

16558_gmvaluator_similarity_based_da. pdf

By CVPR Appeal (with reference)

WORD COUNT

12249

TIME SUBMITTED

29-FEB-2024 12:37AM

PAPER ID

107215663

GMValuator: Similarity-based Data Valuation for Generative Models

Anonymous CVPR submission

Paper ID 16558

Abstract

001 Data valuation plays a crucial role in machine learning.
002 Existing data valuation methods, mainly focused on dis-
003 criminative models, overlook generative models that have
004 gained attention recently. In generative models, data val-
005 uation measures the impact of training data on generated
006 datasets. Very few existing attempts at data valuation meth-
007 ods designed for deep generative models either concentrate
008 on specific models or lack robustness in their outcomes.
009 Moreover, efficiency still reveals vulnerable shortcomings.
010 We formulate the data valuation problem in generative mod-
011 els from a similarity-matching perspective to bridge the
012 gaps. Specifically, we introduce Generative Model Valu-
013 ator (GMVALUATOR), the first training-free and model-
014 agnostic approach to providing data valuation for genera-
015 tion tasks. It empowers efficient data valuation through
016 our innovative similarity matching module, calibrates bi-
017 ased contributions by incorporating image quality assess-
018 ment, and attributes credits to all training samples based
019 on their contributions to the generated samples. Addition-
020 ally, we introduce four evaluation criteria for assessing
021 data valuation methods in generative models, aligning with
022 the principles of plausibility and truthfulness. GMVALU-
023 ATOR is extensively evaluated on benchmark and high-
024 resolution datasets and various generative architectures to
025 demonstrate its effectiveness.

1. Introduction

026 As the driving force behind modern AI, particularly deep
027 learning [27], a substantial volume of data is indispensable
028 for effective machine learning. On one hand, informative
029 data samples relevant to the task at hand play a critical role
030 in the training process. On the other hand, due to data
031 privacy concerns, personal data is safeguarded by various
032 regulations, including the General Data Protection Regula-
033 tion (GDPR) [30], and has become a valuable asset. Con-
034 sequently, data valuation has garnered significant attention
035 from academic and industrial sectors recently.

036 The intricate relationship between data and model pa-

rameters presents a significant challenge in the contribution
038 measurement of each training sample, thus making data val-
039 uation a difficult task. Most existing data valuation studies
040 focus on supervised learning for discriminative models
041 (e.g. classification and regression). These methods can be
042 categorized as (1) Metric-based methods: Methods such
043 as Shapley Value (SV) and Banzhaf Index(BI) [9, 37] pro-
044 vide the assessment of data value by calculating marginal
045 contribution on performance metrics (e.g. accuracy or loss)
046 through retraining the model¹. (2) Influence-based meth-
047 ods: This line of methods measures data value by evalua-
048 ting influence on model parameters of data points [19, 26].
049 (3) Data-driven methods: These techniques avoid retraining
050 by leveraging data characteristics (like data diversity, gen-
051 eralization bound estimation, class-wise distance), though
052 generally necessitate data labels. [20, 39, 40].

053 Data valuation in the context of generative models has
054 NOT been well-investigated in the current literature. More-
055 over, there exist significant challenges in directly adapting
056 the aforementioned data valuation methods for discrimina-
057 tive models listed to generative models: *Firstly*, the chal-
058 lenge of applying metric-based methods arises from lacking
059 robust performance metrics in generative models, in con-
060 trast to the existence of commonly used metrics (e.g. accu-
061 racy or loss) in discriminative models. In addition, the ex-
062 pensive cost of retraining requirements is another obstacle
063 to using metric-based methods. *Secondly*, influence-based
064 methods may not perform well on non-convex objective
065 function of generative models [1]. Besides, estimating the
066 influence function in a deep generative model is expensive,
067 as it requires computing (or approximation) inverse Hes-
068 sian. *Thirdly*, for data-driven methods, they mainly focus
069 on supervised learning, which requires knowing data labels
070 to quantify data values.

071 To the best of our knowledge, limited studies have ex-
072 plored model-dependent data evaluation using influence
073 functions for specific generative models. IF4GAN [35] con-
074 sider multiple evaluation metrics, such as log-likelihood,
075 inception score (IS), and Frechet inception distance
076 (FID) [13] to identify the most responsible training samples

¹Exception for KNN-Shap [18]

078 for the overall performance of generative adversarial networks (GAN)-based models. However, selecting appropriate
 079 metrics is critical, as the results are inconsistent across
 080 metrics. VAE-TracIn [22] finds the most significant contributors
 081 for generating a particular generated sample for Variational
 082 Autoencoders (VAE) using the influence function.
 083 However, viable Hessian estimation in influence function
 084 calculations incur high computational costs and this method
 085 cannot be easily generalized to other generative models.
 086

087 Considering the challenges posed by the selection of performance
 088 metrics and computational efficiency and the need for broader applicability across various generative models,
 089 we aim to propose a unified and efficient data valuation
 090 method. We expect the unified method to satisfy the following
 091 key properties: 1) *Model-agnostic*: the method should
 092 be versatile, capable of being employed across diverse generative
 093 model architectures and algorithms, regardless of specific
 094 design choices; 2) *Computation Efficiency*: the method does not require retraining the model and minimize
 095 computational overhead while maintaining a satisfactory
 096 level of accuracy and reliability in evaluating the value
 097 of data points; 3) *Plausibility*: the method should evaluate
 098 the value of data based on its alignment with human prior
 099 knowledge on the task, enhancing its credibility and reliability;
 100 4) *Truthfulness*: the method should strive to provide an unbiased and accurate assessment of the value associated
 101 with individual data points. In this work, we evaluate the contribution of training data given a set of (good) generated
 102 data from a certain well-trained generative algorithm. To meet the design purposes, we propose a similarity-based
 103 data valuation approach for the generative model, called
 104 GMVALUATOR, which is *model-agnostic* and *efficient*.
 105

106 In summary, our work here provides the following specific
 107 novel contributions:
 108

- To the best of our knowledge, GMVALUATOR is the first modal-agnostic and retraining-free data valuation method for generative models.
- We formulate data valuation for generative models as an efficient similarity-matching problem. We further eliminate the biased contribution measurement by introducing image quality assessment for calibration.
- We propose four evaluation methods to assess the truthfulness of data valuation and evaluate GMVALUATOR on different datasets (including benchmark datasets and high-resolution large-scale datasets) and various deep generative models to verify GMVALUATOR’s validity.

109 **Related Work:** Different from discriminative models for regression or classification tasks, generative models in various forms (*i.e.* Variational auto-encoders (VAEs) [31], Variational auto-encoders (VAEs) [31], Diffusion Model [16, 33]) aim to learn data distribution for generation tasks. Consequently, the approach to data valuation in our work, which focuses on generative models, is distinct and orthogonal to

110 that for discriminative models. The primary limitation of both *metric-based methods* and *influence-based methods* is
 111 the expensive computation cost, a drawback that is further magnified in generative models. Conversely, while data-driven methods are training-free, they predominantly concentrate on supervised learning and data itself [20, 39, 40].
 112 Apart from the expensive computational cost, the limited existing *influence-based methods* for generative models are
 113 model-specific, which can not adapt to the diverse and evolving trends in generative model development [22, 35].
 114 A more detailed related work is discussed in the appendix.
 115

2. Motivation and Problem Formulation

In this section, we present the background and motivation behind our proposed method GMVALUATOR, which formulates the data valuation for generative models as a similarity matching problem.

2.1. Motivations

Let us denote $X = \{x_1, \dots, x_n | x \sim \mathcal{X}\}$ as the training dataset without duplicated data for a generation task, and n is regarded as the size of the dataset. Denote the generated dataset as $\hat{X} = \{\hat{x}_1, \dots, \hat{x}_m\}$ by a well trained generative model G (*e.g.* Generative Adversarial Network (GAN), Variational Autoencoder (VAE), Diffusion Model) on X . The core of these generative models is an intractable probability distribution learning process to train a generator G , that maps each z sampled from the noise distribution \mathcal{Z} to a generated sample \hat{x} in estimated distribution by maximizing the likelihood $p_{\mathcal{X}}(x) \approx \int p_G(x|z)p_{\mathcal{Z}}(z)dz$. This can be achieved by minimizing the distance $d(\cdot, \cdot)$ (*i.e.*, Wasserstein Distance) between the estimated distribution and training data distribution:

$$G^* = \arg \min_G d(G_{z \sim \mathcal{Z}}(z), \mathcal{X}). \quad (1)$$

Therefore, the data generated by an optimal generator, which closely approximates the training distribution, can be considered as drawing from a subdistribution of \mathcal{X} .

2.1.1 Theoretical Justification

Let $T \subset X$ denote the K contributors found by a data valuator (*e.g.*, GMVALUATOR) for \hat{X} , and S is a subset of the training dataset X . We focus on the data type with describable attributes (*e.g.*, image data with semantic attributes). We first provide definitions of data and contributors.

Definition 2.1 (*Data with Describable Attributes.*) We assume an image can be characterized by V attributes, and there is a labeling function f mapping S, T and \hat{X} to the same attribute space $\mathcal{A} = \{0, 1\}^V$.

176 **Definition 2.2** (*Contributors in Generative Models.*) Let
 177 S^* denotes the real K contributors, defined as follows:

$$178 \quad S^* = \arg \min_{S \subset X} d(\mathcal{X}_{(\hat{X}|A)}, \mathcal{X}_{(G(S)|A)}), \quad (2)$$

179 where $|S^*| = K$ and K is a reasonable number for
 180 generative model training. $\mathcal{X}_{(\hat{X}|A)}$ is the data distribution of
 181 generated data \hat{X} on attributes $A \in \mathcal{A}$, and $\mathcal{X}_{(G(S)|A)}$ is the
 182 distribution of data with attribute A that are generated by
 183 the optimal generator trained by contributors S . According
 184 to the objective of generative models and Eq. (2), we have
 185 $\mathcal{X}_{(\hat{X}|A)} \sim \mathcal{X}_{(G(S^*)|A)}$. We follow [20] on assumptions and
 186 lemmas that will be used to obtain the theorem.

