



**HAL**  
open science

## Deep semantic-visual embedding with localization

Martin Engilberge, Louis Chevallier, Patrick Pérez, Matthieu Cord

► **To cite this version:**

Martin Engilberge, Louis Chevallier, Patrick Pérez, Matthieu Cord. Deep semantic-visual embedding with localization. RFIAP 2018 - Congrès Reconnaissance des Formes, Image, Apprentissage et Perception, Jun 2018, Marne-la-Vallée, France. hal-02171880

**HAL Id: hal-02171880**

**<https://hal.science/hal-02171880>**

Submitted on 3 Jul 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Deep semantic-visual embedding with localization

Martin Engilberge<sup>1,2</sup>

Louis Chevallier<sup>2</sup>

Patrick Pérez<sup>3</sup>

Matthieu Cord<sup>1</sup>

<sup>1</sup>Sorbonne Université, Paris, France, <sup>2</sup>Technicolor, Cesson Sévigné, France, <sup>3</sup>Valeo.ai, Paris, France

{martin.engilberge, matthieu.cord}@lip6.fr {louis.chevallier, patrick.perez}@technicolor.com

## Résumé

Nous proposons dans ce papier un réseau de neurones profond pour apprendre un alignement entre des images et leurs descriptions textuelles. Notre architecture est basée sur un réseau à deux branches, l'une visuelle, bénéficiant des mécanismes d'agrégation (pooling) récents, et l'autre encodant l'information textuelle. L'ensemble du réseau est appris de bout en bout dans un schéma supervisé par des paires (image, légende textuelle), fournissant alors une représentation sémantique exploitable dans différents contextes. Notre système obtient des résultats état-de-l'art sur une tâche importante de recherche d'information croisée image-texte. Nous montrons également sa capacité à découvrir la position des concepts de l'espace sémantique dans les images, permettant ainsi d'ancrer des phrases sur des parties d'images.

## Mots Clef

Alignement multimodal, Recherche d'information multimodale, Localisation d'information visuelle.

## Abstract

In this paper, we introduce a deep network to learn a cross-modal mapping between images and texts. It is based on two-path neural network combining a visual path that leverages recent space-aware pooling mechanisms with a textual path. Jointly trained from scratch, our semantic-visual embedding offers a versatile model. Once trained under the supervision of captioned images, it yields new state-of-the-art performance on cross-modal retrieval. It also allows the localization of new concepts from the embedding space into any input image, delivering state-of-the-art result on the visual grounding of phrases.

## Keywords

Multimodal embedding, Cross-modal retrieval, Visual grounding.

## 1 Introduction

Vision and Language understanding has motivated a lot of recent works from Machine Learning and Computer Vision

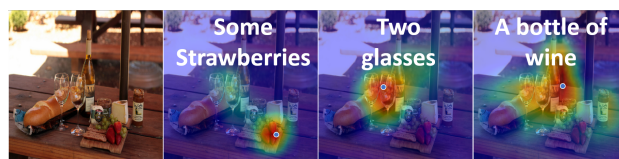


FIGURE 1 – **Concept localization with proposed semantic-visual embedding.** Not only does our deep embedding allows cross-modal retrieval with state-of-the-art performance, but it can also associate to an image, *e.g.*, the picnic table on the left, a localization heatmap for any text query, as shown with overlays for three text examples. The circled blue dot indicates the highest peak in the heatmap.

communities. Several works have demonstrated how deep representations of images and texts can be jointly leveraged to build *visual-semantic embeddings* [10, 16]. The ability to map natural images and texts in a shared representation space not only does permit to revisit visual recognition and captioning tasks, but also opens up new usages, such as cross-modal content search or generation.

One popular approach to semantic-visual joint embedding is to connect two mono-modal paths with one or multiple fully connected layers [19, 37, 9, 2] : A visual path based on a pre-trained convolutional neural network (CNN) and a text path based on a pre-trained recurrent neural network (RNN) operating on a given word embedding. Using aligned text-image data, such as images with multiple captions from MS-COCO dataset [25], final mapping layers can be trained, along with the optional fine-tuning of the two branches.

In this paper, we investigate new pooling mechanisms in the visual path. Inspired by recent work on weakly supervised object localization [43, 7], we propose in particular to leverage selective spatial pooling with negative evidence proposed in [7] to improve visual feature extraction without resorting, *e.g.*, to expensive region proposal strategies. Another important benefit of the proposed joint architecture is that, once trained, it allows localization of arbitrary concepts within arbitrary images : Given an image

and the embedding of a text (or any point of the embedding space), we propose a mechanism to compute a localization map, as demonstrated in Fig. 1.

We discuss in Section 2 the related works, on semantic-visual embedding and on weak supervised localization, and position our work. Section 3 is dedicated to the presentation of our own system, which couples selective spatial pooling with recent architectures and which relies on a triplet ranking loss based on hard negatives. We also show how it can be equipped with a concept localization module by exploiting without pooling the last feature maps in the visual path. The performance is assessed on two very different tasks in Section 4. We first establish new state-of-the-art performance on cross-modal matching, effectively composed of two symmetric sub-tasks : Retrieving captions from query images and vice-versa. Without additional fine-tuning, our model with its built-in concept localization mechanism also outperforms existing work on the "pointing game" sentence-grounding task.

