

# An Economic Solution to Copyright Challenges of Generative AI

Jiachen T. Wang<sup>1</sup>    Zhun Deng<sup>2</sup>    Hiroaki Chiba-Okabe<sup>4</sup>    Boaz Barak<sup>3\*</sup>  
Weijie J. Su<sup>4</sup>

<sup>1</sup> Princeton University

<sup>2</sup> Columbia University

<sup>3</sup> Harvard University

<sup>4</sup> University of Pennsylvania

April 24, 2024

## Abstract

Generative artificial intelligence (AI) systems are trained on large data corpora to generate new pieces of text, images, videos, and other media. There is growing concern that such systems may infringe on the copyright interests of training data contributors. To address the copyright challenges of generative AI, we propose a framework that compensates copyright owners proportionally to their contributions to the creation of AI-generated content. The metric for contributions is quantitatively determined by leveraging the probabilistic nature of modern generative AI models and using techniques from cooperative game theory in economics. This framework enables a platform where AI developers benefit from access to high-quality training data, thus improving model performance. Meanwhile, copyright owners receive fair compensation, driving the continued provision of relevant data for generative model training. Experiments demonstrate that our framework successfully identifies the most relevant data sources used in artwork generation, ensuring a fair and interpretable distribution of revenues among copyright owners.

## 1 Introduction

Recent advancements in generative artificial intelligence (AI) have profoundly impacted the creative industries, ushering in an era of AI-generated content in literature, visual arts, and music. Trained on vast datasets of human-generated material, generative AI models such as large language models and diffusion models can now produce content with a sophistication that rivals—and may potentially displace—the works of human artists [28, 2, 13]. This burgeoning capability raises crucial questions about the legal and ethical boundaries of creative authorship, particularly concerning copyright infringement by generative models [30, 32]. Consequently, several AI companies are currently involved in lawsuits over allegations of producing content that potentially infringes on copyrights [32, 11].

Efforts to mitigate the tension between owners of copyright in the training data and AI developers have emerged, mostly involving modifications to generative model training or inference to reduce the likelihood of generating infringing outputs [35, 4, 33]. However, these modifications

---

\*Currently also at OpenAI. Work done while at Harvard.

may compromise model performance due to either the exclusion of high-quality, copyrighted training data from training or restrictions on content generation [19]. The complexity and ambiguity of copyright law add another layer of difficulty, blurring the line between infringing and non-infringing outputs. The resulting uncertainty could lead to a significant waste of resources on both sides while these issues are debated in courts [32].

Rather than restricting AI developers’ use of copyrighted data, we propose establishing a mutually beneficial revenue-sharing agreement between AI developers and copyright owners. This proposal echoes an argument recently advocated in economics [1]. However, a major challenge in developing a revenue-sharing model for generative AI, in contrast to conventional cases of sharing between digital platforms and independent content creators [6], lies in the complexity of training generative models on diverse data sources. This results in the “black-box” nature of model training and content generation, making the traditional, straightforward pro rata methods unsuitable [21].

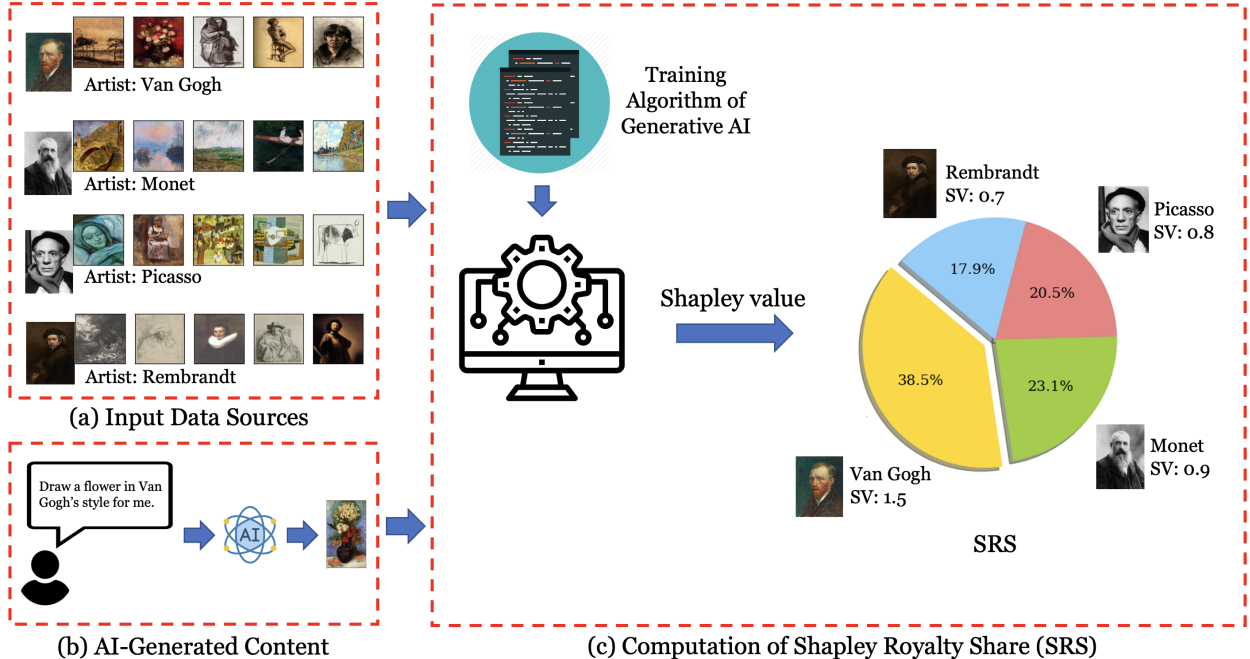
In this paper, we introduce a simple framework that appropriately compensates copyright owners for using their copyrighted data in training generative AI systems based on the cooperative game theory, thereby directly addressing the intricacies of copyright challenges. Our framework does not require modifying the inference process and preserves the full capabilities of generative models. We propose the royalty distribution model for sharing revenue with copyright owners by leveraging the probabilistic nature of generative models: the log-likelihood of generating the user-chosen content is used to measure the utility of the training data. This utility measure captures the capabilities of the model in satisfying users’ needs. Royalties are subsequently distributed among the copyright owners in accordance with their contributions, which are analytically determined using the theory of Shapley value [34]. By aligning compensation with these quantifiable contributions, our framework ensures interpretability in the distribution of royalties, thereby fostering innovation in AI while guaranteeing a fair share of benefits to all copyright holders.

