
Attribute Based Interpretable Evaluation Metrics for Generative Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 While generative models continue to evolve, the field of evaluation metrics has
2 largely remained stagnant. Despite the annual publication of metric papers, the
3 majority of these metrics share a common characteristic: they measure distributional
4 distance using pre-trained embeddings without considering the interpretability of
5 the underlying information. This limits their usefulness and makes it difficult to gain
6 a comprehensive understanding of the data. To address this issue, we propose using
7 a new type of interpretable embedding. We demonstrate how we can transform
8 deeply encoded embeddings into interpretable embeddings by measuring their
9 correspondence with text attributes. With this new type of embedding, we introduce
10 two novel metrics that measure and explain the diversity of the generator: the first
11 metric compares the frequency of appearance of the training set and the attribute,
12 and the second metric evaluates whether the relationships between attributes in the
13 training set are preserved. By introducing these new metrics, we hope to enhance
14 the interpretability and usefulness of evaluation metrics in the field of generative
15 models.

16 1 Introduction

17 Significant advancements have been achieved in the image generation field, from the pioneering
18 introduction of generative adversarial networks (GANs) to the more recent emergence of diffusion
19 models (DMs). [5, 10, 27] In recent years, generated images are hardly distinguishable from real
20 images. In this context, evaluating the generated images for a given training dataset has played a
21 critical role in the development.

22 Envision an evaluation scenario where the outputs of two generative models are compared against
23 a common training dataset. What would be the underlying factors for judging a set as superior to
24 another set? As the goal of generative models is mimicking the real data distribution, various metrics
25 have been designed to assess the similarity between the generated images and the training dataset, e.g.,
26 Fréchet Inception Distance (FID)[9], Precision and Recall[25][17], and Density and Coverage[22].

27 Most of these evaluation metrics capture the disparity between the training data distribution and the
28 distribution of generated images by examining the differences in feature representations within the
29 embedding space of a pre-trained network[26, 28]. FID is a widely used metric that quantifies the
30 dissimilarity in visual features to assess the quality and diversity of the generated images. Specifically,
31 it measures the distance between the real and fake distributions in the embedding space of Inception-
32 V3[28].

33 An important question arises regarding the suitability of the embedding space employed for evaluating
34 generated images. The embedding space of the pre-trained model may vary depending on the dataset
35 and task it was trained on. For instance, Inception V3 was trained for image classification on

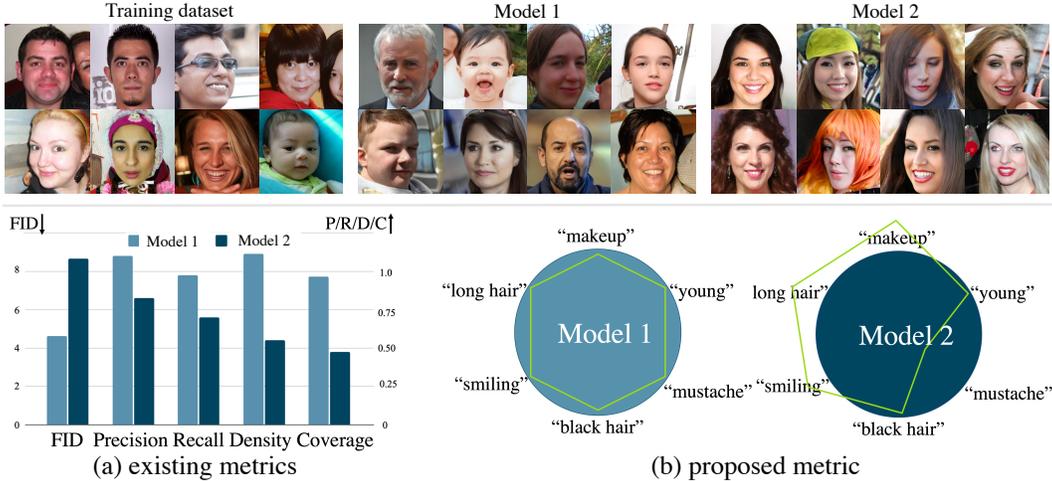


Figure 1: **Conceptual illustration of our method.** We design the scenario, Model 2 lacks diversity. (a) Although existing metrics distinguish the inferiority of Model 2, they provide no explanation about judgment. (b) Our attribute-based proposed metric has interpretation; Model 2 is biased with ‘long hair’ and ‘makeup’.

36 ImageNet[3], suggesting that its embedding space is designed to compress image information and
 37 discern essential patterns for classification. Consequently, the appropriateness of employing this
 38 embedding space for evaluating generated images remains an open question.

39 Returning to the fundamental question at hand, Figure 1 makes an evaluation scenario a little bit more
 40 specific. Suppose there are realistic generated images from two distinct models. As shown in the
 41 example images, it is evident that Model 2 generates biased images, i.e., there are only women, while
 42 Model 1 successfully generates various images that are close to training data. Fortunately, although
 43 there remains an open question about embedding space, the values of various metrics in Figure 1 (b)
 44 align reasonably well with our interpretation; Model 1 is perceived as superior.