## 2 Related Work and Paper Positioning

Deep learning nets are now routinely used to extract versatile deep features from images [22, 34, 12] as well as from words and sentences [28, 31, 4, 24]. In the following, we review learning methods to handle such mono/cross-modal representations, and we also highlight approaches dealing with spatial localization in this context.

**Metric learning for semantic embedding** Several methods have been proposed to learn visual metrics. [41] minimizes the distance within pairs of similar training examples with a constraint on the distance between dissimilar ones. This learning process has been extended to kernel functions as in [27]. Other methods consider triplets or quadruplets of images, to express richer relative constraints among groups of similar and dissimilar examples [38, 11, 3]. This kind of learning strategies has been also considered for deep (Siamese) architecture embeddings in the pairwise framework [35], and recently extended to triplets [14]. To embed words in a continuous space as vector representations, Mikolov *et al.*'s "word2vec" is definitively the leading technique [28]. In recent years, several approaches have been developed for learning operators that map sequences of word vectors to sentence vectors including recurrent networks [13, 4, 24] and convolutional networks [17]. Using word vector learning as inspiration, [20] proposes an objective function that abstracts the skip-gram word model to the sentence level, by encoding a sentence to predict the sentences around it. In our work, we adopt most recent and effective deep architectures on both sides, using a deep convolutional network (ResNet) for images [12] and a simple recurrent unit (SRU) network [24] to encode the textual information. Our learning scheme is based on fine-tuning (on the visual side) and triplet-based optimization, in the context of cross-modal alignment that we

describe now.

**Cross-modal embedding and localization** The Canonical Correlation Analysis (CCA) method uses linear projections of two views of heterogeneous data in a common space [15], which are optimized in order to maximize the cross correlation. Non-linear extensions using kernel [23] or deep net [1] have been proposed. Recently, for the more advanced task of textual image description (captioning), [19, 16] propose a joint embedding encoder with a decoder architecture. Other works focus on the sole building of such a joint embedding, to perform image-text matching and cross-modal retrieval [10, 9, 26]. Our work stems from this latter class. We aim at generating a joint embedding space that offers rich descriptors for both images and texts. We adopt the contrastive triplet loss that follows the margin-based principle to separate the positive pairs from the negative ones with at least a fixed margin. The training strategy with stochastic gradient descent has to be carefully adapted to the cross-modality of the triplets. Following [9], we resort to batch-based hard mining, but we depart from this work, and from other related approaches, in the way we handle localization information.

Existing works that combine localization and multimodal embedding rely on a two-step process. First, regions are extracted either by a dedicated model, *e.g.*, EdgeBox in [37], or by a module in the architecture. Then the embedding space is used to measure the similarity between these regions and textual data. [30, 16] use this approach on the dense captioning task to produce region annotations. It is also used for phrase localization by [37] where the region with the highest similarity with the phrase is picked. To address this specific problem of phrase grounding, Xiao *et al.* [40] recently proposed to learn jointly a similarity score and an attention mask. The model is trained using a structural loss, leveraging the syntactic structure of the textual data to enforce corresponding structure in the attention mask. In contrast to these works, our approach to spatial localization in semantic-visual embedding is weakly supervised and does not rely on a region extraction model. Instead, we take inspiration from other works on weakly supervised visual localization to design our architecture, with no need for a location-dependent loss. A number of weakly supervised object localization approaches extrapolate localization features while training an image classifier, *e.g.*, [43, 7, 5]. The main strategy consists in using a fully convolutional deep architecture that postpones the spatial aggregation (pooling) at the very last layer of the net. We follow the same strategy, but in the context of multimodal embedding learning, hence with a different goal. In particular, richer semantics is sought (and used for training) in the form of visual description, whether at the scene or at the object level.

## 3 Approach

The overall structure of the proposed approach, shown in Fig. 2, follows the dual-path encoding architecture of Kiros

*et al.* [19]. We first explain its specifics before turning to its training with a cross-modal triplet ranking loss.

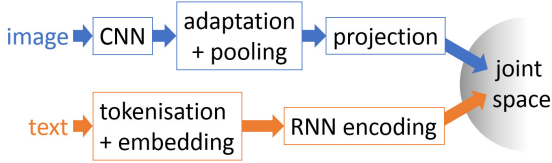


FIGURE 2 – **Two-path multimodal embedding architecture.** Images of arbitrary size and text of arbitrary length pass through dedicated neural networks to be mapped into a shared representation vector space. The visual path (blue) is composed of a fully convolutional neural network (ResNet in experiments), followed by a convolutional adaptation layer, a pooling layer that aggregates previous feature maps into a vector and a final projection to the final output space; The textual path (orange) is composed of a recurrent net running on sequences of text tokens individually embedded with an off-the-shelf map (word2vec in experiments).