## 2 The Shapley Royalty Share Framework

Our framework takes two steps to tackle copyright issues associated with generative AI models. The first step is to evaluate the utility of the model trained on every possible subset of the entire dataset. Intuitively, the utility of the data subset would be large if this model could generate with a great chance the same AI-generated content (e.g., an artwork) as the deployed model, which is trained on the entire dataset. The second step is to determine any participating copyright owners’ rightful share based on the utilities from the first step, using tools from cooperative game theory. Loosely speaking, a copyright owner’s share would be large if the utility tends to increase by including its data in the model training.

**Utilities of different data source combinations** Let there be  $n$  copyright owners and the  $i$ th owns the copyright of training data  $D^{(i)}$ , where  $i \in N := \{1, \dots, n\}$ . The deployed model is trained on the entire dataset  $D := D^{(1)} \cup \dots \cup D^{(n)}$  and generates a content  $x^{(\text{gen})}$ . Consider a counterfactual model that is trained on a subset of training data,  $\cup_{i \in S} D^{(i)}$ , where  $S \subset N$  denotes a subset of data owners. The utility of the counterfactual model could be best reflected by its likelihood of generating the same content  $x^{(\text{gen})}$  as the deployed model. Let  $p_S(\cdot)$  denote the probability density function of this counterfactual model. We define the utility of this model for content  $x^{(\text{gen})}$  as

$$v(S; x^{(\text{gen})}) := \log p_S(x^{(\text{gen})}). \tag{2.1}$$



**Figure 1:** Overview of our method. (a) The artists provide their copyrighted artworks as (part of) the training data for the generative AI model. (b) A user prompts the generative AI and obtains a new artwork. (c) We assess the contribution of each artist to the AI-generated artwork using the Shapley Royalty Share, which determines their compensation.

The utility offers a way to measure the extent to which the data sources from  $S$  are responsible for generating the content. It is small if the counterfactual model is unlikely to generate the same content as the deployed model, and vice versa.

In practice, the generation of the content involves prompts and human interactions, from which we can write the density conditional on an event  $Q$ . The utility definition becomes  $v(S; x^{(\text{gen})}) = \log p_S(x^{(\text{gen})}|Q)$ . More generally, the utility can be defined relative to a baseline model, which, for example, is trained only on data in the public domain (that is,  $S$  is the empty set  $\emptyset$ ). The relative utility is defined as

$$v(S; x^{(\text{gen})}) = \log \frac{p_S(x^{(\text{gen})}|Q)}{p_{\emptyset}(x^{(\text{gen})}|Q)}. \quad (2.2)$$

This formulation can be viewed as the additional information about  $x^{(\text{gen})}$ , in bits, contributed by the data of  $S$  beyond what is available in the public domain dataset.

**Royalty sharing among copyright owners** The utility (2.1) or (2.2) can be interpreted as the total compensation all members of  $S$  collectively deserve for providing their data to train the generative AI model. The next step is to determine the payoff for each individual copyright owner, based on the utilities of all possible combinations of data sources. We propose using the *Shapley value* [34] for this task. The Shapley value is a solution concept in cooperative game theory that offers a principled approach to distributing gains depending on the utility of every combination of players as a coalition. It is the only payment rule satisfying several important economic properties

(see the supplementary materials for details) [34, 29] and has gained popularity in data valuation for machine learning models [9, 16].

Given the utility  $v(S)$  defined in (2.1) or (2.2), the Shapley value of the  $i$ th copyright owner is defined as

$$\phi_i := \frac{1}{n} \sum_{k=1}^n \binom{n-1}{k-1}^{-1} \sum_{\substack{S \subseteq N \setminus \{i\} \\ |S|=k-1}} [v(S \cup \{i\}) - v(S)]. \quad (2.3)$$

The Shapley value remains the same regardless of whether the absolute utility (2.1) or its relative counterpart (2.2) is being used. At a high level, it rewards a copyright owner based on the weighted average of utility changes caused by adding this contributor’s data to all possible coalitions. The Shapley value is large if the addition of the data source enhances the likelihood of generating the artwork for many combinations of contributors’ data in training. In particular, it equals zero if the contributor’s data does not impact the likelihood of generating the content  $x^{(\text{gen})}$  for any combination.<sup>1</sup>

In our framework, the  $i$ th copyright owner receives a payoff proportional to the following Shapley Royalty Share (SRS) for the AI-generated content  $x^{(\text{gen})}$ :

$$\frac{\phi_i(v(\cdot; x^{(\text{gen})}))}{\sum_{j=1}^n \phi_j(v(\cdot; x^{(\text{gen})}))}, \quad (2.4)$$

where the denominator equals the total relative utility  $v(N; x^{(\text{gen})})$  defined in (2.2) due to the efficiency property of the Shapley value [34]. When  $\phi_i(v(\cdot; x^{(\text{gen})}))$  is negative, we replace it by zero in both the numerator and denominator of (2.4). For instance, if the user pays one dollar for generating  $x^{(\text{gen})}$ , a copyright owner would receive a payoff in the amount of its SRS. However, in practice, it is reasonable for the AI developers to retain a fraction of the revenue since it costs considerable resources to train the model. We defer the discussion of this point to Section 5.

From the definition of the SRS in (2.4), if a contributor’s data has a relatively large Shapley value, this contributor would receive a large royalty share, and vice versa. As the Shapley value is a fair metric of each party’s contribution to the coalition [34], the SRS offers a principled approach to assigning royalty shares to copyright owners. A related approach is called leave-one-out (LOO) score [5, 6], which examines the effect of removing a single data point or source from the entire training dataset. However, it may not capture the complex interactions among various data sources. This shortcoming becomes especially pronounced with data duplication across various copyright owners, which is common in machine learning applications [20]; see the supplementary materials for detailed discussion.

**Computational considerations** A main challenge in applying the framework of SRS lies in its substantial computational cost. The evaluation of the utility functions on different combinations of data sources requires retraining the model multiple times. In some applications where the number of

---

<sup>1</sup>We note that, in practice, due to learning stochasticity, utility functions are randomized, rendering the Shapley value a random variable. While previous research has demonstrated that such stochasticity can significantly affect the estimation of Shapley values when each player contributes only a single data point, this paper focuses on scenarios where each player possesses a data source [36]. In such settings, the impact of learning stochasticity on Shapley value estimation is minimal.

copyright owners is small, the computational challenge might not be as severe as it seems. Indeed, we envision that this contract-based framework works best when the entire copyrighted data is partitioned among a handful of copyright owners so that each source has enough data to impact the training outcome. If the data source is very small in size, the royalty share of the owner would be mostly insignificant and, worse, noisy due to the stochastic nature of training AI models [36].