187 **Assumption 2.3** Assume that the function f is ϵ -Lipschitz
 188 and the loss function $\mathcal{L} : \{0, 1\}^V \times \{0, 1\}^V \rightarrow \mathbb{R}^+$ is k -
 189 Lipschitz in both inputs and attributes. We have labeling
 190 functions that are all bounded by V as $\|f\| \leq V$.

191 Then, we introduce the error bound between T and S^* .

192 **Theorem 2.4** (*Bounded Attributes Classification Error on*
 193 S^* *to T.*) Let $f'_{S^*} : \mu \rightarrow \mathcal{A} = \{0, 1\}^V$ be the model trained
 194 on the optimal contributor dataset S^* . Following Assump-
 195 tion 2.3, if the contributors are corresponding to the given
 196 generated data \hat{X} , we have:

$$197 \quad \mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] - \mathbb{E}_{x \sim \mu_S} [\mathcal{L}(f(x), f'_{S^*}(x))] \\ 198 \leq \xi \epsilon \cdot [d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)}) + d_W(\mathcal{X}_{(S^*|f)}, \mathcal{X}_{(\hat{X}|f)})] + B_2 \quad (3)$$

where B_2 is in order $\mathcal{O}(\xi V)$.

199 As shown in the Theorem 2.4, given the optimal contribu-
 200 tor S^* , $d_W(\mathcal{X}_{(S^*|f)}, \mathcal{X}_{(\hat{X}|f)})$ is a deterministic term that
 201 approaching zero. Then, approximate S^* can be achieved
 202 by reducing the distance term $d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)})$. Please
 203 see the appendix for the proof.

204 2.1.2 Empirical Validation

According to theoretical justification, the generated data
 should show more similarity to the data samples used for
 training. In this context, we examine the value of training
 data being more valuable compared to data that wasn't used
 for training, despite originating from a similar distribution.
 We support this by partitioning a class of CIFAR-10 (the
 class is plane here) into two non-overlapped subsets, de-
 noted as X_{v1} and X_{v2} .² Next, we keep X_{v1} as non-training
 data and use X_{v2} as training data to train a BigGAN [2]
 and generate dataset \hat{X} . If our assumption holds, the
 generated data will be more similar to the training data X_{v2} .

²In practice, we perform T-SNE first to randomly split the non-overlapped samples from the embedding space.

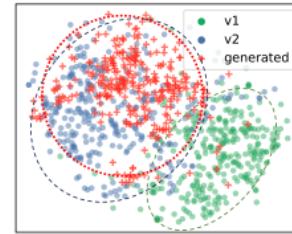


Figure 1. Data distribution for X_{v1}, X_{v2} and \hat{X} for CIFAR-10, and X_{v1}, X_{v2} are both airplane dataset.

As presented in the T-SNE [36] plot (Figure 1), the generated data demonstrate a more substantial overlap in the shared feature space with the training dataset X_{v2} than the non-training ones X_{v1} .³ This observation motivates us to address the data valuation for the generative model from a similarity matching perspective.

2.2 Problem Formulation and Challenges

The primary objective of our research is to tackle the issue of data valuation in generative models using *black-box* access. In essence, given a fixed (good) set of data samples \hat{X} with a size of m generated from the well-trained deep generative model G^* , our aim is to determine the value $\phi_i(x_i, \hat{X}, G^*)$ associated with each data point $x_i \in X$ for $i \in [n]$ in the *deduplicated* training dataset, which contributes to the generated dataset. We denote the value for training data x_i as ϕ_i in the rest of the paper for simplicity.

Following the motivation stated in Sec. 2.1, ϕ_i should be a function of the distance between the data points. We denote the distance between training data x_i and generated data \hat{x}_j as d_{ij} , for $i \in [n]$ and $j \in [m]$.

Definition 2.5 (*Primary Contribution Score.*) The contribution score of x_i to \hat{x}_j is denoted as $\mathcal{V}(x_i, \hat{x}_j) \propto d_{ij}^{-1}$, which is inversely proportional to distance since maximizing the log-likelihood of a generative model is equivalent to minimizing the dissimilarity of real and generated data distribution. To link the dissimilarity $-d_{ij}$ to likelihood, we choose $\exp(-d_{ij})$ to likelihood

$$243 \quad \mathcal{V}(x_i, \hat{x}_j) = \frac{\exp(-d_{ij})}{\sum_i \exp(-d_{ij})}. \quad (4)$$

Therefore, an intuitive definition for the value of data sample x_i can be written as below.

Definition 2.6 (*Data Value.*) The contribution of each training data point $x_i \in S^*$ for the generation of dataset \hat{X} equals to its the sum of contributions to each generated data point \hat{x}_j in \hat{X} .

$$250 \quad \phi_i = \sum_{j=1}^m \mathcal{V}(x_i, \hat{x}_j), x_i \in X, \hat{x}_j \in S^*. \quad (5)$$

³Quantitative statistical testing is provided in Appendix.

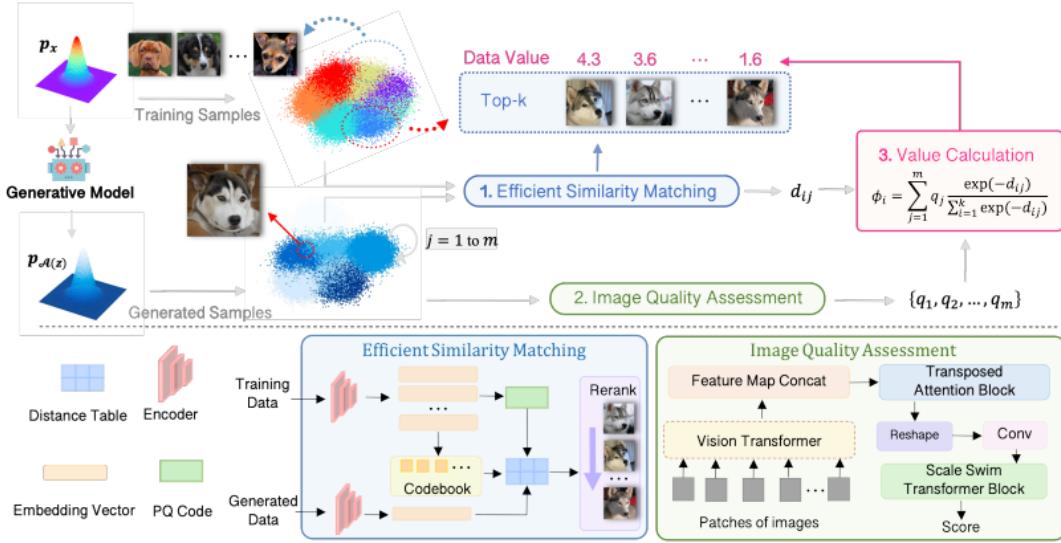


Figure 2. Overview of GMVALUATOR, a unified and training-free data valuation approach for any generative models. GMVALUATOR contains three important modules – (1) Efficient Similarity Matching (ESM), (2) Image Quality Assessment, and (3)Value Calculation. Each generated data \hat{x}_j is matched with training data through ESM approach, resulting in the distances with its top k contributors. The normalized contribution score from training sample x_i to \hat{x}_j , defined as $\exp(-d_{ij})/\sum_i^k \exp(-d_{ij})$, is adjusted based on the quality of the associated generated samples q_j . We compute the data value ϕ_i of each training sample x_i by summing its contributions to the generated samples, where it ranks among the top k contributors.

251 The high-value data will achieve a smaller distance to
 252 the target distributions, thus, a better approximation. Therefore,
 253 data valuation in the above problem formulation contains two steps: *Step 1*, calculating all of the score value
 254 $\mathcal{V}(x_i, \hat{x}_j)$; *Step 2*, mapping the contribution from training
 255 data to generative data based on the scores \mathcal{V} .

256 However, there are several open questions and challenges in performing the above two steps for calculating
 257 Eq. (5):

258 *Challenge 1: Efficiency.* In step 1, considering n training
 259 samples and m generated samples, where $\mathcal{O}(C)$ represents
 260 the complexity of the selected pair-wise distance calcula-
 261 tion, the total complexity of this step amounts to $\mathcal{O}(mnC)$.
 262 In practical scenarios with large training datasets (e.g., $n >$
 263 $10K$), the computation cost becomes prohibitively expen-
 264 sive. Additionally, fitting such a large collection of high-
 265 dimensional data for distance calculation can pose signifi-
 266 cant challenges in system memory.

267 *Challenge 2: Contribution plausibility.* To ensure that a
 268 training data point contributes more if it is similar to high-
 269 quality generated data and less if it is similar to low-quality
 270 generated data, the contribution scores should be adjusted
 271 based on the quality of the generated data.

272 *Challenge 3: Non-zero scores.* In practical scenarios, the
 273 distance between training data and the least similar gen-
 274 erated data is not infinite, which may result in false non-zero
 275 contribution scores. With a large dataset size, the accumu-
 276 lation of these noisy scores yields biased data valuation.

In light of this, we present a novel and efficient data valuation approach suitable for agnostic generative models, termed as GMVALUATOR as elaborated in Section 3.

3. The Proposed Data Valuation Methods

The crucial idea behind GMVALUATOR is to transform the data valuation problem into a similarity matching problem between generated and training data. The overview of GMVALUATOR is presented in Figure 2. To tackle *challenge 1*, we propose to employ efficient similarity matching (ESM), where each generated data point can be linked to multiple contributors from the training dataset (Section 3.1), where each generated data firstly link with top k contributors via recall phase and then re-rank by a refined similarity for effectiveness. After that, the image quality of the generated sample is assessed to weigh the valuation (Sec. 3.2). Finally, the value computation function combines both the quality score and the image-space similarity score to measure the value of the training data (Sec. 3.3).

