# How to Read Paintings: Semantic Art Understanding with Multi-Modal Retrieval

Noa Garcia and George Vogiatzis

Aston University, United Kingdom
{garciadn,g.vogiatzis}@aston.ac.uk

**Abstract.** Automatic art analysis has been mostly focused on classifying artworks into different artistic styles. However, understanding an artistic representation involves more complex processes, such as identifying the elements in the scene or recognizing author influences. We present SemArt, a multi-modal dataset for semantic art understanding. SemArt is a collection of fine-art painting images in which each image is associated to a number of attributes and a textual artistic comment, such as those that appear in art catalogues or museum collections. To evaluate semantic art understanding, we envisage the Text2Art challenge, a multi-modal retrieval task where relevant paintings are retrieved according to an artistic text, and vice versa. We also propose several models for encoding visual and textual artistic representations into a common semantic space. Our best approach is able to find the correct image within the top 10 ranked images in the 45.5% of the test samples. Moreover, our models show remarkable levels of art understanding when compared against human evaluation.

**Keywords:** semantic art understanding · art analysis · image-text retrieval · multi-modal retrieval

## 1 Introduction

The ultimate aim of computer vision has always been to enable computers to understand images the way humans do. With the latest advances in deep learning technologies, the availability of large volumes of training data and the use of powerful graphic processing units, computer vision systems are now able to locate and classify objects in natural images with high accuracy, surpassing human performance in some specific tasks. However, we are still a long way from human-like analysis and extraction of high-level semantics from images. This work aims to push high-level image recognition by enabling machines to interpret art.

To study automatic interpretation of art, we introduce SemArt[1], a dataset for semantic art understanding. We build SemArt by gathering a collection of fine-art images, each with its respective attributes (author, type, school, etc.) as well as a short artistic comment or description, such as those that commonly appear

---

[1] http://noagarciad.com/SemArt/

**Fig. 1. SemArt dataset samples**. Each sample is a triplet of image, attributes and artistic comment.

in art catalogues or museum collections. Artistic comments involve not only descriptions of the visual elements that appear in the scene but also references to its technique, author or context. Some examples of the dataset are shown in Figure 1.

We address semantic art understanding by proposing a number of models that map paintings and artistic comments into a common semantic space, thus enabling comparison in terms of semantic similarity. To evaluate and benchmark the proposed models, we design the Text2Art challenge as a multi-modal retrieval task. The aim of the challenge is to evaluate whether the models capture enough of the insights and clues provided by the artistic description to be able to match it to the correct painting.

A key difference with previously proposed methods for semantic understanding of natural images (e.g. MS-COCO dataset [15]) is that our system relies on background information on art history and artistic styles. As already noted in previous work [3,5,4], paintings are substantially different from natural images in several aspects. Firstly, paintings, unlike natural images, are figurative representations of people, objects, places or situations which may or may not correspond to the real world. Secondly, the study of fine-art paintings usually requires previous knowledge about history of art, different artistic styles as well as contextual information about the subjects represented. Thirdly, paintings commonly exhibit one or more layers of abstraction and symbolism which creates ambiguity in interpretation.

In this work, we harness existing prior knowledge about art and deep neural networks to model understanding of fine-art paintings. Specifically, our contributions are:

1. to introduce the first dataset for semantic art understanding in which each sample is a triplet of images, attributes and artistic comments,
2. to propose models to map fine-art paintings and their high-level artistic descriptions onto a joint semantic space,
3. to design an evaluation protocol based on multi-modal retrieval for semantic art understanding, so that future research can be benchmarked under a common, public framework.

**Table 1. Datasets for art analysis**. *Meta* and *Text* columns state if image metadata and textual information are provided, respectively.

| Dataset | #Paintings | Meta | Text | Task |
|---|---|---|---|---|
| PRINTART [2] | 988 | ✓ | ✗ | Classification and Retrieval |
| Painting-91 [12] | 4,266 | ✓ | ✗ | Classification |
| Rijksmuseum [19] | 3,593 | ✓ | ✗ | Classification |
| Wikipaintings [11] | 85,000 | ✓ | ✗ | Classification |
| Paintings [3] | 8,629 | ✗ | ✗ | Object Recognition |
| Face Paintings [4] | 14,000 | ✗ | ✗ | Face Retrieval |
| VisualLink [22] | 38,500 | ✓ | ✗ | Instance Retrieval |
| Art500k [18] | 554,198 | ✓ | ✗ | Classification |
| SemArt | 21,383 | ✓ | ✓ | Semantic Retrieval |

## 2  Related Work

With the digitalization of large collections of fine-art paintings and the emergence of publicly available online art catalogs such as WikiArt[2] or the Web Gallery of Art[3], computer vision researchers become interested in analyzing fine-art paintings automatically. Early work [10,23,2,12] proposes methods based on handcrafted visual features to identify an author and/or a specific style in a piece of art. Datasets used in these kinds of approaches, such as PRINTART [2] and Painting-91 [12], are rather small, with 988 and 4,266 painting images, respectively. Mensink and Van Gemert introduce in [19] the large-scale Rijksmuseum dataset for multi-class prediction, consisting on 112,039 images from artistic objects, although only 3,593 are from fine-art paintings. With the success of convolutional neural networks (CNN) in large-scale image classification [14], deep features from CNNs replace handcrafted image representations in many computer vision applications, including painting image classification [1,11,21,25,16,18], and larger datasets are made publicly available [11,18]. In these methods, paintings are fed into a CNN to predict its artistic style or author by studying its visual aesthetics.