### 3.1 Semantic-visual embedding architecture

**Visual path** In order to accommodate variable size images and to benefit from the performance of very deep architectures, we rely on fully convolutional residual ResNet-152 [12] as our base visual network. Its penultimate layer outputs a stack of  $D = 2048$  feature maps of size  $(w, h) = (\frac{W}{32}, \frac{H}{32})$ , where  $(W, H)$  is the spatial size of the input image. These feature maps retain coarse spatial information that lends itself to spatial reasoning in subsequent layers. Following the weakly supervised learning framework proposed by Durand *et al.* [7, 6], we first transform this stack through a linear adaptation layer of  $1 \times 1$  convolutions. While in WELDON [7] and in WILDCAT [6] the resulting maps are class-related (one map per class in the former, a fixed number of maps per class in the latter), we do not address classification or class detection here. Hence we empirically set the number  $D'$  of these new maps to a large value, 2400 in our experiments. A pooling à la WELDON is then used, but again in the absence of classes, to turn these maps into vector representations of dimension  $D'$ . A linear projection with bias, followed by  $\ell_2$  normalization accomplishes the last step to the embedding space of dimension  $d$ . More formally, the visual embedding path is defined as follows :

$$\mathbf{I} \xrightarrow{f_{\theta_0}} \mathbf{F} \xrightarrow{g_{\theta_1}} \mathbf{G} \xrightarrow{\text{sPool}} \mathbf{h} \in \mathbb{R}^{D'} \xrightarrow{p_{\theta_2}} \mathbf{x} \in \mathbb{R}^d, \quad (1)$$

where  $\mathbf{I} \in (0, 255)^{W \times H \times 3}$  is the input color image,  $f_{\theta_0}(\mathbf{I}) \in \mathbb{R}^{w \times h \times D}$  is the output of ResNet’s conv5 parametrized by weights in  $\theta_0$ ,  $g_{\theta_1}$  is a convolution layer with  $|\theta_1| = D \times D'$  weights and with activation in  $\mathbb{R}^{w \times h \times D'}$ ,

sPool is the selective spatial pooling with negative evidence defined in [7] :

$$\mathbf{h}[k] = \max \mathbf{G}[:, :, k] + \min \mathbf{G}[:, :, k], \quad k = 1 \cdots D', \quad (2)$$

and  $p_{\theta_2}$  is an  $\ell_2$ -normalized affine function

$$p_{\theta_2}(\mathbf{h}) = \frac{A\mathbf{h} + \mathbf{b}}{\|A\mathbf{h} + \mathbf{b}\|_2}, \quad (3)$$

where  $\theta_2 = (A, \mathbf{b})$  is of size  $d \times (D' + 1)$ . We shall denote  $\mathbf{x} = F(\mathbf{I}; \theta_{0:2})$  for short this visual embedding.

**Textual path** The inputs to this path are tokenized sentences (captions), *i.e.*, variable length sequences of tokens  $S = (s_1 \cdots s_T)$ . Each token  $s_t$  is turned into a vector representation  $s_t \in \mathbb{R}^K$  by the pre-trained word2vec embedding [28] of size  $K = 620$  used in [20]. Several RNNs have been proposed in the literature to turn such variable length sequences of (vectorized) words into meaningful, fixed-sized representations. In the specific context of semantic-visual embedding, [19, 9] use for instance gated recurrent unit (GRU) [4] networks as text encoders. Based on experimental comparisons, we chose to encode sentences with the simple recurrent unit (SRU) architecture recently proposed in [24]. Since we train this network from scratch, we take its output, up to  $\ell_2$  normalization, as the final embedding of the input sentence. There is no need here for an additional trainable projection layer. Formally, the textual path reads :

$$S \xrightarrow{\text{w2v}} \mathbf{S} \xrightarrow{\text{normSRU}_{\phi}} \mathbf{v} \in \mathbb{R}^d, \quad (4)$$

where  $\mathbf{S} = \text{w2v}(S) = \mathbb{R}^{K \times T}$  is an input sequence of text tokens vectorized with word2vec and  $\mathbf{v}$  is the final sentence embedding in the joint semantic-visual space, obtained after  $\ell_2$ -normalizing the output of SRU with parameters  $\phi$ .

### 3.2 Training

The full architecture is summarized in Fig. 3. The aim of training it is to learn the parameters  $\theta_{0:2}$  of the visual path, as well as all parameters  $\phi$  of the SRU text encoder. The goal is to create a joint embedding space for images and sentences such that closeness in this space can be interpreted as semantic similarity. This requires cross-modal supervision such that image-to-text semantic similarities are indeed enforced.<sup>1</sup>

**Contrastive triplet ranking loss** Following [19], we resort to a contrastive triplet ranking loss. Given a training set  $\mathcal{T} = \{(\mathbf{I}_n, S_n)\}_{n=1}^N$  of aligned image-sentence pairs – the

1. Note that mono-modal supervision can also be useful and relatively easier to get in the form, *e.g.*, of categorized images or of categorized sentences. Both are indeed used implicitly when relying on pre-trained CNNs and pre-trained text encoders. It is our case as well as far as the visual path is concerned. However, since our text encoder is trained from scratch, the only pure text (self-)supervision we implicitly use lies in the pre-training of word2vec.