To alleviate such computational burdens, two approaches can be applied here. The first is to use the Monte Carlo method to approximate the Shapley value [16, 15, 26, 38, 3, 25, 23, 37]. This technique is specially tailored to the case of a large number of copyright owners in the coalition. The second approach is to train a model by fine-tuning it from another model that is trained on a smaller subset of data. Hence, one can approximate models trained on different subsets of data sources by training the model with only one pass through the entire training data. Specifically, for a randomly sampled permutation of copyright owners, we can first train on the first copyright owner, then the second, and all the way up to the last copyright owner. This technique can be used together with the famous permutation sampling estimator for the Shapley value [24].

In practice, a commercial AI model could undergo millions of transactions on a daily basis. It suffices to estimate the aggregated payoffs each copyright owner deserves instead of calculating the payoff as specified in (2.4) for each AI-generated content. To save computational cost, we can evaluate the SRS for only a small fraction of all transactions and scale back to obtain estimates of the revenue distributions from all transactions; see a detailed discussion in the supplementary materials.

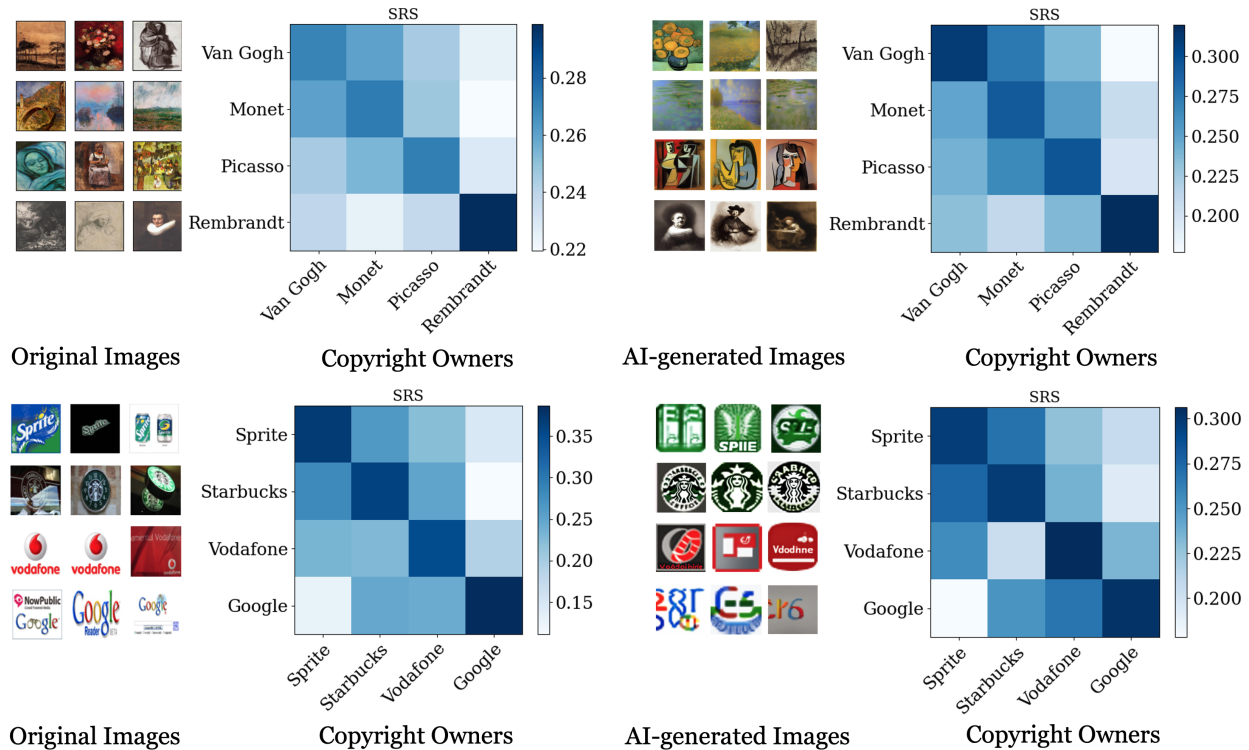
**Beyond copyright considerations** Our framework not only tackles copyright disputes but also addresses scenarios where multiple entities, each holding a private dataset, seek to jointly train a generative AI model with the objective of generating revenue from its application. Though initially driven by copyright concerns, our SRS framework adapts seamlessly to these new scenarios, ensuring fair revenue sharing among private data owners. Crucially, this approach addresses potential financial disagreements and facilitates decentralized model development.

### 3 Results

We assessed the proposed framework’s effectiveness in distributing royalties for AI-generated content using experiments, with a focus on creative art and logos in the image domain. Our evaluation utilized publicly available datasets: WikiArt [31] and FlickrLogo-27 [17], with detailed dataset and training algorithm settings provided in the supplementary materials.

**Evaluation protocol** For the WikiArt dataset, we selected four disjoint subsets of paintings from four renowned artists. A model, initially trained on a broader set of training images (excluding those belonging to the four artists), served as the base model. The SRS is computed by further fine-tuning the base model on various combinations of the four painting sets belonging to the selected artists. Similarly, for the FlickrLogo-27 dataset, we selected four disjoint subsets of logo designs from four brands, and computed the SRS using a base model trained on logo images from other brands. Our goal was to assess whether the SRS can reflect each copyright owner’s contribution to the generation of images.

**Identifying relevant copyright owners** Figure 2 shows the computed SRS for different kinds of  $x^{(\text{gen})}$ 's that are either the original or the AI-generated painting that is in the style of different artists. The results indicate that the SRS has the highest values when the  $x^{(\text{gen})}$ 's closely resembles the training data source in style. This relationship underscores the SRS framework's ability to accurately attribute contributions to the creation of AI-generated images.

