

# Multi-Modal Reasoning Graph for Scene-Text Based Fine-Grained Image Classification and Retrieval

Andres Mafla Sounak Dey Ali Furkan Biten Lluís Gomez Dimosthenis Karatzas  
Computer Vision Center, UAB, Spain

{andres.mafla, sdey, abiten, lgomez, dimos}@cvc.uab.es

## Abstract

Scene text instances found in natural images carry explicit semantic information that can provide important cues to solve a wide array of computer vision problems. In this paper, we focus on leveraging multi-modal content in the form of visual and textual cues to tackle the task of fine-grained image classification and retrieval. First, we obtain the text instances from images by employing a text reading system. Then, we combine textual features with salient image regions to exploit the complementary information carried by the two sources. Specifically, we employ a Graph Convolutional Network to perform multi-modal reasoning and obtain relationship-enhanced features by learning a common semantic space between salient objects and text found in an image. By obtaining an enhanced set of visual and textual features, the proposed model greatly outperforms previous state-of-the-art in two different tasks, fine-grained classification and image retrieval in the Con-Text[23] and Drink Bottle[4] datasets.

## 1. Introduction

Since the advent of written text to represent ideas, humans have employed it to communicate non-trivial and semantically rich information. Nowadays, text can be found in a ubiquitous manner in images and video, especially in urban and man-made environments[52, 24]. Extracting and analyzing such textual information in images jointly with the visual content is indispensable to achieve full scene understanding. In this work, we explore the role of such multi-modal cues, specifically in the form of visual and textual features to solve the task of fine-grained image classification and retrieval.

The task of fine-grained image classification (FGIC) consists of labeling a set of images that are visually alike. A lot of research on this problem has been oriented to differentiate visually similar objects such as birds[15], aircrafts[40], and dog breeds[26] among others, which more often

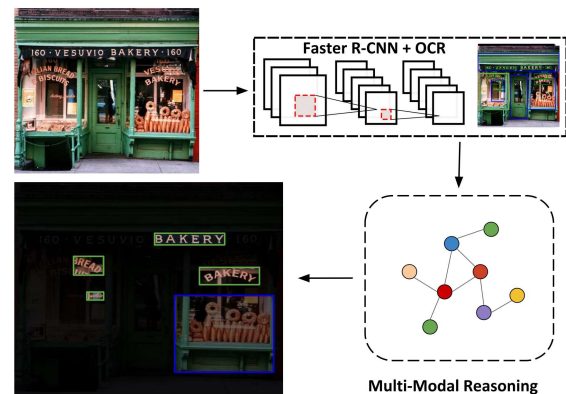


Figure 1. The proposed model uses a Graph-based Multi-Modal Reasoning (MMR) module to enrich location-based visual and textual features in a combined semantic representation. The network learns at the output of the MMR to map strong complementary regions of visual (blue) and text (green) instances to obtain discriminative features to perform fine-grained image classification and retrieval.

than not require domain specific knowledge. However, differentiating objects by leveraging available textual instances in the scene is an omnipresent practice in daily life. In this work, we focus on exploiting scene-text as the main discriminatory feature to perform FGIC. A seminal work on leveraging textual cues was presented by Movshovitz *et al.* [42], who showcased that in order to classify storefronts, a trained Convolutional Neural Network (CNN) had automatically learned to focus on scene text instances as the sole way to solve the given task. In the case of blurred or occluded text instances, the classification task is extremely challenging for humans as well. Consequently, scene text found in an image serves as an additional discriminative signal that a model should incorporate into its design. Further research has been devoted to explicitly leveraging textual cues in the task of FGIC. Similar to our work, Karaoglu *et al.* [23, 22] introduces a simple pipeline to perform fine-grained classification using scene text and extending the previous work, an attention mechanism is proposed by Bai

*et al.* [4] to learn a common semantic space. In a different approach, Mafla *et al.* [38] learns a morphological space by using textual instances as discriminative features rather than semantics to solve this task.