Besides painting classification, other work is focused on exploring image retrieval in artistic paintings. For example, in [2], monochromatic painting images are retrieved by using artistic-related keywords, whereas in [22] a pre-trained CNN is fine-tuned to find paintings with similar artistic motifs. Crowley and Zisserman [4] explore domain transfer to retrieve image of portraits from real faces, in the same way as [3] and [6] explore domain transfer to perform object recognition in paintings.

A summary of the existing datasets for fine-art understanding is shown in Table 1. In essence, previous work studies art from an aesthetics point of view to classify paintings according to author and style [2,12,19,11,18], to find relevant paintings according to a query input [2,4,22] or to identify objects in

---

[2] http://www.wikiart.org
[3] https://www.wga.hu/

artistic representations [3]. However, understanding art involves also identifying the symbolism of the elements, the artistic influences or the historical context of the work. To study such complex processes, we propose to interpret fine-art paintings in a semantic way by introducing SemArt, a multi-modal dataset for semantic art understanding. To the best of our knowledge, SemArt is the first corpus that provides not only fine-art images and their attributes, but also artistic comments for the semantic understanding of fine-art paintings.
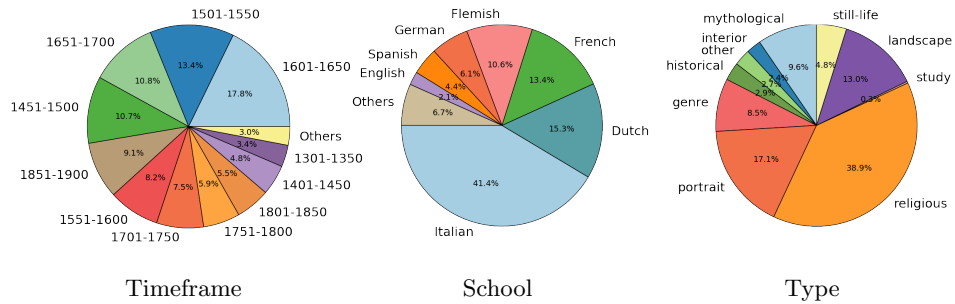
## 3   SemArt Dataset

### 3.1   Data Collection

To create the SemArt dataset, we collect artistic data from the Web Gallery of Art (WGA), a website with more than 44,809 images of European fine-art reproductions between the 8th and the 19th century. WGA provides links to all their images in a downloadable comma separated values file (CSV). In the CSV file, each image is associated with some attributes or metadata: author, author's birth and death, title, date, technique, current location, form, type, school and time-line. Following the links provided in the CSV file, we only collect images from artworks whose field *form* is set as painting, as opposite to images of other forms of art such as sculpture or architecture.

We create a script to collect artistic comments for each painting image, as they are not provided in the aforementioned CSV file. We omit images that are not associated to any comment and we remove irrelevant metadata fields, such as author's birth and death and current location. The final size of the cleaned collection is downsampled to 21,384 triplets, where each triplet is formed by an image, a text and a number of attributes.

### 3.2   Data Analysis

For each sample, the metadata is provided as a set of seven fields, which describe the basic attributes of its associated painting: *Author*, *Title*, *Date*, *Technique*, *Type*, *School* and *Timeframe*. In total, there are 3,281 different authors, the most frequent one being Vincent van Gogh with 327 paintings. There are 14,902 different titles in the dataset, with 38.8% of the paintings presenting a non-unique title. Among all the titles, Still-Life and Self-Portrait are the most common ones. *Technique* and *Date* fields are not available for all samples, but provided for completeness. *Type* field classifies paintings according to ten different genres, such as religious, landscape or portrait. There are 26 artistic schools in the collection, Italian being the most common, with 8,860 paintings and Finnish the least frequent with just 5 samples. Also, there are 22 different timeframes, which are periods of 50 years evenly distributed between 801 and 1900. The distribution of values over the fields *Type*, *School* and *Timeframe* is shown in Figure 2. With respect to artistic comments, the vocabulary set follows the Zipf's law [17]. Most of the comments are relatively short, with almost 70% of the them containing

**Fig. 2. Metadata distribution**. Distribution of samples within the SemArt dataset in Timeframe, School and Type attributes.

100 words or less. Images are provided in different aspect ratios and sizes. The dataset is randomly split into training, validation and test sets with 19,244, 1,069 and 1,069 triplets, respectively.

## 4   Text2Art Challenge

In what follows, we use bold style to refer to vectors and matrices (e.g $\boldsymbol{x}$ and $\boldsymbol{W}$). Given a collection of artistic samples $K$, the $k$-th sample in $K$ is given by the triplet $(img_k, com_k, att_k)$, being $img_k$ the artistic image, $com_k$ the artistic comment and $att_k$ the artistic attributes. Images, comments and attributes are input into specific encoding functions, $f_{img}$, $f_{com}$, $f_{att}$, to map raw data from the corpus into vector representations, $\boldsymbol{i}_k$, $\boldsymbol{c}_k$, $\boldsymbol{a}_k$, as:

$$\boldsymbol{i}_k = f_{img}(img_k; \phi_{img}) \tag{1}$$

$$\boldsymbol{c}_k = f_{com}(com_k; \phi_{com}) \tag{2}$$

$$\boldsymbol{a}_k = f_{att}(att_k; \phi_{att}) \tag{3}$$

where $\phi_{img}$, $\phi_{com}$ and $\phi_{att}$ are the parameters of each encoding function.

