

Text-Guided Synthesis of Artistic Images with Retrieval-Augmented Diffusion Models

Robin Rombach*, Andreas Blattmann*, and Björn Ommer

Ludwig-Maximilian University Munich, Germany

**the first two authors contributed equally to this work*

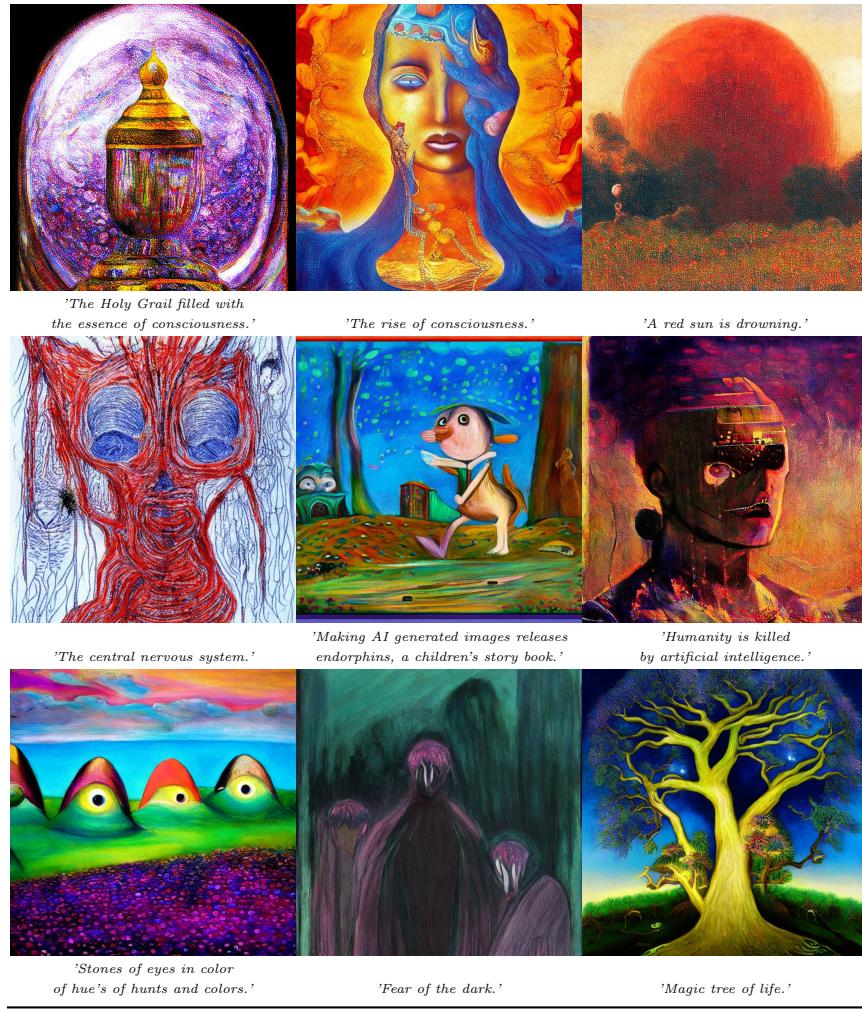


Fig. 1: Selected 768×768 samples for various text inputs obtained with our “stylized” LAION-RDM (see Sec. 2.2). Note that the model was only trained on images.

Abstract. Novel architectures have recently improved generative image synthesis leading to excellent visual quality in various tasks. Of particular note is the field of “AI-Art”, which has seen unprecedented growth with the emergence of powerful multimodal models such as CLIP. By combining speech and image synthesis models, so-called ”prompt-engineering” has become established, in which carefully selected and composed sentences are used to achieve a certain visual style in the synthesized image. In this note, we present an alternative approach based on retrieval-augmented diffusion models (RDMs). In RDMs, a set of nearest neighbors is retrieved from an external database during training for each training instance, and the diffusion model is conditioned on these informative samples. During inference (sampling), we replace the retrieval database with a more specialized database that contains, for example, only images of a particular visual style. This provides a novel way to “prompt” a general trained model after training and thereby specify a particular visual style. As shown by our experiments, this approach is superior to specifying the visual style within the text prompt. We open-source code and model weights at <https://github.com/CompVis/latent-diffusion>.

