

---

# Rethinking Artistic Copyright Infringements in the Era of Text-to-Image Generative Models

---

**Anonymous Author(s)**

Affiliation

Address

email

**ArtSavant Report for Canaletto**

 **You have a unique and recognizable style!**  
We can identify your style (over the style of 372 other artists) in **88.37%** of your works. This puts you in the top 83.6% percentile of artists in recognizability.

 **Your style is detected in works generated by Stable Diffusion.**  
When prompting a gen AI model to copy you, the resultant images exhibit your style more than 372 other artists **70.34%** of the time.

 **We find stylistic elements unique to you that reappear in generated images.**  
We identify some tag signatures (set of stylistic elements that frequently co-occur only in your work) that also appear in generated images. Here's an example; click to see more.  
Matched Tag for Canaletto: oil painting, broad brushwork, geographical symbolism  
0 other artists also have this signature



17 generated images with this signature



8 real images with this signature

Figure 1: We propose a new way to define artistic style and argue style copying, so to be accessible to artists and lawyers. Given artworks by Canaletto, our tool ArtSavant identifies a unique style and recognizes said style in generated art, and produces an easy to understand yet quantitative report.

- 1 Text-to-image generative models [27, 28, 2, 23] have captured widespread attention and at times
- 2 concern, for they may make infringing copyrighted material far easier. While direct copying of
- 3 individual training images seems to be generally rare in diffusion models [5, 30, 31], the degree
- 4 to which image generative models can replicate art *styles* as opposed to art works remains unclear.
- 5 This issue has human and material consequences (potentially unfairly undermining the value of
- 6 original art), and is fundamentally interdisciplinary, engaging artistic and legal communities. There are
- 7 currently no laws to identify and protect an artist's style - mainly due to challenges in definition and a
- 8 previous lack of necessity. However, at least one major actor has proposed such legislation [3], raising
- 9 the issues of how well individual artistic style can be defined, and how much artists should be worried
- 10 that their style can be effectively mimicked. To this end, we seek to tackle the problem of defining and
- 11 identifying artistic styles, as well as building a practical tool to detect instances of style infringement.
- 12 Our tool, ArtSavant, prioritizes accessibility and transparency so that it is useful to a broad audience:
- 13 we make it simple and fast enough for an end-user (e.g., artist or lawyer) to run, and interpretable
- 14 enough so that the user can understand and convey the results to another party (e.g., judge or jury).
- 15 We frame artistic style as characterized by a set of elements that co-occur frequently across an artist's
- 16 *body of work*, which makes it challenging to determine style by inspecting individual works (a la



Figure 2: We define artistic style as a set of elements (or signature) that appear frequently over a body of work, and reduce the problem of style copy detection to classification of *sets* of images to artists. (**left**) We present two ways to recognize artistic styles, including a novel interpretable and attributable method. (**right**) We find gen. models potentially copy artistic styles for 20.2% of 372 prolific artists.

17 previous image-wise copy studies). For e.g., Vincent Van Gogh’s style comprised of expressive wavy  
 18 lines, bright unblended coloring, post-impressionism, choppy textured brushwork, etc. In Figure 3, we  
 19 illustrate that while generative models seldom reproduce Van Gogh’s artworks exactly, they frequently  
 20 capture and replicate elements of his style. While describing his (or any style) can be challenging, and  
 21 making a case for distinctiveness between two styles is even more so, as artists draw inspiration from  
 22 each other, we can still recognize Van Gogh’s style. Building on this intuition, our approach to proving  
 23 the uniqueness of a style is to show that from a collection of artworks, one can identify the artist who  
 24 created them. That is, if an artist’s work can consistently be attributed to its creator, this entails a  
 25 uniqueness to that artist’s style. Therefore, the task of showing the existence and distinctiveness of  
 26 artistic styles can be reduced to **classification** over image *sets*. To empirically study style copying in  
 27 generative models and to build a corpus of artistic styles, we collect a dataset of works from 372 artists,  
 28 and develop two complementary methods to classify artistic style over a body of works, strongly moti-  
 29 vated by notions of ‘holistic’ and ‘analytic’ comparisons from the copyright legal literature [13, 20].

30 The first method – **DeepMatch** – is a neural network that classifies artwork to artists. DeepMatch  
 31 implicitly maps each artist to a vector (via the classification head) during training, which can be inter-  
 32 preted as a *neural signature* representing an artist. Aggregating its predictions over a set of artworks  
 33 via majority voting, we find that DeepMatch achieves 89.3% test accuracy, indicating that **unique**  
 34 **artistic styles indeed exist for a large fraction of artists**. Since deep features are not very interpretable,  
 35 DeepMatch is not suited for articulating the elements that comprise each artistic style. Thus, we com-  
 36plement DeepMatch with a novel inherently *interpretable* and *attributable* method called **TagMatch**.

37 TagMatch first tags individual artworks using a novel method, validated with an MTurk study, based  
 38 on *zero-shot, selective, multilabel classification* with CLIP [24], resulting in tags spanning diverse  
 39 aspects of artistic style. Individual tags are common across artists and thus cannot define unique  
 40 styles alone, but, by efficiently searching the space of tag combinations, we surface *tag signatures*,  
 41 where a set of tags frequently co-occur only over the set of works from a single artist. To map a set of  
 42 unseen works to an artist, we employ a look-up scheme, where we predict the artist who’s works share  
 43 the most unique tag composition with the test set of works. We find tag signatures for *all* artists in our  
 44 dataset, and observe them to be reliable enough to detect the style of the artists in our dataset (on a  
 45 held out set) with 61.6% top-1 and 82.5% top-5 accuracy. Crucially, TagMatch articulates the stylistic  
 46 elements that were uniquely present in the test set of images and the matched reference set, and offers  
 47 as attribution, by way of the subset of images from both sets that contain the matched tag signature.

48 Given a set of works by a concerned artist, ArtSavant applies DeepMatch and TagMatch to generate  
 49 report like Figure 1 in minutes, offering quantitative evidence (if present) of the existence of the  
 50 artist’s unique style and copying by a generative model. To better understand style copying at  
 51 scale, we employ ArtSavant on images generated in the style of artists in our dataset via simple  
 52 prompting of 3 popular text-to-image models. We find 20% of the artists we study to be at risk of  
 53 style copying, though this number may rise as models and prompting schemes grow in sophistication.  
 54 We hope ArtSavant can continue to offer quantitative insight on the prevalence of style copying,  
 55 while also being accessible and practically useful to the broad range of relevant stakeholders.

56 **References**

- 57 [1] Generative artificial intelligence and copyright law, Sep 2023. URL <https://crsreports.congress.gov/product/pdf/LSB/LSB10922>.
- 58 [2] Deepfloyd, Apr 2023. URL <https://github.com/deep-floyd/IF>.
- 59 [3] Adobe. Fair act to protect artists in age of ai, September 12 2023. URL <https://blog.adobe.com/en/publish/2023/09/12/fair-act-to-protect-artists-in-age-of-ai#:~:text=The%20right%20requires%20intent%20to,independent%20creation%20is%20a%20defense>. Accessed: 2024-05-22.
- 60 [4] Samyadeep Basu, Nanxuan Zhao, Vlad Morariu, Soheil Feizi, and Varun Manjunatha. Localizing and editing knowledge in text-to-image generative models, 2023.
- 61 [5] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models, 2023.
- 62 [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *CoRR*, abs/2104.14294, 2021. URL <https://arxiv.org/abs/2104.14294>.
- 63 [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- 64 [8] Stephen Casper, Zifan Guo, Shreya Mogulothu, Zachary Marinov, Chinmay Deshpande, Rui-Jie Yew, Zheng Dai, and Dylan Hadfield-Menell. Measuring the success of diffusion models at imitating human artists, 2023.
- 65 [9] Yingqian Cui, Jie Ren, Han Xu, Pengfei He, Hui Liu, Lichao Sun, Yue Xing, and Jiliang Tang. Diffusionshield: A watermark for copyright protection against generative diffusion models, 2023.
- 66 [10] Yingqian Cui, Jie Ren, Yuping Lin, Han Xu, Pengfei He, Yue Xing, Wenqi Fan, Hui Liu, and Jiliang Tang. FT-SHIELD: A watermark against unauthorized fine-tuning in text-to-image diffusion models, 2024. URL <https://openreview.net/forum?id=0QccFglTb5>.
- 67 [11] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept editing in diffusion models, 2023.
- 68 [12] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016.
- 69 [13] Paul Goldstein. *Goldstein on Copyright, 3rd edition*. Wolters Kluwer Legal & Regulatory U.S., 2014.
- 70 [14] Xinyu Huang, Yi-Jie Huang, Youcai Zhang, Weiwei Tian, Rui Feng, Yuejie Zhang, Yanchun Xie, Yaqian Li, and Lei Zhang. Open-set image tagging with multi-grained text supervision. *arXiv e-prints*, pages arXiv–2310, 2023.
- 71 [15] C. Richard Johnson, Ella Hendriks, Igor J. Berezhnoy, Eugene Brevdo, Shannon M. Hughes, Ingrid Daubechies, Jia Li, Eric Postma, and James Z. Wang. Image processing for artist identification. *IEEE Signal Processing Magazine*, 25(4):37–48, 2008. doi: 10.1109/MSP.2008.923513.
- 72 [16] Sergey Karayev, Aaron Hertzmann, Matthew Trentacoste, Helen Han, Holger Winnemoeller, Aseem Agarwala, and Trevor Darrell. Recognizing image style. In *Proceedings of the British Machine Vision Conference 2014*, BMVC 2014. British Machine Vision Association, 2014. doi: 10.5244/c.28.122. URL <http://dx.doi.org/10.5244/C.28.122>.
- 73 [17] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models, 2020.
- 74 [18] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu. Ablating concepts in text-to-image diffusion models, 2023.
- 75 [19] Mazda Moayeri, Keivan Rezaei, Maziar Sanjabi, and Soheil Feizi. Text-to-concept (and back) via cross-model alignment, 2023.

- 108 [20] Oakes, Calebris, and Sotomayor. Tufenkian import export ventures inc v. einstein moomjy inc,  
 109 2003. URL <https://caselaw.findlaw.com/court/us-2nd-circuit/1455682.html>.
- 110 [21] Maxime Oquab, Timothée Darctet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
 111 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran,  
 112 Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra,  
 113 Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick  
 114 Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without  
 115 supervision, 2024.
- 116 [22] Ed Pizzi, Sreya Dutta Roy, Sugosh Nagavara Ravindra, Priya Goyal, and Matthijs Douze. A  
 117 self-supervised descriptor for image copy detection, 2022.
- 118 [23] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller,  
 119 Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution  
 120 image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.  
 121 URL <https://openreview.net/forum?id=di52zR8xgf>.
- 122 [24] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
 123 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
 124 Sutskever. Learning transferable visual models from natural language supervision, 2021.
- 125 [25] Jie Ren, Han Xu, Pengfei He, Yingqian Cui, Shenglai Zeng, Jiankun Zhang, Hongzhi Wen,  
 126 Jiayuan Ding, Hui Liu, Yi Chang, and Jiliang Tang. Copyright protection in generative ai: A  
 127 technical perspective, 2024.
- 128 [26] Keivan Rezaei, Mehrdad Saberi, Mazda Moayeri, and Soheil Feizi. Prime: Prioritizing inter-  
 129 pretability in failure mode extraction, 2023.
- 130 [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.  
 131 High-resolution image synthesis with latent diffusion models, 2021.
- 132 [28] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed  
 133 Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim  
 134 Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image  
 135 diffusion models with deep language understanding, 2022.
- 136 [29] Shawn Shan, Jenna Cryan, Emily Wenger, Haitao Zheng, Rana Hanocka, and Ben Y. Zhao.  
 137 Glaze: Protecting artists from style mimicry by text-to-image models, 2023.
- 138 [30] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein.  
 139 Diffusion art or digital forgery? investigating data replication in diffusion models, 2022.
- 140 [31] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Gold-  
 141 stein. Understanding and mitigating copying in diffusion models. In A. Oh, T. Neu-  
 142 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neu-  
 143 ral Information Processing Systems*, volume 36, pages 47783–47803. Curran Associates,  
 144 Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/9521b6e7f33e039e7d92e23f5e37bbf4-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/9521b6e7f33e039e7d92e23f5e37bbf4-Paper-Conference.pdf).
- 146 [32] Gowthami Somepalli, Anubhav Gupta, Kamal Gupta, Shramay Palta, Micah Goldblum, Jonas  
 147 Geiping, Abhinav Shrivastava, and Tom Goldstein. Measuring style similarity in diffusion  
 148 models, 2024.
- 149 [33] Wei Ren Tan, Chee Seng Chan, Hernán E. Aguirre, and Kiyoshi Tanaka. Artgan: Artwork  
 150 synthesis with conditional categorial gans. *CoRR*, abs/1702.03410, 2017. URL <http://arxiv.org/abs/1702.03410>.
- 152 [34] Nanne van Noord, Ella Hendriks, and Eric O. Postma. Toward discovery of the artist’s style:  
 153 Learning to recognize artists by their artworks. *IEEE Signal Processing Magazine*, 32:46–54,  
 154 2015. URL <https://api.semanticscholar.org/CorpusID:15774940>.
- 155 [35] Zhenting Wang, Chen Chen, Lingjuan Lyu, Dimitris N. Metaxas, and Shiqing Ma. DIAGNO-  
 156 SIS: Detecting unauthorized data usages in text-to-image diffusion models. In *The Twelfth  
 157 International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=f8S3aLm0Vp>.
- 159 [36] Haotian Xue, Chumeng Liang, Xiaoyu Wu, and Yongxin Chen. Toward effective protection  
 160 against diffusion based mimicry through score distillation, 2024.

- 161 [37] Mert Yuksekgonul, Maggie Wang, and James Zou. Post-hoc concept bottleneck models, 2023.  
 162 [38] Zhengyue Zhao, Jinhao Duan, Kaidi Xu, Chenan Wang, Rui Zhangp Zidong Dup Qi Guo, and  
 163 Xing Hu. Can protective perturbation safeguard personal data from being exploited by stable  
 164 diffusion?, 2023.