3.1. Efficient Similarity Matching

Considering the complexity of calculating $\mathcal{V}(x_i, \hat{x}_j)$, we formulate it as an ESM problem between generative data and training data. Each generated data sample is matched to several training data samples based on their similarity. We denote $\mathcal{P}_j = \{x_1, x_2, \dots, x_k\} = f(X, \hat{x}_j)$ as the subset of training data that contains the $k \ll n$ most similar data samples. Here, f represents the similarity-matching strat-

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304



Figure 3. The value without generated image quality calibration for q high-quality image (top row) and a low-quality image (bottom row). Column 1: generated images. Column 2-5:their top 4 contributors.

305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
egy, encompassing the *recall* and *re-ranking* phases, which will be introduced subsequently.

Recall Phase: The main aim of the *recall* phase is to rapidly identify a subset of training samples that are similar to a generated sample. To achieve this, an initial step involves encoding all original training images and generated samples from the image space R^C to a lower-dimensional embedding space R^D using a pre-trained encoder f_e (such as CLIP [29]) to reduce computational complexity. Subsequently, the technique of *Product Quantization* (PQ) [17] is employed to further decrease the computational burden. Specifically, PQ divides embedding vectors into subvectors and independently quantizes each subvector through Q -means clustering. This process generates compact PQ codes that serve as representations of the original vectors. This representation significantly reduces the vector sizes, allowing for efficient estimation of Euclidean Distance between two samples. By incorporating the *recall* method into the similarity matching process, the computational complexity is lowered from $\mathcal{O}(mnC)$ to $\mathcal{O}(mQD)$, where $D \ll C$ and $Q \ll n$. Consequently, generated images can quickly identify their top- k most similar training data samples.

Re-Ranking Phase: Following the PQ-based efficient recall process in GMVALUATOR, we further improve the precision of the results by utilizing perceptual similarity [7] for precision ranking. Once we have extracted perceptual features from the top k recalled training samples, we proceed to calculate the distance for each pair of items. To obtain precise distance measurements based on their perceptual content, we propose to use Learned Perceptual Image Patch Similarity (LPIPS) [42] or DreamSim [7] as the distance measurement d to gain insights into the perceptual dissimilarity between the generated sample and different training samples. These metrics enable us to precisely measure the most significant contributors according to their perceived similarity, more importantly, the obtained distance d will be employed to compute data valuation in Eq. (6). We do not use Wasserstein distance as it is more proper to measure dissimilarity between two probability distributions rather than a pair of image instances with semantic characteristics.

3.2. Image Quality Assessment

345

To tackle *challenge 2*, we have to establish a connection between the quality of the generated sample that the training samples contribute to and the contribution score. This is necessary before assigning contribution scores to the ranked training samples for the generated sample. For low-quality generated samples \hat{x}_{low} , we expect their total contribution score from contributors to be lower compared to high-quality samples \hat{x}_{high} , namely $\sum_{i \in [k]} \mathcal{V}(x_i, \hat{x}_{\text{low}}) < \sum_{i \in [k]} \mathcal{V}(x_i, \hat{x}_{\text{high}})$. This motivation arises from the observations in Figure 3, where the first and second rows show generated samples using z from a normal distribution and uniform distribution, respectively. The values denoted in Figure 3 directly employs Eq (4), so we have $\sum_{x_i \in \mathcal{P}} \mathcal{V}(x_i, \hat{x}) = 1$ for all $\hat{x} \in \hat{\mathcal{X}}$. The real samples are noticeably dissimilar to the generated data, yet they are still assigned a high value. Therefore, we propose to calibrate the contribution scores with the generated data quality. Specifically, we obtain a comprehensive quality score $q_j \in [0, 1]$ for each generated image integrated using MANIQA [41] model into our evaluation process. A higher q_j indicates better data quality. This score of MANIQA considers various factors such as sharpness, color accuracy, composition, and overall visual appeal. Incorporating the image quality evaluation provided by MANIQA allows us to more accurately assess the generated images and take into account their perceptual fidelity and aesthetic qualities.

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

Algorithm 1 GMVALUATOR

Input: Training dataset $X = \{x_i\}_{i=1}^n$, a well-trained model G^* , random distribution \mathcal{Z} .

Output: Generated dataset \hat{X} , the value of training data points $\Phi = \{\phi_1, \phi_2, \dots, \phi_n\}$

```

1: // Generate the synthetic dataset
2:  $\hat{X} = \{\hat{x}_j\}_{j=1}^m \leftarrow G^*(z_j)$ , for  $z_j \in \mathcal{Z}$ 
3: for  $\hat{x}_j$  in  $\hat{X}$  do
4:   // Matching process (see Sec. 3.1)
5:    $\mathcal{P}_j = f(X, \hat{x}_j)$  // Including two phases
6:   for  $x_i$  in  $\mathcal{P}_j$  do
7:      $d_{ij} \leftarrow \text{DreamSim}(x_i, \hat{x}_j)$  or others
8:   end for
9:   // Image Quality Assessment (see Sec. 3.2)
10:   $q_j = \text{MANIQA}(\hat{x}_j)$ 
11:  // Contribution Score Calculation (see Sec. 3.3)
12:  Calculate score  $\mathcal{V}(x_i, \hat{x}_j, d_{ij}, q_j)$  using Eq. (6)
13: end for
14: // Calculation of data value and return the result  $\Phi$ 
15: for  $x_i$  in  $X$  do
16:   Calculate  $x_i$ 's value  $\phi_i$  using Eq.(5)
17: end for
18: return  $\Phi = \{\phi_1, \phi_2, \dots, \phi_n\}$ 

```

371

372 3.3. Value Calculation

373 Finally, we utilize the image quality q to calibrate the contribution scores. Notably, during the recall phase of similarity matching, we select the top k contributors \mathcal{P}_j for the generated data \hat{x}_j . Consequently, we only consider the contribution score between \mathcal{P}_j and \hat{x}_j by setting $\mathcal{V}(x_i, \hat{x}_j) = 0$ for $x_i \notin \mathcal{P}_j$. This strategy addresses *challenge 3* by assigning zero scores to irrelevant samples, effectively reducing bias and noise in value estimation when dealing with a large n . Hence, we define the contribution score of each training data point for a specific synthetic data point as follows:

$$383 \quad \mathcal{V}(x_i, \hat{x}_j, d_{ij}, q_j) = \begin{cases} q_j \cdot \frac{\exp(-d_{ij})}{\sum_{i=1}^k \exp(-d_{ij})} & x_i \in \mathcal{P}_j \\ 0 & x_i \notin \mathcal{P}_j \end{cases} \quad (6)$$

384 which is used to give different scores to training data points
385 according to their ranking distance. The final data value is
386 obtained by plugging Eq. (6) into Eq. (5).

387 4. Experiments

388 4.1. Experiment Setup

389 **Datasets.** Our experiment settings are listed in Table 1.
390 Particularly, we consider different types of generative models,
391 including GAN based models, β -VAE [14], and diffusion
392 models [15, 38]. The generation tasks are conducted on
393 benchmark datasets (*i.e.* MNIST [24] and CIFAR [23]), face
394 recognition dataset (*i.e.* CelebA [25]), high-resolution
395 image dataset with size 512×512 , and $1024 \times$
396 1024 (*i.e.* AFHQ [3], FFHQ [21]), and also the large-scale
397 dataset with 1,000 classes and 14,197,122 images (*i.e.* ImageNet [5]).

398 **Generative Models.** Note that GMVALUATOR is a model-
399 agnostic method. To showcase its efficacy, we compare
400 it with two baseline methods: VAE-TracIn [22] and
401 IF4GAN [35]. VAE-TracIn identifies the most significant
402 contributors for a generated sample, which is the same as the
403 similarity matching process of GMVALUATOR (Sec. 3.1). Thus,
404 we compare GMVALUATOR with VAE-TracIn in Sec. 4.3 and Sec. 4.4. IF4GAN finds the
405 high-valued training samples for generative model training,
406 which we regard it as the baseline method in section 4.5.
407 We also examine our approach using DreamSim, LPIPS and
408 l_2 -distance in the re-ranking phase, referred to as GMValuator
409 (DreamSim), GMValuator (LPIPS) and GMValuator (l_2 -
410 distance) respectively, while the GMValuator without the
411 re-ranking phase is referred to as GMValuator (No-Rerank).

414 4.2. Metrics

415 Given the absence of a definitive benchmark to evaluate data
416 valuation methods in the context of generative models, we
417 propose to evaluate the truthfulness of data valuation outcomes
418 through the examination of four criteria (C):

- **C1: Identical Class Test.** Following the concept of an identical class [12], we posit that the most significant contributors among the training samples should belong to the same class as the generated sample produced by a well-trained generative model, which is also used for evaluation in VAE-TracIn [22].

- **C2: Identical Attributes Test.** Following C1, we extend the concept of identical class test and consider the attributes of samples as the ground truth to examine the most significant contributors by our approach on datasets without class labels. We examine the overlap level of attributes (*i.e.* the number of identical attributes) between the most significant training samples with the generated sample.

- **C3: Out of Distribution Detection.** Under the assumption, the out of distribution training data samples (*e.g.* noisy samples) can be identified as low-contribution training samples for the generated dataset. To evaluate the performance of data valuation approaches, we consider the contribution level of noisy data samples as the performance metric across various approaches.

- **C4: Efficiency.** As an efficient training-free approach, GMVALUATOR should measure the data value in a limited time and be much more efficient than the previous work while obtaining truthful results.

Table 1. Evaluation Setup.

