C	NTENTS						
1	Introduction						
-							
2	Motivation and Pr	roblem Formulation					
	2.1 Motivations						
	2.1.1 Theor	retical Justification					
	2.1.2 Empe	erical Validation					
	2.2 Problem Form	nulation and Challenges					
3	The Proposed Dat	a Valuation Methods					
-	3.1 Efficient Simi	ilarity Matching					
	3.2 Image Quality						
	3.2 Image Quant	y Assessment					
	5.5 Value Calcula						
4	Experiments						
	4.1 Experiment S	Setup					
	4.2 Metrics						
	4.3 Identical Clas	ss Test (C1)					
	4.4 Identical Attr	ibutes Test (C2)					
	4.5 Out of Distrib	pution Detection (C3)					
	4.6 Efficiency (C	4)					
5	Conclusion	Conclusion 10					
Aŗ	pendix						
A	Notation						
B	Algorithm						
С	Related Work						
	C.1 Data Valuatio	m					
	C.2 Generative M	lodel					
D	Statistical Results	for Figure 1					
Ε	Additional Results	s on C1 and C2					
	E.1 (C1) Identical	l Class Test on Other Generative Models					
	E.2 (C2) Identical	l Attributes Test on CelebA					
	E.3 (C2) Identical	Attributes Test on Other Datasets					
F	Different Generat	ed Data Sizes					

756	G	Alternative Embedding Approaches		
758				
759	Н	Alte	rnative Distance Metric	21
760				
761	Ι	Nec	essity for calibration	21
762			•	
763	J	Add	litional Experimental Details	21
764 765		J.1	Justification of Experiment Setup	21
766		J.2	Datasets	22
767		13	Architecture of Generative models	23
768		J .J		23
769	17	D!-	and an an the Description Annulise theme	
770	ĸ	Disc	sussion on the Possible Applications	23
771	-	~		
772	L	Cur	rent Limitation and Future Directions	24
774				
775	Μ	Om	itted proofs	24
776				
777	Ν	Rep	roducibility	25
778				

Notations	Notationa Description		
INOLALIOIIS	Description		
μ	the distribution of the subset of the training data		
μ_T	a subset of the contributors		
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		
${\cal L}$	loss function		
f, f'_{C*}	labeling function, model trained on S^*		

Roadmap of Appendix In this appendix, we provide a concise summary of key notations and Algorithm 1, detailing the pipeline of GMValuator in Sec. A and Sec. B. A more comprehensive review of related work is also provided in Sec. C. Additionally, Sec. D 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. E, 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 embedding approaches, alternative distance metrics, calibration and further experimental

details in Sec. F, H, I and J, respectively. The potential applications, limitations and future directions are also provided in Sec. K and Sec. L. Lastly, we present the proof of our theorem and ensure the reproducibility of our experiments.

A NOTATION

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

817 818 819

820 821

822

844 845

846 847

848

814

815

816

B ALGORITHM

To better introduce GMVALUTOR, the pseudocode is shown in Algorithm 1.

Algorithm 1 GMVALUATOR 823 824 **Input**: Training dataset $X = \{x_i\}_{i=1}^n$, a well-trained 825 model G^* , random distribution \mathcal{Z} . **Output**: Generated dataset \hat{X} , the value of training data 826 points $\Phi = \{\phi_1, \phi_2, ..., \phi_n\}$ 827 828 1: $\hat{X} = {\hat{x}_j}_{j=1}^m \leftarrow G^*(z_j)$, for $z_j \in \mathbb{Z}$ // Generate the synthetic dataset 829 2: for \hat{x}_i in \hat{X} do 830 3: // Matching process (see Sec. 3.1) 831 4: $\mathcal{P}_j = f(X, \hat{x}_j)$ // Including two phases 832 5: for x_i in \mathcal{P}_j do $d_{ij} \leftarrow \text{DreamSim}(x_i, \hat{x}_j) \text{ or others}$ 833 6: 7: end for 834 $q_i = MANIQA(\hat{x}_i)$ // Image Quality Assessment (see Sec. 3.2) 8: 835 Calculate score $\mathcal{V}(x_i, \hat{x}_j, d_{ij}, q_j)$ using Eq. equation 6 // Contribution Score Calculation (see 9: 836 Sec. 3.3) 837 10: end for 838 11: // Calculation of data value and return the result Φ 839 12: for x_i in X do 840 Calculate x_i 's value ϕ_i using Eq.equation 5 13: 841 14: end for 842 15: return $\Phi = \{\phi_1, \phi_2, ..., \phi_n\}$ 843