Keywords: Synthesis, Diffusion Models, Retrieval, CLIP

Diffusion models have recently set the state of the art in image generation and controllable synthesis [9,22]. In text-to-image synthesis in particular, we have seen impressive results [21,23] that can also be used to create artistic images. Such models thus have the potential to help artists create new content and have contributed to the tremendous growth of the field of AI generated art [7]. However, these models are very compute intensive and so far cannot be reused for tasks other than those for which they were trained. For this reason, in the present work we build on the recently introduced retrieval-augmented diffusion models (RDMs) [3,2], which also have the potential to significantly reduce the computational complexity required in training by providing a comparatively small generative model with a large image database: While the retrieval approach provides the (local) content, the model can now focus on learning the composition of scenes based on this content. In this extended abstract, we scale RDMs and show their capability to generate artistic images as those shown in Fig. 1. Moreover, we can control the synthesis process with natural language by using the joint text-image representation space of CLIP [20] and demonstrate that we obtain fine-grained control over the visual style of the output by retrieving neighbors from highly specialized databases built from WikiArt [24] and ArtBench [16]. Finally, we also consider the release of our model weights as a contribution that allows artists to complement, extend, and evaluate their work and also to investigate the inherent biases of these models.

1 Recap on Retrieval-Augmented Diffusion Models

Following [3,2], a retrieval-augmented diffusion model (*RDM*) is a combination of a conditional latent diffusion model ϵ_θ [12,22], a database of images $\mathcal{D}_{\text{train}}$, which is considered to be an explicit part of the model, and a (non-trainable)

sampling strategy ξ_k to obtain a subset of $\mathcal{D}_{\text{train}}$ based on a query x as introduced in [3]. The model is trained by implementing ξ_k as a nearest neighbor algorithm, such that for each query (i.e., training example) its k nearest neighbors are returned as a set, where the distance is measured in CLIP [20] image embedding space. The CLIP embeddings of these nearest neighbors are then fed to the model via the cross attention mechanism [28,22]. The training objective reads

$$\min_{\theta} \mathcal{L} = \mathbb{E}_{p(x), z \sim E(x), \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \{\phi_{\text{CLIP}}(y) | y \in \xi_k(x, \mathcal{D}_{\text{train}})\})\|_2^2 \right], \quad (1)$$

where ϕ_{CLIP} is the CLIP image encoder and $E(x)$ is the encoder of an autoencoding model as deployed in [22,3]. After training, we replace $\mathcal{D}_{\text{train}}$ of the original *RDM* with alternate databases $\mathcal{D}_{\text{style}}$, derived from art datasets [24,16] to obtain a post-hoc model modification and thereby zero-shot stylization. Furthermore, we can guide the synthesis process with text-prompts by using the shared text-image feature space of CLIP [20] as proposed in [3]. Thus, we obtain a *controllable* synthesis model which is only trained on image data.

2 Text-Guided Synthesis of Artistic Images with RDMs

2.1 General Setting

We conduct experiments for two models: To show the general zero-shot stylization potential of *RDM*, we train an exact replica of the *RDM* on ImageNet [8] as proposed in [3], i.e., we build $\mathcal{D}_{\text{train}}$ from OpenImages [15]. For inference, we achieve stylization by using a database $\mathcal{D}_{\text{style}}$ based on the WikiArt [24] dataset *cf.* Sec. 2.2. In Sec. 2.2 we present a larger model, trained on 100M examples from LAION-2B-en [25,1] with a more diverse database $\mathcal{D}_{\text{train}}$, which contains the remaining 1.9B samples from that dataset. Samples from this model are shown in Fig. 1. By exchanging this database with distinct, style-specific subsets of the ArtBench dataset [16] during inference, we show that *RDM* can further be used for fine-grained stylization, without being trained for this task. Details on training and inference are provided in Appendix B.