## 165 A Related Works

166 The rapid advance of image generative models has made the possibility of mimicking artists' personal  
 167 styles a topic of discussion in the literature [25]. Some works describe ways to either detect direct  
 168 image copying in generated images, or to foil any future copying attempts by imperceptibly altering  
 169 the artists' works to prevent effective training by the generative models. These include techniques  
 170 like adding imperceptible watermarks to copyrighted artworks [35, 9, 10], and crafting "un-learnable"  
 171 examples on which models struggle to learn the style-relevant information [29, 36, 38]. Others  
 172 have suggested methods to mitigate this issue from the model owner's perspective - to either de-  
 173 dupe the dataset before training [5, 30, 31], or to remove concepts from the model after training  
 174 ("unlearning") [18, 11, 4]. Methods like [5, 30, 31] are also more focused on analyzing direct image  
 175 copying from the training data, and thus may not be applicable to preventing style copying.  
 176 None of these works tackle the problem of *detecting* potentially copied art *styles* in generated art,  
 177 especially in a manner which may be relevant to legal standards of copyright infringement. According  
 178 to current US legal standards [1], an artwork has to meet the "substantial similarity" test for it to be  
 179 infringing on copyright. This similarity has to be established on *analytic* and *holistic* terms [20, 13].  
 180 Analytic here refers to explaining an artwork by breaking it down into its constituents using a concrete  
 181 and objective technical vocabulary, while holistic refers to the overall "look and feel" of the artwork.  
 182 So to be relevant to the legal community (who ultimately decides on alleged cases of style copying),  
 183 we design our tool to reflect this dichotomy in its working, while also emphasizing ease of use and  
 184 interpretability, to make our tool practically useful for a concerned artist hoping to protect themselves.  
 185 These priorities manifest in our reformulation of detecting style copying as classification in §C. But  
 186 first, we discuss limitations in applying the typical copy detection approach to artistic styles.



Figure 3: Example generations from Stable Diffusion 2 when prompted to produce specific paintings by Vincent Van Gogh, along with the histogram of similarities between the generated image and corresponding real image. Even for a famous artist like Van Gogh, generative models rarely produce near-exact duplicates. However, Van Gogh's *style* appears consistently, even when similarity is low.

187 **B Motivation: Image-wise similarity may be limited for Style Copying**

188 A prevailing approach to investigating copying involves representing images in a deep embedding  
189 space via models like SSCD [22] or DINO [6], and computing image-to-image similarities across  
190 generated and real images. Such an approach has been employed by [30, 31, 5] to show that generative  
191 models can (though rarely do) create exact replicas of training images. Inspired by these results  
192 and the consequent concerns from artists, we first explore if generative models can recreate famous  
193 artworks, e.g., by Vincent Van Gogh. Specifically, we generate images by prompting “*{artwork title}*  
194 by Vincent Van Gogh” for 1500 Van Gogh works, and compute the DINO similarity between pairs of  
195 a real and corresponding generated image. Figure 3 visualizes the distribution of similarities, as well  
196 as examples at each similarity level. We find that the vast majority of similarities are lower than 0.75,  
197 which amounts to pairs that are far from duplicates. However, even when the generated image differs  
198 significantly from the source real image, certain stylistic elements associated with Van Gogh seem  
199 to appear consistently in the generated works. Thus, while instance-wise copying of artwork appears  
200 rare for even the ultra famous Van Gogh, style copying may require going beyond image-to-image  
201 comparisons, as artists may still have their personal styles, developed over a long career/many artworks  
202 and at significant personal cost, infringed upon in ways that searching for exact replicas would miss.  
203 A concurrent work finetunes embeddings so that cosine similarity better proxies style similarity [32],  
204 though even in this case, the utility of such a tool in court is limited by its lack of interpretability.

205 **C Reformulating Artistic Style Copying as Classification over Image Sets**

206 Having established that style is comprised over a body of work (instead of a single image) and that  
207 copy detection must be interpretable to hold weight in court, we now present an alternate framework  
208 for arguing style infringement, with the following intuition: if an artist’s work can consistently be  
209 distinguished from that of other artists, then there must exist something unique that is present across  
210 that artist’s portfolio. Thus, we can use classification over image sets to demonstrate a unique style  
211 exists given an artist. Then, style infringement can be argued by showing the copied artist can again  
212 be predicted (over many others) given a set of generated works. We now detail DeepMatch and  
213 TagMatch, two complementary methods (w.r.t. accuracy and interpretability) that classify artistic  
214 styles over image sets, in holistic and analytic manners respectively.

215 **A necessary preliminary: WikiArt Dataset.** To distinguish one artist’s style from that of others,  
216 we need a corpus of artistic styles (i.e. portfolios from many artists) to compare against. To this end,  
217 we curate a dataset  $\mathcal{D}$  consisting of artworks from WikiArt<sup>1</sup> (like others [33, 16]) to serve as (i) a  
218 reference set of artistic styles, (ii) a validation set of real art to show (most) artists have unique styles  
219 and our methods can recognize them on held-out sets of their works, and (iii) a test-bed to explore  
220 if text-to-image models replicate the styles of the artists in our dataset in their generated images.  
221 We include  $\sim 91k$  artworks from 372 artists  $\mathcal{A}$  spanning diverse eras and art movements, including  
222 any artist with at least 100 works on WikiArt. Each work is labeled with its genre (e.g., *landscape*)  
223 and style (e.g., *Impressionism*), though we primarily use the artist and title labels. We provide an  
224 easy-to-execute script to enable others to scrape newer versions of this dataset if desired. We now  
225 detail DeepMatch and TagMatch, which each compare a test set of images to our reference corpus.

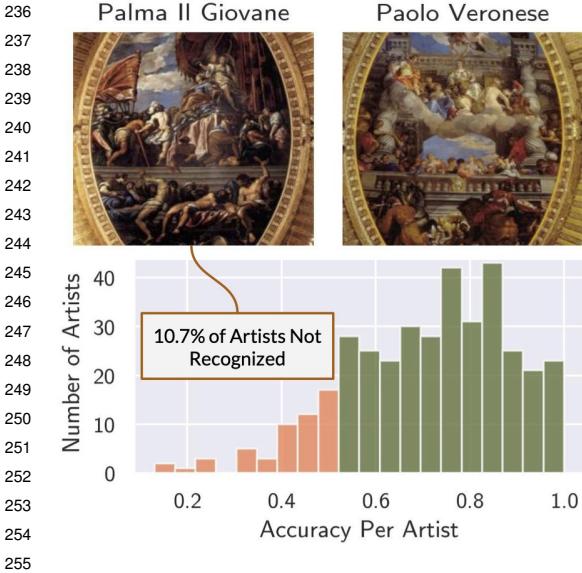
226 **C.1 DeepMatch: Black-Box Detector**

227 DeepMatch consists of a light-weight artist classifier<sup>2</sup> (on images) and a majority voting aggregation  
228 scheme to obtain one prediction for a *set* of images. Majority voting requires that at least half the  
229 images in a test set  $\hat{\mathcal{D}}_a$  are predicted to  $a$  for DeepMatch to predict  $a$ , allowing for abstention in  
230 case no specific style is recognized with sufficient confidence. For our classifier, we train a two layer  
231 MLP on top of embeddings from a frozen CLIP ViT-B\16 vision encoder [24], using a train split  
232 containing 80% of our dataset. We employ weighted sampling to account for class imbalance. Since  
233 we utilize frozen embeddings, training takes only a few minutes on one RTX2080 GPU. Thus, a  
234 new artist could easily retrain a detector to include their works (and thus encode their artistic style).

235 **Validation of the Detector.** We apply DeepMatch on the held-out test split of our dataset and

<sup>1</sup><https://www.wikiart.org/>; note that we only include Public domain or fair use images.

<sup>2</sup>Others have trained art classifiers [16, 15, 34], but they do not operationalize them for style infringement.



256      Figure 4: DeepMatch on held-out real art: 89.3%  
 257      of artists can be recognized. The remaining 10.7%  
 258      of artists have very similar styles to other artists:  
 259      e.g., Palma Il Giovane’s work differs marginally  
 260      from other Italian renaissance painters.

261      artist’s unique style. We do so by tagging images with descriptors (called atomic tags) drawn from  
 262      a vocabulary of stylistic elements. Then, we *compose* tags efficiently to go from atomic tags that are  
 263      common across artists to longer tag compositions that are unique to each artist (i.e. *tag signatures*).  
 264      We detail these steps now, before explaining how tag signatures can be used to classify an image  
 set to an artist in the following section.



Figure 5: Example atomic tags assigned via our proposed CLIP-based zero-shot method. We perform selective multilabel classification along various aspects of art (e.g. medium, colors, shapes, etc), so that atomic tags span diverse categories. Details in section C.2.

265      **Zero-shot Art Tagging** We utilize the zero-shot open-vocabulary recognition abilities of CLIP to  
 266      tag images with descriptors of stylistic elements. First, we construct a concept vocabulary  $\mathcal{V}$  with  
 267      help from LLMs. Namely, we prompt Vicuna-13b and ChatGPT to generate a dictionary of concepts  
 268      along various aspects of art. We manually consolidate and amend the concept dictionary, resulting in  
 269      a vocabulary of 260 concepts over 16 aspects (see Appendix J.1).

270      To assign concepts to images, we design a novel scheme that consists of selective multilabel  
 271      classification per-aspect. Namely, for an image, we compute CLIP similarities to all concepts, and  
 272      normalize similarities *within each aspect*. Then, we only assign a concept its normalized similarity  
 273      (i.e. z-score) exceeds a threshold of 1.75. This means that a concept is only assigned for an aspect if  
 274      the image is substantially more similar to this concept than other concepts describing the same aspect.  
 275      Classifying per-aspect allows for a diversity of descriptors to emerge, as global thresholding results in  
 276      a biased tag description, as concepts for certain aspects (e.g. subject matter) consistently have higher

277 CLIP similarity than those for more nuanced aspects (e.g. brushwork). We call the assigned concepts  
 278 *atomic tags*; figure 5 shows atomic tags assigned for a few examples.

279 **Validation of Quality of Tags Using Human-Study.** We validate the effectiveness of our tagging via  
 280 a human-study involving MTurk workers. In particular, given an image of an artwork and an assigned  
 281 atomic tag  $v_{predict}$  from the vocabulary  $\mathcal{V}$  – MTurk workers are asked “*Does the term  $v_{predict}$  match*  
 282 *(i.e. the concept  $v_{predict}$  present) the artwork below?*”. The workers are then asked to select between  
 283 {Yes, No, Unsure}. We collect responses for 1000 images with 3 annotators each. We find that in  
 284 only 17% cases, a majority of workers disagree with the provided tag, suggesting our tagging results  
 285 in a low false positive rate. We also observe all three annotators agree in only 51% of cases, reflecting  
 286 that describing artistic style can be subjective. While our tagging is not perfect, it is a deterministic and  
 287 automatic method of articulating artistic style elements, and that our tagging method will improve as  
 288 underlying VLMs improve too. See the appendix for more details and discussion on the human study.

289 **Tag Composition for Artists.** Using the atomic tags in the artwork specific vocabulary  $\mathcal{V}$ , in this  
 290 section we design a simple and easy-to-understand iterative algorithm to obtain a set of *tag signatures*  
 291  $\mathcal{S}_a$  for each artist  $a \in \mathcal{A}$ . These signatures are a composition of a subset of tags in  $\mathcal{V}$ . In particular, our  
 292 algorithm efficiently searches the space of tag compositions to go from atomic tags to composition of  
 293 tags which become more unique as the length of the tag composition grows. For e.g., while 40% of  
 294 the artists may use simple colors, *only* 15% may use both simple colors and impressionism style.

295 To efficiently search the space of tag compositions per artist  $a \in \mathcal{A}$ , we first assign a set of tags to  
 296 each of their images  $x \in \mathcal{D}_a$  via the zero-shot *selective multi-label classification* method described  
 297 above. For each image  $x$ , let  $\text{tag}(x)$  denote the set of predicted atomic tags. To get atomic tags *for an*  
 298 *artist*, we aggregate all atomic tags over images, and keep only the tags occurring in at least 3 works.  
 299 We denote this aggregate set of atomic tags as the “Common Atomic Tags Per Artist” and denote it  
 300 as  $\mathcal{C}_a$ . Then, we iterate through all the images  $x \in \mathcal{D}_a$  for a given artist  $a$ , to find the intersection  
 301  $I(x) = \text{tag}(x) \cap \mathcal{C}_a$ . We then compute a powerset  $\mathcal{P}(I(x))$  of the tags occurring in the intersection  
 302  $I(x)$  and increment the count of each occurrence of the tag composition from the powerset in  $\mathcal{S}_a$ .  
 303 Note that the size of  $I(x)$  is much smaller than that of  $\mathcal{C}_a$ , and thus, iterating through  $\mathcal{P}(I(x))$  for  
 304 each image  $x$  is much, much faster than iterating through  $\mathcal{P}(\mathcal{C}_a)$ . Finally, we again filter the tag  
 305 compositions in  $\mathcal{S}_a$ , only including those that occur  
 306 in at least 3 works. We provide the details of this  
 307 tag composition algorithm in 1 and Appendix J.3.

308 **Do Unique Signatures Exist for Artists?** Using  
 309 our tag composition method on the curated dataset  
 310 from WikiArt, we find that *artistic signatures* in  
 311 the form of an unique tag composition exists per  
 312 artist. In Figure 6, we show that our tag composi-  
 313 tion algorithm is able to select unique tag compo-  
 314 sitions such that *only* a very few artists exhibit such  
 315 compositions in their paintings as the tag length  
 316 increases. This shows that artists exhibit *unique*  
 317 *style* which can effectively be captured by our iter-  
 318 ative algorithm. Leveraging these observations, in  
 319 the next section, we describe TagMatch, which can  
 320 classify a set of artworks to an artist by uniquely  
 321 matching such tags (or tag signatures).