	Dataset	Network
C1	MNIST, CIFAR-10	GAN, Diffusion, β -VAE
	ImageNet	Masked Diffusion Transformer [8]
C2	CelebA	Diffusion-StyleGAN
	AFHQ, FFHQ	StyleGAN
C3	MNIST	DCGAN
C4	Hardware Environment	
CPU	One RTX 3080 (10GB) GPU	
GPU	12 vCPU Intel(R) Xeon(R), Platinum 8255C CPU @ 2.50GHz	

414 4.3. Identical Class Test (C1)

415 **Methodology.** In this subsection, we follow setup of
416 VAE-TracIn [22] by training separate β -VAE models on
417 MNIST [23], CIFAR [24]. We then attribute the most significant
418 contributors by VAE-TracIn and our approach over
419 generated samples. By the concept of identical class test,
420 we expect a perfect data valuation approach should indicate
421 training data in the same subclass of the generated
422 sample \hat{x}_j contribute more to \hat{x}_j . We examine GMVALUATOR
423 (DreamSim), GMVALUATOR (LPIPS), GMValuator
424 (l_2 -distance) and GMVALUATOR (No-Rerank) separately.
425 In addition to using VAE and comparing it with the baseline
426 method, we also perform experiments using masked diffusion
427 transformer on the large-scale dataset (ImageNet) and other
428 various generative models (in the appendix) to demonstrate
429 the model-agnostic property of GMVALUATOR.

430 **Results.** For a given generated data, we examine the
431 class(es) of is top k contributors in the training data. We
432

Table 2. Performance comparison of Identical Class Test (C1).

MNIST (%)	$k=30$	$k=50$	$k=100$
VAE-TracIn	72.00	71.11	68.58
GMValuator (No-Rerank)	86.41	85.13	85.95
GMValuator (l_2 -distance)	87.76	86.95	85.92
GMValuator (LPIPS)	88.78	88.19	86.69
GMValuator (DreamSim)	88.78	88.05	86.84
CIFAR-10 (%)	$k=30$	$k=50$	$k=100$
VAE-TracIn	6.28	3.77	1.88
GMValuator (No-Rerank)	72.66	72.25	70.71
GMValuator (l_2 -distance)	72.66	72.25	70.71
GMValuator (LPIPS)	72.60	71.75	70.84
GMValuator (DreamSim)	77.94	76.41	73.47
ImageNet (%)	$k=30$	$k=50$	$k=100$
GMValuator (No-Rerank)	56.18	53.6	48.17
GMValuator (l_2 -distance)	42.44	40.51	39.0
GMValuator (LPIPS)	50.30	47.51	45.13
GMValuator (DreamSim)	77.27	70.90	58.89

462 count the number of training Q samples in the top k contributors, which have the identical class as the generated data.
 463 The identical class ratio $\rho = Q/k$. We report the averaged ρ
 464 over the generated datasets (the data size $m=100$) on different choices of k in Table 2. GMVALUATOR has the most
 465 contributors that belong to the same class with the generated sample among on MNIST and CIFAR-10, while GMVALUATOR (DreamSim) has the best performance. The results for
 466 CIFAR-10 using VAE-TracIn are extremely bad. This could
 467 be attributed to underfitting, as the authors mentioned in
 468 their study [22], training a good β -VAE model on CIFAR-
 469 10 is challenging. In addition, GMVALUATOR (DreamSim)
 470 shows a significant advantage on ImageNet since it considers semantic characteristics while l_2 -distance or LPIPS de-
 471 stroy the performance. In the appendix, further analysis and
 472 explanation for this phenomenon are discussed and the re-
 473 sults from various generative models are presented as well.
 474

4.4. Identical Attributes Test (C2)

480 **Methodology.** In extension to C1, we focus on certain im-
 481 age attributes instead of class labels as the ground truth.
 482 We posit that the most significant contributors to a gener-
 483 ated sample should share similar attributes with it. We train
 484 Diffusion-StyleGAN [38] on CelebA, and leverage our ap-
 485 proach to identify the most significant contributors in train-
 486 ing samples and use some attributes including hat, gender,
 487 and eyeglasses as the ground truth, to check the correctness
 488 of our approach for identifying the most significant con-
 489 tributors for a generated sample. We also examine this on
 490 AFHQ, FFHQ dataset by StyleGAN [3].

491 **Results.** The results in Table 3 indicate that GMVALUATOR
 492 can find the most significant contributors with the same at-
 493 tributes to the generated sample. Besides, using DreamSim
 494 in the similarity re-ranking phase obtains the best perfor-
 495 mance over these three attributes. We also visualize some

Table 3. Performance of Identical Attributes Test (C2) of some attributes including Hat, Gender, and Eyeglasses on CelebA.

Top K contributors:	$k=5$	$k=10$	$k=15$
Attribute: Hat (%)			
GMValuator (No-Rerank)	92.52	92.12	91.92
GMValuator (l_2 -distance)	90.91	89.90	89.83
GMValuator (LPIPS)	96.77	96.57	96.90
GMValuator (DreamSim)	97.78	97.07	96.90
Attribute: Gender (%)			
GMValuator (No-Rerank)	75.15	75.25	75.82
GMValuator (l_2 -distance)	63.64	61.01	60.00
GMValuator (LPIPS)	97.98	97.48	97.37
GMValuator (DreamSim)	99.19	98.99	98.79
Attribute: Eyeglasses (%)			
GMValuator (No-Rerank)	91.52	91.52	91.45
GMValuator (l_2 -distance)	94.95	94.65	93.80
GMValuator (LPIPS)	94.95	94.65	93.80
GMValuator (DreamSim)	96.77	96.26	95.96

visible results in Figures on CelebA, AFHQ and FFHQ in the appendix. As we can see, the most significant contributors have similar attributes (e.g. skin color, hair) to the generated sample. Similarly, the visualized results on AFHQ and FFHQ shown in Figure 6 demonstrate that the top k contributors have similar attributes to the generated sample such as the fur color of cats or dogs in AFHQ or hair color in FFHQ.

4.5. Out of Distribution Detection (C3)

Methodology In this experiment, we introduce 100 noisy images into the MNIST dataset to create a new corrupted training dataset and each noisy image is generated by adding Gaussian noise to the clean data. These contaminated samples can diminish the performance of generative models, and as a result, they are anticipated to possess lower values compared to the rest of the training dataset. Our objective is to assess the performance of both GMVALUATOR and its baseline method, IF4GAN [35], by analyzing the value rankings of noisy samples. An effective approach should place noisy data in lower positions within the value ranking. We conducted this evaluation using 10,000 generated images and referred to the study [35] for the baseline method (IF4GAN).

Results. The results depicted in Figure 5 demonstrate that the values of all noisy training samples, as calculated by GMValuator, are lower compared to the values calculated by IF4GAN. This observation suggests that the performance of GMVALUATOR is significantly better than that of IF4GAN. Besides, the performance of GMVALUATOR with the re-ranking phase is better than GMVALUATOR (No-rerank).



Figure 4. Visualization of Identical Attributes Test on AFHQ and FFHQ. The results shown in the first and second subfigures on the left are conducted on AFHQ-Cat and AFHQ-Dog, respectively. The subfigure on the right presents the results on FFHQ. In each subfigure, the generated samples are on the left, and the top k contributors in the training dataset are on the right.

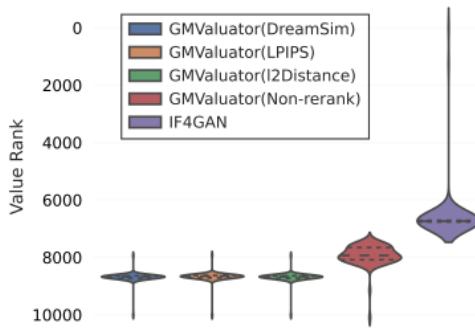


Figure 5. The y -axis represents the ranking of values from high to low, with the top being the highest value and the bottom being the lowest value. The x -axis represents the index of each noisy data.

527 4.6. Efficiency (C4)

528 **Methodology.** We thoroughly evaluate and compare the
529 efficiency of our approach against existing data valuation
530 methods for generative models discussed above: VAE-
531 TracIn and IF4GAN. Specifically, we measure the attribute
532 time for one generated sample on an average of 100 test
533 samples using both VAE-TracIn and our proposed approach
534 on MNIST and CIFAR-10. Additionally, under the same
535 setting in C3, we assess the data valuation time when uti-
536 lizing IF4GAN and our approaches on noise MNIST which
537 we mentioned in C3.

538 **Results.** The average time taken to attribute the most sig-
539 nificant contributors for one generated sample is calculated
540 and compared with the baseline method (VAE-TracIn), as
541 demonstrated in Table 4. Our approaches are all greatly
542 more efficient than the baseline methods on both datasets.
543 When it comes to data valuation for GAN, GMVALUATOR
544 (No-rerank) and GMVALUATOR (LPIPS) are significantly
545 better than IF4GAN. The aforementioned phenomenon is
546 attributed to the costly nature of the Hessian estimation pro-
547 cess, despite the utilization of certain acceleration methods.
548 It is noticeable that GMVALUATOR (DreamSim) is lightly
549 time-consuming when compared to IF4GAN. This lead to a
550 trade-off problem due to the GMVALUATOR using Dream-
551 Sim in the re-ranking phase obtained the better performance

from C1 to C3.

Overall, through the above experiments, GMVALUATOR outperforms than baseline methods, while GMVALUATOR (DreamSim) obtains the best performance. This is due to the fact that DreamSim captures mid-level similarities in image semantic content and layout compared to LPIPS. And the selection of different GMVALUATORS can be combined with the consideration of their efficiency.