As comment encodings, $\boldsymbol{c}_k$, and attribute encodings, $\boldsymbol{a}_k$, are both from textual data, a joint textual vector, $\boldsymbol{t}_k$ can be obtained as:

$$\boldsymbol{t}_k = \boldsymbol{c}_k \oplus \boldsymbol{a}_k \tag{4}$$

where $\oplus$ is vector concatenation.

The transformation functions, $g_{vis}$ and $g_{text}$, can be defined as the functions that project the visual and the textual encodings into a common multi-modal space. The projected vectors $\boldsymbol{p}_k^{vis}$ and $\boldsymbol{p}_k^{text}$ are then obtained as:

$$\boldsymbol{p}_k^{vis} = g_{vis}(\boldsymbol{i}_k; \theta_{vis}) \tag{5}$$

$$\boldsymbol{p}_k^{text} = g_{text}(\boldsymbol{t}_k; \theta_{text}) \tag{6}$$

being $\theta_{vis}$ and $\theta_{text}$ the parameters of each transformation function.

For a given similarity function $d$, the similarity between any text (i.e. pair of comments and attributes) and any image in $K$ is measured as the distance between their projections:

$$d(\boldsymbol{p}_k^{text}, \boldsymbol{p}_j^{vis}) = d(g_{text}(\boldsymbol{t}_k; \theta_{text}), g_{vis}(\boldsymbol{i}_j; \theta_{vis})) \tag{7}$$

In semantic art understanding, the aim is to learn $f_{img}$, $f_{com}$, $f_{att}$, $g_{vis}$ and $g_{text}$ such that images, comments and attributes from the same sample are mapped closer in terms of $d$ than images, texts and attributes from different samples:

$$d(\boldsymbol{p}_k^{text}, \boldsymbol{p}_k^{vis}) < d(\boldsymbol{p}_k^{text}, \boldsymbol{p}_j^{vis}) \text{ for all } k, j \leq |K| \tag{8}$$

and

$$d(\boldsymbol{p}_k^{text}, \boldsymbol{p}_k^{vis}) < d(\boldsymbol{p}_j^{text}, \boldsymbol{p}_k^{vis}) \text{ for all } k, j \leq |K| \tag{9}$$

To evaluate semantic art understanding, we propose the Text2Art challenge as a multi-modal retrieval problem. Within Text2Art, we define two tasks: text-to-image retrieval and image-to-text retrieval. In text-to-image retrieval, the aim is to find the most relevant painting in the collection, $img^* \in K$, given a query comment and its attributes:

$$img^* = \underset{img_j \in K}{\arg\min} \, d(\boldsymbol{p}_k^{text}, \boldsymbol{p}_j^{vis}) \tag{10}$$

Similarly, in the image-to-text retrieval task, when a painting image is given, the aim is to find the comment and the attributes, $com^* \in K$ and $att^* \in K$ , that are more relevant to the visual query:

$$com^*, att^* = \underset{com_j, att_j \in K}{\arg\min} \, d(\boldsymbol{p}_j^{text}, \boldsymbol{p}_k^{vis}) \tag{11}$$

## 5   Models for Semantic Art Understanding

We propose several models to learn meaningful textual and visual encodings and transformations for semantic art understanding. First, images, comments and attributes are encoded into visual and textual vectors. Then, a multi-modal transformation model is used to map these visual and textual vectors into a common multi-modal space where a similarity function is applied.

### 5.1   Visual Encoding

We represent each painting image as a visual vector, $\boldsymbol{i}_k$, using convolutional neural networks (CNNs). We use different CNN architectures, such as VGG16 [24], different versions of ResNet [8] and RMAC [26].

**VGG16** [24] contains 13 3x3 convolutional layers and three fully-connected layers stacked on top of each other. We use the output of one of the fully connected layers as the visual encoding.

**ResNet** [8] uses shortcut connections to connect the input of a layer to the output of a deeper layer. There exist many versions depending on the number of layers, such as ResNet50 and ResNet152 with 50 and 152 layers, respectively. We use the output of the last layer as the visual encoding.

**RMAC** is a visual descriptor introduced by Tolias et al. in [26] for image retrieval. The activation map from the last convolutional layer from a CNN model is max-pooled over several regions to obtain a set of regional features. The regional features are post-processed, sum-up together and normalized to obtain the final visual representation.

## 5.2 Textual Encoding

With respect to the textual information, comments are encoded into a comment vector, $c_k$, and attributes are encoded into an attribute vector, $a_k$. To get the joint textual encoding, $t_k$, both vectors are concatenated.

**Comment Encoding** To encode comments into a comment vector, $c_k$, we first build a comment vocabulary, $V_C$. $V_C$ contains all the alphabetic words that appear at least ten times in the training set. The comment vector is obtained using three different techniques: a comment bag-of-words (BOWc), a comment multi-layer perceptron (MLPc) and a comment recurrent model (LSTMc).