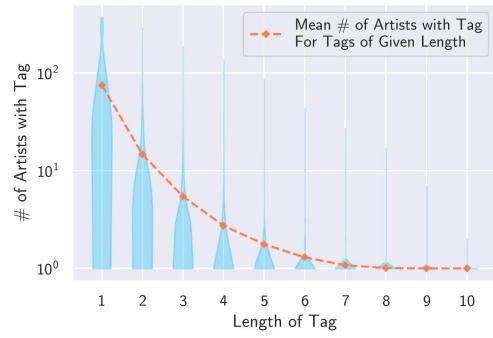


Figure 6: Composing atomic tags results in more unique tags, towards artistic *tag signatures*.

### 322 C.3 TagMatch: Interpretable and Attributable Style Detection

323 In C.1, we outlined a holistic approach to accurately detect artistic styles. While DeepMatch obtains  
 324 high accuracy (recognizing styles for 89.3% of artists), the neural signatures it relies upon lack  
 325 interpretability. For a copyright detection tool to be useful in practice (e.g., to be used as assistive  
 326 technologies), providing explanations of the classification decisions can tremendously benefit the  
 327 end-user. To this end, we leverage our efficient tag composition algorithm as defined in C.2 to develop  
 328 TagMatch - an interpretable classification and attribution method which can effectively classify a set of  
 329 artworks to an artist, as well provide reasoning behind the classification and example images from both  
 330 sets that present the matched tag signature. TagMatch follows the intuition of matching a test portfolio  
 331 to a reference artist who’s portfolio shares the most unique tag signatures. Given a set of  $N$  test images

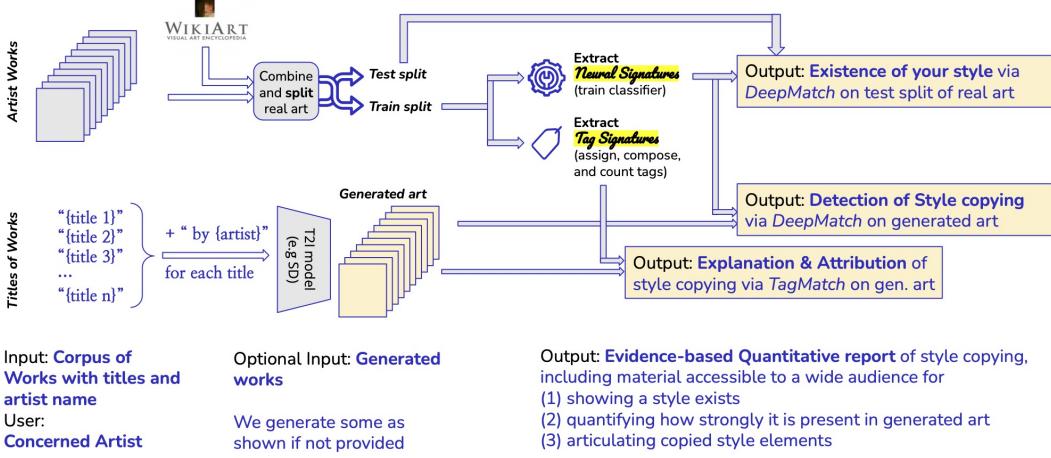


Figure 7: ArtSavant flow. We design our tool with a concerned artist in mind, who wishes to quickly investigate the degree to which they may be at risk of style copying by generative models.

332  $\mathcal{T} = \{x_i\}_{i=1}^N$ , we first obtain a number of tag compositions for them using our iterative algorithm  
 333 in C.2. These tag compositions are then compared with the tag compositions of the artists in the  
 334 reference corpus in order of uniqueness (i.e. we first consider tag signatures present in the test portfolio  
 335 that occur for the fewest number of reference artists). We can then rank reference artists by how  
 336 unique the shared tags are with the test portfolio. Detailed steps of the algorithm is in Appendix J.3.  
 337 Also, TagMatch is fast, taking only about a minute, after caching embeddings of all images.

338 **Validation of TagMatch.** We again utilize the test split of our WikiArt Dataset to validate the  
 339 proposed style detection method. TagMatch predicts the correct artist with top-1 accuracy of 61.6%,  
 340 with top-5 and top-10 accuracies rising to 82.5% and 88.4% respectively. While less accurate than  
 341 DeepMatch, the *tag signatures* provided by TagMatch allow for analytic arguments to be made  
 342 regarding style copying, as the exact tag signatures used in matching can be inspected. Moreover,  
 343 the subset of images in both the test portfolio and matched reference portfolio can be easily retrieved,  
 344 offering direct attribution of the method; examples can be seen in the next section, where we match  
 345 generated images to our reference artists. Overall, we hope TagMatch and DeepMatch can serve  
 346 as automatic and objective tools to navigate the subtle problem of identifying artistic styles, towards  
 347 detecting style copying and helping artists argue their case (i.e. in a court of law) in such instances.

## 348 D ArtSavant: A Practical Tool for Concerned Artists

349 We package DeepMatch and TagMatch into ArtSavant, a practical tool designed with a concerned  
 350 artist in mind. Given a set of works by the concerned artist, ArtSavant would create an easy-to-  
 351 understand report characterizing the degree to which generative models copy the styles of the artist.  
 352 As shown in Figure 7, the artist can present a set of generated images, or we can generate them by  
 353 prompting text-to-image models with prompts of the form “{title of work} by {name of artist}”. The  
 354 provided works are then combined with our existing art repository and split into train/test sets. Using  
 355 the train split, we (a) train a classifier over the 372 + 1 artists, and (b) tag all images, compose tags  
 356 within artists, and store extracted tag compositions per artist, resulting in neural and tag signatures.  
 357 With these, we can apply DeepMatch and TagMatch respectively. Applying DeepMatch to the held-out  
 358 art provides a measure of recognizability, establishing that the artist has an identifiable style to begin  
 359 with. Then, running DeepMatch on generated images provides a quantitative manner to understand  
 360 if (and to what degree) the artist’s style appears consistently in generated works. Finally, running  
 361 TagMatch on the generated images helps articulate the particular style signatures that are copied,  
 362 enabling an analytic way to argue infringement, while also surfacing stylistically similar examples.

363 Figure 1 shows an example report outputted by ArtSavant when presented with art from an artist  
 364 named Canaletto, who we observed was at risk of style infringement. We design the report to be easy  
 365 to read and understand, as well as being evidence-based. Moreover, the report can be generated very

366 quickly. Because all steps operate on embeddings from a frozen CLIP encoder, the process takes  
 367 about 1-2 minutes, as we can simply compute embeddings once (and offline for the WikiArt corpus).

### 368 D.1 Analysis with ArtSavant: Quantifying Style Copying of 372 Prolific Artists

369 While enough anecdotal instances of style mimicry have been observed to raise concern [29, 25],  
 370 the prevalence and nature of such instances remains nebulous. To shed quantitative insight on style  
 371 copying, we now leverage ArtSavant on the  
 372 artists from our WikiArt dataset, generating  
 373 images with three popular text-to-image models:  
 374 (i) Stable-Diffusion-v1.4; (ii) Stable-Diffusion-  
 375 v2.0; and (iii) OpenJourney from PromptHero.  
 376 Following figure 7, we employ a simple prompt-  
 377 ing strategy of augmenting painting titles with  
 378 the name of the artist; we explore alternate  
 379 prompts in I.

380 We first apply **DeepMatch** to see what fraction  
 381 of artists’ styles can be recognized consistently  
 382 over generated images. Namely, each generated  
 383 image is classified to one of 372 artists, and per  
 384 artist, predictions are aggregated via majority  
 385 voting. Figure 8 shows the ‘accuracy’ on  
 386 generated images per artist, where accuracy  
 387 is now interpreted as the rate which images  
 388 generated to copy an artist are classified as  
 389 that artist. In red, the fraction of artists who  
 390 see accuracies of at least 50% (i.e. so that the  
 391 generated image set is classified to the original artist)  
 392 are denoted per model, which we call the  
 393 match rate. We observe an average match rate of 20.2%, indicating that for the vast majority of  
 394 artists in our study, *simple prompting of generative models does not reproduce their styles* in a way  
 395 recognizable to DeepMatch, which has an 89% match rate on real art. For all three models, over half  
 396 the artists see accuracies below 20%, with 26% of artists seeing an average accuracy below 5% for  
 397 generated images. On the other hand, a handful of artists’ styles are matched with high confidence:  
 398 16 artists see average accuracies over 75%. These include ultra famous artists like Van Gogh, Claude  
 399 Monet, Renoir, which we’d expect generative models to do well in emulating. However, a few  
 400 relatively lesser known artists are also present, like Jacek Yerka, who are still alive, and thus could  
 401 be negatively affected by generative models reproducing their styles.

401 With **TagMatch**, in addition to predicting an artistic style, we can also articulate the specific tag  
 402 signature shared between the test set of images and the reference set of images for the predicted  
 403 style. Thus, we can inspect the shared signature, as well as instances from both sets where the  
 404 signature is present, providing direct evidence of the potential style infringement a broader audience  
 405 to independently verify. Inspecting some examples in figure 9 (more in fig. 15), we observe that while  
 406 pixel level differences are common across retrieved image subsets, stylistic elements are consistent in  
 407 both sets with the labeled tags, echoing our motivating claim that style copying goes beyond image  
 408 or pixel-wise similarity. Lastly, TagMatch also allows for understanding image distributions from the  
 409 perspective of interpretable tags. We explore this direction in appendix J.2, finding differences in the  
 410 uniqueness of the tags present in generated art vs real art.

## 411 E Conclusion

412 In our paper, we rethink the problem of copyright infringement in the context of artistic styles. We  
 413 first argue that image-similarity approaches to copy detection may not fully capture the nuance of  
 414 artistic style copying. After reformulating the task to a classification problem over image sets, we  
 415 develop a novel tool – ArtSavant, to reliably and interpretably (via a novel attributable method)  
 416 extract and detect artistic style *signatures* in a way a broader audience can understand. We find  
 417 evidence of the existence of artistic styles, and in an empirical study, quantify the degree to which  
 418 styles are potentially infringed, validating our framework. We hope ArtSavant can be of use to the

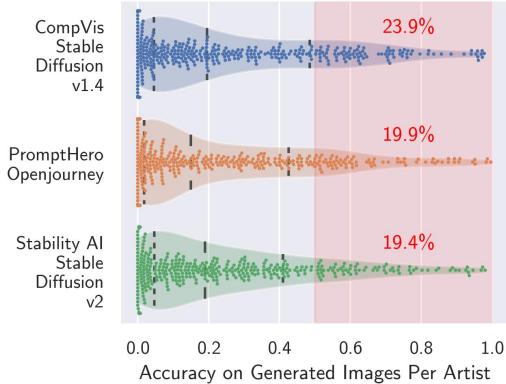


Figure 8: DeepMatch on generated art. In red: the fraction of artists with their styles recognized in at least half of their respective generated images.

match rate. We observe an average match rate of 20.2%, indicating that for the vast majority of artists in our study, *simple prompting of generative models does not reproduce their styles* in a way recognizable to DeepMatch, which has an 89% match rate on real art. For all three models, over half the artists see accuracies below 20%, with 26% of artists seeing an average accuracy below 5% for generated images. On the other hand, a handful of artists’ styles are matched with high confidence: 16 artists see average accuracies over 75%. These include ultra famous artists like Van Gogh, Claude Monet, Renoir, which we’d expect generative models to do well in emulating. However, a few relatively lesser known artists are also present, like Jacek Yerka, who are still alive, and thus could be negatively affected by generative models reproducing their styles.

With **TagMatch**, in addition to predicting an artistic style, we can also articulate the specific tag signature shared between the test set of images and the reference set of images for the predicted style. Thus, we can inspect the shared signature, as well as instances from both sets where the signature is present, providing direct evidence of the potential style infringement a broader audience to independently verify. Inspecting some examples in figure 9 (more in fig. 15), we observe that while pixel level differences are common across retrieved image subsets, stylistic elements are consistent in both sets with the labeled tags, echoing our motivating claim that style copying goes beyond image or pixel-wise similarity. Lastly, TagMatch also allows for understanding image distributions from the perspective of interpretable tags. We explore this direction in appendix J.2, finding differences in the uniqueness of the tags present in generated art vs real art.



Figure 9: Examples of applying TagMatch to generated images. TagMatch is inherently interpretable with respect to tags, as each inference comes with the exact set of tags that are (i) shared between the sets of test art and art from the predicted artist, and (ii) used to predict the artist.

419 broader community who this problem affects, and serve as an accessible framework to quantitatively  
420 examine the nuanced issue of artistic style infringements.

## 421 F Limitations

422 Our work tackles a novel problem of artistic *style* infringements. Style, however, is qualitative. We  
423 merely put forward one definition for artistic style, along with two implementations for demonstrating  
424 the existence of a style given example works from an artist and recognizing the identified style in  
425 other works.

426 Importantly, we argue that an artist’s style is unique if we can consistently distinguish their work from  
427 that of other artists. However, we can only proxy the entire space of artists. We construct a dataset  
428 consisting of works from 372 artists spanning diverse schools of art and time periods in attempt to  
429 represent the space of existing artists, though of course we will always fall short in capturing all kinds  
430 of art. We provide tools to allow for this dataset to grow with time, and we caution that if only one  
431 artist for some broader artistic style is not present in our reference set, the uniqueness of that artist’s  
432 style may be overestimated, and as such, generated images may be matched to this artist with an  
433 overestimated confidence. However, if only one out of 372 artists exhibits some style, than one could  
434 argue that that alone reflects a notable uniqueness of that artist. To employ a stricter criterion for  
435 alleging style copying, we’d recommend augmenting the reference set to include more artists with  
436 very similar styles to the artist in question. Nonetheless, we believe our reference dataset does well in  
437 representing all art, to where analysis based on this reference set is still informative.

438 We also note that our atomic tagging leverages an existing foundation model (CLIP) with no additional  
439 training. While we verify the precision of our tags, CLIP is known to have issues with complex  
440 concepts. Further, we do not claim our tags achieve perfect recall (most image taggers do not). We  
441 advise users to interpret the assignment of a tag to indicate a strong presence of that concept, relative  
442 to similar concepts (i.e. from the same aspect of artistic style). While our tagger is not perfect, it is  
443 objective and automatic, enabling interpretable style articulation and detection. Also, we note that the  
444 field of image tagging in general has seen rapid improvement in the past year [14], and an improved  
445 tagger could easily be swapped into our pipeline.

