# **Understanding Art through Multi-Modal Retrieval in Paintings**

## Noa Garcia Benjamin Renoust Yuta Nakashima Institute for Datability Science Osaka University, Japan

{noagarcia, renoust, n-yuta}@ids.osaka-u.ac.jp

## **Abstract**

In computer vision, visual arts are often studied from a purely aesthetics perspective, mostly by analysing the visual appearance of an artistic reproduction to infer its style, its author, or its representative features. In this work, however, we explore art from both a visual and a language perspective. Our aim is to bridge the gap between the visual appearance of an artwork and its underlying meaning, by jointly analysing its aesthetics and its semantics. We introduce the use of multi-modal techniques in the field of automatic art analysis by 1) collecting a multi-modal dataset with fine-art paintings and comments, and 2) exploring robust visual and textual representations in artistic images.

## 1. Introduction

The large-scale digitisation of artworks from collections all over the world has opened the opportunity to study art from a computer vision perspective, by building tools to help in the conservation and dissemination of cultural heritage. Some of the most promising work on this direction involves the automatic analysis of paintings, in which computer vision techniques are applied to study the content [5], the style [4, 18], or to classify the attributes [16, 15] of a specific piece of art. In this way, art has been mostly studied from a visual perspective [2, 10, 17, 19, 14], and less attention has been paid to automatically analyse the underlying meaning of each painting. In this work, we aim to bridge the gap between the visual analysis and the high-level understanding of art, by proposing robust language and vision representations for multi-modal retrieval in paintings.

We first introduce a multi-modal dataset for visual arts, in which each image of a painting is associated with an artistic comment (Figure 1). Differently from multi-modal datasets in natural images, such as VQA [1], Visual Genome [12], and MS-COCO [13], the interpretation of art is strongly related to the artistic context of each artwork. This peculiarity is observed in the proposed dataset both in terms of images, through the use of style and composition,

#### View of Florence from Villa San Firenze, near San Miniato



This view of Florence is one of a number of views by Lear based upon on the spot sketches he produced in 1861.

#### Water Carriers



This painting was inspired by the painter's travels in Italy. The costume of the two girls and the landscape suggests the Amalfi coast, or Capri as the setting of the scene.

#### Ships Moored Off a Rocky Coastline



This landscape depicts ships moored off a rocky coastline with fishermen unloading their catch.

#### Still-Life



This painting depicts a still-life of grapes, cherries, peaches and other fruit in a basket, with a rose and a dragonfly on a stone ledge.

Figure 1. Examples of paintings and comments in SemArt dataset.

and language, through the use of references.

To leverage these differences and study art from a semantics perspective, we propose to enhance robust visual and language representations with artistic attributes. The enhanced representations are projected into a multi-modal artistic space in which image and text coexist. By fine-tuning the multi-modal representations in the art domain, paintings and comments that are semantically similar are represented closer than dissimilar samples.

The quality of the proposed multi-modal artistic space is evaluated as a retrieval task, in which given a painting image, the most representative comment from the collection must be found, and vice-versa. Multi-modal retrieval allows us to discriminate whether the language and visual representations capture the sufficient artistic insights to match corresponding paintings and comments together. In the evaluation, our method achieves results only 0.059 below human accuracy.

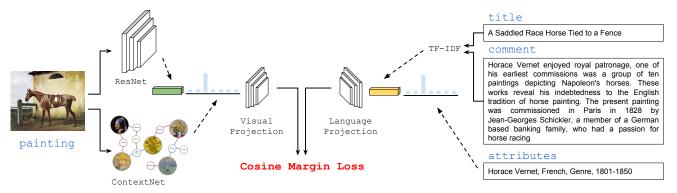


Figure 2. Proposed visual and language representations for multi-modal retrieval in art.

## 2. SemArt Dataset

Existing datasets in art analysis, such as PRINTART [3], Painting-91 [11], Rijksmuseum [16] or Art500k [15], are mostly annotated with attribute labels, such as author, style, or timeframe. Although this information is crucial in the analysis of visual arts, it does not provide enough insights for understanding the high-level semantics of fine-art paintings. To jointly study language and vision in art, we introduce SemArt<sup>1</sup>, a dataset for semantic art understanding.

