g., ears, eyes, eyebrows). Furthermore, we form one mouth part consisting of upper/lower  
 304 lips and mouth. Finally, we collapse all accessories and clothing into background.

305 For the Car-Parts-Segmentation dataset, we remove any distinction between left/right and front/back  
 306 for parts such as doors, lights, bumpers, and mirrors. We also merge trunks and tailgates to be the  
 307 same class.

## 308 B Examples of failure cases of generated images and labels

309 In this section, we highlight failure cases of generated labeled images. These images typically include  
 310 components not included in the training data, such as humans with hats/hoods, multiple humans, and  
 311 babies. Exploration of how including more edge cases in the training data affects generated label  
 312 quality is an interesting direction of future work.

## 313 C Number of MLPs ablation

314 Because training an ensemble of classifiers is timely, we experiment with utilizing fewer MLPs. We  
 315 train the 10 MLPs, then from the 10 trained MLPs, we use 1, 3, 5, 7, and 10 MLPs to generate labels.  
 316 As seen in Figure 7, using only 1 network results in a performance drop, but using anywhere from 3  
 317 to 7 MLPs results in performance meeting or even exceeding the performance of using all 10 MLPs.  
 318 This implies that there is diminishing returns in using more MLPs, in that simply using any  $M > 1$   
 319 achieves (or even surpasses) the robustness benefits of using 10 MLPs.