**BOWc** each comment is encoded as a *term frequency - inverse document frequency* (tf-idf) vector by weighting each word in the comment by its relevance within the corpus.

**MLPc** comments are encoded as a tf-idf vectors and fed into a fully connected layer with tanh activation[4] and $\ell_2$-normalization. The output of the normalization layer is used as the comment encoding.
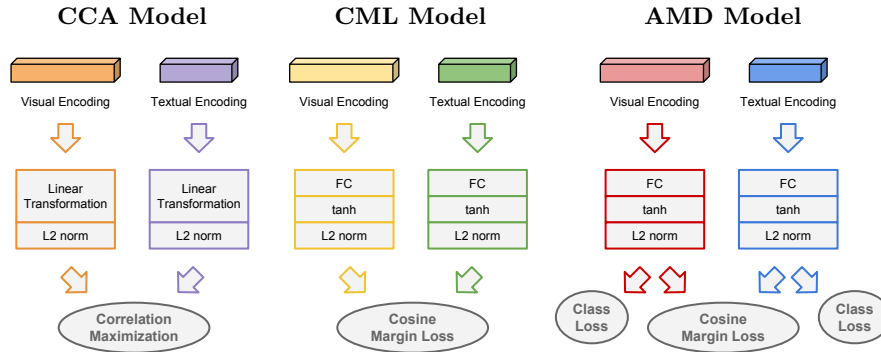
**LSTMc** each sentence in a comment is encoded into a sentence vector using a 2,400 dimensional pre-trained skip-thought model [13]. Sentence vectors are input into a long short-term memory network (LSTM) [9]. The last state of the LSTM is $\ell_2$-normalized and used as the comment encoding.

**Attribute Encoding** We use the attribute field *Title* in the metadata to provide an extra textual information to our model. We propose three different techniques to encode titles into attribute encodings, $a_k$: an attribute bag-of-words (BOWa) an attribute multi-layer perceptron (MLPa) and an attribute recurrent model (LSTMa).

**BOWa** as in comments, titles are encoded as a tf-idf-weighted vector using a title vocabulary, $V_T$. $V_T$ is built with all the alphabetic words in the titles of the training set.

---

[4] $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

**Fig. 3. Multi-modal transformation models**. Models for mapping textual and visual representations into a common multi-modal space.

**MLP$_a$** also as in comments, tf-idf encoded titles are fed into a fully connected layer with tanh activation and a $\ell_2$-normalization. The output of the normalization layer is used as the attribute vector.

**LSTM$_a$** in this case, each word in a title is fed into an embedding layer followed by a LSTM network. The output of the last state of the LSTM is $\ell_2$-normalized and used as the attribute encoding.

### 5.3   Multi-Modal Transformation

The visual and textual encodings, $\boldsymbol{i}_k$ and $\boldsymbol{t}_k$ respectively, encode visual and textual data into two different spaces. We use a multi-modal transformation model to map the visual and textual representations into a common multi-modal space. In this common space, textual and visual information can be compared in terms of the similarity function $d$. We propose three different models, which are illustrated in Figure 3.

**CCA** Canonical Correlation Analysis (CCA) [7] is a linear approach for projecting data from two different sources into a common space by maximizing the normalized correlation between the projected data. The CCA projection matrices are learnt by using training pairs of samples from the corpus. At test time, the textual and visual encodings from a test sample are projected using these CCA matrices.

**CML** Cosine Margin Loss (CML) is a deep learning architecture trained end-to-end to learn the visual and textual encodings and their projections all at once. Each image encoding is fed into a fully connected layer followed by a tanh activation function and a $\ell_2$-normalization layer to project the visual feature, $\boldsymbol{i}_j$, into a $D$-dimensional space, obtaining the projected visual vector $\boldsymbol{p}_j^{vis}$. Similarly, each textual vector $\boldsymbol{t}_k$, is input into another network with identical layer structure (fully connected layer with tanh activation and $\ell_2$-normalization) to map the textual feature into the same $D$-dimensional

space, obtaining the projected textual vector $\boldsymbol{p}_k^{text}$. We train the CML model with both positive ($k = j$) and negative ($k \neq j$) pairs of textual and visual data and cosine similarity with margin as the loss function:

$$L_{CML}(\boldsymbol{p}_k^{vis}, \boldsymbol{p}_j^{text}) = \begin{cases} 1 - \cos(\boldsymbol{p}_k^{vis}, \boldsymbol{p}_j^{text}), & \text{if } k = j \\ \max(0, \cos(\boldsymbol{p}_k^{vis}, \boldsymbol{p}_j^{text}) - m), & \text{if } k \neq j \end{cases} \qquad (12)$$

where cos is the cosine similarity between two normalized vectors and $m$ is the margin hyperparameter.