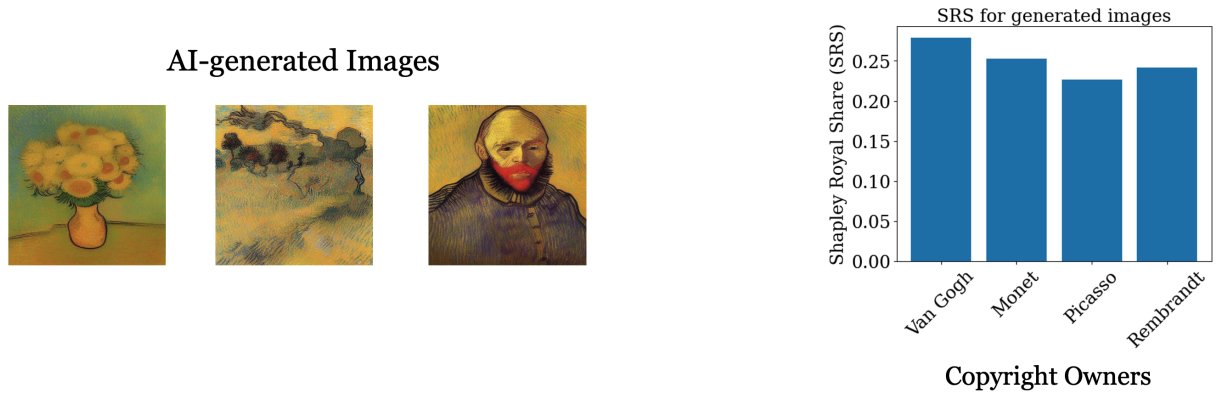


**Figure 2:** Evaluation of the SRS using the WikiArt (upper) and FlickrLogo-27 datasets (lower): Each row displays example target images ( $x^{(\text{gen})}$ 's) for which the SRS is assessed. Left: The heatmap of the SRS of copyright owners in producing the original paintings from different artists (or original logo designs from different brands). Right: The heatmap of the SRS of copyright owners in producing AI-generated paintings in the style of different artists (or AI-generated logo designs of different brands).

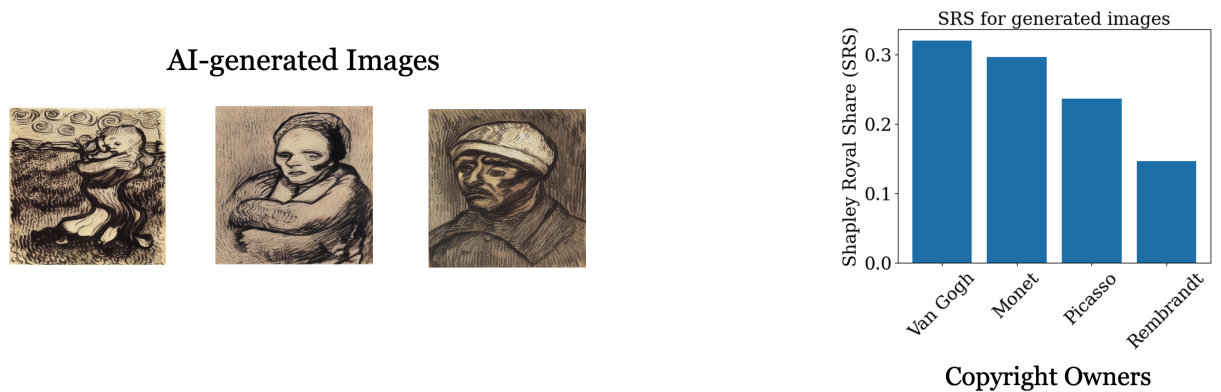
**Assessing mixed-style generation** In Figure 3, we explored the SRS distribution for prompts requesting content generation from multiple data sources. Notably, for the WikiArt dataset, prompts asked the generative model to blend styles from multiple artists. The SRS effectively recognized and rewarded the contributions of data sources integrated into the generated artworks, showcasing the framework's capability to discern and value diverse data source inputs to generate content.

**Non-copyrighted data** We further explored the SRS framework's response to prompts requesting content generation from non-copyrighted data sources, as shown in Figure 4. In these scenarios, the SRS distribution was observed to be nearly uniform across all copyright owners. This outcome aligns with expectations, as the generated content lacks direct ties to any of the copyrighted data

Prompt: “I want a painting in the style of the combination of Van Gogh and Monet”



Prompt: “I want a painting in the style of the combination of Van Gogh, Monet, and Picasso”



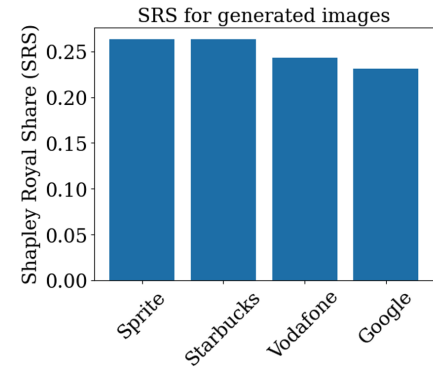
**Figure 3:** Results on the WikiArt dataset when prompting generative AI to produce a painting based on multiple copyright owners’ style. Left: The generated images. Right: The histogram of the SRS of different copyright owners when using AI-generated images.

sources. This uniformity demonstrates the SRS framework’s ability to avoid disproportionate revenue distribution.

**Ranking of contributions via SRS** In many applications, it is essential to understand the hierarchy of contributions from data sources. To validate the SRS framework’s capability to do so, we conducted experiments using the CIFAR100 dataset [18], focusing on four distinct categories: Aquarium Fish, Other Fish, Aquatic Mammals, and Furniture. With “Aquarium Fish” images as the baseline for generation, it is natural to expect the following relevance order: Aquarium Fish > Other Fish > Aquatic Mammals > Furniture. Figure 5 shows that the SRS framework accurately reflects this expected ordering, demonstrating its robustness in discerning the relative significance of contributions from diverse data sources.

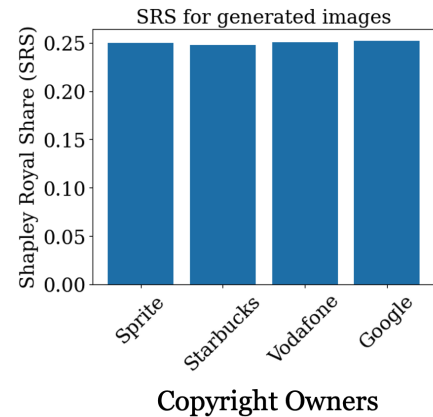
Prompt: “A logo by Coca-Cola”

AI-generated Images



Prompt: “A logo by DHL”