Departing from previous approaches, we exploit a structural representation between the studied modalities. Our work, summarized in Figure 1 with publicly available code at <sup>1</sup>, focuses on learning an enhanced visual representation that incorporates reasoning between salient regions of an image and scene text to construct a semantic space over which fine-grained classification is performed. In this example, we can observe that relevant regions such as the text "Bakery" and "Bread" are associated with a visual region that depicts pastry, both important cues to classify the given image. Additionally, we show experiments of fine-grained image retrieval, using the same multi-modal representation, in the two evaluated datasets. Overall, the main contributions can be summarized as follows:

- We propose a novel architecture that greatly surpasses previous state-of-the-art results in two datasets by more than 5% on fine-grained classification and 10% on image retrieval by considering text and visual features of an image.
- We design a fully end-to-end trainable pipeline that incorporates a Multi-Modal Reasoning module that combines textual and visual features that do not rely on ensemble models or pre-computed features.
- We provide exhaustive experiments in which we analyze the effectiveness of different modules in our model architecture and the importance of scene text towards comprehensive models of image understanding.

## 2. Related Work

### 2.1. Scene Text Detection and Recognition

Localizing and recognizing text instances found in a natural image is a challenging problem due to the variability, orientation, occlusion, and background noise among other factors [10]. Deep learning-based methods began with the work proposed by [21] which focused on a sliding window and a CNN to filter the proposals. The proposals were used as input into another CNN that posed the task as a classification problem over a large fixed dictionary of words. Later works take object detection pipelines such as YOLO [46] used by [18] to obtain a Fully Convolutional Neural Network along with a focus on generating synthetic training data, which later became the go-to data to train text detectors and recognizers. Along these lines, a variation of SSD [36] is presented by [33, 32] to develop a text detector which easily integrates with a module trained for recognition. Methods that focus on an end-to-end recognition

---

<sup>1</sup>[https://github.com/AndresPMD/GCN\\_classification](https://github.com/AndresPMD/GCN_classification)

have been explored by [7] based on Faster R-CNN [47], which performs text detection and incorporates a Connectionist Temporal Classification (CTC) [17] to recognize a given text instance. Similarly, [20] presents a CNN as a region-based feature extractor, features that are fed to two attention-based Long-Short Term Memories (LSTM) to predict bounding boxes and recognize the textual proposals. Multi-lingual models have been proposed as in the case of [8], work that uses a CNN as an encoder and a CTC to decode the characters from a set of different languages.

On a different approach, the Pyramidal Histogram of Characters (PHOC) [1] is used to represent words and it has been amply used in text spotting in documents [50] and text retrieval in natural images [16, 39]. Despite all the progress done in scene text detection and recognition, it remains an open problem in the computer vision community, with a special focus placed lately on multi-oriented text localization and recognition.

### 2.2. Fine-Grained Classification

The task of Fine-Grained Image Classification (FGIC) focuses on finding discriminative visual regions that often require domain-specific knowledge to correctly perform the labeling task [53]. Different to solely visual-based FGIC methods, there has been growing interest to use textual cues to achieve this task by incorporating two modalities.

Closely related to this work, the initial approach taken by [22] was to extract scene text and construct a bag of words, while the visual features were obtained by employing a pre-trained GoogLeNet [51]. Soon after, [4] proposes the usage of Textboxes [33] to read scene-text in an image, a CNN to obtain visual features along with an attention mechanism and a concatenation of the final features to learn a semantic space suitable for scene-text based FGIC. Later work performed by [38] employs a CNN as a visual feature extractor and uses the PHOC representation of a word along with the Fisher Vector [45] to learn a space based on the morphology of text instances to overcome Optical Character Recognition (OCR) errors. Several fusion methods are explored in the work by [38] but finally a concatenation of features is performed to solve the task of image classification and retrieval.