446 Lastly, we only analyze generated images using off-the-shelf text-to-image models. It is possible that  
447 particularly determined and AI-adept style thieves fine-tune a model to more closely replicate specific

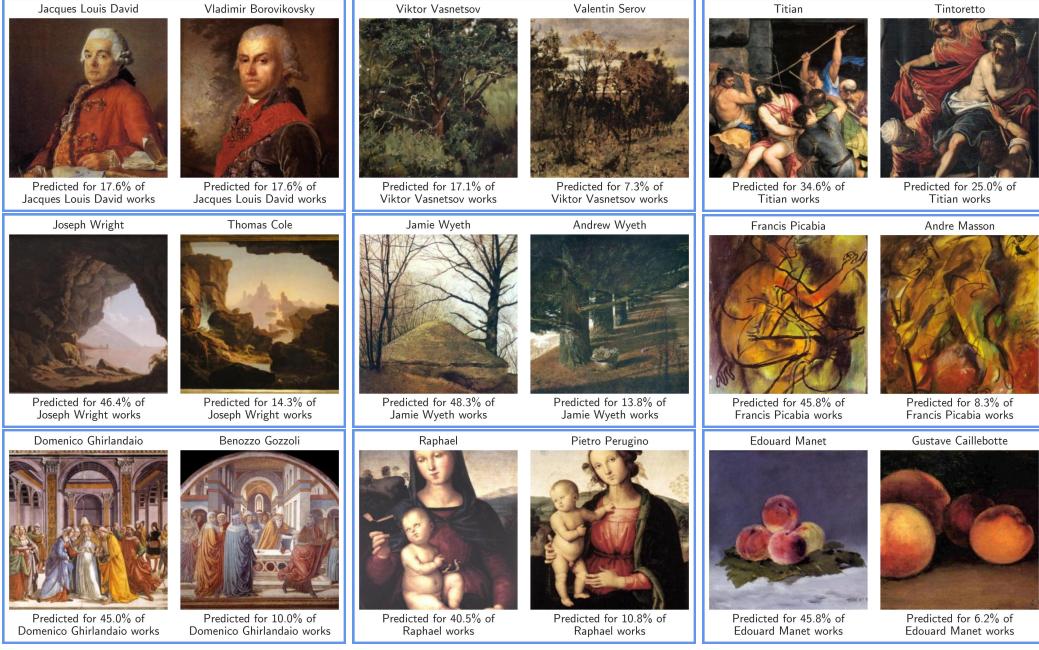


Figure 10: Examples of artists who’s styles were not recognized by DeepMatch (i.e. less than half of their held-out works were predicted to the artist). Each panel shows an example work from (left) the unrecognized artist and (right) the artist that is incorrectly predicted most frequently over works from the unrecognized artist. We see that artists can use very similar, at times arguably indistinguishable, styles.

448 artistic styles. This is a much more threatening scenario, though requires greater effort and ability by  
 449 the style thief. We elect to demonstrate the feasibility of our approach in the more broadly accessible  
 450 setting of using models off-the-shelf, and note that our method can flexibly accept generated images  
 451 produced in a different way (or perhaps discovered on the internet); notice generated images are an  
 452 optional input in figure 7. We look forward to explorations of more threatening scenarios in future  
 453 work, and hope both our formulation and methods for measuring style copying prove to be of use.

## 454 G A nuance in artistic style infringements: Existing Artists can have very 455 similar styles

456 A crucial step in arguing that an artist’s style has been infringed is to first demonstrate the existence  
 457 of the given artist’s *unique* style. We note that doing so objectively is non-trivial, as a style may not  
 458 have a clear definition, and thus, it can be challenging to systematically compare to all other artistic  
 459 styles, so to show uniqueness. In our work, we utilized classification, claiming that if an artist’s works  
 460 can consistently be mapped (i.e. at least half the time) to that artist (over a large set of other artists),  
 461 than that artist must have some underlying unique style (parameterized by a neural signature).

462 In doing so, we found that 89.3% of artists could be recognized based of a set of (at least 20 of) their  
 463 works (held-out in training the classifier). What about the remaining 10.7% of artists? We now take  
 464 a closer look at these artists, and also introduce a second, stricter style copying criterion. Namely,  
 465 we consider the notion that it may be unfair to claim a generative model is copying the style of an  
 466 artist, if another existing artist seems to also be copying that artist. That is, we propose a way to  
 467 verify that the generative model not only shows a substantial similarity to the copied artist, but also  
 468 an *unprecedented* similarity.

### 469 G.1 Artists who’s styles were not recognized

470 First, we inspect more examples from artists who were not recognized using our majority voting  
 471 threshold in DeepMatch. That is, less than half of their held-out works were predicted to them. Figure

472 10 shows a number of examples, from which we can make some qualitative observations. First,  
473 the styles of artists who operate in the same broader genre (e.g. portraiture, landscapes, narrative  
474 scenes in renaissance styles, etc) can be extremely similar. We even see an instance where an artist's  
475 son's style is indistinguishable from his father's (Jamie and Andrew Wyeth). Lastly, we note that in  
476 most cases, the artists only marginally fall short of our recognition threshold (i.e. accuracy for their  
477 held-out works is only a bit below 50%). We utilize majority voting because (i) it is intuitive, (ii) it  
478 requires *consistent* appearance of the neural signature across works, and (iii) it allows for abstention  
479 when no particular style is strongly present. However, the exact threshold of 50% can be altered as  
480 desired. In summary, as in Figure 4, we see artistic styles can be very similar, making the existence  
481 of unique artistic styles for the vast majority of artists a non-trivial observation.

482 If an artist's style cannot be recognized over their own held-out works, arguing that a generative model  
483 copies that style is strenuous, as the style itself is ill-defined. Notably, in these cases, the classifier  
484 had an option to predict the correct artist. However, in applying DeepMatch to generated images,  
485 there is no direct option for the classifier to abstain from predicting anyone, under that generated  
486 art comes from a "new artist", which takes inspiration from existing artists. Note that abstention is  
487 still possible (due to the majority voting in DeepMatch), and occurs when a match confidence falls  
488 below 50%. To make comparisons fairer to generative models, we now discuss a stricter criterion of  
489 *unprecedented similarity*.

490 **G.2 Unprecedented Similarity: Do generative models copy styles more than existing artists  
491 already do?**

492 A nuance that requires consideration when studying artistic style copying is that it is possible for  
493 two artists to have very similar styles. Thus, it may be unfair to allege that a generative model is  
494 copying an artist  $a$  if there exists another artist  $b$  who's style is just as or in fact even more similar to  
495 artist  $a$ . Towards this end, we introduce *unprecedented similarity*, which requires that the similarity  
496 between works of a generative model  $A'$  and works of the artist intended to be copied  $A$  is higher than  
497 the similarity of any existing artist with  $A$ . That is,  $\text{sim}(A, A') \geq \text{sim}(A, B)$  for works  $B$  from all  
498 other existing artists  $b$ .

499 Note that this is a stricter criterion than our previous threshold. In DeepMatch, we required that at  
500 least half of the works in a given set of test images were predicted to a single artist in order for us  
501 to flag the test images as a potential style infringement. In other words, that threshold required that  
502  $\text{sim}(A, A') \geq 0.5$ , which in turn implies that  $\text{sim}(A, A') \geq \text{sim}(A', B)$  for all  $B$  (with room to  
503 spare; here we use match confidence to denote similarity).

504 Now, however, instead of just comparing  $A'$  to all  $B$ , we must also compare all  $B$  to  $A$ . Instead of  
505 comparing all other artists, we inspect the most similar artist  $b^*$  to  $a$ , identified by taking the artist  
506  $b$  with the highest rate of false positive predictions to artist  $a$ . Then, we hold out  $b$ , and train a new  
507 classifier on the remaining 371 artists. Finally, we check for style matches of for the set of generated  
508 images  $A'$  and the works  $B^*$  from the most similar artist  $b^*$ .

509 Figure 11 summarizes our result for OpenJourney (all three models studied show consistent results).  
510 We find that only in three cases do we see a held-out artist's work flagged as potential style copying.  
511 Notably, in all instances where generated work is flagged as potential style copying, the corresponding  
512 held-out artist's work is either not flagged or is flagged with lower confidence, indicating that the  
513 instances of style copying of generative models that we observe always also satisfy the criterion of  
514 unprecedented similarity.

515 Taking a closer look at instances where held-out art is flagged for style copying (or perhaps style  
516 emulation?), we again see just how similar the works of different artists can be. Namely, we see  
517 that some artists works seem to fall into a broader genre of art that many artists utilize (e.g. ukiyo-e  
518 or impressionism). In summary, while generative models can very closely resemble the style of a  
519 given artist, contextualizing copying by generative models with respect to copying (or perhaps, 'style  
520 emulation') already done by existing artists is crucial in order to afford the same artistic liberties to  
521 generative models as have been provided to other artists in the past.

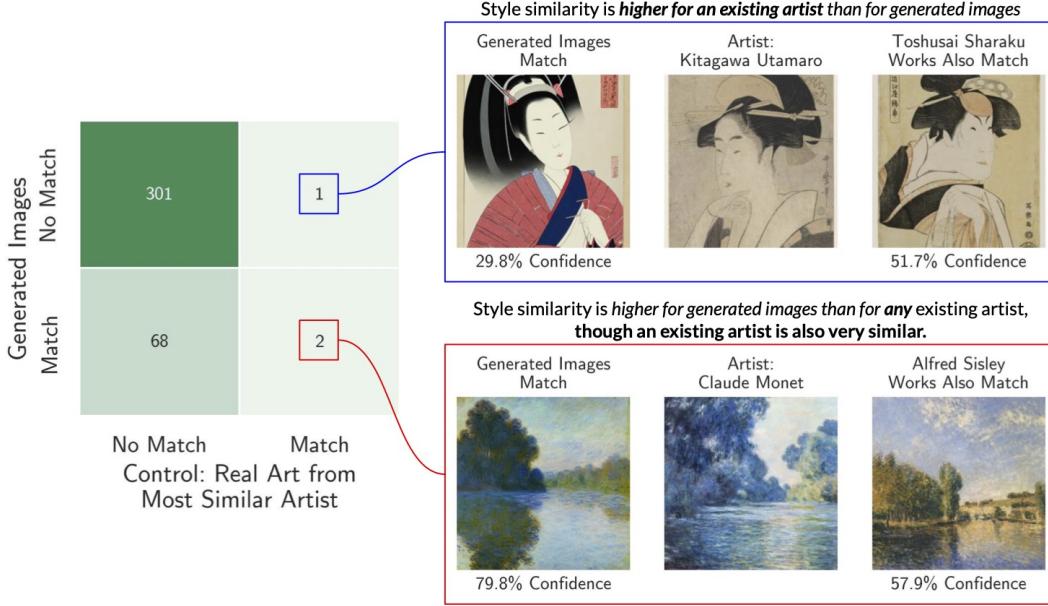


Figure 11: We verify the stricter criterion of *unprecedented similarity* by holding out the real artist with highest similarity to a given artist, and checking if the held-out real artist’s works are flagged as potential style copying by DeepMatch. (**left**) We observe only three artists where the most similar held-out artist has their work flagged as a style match, and in all cases, when generated images are flagged, the match confidence of the generated images exceeds that of the held-out real artist’s works (i.e., **the generated images flagged by our method reflect unprecedented similarity to the given artist’s style**). (**right**) Inspecting the flagged held-out artists further show that style copying is very nuanced, as artists take inspiration from one another, and as such, they may already have very similar styles. While we always observe unprecedented similarity, a potential solution to style copying may be for generative models to ensure that they do not copy any more than what already exists; that is, they may exhibit some copying, but no more than for which precedent already exists.

## 522 H Baselines

523 We now present some alternate implementations to the methods we present, so to serve as base-  
 524 lines. We note that a key contribution of our work is reformulating the problem of detecting style  
 525 infringements from computing image-wise similarity to performing classification over image sets,  
 526 and building a tool around this idea. Thus, it is rather challenging to perform apples-to-apples  
 527 comparisons to prior copy detection works, as our methods implement a different task. We include  
 528 substantial qualitative discussion comparing our approach to image-similarity techniques (and thus  
 529 motivating our framework) in section B, and we add to that discussion here.

530 We further stress that there is not a singular numerical objective that we can use as a way to compare  
 531 methods. For example, we report the accuracy of matching artists (i.e. aggregating classification  
 532 predictions with majority voting), but since it is not necessarily true that all artists are distinguishable,  
 533 it would be imprudent to strictly prefer a higher accuracy, as there is no strict groundtruth; that  
 534 is, there is no completely definitive way to say if an artist has a unique style or not, due to the  
 535 subjective/qualitative nature of style. Nonetheless, for lack of other quantitative metrics, we inspect  
 536 accuracy on real and generated images for a few lightweight approaches to artist classifications, and  
 537 compare them below.

### 538 H.1 DeepMatch

539 Figure 12 shows the performance of different classifiers, where we vary the frozen backbone and the  
 540 number of hidden layers. We find that classifiers trained on CLIP yield higher match-rates for both  
 541 real and generated art than classifiers train on DINOv2 [21] embeddings. Interestingly, zero-shot

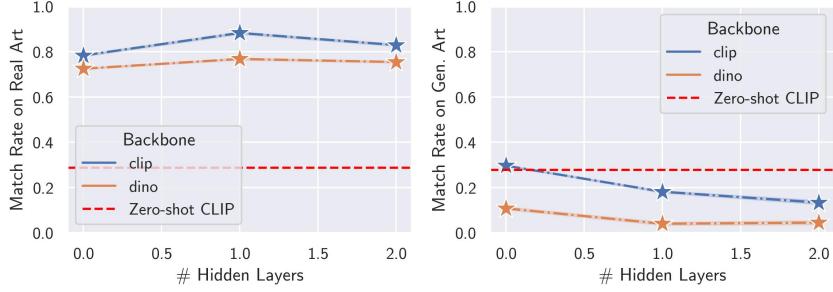


Figure 12: Alternate implementations of DeepMatch, using DINOV2 and CLIP backbones, and varying the number of hidden layers. We also present performance of zero-shot CLIP. Numbers are averaged over five trials, except for zero-shot CLIP, which is deterministic.