45 However, what are the underlying factors that contribute to such judgment? Although the results
 46 are consistent with a person’s conclusion, it far fails to provide a comprehensive explanation. The
 47 interpretation of distances within the embedding space from a pre-trained classification model remains
 48 elusive, posing challenges in evaluation. On the contrary, humans readily discern certain factors for
 49 judgment; individuals easily recognize the bias of Model 2. These factors suggest more information
 50 and a direction beyond simple ranking. In this paper, we propose an evaluation metric that aims to
 51 interpret the underlying factors behind such judgments.

52 To address this objective, we begin by examining attribute comparison methods in human judgment.
 53 When evaluating two generated image distributions, humans compare the attributes present in the
 54 training dataset with those exhibited by the generated images. Key attributes under consideration
 55 include gender, facial representation, and age distribution. Ideally, with well-defined training data, we
 56 anticipate the attributes in the generated images to align with those in the training data. If the model
 57 lacks essential attributes (e.g., gender, age, glasses, or hats), it is insufficient to generate visually
 58 realistic images. Incorporating these attributes into the evaluation process may enable a more explicit
 59 and comprehensive assessment.

60 This paper presents a novel approach for evaluating generative models by leveraging a newly proposed
 61 embedding space that incorporates attribute-specific information. Similar to human visual judgment,
 62 our metrics evaluate images in terms of various characteristic attributes. Figure 1 (b) illustrates the
 63 concept of our metric; it captures the distribution differences of attributes. We use pre-trained CLIP
 64 [24], a language-image model trained on a huge dataset, to define a new embedding space that can
 65 quantify images for multiple attributes.

66 To facilitate our embedding space, we introduce the "Directional CLIPScore" (DCS), a method for
 67 quantifying each attribute based on the training data. Within our proposed embedding space, each
 68 channel comprises DCS values that explicitly indicate the relevance of an image to specific attributes.
 69 The use of a perceptible embedding space offers the advantage of interpretability.

70 We introduce two novel evaluation metrics to use the newly proposed embedding space. Firstly,
71 the "Single attribute KL Divergence (**SaKLD**)" compares attribute distributions between training
72 data and the generated images, providing a quantitative measure of the similarity between attribute
73 distributions. It quantifies how closely the attributes of generated images align with the attribute
74 distribution in training data. Secondly, we introduce the "Paired attribute KL divergence (**PaKLD**)"
75 that considers correlations among multiple attributes. This metric accounts for the relationship
76 between attributes, such as the presence of a beard in an image of a woman. PaKLD successfully
77 evaluates the generated images while taking into consideration attribute relationships.

78 We validate our metrics through a series of carefully designed experiments, demonstrating their
79 effectiveness and interpretability. By employing our metric, we conduct a comprehensive analysis of
80 prominent generation models currently considered state-of-the-art [11, 13, 12, 14, 23]. Interestingly,
81 our findings reveal variations in performance across different datasets. For instance, diffusion models
82 exhibit superior performance on datasets with a large number of samples, such as FFHQ. In contrast,
83 GANs outperform diffusion models on datasets with relatively smaller sample size, such as MetFaces.

84 In summary, this paper presents a novel approach for evaluating generative models using a new
85 embedding space that incorporates attribute-specific information. Our proposed method, along with
86 the introduced evaluation metrics, allows for a comprehensive assessment of generated images by
87 considering attribute distributions and correlations. Our findings contribute to the research field by
88 advancing the understanding and evaluation of generative models, offering insights into their strengths
89 and limitations. Moreover, our work opens avenues for future research and potential improvements in
90 the field of generative image synthesis by comprehensive evaluation metrics.

91 **2 Related Work**

92 **Fréchet Inception Distance** Fréchet Inception Distance (FID) [9] measures the distance between
93 the estimated Gaussian distributions of two datasets by passing them through a pre-trained Inception-
94 v3[28] model. However, Kynkäänniemi et al. [18] revealed that when generated images are far from
95 training data, the embeddings may incorrectly highlight irrelevant parts of images. To address this
96 issue, the researchers proposed using the CLIP [24] image encoder instead of Inception-v3 to calculate
97 the 2-Wasserstein distance, which provides reliable results regardless of the dataset being measured.

98 **Fidelity and diversity** Sajjadi et al. [25] introduced precision and recall for evaluating generative
99 model, and subsequent studies by Kynkäänniemi et al. [17] and Naeem et al. [22] have further refined
100 this approach. Most of these methods use a pre-trained network to examine whether the embedding
101 of generated images falls within the boundary of real image embedding (precision) and whether
102 the embedding of real images falls within the boundary of generated image embedding (recall) for
103 assessing fidelity and diversity.