### 2.3. Multi-Modal Fusion and Reasoning

Several fusion-based techniques such as Multimodal Compact Bilinear Pooling (MCB) [14, 12], Low-rank Bilinear Attention Network (MLB) [27] and Block [5] have been explored to model relationships between language and vision. To model this interaction, attention-based [3] approaches also have been proposed [2, 54, 25]. With the aim of designing models capable of reasoning, the intrinsic synergy between visual and textual features has been explored. Work such as [57, 30] employ variations of an

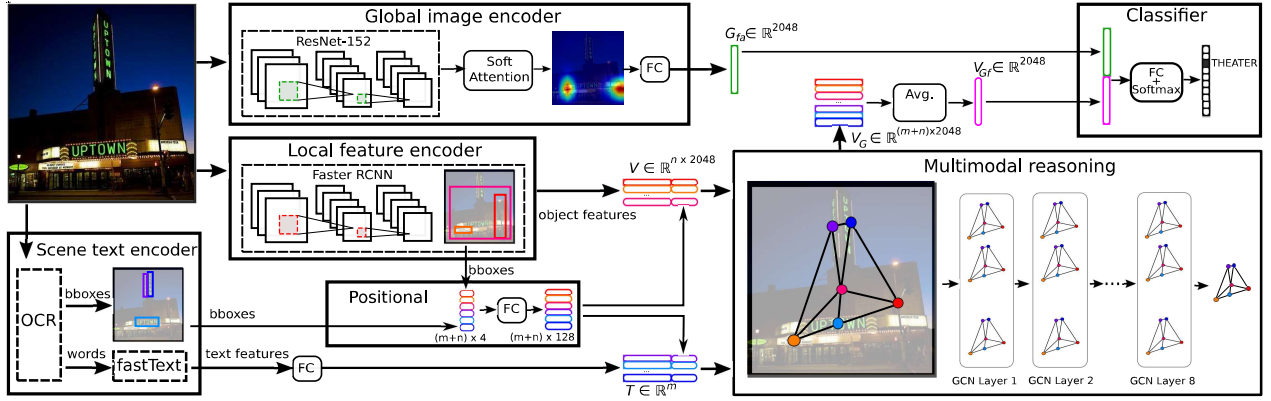


Figure 2. Detailed model architecture. The proposed model combines features of regions of scene text and visual salient objects by employing a graph-based Multi-Modal Reasoning (MMR) module. The MMR module enhances semantic relations between the visual regions and uses the enriched nodes along with features from the Global Encoder to obtain a set of discriminatory signals for fine-grained classification and retrieval.

LSTM and a Gated Recurrent Unit (GRU) to perform reasoning in a sequential manner. However, significant advances have been made by the usage of Graph Convolutional Networks (GCN) [28], due to the proven capability of modeling relationships [48] between nodes in a given graph. Along this road, GCNs have been successfully used in tasks that require reasoning such as VQA [43, 49, 13], image captioning [31, 55] and image-sentence retrieval [30, 34].

In this work, we propose a method to learn a richer set of visual features and model a more discriminative semantic space by employing a GCN. To the best of our knowledge, this is the first approach that integrates multimodal sources that come in the form of visual along with textual features jointly with positional encoding into a GCN pipeline that performs reasoning for the task of scene-text based fine-grained image classification and retrieval.

### 3. Method

In this section, we detail each of the components that comprise the proposed architecture. Figure 2 depicts the overall scheme of the proposed model, which is formed by 6 different modules: global image encoder, local feature encoder, text encoder, positional encoder, multi-modal reasoning graph and classification module. The local feature encoder employs features extracted based on the regions of interest obtained by a Faster R-CNN [47] in a similar manner as the bottom-up attention model [2]. The scene text encoder uses an OCR model to obtain scene text and further embed it into a common space. The goal is to obtain multi-modal node representations that leverage the semantic relationships found between salient objects and text instances within an image that are discriminative enough to perform fine-grained classification.