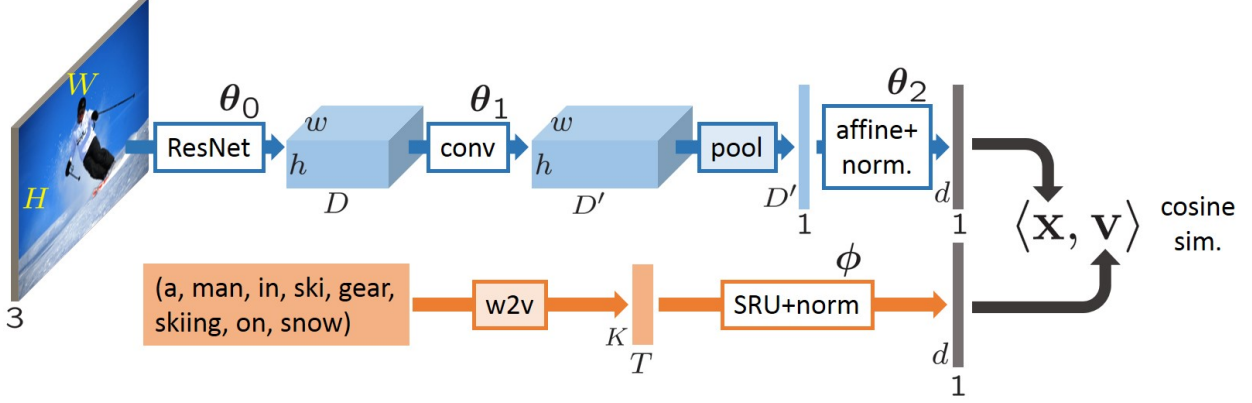


FIGURE 3 – **Details of the proposed semantic-visual embedding architecture.** An image of size  $3 \times W \times H$  is transformed into a unit norm representation  $\mathbf{x} \in \mathbb{R}^d$ ; likewise, a sequence of  $T$  tokenized words is mapped to a normalized representation  $\mathbf{v} \in \mathbb{R}^d$ . Training will aim to learn parameters  $(\theta_0, \theta_1, \theta_2, \phi)$  such that cross-modal semantic proximity translates into high cosine similarity  $\langle \mathbf{x}, \mathbf{v} \rangle$  in the joint embedded space. Boxes with white background correspond to trainable modules, with parameters indicated on top. In our experiments, the dimensions are  $K = 620$ ,  $D = 2048$  and  $D' = d = 2400$ .

sentence describes (part of) the visual scene – the empirical loss to be minimized takes the form :

$$\mathcal{L}(\Theta; \mathcal{T}) = \frac{1}{N} \sum_{n=1}^N \left( \sum_{m \in C_n} \text{loss}(\mathbf{x}_n, \mathbf{v}_n, \mathbf{v}_m) + \sum_{m \in D_n} \text{loss}(\mathbf{v}_n, \mathbf{x}_n, \mathbf{x}_m) \right), \quad (5)$$

where  $\Theta = (\theta_0, \theta_1, \theta_2, \phi)$  are the parameters to learn,  $\mathbf{x}_n = F(\mathbf{I}_n; \theta_{0:2})$  is the embedding of image  $n$ ,  $\mathbf{v}_n = \text{normSRU}_\phi(\text{w2v}(S_n))$  is the embedding of sentence  $n$ ,  $\{S_m\}_{m \in C_n}$  is a set of sentences unrelated to  $n$ -th image,  $\{\mathbf{I}_m\}_{m \in D_n}$  is a set of images unrelated to  $n$ -th sentence. The two latter sets are composed of negative (“contrastive”) examples. The triplet loss is defined as :

$$\text{loss}(\mathbf{y}, \mathbf{z}, \mathbf{z}') = \max \{0, \alpha - \langle \mathbf{y}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z}' \rangle\}, \quad (6)$$

with  $\alpha > 0$  a margin. It derives from triplet ranking losses used to learn metrics and to train retrieval/ranking systems. The first argument is a “query”, while the second and third ones stand respectively for a relevant (positive) answer and an irrelevant (negative) one. The loss is used here in a similar way, but with a multimodal triplet. In the first sum of Eq. 5, this loss encourages the similarity, in the embedding space, of an image with a related sentence to be larger by a margin to its similarity with irrelevant sentences. The second sum is analogous, but centered on sentences.

**Mining hard negatives** In [19, 16], contrastive examples are sampled at random among all images (resp. sentences) in the mini-batch that are unrelated to the query sentence (resp. image). Faghri *et al.* [9] propose instead to focus only on the hardest negatives. We follow the same strategy : For each positive pair in the batch, a single contrastive example

is selected in this batch as the one that has the highest similarity with the query image/sentence while not being associated with it. This amounts to considering the following loss for the current batch  $\mathcal{B} = \{(\mathbf{I}_n, S_n)\}_{n \in \mathcal{B}}$  :

$$\mathcal{L}(\Theta; \mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{n \in \mathcal{B}} \left( \max_{m \in C_n \cap \mathcal{B}} \text{loss}(\mathbf{x}_n, \mathbf{v}_n, \mathbf{v}_m) + \max_{m \in D_n \cap \mathcal{B}} \text{loss}(\mathbf{v}_n, \mathbf{x}_n, \mathbf{x}_m) \right). \quad (7)$$

Beyond its practical interest, this mining strategy limits the amount of gradient averaging, making the training more discerning.

### 3.3 Localization from embedding

As described in Section 2, several works on weak supervised localization [43, 7] combine fully convolutional architectures with specific pooling mechanisms such that the unknown object positions in the training images can be hypothesized. This localization ability derives from the activation maps of the last convolutional layer. Suitable linear combinations of these maps can indeed provide one heatmap per class.