Table 4. Efficiency Comparison

	Attribute for VAE	Time(s)
MNIST	VAE-TracIn	47.945
	GMValuator (No-Rerank)	0.250
	GMValuator (l_2 -distance)	0.339
	GMValuator (LPIPS)	0.477
	GMValuator (DreamSim)	1.709
CIFAR-10	VAE-TracIn	66.178
	GMValuator (No-Rerank)	0.755
	GMValuator (l_2 -distance)	1.226
	GMValuator (LPIPS)	2.412
	GMValuator (DreamSim)	15.491
	Data Valuation for GAN	Time(s)
Noise MNIST	IF4GAN	14.543
	GMValuator (No-Rerank)	2.137
	GMValuator (l_2 -distance)	3.388
	GMValuator (LPIPS)	4.771
	GMValuator (DreamSim)	17.086

5. Conclusion

To measure the contribution of each training data sample, we propose an efficient approach, GMVALUATOR, for generative model. As far as we are aware, there is no prior model-agnostic and training-free data valuation approach for generative models util GMVALUATOR. Our approach is based on efficient similarity matching, and it enables us to calculate the final value of each training data point, aligning with plausible assumptions. The proposed method is validated through a series of comprehensive experiments to showcase its truthfulness and efficacy on four criteria. In the future, we will validate the proposed methods on other data modalities.

573 **References**

- [1] Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. If influence functions are the answer, then what is the question? *Advances in Neural Information Processing Systems*, 35:17953–17967, 2022. 1
- [2] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018. 3, 2
- [3] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020. 6, 7, 3
- [4] R Dennis Cook. Detection of influential observation in linear regression. *Technometrics*, 19(1):15–18, 1977. 1
- [5] 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. 6
- [6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3
- [7] Stephanie Fu, Netanel Tamir, Shobhit Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic data. *arXiv preprint arXiv:2306.09344*, 2023. 5, 3
- [8] Shanghua Gao, Pan Zhou, Ming-Ming Cheng, and Shuicheng Yan. Masked diffusion transformer is a strong image synthesizer. *arXiv preprint arXiv:2303.14389*, 2023. 6
- [9] Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In *International Conference on Machine Learning*, pages 2242–2251. PMLR, 2019. 1
- [10] Amirata Ghorbani, Michael Kim, and James Zou. A distributional framework for data valuation. In *International Conference on Machine Learning*, pages 3535–3544. PMLR, 2020. 1
- [11] 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. 3
- [12] Kazuaki Hanawa, Sho Yokoi, Satoshi Hara, and Kentaro Inui. Evaluation of similarity-based explanations. *arXiv preprint arXiv:2006.04528*, 2020. 6
- [13] 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. 1
- [14] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International conference on learning representations*, 2016. 6
- [15] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 6, 2
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2
- [17] Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128, 2010. 5
- [18] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nezihe Merve Gürel, Bo Li, Ce Zhang, Costas J Spanos, and Dawn Song. Efficient task-specific data valuation for nearest neighbor algorithms. *arXiv preprint arXiv:1908.08619*, 2019. 1
- [19] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li, Ce Zhang, Dawn Song, and Costas J Spanos. Towards efficient data valuation based on the shapley value. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1167–1176. PMLR, 2019. 1
- [20] Hoang Anh Just, Feiyang Kang, Jiachen T Wang, Yi Zeng, Myeongseob Ko, Ming Jin, and Ruoxi Jia. Lava: Data valuation without pre-specified learning algorithms. *arXiv preprint arXiv:2305.00054*, 2023. 1, 2, 3
- [21] 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. 6, 3
- [22] Zhifeng Kong and Kamalika Chaudhuri. Understanding instance-based interpretability of variational auto-encoders. *Advances in Neural Information Processing Systems*, 34:2400–2412, 2021. 2, 6, 7, 1
- [23] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 6, 2, 3
- [24] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 6, 2, 3
- [25] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba) dataset. *Retrieved August, 15(2018):11*, 2018. 6, 3
- [26] Ki Nohyun, Hoyong Choi, and Hye Won Chung. Data valuation without training of a model. In *The Eleventh International Conference on Learning Representations*, 2022. 1
- [27] Jian Pei. A survey on data pricing: from economics to data science. *IEEE Transactions on knowledge and Data Engineering*, 34(10):4586–4608, 2020. 1
- [28] Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*, 33:19920–19930, 2020. 1
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5

- 686 [30] Protection Regulation. General data protection regulation.
687 *Intouch*, 25, 2018. 1
- 688 [31] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wier-
689 stra. Stochastic backpropagation and approximate inference
690 in deep generative models. In *International conference on
691 machine learning*, pages 1278–1286. PMLR, 2014. 2
- 692 [32] Adam Richardson, Aris Filos-Ratsikas, and Boi Faltings.
693 Rewarding high-quality data via influence functions. *arXiv
694 preprint arXiv:1908.11598*, 2019. 1
- 695 [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
696 Patrick Esser, and Björn Ommer. High-resolution image
697 synthesis with latent diffusion models. In *Proceedings of
698 the IEEE/CVF Conference on Computer Vision and Pattern
699 Recognition*, pages 10684–10695, 2022. 2
- 700 [34] Nikunj Saunshi, Arushi Gupta, Mark Braverman, and San-
701 jeev Arora. Understanding influence functions and datamod-
702 els via harmonic analysis. *arXiv preprint arXiv:2210.01072*,
703 2022. 1
- 704 [35] Naoyuki Terashita, Hiroki Ohashi, Yuichi Nonaka, and
705 Takashi Kanemaru. Influence estimation for generative ad-
706 versarial networks. *arXiv preprint arXiv:2101.08367*, 2021.
707 1, 2, 6, 7
- 708 [36] Laurens Van der Maaten and Geoffrey Hinton. Visualizing
709 data using t-sne. *Journal of machine learning research*, 9
710 (11), 2008. 3
- 711 [37] Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust
712 data valuation framework for machine learning. In *Inter-
713 national Conference on Artificial Intelligence and Statistics*,
714 pages 6388–6421. PMLR, 2023. 1
- 715 [38] Zhendong Wang, Huangjie Zheng, Pengcheng He, Weizhu
716 Chen, and Mingyuan Zhou. Diffusion-gan: Training gans
717 with diffusion. *arXiv preprint arXiv:2206.02262*, 2022. 6, 7,
718 3
- 719 [39] Zhaoxuan Wu, Yao Shu, and Bryan Kian Hsiang Low.
720 Davinz: Data valuation using deep neural networks at initial-
721 ization. In *International Conference on Machine Learning*,
722 pages 24150–24176. PMLR, 2022. 1, 2
- 723 [40] Xinyi Xu, Zhaoxuan Wu, Chuan Sheng Foo, and Bryan
724 Kian Hsiang Low. Validation free and replication robust
725 volume-based data valuation. *Advances in Neural Infor-
726 mation Processing Systems*, 34:10837–10848, 2021. 1, 2
- 727 [41] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan
728 Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang.
729 Maniq: Multi-dimension attention network for no-reference
730 image quality assessment. In *Proceedings of the IEEE/CVF
731 Conference on Computer Vision and Pattern Recognition*,
732 pages 1191–1200, 2022. 5
- 733 [42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shecht-
734 man, and Oliver Wang. The unreasonable effectiveness of
735 deep features as a perceptual metric. In *Proceedings of the
736 IEEE conference on computer vision and pattern recogni-
737 tion*, pages 586–595, 2018. 5, 3

GMValuator: Similarity-based Data Valuation for Generative Models

Supplementary Material

In this appendix, we provide a concise summary of key notations in Sec. A and offer a more comprehensive review of related work in Sec. B. Additionally, Sec. C contains detailed statistical analysis for the results depicted in Figure 1. Our main text focuses on experiments conducted on benchmark datasets such as MNIST and CIFAR-10, along with a large-scale dataset, ImageNet, to evaluate the most significant contributors found by GMVALUATOR are in the same class as the generated sample (**C1**). In Sec. D, we extend the evaluation of GMVALUATOR to other generative models for **C1**. Additionally, we expand our evaluation using **C2** by considering the ground truth of multiple combined attributes. We operate under the assumption that the most significant contributors to a generated sample should exhibit similar attributes in the images. For a more thorough validation of the effectiveness of GMVALUATOR, high-resolution datasets like AFHQ and FFHQ are also employed for **C2**. The appendix includes additional experimental insights such as ablation studies, alternative distance metrics, and further experimental details in Sec. E, F, and H, respectively. Lastly, we ensure the reproducibility of our experiments and present the proof of our theorem.

A. Notation

To make motivation and problem formulation clear, we list some significant notations from Sec. 2 in Table 5.

Table 5. Some important notations in Sec. 2

Notations	Description
G	generative model
G^*	well-trained generative model
\mathcal{Z}	noise distribution
z	a latent sample in the latent space
$d(\cdot, \cdot)$	distance function
X, x	training dataset, training sample
\hat{X}, \hat{x}	generated dataset, generated sample
S	a subset of the training data
S^*	the real K contributors
K	$K = S^* $
T	K contributors found by a data valuator
\mathcal{A}	attribute space
\mathcal{X}	data distribution
\mathcal{L}	loss function
f, f'_{S^*}	labeling function, model trained on S^*

771

B. Related Work

B.1. Data Valuation

There are three lines of methods on data valuation: *metric-based methods*, *influence-based methods* and *data-driven methods*. In terms of *metric-based methods*, the commonly-used approach is to calculate its *marginal contribution* (MC) based on performance metrics (e.g. accuracy, loss). As the basic method depending on performance metrics for data valuation, LOO (Leave-One-Out) [4] is used to evaluate the value of the training sample by observational change of model performance when leaving out that data point from the training dataset. To overcome inaccuracy and strict desirability of LOO, SV [9] and BI [39] originated from *Cooperative Game Theory* are widely used to measure the contribution of data [10, 19]. Considering the joining sequence of each training data point, SV needs to calculate the marginal performance of all possible subsets in which the time complexity is exponential. Despite the introduction of techniques such as Monte-Carlo and gradient-based methods, as well as others proposed in the literature, approximating data significance value (SV) is computationally expensive and it typically requires retraining [9, 18]. The computational cost and need for unconventional performance metrics present difficulties in adapting the methods to generative models. As for *influence-based methods*, they evaluate the influence of data points on model parameters by computing the inverse Hessian for data valuation [18, 34, 36]. Due to the high computational cost, some approximation methods have also been proposed [30]. In addition, the use of influence function for data valuation is not limited to discriminative models, but can also be applied to specific generative models such as GAN and VAE [23, 37]. When it comes to data-driven methods, most of them are training-free methods that focus on the data itself [20, 41, 42].