#### 3.1. Global Image Encoder

We employ a CNN as an encoder, which in our case is a ResNet-152 [19] pre-trained on ImageNet [11] to acquire global image features. Particularly, given an image  $I$  we take the output features before the last average pooling layer, which output is denoted as  $G_f = \psi(I)$ . In order to obtain a more descriptive set of global features and due to its differentiable properties, we compute a soft attention mechanism on top of the global features. This self-attention mechanism yields an attention mask,  $attn_{mask}$ , that assigns weights on different regions of the input image. The attention weights are learned in an end-to-end manner by convolving  $1 \times 1$  kernels projected into a single-dimensional filter and later followed by a Softmax function. In order to obtain the final attended global features, the attention mask is broadcasted and multiplied with the global features, which result is added to the global features  $G_f$  to later be used as input of a Fully-Connected layer,  $FC$ , in the form of:

$$G_{fa} = FC(G_f + (G_f \times attn_{mask})) \quad (1)$$

where  $G_{fa} \in \mathbb{R}^{1 \times D}$ ,  $G_f \in \mathbb{R}^{H \times W \times D}$ ,  $attn_{mask} \in \mathbb{R}^{H \times W}$  stands for the final encoded global features, where  $D = 2048$ ,  $H = 7$  and  $W = 7$ .

#### 3.2. Local Feature Encoder

Following [2], we employ a Faster R-CNN [47] pre-trained on Visual Genome [29] as the extractor of local visual features. This approach allows us to obtain salient image regions that are potentially discriminative for our task. We use an IoU threshold of 0.7 and a confidence threshold of 0.3, and sort the obtained predictions before the last average pooling layer to use the top  $n$  most confident regions of interest. Thus, we can represent the output of an image  $I$  with a set of region features  $R_f =$

$\{(r_1, bbox_{r_1}), \dots, (r_n, bbox_{r_n})\}$ ,  $r_i \in \mathbb{R}^d$ , where  $r_i$  is the  $i^{th}$  region of interest and  $bbox_{r_i}$  is the  $r_i$ 's corresponding bounding box coordinates normalized with respect to the image. In our experiments, we set  $n = 36$  and the obtained features have a dimension of  $d = 2048$ . In order to encode the local visual features, we project the features through a fully-connected layer.

In this manner, we obtain the final encoded local features that will serve as input to the multi-modal GCN in the form of  $V_f = \{v_1, \dots, v_n\}$ ,  $v_i \in \mathbb{R}^D$ , where  $D = 1920$  is the dimension of the final embedding space. We use  $D = 1920$  to further add positional encoding information  $D = 128$  to have a final feature representation of  $D = 2048$ . The bounding boxes obtained to represent these regions are later used as input into the positional encoder module. If there are less than  $n = 36$  regions in an image, a zero padding scheme is adopted.

### 3.3. Text Encoder

To extract text contained in an image, we ran several public state of the art text recognizers as well as a commercial OCR model provided by Google<sup>2</sup>. We extract the transcriptions of each word, denoted as  $w_i$ , as well as the corresponding bounding boxes,  $bbox_{w_i}$ . In particular, we extract the top  $m$  most confident textual instances found in an image. The transcriptions are embedded using fastText [6] and the bounding boxes will be used as input in the positional encoder branch. We employ the fastText embedding due to its capability of encoding word morphology in the form of n-grams as well as preserving a semantic space similar to Word2Vec [41] while at the same time dealing with out of vocabulary words. Analogously to the case of local features, we project the obtained embedded textual features by passing them through a fully-connected layer. The final textual features are represented by  $T_f = \{t_1, \dots, t_m\}$ ,  $t_i \in \mathbb{R}^D$ , where  $D = 1920$  is the dimension of the final embedding space and  $m = 15$  is the number of text proposals extracted from an image. In the case that there is no text found in a given image, similarly to the local encoder module, zero padding is employed.