Based on the pooling architecture of [7] which is included in our system, we derive the localization mechanism for our semantic-visual embedding. Let’s remind that in our case, the number of feature maps is arbitrary since we are not training on a classification task but on a cross-modal matching one. Yet, one can imagine several ways to leverage these maps to try and map an arbitrary vector of the joint embedding space into an arbitrary input image. When this vector is the actual embedding of a word or sentence, this spatial mapping should allow localizing the associated concept(s) in the image, if present. Ideally, a well-trained joint embedding should allow such localization even for

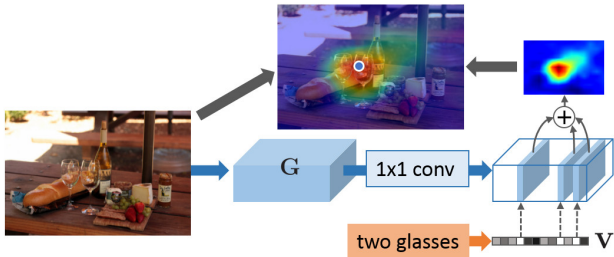


FIGURE 4 – **From text embedding to visual localization.** Given the feature maps  $\mathbf{G}$  associated to an image by our semantic-visual architecture and the embedding of a sentence, a heatmap can be constructed: Learned projection matrix  $A$  serves as a  $1 \times 1$  convolution; Among the  $d$  maps thus generated, the  $k$  ones associated with the largest values among the  $d$  entries of  $\mathbf{v}$  are linearly combined. If the sentence relates to a part of the visual scene, like “two glasses” in this example, the constructed heatmap should highlight the corresponding location. Blue dot indicates the heat maximum.

concepts that are absent from the training captions. To this end, we propose the following localization process (Fig. 4). Let  $\mathbf{I}$  be an image and  $\mathbf{G}$  its associated  $D'$  feature maps (Eq. 1). This stack is turned into a stack  $\mathbf{G}' \in \mathbb{R}^{w \times h \times d}$  of  $d$  heatmaps using the linear part of the projection layer  $p_{\theta_2}$ :<sup>2</sup>

$$\mathbf{G}'[i, j, :] = A\mathbf{G}[i, j, :], \forall (i, j) \in \llbracket 1, w \rrbracket \times \llbracket 1, h \rrbracket, \quad (8)$$

which is a  $1 \times 1$  convolution. Given  $\mathbf{v} \in \mathbb{R}^d$  the embedding of a word or sentence (or any unit vector in the embedded space) and  $K(\mathbf{v})$  the set of the indices of its  $k$  largest entries, the 2D heatmap  $\mathbf{H} \in \mathbb{R}^{w \times h}$  associated with the embedded text  $\mathbf{v}$  in image  $\mathbf{I}$  is defined as:

$$\mathbf{H} = \sum_{u \in K(\mathbf{v})} |\mathbf{v}[u]| \times \mathbf{G}'[:, :, u]. \quad (9)$$

In the next section, such heatmaps will be shown in false colors, overlaid on the input image after suitable resizing, as illustrated in Figs. 1 and 4. Note that [33] also proposes to build semantic heatmaps as weighted combinations of feature maps, but with weights obtained by back-propagating the loss in their task-specific network (classification or captioning net). Such heatmaps help visualize which image regions explain the decision of the network for this task.

## 4 Experiments

Using the MS-COCO dataset, we evaluate the overall quality of our model for cross-modal retrieval and visual grounding of phrases.

2. In other words, the pooling is removed. Bias and normalization being of no incidence on the location of the peaks, they are ignored.

## 4.1 Training

**Datasets** To train our model, we used the MS-COCO dataset [25]<sup>3</sup>. This dataset contains 123,287 images (train+val), each of them annotated with 5 captions. It is originally split into a training set of 82,783 images and a validation set of 40,504 images. The authors of [16] proposed another split (called rVal in the rest of the paper) keeping from the original validation set 5,000 images for validation and 5,000 for testing and using the remaining 30,504 as additional training data. To make our results comparable, we trained a model using each split. For evaluation, we also use the MS-COCO dataset, complemented with the annotations from Visual Genome dataset [21]<sup>4</sup> to get localization ground-truth when needed.

**Image pipeline** The image pipeline is pre-trained on its own in two stages. We start from original ResNet-152 [12] pre-trained on ImageNet classification task. Then, to initialize the convolutional adaptation layer  $g_{\theta_1}$ , we consider temporarily that the post-pooling projection is of size 1000 such that we can train both on ImageNet as well. Once this pre-training is complete, the actual projection layer  $p_{\theta_2}$  onto the joint space is put in place with random initialization, and combined with a 0.5-probability dropout layer. As done in [9], random rectangular crops are taken from training images and resized to a fixed-size square (of size  $256 \times 256$ ).

**Text pipeline** To represent individual word tokens as vectors, we used pre-trained word2vec with no further fine-tuning. The SRU text encoder [24] is trained from scratch jointly with the image pipeline. It has four stacked hidden layers of dimension 2400. Following [24], 0.25-probability dropout is applied on the linear transformation from input to hidden state and between the layers.

