

# Self-supervised disentangled representation learning of artistic style through Neural Style Transfer

Dan Ruta<sup>1</sup>, Gemma Canet Tarrés<sup>1</sup>, Alexander Black<sup>1</sup>, Andrew Gilbert<sup>1</sup>, and John Collomosse<sup>1,2</sup>

<sup>1</sup> University of Surrey

<sup>2</sup> Adobe Research

**Abstract.** We present a new method for learning a fine-grained representation of visual style. Representation learning aims to discover individual salient features of a domain in a compact and descriptive form that strongly identifies the unique characteristics of that domain. Prior visual style representation works attempt to disentangle style (*i.e.* appearance) from content (*i.e.* semantics) yet a complete separation has yet to be achieved. We present a technique to learn a representation of visual style more strongly disentangled from the semantic content depicted in an image. We use Neural Style Transfer (NST) to measure and drive the learning signal and achieve state-of-the-art representation learning on explicitly disentangled metrics. We show that strongly addressing the disentanglement of style and content leads to large gains in style-specific metrics, encoding far less semantic information and achieving state-of-the-art accuracy in downstream style matching (retrieval) and zero-shot style tagging tasks.

## 1 Introduction

Visual style refers to the depiction or ‘appearance’ of subject matter within an image, as opposed to the subject matter or ‘content’ depicted. Describing the visual style of image through a learned representation remains an open research challenge for computer vision. This is, in part, due to scarcity of labelled style data, and in part due to the subjectivity of defining a descriptive ontology for style beyond coarse-grained categorical labels. Training is therefore often performed in a comparative (*e.g.* contrastive) setting, where similarities and differences between two stylistically similar images can hint at common properties. A recurring goal in style representation learning is therefore separating and disentangling visual style from the depicted subject matter.

Use cases for learned representations of visual style include style conditioned image retrieval [19] and generation [6,24], stylization [21], automatic style tagging [19], and image translation [17]. Disentanglement in embeddings is critical in such multimodal applications, where a clean disentangled signal of the style modality is needed to independently control or describe style. In this work, we show that

style-content entanglement is still present in state-of-the-art representations. We propose a novel learning methodology for a fully disentangled learning of style. We show that this explicitly disentangled representation benefits downstream tasks like style-based image retrieval and tagging. Our contributions are:

1. Novel methodology for training a style representation model without content-style entanglement in the data, trained over the BBST-4M dataset [20].
2. State-of-the-art in style representation learning with enforced disentanglement, with a new benchmark dataset.
3. New state-of-the-art multimodal vision/language learning in the context of artistic style for zero-shot style tagging.



**Fig. 1:** Please zoom in for details. (Left) Example style groups from the BAM-FG dataset. The images in each group are style consistent, but they are also semantically consistent. For example, the top left style group has a consistent *weathered paper* style but is also consistent in the subject matter of character design. The top right has consistent *pastel* style but is consistently interiors. The bottom left is consistent *moody vignette dark photography* style, but all images are of landscapes. Bottom right *vector art* images all contain faces. (Right) Example synthetic style consistent images, as used in our work (via NeAT). The left-most images in each style group are the reference style image. The BAM-FG data (left) shows style consistency at the cost of entanglement with semantic consistency, unlike the synthetic data (right).

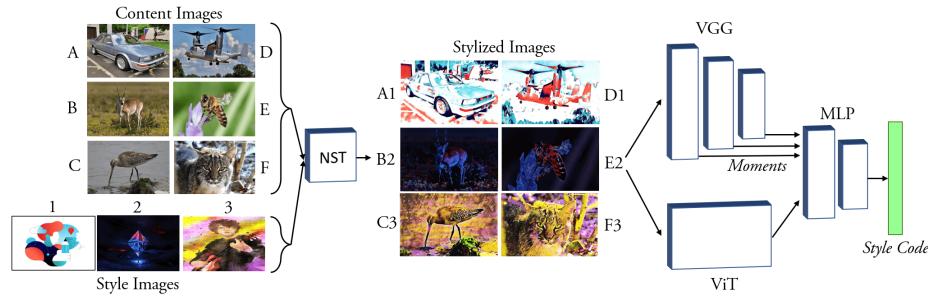
## 2 Related Work

The seminal work of Gatys’s neural style transfer [5] introduced the concept of using neural, learning based methods for re-rendering a given content image, such as a photograph, to match the visual artistic style of a second stylistic image, typically an artwork. Other works extended neural style transfer to multiple, and eventually arbitrary styles per model [10, 12–14, 25].

AvatarNet [23], and later SANet [16] explore the application of self attention modules in performing style transfer in a feed-forward manner, aligning

feature statistics between the content and style images. PAMA [15] proposes an alternative attention mechanism based on iterative refinement of feature alignment, progressively adjusting content features to match style features in a more spatially consistent manner. ContraAST [1] uses self attention as per SANet in conjunction with domain-level adversarial losses and contrastive losses to push the stylized images to resemble distributions of real images better - thereby creating more convincingly real looking images regardless of style. CAST [27] primarily improve this process by including ground truth style images in the loss.

NeAT [20] further build on the work in ContraAST and CAST, expanding on the attention approach in PAMA, and using other robustness and quality improvements. They perform stylization as an image editing process rather than an image re-generation process by predicting deltas over a partially corrupted version of the reference content image.



**Fig. 2:** Visualization of our NST-driven style representation learning method. We show a training iteration with batch size 6, with 6 content images and 3 style images (in our experiments, we use much larger batch sizes but use 6 here for clarity). The content images are stylized with a pre-trained and frozen Neural Style Transfer method using two copies of the 3 style images. We extract a style embedding using layer-wise global moment statistics and the logits from a more localized vision transformer.

In a similar branch of research, image translation works like MUNIT [9] and Swapping Autoencoders (SAE) [17] decompose a pair of images into structural information and global unlocalized latents, which can be mixed during inference to render an image with mixed properties. These works demonstrate how an embedding optimized to capture global information (such as style) can be used in a generative setting.