**AMD** Augmented Metadata (AMD) is a model in which the network is informed with attribute data for an extra alignment between the visual and the textual encodings. The AMD model consists on a deep learning architecture that projects both visual and textual vectors into the common multi-modal space whereas, at the same time, ensures that the projected encodings are meaningful in the art domain. As in the CML model, image and textual encodings are projected into $D$-dimensional vectors using fully connected layers, and the loss between the multi-modal transformations is computed using a cosine margin loss. Attribute metadata is used to train a pair of classifiers on top of the projected data (Figure 3, AMD Model), each classifier consisting of a fully connected layer without activation. Metadata classifiers are trained using a standard cross entropy classification loss function:

$$L_{META}(\boldsymbol{x}, class) = -\log\left(\frac{\exp(\boldsymbol{x}[class])}{\sum_j \exp(\boldsymbol{x}[j])}\right) \qquad (13)$$

which contribute to the total loss of the model in addition to the cosine margin loss. The total loss of the model is then computed as:

$$
\begin{aligned}
L_{AMD}(\boldsymbol{p}_k^{text}, \boldsymbol{p}_j^{vis}, l_{p_k^{text}}, l_{p_j^{vis}}) = {} & (1 - 2\alpha)L_{CML}(\boldsymbol{p}_k^{text}, \boldsymbol{p}_j^{vis}) \\
& + \alpha L_{META}(\boldsymbol{p}_k^{text}, l_{p_k^{text}}) \qquad (14) \\
& + \alpha L_{META}(\boldsymbol{p}_j^{vis}, l_{p_j^{vis}})
\end{aligned}
$$

where $l_{p_k^{text}}$ and $l_{p_j^{vis}}$ are the class labels of the $k$-th text and the $j$-th image, respectively, and $\alpha$ is the weight of the classifier loss.

## 6    Experiments

**Experimental Details.** In the image encoding part, each network is initialized with its standard pre-trained weights for image classification. Images are scaled down to 256 pixels per side and randomly cropped into $224 \times 224$ patches. Visual data is augmented by randomly flipping images horizontally. In the textual encoding part, the dimensionality of LSTM hidden state for comments is 1,024, whereas in the LSTM for titles is 300. The title vocabulary size is 9,092. Skip thoughts dimensionality is set to 2,400. In the multi-modal transformation part,

**Table 2. Visual Domain Adaptation.** Transferability of visual features from the natural image classification domain to the Text2Art challenge.

| Encoding | | Text-to-Image | | | | Image-to-Text | | | |
|---|---|---|---|---|---|---|---|---|---|
| Img | Dim | R@1 | R@5 | R@10 | MR | R@1 | R@5 | R@10 | MR |
| VGG16 FC1 | 4,096 | 0.069 | 0.129 | 0.174 | 115 | 0.061 | 0.129 | 0.180 | 121 |
| VGG16 FC2 | 4,096 | 0.051 | 0.097 | 0.109 | 278 | 0.051 | 0.085 | 0.103 | 275 |
| VGG16 FC3 | 1,000 | 0.101 | 0.211 | 0.285 | 44 | 0.094 | 0.217 | 0.283 | 51 |
| ResNet50 | 1,000 | 0.114 | 0.231 | 0.304 | 42 | 0.114 | 0.242 | 0.318 | 44 |
| ResNet152 | 1,000 | 0.108 | **0.254** | **0.343** | **36** | **0.118** | **0.250** | **0.321** | **36** |
| RMAC VGG16 | 512 | 0.092 | 0.206 | 0.286 | 41 | 0.084 | 0.202 | 0.293 | 44 |
| RMAC Res50 | 2,048 | 0.084 | 0.202 | 0.293 | 48 | 0.097 | 0.215 | 0.288 | 49 |
| RMAC Res152 | 2,048 | **0.115** | 0.233 | 0.306 | 44 | 0.103 | 0.238 | 0.305 | 44 |

the CCA matrices are learnt using scikit-learn [20]. For the deep learning architectures, we use Adam optimizer and the learning rate is set to 0.0001, $m$ to 0.1 and $\alpha$ to 0.01. Training is conducted in mini batches of 32 samples. Cosine similarity is used as the similarity function $d$ in all of our models.

**Text2Art Challenge Evaluation.** Painting images are ranked according to their similarity to a given text, and vice versa. The ranking is computed on the whole set of test samples and results are reported as median rank (MR) and recall rate at K (R@K), with K being 1, 5 and 10. MR is the value separating the higher half of the relevant ranking position amount all samples, so the lower the better. Recall at rate K is the rate of samples for which its relevant image is in the top K positions of the ranking, so the higher the better.

## 6.1   Visual Domain Adaptation

We first evaluate the transferability of visual features from the natural image domain to the artistic domain. In this experiment, texts are encoded with the BOW$_c$ approach with $V_C = 3{,}000$. As multi-modal transformation model, a 128-dimensional CCA is used. We extract visual encodings from networks pre-trained for classification of natural images without further fine-tunning or refinement. For the VGG16 model, we extract features from the first, second and third fully connected layer (VGG16$_{FC1}$ , VGG16$_{FC2}$ and VGG16$_{FC3}$). For the ResNet models, we consider the visual features from the output of the networks (ResNet50 and ResNet152). Finally, RMAC representation is computed using a VGG16, a ResNet50 and a ResNet152 (RMAC$_{VGG16}$ , RMAC$_{Res50}$ and RMAC$_{Res152}$). Results are detailed in Table 2. As semantic art understanding is a high-level task, it is expected that representations acquired from deeper layers perform better, as in the VGG16 models, where the deepest layer of the network obtains the best performance. RMAC features respond well when transferring from natural images to art, although ResNet models obtain the best performance. Considering these results, we use ResNets as visual encoders in the following experiments.