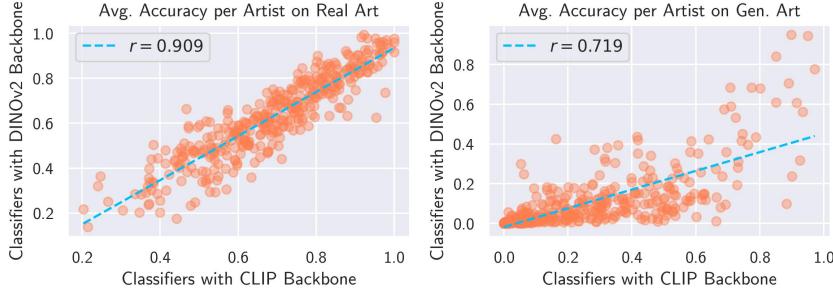


Figure 13: Per-artist accuracy for classifiers using CLIP and DINOV2 backbones are highly correlated. While each classifier may yield different overall accuracy, the *relative* notions of (i) how recognizable the artist’s real art is and (ii) how much so the artist’s style appears in generated works appear to be classifier agnostic.

542 CLIP does poorly on real art, but well on generated art, perhaps because many generative models  
543 optimize using CLIP-score, which applies the same mechanism as zero-shot CLIP classification,  
544 perhaps explaining the assertion that generative models are highly capable of imitating humans found  
545 in this brief work [8]. The number of hidden layers does not have a very strong affect on recognizing  
546 real art, but it does appear inversely related to the ability of the model to recognize generated art. It is  
547 possible that having two many hidden layers can overfit the model to the distribution of real images,  
548 creating a distribution shift when applied on the generated images.

549 While exact numbers seem to vary, we note that relative trends (i.e. between artists) appear agnostic  
550 to the underlying classifier. Figure 13 shows accuracy per artist for classifiers trained on CLIP vs  
551 DINOV2 embeddings. For both real and generated art, the per-artist accuracies are strongly correlated,  
552 which could motivate using relative metrics in addition to absolute values dependent on exact accuracy  
553 values; note that we include relative numbers in our ArtSavant report (see Figure 1; e.g., ‘percentile  
554 of recognizability’).

555 We ultimately choose something in the middle of the round: a 1-hidden layer MLP on CLIP  
556 embeddings, which has the strongest performance recognizing real art, and appears to have some  
557 ability to recognize generated art. We note the majority aggregation that we apply is just one way  
558 to summarize the classification output across an image set. We opt for it because it is intuitive  
559 and it provides a natural avenue for abstention, though this threshold can be modified as desired,  
560 and inspecting relative accuracies could be most informative. We again stress that our current  
561 implementation serves as a proof of concept of our framework, which is our primary contribution.

## 562 H.2 TagMatch

563 We now present baselines for TagMatch. Like above, and indeed more so, accuracy is not exactly an  
564 objective to maximize. In fact, what is most important with TagMatch is interpretability, and ease with

	CBM	CBM + sparsity	Ours
Accuracy on real art	62.8%	58.7%	61.5%

Table 1: Baselines for TagMatch

565 which the output of TagMatch can be used in arguments to a broader, non-technical audience. Thus,  
 566 we consider a popular framework from the interpretable classification literature: concept bottleneck  
 567 models (CBM) [17]. Namely, we train a linear layer atop concept predictions extracted from CLIP,  
 568 so to create a CBM without direct concept supervision, as in [19, 37].

569 As shown in 1, accuracy values are roughly similar. We note the interpretability provided by the  
 570 methods are markedly different. CBM allows one to inspect the final linear layer to discern which  
 571 concepts are important to which class, but this results in requiring users to inspect a coefficient for  
 572 every concept. Adding sparsity by way of an  $\ell_1$  penalty can help, but the problem persists. Our  
 573 version of TagMatch, on the other hand, affords concise articulations of tag signatures, as well as a  
 574 number of how many other artists share a given signature. Perhaps most crucially, our implementation  
 575 also yields faithful attribution, which can be critical in gathering evidence to present to a judge or  
 576 jury.

### 577 H.3 Stability

578 We also explore the stability of our method to using different data splits. We perform five different  
 579 random train / test splits, and inspect the accuracy of our implementations of DeepMatch and  
 580 TagMatch. DeepMatch per-image accuracies are very stable, with a standard deviation of 0.1%.  
 581 TagMatch is also stable, though less so, with a standard deviation 1.1%.

## 582 I On alternate prompts

583 We briefly explore using alternate prompts to generate images. Namely, we create 120 prompts of  
 584 the form “{an object} in {location} in the style of {artist}” (e.g. “A bottle in forest in the style of  
 585 Jeff Koons”, which are by nature no longer artist-specific (like the titles we originally use). Using  
 586 DeepMatch, average match rate drops considerably in this less specific case, from 20% to 8%. This  
 587 is in line with existing wisdom that prompting can significantly affect the behavior of a model, and  
 588 also echoes our overall empirical observation that current style copying does not appear to be very  
 589 prevalent. We hope that our framework can be useful in examining which prompts induce greatest  
 590 copying going forward, especially as prompt and model sophistication grows.

## 591 J Details on TagMatch

592 We now provide greater details regarding the implementation of TagMatch, a central technical  
 593 contribution of our work. TagMatch is a method to classify a set of images to a class; specifically,  
 594 we map a set of artworks to one artist, selected over 372 choices. TagMatch is not as accurate as  
 595 DeepMatch, as it maps held-out works of each artist in our WikiArt dataset to the correct artist about  
 596 61% of the time (compared to 89% top-1 accuracy for DeepMatch). However, top-5 accuracy is  
 597 more reasonable, achieving above 80%. Most notably, **TagMatch is inherently interpretable and**  
 598 **attributable**. It consists of three steps: (i) assigning atomic tags to images, (ii) efficiently composing  
 599 tags to obtain more unique *tag signatures*, and (iii) matching a test set of images to a reference  
 600 artist based on the uniqueness of the tags shared between the test set and works from the predicted  
 601 reference artist.

602 Our method is fast and flexible: after caching image embeddings, the whole thing only takes minutes,  
 603 and it is easy to modify the concept vocabulary as desired, as the tagging is done in a zero-shot  
 604 manner. Through MTurk studies, we verify that the atomic tags we assign are mostly precise, though  
 605 we recognize that these descriptors can be subjective. Thus, while we do not claim perfect tagging,  
 606 we stress that our method is easy to understand, and crucially, is deterministic per image. Therefore,  
 607 ideally our tagging may be more reliable than human judgements, particularly when the

608 humans involved may be biased (e.g. an artist alleging copying and a lawyer defending a generative  
609 model would have strong and opposing stakes).

610 Below, we provide details for image tagging (§J.1), artist tagging (§J.2), artistic style inference via  
611 tag matching (§J.3), effect of hyperparameters (§J.4), details on efficiency (§J.5), and a review of  
612 validation (§J.6).

## 613 J.1 Image Tagging

614 As explained in §C.2, we utilize CLIP to attain a diverse set of atomic tags per image in a zero-shot  
615 manner. Specifically, we first define a vocabulary of descriptors along various aspects of artistic  
616 style. Then, given an image, we do selective multi-label zero-shot classification *for each aspect*.  
617 Performing zero-shot classification per aspect proves to be critical in order to achieve a diversity of  
618 tags and a similar number of tags per image. We find that some descriptors always lead to higher  
619 CLIP similarities than others. Specifically, descriptors for simple aspects, like colors and shapes,  
620 yield higher similarities than more complex aspects like brushwork and style. Thus, using a global  
621 threshold across descriptors would lead to a less diverse descriptor set. Moreover, we observe  
622 some images have higher similarities across the board than others, which again would lead global  
623 thresholding to result in a disparate number of tags per image. Our per-aspect scheme requires that  
624 the descriptors within each aspect are mostly mutually exclusive; we prioritize this in the construction  
625 of the concept vocabulary, via the prompt we present the LLM assistants and our manual verification.

626 Namely, we prompt both Vicuna-33b and ChatGPT with “*I want to build a vocabulary of tags to be  
627 able to describe art. First, consider different aspects of art, and then for each aspect, list about 20  
628 distinct descriptors that could describe that aspect of art. Please return your answer in the form of a  
629 python dictionary.*”. We then perform a filtering step with a human in the loop, where we manually  
630 remove tags that are difficult to recognize or redundant. After this filtering step, we add in a few new  
631 aspects. First, we incorporate the 20 *styles* (e.g., “impressionism”) and *genres* (e.g., “portrait”) that  
632 are most common amongst works in our WikiArt dataset; note that all WikiArt images also contain  
633 metadata for these categories. Finally, we add some easy to understand tags such as *color* and *shape*  
634 which can be important characteristics describing a given painting. The concept vocabulary we use is  
635 contains shown below:

- 636 • **Style**, caption template: *{}* *style*. Descriptors:
  - 637 – *realism, impressionism, romanticism, expressionism, post impressionism, art nouveau  
638 modern, baroque, symbolism, surrealism, neoclassicism, naïve art primitivism, north-  
639 ern renaissance, rococo, cubism, ukiyo e, abstract expressionism, mannerism late  
640 renaissance, high renaissance, magic realism, neo impressionism*
- 641 • **Genre**, caption template: *the genre of {}*. Descriptors:
  - 642 – *portrait, landscape, genre painting, religious painting, cityscape, sketch and study,  
643 illustration, abstract art, figurative, nude painting, design, still life, symbolic painting,  
644 marina, mythological painting, flower painting, self portrait, animal painting, photo,  
645 history painting, digital art*
- 646 • **Colors**, caption template: *{}* *colors*. Descriptors:
  - 647 – *pale red, pale blue, pale green, pale brown, pale yellow, pale purple, pale gray, black  
648 and white, dark red, dark blue, dark green, dark brown, dark yellow, dark purple, dark  
649 gray*
- 650 • **Shapes**, caption template: *{}*. Descriptors:
  - 651 – *circles, squares, straight lines, rectangles, triangles, curves, sharp angles, curved an-  
652 gles, cubes, spheres, cylinders, diagonal lines, spirals, swirling lines, radial symmetry,  
653 grid patterns*
- 654 • **Common Objects**, caption template: *{}*. Descriptors:
  - 655 – *male figures, female figures, children, farm animals, pet animals, wild animals, geo-  
656 metric shapes, fruit, vegetables, instruments, flowers, boats, waves, roads, household  
657 items, the moon, the sun, saints, angels, demons*
- 658 • **Backgrounds**, caption template: *{}* *in the background*. Descriptors:

- 659        – *fields, blue sky, night sky, sunset or sunrise, forest, rolling hills, simple colors, beach,*  
 660        *port, river, starry night, clouds, shadows, living room, bedroom, trees, buildings,*  
 661        *chapels, heaven, hell, houses, streets*
- 662     • **Color Palette**, caption template: *{}* color palette. Descriptors:
    - 663        – *vibrant, muted, monochromatic, complementary, pastel, bright, dull, earthy, bold,*  
 664        *subdued, rich, simple, complex, varying, minimal, contrasting*
  - 665     • **Medium**, caption template: *the medium of {}*. Descriptors:
    - 666        – *oil painting, watercolor, acrylic, ink, pencil, charcoal, etching, screen printing, relief,*  
 667        *intaglio, collage, montage, photography, sculpture, ceramics, glass*
  - 668     • **Cultural Influence**, caption template: *{}* influences. Descriptors:
    - 669        – *Indigenous, European, American, East Asian, Indian, Middle Eastern, Hispanic, Aztec,*  
 670        *Contemporary, Greek, Roman, Byzantine, Russian, African, Egyptian, Tahitian, Polynesian, Dutch*
  - 672     • **Texture**, caption template: *{}* texture. Descriptors:
    - 673        – *rough, smooth, bumpy, glossy, matte, roughened, polished, textured, smoothed, brush-*  
 674        *stroked, layered, scraped, glazed, streaked, blended, uneven, smudged*
  - 675     • **Other Elements**, caption template: *{}*. Descriptors:
    - 676        – *stippled brushwork, chiaroscuro lighting, pointillist brushwork, multimedia compo-*  
 677        *sition, impasto technique, repetitive, pop culture references, written words, chinese*  
 678        *characters, japanese characters*

679 Now, we detail the implementation of our modified zero-shot classification. Recall that in zero-shot  
 680 classification, one computes a text embedding per class, which amounts to the classification head,  
 681 and computes an image embedding for the test input, so that the prediction is the class who’s text  
 682 embedding has the highest cosine similarity to the test image embedding. In computing the text  
 683 embeddings, we take each descriptor (e.g. *Dutch*) and place it an aspect-specific caption template (e.g.  
 684 *Dutch → Dutch influences*), and then average embedddings over multiple prompts (e.g. “artwork  
 685 containing *Dutch influences*”, “a piece of art with *Dutch influences*”, etc), as done in [24]. We  
 686 modify standard zero-shot classification to allow for the fact that more than one descriptor (or perhaps  
 687 none) from a given aspect may be present. Namely, instead of assigning the most similar descriptor  
 688 per-aspect, we assign an atomic tag for any descriptor who’s similarity is significantly higher than  
 689 other descriptors for that aspect. We achieve this via z-score thresholding: per-aspect, we convert  
 690 similarities to z-scores by subtracting away the mean and dividing by the standard deviation, and then  
 691 admit atomic tags who’s z-score is at least 1.5.