Using a triplet loss, [3] learn a coarse metric style representation for 7 styles, using the style-labeled subset of the BAM dataset [26]. ALADIN [22] first explored a *fine-grained* style representation using their newly labeled BAM-FG dataset. Depending on the chosen similarity strength, this larger dataset contains up to 135k style groups. They design their model for the disentangled representation of content and style by extracting features as global AdaIN statistics

from each encoder layer. The BAM-FG dataset was curated via crowd annotation to select style-coherent images in existing weakly labeled style coherent groups of images from *Behance.net*. The labeling process removed anomalies and cleaned the coherent style groups such that the remaining images in a group, at several thresholds, was human verified to be style consistent. This helped to drive ALADIN to be state-of-the-art in style representation capabilities, further to the training methodology.

However, this labeling process needs to be revised. The dataset is indeed style coherent, but the labeling process only improves the style coherency of any given small set of images - it does not help avoid content and style entanglement. If all the images in a style group have the same content depicted, the BAM-FG cleaning process only helps ensure they are also style consistent. But resulting style groups are also consistent in the content information.

<i>Anchor</i>	<i>Positive</i>	<i>Negatives</i>
A 1	D 1	B 2 C 3 E 2 F 3
B 2	E 2	A 1 C 3 D 1 F 3
C 3	F 3	A 1 B 2 D 1 E 2
D 1	A 1	B 2 C 3 E 2 F 3
E 2	B 2	A 1 C 3 D 1 F 3
F 3	C 3	A 1 B 2 D 1 E 2

Contrastive Losses (stylized)

<i>Anchor</i>	<i>Positive</i>	<i>Negatives</i>
A 1	1	B 2 C 3 2 3
D 1	1	B 2 C 3 3 E 2 F 3
B 2	2	A 1 C 3 1 3
E 2	2	A 1 C 3 3 D 1 F 3
C 3	3	A 1 2 B 2 1
F 3	3	B 1 2 D 1 E 2

Contrastive Losses (ground truth)

**Fig. 3:** A set of contrastive losses are computed for each stylized image in the batch. The positive sample is the other sample in the batch where the same original style image was used as a style reference during stylization. As half the number of style images are selected per batch, there will always be two images with the same style. The negative samples in the contrastive losses are thus the remaining images in the batch, which are stylized with other randomly sampled style images in the batch. Additional sets of contrastive losses compare the stylized images' embeddings to the embeddings of the source style images, as per CAST [27] and NeAT. Red squares represent ground truth images' embeddings.

As artists develop their skills, they likely specialize in specific subsets of subject matter, such as faces or character design. Alternatively, the work they publish can showcase a project they worked on where the subject matter was constrained to some requirements. This effect is visualized in Figure 1 (left), showing a few style groups from the BAM-FG dataset. The images therein are indeed style consistent, but they also share semantic features.

Model	NST learning signal			NeAT test set		PAMA test set		SANet test set		Average values		CAST test set	
	NeAT	PAMA	SANet	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1
ALADIN-ViT	-	-	-	16.823	0.270	9.9964	0.0575	12.493	0.0599	13.104	0.129	39.541	0.960
Ours (ViT)	✓			85.306	66.308	51.012	15.303	67.525	28.423	67.948	36.678	39.472	<b>10.765</b>
Ours (ViT)		✓		69.226	23.415	62.886	20.628	51.934	6.393	61.349	16.812	24.563	3.608
Ours (ViT)			✓	80.230	46.466	56.215	18.738	74.621	35.693	70.355	33.632	38.171	8.895
Ours (ViT)	✓	✓		84.997	59.650	68.468	30.563	67.021	23.443	73.495	37.885	38.864	8.995
Ours (ViT)	✓		✓	85.052	64.993	53.657	17.688	77.410	46.408	72.040	<b>43.030</b>	<b>39.880</b>	<b>10.715</b>
Ours (ViT)		✓	✓	77.056	36.413	64.596	23.048	67.229	22.835	69.627	27.432	32.523	5.243
Ours (ViT)	✓	✓	✓	83.915	58.900	67.484	29.745	74.755	34.460	<b>75.385</b>	<u>41.035</u>	<u>39.688</u>	9.010

**Table 1:** Style representation learning metrics (IR-topk and mAP) of our model with different NST learning signals. We also compare against ALADIN-ViT [19], keeping the same backbone. We measure the representation learning using multiple test sets compiled with different style transfer works from literature: NeAT [20], PAMA [15], and SANet [16]. We also measure using CAST [27], which we reserve as an NST method not seen during training, only for use during evaluation.

### 3 Methodology

In our work, we set out to create a model to learn disentangled representations of style without being affected by semantics data biases. We seek to train a model on data that has high variance in the semantic content depicted but has a consistent style. As discussed in previous sections, real data with such properties is rare or impractical to create through human artists. Instead, we use the current state-of-the-art neural style transfer methods to create synthetic datasets of stylized images where the style is consistent, but where the content varies depending on our source content images. Fig. 1 (right) visualizes synthetic stylized data used in our work. Given a style image, we can generate images with the same style but completely random and arbitrary semantic content.

Given a batch of content and style images, we know the synthesized data’s ground truth style and content relations. We dynamically use fast, feed-forward NST methods during training to maximize the number of styles we can use without the impractical storage space needed to pre-compute the images. We induce the style learning signal through contrastive losses [2], computed amongst the images generated by the NST method and the reference style image. We sample only half the number of style images in a batch to synthesize two images with the same style in each batch, for each style. For each synthetic stylized image, we use the other image stylized with the same style in this batch as the positive

and the remaining images in the batch (stylized with the different style images) as negatives. This encourages our embedding to represent the style information shared in the stylized pairs, regardless of the semantic content depicted, which is random. We use standard contrastive losses to drive the learning signal using this self-supervised approach of data labelling, as visualized in Fig 2.