B.2. Generative Model

Generative models are a type of unsupervised learning that can learn data distributions. Recently, there has been significant interest in combining generative models with neural networks to create *Deep Generative Models*, which are particularly useful for complex, high-dimensional data distributions. They can approximate the likelihood of each observation and generate new synthetic data by incorporating variations. Variational auto-encoders (VAEs) [33] optimize the log-likelihood of data by maximizing the evidence lower bound (ELBO), while generative adversarial networks (GANs) [11, 22] involves a generator and discriminator that compete with each other, resulting in

819 strong image generation. Recently proposed diffusion models [16, 35] add Gaussian noise to training data and learn to
 820 recover the original data. These models use variational inference and have a fixed procedure with a high-dimensional
 821 latent space.
 822

824 C. Statistical Results for Figure 1

825 It is evident by visualization in Figure 1 that the data points
 826 in X_{v2} (used for training) are more overlapped with generated data than data points in X_{v1} (not used for training).
 827 We perform statistic testing on data values obtained by GM-
 828 VALUATOR, to examine if data points X_{v2} (used for training)
 829 have significantly higher values than those of the data
 830 points in X_{v1} (not used for training).
 831

Table 6. The statistic test of data values of X_{v1} versus X_{v2} using different generative models. X_{v2} is supposed to have higher value than X_{v1} , given the generated data.

$H_0: \phi(D_i, S, \mu_i) \geq \phi(D_j, S, \mu_i)$
$H_1: \phi(X_i, S, \mu_i) < \phi(X_j, S, \mu_i), i \in X_{v1}, j \in X_{v2}$
BigGAN
Average value (v1) 0.319654
Average value (v2) 1.632352
P-value 6.937027×10^{-68}
T-statistic 17.924512
Significance level 0.01
Result p-value less than 0.01, reject H_0 , value of v2 less than v1 averagely
Classifier-free Guidance Diffusion
Average value (v1) 0.030434
Average value (v2) 0.369565
P-value 8.053195×10^{-55}
T-statistic 15.947860
Significance level 0.01
Result p-value less than 0.01, reject H_0

832 To this end, we use a t-test with the null hypothesis that
 833 data values in X_{v1} should not be smaller than those of X_{v2} .
 834 We compute p -value, which is the probability of getting a
 835 difference as large as we observed, or larger, under the null
 836 hypothesis. If the p -value is very low, we reject the null hy-
 837 pothesis and consider our approach, GMVALUATOR, to be
 838 verified with a high level of confidence ($1-p$). Typically, a p -
 839 value smaller than significance level 0.01 is used as a thresh-
 840 old for rejecting the null hypothesis. Table 6 showcases the
 841 outcomes of X_{v1} and X_{v2} in CIFAR-10 with $p \ll 0.01$ for
 842 both BigGAN and diffusion model, indicating that the data
 843 points in X_{v2} have significantly more value than those in
 844 X_{v1} . Consequently, these findings align with the presumption
 845 that the trained dataset X_{v2} has a higher value than the
 846 untrained dataset X_{v1} and verify our approach.

847 D. Additional Results on C1 and C2

Table 7. Performance comparison of Identical Class Test.

MNIST			
GAN (%)	$k=30$	$k=50$	$k=100$
GMValuator (No-Rerank)	96.27	96.26	95.86
GMValuator (l_2 -distance)	97.73	97.58	96.03
GMValuator (LPIPS)	97.77	97.72	97.38
GMValuator (DreamSim)	97.43	97.44	97.40
Diffusion (%)	$k=30$	$k=50$	$k=100$
GMValuator (No-Rerank)	92.40	91.82	91.26
GMValuator (l_2 -distance)	92.90	92.66	91.88
GMValuator (LPIPS)	93.73	97.72	92.42
GMValuator (DreamSim)	93.90	93.44	92.55
CIFAR-10			
BigGAN (%)	$k=30$	$k=50$	$k=100$
GMValuator (No-Rerank)	64.70	63.80	62.14
GMValuator (l_2 -distance)	64.70	63.80	62.14
GMValuator (LPIPS)	63.67	62.80	61.51
GMValuator (DreamSim)	70.33	68.74	65.18
Class-free Guidance Diffusion (%)	$k=30$	$k=50$	$k=100$
GMValuator (No-Rerank)	72.67	72.00	71.00
GMValuator (l_2 -distance)	72.67	72.00	71.00
GMValuator (LPIPS)	72.53	72.28	71.06
GMValuator (DreamSim)	79.37	78.08	74.61

848 D.1. (C1) Identical Class Test on Other Generative 849 Models

We have presented Identical Class Test (C1) on β -VAE and
 850 MNIST [25], CIFAR-10 [24] in Sec. 4.3 in our main
 851 context. Since GMVALUATOR is model-agnostic, we further
 852 validate our method of C1 on other generative models.
 853

Here, we conduct the experiments using a GAN and a
 854 Diffusion Model on MNIST. The architectural details of the
 855 used generative models are described in Sec. G.3 in the
 856 appendix. We also conduct the experiment on BigGAN [2]
 857 and Class-free Guidance Diffusion [15] with CIFAR-10.
 858 We used the same number of generated samples $m = 100$
 859 as the experiments presented in Sec. 4.
 860

Following the similar settings in Sec. 4.3 (C1), we ex-
 861 amine the class(es) of top k contributors for a given gen-
 862 erated data in the training data. We calculate the number of
 863 training samples, denoted as Q , from the top k contributors
 864 that have the same class as the generated data. The identi-
 865 cal class ratio, denoted as ρ , is calculated as $\rho = Q/k$. We re-
 866 port the average value of ρ across the generated datasets for
 867 different choices of k in Table 7. GMVALUATOR (Dream-
 868 Sim) has the highest ratio of contributors that belong to the
 869 same class as the generated sample among most of the mod-
 870 els evaluated on MNIST and CIFAR-10 datasets for differ-
 871 ent values of k . And the ratio improves as the value of k
 872 decreases, which is consistent with the top k assumption
 873 and validates our method.
 874

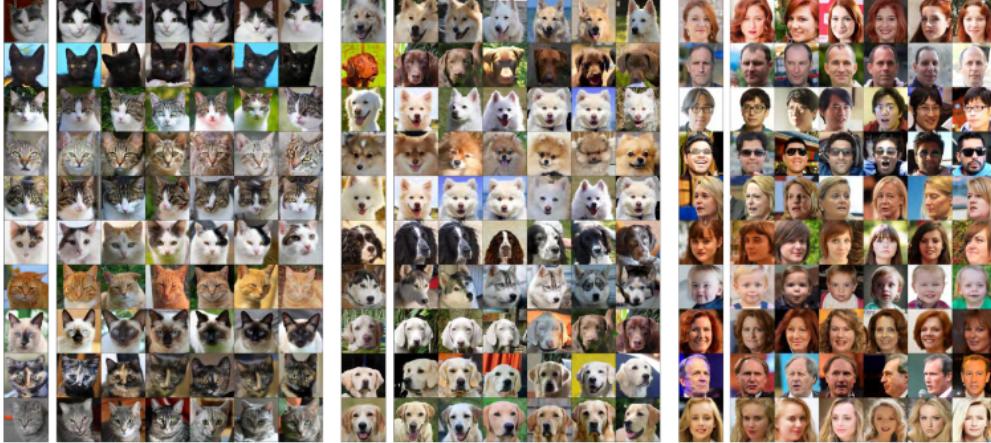


Figure 6. Visualization of Identical Attributes Test on AFHQ and FFHQ. The results shown in the first and second subfigures on the left are conducted on AFHQ-Cat and AFHQ-Dog, respectively. The subfigure on the right presents the results of FFHQ. In each subfigure, the generated samples are on the left, and the top k contributors in the training dataset are on the right.

875 D.2. (C2) Identical Attributes Test on CelebA

876 We extend C1 to focus on the attributes present in the images, treating them as ground truth rather than class labels, 877 as discussed in Sec. 4.4. Our experiments now incorporate 878 multiple attributes simultaneously, rather than just a single 879 attribute, to verify the performance of GMVALUATOR. For 880 instance, when evaluating a generated image with both a hat 881 and eyeglasses, the most significant contributors identified 882 should also include both of these attributes. Table 8 shows 883 some outcomes of identical attributes test when regard- 884 ing multiple combined attributes as ground truth, which 885 validate the effectiveness of our methods. 886

Table 8. Performance of Identical Attributes Test (C2) of multiple combined attributes.