104 **Rarity score** Han et al. [6] proposed a metric for measuring the rarity of generated images. They
105 quantified how rare the generated images are within a k-NN sphere to assess their rarity. The key
106 difference between the rarity score and diversity in precision and recall is that the rarity score
107 considers only the generated samples that fall within the manifold of real samples. In other words, it
108 focuses on how well the generated images fit within the distribution of real images in terms of rarity,
109 rather than capturing the overall diversity of generated samples.

110 However, we note that the concept of using raw embeddings from a pre-trained classifier remains
111 consistent among all these metrics.

112 **A call for explainable evaluation** Existing evaluation metrics in the field of generative models lack
113 the ability to provide detailed insights into the diversity of generated images. As shown in Figure 1,
114 even though metrics like FID, Precision and Recall indicate poor performance for a biased generator
115 towards specific attributes (e.g., "makeup" and "long hair"), they do not provide an explanation
116 for judgment factors. Therefore, researchers manually identified the underlying factors by visual
117 inspection but it becomes increasingly challenging with larger sample sizes. To address this issue,
118 we propose novel explainable evaluation metrics that provide in-depth analysis and insights into the
119 diverse generation abilities of models.

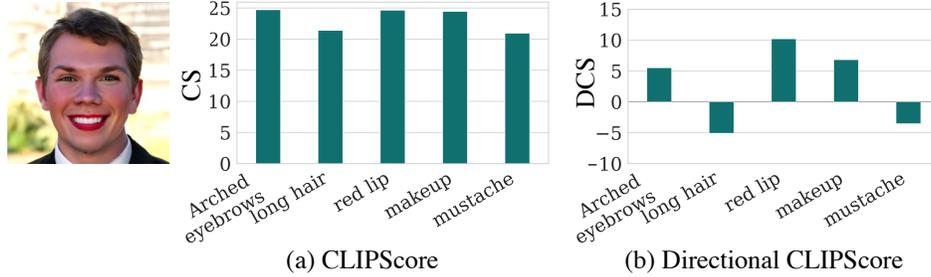


Figure 2: **Difference between CS and DCS.** (a) CLIPScore[8] exhibits similar values, making it difficult to discern. (b) Directional CLIPScore has an intuitive value based on zero. We design a new embedding space; each channel represents the intensity of a specific attribute by DCS, informing explanations about the single image.

120 3 Attribute-Driven Embedding

121 Existing metrics for evaluating generated images commonly utilize embeddings before FCN, from
 122 Inception-V3 or CLIP image encoder[7, 4]. However, these approaches lack interpretability as the
 123 meaning of each channel in the embedding. Additionally, Kynkäänniemi et al. [18] have shown the
 124 FID scores improve significantly when the classification distribution matches that of the training
 125 set, irrespective of the quality, highlighting another limitation of the existing embedding. To address
 126 these issues and develop an explainable evaluation metric, we design each embedding of images to
 127 possess an 'interpretation'. Section 3.1 presents the process of generating explainable embedding for
 128 individual images using the CLIP encoder, and Section 3.2 introduces the Directional CLIPScore, a
 129 novel embedding approach that enhances interpretability and accuracy.

130 3.1 Attribute-driven embeddings for better representations

131 To achieve an interpretable embedding, we utilized each channel of the embedding as a measure of
 132 the attribute's prominence in the image. A straightforward approach to quantify attribute strength is
 133 by employing CLIPScore;

$$\text{CLIPScore}(x, a) = 100 * \text{sim}(\mathbf{E}_I(x), \mathbf{E}_T(a)), \quad (1)$$

134 where x is a single image, a is a given text of attribute, $\text{sim}(*, *)$ is cosine similarity, and \mathbf{E}_I and \mathbf{E}_T
 135 are CLIP image encoder and text encoder respectively. We selected multiple attributes that effectively
 136 represent image characteristics as textual descriptions and measured CLIPScore with individual
 137 images and selected attributes. The way to select attributes will refer to Section 3.3. By assigning
 138 these CLIPScores as the values for each channel in the embedding, we obtained an interpretable
 139 representation. However, relying solely on CLIPScore has challenges as the cosine similarity values
 140 tend to be similar, making it difficult to discern the relative differences between attribute strengths.
 141 Intuitively, selected human-related attributes tend to cluster closely in the CLIP embedding, resulting
 142 in smaller variations in cosine similarity. To address this limitation, subsequent subsections introduce
 143 the Directional CLIPScore, which offers a more precise scoring approach.

144 3.2 Directional CLIPScore

145 As discussed, CLIPScore exhibits a narrow distribution of values, which can be attributed to measuring
 146 similarity between human-related attributes, resulting in their dense clustering on the CLIP embedding.
 147 Figure 3 (a) visualizes it. To address this issue, we propose Directional CLIPScore (DCS), which
 148 leverages the centers of training images and predefined attribute texts on the CLIP embedding.