**Full model training** Both pipelines are trained together with pairs of images and captions, using Adam optimizer [18]. Not every part of the model is updated from the beginning. For the first 8 epochs only the SRU (parameters  $\phi$ ) and the last linear layer of the image pipeline ( $\theta_2$ ) are updated. After that, the rest of the image pipeline ( $\theta_{0:1}$ ) is also fine-tuned. The training starts with a learning rate of 0.001 which is then divided by two at every epoch until the seventh and kept fixed after that. Regarding mini-batches, we found in contrast to [9] that their size has an important impact on the performance of our system. After parameter searching, we set this size to 160. Smaller batches result in weaker performance while too large ones prevent the model from converging.

## 4.2 Results

**MS-COCO retrieval task** Our model is quantitatively evaluated on a cross-modal retrieval task. Given a query

3. <http://cocodataset.org>

4. <http://visualgenome.org/>

model	caption retrieval				image retrieval			
	R@1	R@5	R@10	Med. r	R@1	R@5	R@10	Med. r
Embedding network [37]	50.4	79.3	89.4	-	39.8	75.3	86.6	-
2-Way Net [8]	55.8	75.2	-	-	39.7	63.3	-	-
LayerNorm [2]	48.5	80.6	89.8	5.1	38.9	74.3	86.3	7.6
VSE++ [9]	64.6	-	95.7	1	52.0	-	92.0	1
Ours	<b>69.8</b>	<b>91.9</b>	<b>96.6</b>	1	<b>55.9</b>	<b>86.9</b>	<b>94.0</b>	1

TABLE 1 – **Cross-modal retrieval results on MS-COCO.** On both caption retrieval from images and image retrieval from captions, the proposed architecture outperforms the state-of-the-art systems. It yields an R@1 relative gain of 38% (resp. 40%) with respect to best published results [37] on cross-modal caption retrieval (resp. image retrieval), and 8% (resp 7.5%) with respect to best online results [9].

image (resp. a caption), the aim is to retrieve the corresponding captions (resp. image). Since MS-COCO contains 5 captions per image, recall at  $r$  (“R@ $r$ ”) for caption retrieval is computed based on whether at least one of the correct captions is among the first  $r$  retrieved ones. The task is performed 5 times on 1000-image subsets of the test set and the results are averaged. All the results are reported on Tab. 1. We compare our model with recent leading methods. As far as we know, the best published results on this task are obtained by the Embedding Network [37]. For caption retrieval, we surpass it by (19.4%,12.6%,7.2%) on (R@1,R@5,R@10) in absolute, and by (16.1%,11.6%,7.4%) for image retrieval. Three other methods are also available online, 2-Way Net [8], LayerNorm [2] and VSE++ [9]. The first two are on the par with Embedding Network while VSE++ reports much stronger performance. We consistently outperform the latter, especially in terms of R@1. Note that in [9], the test images are scaled such that the smaller dimension is 256 and centrally cropped to  $224 \times 224$ . Our best results are obtained with a different strategy : Images are resized to  $400 \times 400$  irrespective of their size and aspect ratio, which our fully convolutional visual pipeline allows. When using the scale-and-crop protocol instead, the recalls of our system are reduced by approximately 1.4% on average on the two tasks, remaining above VSE++ but less so. For completeness we tried our strategy with VSE++, but it proved counterproductive in this case. One of the key elements of the proposed architecture is the final pooling layer, adapted from WELDON [7]. To see how much this choice contributes to the performance of the model, we tried instead the Global Average Pooling (GAP) [43] approach. With this single modification, the model is trained following the exact same procedure as the original one. This results in less good results : For caption retrieval (resp. image retrieval), it incurs a loss of 5.3% for R@1 (resp. 4.7%) for instance, and a loss of 1.1% in accuracy in the pointing game.

**Visual grounding of phrases** We evaluate quantitatively our localization module with the pointing game defined by [40]. This task relies on images that are present both in MS-COCO val 2014 dataset and in Visual Genome dataset. The data contains 17,471 images with 86,5582 text

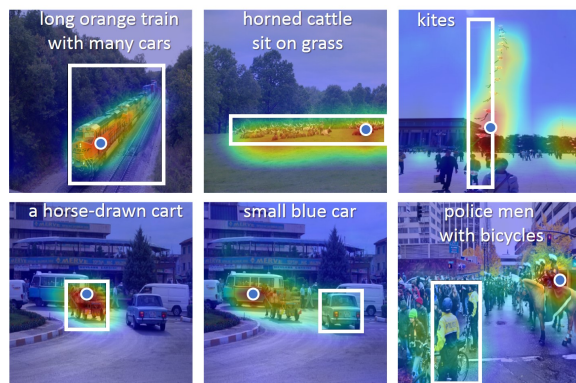


FIGURE 5 – **Pointing game examples.** Images from the Visual Genome dataset overlaid with the heatmap localizing the input text according to our system. The white box is the ground-truth localization of the text and the blue dot marks the location predicted by our model for this text. The first four predictions are correct, unlike the last two ones. In the last ones, the heatmap is nonetheless active inside the ground-truth box.

region annotations (a bounding box associated with a caption). The task consists in “pointing” the region annotation in the associated image. If the returned location lies inside the ground-truth bounding box, it is considered as a correct detection, a negative one otherwise. Since our system produces a localization map, the location of its maximum is used as output for the evaluation. For this evaluation, the number of feature maps from  $G'$  that are used to produce the localization map was set through cross-validation to  $k = 180$  (out of 2400). We keep this parameter fixed for all presented visualizations. The quantitative results are reported in Tab. 2 and some visual examples are shown in Fig. 5. We add to the comparison a baseline that always outputs the center of the image as localization, leading to a surprisingly high accuracy of 19.5%. Our model, with an accuracy of 33.8%, offers absolute (resp. relative) gains of 9.4% (resp. 38%) over [40] and of 14% (resp. 73%) over the trivial baseline.