This approach also benefits from using style images from datasets where style-consistent labeling is not required in a self-supervised manner. We thus use the BBST-4M dataset [20], as it has one of the highest diversity of style images - 2 million images in the style subset. The style subset in BBST-4M is also filtered to only contain stylistic (artistic) data, unlike BAM-FG, which includes style groups of non-stylistic images such as photos. Artistic images are better suited for NST - processes specifically designed to transfer such style.

### 3.1 Moments

Global feature statistics have been demonstrated in literature [8, 14] to capture global style in an image. Standard statistics used are mean and variance. In moments, these represent the first and second moment, though higher order moments have been used to drive NST through moment matching [11], with higher quality. We thus use the first four moments in our work, further extracting skewness and kurtosis from feature statistics in the VGG branch.

The skewness formula is shown in Eq 2, calculated via the z scores (Eq 1), and the kurtosis is shown in Eq 3, where a positive value indicates leptokurtic data distribution, and a negative value indicates a platykurtic distribution - measures of the tails of the data distribution.

$$z_{scores} = \frac{X - \mu}{\sigma} \quad (1)$$

$$m_3 = \frac{\sum z_{scores}^3}{n} \quad (2)$$

$$m_4 = \frac{\sum z_{scores}^4}{n} - 3 \quad (3)$$

We also include highly expressive features extracted from a vision transformer [4] model to capture more strongly capture features from an image. We concatenate these embeddings and project them into a 1024-dim style code, shown in Figure 2.

### 3.2 Loss

The loss objective is a standard contrastive loss, shown in Eq 6, where  $\mathcal{A}$  represents our model,  $x_s$  and  $x_c$  represent style and content images respectively, and  $NST$  represents a randomly sampled NST method from the methods used, to stylize  $x_s$  and  $x_c$  into  $\mathcal{S}_{sc}$ :

$$\mathcal{S}_{sc} = NST(x_s, x_c) \quad (4)$$

$$pos = \mathcal{A}(\mathcal{S}_{sc})_a^T \mathcal{A}(\mathcal{S}_{sc})_p / \tau \quad (5)$$

$$\mathcal{L} := -\log \left( \frac{\exp(pos)}{\exp(pos) + \sum \exp(\mathcal{A}(\mathcal{S}_{sc})_a^T \mathcal{A}(\mathcal{S}_{sc})_n / \tau)} \right) \quad (6)$$

## 4 Experiments

Model	Dataset	NeAT test set				PAMA test set				SANet test set				<i>Average values</i>		CAST test set	
		mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1
ALADIN	BAM-FG	59.549	8.085	38.423	1.7025	48.712	3.560	48.895	4.449	14.066	0.145						
→ Fused	BAM-FG	53.941	4.485	32.686	0.550	42.592	2.395	43.073	2.477	12.716	0.103						
ALADIN-ViT	BAM-FG	16.823	0.270	9.996	0.058	12.493	0.060	13.104	0.129	39.541	0.960						
SAE	BAM-FG	51.600	16.100	28.500	4.000	28.814	4.643	36.305	8.248	24.001	3.622						
Ours	BAM-FG	85.955	58.108	67.699	24.967	74.355	27.154	76.003	36.743	45.352	8.963						
Ours	BBST-4M	90.965	69.523	80.861	42.803	84.953	45.258	<b>85.593</b>	<b>52.528</b>	49.336	9.003						
Ours + SS	BBST-4M	90.224	69.950	79.053	39.638	82.308	38.403	83.862	49.330	<b>60.811</b>	<b>14.640</b>						
Ours + SS, no NST	BBST-4M	11.801	0.665	5.080	0.070	8.554	0.475	8.478	0.403	4.097	0.098						

**Table 2:** Style representation strength, compared to baselines. Higher values are better. Computed over the ALADIN (and its fused variant) [22], ALADIN-ViT [19], and Swapping Autoencoders [17] baselines.

Model	Dataset	NeAT test set				PAMA test set				SANet test set				<i>Average values</i>		CAST test set	
		mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1	mAP	IR-1
ALADIN	BAM-FG	5.547	0	8.523	0	4.642	0	6.237	<b>0</b>	23.279	0.045						
→ Fused	BAM-FG	11.008	0	19.239	0.013	8.920	0	13.056	0.004	41.015	0.183						
ALADIN-ViT	BAM-FG	15.058	0.023	10.081	0.028	8.097	0.003	11.079	0.018	7.485	<b>0.010</b>						
SAE	BAM-FG	2.198	0.003	3.815	0.005	2.758	0	2.924	0.003	5.354	<b>0.010</b>						
Ours	BAM-FG	1.523	0	1.575	0	1.630	0	1.576	<b>0</b>	6.974	0.063						
Ours	BBST-4M	1.491	0	1.427	0	1.652	0	1.523	<b>0</b>	7.463	0.070						
Ours + SS	BBST-4M	1.381	0	1.461	0	1.629	0	<b>1.490</b>	<b>0</b>	<b>4.626</b>	0.018						
Ours + SS, no NST	BBST-4M	10.789	0.355	12.034	0.345	5.969	0.010	9.597	0.237	31.687	3.443						

**Table 3:** Comparisons with the same baselines as Table 2, measuring similarity between content images. In this table, a higher value is worse, as it represents higher content entanglement. Computed over the ALADIN (and its fused variant) [22], ALADIN-ViT [19], and Swapping Autoencoders [17] baselines.

Due to the cross-NST approach in our style representation learning signal, the 2 million images in each of the content and style splits of BBST-4M lead to a synthetic dataset of an effective 4 *trillion* images. This creates a practically limitless combination of style and content during training.



**Fig. 4:** High resolution real style image from the test set (left), randomly cropped and rotated, to gather 25 stylistically similar crops (right) with minimal semantic entanglement. Representative example.