2.2 Zero-Shot Text-Guided Stylization by Exchanging the Database

By replacing the train database $\mathcal{D}_{\text{train}}$ with WikiArt [24] we show the zero-shot stylization capabilities of the ImageNet-*RDM* from in Sec 2.1 in Fig. 2. Our model, though only trained on ImageNet, generalizes to this new database and is capable of generating artwork-like images which depict the content defined by the text prompts. To further emphasize the effects of this post-hoc exchange of the database, we show samples obtained with the same procedure but using $\mathcal{D}_{\text{train}}$ (bottom row).

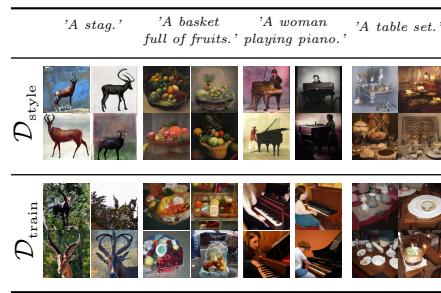


Fig. 2: Zero-shot text-guided stylization with our ImageNet-*RDM*. Best viewed when zoomed in.

Fine-Grained Stylization with ArtBench Many powerful models emulate text-driven stylization by adding the postfix "... *in the style of ...*" to a given prompt [22,19,21,23,30]. By using style specific databases obtained from the ArtBench dataset [16] during inference, we here present an alternative approach. Fig. 3 presents results for the prompt "*Day and night fighting for the domination of time.*" and the LAION-RDM. Each column contains samples generated by replacing $\mathcal{D}_{\text{train}}$ with a style-specific ArtBench-subset. For a quantitative

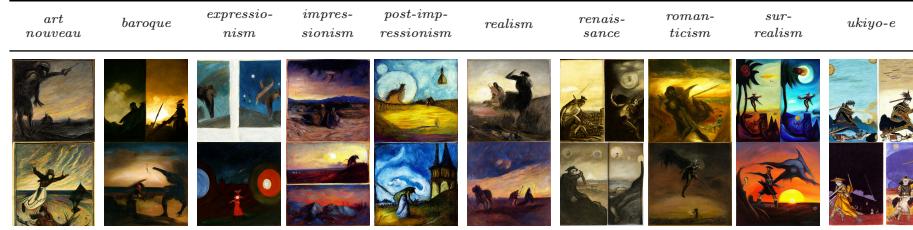
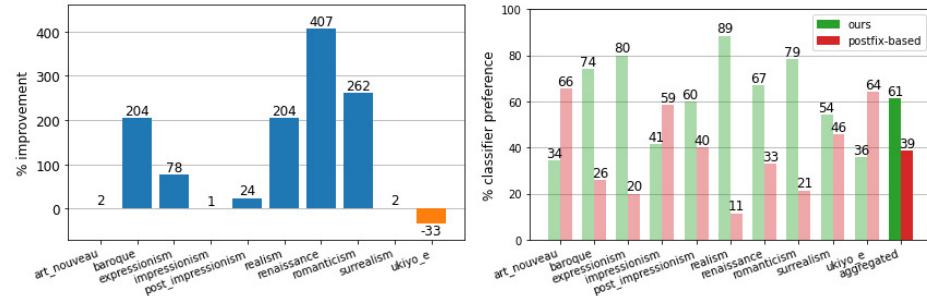


Fig. 3: Visual Examples of fine-grained zero-shot text based stylization with LAION model. The prompt used to generate these samples is '*Day and night fighting for the domination of time.*' More samples provided in the Appendix D.

evaluation, we generate 70 samples per style for both approaches, which we then classify with a style-classifier (details in Appendix C) and compare both the relative improvement in accuracy (Fig. 4a) and the classifier logits (Fig. 4b). The retrieval-based approach almost always outperforms postfix-based stylization.



(a) Relative improvement of our proposed approach over postfix-based stylization. (b) Classifier voting results: Our method outperforms the postfix approach for nearly every style.

Fig. 4: Quantitative comparison of our retrieval and postfix-based prompt stylization.