**Table 3. Text Encoding in Art.** Comparison between different text encodings in the Text2Art challenge.

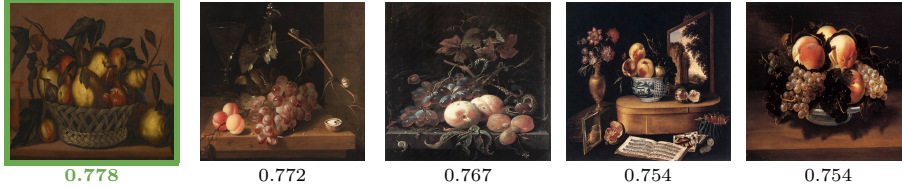| Encoding | | Text-to-Image | | | | Image-to-Text | | | |
|---|---|---|---|---|---|---|---|---|---|
| **Com** | **Att** | **R@1** | **R@5** | **R@10** | **MR** | **R@1** | **R@5** | **R@10** | **MR** |
| LSTMc | LSTMa | 0.053 | 0.162 | 0.256 | 33 | 0.053 | 0.180 | 0.268 | 33 |
| MLPc | LSTMa | 0.089 | 0.260 | 0.376 | 21 | 0.093 | 0.249 | 0.363 | 21 |
| MLPc | MLPa | 0.137 | 0.306 | 0.432 | 16 | **0.140** | 0.317 | 0.436 | 15 |
| BOWc | BOWa | **0.144** | **0.332** | **0.454** | **14** | 0.138 | **0.327** | **0.457** | **14** |

## 6.2   Text Encoding in Art

We then compare the performance between the different text encoding models in the Text2Art challenge. In this experiment, images are encoded with a ResNet50 network and the CML model is used to learn the mapping of the visual and the textual encodings into a common 128-dimensional space. The different encoding methods are compared in Table 3. The best performance is obtained when using the simplest bag-of-words approach both for comments and titles (BOWc and BOWa), although the multi-layer perceptron model (MLPc and MLPa) obtain similar results. Models based on recurrent networks (LSTMc and LSTMa) are not able to capture the insights of semantic art understanding. These results are consistent with previous work [27], which shows that text recurrent models perform worse than non-recurrent methods for multi-modal tasks that do not require text generation.
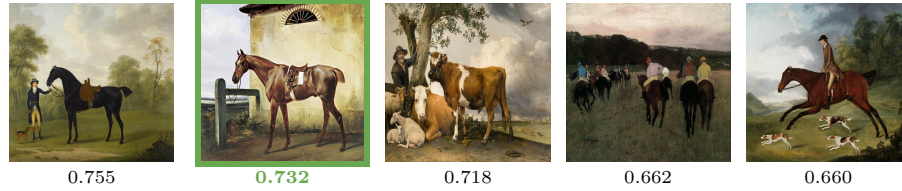
## 6.3   Multi-Modal Models for Art Understanding

Finally, we compare the three proposed multi-modal transformation models in the Text2Art challenge: CCA, CML and AMD. For the AMD approach, we use four different attributes to inform the model: Type (AMDt), TimeFrame (AMDtf), School (AMDs) and Author (AMDa). ResNet50 is used to encode visual features. Results are shown in Table 4. Random ranking results are provided as reference. Overall, the best performance is achieved with the CML model and bag-of-words encodings. CCA achieves the worst results among all the models, which suggests that linear transformations are not able to adjust properly to the task. Surprisingly, adding extra information in the AMD models does not lead to further improvement over the CML approach. We suspect that this might be due to the unbalanced number of samples within the classes of the dataset. Qualitative results of the CML model with ResNet50 and bag-of-words encodings are shown in Figures 4 and 5. In the positive examples (Figure 4), not only the ground truth painting is ranked within the top five returned images, but also all the images within the top five are semantically similar to the query text. In the unsuccessful examples (Figure 5), although the ground truth image is not ranked in the top positions of the list, the algorithm returns images that are semantically meaningful to fragments of the text, which indicates how challenging the task is.

**Title**: Still-Life of Apples, Pears and Figs in a Wicker Basket on a Stone Ledge
**Comment**: The large dark vine leaves and fruit are back-lit and are sharply silhouetted against the luminous background, to quite dramatic effect. Ponce's use of this effect strongly indicates the indirect influence of Caravaggio's Basket of Fruit in the Pinacoteca Ambrosiana, Milan, almost 50 years after it was created.
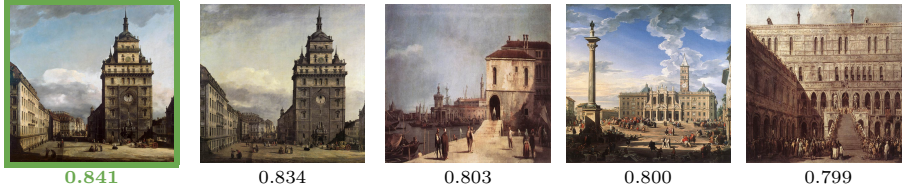


| **0.778** | 0.772 | 0.767 | 0.754 | 0.754 |

**Title**: A Saddled Race Horse Tied to a Fence
**Comment**: Horace Vernet enjoyed royal patronage, one of his earliest commissions was a group of ten paintings depicting Napoleon's horses. These works reveal his indebtedness to the English tradition of horse painting. The present painting was commissioned in Paris in 1828 by Jean Georges Schickler, a member of a German based banking family, who had a passion for horse racing.