**Towards zero-shot localization** The good performance

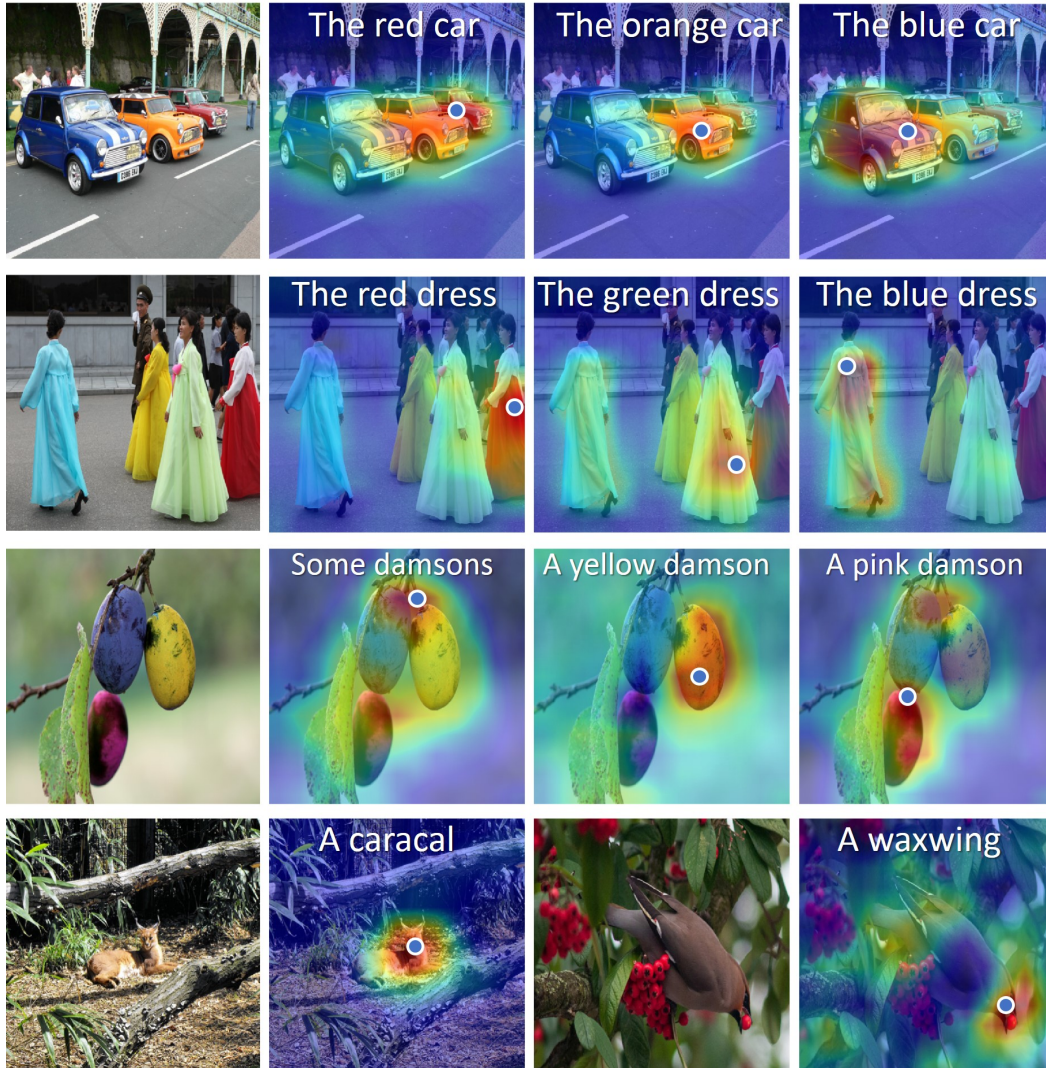


FIGURE 6 – **Toward zero-shot localization.** The first three rows show the ability to differentiate items according to their colors, even if, as in third example, the colors are unnatural and the concept has not been seen at training. This example, and the two last ones could qualify as “zero-shot localization” as damson, caracal, and waxwing are not present in MS-COCO train set.

Model	Accuracy
“center” baseline	19.5
Linguistic structure [40]	24.4
Ours (train 2017)	33.5
Ours (rVal)	33.8

TABLE 2 – **Pointing game results.** Our architecture outperforms the state-of-the-art system [40] by more than 9% in accuracy, when trained with either train or rVal split from MS-COCO.

we obtain in the pointing game highlights the ability of our system to localize visual concepts based on their embedding in the learned joint space. We illustrate further this

strength of the system with additional examples like the one already presented in Fig. 1. Going one step further, we conducted similar experiments with images from the web and concepts that were checked *not to appear in any of the training captions*, see Fig. 6.

## 5 Conclusion

We propose in this paper a new cross-modal text-image embedding pipeline. The use of a selective spatial pooling at the very end of the fully convolutional visual pipeline allows us to equip our system with a powerful mechanism to locate in images the regions corresponding to any text. Extensive experiments show that our model achieves high performance on cross-modal retrieval tasks as well as on phrase localization.