C RELATED WORK

C.1 DATA VALUATION

849 There are three lines of methods on data valuation: *metric-based methods*, *influence-based methods* 850 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). 851 As the basic method depending on performance metrics for data valuation, LOO (Leave-One-852 Out) Cook (1977) is used to evaluate the value of the training sample by observational change 853 of model performance when leaving out that data point from the training dataset. To overcome 854 inaccuracy and strict desirability of LOO, SV Ghorbani & Zou (2019) and BI Wang & Jia (2023) 855 originated from *Cooperative Game Theory* are widely used to measure the contribution of data Jia 856 et al. (2019b); Ghorbani et al. (2020). 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 858 is exponential. Despite the introduction of techniques such as Monte-Carlo and gradient-based 859 methods, as well as others proposed in the literature, approximating data significance value (SV) 860 is computationally expensive and it typically requires retraining Ghorbani & Zou (2019); Jia et al. (2019a). The computational cost and need for unconventional performance metrics present difficulties 861 in adapting the methods to generative models. As for influence-based methods, they evaluate the 862 influence of data points on model parameters by computing the inverse Hessian for data valuation Jia 863 et al. (2019a); Richardson et al. (2019); Saunshi et al. (2022). Due to the high computational cost,

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

$H_0: \phi(D_i, S, \mu_i) \ge \phi(D_j, S, \mu_i)$	
$H_1: \phi(X_i, S, \mu_i) < \phi(X_i, S, \mu_i), i \in X_{v1}, j \in X_{v2}$	2

	BigGAN	Classifier-free Guidance Diffusion
Average value (v1)	0.319654	0.030434
Average value (v2)	1.632352	0.369565
P-value	6.937027×10^{-68}	8.053195×10^{-55}
T-statistic	17.924512	15.947860
Significance level	0.01	0.01
Result	p-value less than 0.01, reject H_0 , value of v2 less than v1 averagely	p-value less than 0.01, reject H_0 , value of v2 less than v1 averagely

some approximation methods have also been proposed Pruthi et al. (2020). 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 Terashita et al. (2021); Kong & Chaudhuri (2021). When it comes to data-driven methods, most of them are training-free methods that focus on the data itself Xu et al. (2021); Wu et al. (2022); Just et al. (2023).

C.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) Rezende et al. (2014) optimize the log-likelihood of data by maximizing the evidence lower bound (ELBO), while generative adversarial networks (GANs) Goodfellow et al. (2020); Karras et al. (2020) involves a generator and discriminator that compete with each other, resulting in strong image generation. Recently proposed diffusion models Ho et al. (2020); Rombach et al. (2022) add Gaussian noise to training data and learn to recover the original data. These models use variational inference and have a fixed procedure with a high-dimensional latent space.

D STATISTICAL RESULTS FOR FIGURE 1

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

To this end, we use a t-test with the null hypothesis that data values in X_{v1} should not be smaller than those of X_{v12} . We compute p-value, which is the probability of getting a difference as large as we observed, or larger, under the null hypothesis. If the *p*-value is very low, we reject the null hypothesis and consider our approach, GMVALUATOR, to be verified with a high level of confidence (1-p). Typically, a p-value smaller than significance level 0.01 is used as a threshold for rejecting the null hypothesis. Table 6 showcases the outcomes of X_{v1} and X_{v2} in CIFAR-10 with $p \ll 0.01$ for both BigGAN and diffusion model, indicating that the data points in X_{v2} have significantly more value than those in X_{v1} . Consequently, these findings align with the presumption that the trained dataset X_{v2} has a higher value than the untrained dataset X_{v1} and verify our approach.

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.