692 The template prompts we utilize for embedding each concept caption are as follows:

- 693        • art with
- 694        • a painting with
- 695        • an image of art with
- 696        • artwork containing
- 697        • a piece of art with
- 698        • artwork that has
- 699        • a work of art with
- 700        • famous art that has
- 701        • a cropped image of art with

## 702 J.2 From Image Tags to *unique* Artist Tags

703 Recall that we define styles not per-image, but over a set of images. Namely, we seek to surface  
 704 tags that occur frequently. The best way to do so is to simply count the occurrences of each tag, and  
 705 discard the ones that rarely appear. However, each atomic tag is not particularly unique with respect  
 706 to artists. We utilized *efficient composition* of atomic tags to arrive at more unique tag signatures, as

---

**Algorithm 1** Iterative Algorithm to Obtain Tag Composition Per Artist  $a \in \mathcal{A}$ 

---

**Require:**  $\mathcal{D}_a$  (Images for artist  $a$ ),  $\mathcal{C}_a$  (Common tags for artist  $a$ )

```

 $\mathcal{S}_a = \{\}$                                  $\triangleright$  Stores the tag compositions with their associated counts
for  $x \in \mathcal{D}_a$  do
     $I(x) = \text{tag}(x) \cap \mathcal{C}_a$            $\triangleright$  Compute the intersection with common atomic tags
     $\mathcal{P}(I(x)) = \text{ComputePowerSet}(I(x))$        $\triangleright$  Compute power-set of the tags
     $\text{UpdateCount}(\mathcal{S}_a, \mathcal{P}(I(x)))$          $\triangleright$  Update the count of each tag composition
end for
 $\text{Filter}(\mathcal{S}_a)$                           $\triangleright$  Keep tag compositions which occur above a count threshold of 3
```

---

707 shown in figure 6 and detailed in algorithm 1. Importantly, we utilize a threshold here to differentiate  
708 what a common tag is; we require a tag to appear in at least three works for an artist in order for the  
709 tag to count as a frequently used tag by the artist. We note that tag composition can be done efficiently  
710 because we have a relatively low number of tags per image: on average, there are 6.2 atomic tags  
711 per image. Moreover, because the number of occurrences for a composed tag is bound below by the  
712 number of occurrences of each atomic tag in the composition, we can ignore all non-frequent atomic  
713 tags. Thus, we can iterate over the powerset of common atomic tags per image without it taking  
714 exorbitantly long. We include one fail safe, which is that in the rare instance where an image has a  
715 very high number of common atomic tags, we truncate the tag list to include only 25 tags. Over the  
716 91k images that we encounter, this happens only once. We highlight that our tag composition takes  
717 inspiration from [26].

718 **J.3 Predicting Artistic Styles based on Matched Tags**

719 Once we have converted tags per image to tags per artist, we can then utilize these artist tags to perform  
720 inference over a set of images. Namely, given a test set of images, we extract common tags (including  
721 tag compositions) for the test set and compare them to tags extracted for each artist in our reference  
722 corpus. Then, we predict the reference artist who shares the most unique tags with the test set.

723 Figure 14 best explains our method, as it shows the documented code. We note that all code will be  
724 released upon acceptance. We’ll now explain it step by step. First, for each artist and for the test set of  
725 images, we find common tags via (i) assigning atomic tags to each image, (ii) finding the commonly  
726 occurring atomic tags, (iii) counting compositions of the commonly occurring atomic tags, and (iv) dis-  
727 carding tags (including compositions) that do not occur frequently enough. The code shows this done  
728 for the test set of images; we perform this per reference artist when the `TagMatcher` object (for which  
729 `tag_match` is function) is initialized; notice fields like `self.ref_tags_w_counts_by_artist`,  
730 which contain useful information about the reference artists, computed once and re-used for each  
731 inference.

732 Then, we loop through the set of ‘matched’ tags (i.e. those that occur for both the test set of images  
733 and at least one reference artist), starting with the most unique ones. Here, uniqueness refers to the  
734 number of reference artists that frequently use a tag. For each tag, we loop through all artists that also  
735 use that tag. For the first  $k$  (denoted by `self.matches_per_artist_to_consider` in the code)  
736 matched tags per artist, we add a score to a list of scores for the artist, which ultimately are averaged.  
737 The score contains an integer and a decimal component. The integer component is the number of  
738 reference artists that share the matched tag. The decimal component is the absolute value of the  
739 difference in frequency with which the tag appears, over the reference artist’s works and the test set  
740 of images; note that this is always less than one. This way, when comparing two matched tags, a  
741 lower score is assigned to a more unique one, and if there is a tie in uniqueness, we break the tie  
742 based on how similar the frequency of the matched tag is for the test artist and reference artist.

743 Finally, we average the list of scores per artist to get a single score per reference  
744 artist, analogous to a logit. We assign a score of `inf` for any artist with less than  
745 `self.matches_per_artist_to_consider` (which we set to 10) matched tags. This hyperpa-  
746 rameter makes our tag matching less sensitive to individual matched tags, and empirically results in a  
747 substantial improvement in top-1 accuracy on held-out art from WikiArt artists (see next section).

```

def tag_match(self, test_img_paths: List[str], test_artist: str):
    dset = BasicDsetFromImgPaths(test_img_paths, self.vlm.transform, dsetname=test_artist)

    tags_by_path = self.tag_images(dset)
    common_tags = self.find_common_tags(tags_by_path)
    composed_tags_w_counts = self.compose_tags(common_tags, tags_by_path)

    # Now we cross-reference the found tags w/ tags for reference artist
    counts_over_ref_artists_by_tag = dict({
        t:len(self.ref_artists_by_tag[t])
        for t in composed_tags_w_counts if t in self.ref_artists_by_tag
    })
    # We sort the tags by uniqueness: we first inspect tags that occur for the lowest number of reference artists
    counts_over_ref_artists_by_tag = dict(sorted(counts_over_ref_artists_by_tag.items(), key=lambda x:x[1]))

    # We will return a score per artist to resemble the typical output of a classifier
    scores_by_artist = dict((artist: [] for artist in self.ref_dset.artists))
    # We will also keep track of the tags used in computing the score per artist -- this provides faithful interpretations
    matched_tags_by_artist = dict((artist: [] for artist in self.ref_dset.artists))
    # Now we loop through each tag that also occurs for reference artists
    for t, num_ref_artists_w_tag in counts_over_ref_artists_by_tag.items():
        # For each tag, we loop through all matches (i.e. any reference artist that also has the tag)
        for ref_artist in self.ref_artists_by_tag[t]:
            # We only consider the top k most unique matched tags per artist (k = self.matches_per_artist_to_consider)
            if len(scores_by_artist[ref_artist]) < self.matches_per_artist_to_consider:
                # Compute frequency of matched tag over works from the reference artist
                num_works_of_ref_artist_w_tag = self.ref_tags_w_counts_by_artist[ref_artist][t]
                freq_for_ref_artist = num_works_of_ref_artist_w_tag / self.num_works_by_ref_artist[ref_artist]
                # Compute frequency of matched tag over works from the test artist
                freq_for_test_artist = composed_tags_w_counts[t] / len(tags_by_path)
                # Our score is the uniqueness of the matched tag + |diff in frequencies of tag for ref artist and test artist|
                scores_by_artist[ref_artist].append(num_ref_artists_w_tag + np.abs(freq_for_ref_artist - freq_for_test_artist))
                matched_tags_by_artist[ref_artist].append(t)

    # We set the score to inf for any artists that did not have enough matched tags
    scores = np.array([np.mean(scores_by_artist[artist][:self.matches_per_artist_to_consider])
        if len(scores_by_artist[artist]) >= self.matches_per_artist_to_consider else np.inf for artist in self.ref_dset.artists])

    # Finally, we return scores along with explanations for each artist
    return scores, matched_tags_by_artist

```

Figure 14: Code for predicting artistic styles via matched tags.

#### 748 J.4 Choosing Hyperparameters

749 Overall, there are three hyperparameters to our method: the z-score threshold, the tag count threshold,  
 750 and the number of matches to consider per artist. Here is quick refresher on what they each do:

- 751 • The z-score threshold determines how much more similar a descriptor needs to be to an  
 752 image compared to other descriptors for the same aspect in order for the descriptor to be  
 753 assigned as an atomic tag of the image. The value we use is 1.75.
- 754 • The tag count threshold is the minimum number of an artist's works that a tag needs to be  
 755 present in order for a the tag to be deemed common for the artist. The value we use is 3.
- 756 • The number of matches to consider per artist pertains to how many matched tags are  
 757 considered when computing the final score per artist in tag match. That is, the final score for  
 758 an artist is the average of the top-k most unique tags that the artist shares with the test set of  
 759 images, where  $k$  corresponds to this hyperparameter. The value we use is 10.

760 Now that the role of each hyperparameter is clear, let's discuss how hyperparameters can be adjusted  
 761 towards particular ends, along with the potential consequence of each action:

- 762 • To increase the number of atomic tags, lower the z-score threshold. Risk: atomic tags may be  
 763 less precise, and the method will take longer to run, as there will atomic tags and composed  
 764 tags.
- 765 • To get more tags per artist, lower the tag count threshold. Risk: some tags will become  
 766 less unique. Other tags will be introduced, and may be very unique, which could skew tag  
 767 matching. Also, the method may take longer to run, as there will be more tags.
- 768 • To make inference less sensitive to a low number of matched tags, increase the number of  
 769 matches to consider per artist. Risk: when you consider more matches, interpretation is a

Matched Tag for Antoine Blanchard: simple colors, streets, Contemporary influences, social symbolism  
0 other artists also have this signature



8 generated images with this signature



5 real images with this signature

Matched Tag for Franz Xaver Winterhalter: broad brushwork, female figures, historical symbolism  
0 other artists also have this signature



19 generated images with this signature



7 real images with this signature

]

Matched Tag for Arthur Rackham: illustration, children, fantastical subject matter  
0 other artists also have this signature

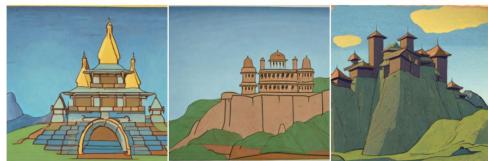


12 generated images with this signature



5 real images with this signature

Matched Tag for Nicholas Roerich: geometric shapes, simple colors, geographical symbolism  
0 other artists also have this signature



47 generated images with this signature



22 real images with this signature

Figure 15: Additional examples of applying TagMatch to generated images.

770 little more difficult, as you have more reasons for each inference, and it will take longer to  
771 view them all.

772 To choose hyperparameters, we selected a small range of reasonable values and swept each hyperpa-  
773 rameter individually. While a combined search would likely yield better accuracy numbers, we opt  
774 out of hyper-tuning TagMatch for accuracy, as its main objective is to provide and interpretable and  
775 attributable complement to DeepMatch. We find the (relatively strong, considering the high number of  
776 artists considered) accuracy numbers encouraging, but do not find it a priority, as DeepMatch arguably  
777 provides a stronger and easier to understand signal of *if* style copying is happening. TagMatch, on  
778 the other hand, tells us *how* and *where* it is happening (if observed with DeepMatch).

779 We also include a hyperparameter sweep, of the z-score threshold and tag count threshold jointly,  
780 and of the number of matches to consider separately afterwards. Figure 16 visualizes the results.  
781 Choosing a lower z-score threshold results in higher TagMatch accuracies. However, a lower z-score  
782 threshold would admit a greater number of false positive tags, and also incurs a longer time of  
783 computation, as there are more tags to compose (we empirically observe an increase of about 50%

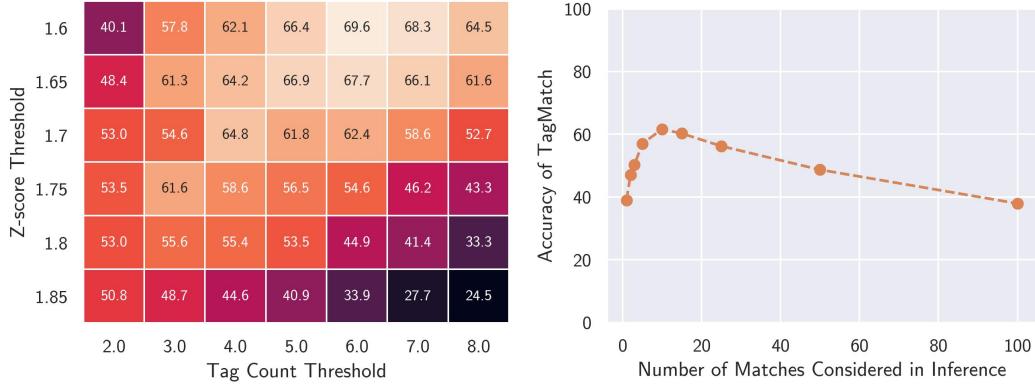


Figure 16: Sweep of hyperparameters associated with TagMatch. **(left)** We jointly sweep the z-score threshold and the tag count threshold. **(right)** Having fixed the first two parameters, we sweep the last one: the number of matches considered in inference. See detailed discussion in §J.4.

784 in run time using our 372 artist reference corpus). Increasing the tag count threshold can reduce  
 785 the time of computation and also increase sensitivity to false positive tags (on individual images),  
 786 resulting in higher TagMatch accuracies. Interestingly, considering more matches improves accuracy  
 787 considerably, but eventually saturates and reduces accuracy. Essentially, by considering more matches  
 788 per artist, inference becomes less sensitive to the most unique matched tag between the artist and the  
 789 test set. The smoothed predictions are more accurate up to a point (i.e. 10 matches), but then hinder  
 790 accuracy. Also, choosing too high a number here can make faithful interpretation more cumbersome,  
 791 as there are more matches to inspect afterwards.

792 We reiterate that the main goal of TagMatch is not to be super accurate, but to complement DeepMatch  
 793 with interpretations (via matched tag signatures) and attributions (via works from the test set and  
 794 from the reference artist that present the matched tags). We ultimately first choose a high z-score  
 795 threshold of 1.75, as a preliminary check revealed this threshold to have considerably higher precision  
 796 in its atomic tags (which we validate with a human study), and since it speeds up the analysis. Then,  
 797 we choose the best tag count threshold (3) and number of matches to consider (10), in that order. We  
 798 hope our discussion of the impact of each hyperparameter can enable practitioners to modify these  
 799 choices as they please. Furthermore, as base VLMs and tagging methods improve, our framework  
 800 can modularly swap out our zero-shot tagging (and thus also the z-score threshold) for a stronger  
 801 method, while retaining the other structure of TagMatch.