## Références

- [1] G. Andrew, R. Arora, J. Bilmes, and K. Livescu. Deep canonical correlation analysis. In *ICML*, 2013.
- [2] J. L. Ba, J. R. Kiros, and G. Hinton. Layer normalization. *arXiv preprint arXiv :1607.06450*, 2016.
- [3] G. Chechik, V. Sharma, U. Shalit, and S. Bengio. Large scale online learning of image similarity through ranking. *JMLR*, 11 :1109–1135, 2010.
- [4] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS w. on Deep Learning*, 2014.
- [5] J. Dai, Y. Li, K. He, and J. Sun. R-FCN : Object detection via region-based fully convolutional networks. In *NIPS*, 2016.
- [6] T. Durand, T. Mordan, N. Thome, and M. Cord. Wildcat : Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation. In *CVPR*, 2017.
- [7] T. Durand, N. Thome, and M. Cord. Weldon : Weakly supervised learning of deep convolutional neural networks. In *CVPR*, 2016.
- [8] A. Eisenschlat and L. Wolf. Linking image and text with 2-way nets. In *CVPR*, 2017.
- [9] F. Faghri, D. Fleet, J. R. Kiros, and S. Fidler. VSE++ : Improved visual-semantic embeddings. *arXiv preprint arXiv :1707.05612*, 2017.
- [10] A. Frome, G. Corrado, J. Shlens, S. Bengio, J. Dean, and T. Mikolov. DeViSE : A deep visual-semantic embedding model. In *NIPS*, 2013.
- [11] A. Frome, Y. Singer, F. Sha, and J. Malik. Learning globally-consistent local distance functions for shape-based image retrieval and classification. In *ICCV*, 2007.
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [13] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8) :1735–1780, 1997.
- [14] E. Hoffer and N. Ailon. Deep metric learning using triplet network. In *ICLRw*, 2015.
- [15] H. Hotelling. Relations between two sets of variates. *Biometrika*, 1936.
- [16] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015.
- [17] Y. Kim. Convolutional neural networks for sentence classification. In *EMNLP*, 2014.
- [18] D. Kingma and J. Ba. Adam : A method for stochastic optimization. In *ICLR*, 2014.
- [19] R. Kiros, R. Salakhutdinov, and R. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *arXiv preprint arXiv :1411.2539*, 2014.
- [20] R. Kiros, Y. Zhu, R. R. Salakhutdinov, R. Zemel, R. Urtasun, A. Torralba, and S. Fidler. Skip-thought vectors. In *NIPS*, 2015.
- [21] 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. *arXiv preprint arXiv :1602.07332*, 2016.
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [23] P. L. Lai and C. Fyfe. Kernel and nonlinear canonical correlation analysis. *Int. J. Neural Syst.*, 2000.
- [24] T. Lei and Y. Zhang. Training RNNs as fast as CNNs. *arXiv preprint arXiv :1709.02755*, 2017.
- [25] 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 *ECCV*, 2014.
- [26] L. Ma, Z. Lu, L. Shang, and H. Li. Multimodal convolutional neural networks for matching image and sentence. In *ICCV*, 2015.
- [27] A. Mignon and F. Jurie. PCCA : A new approach for distance learning from sparse pairwise constraints. In *CVPR*, 2012.
- [28] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv :1301.3781*, 2013.
- [29] H. Nam, J.-W. Ha, and J. Kim. Dual attention networks for multimodal reasoning and matching. *arXiv preprint arXiv :1611.00471*, 2017.
- [30] Z. Niu, M. Zhou, L. Wang, X. Gao, and G. Hua. Hierarchical multimodal LSTM for dense visual-semantic embedding. In *CVPR*, 2017.
- [31] J. Pennington, R. Socher, and C. Manning. GloVe : Global vectors for word representation. In *EMNLP*, 2014.
- [32] Z. Ren, H. Jin, Z. Lin, C. Fang, and A. Yuille. Joint image-text representation by gaussian visual-semantic embedding. In *ACMMM*, 2016.
- [33] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-CAM : Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017.
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv :1409.1556*, 2014.
- [35] Y. L. Sumit Chopra, Raia Hadsell. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, 2005.
- [36] L. Wang, Y. Li, and S. Lazebnik. Learning deep structure-preserving image-text embeddings. In *CVPR*, 2016.
- [37] L. Wang, Y. Li, and S. Lazebnik. Learning two-branch neural networks for image-text matching tasks. *T-PAMI*, 2017.
- [38] K. Q. Weinberger and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. *JMLR*, 10 :207–244, 2009.
- [39] J. Weston, S. Bengio, and N. Usunier. Wsabie : Scaling up to large vocabulary image annotation. In *IJCAI*, 2011.
- [40] F. Xiao, L. Sigal, and Y. Jae Lee. Weakly-supervised visual grounding of phrases with linguistic structures. In *CVPR*, 2017.
- [41] E. P. Xing, M. I. Jordan, S. J. Russell, and A. Y. Ng. Distance metric learning, with application to clustering with side-information. In *NIPS*, 2002.
- [42] F. Yan and K. Mikolajczyk. Deep correlation for matching images and text. In *CVPR*, 2015.
- [43] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016.