149 Given training data, denoted as $\{x_1, x_2, x_3, \dots\} \in \mathcal{X}$, we define $C_{\mathcal{X}}$ as the center of images and $C_{\mathcal{T}}$
 150 as another center of images for text attributes on the CLIP embedding, respectively. By using the
 151 image captioning model, BLIP[19], we define $C_{\mathcal{T}}$ as the center of images in text respect;

$$C_{\mathcal{X}} = \frac{1}{N} \sum_{i=1}^N \mathbf{E}_I(x_i), \quad C_{\mathcal{T}} = \frac{1}{N} \sum_{i=1}^N \mathbf{E}_T(\text{BLIP}(x_i)). \quad (2)$$

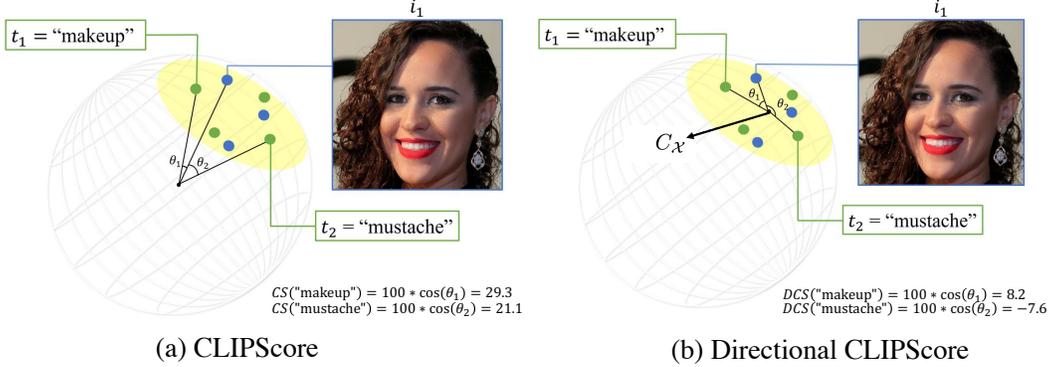


Figure 3: **Illustration of CLIPScore and Directional CLIPScore.** (a) CLIPScore measures the similarity between vectors with coordinate origin. (b) Directional CLIPScore measures the similarity between vectors with a defined mean of the images, $C_{\mathcal{X}}$, as the origin. In the figure, we illustrate $C_{\mathcal{X}}$ and $C_{\mathcal{T}}$ as the same point for ease of clarity and comprehension.

Table 1: **CLIPScore and Directional CLIPScore’s mean accuracy on CelebA dataset.**

	All attributes		Refined attributes	
	CLIPScore	Directional CLIPScore	CLIPScore	Directional CLIPScore
mean accuracy	0.395	0.409	0.501	0.530

152 These centers serve as reference points in the embedding space and aid more accurate attribute
 153 scores. We define DCS as the measure of similarity between two directions, V_x and V_a where a set
 154 of attributes defined as $\{a_1, a_2, a_3, \dots\} \in \mathcal{A}$. The first direction spans from the center of the image
 155 to the image itself, and the second direction extends from the center of the attributes to the desired
 156 attribute.

$$V_x = \mathbf{E}_{\mathbf{I}}(x) - C_{\mathcal{X}}, \quad V_a = \mathbf{E}_{\mathbf{T}}(a) - C_{\mathcal{T}}, \quad (3)$$

157

$$\text{DCS}(x, a) = 100 * \text{sim}(V_x, V_a), \quad (4)$$

158 where $\text{sim}(*, *)$ is cosine similarity. For extending DCS from a single sample to data we denote the
 159 probability density function (PDF) of $\text{DCS}(x_i, a_i)$ for all $x_i \in \mathcal{X}$ as $\text{DCS}_{\mathcal{X}}(a_i)$ for brevity.

160 Figure 3 visually illustrates the distinction between DCS (Directional CLIPScore) and CS (CLIP-
 161 Score). Unlike CS, which lacks a clear reference point, DCS is based on the center, enabling the
 162 determination of attribute magnitudes relative to a zero point. Furthermore, DCS exhibits superior
 163 accuracy compared to CS, as demonstrated in Table 1. The table presents the accuracy results of CS
 164 and DCS for annotated attributes in CelebA[20]. By evaluating how well positive samples with the
 165 highest score align with positive samples for a given attribute, DCS consistently outperforms CS
 166 in accuracy. Notably, this trend remains consistent across refined attributes, which are removed for
 167 subjective attributes such as "Attractive" or "Blurry".

168 3.3 attribute selection methodologies

169 Our evaluation metric for measuring the performance of the generator is dependent on the attributes we
 170 choose to measure. To explore how to choose attributes that accurately reflect generator performance,
 171 we introduce three methods for attribute selection.

172 **BLIP extracted attribute** We aim to identify and quantify the attributes present in the training
 173 data from image descriptions. We can determine which attributes are most commonly occurring in
 174 the training data by counting attributes that appear in the training data. We use the image captioning
 175 model, BLIP[19], to extract attribute-related words from training data. We use N attributes that
 176 appear frequently in the training data as a set of attributes \mathcal{A} for our proposed metric.