Top K contributors:	$k=5$	$k=10$	$k=15$
Attribute: Eyeglasses & Gender (%)			
GMValuator (No-Rerank)	50.10	48.48	48.42
GMValuator (l_2 -distance)	59.19	56.67	56.23
GMValuator (LPIPS)	78.18	74.24	73.87
GMValuator (DreamSim)	92.53	90.61	90.30
Attribute: Eyeglasses & Hat (%)			
GMValuator (No-Rerank)	78.79	78.38	85.86
GMValuator (l_2 -distance)	84.44	82.93	83.23
GMValuator (LPIPS)	86.87	86.16	86.33
GMValuator (DreamSim)	89.49	87.58	87.41
Attribute: Gender & Hat (%)			
GMValuator (No-Rerank)	58.59	57.47	57.71
GMValuator (l_2 -distance)	61.62	60.51	60.40
GMValuator (LPIPS)	63.84	63.23	62.96
GMValuator (DreamSim)	65.25	64.44	64.18

887 D.3. (C2) Identical Attributes Test on Other 888 Datasets

To further validate GMVALUATOR, we also conducted 889 experiments on high-resolution datasets: AFHQ [3] and 890 FFHQ [21], following the same settings in Sec 4.4. The 891 results are shown in Figure 6, which demonstrates the 892 effectiveness of our methods in C2. The results show that the 893 top k contributors have similar attributes with the generated 894 sample such as fur color of cats or dogs in AFHQ. For the 895 experiment conducted on FFHQ, human faces attributes of 896 the most significant contributors are also similar to the 897 attributes in generated images. 898

899 E. Different Generated Data Sizes

Since our value function ϕ_i (Eq. (5)) for training data x_i 900 is computed by averaging over generated samples, it is ex- 901 pected that the sensitivity of ϕ_i is connected to the size m of 902 the investigated generated sample. To explore the influence 903 of generated data size m on the utilization of GMVALUA- 904 TOR, we perform sensitivity testing on MNIST, CIFAR10 905 using generative models GAN [11], and Diffusion mod- 906 els [6] as depicted below. The dataset, model and used k are 907 denoted under each subfigure of Figure 7. First, we generate 908 a varying number of samples from the same class. Specif- 909 ically, we consider four different sample sizes, denoted by 910 m , which are given by 1, 10, 30, and 50. Next, we evaluate 911 the GMVALUATOR using parameter C1, and this evalua- 912 tion is performed for each of the aforementioned values of 913 m . Subsequently, we conduct the experiment 10 times us- 914 ing GMVALUATOR (No-Rerank), each time with different 915 m -sized generated data samples from the same class. The 916

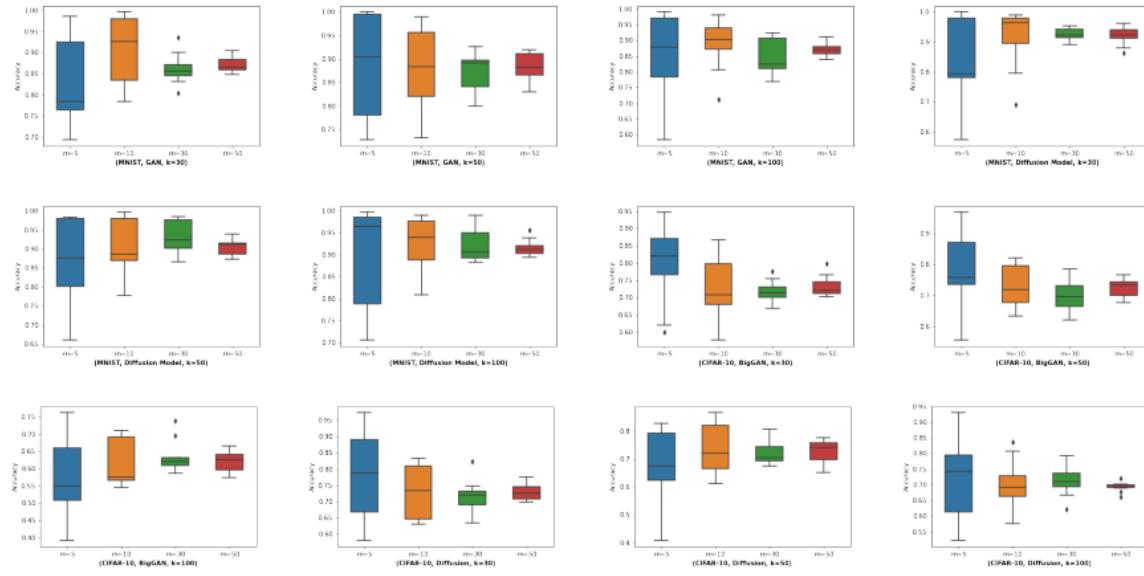


Figure 7. The change of ρ with the different number of generated samples m on MNIST and CIFAR-10 by diverse generative models.

917 results are presented as the mean and standard deviation (ρ)
 918 for accuracy, taken over these 10 runs. The results shown in
 919 Figure 7 imply that varying m does not yield notable differ-
 920 ences in mean accuracy and increasing the number of
 921 generated samples m leads to more stable and consistent
 922 results.

923 F. Alternative Distance Metric

924 We suggest utilizing Learned Perceptual Image Patch Sim-
 925 ilarity (LPIPS) [44] or DreamSim [7] as the distance metric
 926 d during the re-ranking phase. This enables to understand
 927 the perceptual dissimilarity between generated and real data
 928 points. These metrics are applied in the image embedding
 929 space derived from pre-trained models. Alternatively,
 930 distance measurement can be based on pixel space, such as
 931 employing l_2 -distance. Notably, we find that distances cal-
 932 culated in the input pixel space yield comparable outcomes
 933 as shown in Table 7. This implies that our method GMVAL-
 934 UATOR could be flexible to the different choices of distance
 935 metrics and the selection can depend on data prior.

936 G. Additional Experimental Details

937 G.1. Justification of Experiment Setup

938 We provide a detailed justification in Table 9 for our exper-
 939 iment setup from C1 to C4.

- 940 • C1. In Identical Class Test, the baseline method is VAE-
 941 TracIn [23], which can find the most influenced instances

942 in the training dataset. Since VAE-TracIn is the model-
 943 specific method, we only need to compare it with GM-
 944 VALUATOR when VAE model is used. Besides, all the
 945 datasets used in C1 should have the class labels. Con-
 946 sidering the computational demands detailed in VAE-
 947 TracIn [23], the runtime complexity of VAE-TracIn corre-
 948 lates with the number of network parameters and the size
 949 of the dataset. Therefore, our analysis primarily utilizes
 950 simpler benchmark datasets such as MNIST and CIFAR-
 951 10 for comparative evaluations with VAE-TracIn.

- 952 • C2. In extending C1, we use image attributes to supplant
 953 the concept of class for datasets lacking class labels. For
 954 datasets with attribute labels, quantified experiments are
 955 feasible, as demonstrated in Table 3 and Table 8. For
 956 datasets without attribute labels, we employ visualized
 957 experiments.

- 958 • C3. IF4GAN, a model-specific method for GAN, serves
 959 as the baseline for measuring data value. Adhering to the
 960 settings outlined in [37] and considering the impractical
 961 computational costs, we conduct experiments on the same
 962 dataset (MNIST) in [37] using DCGAN.

- 963 • C4. We present a comparison of the efficiency of GM-
 964 VALUATOR against baseline methods. In accordance with
 965 the settings used for VAE-TracIN and IF4GAN in [23],
 966 we report the efficiency results for C1 on the MNIST and
 967 CIFAR-10 datasets, and for C3 on MNIST.

Table 9. Justification of Experiment Setup. The selection of datasets and models was based on three critical factors that guided the process. Firstly, attribute labels were required to evaluate C2 effectively. Secondly, benchmark datasets were meticulously chosen to ensure a fair comparison with baselines while also taking into account computational costs (C3 and C4). Finally, the selected generative models are powerful enough to generate good-quality data for the datasets.

	Dataset	Model	Baseline	Label Requirements
C1	MNIST, CIFAR-10	VAE	VAE-TracIn [23]	Class labels
		Diffusion, GAN	-	Class labels
	ImageNet	Masked Diffusion Transformer [8]	-	Class labels
C2	CelebA	Diffusion-StyleGAN	-	Attribute labels
	AFHQ, FFHQ	StyleGAN	-	-
C3	MNIST	DCGAN	IF4GAN [23]	Class labels
C4	MNIST, CIFAR-10	VAE	VAE-TracIn [23]	Class labels
	MNIST	DCGAN	IF4GAN [23]	Class labels

968 G.2. Datasets

969 We conduct the generation tasks in the experiments on
 970 benchmark datasets (*i.e.* MNIST [25] and CIFAR [24]), face
 971 recognition dataset (*i.e.* CelebA [26]), high-resolution im-
 972 age dataset AFHQ [3] and FFHQ [21], large-scale image
 973 dataset ImageNet [5].

974 **MNIST.** The MNIST dataset consists of a collection of
 975 grayscale images of handwritten digits (0-9) with a reso-
 976 lution of 28x28 pixels. The dataset contains 60,000 training
 977 images and 10,000 testing images.

978 **CIFAR-10.** CIFAR-10 dataset consists of 60,000 color im-
 979 ages in 10 different classes, with 6,000 images per class.
 980 The classes include objects such as airplanes, cars, birds,
 981 cats, deer, dogs, frogs, horses, ships, and trucks. Each im-
 982 age in the CIFAR-10 dataset has a resolution of 32x32 pix-
 983 els.

984 **CelebA.** The CelebA dataset is a widely used face recog-
 985 nition and attribute analysis dataset, which contains a large
 986 collection of celebrity images with various facial attributes
 987 and annotations. The dataset consists of more than 200,000
 988 celebrity images, with each image labeled with 40 binary
 989 attribute annotations such as gender, age, facial hair, and
 990 presence of eyeglasses.

991 **AFHQ.** The AFHQ dataset is a high-resolution image
 992 dataset that focuses on animal faces (*e.g.* dogs, cat), and
 993 it consists of high-resolution images with 512×512 pixels.

994 **FFHQ.** The FFHQ dataset is a high-resolution face dataset
 995 that contains high-quality images (1024×1024 pixels) of hu-
 996 man faces.

997 **ImageNet.** ImageNet is a large-scale image dataset, which
 998 contains over 14 million images and is categorized into
 999 more than 20,000 classes.

1000 G.3. Architecture of Generative models

1001 In our experiments, we leverage different generative models
 1002 in the class of GAN, VAE and diffusion models. We util-
 1003 ize β -VAE for both MNIST and CIFAR-10 datasets while
 1004 a simple GAN is conducted on MNIST. BigGAN and β -

Table 10. The architecture of GAN for MNIST.