Table 7: Performance comparison of Identical Class Test.			
MNIST			
GAN (%)	<i>k</i> =30	<i>k</i> =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 (%)	<i>k</i> =30	<i>k</i> =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 (%)	<i>k</i> =30	<i>k</i> =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 (%)	<i>k</i> =30	<i>k</i> =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

E ADDITIONAL RESULTS ON C1 AND C2

E.1 (C1) IDENTICAL CLASS TEST ON OTHER GENERATIVE MODELS

We have presented Identical Class Test (C1) on β -VAE and MNIST LeCun et al. (1998), CIFAR-10 Krizhevsky et al. (2009) in Sec. 4.3 in our main context. Since GMVALUATOR is model-agnostic, we further validate our method of C1 on other generative models.

Here, we conduct the experiments using a GAN and a Diffusion Model on MNIST. The architectural
 details of the used generative models are described in Sec. J.3 in the appendix. We also conduct
 the experiment on BigGAN Brock et al. (2018) and Class-free Guidance Diffusion Ho & Salimans

(2022) with CIFAR-10. We used the same number of generated samples m = 100 as the experiments presented in Sec. 4.

Following the similar settings in Sec. 4.3 (C1), we examine the class(es) of is top k contributors for a given generated data in the training data. We calculate the number of training samples, denoted as Q, from the top k contributors that have the same class as the generated data. The identical class ratio, denoted as ρ , is calculated as $\rho = Q/k$. We report the average value of ρ across the generated datasets for different choices of k in Table 7. GMVALUATOR (DreamSim) has the highest ratio of contributors that belong to the same class as the generated sample among most of the models evaluated on MNIST and CIFAR-10 datasets for different values of k. And the ratio improves as the value of k decreases, which is consistent with the top k assumption and validates our method.

982 983 984

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

			-
Top K contributors:	<i>k</i> =5	k=10	<i>k</i> =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: Eyegla	sses & H	lat (%)	
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

1001 E.2 (C2) IDENTICAL ATTRIBUTES TEST ON CELEBA

We extend **C1** to focus on the attributes present in the images, treating them as ground truth rather than class labels, as discussed in Sec. 4.4. Our experiments now incorporate multiple attributes simultaneously, rather than just a single attribute, to verify the performance of GMVALUATOR. For instance, when evaluating a generated image with both a hat and eyeglasses, the most significant contributors identified should also include both of these attributes. Table 8 showcases some outcomes of identical attributes test when regarding multiple combined attributes as ground truth, which validate the effectiveness of our methods.

1009 1010 1011

E.3 (C2) IDENTICAL ATTRIBUTES TEST ON OTHER DATASETS

To further validate GMVALUATOR, we also conducted experiments on high-resolution datasets: AFHQ Choi et al. (2020) and FFHQ Karras et al. (2019), following the same settings in Sec 4.4. The results are shown in Figure 6, which demonstrates the effectiveness of our methods in C2. The results show that the top k contributors have similar attributes with the generated sample such as fur color of cats or dogs in AFHQ. For the experiment conducted on FFHQ, human faces attributes of the most significant contributors are also similar to the attributes in generated images.

1018

1020

¹⁰¹⁹ F DIFFERENT GENERATED DATA SIZES

1021 Since our value function ϕ_i (Eq. equation 5) for training data x_i is computed by averaging over 1022 generated samples, it is expected that the sensitivity of ϕ_i is connected to the size m of the in-1023 vestigated generated sample. To explore the influence of generated data size m on the utilization 1024 of GMVALUATOR, we perform sensitivity testing on MNIST, CIFAR10 using generative models 1025 GAN Goodfellow et al. (2020), and Diffusion models Dhariwal & Nichol (2021) as depicted below. 1026 The dataset, model and used k are denoted under each subfigure of Figure 7. First, we generate a



Figure 7: The change of ρ with the different number of generated samples m on MNIST and CIFAR-1056 10 by diverse generative models. 1057

1061

1062 1063

1065

1066

varying number of samples from the same class. Specifically, we consider four different sample sizes, denoted by m, which are given by 1, 10, 30, and 50. Next, we evaluate the GMVALUATOR using parameter C1, and this evaluation is performed for each of the aforementioned values of m. Subsequently, we conduct the experiment 10 times using GMVALUATOR (No-Rerank), each time 1064 with different m-sized generated data samples from the same class. The results are presented as the mean and standard deviation (ρ) for accuracy, taken over these 10 runs. The results shown in Figure 7 imply that varying m does not yield notable differences in mean accuracy and increasing the number of generated samples m leads to more stable and consistent results.