#### 4.1 Data

To evaluate how well the model represents specifically disentangled style information, we need to consider test data that is also wholly disentangled. We also apply our synthetic NST dataset creation methodology for the test set, ensuring no overlap with training data. Using 400 new style images from Behance, and 100 new content images from Flickr, we extend BBST-4M with synthetic stylized images. We create 40k images, stylizing each content image with each style.

We use NeAT, PAMA, and SANet variants of this test set to evaluate the generalization of style representation independent of any systematic signatures specific to any NST method - visible or otherwise. We select these three NST methods given their fast and leading stylization qualities in literature. We manually selected the source style images to ensure a high variety of styles and no duplicates or styles too similar by manually inspecting style-based image retrieval for each style image as a query over the remaining test set style corpus.

We additionally run the same tests but using CAST, an NST method held out, not seen during training. We use this to evaluate generalization to unseen synthetic stylization methods, with unseen model biases.

We'd like to stress that although we are evaluating the disentanglement capabilities of our technique on synthetic data, we show through our other experiments that our representation carries over its strengths to real data, also. The synthetic data is only used for evaluating disentanglement properties, as real disentangled data is not available.

## 4.2 Metrics

We build our evaluation pipeline around image retrieval using these synthetic test sets. The primary metric we measure is mean average precision (mAP). We calculate the mAP by considering the other 99 content images stylized with the same style as positives and those stylized with the different style images as negatives. For each image in the test set, we re-arrange the remaining test set images, sorted by similarity in the style embedding space to this query image. We also use the Instance Retrieval (IR) [22] metric from ALADIN, for which we measure the *top-k* accuracy of retrieving the source style image from the corpus. We remove the other stylized images of the same style from the corpus, leaving just the query and source images to share the style for a given search.

## 4.3 Real data evaluation

The primary difficulty facing this work is the lack of fully disentangled data, as discussed in earlier sections. Therefore, the best way to perform evaluation is through synthetic data, as we have done during training. However, in an effort to include evaluations using non-synthetic data, we devise a secondary evaluation strategy based on image retrieval using crops of real style images. Style is a global property in an image, thus most of the image will generally have uniform style. Semantic content however is generally represented in sub-located areas of the image. Therefore, we can attempt to remove content information from an image by sub-dividing it into small non-overlapping crops, which will also maintain style information. Furthermore, we can randomly rotate the crops to further disassociate the semantics, while still keeping style information, as style tends to be rotation invariant. We use 5x5 image crops for this, dividing our high resolution real style images into 25 non-overlapping, randomly rotated crops, as visualized in Fig 4. Despite these measures, some images may still contain similar semantic cues - but this is likely the best way to use real images for evaluation, with minimal entanglement.

Model	Dataset	Real image evaluation	
		mAP	IR-1
ALADIN	BAM-FG	0.519	0.840
→ Fused	BAM-FG	0.539	0.380
ALADIN-ViT	BAM-FG	0.483	9.450
SAE	BAM-FG	0.503	1.422
Ours	BAM-FG	0.510	8.040
Ours	BBST-4M	0.566	13.470
Ours + SS	BBST-4M	<b>0.575</b>	<b>14.060</b>
Ours + SS, no NST	BBST-4M	0.463	6.63

**Table 4:** Evaluation using real style image crops. Higher is better. Note that despite rotated crops, there is still some inherent content entanglement in this experiment.

We evaluate using image retrieval, but this time we use each crop image as a query, and we measure the mean average precision and IR- $k$  for retrieving the remaining crops.

#### 4.4 Ablations

Table 1 contains ablations where we experiment with the NST methods used for driving the style learning signal during training. Using more than one method is essential for generalizing the style representation. Using only one NST method risks modeling specific artifacts of that method. Future work could incorporate additional new NST methodologies as the field advances. Our final model uses all 3 methods: NeAT, PAMA, and SANet. From our experiments, we note that NeAT and SANet add the most value to the model quality, but we keep PAMA for the added method de-biasing. We explore these ablations using only a ViT backbone, such that we can also draw fair comparisons to literature using the same architecture, ALADIN-ViT.

We additionally explore adding an additional self-supervised loss, using only real image data, without the NST. We denote these experiments as "+ SS" in Tables 2, 3, and 4. This added loss is in addition to the NST-driven learning signal, simply adding a contrastive loss treating two crops of a real style image as positives, with the remaining batch items as negatives, similar to other self-supervised works in literature [2]. This additional loss helps to reduce the entanglement, as shown in Table 3, and it helps with generalization to unseen NST evaluation (CAST - Table 2) and real image evaluation (Table 4). The average scores across the three NST methods is lower, due to the overall training objective having a less strong focus on optimising for any biases introduced by specific models.

Finally, we also train our model with *only* the self-supervised crop training, without our core NST-driven approach, to signify its value. Table 2 shows this approach does not encode style information as strongly as the other methods. Table 3 shows competitive content entanglement with some baselines, but much stronger entanglement than our same method where NST was used - showing the disentanglement strength of our NST approach.

#### 4.5 Baselines

We compare our model against baselines in Table 2, demonstrating also how the BAM-FG dataset needs to be revised as a style-only dataset. We train Swapping Autoencoders (SAE) [17], and our model with BAM-FG for fair comparison.

The accuracy rankings according to the literature are flipped in our measurements as we test with purely disentangled labels. ALADIN scores highest despite being the first relevant model because its features are extracted globally. Their fusion includes ResNet embeddings, which introduce semantic entanglement as a by-product of the higher BAM-FG style scores. ALADIN-ViT scores are even lower on our disentangled test set due to a lack of explicitly global

features, therefore more intensely focusing on localized and, thus, typically more semantic information.

In Table 3, we repeat our evaluation found in Table 2, but instead of measuring the style retrieval in our test sets, we measure *content* retrieval. Like the style evaluation, we compute mAP by using a query image and evaluating retrieval of the corpus concerning all the other 399 images of the same content, but stylized with different style images. So, measuring semantics-based image retrieval, irrespective of style.