### 3.4. Positional Encoder

Encoding the position of objects and text instances within an image can provide important relational information about the scene. For example text found on top of a building often refers to its class in an explicit manner contrary to text found in any other location in the image. To meet this end, we design a positional encoding that takes as input a predicted bounding box of an object or text instance. The input to the positional encoder describes the top left  $(x_1, y_1)$ , and bottom right  $(x_2, y_2)$  coordinates normalized according to the image size, and

<sup>2</sup><https://cloud.google.com/vision/>

is a concatenation of the bounding boxes of the local and text regions of interest. The bbox matrix is given by:  $bboxes_{input} = \{bbox_{r_1}, \dots, bbox_{r_n}, bbox_{t_1}, \dots, bbox_{t_m}\}$  where  $bbox_i = (x_1, y_1, x_2, y_2)$ . In order to encode them, we pass the bounding boxes over a fully-connected in a similar way as the same as previous sections. The final encoded representation can be described as:  $bboxes = \{bbox_{r_1}, \dots, bbox_{r_n}, bbox_{t_1}, \dots, bbox_{t_m}\}$ ,  $bbox_i \in \mathbb{R}^b$ , in which the dimension  $b = 128$  represents the final encoded bounding boxes.

### 3.5. Multi-modal Reasoning Graph

Due to the showcased capability of graphs to describe reasoning between objects [49, 55, 13, 34], we construct a richer set of region-based visual descriptors that exploit the semantic correlation between visual and textual features. In order to do so, we initialize the node features as local visual features and textual features concatenated with their respective positional encoding of bounding boxes. We can describe the node features as:

$$V = \{(v_1, bbox_{r_1}), \dots, (v_n, bbox_{r_n}), (t_1, bbox_{t_1}), \dots, (t_m, bbox_{t_m})\}, V \in \mathbb{R}^{(n+m) \times D}$$

where  $n, m$  is the number of visual and textual features, respectively. In our case,  $n + m = 51$  and  $D = 1920 + 128 = 2048$ . Furthermore, we construct the affinity matrix  $R$  which measures the degree of correlation of between two visual regions. The construction of the affinity matrix is given by:

$$R_{ij} = \phi(k_i)^T \gamma(k_j) \quad (2)$$

where  $k_i, k_j \in V$ ,  $\phi(\cdot)$  and  $\gamma(\cdot)$  are two fully connected layers that are learned end-to-end by back propagation at training time. If we define  $k = n + m$ , then the obtained affinity matrix consists of a shape  $k \times k$ . Once  $R$  is calculated, we can define our graph by  $G = (V, R)$ , in which the nodes are represented by the local and textual features  $V$ , and the edges are described by  $R$ . The obtained graph describes through the affinity matrix  $R$  the degree of semantic and spatial correlation between two nodes. We use the formulation of Graph Convolutional Networks given by [28] to obtain reasoning over the nodes and edges. Particularly, we use residual connections in the GCN formulation as it is presented by [30]. We can write the equation that describes a single Graph Convolution layer performed as:

$$V_g^l = W_r^l (R^l V^{l-1} W_g^l) + V^{l-1} \quad (3)$$

where  $R \in \mathbb{R}^{k \times k}$  is the affinity matrix,  $V \in \mathbb{R}^{k \times D}$  the local visual features,  $W_g \in \mathbb{R}^{D \times D}$  is a learnable weights matrix of the GCN,  $W_r \in \mathbb{R}^{k \times k}$  corresponds to the residual weights matrix and  $l$  is the number of GCN layer. Notice

that passing  $V$  through the GCN layer, a richer set of multi-modal features is obtained. In order to find an enhanced representation of the visual features we apply  $l = 8$  GCN layers in total, which finally yields a set of enriched nodes that represent the visual features  $V_G$  such that:

$$V_G = \{v_{g1}, \dots, v_{gk}\}, V_G \in \mathbb{R}^{k \times D}$$

### 3.6. Classification

In order to combine the global  $G_{fa}$  and the enriched local and textual  $V_G$  visual features, firstly we perform an average pooling of the  $V_G$  tensor. Specifically, we can rewrite the final local feature vector  $V_{Gf}$  as:

$$V_{Gf} = \frac{1}{k} \sum_{n=1}^k V_{gi} \quad (4)$$

Lastly, we simply concatenate the two obtained vectors  $V_{Gf}$  and  $G_{fa}$ , to obtain the final vector  $F$  that is used as input for the final fully-connected layer for classification denoted by:  $F = [G_{fa}, V_{Gf}]$

By applying a softmax to the output of the final layer, we obtain a probability distribution of a class label given an input image. The model is trained in an end-to-end fashion optimized with the cross entropy loss function described by:

$$J(\theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^C y_i^n \log(p_i^n) \quad (5)$$

Where,  $C$  is the number of classes,  $N$  the dataset samples such that each pair contains an annotation  $\{x^{(n)}, y^{(n)}\} | n = 1, 2, \dots, N$ , and  $p^n$  is the predicted output label.

## 4. Experiments and Results

This section presents an introduction to the datasets employed in this work, as well as the implementation details, ablation studies performed, and a thorough analysis of the results obtained in the experiments conducted.

### 4.1. Datasets

The *Con-Text* dataset was introduced by Karaoglu *et al.* [23] and is a subset of ImageNet [11], constructed by selecting the sub-categories of "building" and "place of business". This dataset contains 24,255 images in total divided into three-folds to divide training and testing sets. This dataset introduces 28 visually similar categories of images such as Cafe, Pizzeria, and Pharmacy in which in order to perform fine-grained classification, text is a necessary cue to solve otherwise a very difficult task even for humans. This dataset closely resembles natural circumstances due to the fact that the images are taken without considering scene

text instances, thus some images do not have text present in them.

The *Drink Bottle* dataset was presented by Bai *et al.* [4] and as the *Con-Text* dataset, it is a subset of images of ImageNet [11], specifically taken from the sub-categories of soft drink and alcoholic drink. The dataset is divided in three-folds as well and contains 18,488 images. There are 20 image categories which include visually similar instances such as Coca Cola, Pepsi Cola and Cream Soda. Akin to the *Con-Text* dataset, some images contain scene-text while others do not have it.

### 4.2. Implementation Details

In our experiments in order to extract visual regions of an image, we use the same settings as [2]. We take the top  $n = 36$  ROIs and encode them along with their bounding boxes into a common space of 2048-d. The transcribed text is sorted by confidence score and we take the top  $m = 15$  confident predictions. We embed the textual instances by using a pre-trained fastText model with 1 million 300-d word vectors, trained with sub-word information on Wikipedia2017, UMBC webbase corpus and statmt.org news dataset. The obtained 300-d textual vectors are projected with the corresponding bounding boxes into a 2048-d space. The Faster R-CNN [47] from [2] and the OCR models, both employed as initial feature extractor modules use pre-trained weights and are not updated at training stage. The rest of the weights of each module in the model are learned in an end-to-end manner during training. The graph-based multimodal reasoning module employs 8 multi-modal GCN layers to obtain the final enriched visual features. In the last full-connected layer before classification, we employ a dropout rate of 0.3 to avoid over-fitting on the evaluated datasets. In general, we employ Leaky ReLU as an activation function in all layers except the last one, in which we use a Softmax to compute the class label probabilities. The proposed model is trained for 45 epochs, but an early stop condition is employed. We use a combination of optimizers comprised by RAdam [35] and Lookahead [58]. The batch size employed in all our experiments is 64, with a starting learning rate of 0.001 that decays by a factor of 0.1 on the epochs 15, 30 and 45. The momentum value used on the optimizers is 0.9 and the weight decay is 0.0005.