177 **User annotation** Another option for attribute selection is to use a set of human-annotated attributes.
 178 By explicitly assigning attributes for evaluating generative models, users can fairly compare the

179 impact of each attribute score or focus on specific attributes. Especially, the CelebA dataset provides
 180 40 binary attributes about the human face domain, which can be used to evaluate a wide range of
 181 (generated) human image sets.

182 **GPT attributes** We leveraged the power of GPT-3[1] to extract attributes. Through repetitive
 183 questioning, such as ‘Give me 50 words of useful visual attributes for distinguishing faces in a
 184 photo’ and ‘Give me 50 words of useful visual attributes for discerning variations in facial features to
 185 identify people in images,’ we obtained a set of attributes, which frequently appeared in the responses
 186 across different datasets. The list of questions posed to GPT-3 can be found in the Appendix, and we
 187 followed the questioning methodology outlined in [21].

188 4 Evaluation Metric with Interpretable Attribute-Driven Embedding

189 In this section, by leveraging the knowledge of attribute intensities, we have developed two un-
 190 derstandable metrics. In Section 4.1, we present Single attribute KL Divergence (SaKLD), which
 191 measures the distance of attribute distributions between training data and generated images. In Section
 192 4.2, we introduce Paired attribute KL divergence (PaKLD), a metric that assesses the relationship of
 193 attributes.

194 4.1 Single attribute KL Divergence (SaKLD)

195 We design SaKLD to distinguish a good generative model which produces the same quantity of each
 196 attribute present in the training data. For example, if 50,000 training data contains 3,000 images with
 197 eyeglasses, the model should generate exactly 3,000 images with eyeglasses. Any deviation from this
 198 ideal distribution is considered undesirable. We introduce a new metric that quantifies density of each
 199 attribute in dataset by utilizing interpretable embedding. Our metric, SaKLD, quantifies the difference
 200 in density for each attribute between the training dataset (\mathcal{X}) and the set of generated images (\mathcal{Y}).

201 We define SaKLD as

$$\text{SaKLD}(\mathcal{X}, \mathcal{Y}) = \frac{1}{N} \sum_i^N \text{KL}(\text{DCS}_{\mathcal{X}}(a_i), \text{DCS}_{\mathcal{Y}}(a_i)), \quad (5)$$

202 where i denotes an index for each attribute, N is the number of attributes, KL^* is Kullback-Leibler
 203 Divergence, and note that we denote the PDF of $\text{DCS}(x_i, a_i)$ for all $x_i \in \mathcal{X}$ as $\text{DCS}_{\mathcal{X}}(a_i)$.

204 We compare the PDFs of Directional CLIPScore for each attribute in \mathcal{X} and \mathcal{Y} . The DCS PDF for
 205 each attribute in \mathcal{X} and \mathcal{Y} represent the distribution of the amount of that attribute in the respective
 206 sets. If the distribution of the amount of a specific attribute in \mathcal{X} and \mathcal{Y} is similar, the DCS distri-
 207 bution will also be similar, and the PDFs of the two sets will be close. We used Kullback-Leibler
 208 Divergence(KLD) to compare the each Directional CLIPScore PDFs for their attribute in \mathcal{X} and \mathcal{Y} , to
 209 quantify the extent to which the generator has created too few or too many instances of a specific
 210 attribute. We then calculate the average KLD value between the PDFs of each attribute in \mathcal{X} and \mathcal{Y} to
 211 obtain the final value of SaKLD.

212 4.2 Paired attribute KL Divergence (PaKLD)

213 We design another metric, PaKLD for examining that generated images preserve the attribute re-
 214 lationships present in training data. The model should generate images that adhere to the attribute
 215 relationships observed in the training data. For instance, if all 50,000 male images in the training data
 216 wear glasses, then all generated male images should also wear glasses. To evaluate the preservation
 217 of attribute relationships, we compare the difference in the joint probability density distribution
 218 of attribute pairs between training data. Our proposed metric, Pairwise Attribute KL Divergence
 219 (PaKLD), is defined with joint probability density functions as follows:

$$\text{PaKLD}(\mathcal{X}, \mathcal{Y}) = \frac{1}{M} \sum_{(i,j)}^M \text{KL}(\text{DCS}_{\mathcal{X}}(a_{i,j}), \text{DCS}_{\mathcal{Y}}(a_{i,j})), \quad (6)$$

220 where $M = nP_2$, (i, j) denotes an index pair of attributes, and the pair of attributes’ joint PDF is
 221 denoted as $\text{DCS}_{\mathcal{X}}(a_{i,j})$.