For IR- $k$ , we filter out the other stylized versions of the content image stylized in the query and measure the retrieval of the original un-stylized content image. We run this set of evaluations to measure how strongly the style embeddings capture content/semantic information. Supporting our previous explanations, ALADIN’s fused and ViT variants each capture more semantic information. Our work improves upon this, as our style embeddings perform much more poorly at retrieving images with the same content.

#### 4.6 Style-based image retrieval

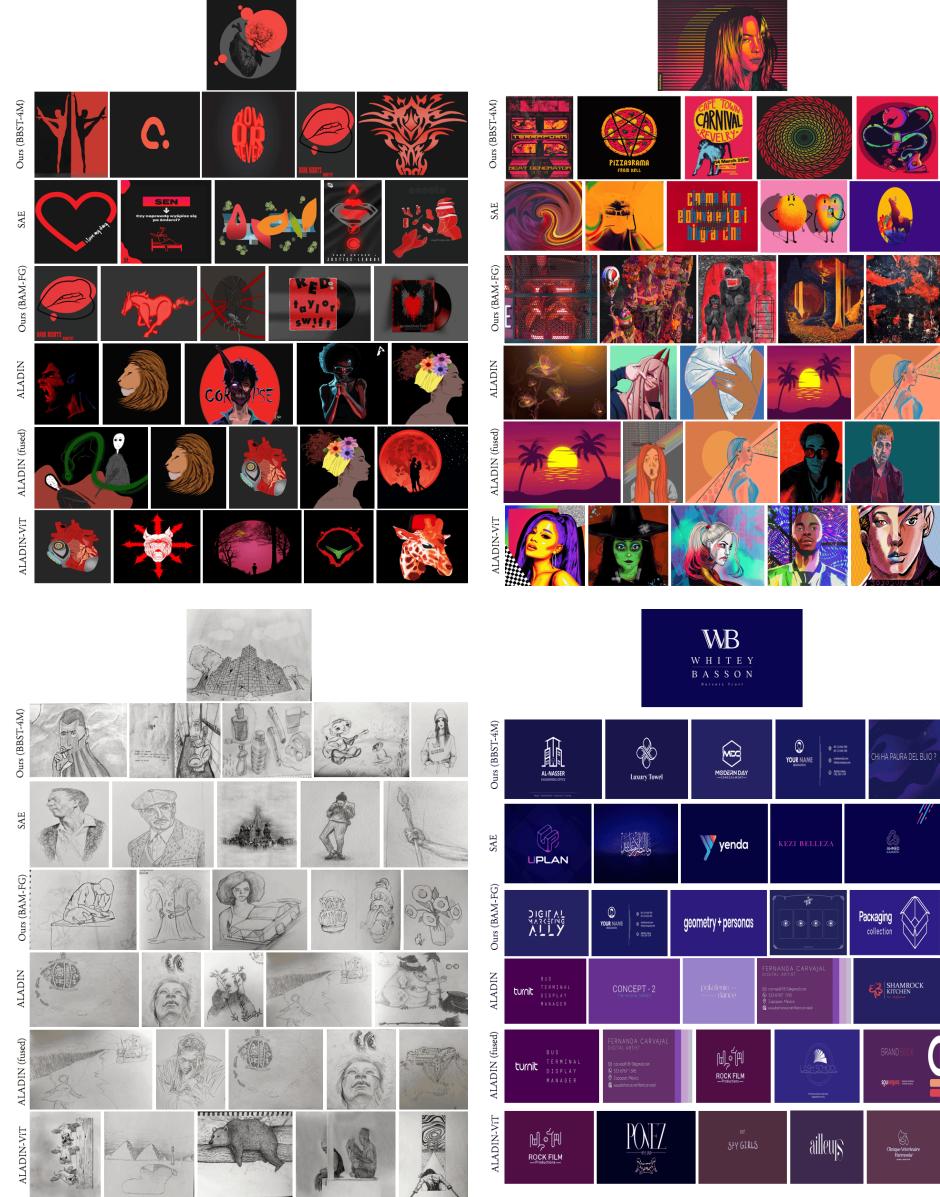
We visualize in Figure 5 a comparison of our method to baseline methods in the literature for style-based image retrieval. We show our approach trained on both BBST-4M and BAM-FG. We perform the retrieval over a corpus of 500k images from BBST-4M.

Our results are comparable regarding visual features, improving slightly on the color consistency of retrieved results (bottom right). However, there is less semantic consistency between our model’s query and outcomes, especially compared to the previous ALADIN models. In the top left of Figure 6, the *heart* in the query image retrieves some other heart-related imagery in baselines. In the top right, baseline recovered results contain faces and character designs, also present in the query.

### 5 Multimodal vision-language learning

We use style embeddings from our proposed model to learn a joint multimodal representation between style and language. We replicate the work in StyleBabel [19], where style tags attributed to images can be used as labels for this task. We replace their ALADIN-ViT vision backbone with ours, and we similarly train an MLP, joining these style embeddings to text embeddings extracted using CLIP [18] through contrastive learning. We aim to measure how our new style embeddings can be used in this multimodal setting.

We measure a WordNet score of 0.329, which beats their baseline CLIP WordNet score of 0.215, but does not beat the ALADIN-ViT WordNet score of 0.352. This may be due to the inherent content/style entanglement of the style tags in StyleBabel, which itself is not strictly disentangled. There exist several tags in StyleBabel, such as *t-shirt design*, *interior design*, *fashion photography*, which do describe the style, but in a context that is also grounded in semantics. By