### 4.3. Comparison with the State-of-the-Art

We show the experimental results of our method compared to previous state-of-the-art on Table 1. We can note that the performance obtained in the *Con-Text* significantly surpasses the previous best performing method by 5.9%. The improvement in the *Drink-Bottle* dataset is more modest, of about 1.98%, however it is still significant.

We believe the improvement is greater in *Con-Text* due to the text instances found in it, which refer mostly to busi-

| Method              | OCR               | Emb.               | Con-Text     | Bottles      |
|---------------------|-------------------|--------------------|--------------|--------------|
| Karao.[23]          | Custom            | BoB <sup>1</sup>   | 39.0         | –            |
| Karao.[22]          | Jaderberg         | Probs <sup>2</sup> | 77.3         | –            |
| Bai[4]              | Textboxes         | GloVe              | 78.9         | –            |
| Bai[4] <sup>†</sup> | Textboxes         | GloVe              | 79.6         | 72.8         |
| Bai[4] <sup>†</sup> | Google OCR        | GloVe              | 80.5         | 74.5         |
| Mafla[38]           | SSTR-PHOC         | FV                 | 80.2         | 77.4         |
| Proposed            | E2E-MLT           | fastText           | 82.36        | 78.14        |
| Proposed            | SSTR-PHOC         | PHOC               | 82.77        | 78.27        |
| Proposed            | SSTR-PHOC         | FV                 | 83.15        | 77.86        |
| <b>Ours</b>         | <b>Google OCR</b> | <b>fastText</b>    | <b>85.81</b> | <b>79.87</b> |

Table 1. Classification performance of state-of-the art methods on the Con-Text and Drink-Bottle datasets. The results depicted with <sup>†</sup> are based on an ensemble model. The embeddings labeled as <sup>1</sup> refer to a Bag of Bigrams, <sup>2</sup> is a probability vector along a dictionary. The acronym FV stands for Fisher Vector. The metric depicted is the mean Average Precision (mAP in %).

ness places without many out of vocabulary words, therefore a semantic space for classification is more discriminative when compared to the Drink-Bottle dataset. To provide further insights, we conducted experiments by employing the final model along with different OCRs and word embeddings in both datasets. It is essential to note that state-of-the-art results are achieved by the usage of other OCRs as well, showing that the proposed pipeline still outperforms previous methods. Results showing the classification scores of each evaluated class and further analysis are shown in the Supplementary Material section.

When comparing to previous methods, it is worth revisiting previous approaches. The results reported by [4] used an ensemble of classifiers to reach the obtained performance. As an additional experiment to showcase the effect of using the same OCR as our proposed model is included, and it shows that our model vastly outperforms the evaluated pipeline not because of the OCR system employed. On the other side, the work done by [38] requires offline pre-computation of the Fisher Vector by training a Gaussian Mixture Model and tuning the hyper-parameters involved. In this manner, the method proposed in this work does not require an ensemble and the features used are learned in an end-to-end manner at training time. We clearly show that the proposed pipeline surpasses other approaches even when employing a set of different scene-text OCRs.

With the aim of offering additional insights, we present in Table 2 the performance of previous state of the art methods compared with our proposed method in a subset of the test set such that the evaluated images either contain scene-text or not. The results show the average performance along with 3 different splits of each dataset. We can observe that our model is able to perform better than previous approaches in both scenarios while a more significant improvement is achieved in images that contain scene-text,

| Method      | Con-Text     |              | Drink Bottle |              |
|-------------|--------------|--------------|--------------|--------------|
|             | I + T        | I - T        | I + T        | I - T        |
| Bai [4]     | 78.92        | 71.63        | 71.61        | 62.25        |
| Mafla [38]  | 80.94        | 72.59        | 78.57        | 68.97        |
| <b>Ours</b> | <b>86.76</b> | <b>74.31</b> | <b>82.75</b> | <b>69.19</b> |

Table 2. Classification performance of the proposed method on the subset of images from the test set of the Con-Text and Drink-Bottle datasets such that the images: contain scene-text (I + T) and do not contain scene-text (I - T). The metric depicted is the mean Average Precision (mAP in %).

which we treat as the major discriminative feature to perform the task of fine-grained classification.