SemArt contains 21,384 reproductions of European paintings collected from the Web Gallery of Art<sup>2</sup>, randomly split into training, validation, and test sets with 19,244, 1,069 and 1,069 samples, respectively. Each image is annotated with its main attributes – author, title, date, technique, type, school and timeframe<sup>3</sup> – and with a natural language comment. Interestingly, comments involve not only a description of the elements in the scene but also references to its technique, author or context. Some examples are shown in Figure 1, and a complete analysis of the dataset can be found in [7].

## 3. Multi-Modal Representations

To jointly represent aesthetics and semantics in art, we propose to project robust visual and language representations enhanced with artistic attributes into a multi-modal artistic space, as depicted in Fig. 2. In total, we combine four different representations, which are described below.

**Language representation** The language representation captures the insights of the high-level semantics of artworks by encoding both titles and artistic comments. Titles are encoded as a term frequency - inverse document frequency (tf-idf) vector,  $\mathbf{v}_{\text{tit}} \in \mathbb{R}^{N_t}$ , with  $N_t = 9,092$  being the size of the title vocabulary built with the alphabetic words in the titles in the training set. Comments are encoded as another

tf-idf vector,  $\mathbf{v}_{\text{com}} \in \mathbb{R}^{N_c}$ , with  $N_c = 9,708$  being the comments vocabulary built with the alphabetic words occurring at least ten times in the training set. The language representation is obtained by  $\mathbf{v}_{\text{lang}} = \mathbf{v}_{\text{tit}} \oplus \mathbf{v}_{\text{com}}$ , where  $\oplus$  is vector concatenation.

**Language attributes** Attributes capture the essential information of a painting, such as its painter or its date of creation. We encode the type, school, timeframe, or author labels in the dataset as a one-hot vector,  $\mathbf{v}_{\text{att}} \in \mathbb{R}^c$ , with c being the number of labels in each attribute.

**Visual representation** The visual representation captures the visual appearance of paintings. Painting images are scaled down to 256 pixels per side, randomly cropped into  $224 \times 224$  patches and fed into a ResNet50 [9], initialised with its standard pre-trained weights. Appearance is then represented by the output of the model as  $\mathbf{v}_{\text{vis}} \in \mathbb{R}^{1000}$ .

Image attributes From the painting image, we use a contextual network (ContextNet) [6] to predict the artistic attributes. ContextNet is composed by two core modules, as depicted in Fig. 3: a ResNet<sup>4</sup> [9], which obtains the visual information of the image, and a knowledge graph, which captures the contextual relationships of the painting. The visual encoding from the ResNet is further input into an attribute classifier<sup>5</sup> for predicting the artistic attributes, and into an encoder module<sup>6</sup> for projecting the visual encoding into the knowledge graph space. The knowledge graph is built by connecting the training paintings in SemArt with their attributes, and its nodes are encoded into a 128-dimensional graph representations using node2vec [8].

At training time, we compute the cross-entropy loss function,  $\ell_c$ , between the predicted attribute and the real attribute of the painting, and the smooth L1 loss function,  $\ell_e$ ,

<sup>&</sup>lt;sup>1</sup>Available at http://noagarciad.com/SemArt/

<sup>&</sup>lt;sup>2</sup>https://www.wga.hu/

<sup>&</sup>lt;sup>3</sup>Periods of 50 years evenly distributed between 801 and 1900.

<sup>&</sup>lt;sup>4</sup>Without the last fully connected layer.

 $<sup>^5</sup>$ A n-dimensional fully-connected layer with ReLU and softmax, where n is the number of classes for the predicted attribute.

<sup>&</sup>lt;sup>6</sup>A 128-dimensional fully-connected layer.

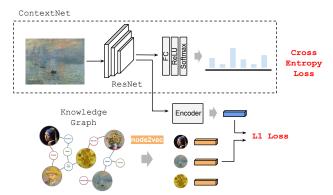


Figure 3. ContextNet predicts the painting attributes, such as type, school, timeframe, or author, by fine-tuning a ResNet model based on the information captured by an artistic Knowlegde Graph.

between the encoder output and the graph embedding from the knowledge graph. The ContextNet weights are learnt by jointly optimising both losses as:

$$\mathcal{L} = \lambda_c \sum_{j=1}^{N} \ell_{c,j} + \lambda_e \sum_{j=1}^{N} \ell_{e,j}$$
 (1)

where  $\lambda_c$  and  $\lambda_e$  are parameters that weight the contribution of the classification and the encoder modules, respectively, and N is the number of training samples. To predict the painting attribute, the graph computation part is removed, and the attribute is predicted as the maximum value of the output of the ContextNet classifier, represented as  $\mathbf{v}_{\text{ctx}} \in \mathbb{R}^c$ , with c being the number of labels in the attribute.