AI-generated Images



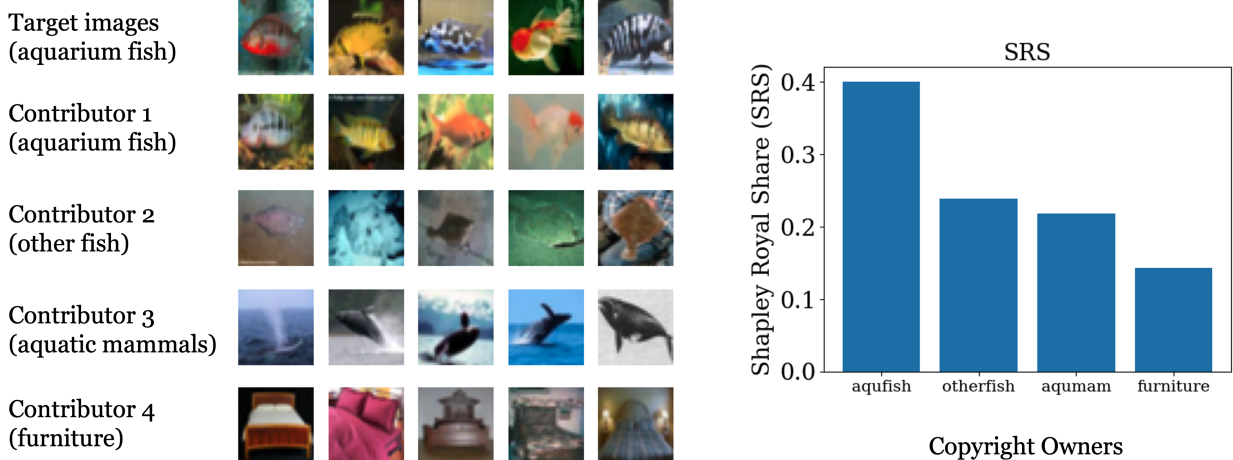
**Figure 4:** Results on FlickrLogo-27 Dataset when prompting generative AI for producing logos for Coca-Cola / DHL, brands whose logo images are not contained in any of the copyright owners’ training set. Left: The generated images. Right: The histogram of the SRS of different copyright owners for the AI-generated images.

## 4 Related Work

Recent efforts in machine learning have primarily focused on minimizing the likelihood of creating copyright-infringing content by generative AI models. One approach involves training an auxiliary generative model on non-copyrighted data and utilizing rejection sampling to reduce the likelihood of reproducing copyrighted material [35]. However, this method is susceptible to adversarial attacks [22]. Alternatively, [4] suggests modifying generative models’ training objectives to avoid generating outputs that closely resemble copyrighted data. Yet another technique focuses on protecting unique artistic styles by incorporating adversarial perturbations into copyrighted images for model fine-tuning [33].

The Shapley value has been suggested as a means to fairly distribute revenue in traditional sectors such as royalty agreements between music copyright holders and radio broadcasters [39]. The Shapley value has been used for data valuation where the utility function is the prediction accuracy of the machine learning model [9, 16]. This differs from our SRS framework, which uses the log-likelihood as the utility since there is no such thing as prediction accuracy for generative





**Figure 5:** We divide the CIFAR100 training set into 4 disjoint groups, where each group has a different level of similarity to the target images, and then evaluate the SRS.

models.

Other data valuation techniques have been developed for generative models. A simple approach utilizes similarity scores between training data and generated content as a valuation metric [40]. Another commonly used approach is the leave-one-out (LOO) score or its approximations. For example, [8] extends the TRAK framework [27] to generative models, and [41] further introduced empirical approaches to improving the performance of [8]. Notably, [6] proposed a revenue-sharing mechanism for AI-generated music based on TRAK, which is closely related to our work. However, the LOO scores neglect the high-order training data interactions, which may result in undesirable attribution scores (see Appendix B for detailed discussion).

## 5 Discussion

The recent rise of generative AI has profoundly challenged traditional copyright laws, driven by its powerful generating capabilities. This is compounded by the intricacies in the interpretation of copyrights for AI-generated content as well as the black-box nature of large AI systems. We have addressed these issues from an economic standpoint by developing a royalty sharing model that permits training on copyrighted data in exchange for revenue distribution among copyright owners. This fosters mutually beneficial cooperation between the AI developers and copyright owners. Our framework has several economic underpinnings that render it fair and interpretable. We demonstrate the effectiveness and feasibility of this framework through numerical experiments.

Our study, however, has limitations and opens avenues for future investigation. One concern is potential strategic behaviors, such as copyright owners merging or splitting their data to maximize their royalty share. The SRS could be manipulated by a malicious copyright owner creating multiple copies of their data. While replication-robust solution concepts have been explored [12], they focused on the impact on Shapley values rather than ratios under replication. Developing a mechanism robust against such manipulation is an important direction for future work. Another open question is handling copyrighted data when owners are unable or unwilling to negotiate agree-

ments, particularly with numerous owners each having small datasets. In such cases, our approach could be combined with methods for generating lawful content [35]. Enhancing our model to determine appropriate revenue division between copyright owners and AI developers, acknowledging the critical role of computational resources, algorithm design, and engineering expertise in developing high-performance AI models, is another avenue for research. We have made preliminary progress toward this by adapting the concept of permission structure from cooperative game theory [10] to model the scenario where the AI developers and copyright owners jointly train a generative AI; see the supplementary materials for details.

From a methodological perspective, a crucial aspect warranting future research is the use of Shapley value ratios for revenue distribution. The key challenge with directly using the Shapley value lies in the unknown total revenue for any coalition of copyright owners' data. The log-likelihood ratio (2.2) serves as a surrogate for this unknown quantity. However, the efficiency property of the Shapley value [34], which ensures the sum of Shapley values equals the grand coalition's utility, loses meaning when considering ratios. In this light, semivalues [7], which are a generalization of the Shapley value that drop the efficiency axiom, could provide a viable alternative. Future work could aim to establish axiomatic justifications to identify the most suitable solution concepts within the semivalue class for royalty distribution in this context.

## Acknowledgements

JTW and ZD conducted this work as independent researchers. HCO, BB, and WJS were supported in part by NSF grant DMS-2310679, Wharton AI for Business, a Simons Investigator Fellowship, the Simons Foundation Math+X Grant to the University of Pennsylvania, NSF grant DMS-2134157, DARPA grant W911NF2010021, and DOE grant DE-SC0022199. We are grateful to Peter Henderson for providing helpful feedback on an early version of this paper.

## References

- [1] Erik Brynjolfsson, Danielle Li, and Lindsey R Raymond. Generative ai at work. Technical report, National Bureau of Economic Research, 2023.
- [2] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [3] Mark Alexander Burgess and Archie C Chapman. Approximating the shapley value using stratified empirical bernstein sampling. In *IJCAI*, pages 73–81, 2021.
- [4] Timothy Chu, Zhao Song, and Chiwun Yang. How to protect copyright data in optimization of large language models? *arXiv preprint arXiv:2308.12247*, 2023.
- [5] R Dennis Cook and Sanford Weisberg. Characterizations of an empirical influence function for detecting influential cases in regression. *Technometrics*, 22(4):495–508, 1980.
- [6] Junwei Deng and Jiaqi Ma. Computational copyright: Towards a royalty model for ai music generation platforms. *arXiv preprint arXiv:2312.06646*, 2023.