Table 2: **Validation of metrics by including correlated images.** The first row shows metric scores between two distinct subsets of the FFHQ dataset (30,000 images each). The rest rows show the correlated-sample-injected-scores where only one of the subsets contains an additional 300 or 600 edited images. We examine the metric performance on ("man"- "makeup") and ("man"- "bangs") correlated images. All results are average values for five random subset pairs.

include edited images to one subset	SaKLD↓			PaKLD↓			FID↓	FID _{CLIP} ↓
	BLIP	USER	GPT	BLIP	USER	GPT		
not included	0.904	0.920	1.095	3.357	3.924	4.438	1.275	0.115
("man"- "makeup") 300	0.985	1.048	1.115	3.676	4.205	4.453	1.282	0.132
("man"- "makeup") 600	1.079	1.368	1.286	3.910	4.819	4.710	1.306	0.162
("man"- "bangs") 300	0.991	1.102	1.171	3.679	4.297	4.496	1.278	0.122
("man"- "bangs") 600	1.201	1.521	1.314	4.031	5.064	4.718	1.288	0.140

222 PaKLD analyzes the performance of the model more comprehensively. For example, if the generator’s
 223 probability density function for the attribute pair ("makeup", "long hair") significantly differs from
 224 that of the training data, we can infer that the generator does not preserve the ("makeup", "long hair")
 225 relationship. PaKLD allows to quantify the degree of preservation of attribute relationships and
 226 measure quantitative entanglements between attributes that have not been considered in previous
 227 researches.

228 5 Experiments

229 **Experimental details** To estimate the probability density function (PDF) of Directional CLIPScore
 230 (DCS) in the training data and generated images, we use Gaussian kernel density estimation. We
 231 sample 10,000 points from each PDF to obtain a discretized distribution and use it to calculate SaKLD
 232 and PaKLD. In all experiments, we use a set of $N = 20$ attributes.

233 5.1 Correlated Image Injection Experiment: Validating the Effectiveness of Our Metric

234 In this subsection, we provide a carefully designed experiment to compare the proposed metrics with
 235 FID; we first create two non-overlapping subsets of 30,000 images from FFHQ and consider them as
 236 training data \mathcal{X} and generated images \mathcal{Y} , respectively. We then compare the scores for all metrics
 237 after including the edited images in set \mathcal{Y} . Specifically, we use DiffuseIT[16] to prepare two sets
 238 of edited images: ‘man’ with ‘makeup’ and ‘man’ with ‘bangs’. We use CelebA attributes for user
 239 annotation method (denoted by USER in Table 2).

240 As shown in Table 2, our metrics and FID show consistent tendency: score increases when more
 241 edited images are included in imageset \mathcal{Y} . Furthermore, thanks to the nature of focusing on the
 242 attributes of the image domain, our metrics show more obvious numerical differences compared to
 243 FID. These results demonstrate that SaKLD successfully captures the attribute distribution difference
 244 and PaKLD captures the joint distribution difference between attribute pairs. Basically, our three
 245 attribute selection scenarios have similar tendencies across the two proposed metrics, but there are
 246 several differences. See supplement material for more details.

247 5.2 Necessity of PaKLD

248 We conducted another toy experiment, a scenario in which the SaKLD metric fails to detect a particular
 249 attribute relationship, while PaKLD metric successfully identified it. We define the curated subsets
 250 of CelebA-HQ as training data and generated images with discrepancies in attribute relationship.
 251 Specifically, for training data, we collect 20,000 ‘smiling men’ images and 20,000 ‘non-smiling
 252 women’ images using ground truth labels of CelebA-HQ. Conversely, the generated images consist
 253 of 20,000 ‘non-smiling men’ and 20,000 ‘non-smiling women’. In this scenario, the PDFs of the
 254 ‘man’, ‘woman’, and ‘smile’ attributes would not differ significantly between the two sets, and thus
 255 the SaKLD score would not capture it well. However, Paired attribute KL divergence would exhibit
 256 significant differences because the relationships between attributes within each set are completely
 257 different.

258 Figure 4 clearly illustrates the disparities in the evaluation results. While SaKLD score remained
 259 relatively unchanged for noteworthy attributes such as ‘man’, ‘woman’, and ‘smile’, the Paired

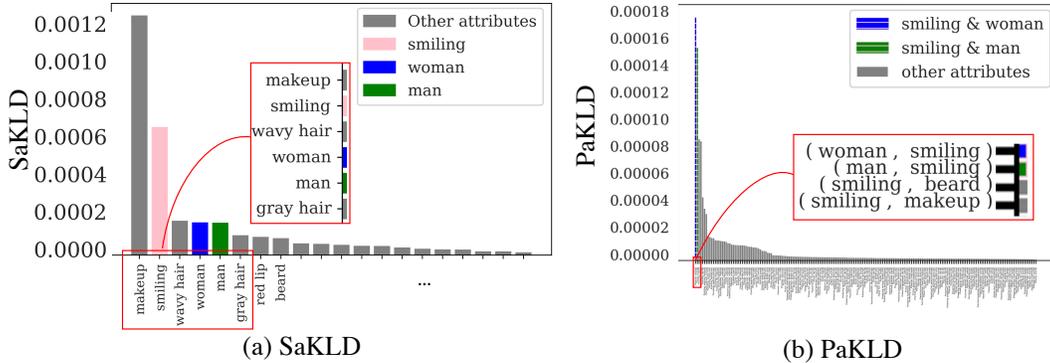


Figure 4: **Superiority of PaKLD.** We define the curated subsets of CelebA-HQ as training data, consisting of smiling men and non-smiling women, and generated images, consisting of non-smiling men and smiling women. (a) The most influential attribute on SaKLD is not the attribute we manipulate. (b) The most influential attributes on PaKLD provides explicit insights into the contributions of attribute pairs, such as (woman, smiling).