## 4. Multi-Modal Projections

To learn the relationship between the visual attributes from the paintings and the semantics from the comments, we project the multi-modal representations from paintings and comments into a multi-modal artistic space. We define the vectors  $\mathbf{p} \in \mathbb{R}^{1000+c}$  and  $\mathbf{q} \in \mathbb{R}^{N_t+N_c+c}$  as the joint representation of visual and image attributes, and language and language attributes, respectively:

$$\begin{aligned} \mathbf{p} &= \mathbf{v}_{vis} \oplus \mathbf{v}_{ctx} \\ \mathbf{q} &= \mathbf{v}_{lang} \oplus \mathbf{v}_{att} \end{aligned}$$

The two joint representation vectors are projected into a multi-modal artistic space using the non-linear functions  $f(\cdot)$  and  $g(\cdot)$ , respectively, which are implemented with a 128-dimensional fully connected layer followed by a tanh activation function and a  $\ell_2$ -normalisation layer. The whole model, except for the ContextNet which is previously finetuned and frozen, is trained end-to-end using both matching and non-matching pairs of samples from the training set. The loss is computed as a cosine margin loss function:

$$\operatorname{Loss}(\mathbf{p}_i, \mathbf{q}_j) = \begin{cases} 1 - \operatorname{sim}(f(\mathbf{p}_i), g(\mathbf{q}_j)), & \text{if } i = j \\ \max(0, \operatorname{sim}(f(\mathbf{p}_i), f_q(\mathbf{q}_j)) - \Delta), & \text{if } k \neq j \end{cases}$$

where the sub-indeces i and j are the representations for the i-th and j-th training sample,  $sim(\cdot, \cdot)$  is the cosine similarity between two vectors, and  $\Delta = 0.1$  is the margin. We use Adam optimiser with learning rate 0.0001.

### 5. Evaluation

To evaluate the quality of language and vision representations in art, we design the Text2Art challenge based on multi-modal retrieval, in which the aim is to find the most representative painting given an artistic comment, and vice versa, by ranking test samples according to their cosine similarity. In this way, the challenge evaluates whether the models capture enough of the insights and clues provided by the artistic comments to be able to match it to the correct painting. Results are reported with standard retrieval metrics: median rank (MR), and recall rate at K (R@K), with K being 1, 5 and 10.

Table 1 reports an ablation study when different combinations of the proposed representations are used. Vis&Lang uses the visual and language representations only. Att uses the vision, language, and language attribute (specified in brackets) as well as the output of a ResNet152 attribute classification network as a simplier image attribute representation. Note that the image attribute representation predicted in this way has not been informed with the graph representation from the knowledge graph. Finally, Att&ContextNet considers the four multi-modal representations from Section 3, including the context-aware classifier.

The best results are obtained when the four proposed representations are used, with attributes from language and image are given by the author. Att&ContextNet (Author) improves results by a 37.24% in average with respect to vision and language only, suggesting the importance of considering context when studying art. When compared against Att, the use of ResNet152 instead of the context-aware classifier performs better with type and school attributes, whereas Att&ContextNet is the best in timeframe and author.

In Table 2, we evaluate the proposed multi-modal art representations against human performance, where human evaluators were asked to choose between 10 paintings according to an artistic comment, title, author, type, school, and timeframe. We performed two evaluations: in the easy setup, the 10 paintings were chosen randomly, whereas in the difficult setup, the 10 paintings shared the same type (i.e. landascape, portrait, etc.). The multi-modal representations using the ContextNet reached values close to human accuracy, outperforming Vis&Lang by a 10.67% in the easy task and a 9.67% in the difficult task.