| 0.755 | **0.732** | 0.718 | 0.662 | 0.660 |

**Title**: Portrait of a Girl
**Comment**: This painting shows a girl in a yellow dress holding a bouquet of flowers. It is a typical portrait of the artist showing the influence of his teacher, Agnolo Bronzino.



| **0.870** | 0.848 | 0.847 | 0.827 | 0.825 |

**Title**: The Kreuzkirche in Dresden
**Comment**: A few years later, during his second stay in Saxony, Bellotto depicted the demolition of this Gothic church. There exists an almost identical version in the Gemldegalerie, Dresden.



| **0.841** | 0.834 | 0.803 | 0.800 | 0.799 |

**Fig. 4. Qualitative positive results**. For each text (i.e. title and comment), the top five ranked images, along with their score, are shown. The ground truth image is highlighted in green.

**Title**: Brunette with Bare Breasts
**Comment**: The 1870s were rich in female models for Manet: the Brunette with Bare Breasts, the Blonde with Bare Breasts and the Sultana testify to it.

**ranked 28, 0.445**



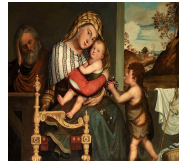| 0.640 | 0.622 | 0.605 | 0.572 | 0.569 |



**Title**: Virgin and Child with the Young St John the Baptist
**Comment**: The stylistic characteristics of this painting, such as rounded faces and narrow, elongated eyes seem to be a general reflection of the foreign presence in Genoese painting at this time.

**ranked 17, 0.690**



| 0.754 | 0.751 | 0.730 | 0.727 | 0.721 |

**Fig. 5. Qualitative negative result**. For each text, the ground truth image is shown next to it, along with its ranking position and score. Below, the five top ranked images.

## 6.4   Human Evaluation

We design a task in Amazon Mechanical Turk[5] for testing human performance in the Text2Art challenge. For a given artistic text, which includes comment, title, author, type, school and timeframe, human evaluators are asked to choose the most appropriate painting from a pool of ten images. The task has two different levels: *easy*, in which the pool of images is chosen randomly from all the paintings in test set, and *difficult*, in which the ten images in the pool share the same attribute type (i.e. portraits, landscapes, etc.). For each level, evaluators are asked to perform the task in 100 artistic texts. Accuracy is measured as the ratio of correct answers over the total number of answers. Results are shown in Table 5. Although human accuracy is considerable high, reaching 88.9% in the easiest set, there is a drop in performance in the difficult level, mostly because images from the same type contain more similar comments than images from

[5] https://www.mturk.com/

**Table 4. Multi-modal transformation models.** Comparison between different multi-modal transformation models in the Text2Art challenge.

| Technique | | | Text-to-Image | | | | Image-to-Text | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Model** | **Com** | **Att** | **R@1** | **R@5** | **R@10** | **MR** | **R@1** | **R@5** | **R@10** | **MR** |
| Random | - | - | 0.0008 | 0.004 | 0.009 | 539 | 0.0008 | 0.004 | 0.009 | 539 |
| CCA | MLPc | MLPa | 0.117 | 0.283 | 0.377 | 25 | 0.131 | 0.279 | 0.355 | 26 |
| CML | BOWc | BOWa | **0.144** | **0.332** | **0.454** | **14** | 0.138 | **0.327** | **0.457** | **14** |
| CML | MLPc | MLPa | 0.137 | 0.306 | 0.432 | 16 | **0.140** | 0.317 | 0.436 | 15 |
| AMDт | MLPc | MLPa | 0.114 | 0.304 | 0.398 | 17 | 0.125 | 0.280 | 0.398 | 16 |
| AMDтF | MLPc | MLPa | 0.117 | 0.297 | 0.389 | 20 | 0.123 | 0.298 | 0.413 | 17 |
| AMDs | MLPc | MLPa | 0.103 | 0.283 | 0.401 | 19 | 0.118 | 0.298 | 0.423 | 16 |
| AMDA | MLPc | MLPa | 0.131 | 0.303 | 0.418 | 17 | 0.120 | 0.302 | 0.428 | 16 |

**Table 5. Human Evaluation.** Evaluation in both the easy and the difficult sets.

| | Technique | | | | Text-to-Image | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | **Model** | **Img** | **Com** | **Att** | **Land** | **Relig** | **Myth** | **Genre** | **Port** | **Total** |
| Easy | CCA | ResNet152 | MLPc | MLPa | 0.708 | 0.609 | 0.571 | 0.714 | 0.615 | 0.650 |
| | CML | ResNet50 | BOWc | BOWa | 0.917 | 0.683 | 0.714 | 1 | 0.538 | 0.750 |
| | Human | - | - | - | 0.918 | 0.795 | 0.864 | 1 | 1 | 0.889 |
| Diff. | CCA | ResNet152 | MLPc | MLPa | 0.600 | 0.525 | 0.400 | 0.300 | 0.400 | 0.470 |
| | CML | ResNet50 | BOWc | BOWa | 0.500 | 0.875 | 0.600 | 0.200 | 0.500 | 0.620 |
| | Human | - | - | - | 0.579 | 0.744 | 0.714 | 0.720 | 0.674 | 0.714 |

different types. We evaluate a CCA and a CML model in the same data split as humans. The CML model with bag-of-words and ResNet50 is able to find the relevant image in the 75% of the samples in the easy set and in the 62% of the cases in the difficult task. There is around ten points of difference between CML model and the human evaluation, which suggests that, although there is still room for improvement, meaningful art representations are being obtained.