3 Conclusion

In this note, we present an approach to train accessible and controllable models for visual art: By building on the recently introduced retrieval-augmented diffusion models, our approach becomes *accessible* as we efficiently store an image database and condition a comparatively small generative model directly on meaningful samples from the database, rather than compressing large training data into increasingly large generative models. Our approach is *controllable* because it allows post-hoc replacement of the external database and thus specification of a desired visual style, which emerges in our experiments as a strong alternative to pure text-based approaches. In future work, we plan to combine this approach with post-hoc finetuning on paired text-image data.

References

1. Laion-5b: A new era of open large-scale multi-modal datasets. <https://laion.ai/blog/laion-5b/>, accessed: 2022-07-07
2. Ashual, O., Sheynin, S., Polyak, A., Singer, U., Gafni, O., Nachmani, E., Taigman, Y.: Knn-diffusion: Image generation via large-scale retrieval. arXiv preprint arXiv:2204.02849 (2022)
3. Blattmann, A., Rombach, R., Oktay, K., Ommer, B.: Retrieval-augmented diffusion models. arXiv (2022)
4. Borgeaud, S., Mensch, A., Hoffmann, J., Cai, T., Rutherford, E., Millican, K., Driessche, G.v.d., Lespiau, J.B., Damoc, B., Clark, A., et al.: Improving language models by retrieving from trillions of tokens. arXiv preprint arXiv:2112.04426 (2021)
5. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems **33**, 1877–1901 (2020)
6. Casanova, A., Careil, M., Verbeek, J., Drozdzal, M., Romero Soriano, A.: Instance-conditioned gan. Advances in Neural Information Processing Systems **34** (2021)
7. Crowson, K., Biderman, S., Kornis, D., Stander, D., Hallahan, E., Castricato, L., Raff, E.: Vqgan-clip: Open domain image generation and editing with natural language guidance (2022). <https://doi.org/10.48550/ARXIV.2204.08583>, <https://arxiv.org/abs/2204.08583>
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
9. Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. Advances in Neural Information Processing Systems **34** (2021)
10. Guo, R., Sun, P., Lindgren, E., Geng, Q., Simcha, D., Chern, F., Kumar, S.: Accelerating large-scale inference with anisotropic vector quantization. In: III, H.D., Singh, A. (eds.) Proceedings of the 37th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 119, pp. 3887–3896. PMLR (13–18 Jul 2020), <https://proceedings.mlr.press/v119/guo20h.html>
11. Guu, K., Lee, K., Tung, Z., Pasupat, P., Chang, M.: Retrieval augmented language model pre-training. In: International Conference on Machine Learning. pp. 3929–3938. PMLR (2020)
12. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems **33**, 6840–6851 (2020)
13. Khandelwal, U., Fan, A., Jurafsky, D., Zettlemoyer, L., Lewis, M.: Nearest neighbor machine translation. arXiv preprint arXiv:2010.00710 (2020)
14. Khandelwal, U., Levy, O., Jurafsky, D., Zettlemoyer, L., Lewis, M.: Generalization through memorization: Nearest neighbor language models. arXiv preprint arXiv:1911.00172 (2019)
15. Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J.R.R., Krasin, I., Pont-Tuset, J., Kamali, S., Popov, S., Malloci, M., Duerig, T., Ferrari, V.: The open images dataset V4: unified image classification, object detection, and visual relationship detection at scale. CoRR **abs/1811.00982** (2018), <http://arxiv.org/abs/1811.00982>
16. Liao, P., Li, X., Liu, X., Keutzer, K.: The artbench dataset: Benchmarking generative models with artworks. arXiv (2022)
17. Long, A., Yin, W., Ajanthan, T., Nguyen, V., Purkait, P., Garg, R., Blair, A., Shen, C., Hengel, A.v.d.: Retrieval augmented classification for long-tail visual recognition. arXiv preprint arXiv:2202.11233 (2022)