- 1067 1068
- 1069 1070
- 1071

G ALTERNATIVE EMBEDDING APPROACHES

1072 1073

1074 In the step of efficient similarity matching, the embedding f_e is exclusively used in the recall phase 1075 for retaining n samples with non-minimal contributions. Apart from using CLIP for embedding, we also investigate the influence for the results when using other embedding methods. We present results in Table 9 using more embedding methods, showing similar performance with CLIP with "Rerank". 1077 The reason for this results is that we implement re-ranking that relies on the image space rather than 1078 the embedding space, which allows us to derive accurate top k contributors from n samples $n \gg k$. 1079 Thus, the embedding could tolerate some noise in the recall phase.

	Table 9: Comparison of different embedding methods.				
-	MNIST (%)	No-Rerank	l_2 -distance	LPIPS	DreamSim
-	CLIP	86.41	87.76	88.78	88.78
	Alexnet	79.77	84.82	86.51	88.72
	Densenet	80.47	86.77	87.97	89.35

1088

1083

1080 1081 1082

H ALTERNATIVE DISTANCE METRIC

1089 We suggest utilizing Learned Perceptual Image Patch Similarity (LPIPS) Zhang et al. (2018) or 1090 DreamSim Fu et al. (2023) as the distance metric d during the re-ranking phase. This enables to 1091 understand the perceptual dissimilarity between generated and real data points. These metrics are 1092 applied in the image embedding space derived from pre-trained models. Alternatively, distance 1093 measurement can be based on pixel space, such as employing l^2 -distance. Notably, we find that 1094 distances calculated in the input pixel space yield comparable outcomes as shown in Table 7. This implies that our method GMVALUATOR could be flexible to the different choices of distance metrics 1095 and the selection can depend on data prior. 1096

I NECESSITY FOR CALIBRATION

1099 1100

1108

1121 1122

1127

1128 1129

1130

1131

1132

1133

1097

For challenges 2 and 3, we utilize image quality assessment and establish a non-zero scores rule for calibration in GMVALUATOR. To better understand the impact and necessity of calibration, we compare the distribution of the top 1, 2, and 3 contributors' scores with (w) and without (w/o) calibration conducted on MNIST in Figure 8. It is evident that scores without calibration are generally higher than with calibration. We also extend C3 and perform an ablation study in Figure 9 on scenarios where the generated samples include low-quality outputs. The results show that the OOD (out-of-distribution) training samples' value rank by GMValuator without (w/o) calibration is smaller, indicating significant bias and poor performance.



Figure 8: Sore Distribution

Figure 9: Value Rank

1123 J ADDITIONAL EXPERIMENTAL DETAILS

1125 1126 J.1 JUSTIFICATION OF EXPERIMENT SETUP

We provide a detailed justification in Table 10 for our experiment setup from C1 to C4.

• **C1.** In Identical Class Test, the baseline method is VAE-TracIn Kong & Chaudhuri (2021), which can find the most influenced instances in the training dataset. Since VAE-TracIn is the model-specific method, we only need to compare it with GMVALUATOR when VAE model is used. Besides, all the datasets used in **C1** should have the class labels. Considering the computational demands detailed in VAE-TracIn Kong & Chaudhuri (2021), the runtime complexity of VAE-TracIn correlates with the number of network parameters and the size of