**Fig. 5:** Style-based image retrieval comparison between our method variants and previous literature.

ALADIN-ViT	neon-like, warm narrative toned lighting, driven, manga, posh, promoting, layered composition, light trail, strong Japanese, car-triangle composition, outline, all cap toon drawing, sition, branding painting, hand-word sound graphic art package made artwork effects, subtle painting, illus-colors trative	regulated layout, mixed media, color splash, marker drawing, watercolor expressionist, outline, all cap toon drawing, sition, branding painting, hand-word sound graphic art package made artwork effects, subtle painting, illus-colors trative	trustful, housing, colorful drawing, storyboarding, chiaroscuro, black vigorous, past, posh, housing abstract and red, red, dark documentary architecture ren-abstract artwork, panel, story picture shot, dark con-der, solarpunk nonobjective, storyboarding trast	glass, trustful, soft, soft color, comic book art, vibe, explosion, high contrast, architectural soft and bright black and white dark space suburban, doc-landscape, pen-color, feminine, art, comic art, umentary shot, tagonal, pointy soft colors documentary	art, interpretive, interpretive, abstract artwork, panel, story cubism
Ours (fused)	dark image, documentary, chiaroscuro, black vigorous, past, posh, housing abstract and red, red, dark documentary architecture ren-abstract artwork, panel, story picture shot, dark con-der, solarpunk nonobjective, storyboarding trast	trustful, housing, colorful drawing, storyboarding, chiaroscuro, black vigorous, past, posh, housing abstract and red, red, dark documentary architecture ren-abstract artwork, panel, story picture shot, dark con-der, solarpunk nonobjective, storyboarding trast	glass, trustful, soft, soft color, comic book art, vibe, explosion, high contrast, architectural soft and bright black and white dark space suburban, doc-landscape, pen-color, feminine, art, comic art, umentary shot, tagonal, pointy soft colors documentary	glass, trustful, soft, soft color, comic book art, vibe, explosion, high contrast, architectural soft and bright black and white dark space suburban, doc-landscape, pen-color, feminine, art, comic art, umentary shot, tagonal, pointy soft colors documentary	art, interpretive, interpretive, abstract artwork, panel, story cubism
Ours	flame, spark, dark contrasted, flame, spark, dark contrasted, vibe, explosion, high contrast, architectural dark space	trustful, soft, soft color, comic book art, vibe, explosion, high contrast, architectural dark space	glass, trustful, soft, soft color, comic book art, vibe, explosion, high contrast, architectural dark space	soft color, comic book art, vibe, explosion, high contrast, architectural dark space	comic, doodle art
ALADIN-ViT	cold hue, product- had material, measuring, tech- user input, 3d expressive, irreg- focused, product watery, bird eye nical sketch, building plan, in-ular angle, clean description, blue- view, nobody design sketch, ternet, gathering stroke, blue glow, based, digital sketch, typogra-formation, da- copy publication	had material, measuring, tech- user input, 3d expressive, irreg- focused, product watery, bird eye nical sketch, building plan, in-ular angle, clean description, blue- view, nobody design sketch, ternet, gathering stroke, blue glow, based, digital sketch, typogra-formation, da- copy publication	material, measuring, tech- user input, 3d expressive, irreg- focused, product watery, bird eye nical sketch, building plan, in-ular angle, clean description, blue- view, nobody design sketch, ternet, gathering stroke, blue glow, based, digital sketch, typogra-formation, da- copy publication	material, measuring, tech- user input, 3d expressive, irreg- focused, product watery, bird eye nical sketch, building plan, in-ular angle, clean description, blue- view, nobody design sketch, ternet, gathering stroke, blue glow, based, digital sketch, typogra-formation, da- copy publication	material, measuring, tech- user input, 3d expressive, irreg- focused, product watery, bird eye nical sketch, building plan, in-ular angle, clean description, blue- view, nobody design sketch, ternet, gathering stroke, blue glow, based, digital sketch, typogra-formation, da- copy publication
Ours (fused)	mystical, fantasy reflective, com- sketch scamp, uxui design, poverty, struggle, concept art, fan- mercial shot, scamp, sketch design interface, slum, evocative, tasy art, aqua, rough texture, work, sketch, user interface, documented fantasy painting geometric shape, sketched line ui instructional geometric line	com- sketch scamp, uxui design, poverty, struggle, concept art, fan- mercial shot, scamp, sketch design interface, slum, evocative, tasy art, aqua, rough texture, work, sketch, user interface, documented fantasy painting geometric shape, sketched line ui instructional geometric line	scamp, uxui design, poverty, struggle, concept art, fan- mercial shot, scamp, sketch design interface, slum, evocative, tasy art, aqua, rough texture, work, sketch, user interface, documented fantasy painting geometric shape, sketched line ui instructional geometric line	scamp, uxui design, poverty, struggle, concept art, fan- mercial shot, scamp, sketch design interface, slum, evocative, tasy art, aqua, rough texture, work, sketch, user interface, documented fantasy painting geometric shape, sketched line ui instructional geometric line	scamp, uxui design, poverty, struggle, concept art, fan- mercial shot, scamp, sketch design interface, slum, evocative, tasy art, aqua, rough texture, work, sketch, user interface, documented fantasy painting geometric shape, sketched line ui instructional geometric line
Ours	layered composi- embossed, small idea, blended, page layout, high-contrast, tion, aqua, mysti- shape, cup, stone pen and pencil, image of article, retouched, noir, cal, fantasy paint- image, light grey sketch, quick ing, digital print	embossed, small idea, blended, page layout, high-contrast, shape, cup, stone pen and pencil, image of article, retouched, noir, image, light grey sketch, quick digital print	blended, page layout, high-contrast, shape, cup, stone pen and pencil, image of article, retouched, noir, image, light grey sketch, quick digital print	blended, page layout, high-contrast, shape, cup, stone pen and pencil, image of article, retouched, noir, image, light grey sketch, quick digital print	blended, page layout, high-contrast, shape, cup, stone pen and pencil, image of article, retouched, noir, image, light grey sketch, quick digital print
ALADIN-ViT	cool hue, mockup, hand drawn, colorful, bright, curved line, col-fading, sans serif magazine book-changing pro-bright bold, flow, orful drawing, and serif, red let layout, blur portion, sketchy, color-heavy mark making, highlighting, themed, graphic saddening, line blue ink drawing, thin letter, type- layout drawing thin stroke face	mockup, hand drawn, colorful, bright, curved line, col-fading, sans serif magazine book-changing pro-bright bold, flow, orful drawing, and serif, red let layout, blur portion, sketchy, color-heavy mark making, highlighting, themed, graphic saddening, line blue ink drawing, thin letter, type- layout drawing thin stroke face	hand drawn, colorful, bright, curved line, col-fading, sans serif magazine book-changing pro-bright bold, flow, orful drawing, and serif, red let layout, blur portion, sketchy, color-heavy mark making, highlighting, themed, graphic saddening, line blue ink drawing, thin letter, type- layout drawing thin stroke face	hand drawn, colorful, bright, curved line, col-fading, sans serif magazine book-changing pro-bright bold, flow, orful drawing, and serif, red let layout, blur portion, sketchy, color-heavy mark making, highlighting, themed, graphic saddening, line blue ink drawing, thin letter, type- layout drawing thin stroke face	hand drawn, colorful, bright, curved line, col-fading, sans serif magazine book-changing pro-bright bold, flow, orful drawing, and serif, red let layout, blur portion, sketchy, color-heavy mark making, highlighting, themed, graphic saddening, line blue ink drawing, thin letter, type- layout drawing thin stroke face
Ours (fused)	editorial design line, ink work, pastel, pastoral, abstract line, spontaneous, editorial, editorial ink, ink drawing, black tea paint-mark-making, beginner, sketch- mockup, editorial outline ing, art appre-line, fine, fluid book, paper work, editorial ciation, colorful line pattern, light mockup drawing pencil work	design line, ink work, pastel, pastoral, abstract line, spontaneous, editorial, editorial ink, ink drawing, black tea paint-mark-making, beginner, sketch- mockup, editorial outline ing, art appre-line, fine, fluid book, paper work, editorial ciation, colorful line pattern, light mockup drawing pencil work	design line, ink work, pastel, pastoral, abstract line, spontaneous, editorial, editorial ink, ink drawing, black tea paint-mark-making, beginner, sketch- mockup, editorial outline ing, art appre-line, fine, fluid book, paper work, editorial ciation, colorful line pattern, light mockup drawing pencil work	design line, ink work, pastel, pastoral, abstract line, spontaneous, editorial, editorial ink, ink drawing, black tea paint-mark-making, beginner, sketch- mockup, editorial outline ing, art appre-line, fine, fluid book, paper work, editorial ciation, colorful line pattern, light mockup drawing pencil work	design line, ink work, pastel, pastoral, abstract line, spontaneous, editorial, editorial ink, ink drawing, black tea paint-mark-making, beginner, sketch- mockup, editorial outline ing, art appre-line, fine, fluid book, paper work, editorial ciation, colorful line pattern, light mockup drawing pencil work
Ours	editorial, editorial animation pastel, art ther-fluid, blue.cut printmaking, del- design, readable, sketch, fine-ap, fine art, tra-paper, plump, icate type, vari- editorial mockup, line drawing, ditional illustra-material sample ous printed pa- professional style figurative draw-tion, child illus- per work, limo, ing, line drawing tration illustration, sketch of cartoon character	pastel, art ther-fluid, blue.cut printmaking, del- design, readable, sketch, fine-ap, fine art, tra-paper, plump, icate type, vari- editorial mockup, line drawing, ditional illustra-material sample ous printed pa- professional style figurative draw-tion, child illus- per work, limo, ing, line drawing tration illustration, sketch of cartoon character	pastel, art ther-fluid, blue.cut printmaking, del- design, readable, sketch, fine-ap, fine art, tra-paper, plump, icate type, vari- editorial mockup, line drawing, ditional illustra-material sample ous printed pa- professional style figurative draw-tion, child illus- per work, limo, ing, line drawing tration illustration, sketch of cartoon character	pastel, art ther-fluid, blue.cut printmaking, del- design, readable, sketch, fine-ap, fine art, tra-paper, plump, icate type, vari- editorial mockup, line drawing, ditional illustra-material sample ous printed pa- professional style figurative draw-tion, child illus- per work, limo, ing, line drawing tration illustration, sketch of cartoon character	pastel, art ther-fluid, blue.cut printmaking, del- design, readable, sketch, fine-ap, fine art, tra-paper, plump, icate type, vari- editorial mockup, line drawing, ditional illustra-material sample ous printed pa- professional style figurative draw-tion, child illus- per work, limo, ing, line drawing tration illustration, sketch of cartoon character