18. Meng, Y., Zong, S., Li, X., Sun, X., Zhang, T., Wu, F., Li, J.: Gnn-lm: Language modeling based on global contexts via gnn. arXiv preprint arXiv:2110.08743 (2021)
19. Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint arXiv:2112.10741 (2021)
20. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning. pp. 8748–8763. PMLR (2021)
21. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 (2022)
22. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. arXiv preprint arXiv:2112.10752 (2021)
23. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S.K.S., Ayan, B.K., Mahdavi, S.S., Lopes, R.G., Salimans, T., Ho, J., Fleet, D.J., Norouzi, M.: Photorealistic text-to-image diffusion models with deep language understanding (2022). <https://doi.org/10.48550/ARXIV.2205.11487>, <https://arxiv.org/abs/2205.11487>
24. Saleh, B., Elgammal, A.M.: Large-scale classification of fine-art paintings: Learning the right metric on the right feature. CoRR **abs/1505.00855** (2015), <http://arxiv.org/abs/1505.00855>
25. Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., Komatsuzaki, A.: Laion-400m: Open dataset of clip-filtered 400 million image-text pairs (2021)
26. Siddiqui, Y., Thies, J., Ma, F., Shan, Q., Nießner, M., Dai, A.: Retrievalfuse: Neural 3d scene reconstruction with a database. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12568–12577 (2021)
27. Tseng, H.Y., Lee, H.Y., Jiang, L., Yang, M.H., Yang, W.: Retrievegan: Image synthesis via differentiable patch retrieval. In: European Conference on Computer Vision. pp. 242–257. Springer (2020)
28. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
29. Xu, R., Guo, M., Wang, J., Li, X., Zhou, B., Loy, C.C.: Texture memory-augmented deep patch-based image inpainting. IEEE Transactions on Image Processing **30**, 9112–9124 (2021)
30. Yu, J., Xu, Y., Koh, J.Y., Luong, T., Baid, G., Wang, Z., Vasudevan, V., Ku, A., Yang, Y., Ayan, B.K., Hutchinson, B., Han, W., Parekh, Z., Li, X., Zhang, H., Baldridge, J., Wu, Y.: Scaling autoregressive models for content-rich text-to-image generation (2022). <https://doi.org/10.48550/ARXIV.2206.10789>, <https://arxiv.org/abs/2206.10789>