	the dataset	. Therefore, our analysis primaril	v utilizes simpler benchmark d	atasets such as
	MNIST an	d CIFAR-10 for comparative eval	uations with VAE-TracIn.	
		1		
	• C2. In exte	ending C1, we use image attribute	s to supplant the concept of cla	ass for datasets
	lacking cla	ss labels. For datasets with attribute	ite labels, quantified experiment	nts are feasible,
	as demons	trated in Table 3 and Table 8. For	datasets without attribute lab	els, we employ
	visualized	experiments.		
	• C3. IF4G.	AN, a model-specific method for	GAN, serves as the baseline	for measuring
	data value.	Adhering to the settings outlined	in Terashita et al. (2021) and o	considering the
	impractica	a st al. (2021) using DCCAN	ct experiments on the same da	taset (MINIST)
	III Terasiin	a et al. (2021) using DCGAN.		
	• C4. We pre	esent a comparison of the efficienc	v of GMVALUATOR against ba	seline methods
	In accorda	nce with the settings used for VA	E-TracIN and IF4GAN in Kon	g & Chaudhuri
	(2021), we	report the efficiency results for C	1 on the MNIST and CIFAR-1	0 datasets, and
	for C3 on 1	MNIST.		
. .	-			
J.2	DATASETS			
xx 7	1 1		1 1 1 1	
We (conduct the gene	eration tasks in the experiments on $Z_{\rm right}$	benchmark datasets (<i>i.e.</i> , MNI	ST LeCun et al. (2018)
199 1194	o) and CIFAR I	Anznevsky et al. (2009)), face fect ge dataset AEHO Choi et al. (202	oginuon dataset (<i>i.e.</i> , CelebA Li	u et al. (2018)),
mag	re dataset Imag	eNet Deng et al. (2009)	(20) and FFIQ Kallas et al. (20)	19), large-scale
inaz				
<i>MN</i>	IST. The MNIS	T dataset consists of a collection of	f grayscale images of handwrit	tten digits (0-9)
with	a resolution of	28x28 pixels. The dataset contai	ns 60,000 training images and	10,000 testing
inag	ges.			
CIF	A R-10. CIFAR	-10 dataset consists of 60,000 co	or images in 10 different class	ses, with 6,000
imag hors	ges per class. T es, ships, and tr	he classes include objects such as ucks. Each image in the CIFAR-1	airplanes, cars, birds, cats, de 0 dataset has a resolution of 32	er, dogs, frogs, x32 pixels.
Colo	ba The Celeb	A dataset is a widely used face rec	contition and attribute analysis	dataset which
cont	ains a large col	lection of celebrity images with	various facial attributes and an	notations. The
data	set consists of	more than 200,000 celebrity ima	ges, with each image labeled	with 40 binary
uttri	bute annotations	s such as gender, age, facial hair, a	ind presence of eyeglasses.	5
FL	IO The ΛFHO	dataset is a high-resolution image	lataset that focuses on animal f	aces (a a doos
at),	and it consists	of high-resolution images with 51	2×512 pixels.	aces (e.g., uogs,
FFE	IQ. The FFHQ	dataset is a high-resolution fac	e dataset that contains high-	quality images
102	4x1024 pixels)	of human faces.	_	
ma	<i>geNet</i> . ImageN	et is a large-scale image dataset	which contains over 14 millior	images and is
cate	gorized into mo	re than 20,000 classes.		
Tabl	e 10: Justificat	ion of Experiment Setup. The se	lection of datasets and models	s was based on
three	e critical factors	s that guided the process. Firstly,	attribute labels were required	to evaluate C2
effe	ctively. Second	ly, benchmark datasets were met	iculously chosen to ensure a fa	air comparison
with	baselines while	e also taking into account computa	tional costs (C3 and C4). Final	lly, the selected
gene	erative models a	re powerful enough to generate go	ood-quality data for the datasets	S.
	MNIST CIEAD 10	VAE	VAE-TracIn Kong & Chaudhuri (2021)	Class labels
C1	ImageNet	Diffusion, GAN Masked Diffusion Transformer Gao et al. (2023)	-	Class labels
C2	CelebA	Diffusion-StyleGAN	-	Attribute labels
C3	AFHQ, FFHQ MNIST	StyleGAN DCGAN	- IF4GAN Terashita et al. (2021)	- Class labels
C4	MNIST, CIFAR-10	VAE	VAE-TracIn Kong & Chaudhuri (2021)	Class labels
	MNIST	DCGAN	IF4GAN Terashita et al. (2021)	Class labels

1188 J.3 ARCHITECTURE OF GENERATIVE MODELS

1190 In our experiments, we leverage different generative models in the class of GAN, VAE and diffusion 1191 models. We utilize β -VAE for both MNIST and CIFAR-10 datasets while a simple GAN is conducted 1192 on MNIST. BigGAN and β -VAE are also conducted on CIFAR-10. We list the architecture details 1193 for these generative models from Table 11 to Table 13. StyleGAN is used for high-resolution datasets 1194 AFHQ and FFHQ. CelebA uses Diffusion-StyleGAN Wang et al. (2022), for which we use the exact 1195 architecture in their open-sourced code. In addition, Masked Diffusion Transformer, as introduced by 1196 Gao et al. (2023), is applied to the ImageNet.