**Fig. 6:** Please zoom for more image detail. Zero-shot automatic style tagging comparison, between ALADIN-ViT, our model, and our fused variant, joining our disentangled embeddings with ALADIN-ViT. We show the top 5 tags for each image.

explicitly not encoding semantic information in our style embeddings, such tags are more difficult to retrieve.

Inspired by the fusing [22] of the complementary ALADIN and ResNet [7] embeddings, we explore a fusion of our disentangled style embeddings with ALADIN-ViT embeddings, which contain some semantic information. We extract and concatenate embeddings from both models and use this dual-model embedding as a style embedding for learning the joint vision+language multimodal embedding against CLIP. We achieve a state-of-the-art WordNet score of **0.415** on StyleBabel tags. We use the same test split for our measurements. In Figure 6, we visualize automatic zero-shot style tagging with ALADIN-ViT (baseline), our model, and our model fused with ALADIN-ViT over images from the StyleBabel test set.

## 6 Training details

We train our model with the NeAT, PAMA, and SANet NST methods for roughly 3 days on a single A100 GPU until convergence. We stylized images for training using 512px resolution, which we downsample to 256x256 for the VGG branch and 224x224 as needed for ViT-B\_16, using the same ViT as ALADIN. We disable the prior blurring in NeAT for speed. We use the Adam optimizer, and a target batch size of 1024 via logit accumulation. We decay the learning rate by 0.999875 every 100 iterations.