Appendix

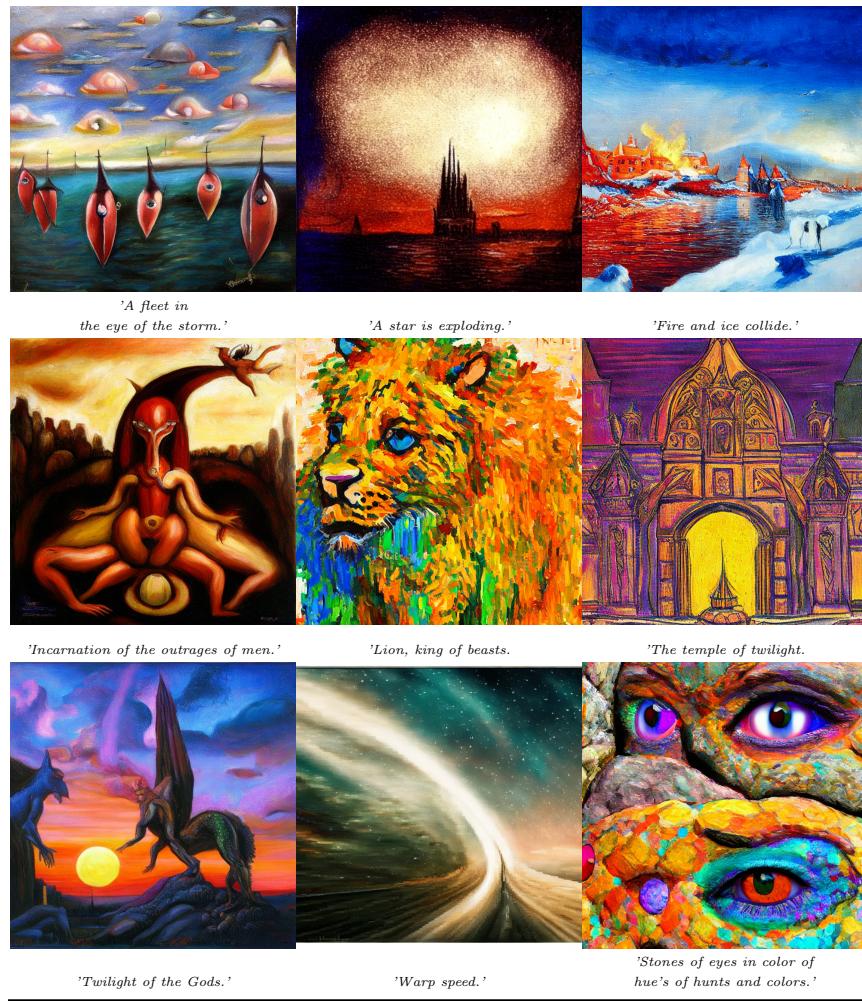


Fig. 5: More 768×768 samples for various text inputs obtained with our “stylized” LAION-RDM (see Sec. 2.2). Note that the model was only trained on images.

A Related Work

Retrieval-Augmented Generative Models. Using external memory to augment traditional models has recently drawn attention in natural language processing (NLP) [14,13,18,11]. For example, RETRO [4] proposes a retrieval-enhanced

transformer for language modeling which performs on par with state-of-the-art models [5] using significantly less parameters and compute resources. These retrieval-augmented models with external memory turn purely parametric deep learning models into semi-parametric ones. Early attempts [17,26,27,29] in retrieval-augmented visual models do not use an external memory and exploit the training data itself for retrieval. In image synthesis, IC-GAN [6] utilizes the neighborhood of training images to train a GAN and generates samples by conditioning on single instances from the training data. However, using training data itself for retrieval potentially limits the generalization capacity, and thus, we favor an external memory in this work. For diffusion models, RDM [3] and kNN-diffusion [2] pioneer retrieval augmentation.

B Details on Presented Models

	<i>RDM</i>	<i>RDM</i>
Train Dataset	ImageNet	LAION-2B-en (100M examples) [†]
image size	256^2	$256^2 / 512^2 / 768^2$
z -shape (smallest image size)	$64 \times 64 \times 3$	$16 \times 16 \times 16$
$ \mathcal{Z} $	8192	KL
Diffusion steps	1000	1000
Noise Schedule	linear	linear
Model Size	400M	1.3B
Channels	192	448
Depth	2	2
Channel Multiplier	1,2,3,5	1,2,3,4
Number of Channels per Head	32	32
Batch Size	1240	$1600 / 800 / 432$
Iterations	112K	$275\text{K}(256^2) / 225\text{K}(512^2) / 270\text{K}(768^2)$
Learning Rate	1.0e-4	1.0e-4
Conditioning	CA	CA
CA-resolutions	32, 16, 8	16, 8, 4
ϕ CLIP model spec	ViT-B/32	ViT-L/14
Embedding Dimension	512	768
Transformer Depth	1	1
size $\mathcal{D}_{\text{train}}$	20M	1.9B
k	4	8
$\mathcal{D}_{\text{train}} \cap \text{Train Dataset} = \emptyset$	✓	✓
image size $\mathcal{D}_{\text{train}}$	256×256	224×224

Table 1: Hyperparameters for the models introduced in Sec. 2.1. [†]: Image size has been successively increased during training, see image sizes

B.1 ImageNet Model

The hyperparameters used for the ImageNet model are presented in Sec. 2.1 and evaluated in Sec. 2.2 are presented in the first column of Tab. 1. We here note again that this model is an exact replica of the ImageNet model presented in [3]. During training, we use a database $\mathcal{D}_{\text{train}}$ of 20M examples from the OpenImages [15] dataset. We extract 2-3 patches per image (more details, see [3]) to use each image at least once. To obtain nearest neighbors we apply the ScaNN search algorithm [10] in the feature space of a pretrained CLIP-ViT-B/32 [20]. Using this setting, retrieving 20 nearest neighbors from the database described above takes approximately 0.95 ms. The model is trained on eight NVIDIA A-100-SXM4 with 80GB RAM per GPU.