- [7] Pradeep Dubey, Abraham Neyman, and Robert James Weber. Value theory without efficiency. *Mathematics of Operations Research*, 6(1):122–128, 1981.
- [8] Kristian Georgiev, Joshua Vendrow, Hadi Salman, Sung Min Park, and Aleksander Madry. The journey, not the destination: How data guides diffusion models. 2023.
- [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.
- [10] Robert P Gilles, Guillermo Owen, and Rene van den Brink. Games with permission structures: the conjunctive approach. *International Journal of Game Theory*, 20(3):277–293, 1992.
- [11] Michael M Grynbaum and Ryan Mac. The times sues openai and microsoft over ai use of copyrighted work. *The New York Times*, 27, 2023.
- [12] Dongge Han, Michael Wooldridge, Alex Rogers, Shruti Tople, Olga Ohrimenko, and Sebastian Tschiatschek. Replication-robust payoff-allocation for machine learning data markets. *arXiv preprint arXiv:2006.14583*, 2020.
- [13] Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A. Lemley, and Percy Liang. Foundation models and fair use. *arXiv preprint arXiv:2303.15715*, 2023.
- [14] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- [15] Ferenc Illés and Péter Kerényi. Estimation of the shapley value by ergodic sampling. *arXiv preprint arXiv:1906.05224*, 2019.
- [16] 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.
- [17] Yannis Kalantidis, Lluís Garcia Pueyo, Michele Trevisiol, Roelof van Zwol, and Yannis Avrithis. Scalable triangulation-based logo recognition. In *Proceedings of the 1st ACM international conference on multimedia retrieval*, pages 1–7, 2011.
- [18] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [19] Katherine Lee, A Feder Cooper, and James Grimmelmann. Talkin’ bout ai generation: Copyright and the generative-ai supply chain. *arXiv preprint arXiv:2309.08133*, 2023.
- [20] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8424–8445, 2022.
- [21] Xiaochang Lei. Pro-rata vs user-centric in the music streaming industry. *Economics Letters*, 226:111111, 2023.

- [22] Xiang Li, Qianli Shen, and Kenji Kawaguchi. Probabilistic copyright protection can fail for text-to-image generative models. *arXiv preprint arXiv:2312.00057*, 2023.
- [23] Jinkun Lin, Anqi Zhang, Mathias Lécuyer, Jinyang Li, Aurojit Panda, and Siddhartha Sen. Measuring the effect of training data on deep learning predictions via randomized experiments. In *International Conference on Machine Learning*, pages 13468–13504. PMLR, 2022.
- [24] Sasan Maleki, Long Tran-Thanh, Greg Hines, Talal Rahwan, and Alex Rogers. Bounding the estimation error of sampling-based shapley value approximation. *arXiv preprint arXiv:1306.4265*, 2013.
- [25] Rory Mitchell, Joshua Cooper, Eibe Frank, and Geoffrey Holmes. Sampling permutations for shapley value estimation. 2022.
- [26] Ramin Okhrati and Aldo Lipani. A multilinear sampling algorithm to estimate shapley values. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 7992–7999. IEEE, 2021.
- [27] Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale. 2023.
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [29] Alvin E Roth. Introduction to the shapley value. *The Shapley value*, pages 1–27, 1988.
- [30] Matthew Sag. Copyright safety for generative ai. *Houston Law Review*, 61(2):295–347, 2023.
- [31] Babak Saleh and Ahmed Elgammal. Large-scale classification of fine-art paintings: Learning the right metric on the right feature. *International Journal for Digital Art History*, (2), 2016.
- [32] Pamela Samuelson. Generative ai meets copyright. *Science*, 381(6654):158–161, 2023.
- [33] Shawn Shan, Jenna Cryan, Emily Wenger, Haitao Zheng, Rana Hanocka, and Ben Y Zhao. Glaze: Protecting artists from style mimicry by text-to-image models. *arXiv preprint arXiv:2302.04222*, 2023.
- [34] Lloyd S Shapley. A value for  $n$ -person games. *Contributions to the Theory of Games*, 2(28):307–317, 1953.
- [35] Nikhil Vyas, Sham Kakade, and Boaz Barak. Provable copyright protection for generative models. *arXiv preprint arXiv:2302.10870*, 2023.
- [36] Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine learning. In *International Conference on Artificial Intelligence and Statistics*, pages 6388–6421. PMLR, 2023.
- [37] Jiachen T Wang and Ruoxi Jia. A note on ” towards efficient data valuation based on the shapley value”. *arXiv preprint arXiv:2302.11431*, 2023.

- [38] Tianhao Wang, Yu Yang, and Ruoxi Jia. Improving cooperative game theory-based data valuation via data utility learning. *ICLR 2022 Workshop on Socially Responsible Machine Learning*, 2022.
- [39] Richard Watt. Fair remuneration for copyright holders and the shapley value. *Handbook on the Economics of Copyright. A Guide for Students and Teachers*, 2014.
- [40] Jiayi Yang, Wenglong Deng, Benlin Liu, Yangsibo Huang, and Xiaoxiao Li. Matching-based data valuation for generative model. *arXiv preprint arXiv:2304.10701*, 2023.
- [41] Xiaosen Zheng, Tianyu Pang, Chao Du, Jing Jiang, and Min Lin. Intriguing properties of data attribution on diffusion models. *arXiv preprint arXiv:2311.00500*, 2023.

## A Supplementary Materials

Recall the setting where we have  $n$  copyright owners where each copyright owner  $i \in N = \{1, \dots, n\}$  owns the copyright of training data  $D^{(i)}$ . The utility function  $v$  is defined as

$$v(S; x^{(\text{gen})}) = \log p_S(x^{(\text{gen})}),$$

where  $p_S(\cdot)$  denotes the probability density function of the model trained on  $\cup_{i \in S} D^{(i)}$ . The Shapley value of the copyright owner  $i$  for generating  $x^{(\text{gen})}$  is

$$\phi_i = \frac{1}{n} \sum_{k=1}^n \binom{n-1}{k-1}^{-1} \sum_{\substack{S \subseteq N \setminus \{i\} \\ |S|=k-1}} [v(S \cup \{i\}) - v(S)].$$

The Shapley value is a concept from cooperative game theory and provides a principled approach to fairly distribute the total gains (or costs) among coalition participants based on their individual contributions. The theoretical foundation of the Shapley value is established through four axioms introduced by Lloyd Shapley in his 1953 paper [34]. These axioms delineate criteria for an equitable and logical distribution of payoffs, where the Shapley value is the *unique* solution concept that satisfies all of them. Recall that  $N = \{1, \dots, n\}$  the set of players.