Table 3: **Comparing the performance of generative models.** We computed each generative model’s performance on our metric with their official pretrained checkpoints. For FFHQ[11] and LSUN Cat[29], we used 50,000 images for both GT and generated set, and we used 1,336 and 50,000 images for GT and generated set for MetFaces[13]. We used BLIP-extracted attributes for this experiment.

	SaKLD↓			PaKLD↓		
	FFHQ	LSUN Cat	MetFaces	FFHQ	LSUN Cat	MetFaces
StyleGAN1[11]	9.902	74.626	-	19.431	119.456	-
StyleGAN2[13]	6.377	63.601	-	12.838	100.896	-
StyleGAN2-ADA[12]	14.118	-	40.769	21.930	-	87.118
StyleGAN3[14]	5.993	-	31.140	12.285	-	58.065
iDDPM [23]	-	110.229	-	-	136.579	-
iDDPM(P2) [2]	12.040	-	129.627	21.507	-	230.720

260 attribute KL divergence score showed significant variations. This can be attributed to the distinct
 261 probability density functions (PDFs) of the ‘woman \cap smiling’. Note that we can easily understand
 262 the judgment factors; top attributes such as ‘woman \cap smiling’ and ‘man \cap smiling’ increase the
 263 score. These findings demonstrate the superior sensitivity and discernment of our proposed metrics,
 264 allowing for a more comprehensive evaluation of the generator’s generation ability.

265 5.3 Comparing generative models including GANs and diffusion models with our methods

266 Leveraging the superior sensitivity and discernment of our proposed metrics, we compare the
 267 performance of GANs and Diffusion Models (DMs) in Tables 3. Interestingly, there are two attractions;
 268 1) StyleGAN2-ADA shows the worst performance and 2) despite the respectable generative capability
 269 of DMs, iDDPM showed worse performance than StyleGAN models in all datasets.

270 The score of StyleGAN2-ADA implies that data augmentation for generative models may ruin
 271 attribute distribution in spite of FID’s superiority. Please refer to Appendix for an analysis. And we
 272 suppose that although there are many advantages of DMs, it is inferior to GANs in attribute-based
 273 analysis.

274 To investigate the reason for the inferiority of DMs, we leverage the flexibility of constructing
 275 attributes to analyze the score changes according to the characteristics of attributes. We constructed
 276 attributes that focus only on color (e.g., ‘yellow fur’, ‘black fur’) and attributes that focus on shape
 277 (e.g., ‘pointy ears’, ‘long tail’) for LSUN Cat.

278 Table 4 shows that iDDPM’s performance was particularly poor for color attributes. This is consistent
 279 with the assumption by Khruikov et al. [15] that the encoder map of DMs coincides with the optimal
 280 transport map for common distributions; which means the pixel-based Euclidean distance corresponds
 281 to high-level texture and color-level similarity regardless of dataset and model. Therefore, the color

Table 4: **Computing performance of models with different attributes for LSUN Cat.** Analyzing the weakness of iDDPM for specific attribute types, such as color or shape. We used BLIP-extracted attributes for this experiment.

	color attributes		shape attributes	
	SaKLD↓	PaKLD↓	SaKLD↓	PaKLD↓
StyleGAN1[11]	36.614	75.884	33.214	72.454
StyleGAN2[13]	36.621	67.518	34.642	68.954
iDDPM [23]	111.302	121.877	72.181	80.511

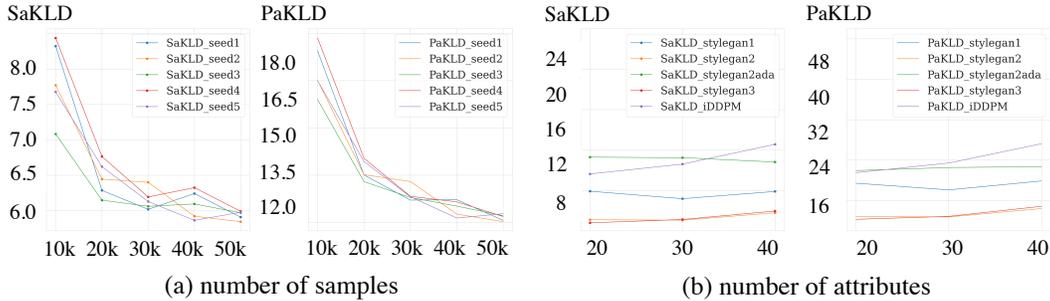


Figure 5: **(a) The effect of sample size on our metric.** Proposed metrics started to stabilize when using more than 50,000 images. **(b) The effect of the attribute counts on our metric.** Although depending on the characteristics of the additional attributes, the ranking of scores between models can vary, the rank of the models mostly remained consistent regardless of the number of attributes.