1198	Table	e 11: The architecture of GAN for MNIST.
1199		Generator
1200		FC(100, 8192), BN(32), ReLU
1201	(Conv2D(128, 64, 4, 2, 1), BN(64), ReLU
1202		Discriminator
1203	Con	v2D(1, 128, 4, 2, 1), BN(128), LeakyReLU
1204	Ī	FC(8192, 1024), BN(1024), LeakyReLU
1205		
1206		
1207	7	Table 12: The architecture of BigGAN.
1208	Innut	$28 \times 28 \times 1$ (MNIST) & $32 \times 32 \times 3$ (CIFAR-10)
1209	mput	$20 \times 20 \times 1$ (MR(51) & $32 \times 32 \times 3$ (CH7(R 10)).
1210		Conv $32 \times 4 \times 4$ (stride 2), $32 \times 4 \times 4$ (stride 2),
1211	Encoder	$64 \times 4 \times 4$ (stride 2), $64 \times 4 \times 4$ (stride 2),
1212		FC 256. ReLU activation.
1213	Latents	32
1214	Decel	Deconv reverse of encoder. ReLu acitvation.
1215	Decoder	Gaussian.
1216		

1217 1218

1219

1225 1226 1227

1107

K DISCUSSION ON THE POSSIBLE APPLICATIONS

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

1228				
1229	β-VAE			
1230	Generator	Discriminator		
1231	$z \in \mathbb{R}^{120} \sim \mathcal{N}(0, I)$			
1232	$\mathbb{E} = \mathbb{E} = $	RGB image $x \in \mathbb{R}^{32 \times 32 \times 3}$		
1233	Linear $(20 + 128) \rightarrow 4 \times 4 \times 16ch$	ResBlock down $ch \rightarrow 2ch$		
1234	ResBlock up $16ch \rightarrow 16ch$	Non-Local Block (64×64)		
1235	ResBlock up $16ch \rightarrow 8ch$	ResBlock down $2ch \rightarrow 4ch$		
1236	ResBlock up $8ch \rightarrow 4ch$	ResBlock down $4ch \rightarrow 8ch$		
1237	ResBlock up $4ch \rightarrow 2ch$	ResBlock down $8ch \rightarrow 16ch$		
1238	Non-Local Block (16×16)	ResBlock down $16ch \rightarrow 16ch$		
1239	ResBlock up $2ch \rightarrow ch$	ResBlock $16ch \rightarrow 16ch$		
1240	BN, ReLU, 3×3 Conv ch $\rightarrow 3$	ReLU, Global sum pooling		
1241	Tanh	Embed $(y) \cdot h + (\text{linear} \to 1)$		

Table 13: The architecture of β -VAE.

1242 Another promising application is material pricing and finding in content creation. For example, when 1243 training generative models for various purposes, such as content recommendation or personalized 1244 advertising, data evaluation can be used to measure the value of reference content.

1245 In addition, GMVALUATOR can play an important role in the development of ensuring the responsi-1246 bility of using synthetic data in safe-sensitive fields, such as healthcare or finance. By assessing the 1247 value of the data used in generative model training, researchers can ensure that the generated data are 1248 robust and reliable.

Last but not least, the applications of GMVALUATOR can promote the recognition of intellectual 1250 property rights. Determining the value of the intellectual property being generated by generative 1251 models is critical. By evaluating the data employed in training generative models, we can develop a 1252 more comprehensive understanding of copyright that may emerge from the generative models. In 1253 essence, such insights can help advance licensing agreements for the utilization of the generative 1254 model and its outputs. 1255

- 1256
- 1257 1258

CURRENT LIMITATION AND FUTURE DIRECTIONS L

The limitation of this work is that it only measures data value for vision-related generative models 1259 and conducts experiments exclusively within the field of computer vision. However, this does not 1260 mean that GMVALUATOR cannot be easily adapted to Natural Language Processing (NLP) fields, 1261 given its core idea of similarity matching. In the future, we should extend GMVALUATOR to NLP 1262 and assess the data value for language-related generative models, such as large language models 1263 (LLMs).