- Dummy player: if  $v(S \cup i) = v(S) + c$  for all  $S \subseteq N \setminus \{i\}$  and a scalar  $c$ , then  $\phi_i = c$ .
- Symmetry: if  $v(S \cup i) = v(S \cup j)$  for all  $S \subseteq N \setminus \{i, j\}$ , then  $\phi_i = \phi_j$ .
- Linearity: for utility functions  $v_1, v_2$  and any scalar values  $\alpha_1, \alpha_2$ ,

$$\phi_i(\alpha_1 v_1 + \alpha_2 v_2) = \alpha_1 \phi_i(v_1) + \alpha_2 \phi_i(v_2).$$

- Efficiency: for every  $v$ ,  $\sum_{i \in N} \phi_i = v(N)$ .

In plain words, the *efficiency* axiom requires the total value to be distributed among individuals. The *symmetry* and *null player* axiom refer to “same contribution, same value” and “no contribution, no value”, respectively. The *linearity* axiom requires the value scores add up when utility functions add up. These principles lay the groundwork for a revenue distribution method that ensures that every participant receives a share of the total value that reflects their contribution to the coalition. The Shapley value’s unique ability to satisfy these conditions makes it a powerful tool for analyzing cooperative scenarios and allocating resources or costs in a manner that is widely considered fair.

### A.1 Efficient SRS Estimation for High-volume Transactions

Considering the potential for millions of daily uses and transactions of a commercial AI model, computing the Shapley Royalty Share (SRS) for each individual transaction may not be feasible due to the high computational cost. To address this challenge, a practical solution is to estimate the average SRS for each copyright owner based on a subset of transactions. For any given transaction involving AI-generated content  $x^{(\text{gen})}$ , let  $p_{x^{(\text{gen})}}$  represent the payment by the user. The share of this payment received by copyright owner  $i$  is determined by:

$$p_{x^{(\text{gen})}} \cdot \text{SRS}_{x^{(\text{gen})}, i},$$

where  $\text{SRS}_{x^{(\text{gen})}} := \frac{\phi_i(v(\cdot; x^{(\text{gen})}))}{\sum_{j=1}^n \phi_j(v(\cdot; x^{(\text{gen})}))}$ . Ideally, the goal is to compute the cumulative payment distribution over all of the daily transactions:

$$\sum_{x^{(\text{gen})}} p_{x^{(\text{gen})}} \cdot \text{SRS}_{x^{(\text{gen})}}.$$

Given the computational cost of SRS calculations, a Monte Carlo method can be applied to randomly sample a batch of  $x^{(\text{gen})}$  transactions and estimate  $\mathbb{E}_{x^{(\text{gen})}}[\text{SRS}_{x^{(\text{gen})}}]$ , and then multiply it by the total amount of the income in a day  $\sum_{x^{(\text{gen})}} p_{x^{(\text{gen})}}$ . This estimation remains unbiased as long as  $\text{SRS}_{x^{(\text{gen})}}$  and  $p_{x^{(\text{gen})}}$  are independent, allowing for:

$$\mathbb{E}_{x^{(\text{gen})}} [p_{x^{(\text{gen})}} \text{SRS}_{x^{(\text{gen})}}] \approx \mathbb{E}_{x^{(\text{gen})}} [p_{x^{(\text{gen})}}] \mathbb{E}_{x^{(\text{gen})}} [\text{SRS}_{x^{(\text{gen})}}].$$

We note that, in practice,  $p_{x^{(\text{gen})}}$  is often a constant number as the charges for different user queries are usually of the same rate. In this case,  $p_{x^{(\text{gen})}}$  and  $\text{SRS}_{x^{(\text{gen})}}$  are clearly independent.

## A.2 The SRS of AI Developers

Recall that within our framework, each copyright owner, denoted by  $i$ , is awarded a payoff directly proportional to their SRS for AI-generated content  $x^{(\text{gen})}$ :

$$\frac{\phi_i(v(\cdot; x^{(\text{gen})}))}{\sum_{j=1}^n \phi_j(v(\cdot; x^{(\text{gen})}))}. \quad (\text{A.1})$$

However, in real-world applications, the AI developing company decides on a collective share,  $\beta_{\text{data}}$ , ranging from 0 to 1, designated for distribution among the copyright owners. The remaining share of the royalty, quantified as  $1 - \beta_{\text{data}}$ , is retained by the company. Consequently, the adjusted royalty share for each copyright owner concerning AI-generated content  $x^{(\text{gen})}$  is calculated as:

$$\beta_{\text{data}} \cdot \frac{\phi_i(v(\cdot; x^{(\text{gen})}))}{\sum_{j=1}^n \phi_j(v(\cdot; x^{(\text{gen})}))}.$$

While the determination of  $\beta_{\text{data}}$  typically falls to the discretion of the AI developer responsible for training the generative models, we propose a principled approach grounded in cooperative game theory to decide  $\beta_{\text{data}}$ . By including the AI developer as a special player in the SRS framework, we can model the scenario as a game with a “permission structure” [10]. Specifically, the permission structure, which is usually represented as a directed graph where nodes symbolize players and directed edges signify permissions, determines valid coalitions. A coalition is considered valid only if every member has the requisite permissions, either inherently or through other coalition members.

In our context, incorporating the AI developer as an additional player expands the player set to  $\mathcal{P} := N \cup \{\text{AI-Dev}\}$ . The presence of the AI developer is indispensable for model training, making it impossible to achieve any utility without their involvement. Consequently, the utility function for  $\mathcal{P}$  is adjusted to  $v^*(S) = v(S \setminus \{\text{AI-Dev}\})$  when **AI-Dev** is included in  $S$ , and  $v^*(S) = 0$  for coalitions excluding **AI-Dev**. To reasonably define  $\beta_{\text{data}}$ , we calculate the Shapley value (and consequently the SRS) of the AI developer within this modified game structure, offering a quantifiable metric of the AI developer’s contribution.

## B Comparison Between Leave-one-out and the Shapley Value

The leave-one-out (LOO) score [5] is a simple, straightforward method for assessing the contribution of data sources. Specifically, the LOO score of a copyright owner is calculated as the model performance change when the data source belonging to the copyright owner is excluded from the full training set:

$$\phi_i^{\text{loo}} := v(N) - v(N \setminus \{i\}).$$

Although intuitive for evaluating the impact of individual data sources, the LOO score has limitations. It solely examines the consequence of removing a data source from the entire dataset. This approach might not accurately reflect the significance of a data point due to potential complex interactions among data sources. Duplicated data points are prevalent across many widely-used machine learning datasets [20]. Consider two copyright owners  $i$  and  $j$  having nearly identical data. The removal of either from the dataset would likely result in minimal change to the model’s content generation likelihood, rendering both LOO scores close to zero. This scenario could unjustly allocate no royalty share to either contributor, despite their datasets’ crucial role in model performance. Moreover, in situations with numerous data sources, the LOO score might diminish to near zero, failing to recognize the nuanced contributions of individual sources.