## 6. Conclusion

We addressed art understanding by introducing a new dataset of paintings with associated comments and explor-

	$\textbf{Text} \rightarrow \textbf{Image}$			$\mathbf{Image} \rightarrow \mathbf{Text}$				
Encoding	R@1	R@5	R@10	MR	R@1	R@5	R@10	MR
Vis⟪	0.164	0.384	0.505	10	0.162	0.366	0.479	12
Att Type	0.178	0.383	0.525	9	0.165	0.364	0.491	11
Att School	0.192	0.386	0.507	10	0.163	0.364	0.484	12
Att Tf	0.127	0.322	0.432	18	0.130	0.336	0.444	16
Att Author	0.236	0.451	0.572	7	0.204	0.440	0.535	8
Att&ContextNet Type	0.152	0.367	0.506	10	0.147	0.367	0.507	10
Att&ContextNet School	0.162	0.371	0.483	12	0.156	0.355	0.483	11
Att&ContextNet Tf	0.175	0.399	0.506	10	0.148	0.360	0.472	12
Att&ContextNet Author	0.247	0.477	0.581	6	0.212	0.446	0.563	7

Table 1. Results on the Text2Art Challenge when using vision and language only (Vis&Lang), when adding attributes (Attributes) and when adding the ContextNet classifier (Att&ContextNet).

Model	Easy	Difficult
Vis⟪	0.750	0.620
Att&Context	0.830	0.680
Human	0.889	0.714

Table 2. Multi-modal representations against humans.

ing multi-modal representations in art. Results showed that robust vision and language representations were able to capture the semantic content of paintings relatively well. However, performance was considerably improved when contextual information in the form of a knowledge graph was used to inform the model, which suggested the existence of a strong relationship between art and context. As a future work, we would like to pursue effort in the use of knowledge graph to connect vision and language. We could enhance ContextNet with more robust graph embedding techniques, such as StarSpace [20], as well as enhance the language representation with the knowledge graph attributes.

**Acknowledgement**: This work was partly supported by JSPS KAKENHI Grant No. 18H03264.

## References

- S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *ICCV*, 2015.
- [2] Y. Bar, N. Levy, and L. Wolf. Classification of artistic styles using binarized features derived from a deep neural network. In *ECCV Workshops*, 2014.
- [3] G. Carneiro, N. P. da Silva, A. Del Bue, and J. P. Costeira. Artistic image classification: An analysis on the printart database. In ECCV, 2012.
- [4] J. Collomosse, T. Bui, M. J. Wilber, C. Fang, and H. Jin. Sketching with style: Visual search with sketches and aesthetic context. In *ICCV*, 2017.
- [5] E. J. Crowley and A. Zisserman. The art of detection. In ECCV. Springer, 2016.
- [6] N. Garcia, B. Renoust, and Y. Nakashima. Context-aware embeddings for automatic art analysis. In *ICMR*, 2019.

- [7] N. Garcia and G. Vogiatzis. How to read paintings: Semantic art understanding with multi-modal retrieval. In *EECV Workshops*, 2018.
- [8] A. Grover and J. Leskovec. node2vec: Scalable feature learning for networks. In *KDD*, pages 855–864. ACM, 2016.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [10] S. Karayev, M. Trentacoste, H. Han, A. Agarwala, T. Darrell, A. Hertzmann, and H. Winnemoeller. Recognizing image style. In *BMVC*, 2014.
- [11] F. S. Khan, S. Beigpour, J. Van de Weijer, and M. Felsberg. Painting-91: a large scale database for computational painting categorization. *Machine vision and applications*, 2014.
- [12] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123(1):32–73, 2017.
- [13] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, editors, ECCV, 2014.
- [14] D. Ma, F. Gao, Y. Bai, Y. Lou, S. Wang, T. Huang, and L.-Y. Duan. From part to whole: Who is behind the painting? In ACMMM, 2017.
- [15] H. Mao, M. Cheung, and J. She. Deepart: Learning joint representations of visual arts. In ACMMM, 2017.
- [16] T. Mensink and J. Van Gemert. The rijksmuseum challenge: Museum-centered visual recognition. In *ICMR*, 2014.
- [17] B. Saleh and A. Elgammal. Large-scale classification of fineart paintings: Learning the right metric on the right feature. *International Journal for Digital Art History*, (2), 2016.
- [18] A. Sanakoyeu, D. Kotovenko, S. Lang, and B. Ommer. A style-aware content loss for real-time hd style transfer. In ECCV, volume 2, 2018.
- [19] W. R. Tan, C. S. Chan, H. E. Aguirre, and K. Tanaka. Ceci n'est pas une pipe: A deep convolutional network for fine-art paintings classification. *ICIP*, 2016.
- [20] L. Y. Wu, A. Fisch, S. Chopra, K. Adams, A. Bordes, and J. Weston. Starspace: Embed all the things! In AAAI, 2018.