## 7   Conclusions

We presented the SemArt dataset, the first collection of fine-art images with attributes and artistic comments for semantic art understanding. In SemArt, comments describe artistic information of the painting, such as content, techniques or context. We designed the Text2Art challenge to evaluate semantic art understanding as a multi-modal retrieval task, whereby given an artistic text (or image), a relevant image (or text) is found. We proposed several models to address the challenge. We showed that for visual encoding, ResNets perform the best. For textual encoding, recurrent models performed worse than multi-layer preceptron or bag-of-words. We projected the visual and textual encodings into a common multi-modal space using several methods, the one with the best results being a neural network trained with cosine margin loss. Experiments with human evaluators showed that current approaches are not able to reach human levels of art understanding yet, although meaningful representations for semantic art understanding are being learnt.

# References

1. Bar, Y., Levy, N., Wolf, L.: Classification of artistic styles using binarized features derived from a deep neural network. In: ECCV Workshops (2014)
2. Carneiro, G., da Silva, N.P., Del Bue, A., Costeira, J.P.: Artistic image classification: An analysis on the printart database. In: ECCV (2012)
3. Crowley, E., Zisserman, A.: The state of the art: Object retrieval in paintings using discriminative regions. In: BMVC (2014)
4. Crowley, E.J., Parkhi, O.M., Zisserman, A.: Face painting: querying art with photos. In: BMVC. pp. 65–1 (2015)
5. Crowley, E.J., Zisserman, A.: In search of art. In: ECCV Workshop (2014)
6. Crowley, E.J., Zisserman, A.: The art of detection. In: ECCV. Springer (2016)
7. Gong, Y., Ke, Q., Isard, M., Lazebnik, S.: A multi-view embedding space for modeling internet images, tags, and their semantics. International journal of computer vision **106**(2) (2014)
8. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition (2016)
9. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8) (1997)
10. Johnson, C.R., Hendriks, E., Berezhnoy, I.J., Brevdo, E., Hughes, S.M., Daubechies, I., Li, J., Postma, E., Wang, J.Z.: Image processing for artist identification. IEEE Signal Processing Magazine **25**(4) (2008)
11. Karayev, S., Trentacoste, M., Han, H., Agarwala, A., Darrell, T., Hertzmann, A., Winnemoeller, H.: Recognizing image style. In: BMVC (2014)
12. Khan, F.S., Beigpour, S., Van de Weijer, J., Felsberg, M.: Painting-91: a large scale database for computational painting categorization. Machine vision and applications (2014)
13. Kiros, R., Zhu, Y., Salakhutdinov, R.R., Zemel, R., Urtasun, R., Torralba, A., Fidler, S.: Skip-thought vectors. In: Advances in neural information processing systems (2015)
14. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NIPS (2012)
15. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) Computer Vision – ECCV 2014. pp. 740–755 (2014)
16. Ma, D., Gao, F., Bai, Y., Lou, Y., Wang, S., Huang, T., Duan, L.Y.: From part to whole: Who is behind the painting? In: Proceedings of the 2017 ACM on Multimedia Conference. ACM (2017)
17. Manning, C.D., Schütze, H.: Foundations of statistical natural language processing. Information Retrieval **4**, 80–81 (2001)
18. Mao, H., Cheung, M., She, J.: Deepart: Learning joint representations of visual arts. In: ACM on Multimedia Conference (2017)
19. Mensink, T., Van Gemert, J.: The rijksmuseum challenge: Museum-centered visual recognition. In: ICMR (2014)
20. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12** (2011)

21. Saleh, B., Elgammal, A.M.: Large-scale classification of fine-art paintings: Learning the right metric on the right feature. CoRR (2015)
22. Seguin, B., Striolo, C., Kaplan, F., et al.: Visual link retrieval in a database of paintings. In: ECCV Workshops (2016)
23. Shamir, L., Macura, T., Orlov, N., Eckley, D.M., Goldberg, I.G.: Impressionism, expressionism, surrealism: Automated recognition of painters and schools of art. ACM Transactions on Applied Perception (2010)
24. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: International Conference on Learning Representations (2015)
25. Tan, W.R., Chan, C.S., Aguirre, H.E., Tanaka, K.: Ceci n'est pas une pipe: A deep convolutional network for fine-art paintings classification. ICIP (2016)
26. Tolias, G., Sicre, R., Jégou, H.: Particular object retrieval with integral max-pooling of cnn activations. International Conference on Learning Representations (2015)
27. Wang, L., Li, Y., Huang, J., Lazebnik, S.: Learning two-branch neural networks for image-text matching tasks. IEEE Transactions on Pattern Analysis and Machine Intelligence (2018)