Generator	
FC(100, 8192), BN(32), ReLU	
Conv2D(128, 64, 4, 2, 1), BN(64), ReLU	
Discriminator	
Conv2D(1, 128, 4, 2, 1), BN(128), LeakyReLU	
FC(8192, 1024), BN(1024), LeakyReLU	

Table 11. The architecture of β -VAE.

Input	$28 \times 28 \times 1$ (MNIST) & $32 \times 32 \times 3$ (CIFAR-10).
Encoder	Conv $32 \times 4 \times 4$ (stride 2), $32 \times 4 \times 4$ (stride 2), $64 \times 4 \times 4$ (stride 2), $64 \times 4 \times 4$ (stride 2), FC 256. ReLU activation.
Latents	32
Decoder	Deconv reverse of encoder. ReLu acitvation. Gaussian.

VAE are also conducted on CIFAR-10. We list the archi-
 1005 tecture details for these generative models from Table 10
 1006 to Table 12. StyleGAN is used for high-resolution datasets
 1007 AFHQ and FFHQ. CelebA uses Diffusion-StyleGAN [40],
 1008 for which we use the exact architecture in their open-
 1009 sourced code. In addition, Masked Diffusion Transformer,
 1010 as introduced by Gao et al. [8], is applied to the ImageNet.
 1011

H. Discussion on the Possible Applications

1012 The application of data valuation within generative models
 1013 offers a wide range of opportunities. A potential use case is
 1014 to quantify privacy risks associated with generative model
 1015 training using specific datasets, since the matching mecha-
 1016 nism GMVALUATOR can help re-identify the training sam-
 1017 ples given the generated data. By doing so, organizations
 1018 and individuals will be able to audit the usage of their data
 1019 more effectively and make informed decisions regarding its
 1020

Table 12. The architecture of β -VAE.

β -VAE	
Generator	Discriminator
$z \in \mathbb{R}^{120} \sim \mathcal{N}(0, I)$	RGB image $x \in \mathbb{R}^{32 \times 32 \times 3}$
Embed(y) $\in \mathbb{R}^{32}$	
Linear ($20 + 128$) $\rightarrow 4 \times 4 \times 16ch$	ResBlock down $ch \rightarrow 2ch$
ResBlock up $16ch \rightarrow 16ch$	Non-Local Block (64×64)
ResBlock up $16ch \rightarrow 8ch$	ResBlock down $2ch \rightarrow 4ch$
ResBlock up $8ch \rightarrow 4ch$	ResBlock down $4ch \rightarrow 8ch$
ResBlock up $4ch \rightarrow 2ch$	ResBlock down $8ch \rightarrow 16ch$
Non-Local Block (16×16)	ResBlock down $16ch \rightarrow 16ch$
ResBlock up $2ch \rightarrow ch$	ResBlock $16ch \rightarrow 16ch$
BN, ReLU, 3×3 Conv $ch \rightarrow 3$	ReLU, Global sum pooling
Tanh	Embed(y) $\cdot h + (\text{linear} \rightarrow 1)$

1021 use.

1022 Another promising application is material pricing and
1023 finding in content creation. For example, when training generative
1024 models for various purposes, such as content recommendation or personalized advertising, data evaluation can
1025 be used to measure the value of reference content.1026 In addition, GMVALUATOR can play an important role
1027 in the development of ensuring the responsibility of using
1028 synthetic data in safe-sensitive fields, such as healthcare or
1029 finance. By assessing the value of the data used in generative
1030 model training, researchers can ensure that the generated data are robust and reliable.1031 Last but not least, the applications of GMVALUATOR
1032 can promote the recognition of intellectual property rights.
1033 Determining the value of the intellectual property being
1034 generated by generative models is critical. By evaluating
1035 the data employed in training generative models, we can
1036 develop a more comprehensive understanding of copyright
1037 that may emerge from the generative models. In essence,
1038 such insights can help advance licensing agreements for the
1039 utilization of the generative model and its outputs.1040

I. Reproducibility

1041 To ensure reproducibility, we make our implementation
1042 available to reviewers through this anonymous link: <https://anonymous.4open.science/r/GMValuator-V2-164D>.

J. Omitted proofs

1046

We follow [20] to prove the theorem. Firstly, we give several assumptions that will be used in later proof.

1047

Assumption J.1 Following Assumption 2.3, given a distance function $d(\cdot, \cdot)$ between $\mathcal{X}_{(T|f)}$ and $\mathcal{X}_{(S^*|f)}$ we defined the coupling between $\mathcal{X}_{(T|f)}$ and $\mathcal{X}_{(S^*|f)}$ as π^* :

1048

1049

$$\pi^* := \arg \inf_{\pi \in \Pi(\mathcal{X}_{(T|f)}, \mathcal{X}_{(S^*|f)})} \mathbb{E}_{(x_T, x_{S^*}) \sim \pi} d(x_T, x_{S^*}) \quad (7) \quad 1050$$

It is easy to see that all joint distributions defined above are couplings between the corresponding distribution pairs. Then, following [20] we prove the main Theorem.

1051

1052

Theorem J.2 (Restated of Theorem 2.4.) Let $f'_{S^*} : \mu \rightarrow \mathcal{A} = \{0, 1\}^V$ be the model trained on the optimal contributor dataset S^* . Following Assumption 2.3, if the contributors are corresponding to the given generated data \hat{X} , we have:

1053

1054

$$\mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] - \mathbb{E}_{x \sim \mu_S} [\mathcal{L}(f(x), f'_{S^*}(x))] \leq k\epsilon \cdot [d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)}) + d_W(\mathcal{X}_{(S^*|f)}, \mathcal{X}_{(\hat{X}|f)})] \quad (8) \quad 1055$$

Proof J.3

$$\mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] = \mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] - \mathbb{E}_{x \sim \mu_S} [\mathcal{L}(f(x), f'_{S^*}(x))] + \mathbb{E}_{x \sim \mu_S} [\mathcal{L}(f(x), f'_{S^*}(x))] \quad 1056$$

$$\leq \mathbb{E}_{x \sim \mu_S} [\mathcal{L}(f(x), f'_{S^*}(x))] + \left| \mathbb{E}_{x \sim \mu_S} [\mathcal{L}(f(x), f'_{S^*}(x))] - \mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] \right| \quad 1057$$

(9)

We bound $\left| \mathbb{E}_{x \sim \mu_{S^*}} [\mathcal{L}(f(x), f'_{S^*}(x))] - \mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] \right|$ as follows:

1058

$$\left| \mathbb{E}_{x \sim \mu_{S^*}} [\mathcal{L}(f(x), f'_{S^*}(x))] - \mathbb{E}_{x \sim \mu_T} [\mathcal{L}(f(x), f'_{S^*}(x))] \right| \quad 1059$$

$$= \left| \int_{\mathcal{X}^2} \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_S)) - \mathcal{L}(f(x_T), f'_{S^*}(x_T)) d\pi^*(x_T, x_{S^*}) \right| \quad 1060$$

$$= \left| \int_{\mathcal{X}^2} \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_{S^*})) - \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_T)) + \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_T)) - \mathcal{L}(f(x_T), f'_{S^*}(x_T)) d\pi^*(x_T, x_{S^*}) \right| \quad 1061$$

$$\leq \int_{\mathcal{X}^2} \left| \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_{S^*})) - \int_{\mathcal{X}^2} \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_T)) \right| d\pi^*(x_T, x_{S^*}) \quad 1062$$

$$+ \int_{\mathcal{X}^2} \left| \mathcal{L}(f(x_{S^*}), f'_{S^*}(x_T)) - \int_{\mathcal{X}^2} \mathcal{L}(f(x_T), f'_{S^*}(x_T)) \right| d\pi^*(x_T, x_{S^*}) \quad (10) \quad 1063$$

Then due to k -Lipschitzness of \mathcal{L} and ϵ -Lipschitzness of f , we can obtain:

$$\text{RHS of Eq.(10)} \leq k \int_{\mathcal{X}^2} \|f'_{S^*}(x_{S^*}) - f'_{S^*}(x_T)\| d\pi^*(x_T, x_{S^*}) + k \int_{\mathcal{X}^2} \|f(x_{S^*}) - f(x_T)\| d\pi^*(x_T, x_{S^*}) \quad 1065$$

$$\leq k\epsilon \int_{\mathcal{X}^2} 2d(x_T, x_{S^*}) d\pi^*(x_T, x_{S^*}) \quad 1066$$

$$= k\epsilon d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(S^*|f)}), \quad 1067$$

where the last step is due to the definition of 1-Wasserstein distance. Then, according to the triangle inequality of Wasserstein distance [29], we can obtain:

1068

1069

$$d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(S^*|f)}) \leq d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)}) + d_W(\mathcal{X}_{(\hat{X}|f)}, \mathcal{X}_{(S^*|f)}) \quad (11) \quad 1070$$

Combining Eq. (9) and Eq. (11) we finished the proof and obtained the Theorem 2.4. By reducing the distance term $d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)})$, we have $\mathcal{X}_{(T|f)} \rightarrow \mathcal{X}_{(S^*|f)}$. As a result, the expected distance

$$\mathbb{E}_{(S^* \sim \mathcal{X}_{(S^*|f)}, T \sim \mathcal{X}_{(T|f)})} \min_{\pi \in \Pi(T, S^*)} \mathbb{E}_{(x_T, x_{S^*}) \sim \pi} d(x_T, x_{S^*}) \rightarrow 0,$$

with randomly sampling S^* and T with K elements.

1071

16558_gmvaluator_similarity_based_da.pdf

ORIGINALITY REPORT

4%

SIMILARITY INDEX

MATCH ALL SOURCES (ONLY SELECTED SOURCE PRINTED)

★workshop.pdf

Comparison Document

4%

EXCLUDE QUOTES ON

EXCLUDE SOURCES OFF

EXCLUDE BIBLIOGRAPHY ON

EXCLUDE MATCHES OFF