B.2 LAION Model

The hyperparameters used for our LAION model are subsumed in the second column of Tab. 1. For this model, we scale the size of the training dataset as well as those of the database $\mathcal{D}_{\text{train}}$. The train dataset is obtained by only using images of the LAION-2B-en dataset [25,1], with a shorter edge length larger than 768 px and filtering images containing watermarks and unsafe content. The database $\mathcal{D}_{\text{train}}$ is constructed based on the remaining 1.9B images. Hence the database $\mathcal{D}_{\text{train}}$ is disjoint from the training set. In contrast to the ImageNet model we use the CLIP-ViT-L/14 model both as a retrieval distance and as nearest neighbor encoder ϕ_{CLIP} . We train the model in three stages and successively increase the image size of the train images from 256^2 px in the first to 512^2 in the second stage before we reach our final resolution of 768^2 px in the final third stage. The model is trained on eight NVIDIA A-100-SXM4 with 80GB RAM per GPU.

B.3 Retrieval Strategy During Inference

To obtain the nearest neighbors for text-based stylization from the inference-time database $\mathcal{D}_{\text{style}}$ we embed the query prompt into the shared CLIP text-image space by using the CLIP text-encoder [20]. We then retrieve the $k = 19$ nearest neighbors from $\mathcal{D}_{\text{style}}$. As a distance measure, we use cosine similarity.

C Details on Style-Classifier

The style classifier is trained on ArtBench and implemented as a two-layer perceptron on top of CLIP image features. Its top-1 accuracy on the validation set is 77%.

D Addtional Qualitative Results

In Fig. 5 we show additional samples from our LAION-model and in Fig. 6 we provide additional examples showing the style-specific stylization capabilities of this model.

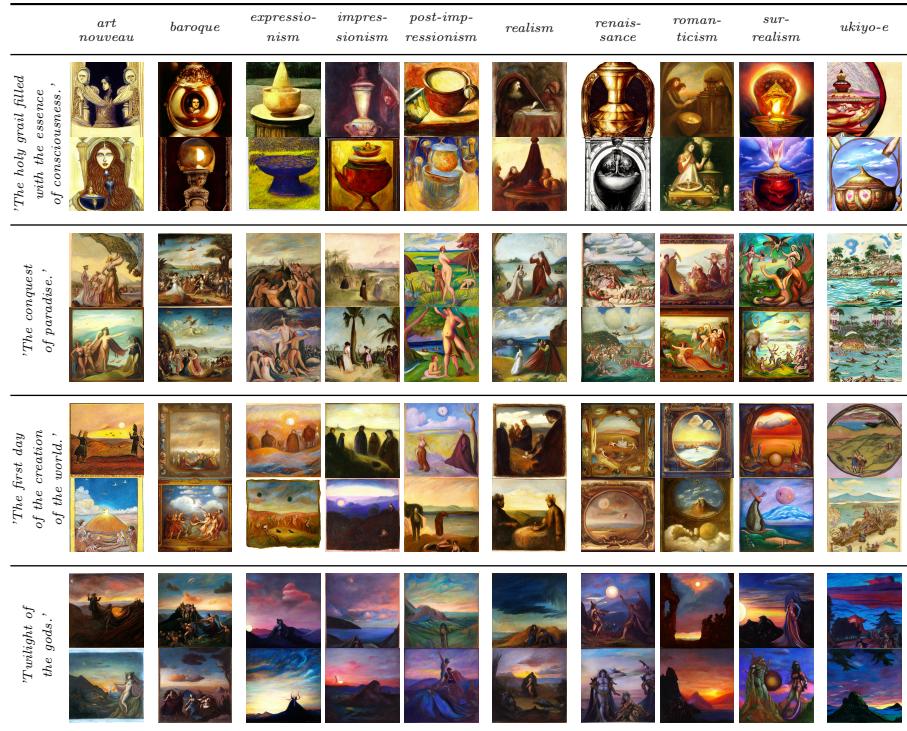


Fig. 6: More visual examples of stylization with our LAION model as in Fig. 3.