## 7 Limitations and Conclusions

We explore a novel learning methodology for artistic style, achieving stronger disentanglement. We demonstrate the value of this by further achieving state-of-the-art multimodal vision+language learning on StyleBabel tags.

Our approach relies on NST as a strong driver of style consistency, thus limiting us to the capabilities of such automated stylization methods. However, as NST methods continue to improve, so can our method in how well it can capture style. The better the artistic stylization process can be modeled by NST models, the better our technique can be trained to capture that style, by including the technique in our pipeline.

For practicality, we can only rely on fast, feed-forward approaches. Optimization and diffusion based techniques are too slow to dynamically synthesize training data during the training loop unless this data is synthesized ahead of time, trading off variety of artistic styles and high storage costs.

It should be noted that we specifically focus on disentanglement of style embeddings. As shown in the multi-modal experiments, some use cases may suffer lower accuracy if a certain degree of entanglement is actually necessary. Downstream tasks need to evaluate this trade-off carefully, depending on the priority of disentanglement or strength of encoded style information. Further work may explore scaling the ViT branch, and explore variants with a more global context.

## References

1. Chen, H., Zhao, L., Wang, Z., Ming, Z.H., Zuo, Z., Li, A., Xing, W., Lu, D.: Artistic style transfer with internal-external learning and contrastive learning. In: Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (eds.) Neural Information Processing Systems (2021)
2. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. arXiv preprint arXiv:2002.05709 (2020)
3. Collomosse, J., Bui, T., Wilber, M., Fang, C., Jin, H.: Sketching with style: Visual search with sketches and aesthetic context. In: Proc. ICCV (2017)
4. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
5. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: Proc. CVPR. pp. 2414–2423 (2016)
6. Ham, C., Tarres, G.C., Bui, T., Hays, J., L., Z., Collomosse, J.: Cogs: Controllable generation and search from sketch and style (2022). <https://doi.org/10.48550/ARXIV.2203.09554>, <https://arxiv.org/abs/2203.09554>
7. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. CoRR **abs/1512.03385** (2015), <http://arxiv.org/abs/1512.03385>
8. Huang, X., Belongie, S.J.: Arbitrary style transfer in real-time with adaptive instance normalization (2017), <http://arxiv.org/abs/1703.06868>
9. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: ECCV (2018)
10. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution (2016). <https://doi.org/10.48550/ARXIV.1603.08155>, <https://arxiv.org/abs/1603.08155>
11. Kalischeck, N., Wegner, J.D., Schindler, K.: In the light of feature distributions: moment matching for neural style transfer. CoRR **abs/2103.07208** (2021), <https://arxiv.org/abs/2103.07208>
12. Karayev, S., Trentacoste, M., Han, H., Agarwala, A., Darrell, T., Hertzmann, A., Winnemoller, H.: Recognizing image style. In: Proc. BMVC (2014)
13. Li, C., Wand, M.: Precomputed real-time texture synthesis with markovian generative adversarial networks (2016). <https://doi.org/10.48550/ARXIV.1604.04382>, <https://arxiv.org/abs/1604.04382>
14. Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.: Universal style transfer via feature transforms. CoRR **abs/1705.08086** (2017), <http://arxiv.org/abs/1705.08086>
15. Luo, X., Han, Z., Yang, L., Zhang, L.: Consistent style transfer. CoRR **abs/2201.02233** (2022), <https://arxiv.org/abs/2201.02233>
16. P., D.Y., L., K.H.: Arbitrary style transfer with style-attentional networks. CoRR (2018)
17. Park, T., Zhu, J.Y., Wang, O., Lu, J., Shechtman, E., Efros, A.A., Zhang, R.: Swapping autoencoder for deep image manipulation. In: Neural Information Processing Systems (2020)
18. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. arXiv preprint arXiv:2103.00020 (2021)

19. Ruta, D., Gilbert, A., Aggarwal, P., Marri, N., Kale, A., Briggs, J., Speed, C., Jin, H., Faieta, B., Filipkowski, A., Lin, Z., Collomosse, J.: Stylebabel: Artistic style tagging and captioning (2022). <https://doi.org/10.48550/ARXIV.2203.05321>, <https://arxiv.org/abs/2203.05321>
20. Ruta, D., Gilbert, A., Collomosse, J., Shechtman, E., Kolkin, N.: Neat: Neural artistic tracing for beautiful style transfer (2023)
21. Ruta, D., Gilbert, A., Motiian, S., Faieta, B., Lin, Z., Collomosse, J.: Hypernst: Hyper-networks for neural style transfer (2022). <https://doi.org/10.48550/ARXIV.2208.04807>, <https://arxiv.org/abs/2208.04807>
22. Ruta, D., Motiian, S., Faieta, B., Lin, Z., Jin, H., Filipkowski, A., Gilbert, A., Collomosse, J.: Aladin: All layer adaptive instance normalization for fine-grained style similarity. In: Proc. ICCV (2021)
23. Sheng, L., Lin, Z., Shao, J., Wang, X.: Avatar-net: Multi-scale zero-shot style transfer by feature decoration. In: Computer Vision and Pattern Recognition (CVPR), 2018 IEEE Conference on (2018)
24. Tarrés, G.C., Ruta, D., Bui, T., Collomosse, J.: Parasol: Parametric style control for diffusion image synthesis (2023)
25. Ulyanov, D., Lebedev, V., Vedaldi, A., Lempitsky, V.: Texture networks: Feed-forward synthesis of textures and stylized images (2016). <https://doi.org/10.48550/ARXIV.1603.03417>, <https://arxiv.org/abs/1603.03417>
26. Wilber, M.J., Fang, C., Jin, H., Hertzmann, A., Collomosse, J., Belongie, S.: Bam! the behance artistic media dataset for recognition beyond photography. In: Proc. ICCV (2017)
27. Zhang, Y., Tang, F., Dong, W., Huang, H., Ma, C., Lee, T.Y., Xu, C.: Domain enhanced arbitrary image style transfer via contrastive learning. In: ACM SIGGRAPH (2022)