In contrast, the Shapley value method accounts for the incremental impact of incorporating a data source alongside all possible combinations of other sources. This comprehensive approach effectively captures the intricate dynamics among data sources, offering a more accurate and fair assessment of each contribution.

## C Experimental Settings

**Datasets** Our study focuses on art painting and logo design, two domains where copyright plays a pivotal role in safeguarding the integrity and commercial value of creative outputs. For art paintings, our research employs the WikiArt dataset [31], which comprises approximately 80,000 artworks spanning the last 400 years. This collection features pieces from over 1,000 renowned artists with a wide variety of styles and genres. For logo design, we use FlickrLogo-27 dataset [17], which consists of images from 27 distinct logo classes or brands, sourced from Flickr.

**Model architectures & training details** The generative models used in our experiment follow from the recent advancements in high-resolution image synthesis utilizing latent diffusion models [28].<sup>2</sup> All images were cropped and resized to  $512 \times 512$  resolution. As training a new generative model from scratch would be prohibitively costly, we use LoRA [14], an efficient fine-tuning method that enables us to scale up the models being used in a tractable manner.<sup>3</sup> The model used for fine-tuning is Stable Diffusion V1-4<sup>4</sup>, which is a latent image diffusion model trained on LAION2B-en<sup>5</sup>. For the WikiArt dataset, we train each model on painting images from the same artist, with the text prompt “A painting in the style of [artist name].” Similarly, for the FlickrLogo-27 dataset, we train each model on logo images from the same artist, with the text prompt “A logo by [company

---

<sup>2</sup>Part of the codebase is adapted from <https://github.com/artem-gorodetskii/WikiArt-Latent-Diffusion> and <https://github.com/VSehwag/minimal-diffusion>.

<sup>3</sup>The implementation is adapted from <https://huggingface.co/blog/lora>.

<sup>4</sup><https://huggingface.co/CompVis/stable-diffusion-v1-4>.

<sup>5</sup><https://huggingface.co/datasets/laion/laion2B-en>.



name].” For each dataset, the model is fine-tuned with an initial learning rate  $10^{-4}$ , 10 epochs, and batch size 4. We train the diffusion models on A100 GPU cluster. For all the SRS results in the maintext, they are being averaged over 20  $x^{(\text{gen})}$ ’s that are randomly selected/generated according to the specifications in the maintext.

**Calculating log-likelihood of AI-generated content** The diffusion models have two main processes: the forward (noise-adding) process and the reverse (noise-removing) process. The *forward process* is a Markov chain that gradually adds Gaussian noise to the data over a series of steps. If we represent the original data as  $x_0$ , the process of adding noise can be expressed as:

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon,$$

where  $x_t$  is the data at step  $t$ ,  $\epsilon$  is a sample from a standard Gaussian distribution  $\mathcal{N}(0, I)$ ,  $\alpha_t$  is a variance schedule that determines how much noise to add at each step, and  $t$  ranges from 0 to  $T$ , with  $T$  being the total number of diffusion steps, and  $x_T$  being almost entirely noise. The *reverse process* aims to learn the distribution of the original data by starting from noise and progressively removing it. This can be modeled as:

$$p(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)),$$

where  $\mu_\theta(x_t, t)$  and  $\Sigma_\theta(x_t, t)$  are the mean and covariance of the Gaussian distribution at step  $t$ , learned by the model with parameters  $\theta$ .<sup>6</sup> The model is trained to minimize the difference between the noisy data and its prediction of the denoised data at each step. This can be formalized as minimizing a loss function, e.g., the mean squared error (MSE), between the original data and the reconstructed data:

$$L(\theta) = \mathbb{E}_{t, x_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2],$$

where  $\epsilon_\theta(x_t, t)$  is the model’s prediction of the noise  $\epsilon$  added at step  $t$ , and the expectation  $\mathbb{E}$  is over different noise levels  $t$ , the original data  $x_0$ , and the noise  $\epsilon$ . To generate new data, we start with a sample from the noise distribution  $x_T \sim \mathcal{N}(0, I)$  and iteratively apply the reverse process to obtain  $x_{T-1}, x_{T-2}, \dots, x_0$ , with  $x_0$  being the final generated sample.

In order to estimate the density of a diffusion model  $p_\theta(x^{(\text{gen})})$  on a generated sample  $x_0 := x^{(\text{gen})}$ , note that

$$\begin{aligned} p_\theta(x^{(\text{gen})}) &= \int p_\theta(x^{(\text{gen})}, x_{1:T}) dx_{1:T} \\ &= \int p_\theta(x^{(\text{gen})}|x_{1:T}) p_\theta(x_{1:T}) dx_{1:T} \\ &= \mathbb{E}_{x_{1:T}} [p_\theta(x^{(\text{gen})}|x_{1:T})] \\ &= \mathbb{E}_{x_{1:T}} [p_\theta(x^{(\text{gen})}|x_1)], \end{aligned} \tag{C.1}$$

where  $x_{1:T} := (x_1, \dots, x_T)$ . Since  $p_\theta(x^{(\text{gen})}|x_1) = \mathcal{N}(x^{(\text{gen})}; \mu_\theta(x_1, 1), \Sigma_\theta(x_1, 1))$  whose probability density can be efficiently computed, we can use Monte Carlo technique to estimate  $p_\theta(x^{(\text{gen})})$  based on Equation (C.1). For all experiments we show in the paper, we use 20 random samples of  $x_T \sim \mathcal{N}(0, I)$ , apply the reverse process to obtain random samples of  $x_{1:T}$  and use the sample average of  $p_\theta(x^{(\text{gen})}|x_1)$  as the estimation for  $p_\theta(x^{(\text{gen})})$ .

<sup>6</sup>We note that for text-to-image diffusion models,  $\mu_\theta(\cdot)$  and  $\Sigma_\theta(\cdot)$  will also depend on the input text encoding, but the core mathematical representation remains the same.

## C.1 The SRS of AI developers

In supplementary materials Section A.2, we explore the SRS when the AI developer is considered a special player within a game characterized by a permission structure. Here in Figure 6, we empirically evaluate the SRS in this setting. Specifically, Figure 6 (a) shows the result of SRS for an AI-generated painting in Van Gogh’s style, and Figure 6 (b) shows the result of SRS for an AI-generated logo for Sprite. Both figures show that the AI developer achieves a markedly higher SRS compared to training data contributors. This observation aligns with the intuitive understanding that the AI developer’s contribution is foundational; without their computational input and expertise, it would be infeasible to generate any valuable content.



**Figure 6:** (a) The SRS results on paintings that are generated by prompting “paint in Van Gogh’s style.” (b) The SRS results on logo designs that are generated by prompting “design a logo similar to Sprite.”