## 802 J.5 Efficiency of TagMatch: Runs in roughly 1 minute

803 TagMatch is surprisingly fast. The longest step by far is computing CLIP embeddings for the reference  
 804 artworks. This takes us about 5 minutes using one rtx2080 GPU with four CPU cores to embed the  
 805 73k training split images using a CLIP ViT-B\16 model. Importantly, this step is done only once,  
 806 and in practice, is done offline. The other steps and approximate time needed for each are as follows:  
 807 embedding concepts (5 seconds), extracting common atomic tags and composing them (45 seconds),  
 808 reorganizing tags and removing non-common tags (3 seconds). Then, inference for a test set of  
 809 100 – 200 works takes about 10 to 15 seconds. Again, we will release all code upon acceptance,  
 810 as we truly hope our tool can be of use to artists who are concerned by generative models potential  
 811 infringing upon their unique styles.

## 812 J.6 Validation

813 Because tag match has multiple steps, we perform multiple validations. First, for image tagging,  
 814 we utilize an MTurk study. We collect 3000 separate human judgements on instances of assigned  
 815 atomic tags. Namely, we show 1000 randomly selected (tag, image) pairs to three annotators each.  
 816 Figure 17 shows an example of the form presented to MTurk workers. MTurkers provide consent  
 817 and are awarded \$0.15 per task, resulting in an estimated hourly pay of \$12 – \$18. For each task,  
 818 they answer ‘yes’, ‘no’, or ‘unsure’ to the question ‘does the term {atomic tag} match the artwork

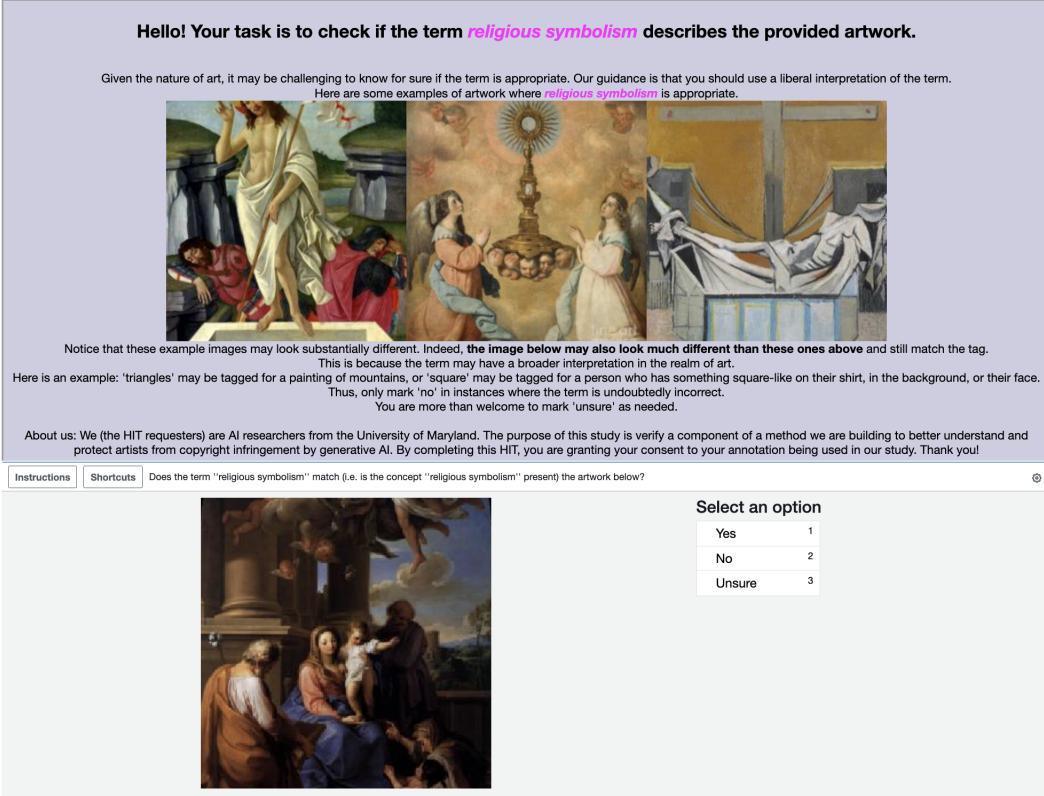


Figure 17: Instructions showed to MTurk workers to validate atomic tags.

819 below?' They are also shown example artworks for each term which were manually verified to be  
 820 correct. Response rates were as follows: 69.89% yes, 8.99% unsure, 21.12% no. In investigating  
 821 inter-annotator agreement, we find that at least 2 annotators agree 92.1% of the time, but all 3 agree  
 822 only 51.52% of the time. This reflects the subjectivity associated with assigning artistic tags, and  
 823 partially motivates the need for a deterministic automated alternative, in order to objectively tag  
 824 images at scale. All three annotators said no only 5.16% of the time, and at least two said no 17.11%  
 825 of the time, suggesting that our zero-shot tagging mechanism achieves reasonable precision.

826 To validate the value of tag composition, we refer to figure 6, which shows how tags become more  
 827 unique as they get longer (i.e. consist of more atomic tags). Moreover, our time analyses show that  
 828 the added benefit of composing tags to find unique tag signatures does not come at the cost of the  
 829 efficiency of our method. Finally, the non-trivial top-1 matching accuracy and strong top-5 matching  
 830 accuracy shows that the extracted tag signatures do indeed capture some unique properties of artistic  
 831 style. Figure 15 reflects a few more examples of successful inference, interpretation, and attribution  
 832 for the task of detecting style copying by generative models.

## 833 K A Sim2Real Gap in Tag Distributions

834 An added advantage of ascribing tags to images is that we can better compare image distributions  
 835 from an interpretable basis (the tags). We briefly explore this direction now.

836 First, we provide complete results from applying TagMatch to generated images from each of the  
 837 three text-to-image models in our study, presented in table 2. Consistent with our DeepMatch results,  
 838 we observe substantially lower matching accuracy for generated images than for real held-out artwork.  
 839 While the primary takeaway is that for many artists, generative models struggle to replicate their  
 840 styles, we can also hypothesize that generative models may output images that follow a different  
 841 distribution than the distribution of real artworks.

			Top 1	Top 5	Top 10
Generated Art	CompVis Stable Diffusion v1.4		10.10	35.49	49.74
	Stability AI Stable Diffusion v2		12.95	37.82	52.59
	PromptHero Openjourney		6.99	31.87	45.08
	Average		10.02	35.06	49.14
Real Art (held out)			61.56	82.53	88.44

Table 2: Match rates using TagMatch for three generative models, as well as on real held out art.

Motivated by this hypothesis, we now compare the distribution of real to generated artworks from the perspective of tags. Because we consider composed tags, the total space of tags is vast and hard to reason over. However, we can look at properties of each tags. Namely, we can inspect the uniqueness of tags. That is, for each tag present in generated images, we inspect the number of reference artists that also present that tag; we do the same for real art as well (subtracting one so to not double count the artist for which a given a tag is being considered). Figure 18 shows a kernel density estimation plot of the distributions of tag commonality, where a tag commonality of 5 means that for each tag assigned to a set of images (either from a real artist or from a generative model emulating an artist), 5 other artists also commonly use that tag. We see tags tend to be rather unique (due to our tag composition), and notably, tags for generated images are more unique.

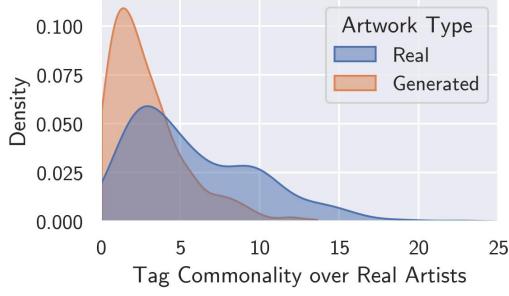


Figure 18: The tags for generated images are less common compared to tags in real art.

## L Patch Match: Generating Additional Visual Evidence of Copying

Detecting artistic style copying in a given art requires analyzing local stylistic elements that manifest across an artist’s body of work. To address this, we employ a patch-based approach that compares small image regions between a given art and original artworks, enabling a fine-grained analysis of stylistic and semantic (e.g. objects) similarities at a local level. We consider three patch matching methods: CLIP-based, DINO-based, and Gram matrix-based.

**Gram Matrix-based Patch Matching [12]:** The Gram matrix is a measure of style similarity introduced in the context of neural style transfer. It captures the correlations between the activations of different feature maps in a convolutional neural network, representing the style of an image. For patch matching, the Gram matrices of patches from the given art and original arts can be computed and compared using a suitable distance metric (e.g., Frobenius norm). The Gram matrix is specifically designed to capture stylistic elements, making it well-suited for detecting style copying.

**CLIP-based Patch Matching [24]:** CLIP (Contrastive Language-Image Pre-training) is a powerful model that can effectively capture the semantic similarity between text and images. In the context of patch matching, CLIP embeddings can be used to measure the similarity between a patch from a given art and patches from original artworks. The patches can be encoded using the CLIP image encoder, and the cosine similarity between their embeddings can be computed to find the closest matches. CLIP may not be as sensitive to low-level stylistic elements, such as brushstrokes, textures, and color palettes, however it focuses more on higher-level semantic concepts, which can be useful to find if the given art pictured the same objects as the selected original patch.

**DINO-based Patch Matching [7]:** DINO is a self-supervised vision transformer that learns robust visual representations by solving a self-distillation task. DINO embeddings can be used for patch matching by computing the cosine similarity between the embeddings of patches from the given art and original artworks. We use DINO to capture higher semantical similarities, and check whether

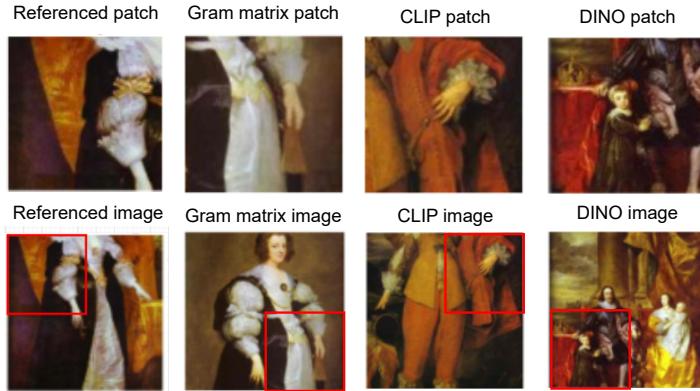


Figure 19: The most similar patches to a referenced patch in an image using Gram-matrix, CLIP, and DINO.

884 the given art pictured similar subjects of interest and high-level visual features as selected original  
 885 artworks.

### 886 **L.1 Experimental setting**

887 For our experiments, we aim to identify the most similar artwork from a pool of 10,000 original  
 888 artworks in the WikiArt dataset given a reference image. The reference image is first resized to a  
 889 resolution of  $512 \times 512$  pixels and normalized. From this normalized image, we select a patch size of  
 890  $128 \times 128$  pixels. This process is repeated for all original artworks in the dataset, resulting in a total  
 891 of 40,000 patches from original artworks for comparison with the reference patch. We then use three  
 892 methods, namely Gram matrix, CLIP, and DINO, to find the most similar patches.

893 Figure 19 showcases the patches that are deemed most similar to the image being referenced. These  
 894 matches are determined using Gram-matrix, CLIP, and DINO methods.

895 We then select an artist and find patches from our original image dataset that closely match this  
 896 artist's style. In Figure 20, we utilize the Gram-matrix method to identify the most similar patches  
 897 to three chosen artworks by Van Gogh. Our dataset includes all paintings by Van Gogh as well as  
 898 works by nine other artists. Gram-matrix selects original artworks that closely resemble the style  
 899 of the reference image, all of which are from Van Gogh. Essentially, this means that Gram-matrix  
 900 predominantly selects Van Gogh's artworks because they are the most stylistically similar to the  
 901 referenced paintings compared to the works of the other nine artists.

### 902 **L.2 Discussion and limitations**

903 Patch matching methods like Gram-matrix, CLIP, and DINO are effective in detecting similarities  
 904 between artworks by examining their local stylistic and semantic elements. Gram-matrix focuses  
 905 on capturing stylistic correlations, CLIP evaluates semantic similarity, and DINO concentrates on  
 906 higher-level features. However, these methods have limitations. They primarily focus on local  
 907 aspects of artworks and may overlook broader artistic characteristics such as texture, composition,  
 908 and brushwork that are crucial to detect copyright infringements. Moreover, the process of finding  
 909 the most similar patches for each given art takes approximately fifteen minutes when considering  
 910 10,000 original artworks, and if we opt to include more original artworks, the duration of the process  
 911 would inevitably increase. Therefore, patch-matching methods are computationally expensive,  
 912 which restricts their practical application. Despite these limitations, patch matching is valuable  
 913 for identifying instances of direct copying in artworks and they aid in the detection of plagiarized  
 914 content.

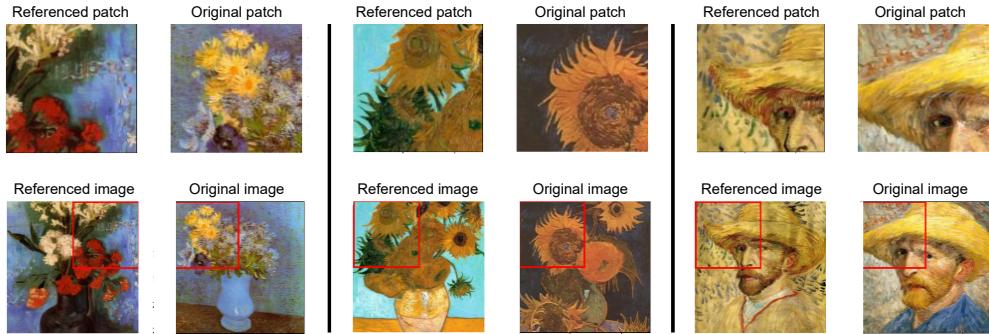


Figure 20: Comparison of patches using the Gram-matrix method, highlighting the closest matches to three selected artworks by Van Gogh. The selected original arts, all from Van Gogh, closely resemble the style of the referenced paintings.