1264 1265

1267

1273

1274 1275

Μ **OMITTED PROOFS**

We follow Just et al. (2023) to prove the theorem. Firstly, we give several assumptions that will be 1268 used in later proof. 1269

1270 **Assumption M.1** Following Assumption 2.3, given a distance function $d(\cdot, \cdot)$ between, we defined 1271 the coupling between $\mathcal{X}_{(T|f)}$ and $\mathcal{X}_{(S^*|f)}$ as π^* : 1272

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

It is easy to see that all joint distributions defined above are couplings between the corresponding 1276 distribution pairs. Then, following Just et al. (2023) we prove the main Theorem. 1277

1278 **Theorem M.2** (Restated of Theorem 2.4.) Let $f'_{S^*}: \mu \to \mathcal{A} = \{0,1\}^V$ be the model trained on the 1279 optimal contributor dataset S^* . Following Assumption 2.3, if the contributors are corresponding to 1280 the given generated data \hat{X} , we have: 1281

1282

1284 1285 1286

1287

1290 1291

1292

1294

$$\mathbb{E}_{x \sim \mu_{T}} \left[\mathcal{L} \left(f(x), f_{S^{*}}^{'}(x) \right) \right] - \mathbb{E}_{x \sim \mu_{S}} \left[\mathcal{L} \left(f(x), f_{S^{*}}^{'}(x) \right) \right]$$

$$\leq k \epsilon \cdot \left[d_{W}(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)}) + d_{W}(\mathcal{X}_{(S^{*}|f)}, \mathcal{X}_{(\hat{X}|f)}) \right]$$
(8)

Proof M.3

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

 $\leq \mathbb{E}_{x \sim \mu_S} \left[\mathcal{L} \left(f(x), f_{S^*}(x) \right) \right]$

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

$$(9)$$

$$1295$$

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

 $= \left| \int_{\mathcal{X}^{2}} \mathcal{L}\left(f(x_{S^{*}}), f_{S^{*}}^{'}(x_{S^{*}}) \right) \right|$

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

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

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

 $-\mathcal{L}\left(f(x_{S^*}), f_{S^*}'(x_T)\right) + \mathcal{L}\left(f(x_{S^*}), f_{S^*}'(x_T)\right) - \mathcal{L}\left(f(x_T), f_{S^*}'(x_T)\right) d\pi^*(x_T, x_{S^*})|$

$$\begin{aligned} \text{RHS of } Eq.equation \ 10 &\leq k \int_{\mathcal{X}^2} ||f'_{S^*}(x_{S^*}) - f'_{S^*}(x_T)||d\pi^*(x_T, x_{S^*}) \\ &+ k \int_{\mathcal{X}^2} ||f(x_{S^*}) - f(x_T)||d\pi^*(x_T, x_{S^*}) \\ &\leq k\epsilon \int_{\mathcal{X}^2} 2d(x_T, x_{S^*})d\pi^*(x_T, x_{S^*}) \\ &= k\epsilon d_W(\mathcal{X}_{(T|f)}, \mathcal{X}_{(S^*|f)}), \end{aligned}$$

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

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

where the last step is due to the definition of 1-Wasserstein distance. Then, according to the triangle
 inequality of Wasserstein distance Peyré et al. (2019), we can obtain:

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

(10)

1325 Combining Eq. equation 9 and Eq. equation 11 we finished the proof and obtained the Theorem 2.4. 1326 By reducing the distance term $d_W\left(\mathcal{X}_{(T|f)}, \mathcal{X}_{(\hat{X}|f)}\right)$, we have $\mathcal{X}_{(T|f)} \to \mathcal{X}_{(S^*|f)}$. As a result, the 1327 expected distance

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

1330 with randomly sampling S^* and T with K elements.

N REPRODUCIBILITY

To ensure reproducibility, we make our implementation available to reviewers through this anonymous link: https://anonymous.4open.science/r/GMValuator-V2-E0BE.