282 of the output images only depends on the initial latent noise x_T , and the Monge optimal transport
 283 map between training data and the standard normal distribution. We conclude that the distribution of
 284 color-related attributes is the inferiority of DMs.

285 5.4 Impact of Sample Size and Attribute Count on Proposed Metric

286 We provide ablation experiments to investigate the effect of a number of samples and attributes in
 287 Figure 5. We obtain generated images by StyleGAN3 from FFHQ with various random seeds. When
 288 the number of samples increases, SaKLD and PaKLD converge, especially more than 50,000 samples
 289 (Figure 5 (a)). We argue that the scores started to stabilize when using more than 50,000 images and
 290 note that we use 50,000 images for Tables 3 and 4. As for the number of attributes, we observe that
 291 the rank of the models mostly remained consistent regardless of the number of attributes. However,
 292 scores of DMs, purple line of Figure 5 (b), is increased as the number of attributes is increased
 293 because of color-related attributes. We argue that 20 attributes are sufficient, but more information
 294 can be obtained by using more diverse cases. Please see Appendix for an analysis of each score.

295 6 Discussion and Conclusion

296 In this paper, we introduce a novel metric that not only assesses the performance of the generator
 297 but also provides explicit explanations. Our proposed method, Directional CLIPScore, quantifies
 298 the attributes captured in an image and aligns them close to human judgment. Leveraging the
 299 interpretability of DCS, we propose two novel metrics, namely the SaKLD and PaKLD, which allow
 300 us to compare attribute appearance frequencies and examine attribute relationships, respectively.

301 While our metrics offer comprehensive explanations, unreliable results may arise when the attributes
 302 present in the images are ambiguous. For instance, in complex modern artworks with intricate color
 303 patterns, extracting appropriate attributes becomes challenging or even impossible, rendering our
 304 metric ineffective. Additionally, if the generative model’s ability is significantly poor, the same
 305 limitation arises: measuring DCS from generated images becomes challenging.

306 Despite these limitations, our research establishes a solid foundation for the development of explain-
 307 able evaluation metrics for generative models and contributes to the advancement of the field.

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11472–11481, 2022.
- [3] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [4] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [6] Jiyeon Han, Hwanil Choi, Yunjey Choi, Junho Kim, Jung-Woo Ha, and Jaesik Choi. Rarity score: A new metric to evaluate the uncommonness of synthesized images. *arXiv preprint arXiv:2206.08549*, 2022.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [8] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- [9] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [11] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [12] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *Advances in neural information processing systems*, 33:12104–12114, 2020.
- [13] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020.
- [14] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021.
- [15] Valentin Khulkov, Gleb Ryzhakov, Andrei Chertkov, and Ivan Oseledets. Understanding ddpn latent codes through optimal transport. *arXiv preprint arXiv:2202.07477*, 2022.
- [16] Gihyun Kwon and Jong Chul Ye. Diffusion-based image translation using disentangled style and content representation. *arXiv preprint arXiv:2209.15264*, 2022.
- [17] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in Neural Information Processing Systems*, 32, 2019.
- [18] Tuomas Kynkäänniemi, Tero Karras, Miika Aittala, Timo Aila, and Jaakko Lehtinen. The role of imagenet classes in frechet inception distance. *arXiv preprint arXiv:2203.06026*, 2022.

- 359 [19] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
360 training for unified vision-language understanding and generation. In *International Conference*
361 *on Machine Learning*, pages 12888–12900. PMLR, 2022.
- 362 [20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the
363 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 364 [21] Sachit Menon and Carl Vondrick. Visual classification via description from large language
365 models. *arXiv preprint arXiv:2210.07183*, 2022.
- 366 [22] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo.
367 Reliable fidelity and diversity metrics for generative models. In *International Conference on*
368 *Machine Learning*, pages 7176–7185. PMLR, 2020.
- 369 [23] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic
370 models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- 371 [24] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
372 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
373 models from natural language supervision. In *International conference on machine learning*,
374 pages 8748–8763. PMLR, 2021.
- 375 [25] Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. As-
376 sessing generative models via precision and recall. *Advances in neural information processing*
377 *systems*, 31, 2018.
- 378 [26] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
379 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 380 [27] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
381 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*
382 *preprint arXiv:2011.13456*, 2020.
- 383 [28] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Re-
384 thinking the inception architecture for computer vision. In *Proceedings of the IEEE conference*
385 *on computer vision and pattern recognition*, pages 2818–2826, 2016.
- 386 [29] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun:
387 Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv*
388 *preprint arXiv:1506.03365*, 2015.