## 915 M Details on WikiArt Scraping

916 WikiArt is a free project intended to collect art from various institutions, like museums and uni-  
 917 versities, to make them readily accessible to a broader audience. We design a scraper to col-  
 918 lect a corpus of reference artists, with which we can define a test artist’s style in contrast to  
 919 the other artists, and to provide a testbed to empirically study copying behavior of generative  
 920 models. Some important landing pages to perform scraping are (i) the works by artist page  
 921 (<https://www.wikiart.org/en/Alphabet/j/text-list>; url shows all artists starting with  
 922 the letter ‘j’, and we loop through all letters), (ii) the page containing information on allowed  
 923 usage (<https://www.wikiart.org/en/terms-of-use>), (iii) an example artist landing page  
 924 (<https://www.wikiart.org/en/vincent-van-gogh>), and (iv) an example painting landing  
 925 page (<https://www.wikiart.org/en/vincent-van-gogh/the-starry-night-1889>). As  
 926 you can see, many pages have standard formats, making scraping particularly feasible. We will  
 927 provide our scraping code, along with all other code, to facilitate easy updating of our dataset as time  
 928 goes by.

929 We obtain artworks only from artists with at least 100 works on WikiArt, so to focus on somewhat  
 930 famous artists who are arguably more likely to be copied. For every work, we also scrape the licensing  
 931 information, and annotation for styles, genres, and title. In total, our dataset has 90,960 artworks over  
 932 372 artists. There are 81 styles with at least 100 works, with the most popular styles being *realism*,  
 933 *impressionism*, *romanticism*, and *expressionism*. There were 37 genres with at least 100 works, with  
 934 the most popular being *portrait*, *landscape*, *religious painting*, *sketch and study*, and *cityscape*. We  
 935 note that we only include images whose license is either public domain or fair use, with the vast  
 936 majority of works being public domain. Nonetheless, we strongly advise against using this dataset  
 937 for commercial purposes, and especially for the purpose of copying artists.

## 938 NeurIPS Paper Checklist

### 939 1. Claims

940 Question: Do the main claims made in the abstract and introduction accurately reflect the  
 941 paper’s contributions and scope?

942 Answer: [Yes]

943 Justification: Yes, the abstract accurately summarizes the paper’s claims, contributions, and  
 944 scope. We do indeed release a tool consisting of two complementary components, including  
 945 a highly interpretable one, and we utilize this tool to conduct an empirical study whose  
 946 results are as stated in the abstract.

947 Guidelines:

- 948           • The answer NA means that the abstract and introduction do not include the claims  
 949            made in the paper.
- 950           • The abstract and/or introduction should clearly state the claims made, including the  
 951            contributions made in the paper and important assumptions and limitations. A No or  
 952            NA answer to this question will not be perceived well by the reviewers.
- 953           • The claims made should match theoretical and experimental results, and reflect how  
 954            much the results can be expected to generalize to other settings.
- 955           • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
 956            are not attained by the paper.

957           **2. Limitations**

958           Question: Does the paper discuss the limitations of the work performed by the authors?

959           Answer: [Yes]

960           Justification: We include a detailed discussion of limitations as the first section in our  
 961           Appendix.

962           Guidelines:

- 963           • The answer NA means that the paper has no limitation while the answer No means that  
 964            the paper has limitations, but those are not discussed in the paper.
- 965           • The authors are encouraged to create a separate "Limitations" section in their paper.
- 966           • The paper should point out any strong assumptions and how robust the results are to  
 967            violations of these assumptions (e.g., independence assumptions, noiseless settings,  
 968            model well-specification, asymptotic approximations only holding locally). The authors  
 969            should reflect on how these assumptions might be violated in practice and what the  
 970            implications would be.
- 971           • The authors should reflect on the scope of the claims made, e.g., if the approach was  
 972            only tested on a few datasets or with a few runs. In general, empirical results often  
 973            depend on implicit assumptions, which should be articulated.
- 974           • The authors should reflect on the factors that influence the performance of the approach.  
 975            For example, a facial recognition algorithm may perform poorly when image resolution  
 976            is low or images are taken in low lighting. Or a speech-to-text system might not be  
 977            used reliably to provide closed captions for online lectures because it fails to handle  
 978            technical jargon.
- 979           • The authors should discuss the computational efficiency of the proposed algorithms  
 980            and how they scale with dataset size.
- 981           • If applicable, the authors should discuss possible limitations of their approach to  
 982            address problems of privacy and fairness.
- 983           • While the authors might fear that complete honesty about limitations might be used by  
 984            reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
 985            limitations that aren't acknowledged in the paper. The authors should use their best  
 986            judgment and recognize that individual actions in favor of transparency play an impor-  
 987            tant role in developing norms that preserve the integrity of the community. Reviewers  
 988            will be specifically instructed to not penalize honesty concerning limitations.

989           **3. Theory Assumptions and Proofs**

990           Question: For each theoretical result, does the paper provide the full set of assumptions and  
 991           a complete (and correct) proof?

992           Answer: [NA]

993           Justification: No theoretical results.

994           Guidelines:

- 995           • The answer NA means that the paper does not include theoretical results.
- 996           • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
 997            referenced.
- 998           • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 999           • The proofs can either appear in the main paper or the supplemental material, but if  
 1000            they appear in the supplemental material, the authors are encouraged to provide a short  
 1001            proof sketch to provide intuition.

- 1002           • Inversely, any informal proof provided in the core of the paper should be complemented  
1003            by formal proofs provided in appendix or supplemental material.  
1004           • Theorems and Lemmas that the proof relies upon should be properly referenced.

1005          **4. Experimental Result Reproducibility**

1006          Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
1007          perimental results of the paper to the extent that it affects the main claims and/or conclusions  
1008          of the paper (regardless of whether the code and data are provided or not)?

1009          Answer: [Yes]

1010          Justification: We explain all methods and experiments in detail, with lots of additional detail  
1011          provided in the appendix. We also provide code in a zip file, and will fully open source all  
1012          code and data if the paper is accepted.

1013          Guidelines:

- 1014           • The answer NA means that the paper does not include experiments.  
1015           • If the paper includes experiments, a No answer to this question will not be perceived  
1016            well by the reviewers: Making the paper reproducible is important, regardless of  
1017            whether the code and data are provided or not.  
1018           • If the contribution is a dataset and/or model, the authors should describe the steps taken  
1019            to make their results reproducible or verifiable.  
1020           • Depending on the contribution, reproducibility can be accomplished in various ways.  
1021            For example, if the contribution is a novel architecture, describing the architecture fully  
1022            might suffice, or if the contribution is a specific model and empirical evaluation, it may  
1023            be necessary to either make it possible for others to replicate the model with the same  
1024            dataset, or provide access to the model. In general, releasing code and data is often  
1025            one good way to accomplish this, but reproducibility can also be provided via detailed  
1026            instructions for how to replicate the results, access to a hosted model (e.g., in the case  
1027            of a large language model), releasing of a model checkpoint, or other means that are  
1028            appropriate to the research performed.  
1029           • While NeurIPS does not require releasing code, the conference does require all submis-  
1030            sions to provide some reasonable avenue for reproducibility, which may depend on the  
1031            nature of the contribution. For example  
1032              (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
1033                to reproduce that algorithm.  
1034              (b) If the contribution is primarily a new model architecture, the paper should describe  
1035                the architecture clearly and fully.  
1036              (c) If the contribution is a new model (e.g., a large language model), then there should  
1037                either be a way to access this model for reproducing the results or a way to reproduce  
1038                the model (e.g., with an open-source dataset or instructions for how to construct  
1039                the dataset).  
1040              (d) We recognize that reproducibility may be tricky in some cases, in which case  
1041                authors are welcome to describe the particular way they provide for reproducibility.  
1042                In the case of closed-source models, it may be that access to the model is limited in  
1043                some way (e.g., to registered users), but it should be possible for other researchers  
1044                to have some path to reproducing or verifying the results.

1045          **5. Open access to data and code**

1046          Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1047          tions to faithfully reproduce the main experimental results, as described in supplemental  
1048          material?

1049          Answer: [Yes]

1050          Justification: Documented code is attached in a zip file, and lots of details are included in  
1051          the appendix, including a code block. We include code to scrape the dataset as well, but  
1052          provide cached embeddings so that experiments can be run without scraping the dataset.

1053          Guidelines:

- 1054           • The answer NA means that paper does not include experiments requiring code.

- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Again, we provide extensive details in the appendix for both of our methods. These details can also be found in the attached code.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: In the appendix, we perform stability analyses where we conduct multiple trials in instances where randomness may be at play, and even try different splitting of our data to confirm that our hyperparameters are not overfit to our test set.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 1106           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1107           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1108           of Normality of errors is not verified.  
1109           • For asymmetric distributions, the authors should be careful not to show in tables or  
1110           figures symmetric error bars that would yield results that are out of range (e.g. negative  
1111           error rates).  
1112           • If error bars are reported in tables or plots, The authors should explain in the text how  
1113           they were calculated and reference the corresponding figures or tables in the text.

1114           **8. Experiments Compute Resources**

1115           Question: For each experiment, does the paper provide sufficient information on the com-  
1116           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1117           the experiments?

1118           Answer: [Yes]

1119           Justification: We conduct all experiments using a single RTX2080 GPU with four cpu  
1120           workers. We discuss the time to run our method as well. In general, this method is efficient  
1121           and does not require much compute.

1122           Guidelines:

- 1123           • The answer NA means that the paper does not include experiments.  
1124           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1125           or cloud provider, including relevant memory and storage.  
1126           • The paper should provide the amount of compute required for each of the individual  
1127           experimental runs as well as estimate the total compute.  
1128           • The paper should disclose whether the full research project required more compute  
1129           than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1130           didn't make it into the paper).

1131           **9. Code Of Ethics**

1132           Question: Does the research conducted in the paper conform, in every respect, with the  
1133           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1134           Answer: [Yes]

1135           Justification: We adhere to the ethical guidelines and discuss the societal implications of our  
1136           work at length.

1137           Guidelines:

- 1138           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
1139           • If the authors answer No, they should explain the special circumstances that require a  
1140           deviation from the Code of Ethics.  
1141           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
1142           eration due to laws or regulations in their jurisdiction).

1143           **10. Broader Impacts**

1144           Question: Does the paper discuss both potential positive societal impacts and negative  
1145           societal impacts of the work performed?

1146           Answer: [Yes]

1147           Justification: This paper is designed to answer a pressing legal and material question around  
1148           how AI ultimately affects people. We attempt to be objective in our analysis, while building  
1149           a tool that will help artists with stylistic infringements, even if they are not being infringed  
1150           upon yet. This tool can also help producers of generative models defend themselves, as  
1151           they now have a way to say that they aren't producing infringing upon unique artistic styles  
1152           (when that is the case).

1153           Guidelines:

- 1154           • The answer NA means that there is no societal impact of the work performed.  
1155           • If the authors answer NA or No, they should explain why their work has no societal  
1156           impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 1176 11. Safeguards

1177 Question: Does the paper describe safeguards that have been put in place for responsible  
 1178 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
 1179 image generators, or scraped datasets)?

1180 Answer: [Yes]

1181 Justification: We discuss the potential risks and safeguards associated with our dataset in the  
 1182 Appendix.

1183 Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 1194 12. Licenses for existing assets

1195 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
 1196 the paper, properly credited and are the license and terms of use explicitly mentioned and  
 1197 properly respected?

1198 Answer: [Yes]

1199 Justification: We mention the licenses of the data we use, and include these licenses in the  
 1200 metadata of our dataset for others to see later.

1201 Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 1209           • If assets are released, the license, copyright information, and terms of use in the  
1210           package should be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets)  
1211           has curated licenses for some datasets. Their licensing guide can help determine the  
1212           license of a dataset.  
1213           • For existing datasets that are re-packaged, both the original license and the license of  
1214           the derived asset (if it has changed) should be provided.  
1215           • If this information is not available online, the authors are encouraged to reach out to  
1216           the asset's creators.

1217           **13. New Assets**

1218           Question: Are new assets introduced in the paper well documented and is the documentation  
1219           provided alongside the assets?

1220           Answer: [NA]

1221           Justification: No new assets.

1222           Guidelines:

- 1223           • The answer NA means that the paper does not release new assets.  
1224           • Researchers should communicate the details of the dataset/code/model as part of their  
1225           submissions via structured templates. This includes details about training, license,  
1226           limitations, etc.  
1227           • The paper should discuss whether and how consent was obtained from people whose  
1228           asset is used.  
1229           • At submission time, remember to anonymize your assets (if applicable). You can either  
1230           create an anonymized URL or include an anonymized zip file.

1231           **14. Crowdsourcing and Research with Human Subjects**

1232           Question: For crowdsourcing experiments and research with human subjects, does the paper  
1233           include the full text of instructions given to participants and screenshots, if applicable, as  
1234           well as details about compensation (if any)?

1235           Answer: [Yes]

1236           Justification: Included in the appendix, with workers receiving pay between \$12 and \$18 an  
1237           hour (USD).

1238           Guidelines:

- 1239           • The answer NA means that the paper does not involve crowdsourcing nor research with  
1240           human subjects.  
1241           • Including this information in the supplemental material is fine, but if the main contribu-  
1242           tion of the paper involves human subjects, then as much detail as possible should be  
1243           included in the main paper.  
1244           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1245           or other labor should be paid at least the minimum wage in the country of the data  
1246           collector.

1247           **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
1248           Subjects**

1249           Question: Does the paper describe potential risks incurred by study participants, whether  
1250           such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1251           approvals (or an equivalent approval/review based on the requirements of your country or  
1252           institution) were obtained?

1253           Answer: [NA]

1254           Justification: We confirm with IRB that our crowdsourced validation does not require IRB  
1255           review.

1256           Guidelines:

- 1257           • The answer NA means that the paper does not involve crowdsourcing nor research with  
1258           human subjects.

- 1259
- Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1260 may be required for any human subjects research. If you obtained IRB approval, you
- 1261 should clearly state this in the paper.
- 1262
- We recognize that the procedures for this may vary significantly between institutions
- 1263 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1264 guidelines for their institution.
- 1265
- For initial submissions, do not include any information that would break anonymity (if
- 1266 applicable), such as the institution conducting the review.