#### 4.4. Importance of Textual Features

In order to assess the importance of the scene text found in images, we follow the previous works [22, 4, 38] by defining two different evaluation baselines, the visual features based and the textual features based. Moreover, due to the fact that the evaluated datasets do not contain text transcriptions as ground truth, we evaluated the effectiveness of the OCR employed in the fine-grained classification task.

The visual only evaluates all the test set images by only employing the global encoder features  $G_f$  in the first scenario and the global encoder along with the self attention features  $G_{fa}$  in the second scenario. In both cases the output of the global encoder, a 2048-d feature vector, is directly passed through a fully connected layer to obtain the final classification prediction. In the textual only, the baselines are evaluated only in the subset of images which contained spotted scene text. The results of each baseline by employing visual only, different OCRs and word embeddings are shown in Table 3.

Following a previous approach [38], we employ  $m = 15$  text instances and pre-trained word embeddings that yield 300-d vectors in the case of Word2Vec [41], GloVe [44] and fastText [6]. The textual tensor obtained is used as input to a fully connected layer, which output is used for classification purposes. In our experiments we evaluate two additional state-of-the-art scene text recognizers, FOTS [37] and the commercially used Google OCR Cloud Vision based on an API. We note that the embedding that performs the best is fastText due to the capability of embedding out of vocabulary words by using character n-grams. Regarding the results, it was found that the best performing standard recognizer is the Google OCR, which employs a more compact (300-d) vector compared to a PHOC or a Fisher Vector. The PHOC embedding employs a 604-d feature vector along with  $m = 15$  and the Fisher Vector is a single 38400-d vector in our experiments. Overall, by using only textual features, the Fisher Vector based on PHOCs remains as the best performing descriptor. However, besides the high dimensional vector employed, extensive offline pre-computation

|         | Model                                   | Con-Text     | Bottles      |
|---------|---|--------------|--------------|
| Visual  | CNN                                     | 62.11        | 65.15        |
|         | CNN + Self Attention                    | 63.78        | 66.62        |
| Textual | Texspotter+w2v <sup>†</sup>             | 35.09        | 50.68        |
|         | Texspotter+glove <sup>†</sup>           | 34.52        | 50.26        |
|         | Texspotter+fastText <sup>†</sup>        | 36.71        | 51.93        |
|         | E2E_MLT+w2v <sup>†</sup>                | 44.36        | 43.98        |
|         | E2E_MLT+glove <sup>†</sup>              | 44.25        | 42.64        |
|         | E2E_MLT+fastText <sup>†</sup>           | 45.07        | 44.31        |
|         | FOTS+w2v                                | 43.22        | 41.33        |
|         | FOTS+glove                              | 43.71        | 41.85        |
|         | FOTS+fastText                           | 44.19        | 42.69        |
|         | Google OCR+w2v                          | 53.87        | 53.47        |
|         | Google OCR+glove                        | 54.48        | 54.39        |
|         | <b>Google OCR+fastText</b>              | <b>55.61</b> | <b>55.16</b> |
|         | PHOC <sup>†</sup>                       | 49.18        | 52.39        |
|         | <b>Fisher Vector (PHOC)<sup>†</sup></b> | <b>63.93</b> | <b>62.41</b> |

Table 3. Visual only and Textual only results. The textual only results were performed on the subset of images that contained spotted text. The results with <sup>†</sup> were reported by [38]. The metric depicted is the mean Average Precision (mAP in %).

is required to obtain such descriptor. Nonetheless, as it can be seen in Table 1, the FV descriptor does not achieve the best results in our final model.

#### 4.5. Ablation studies

In this section, we present the incremental improvements and the effects obtained by the addition of each module that comprises the final architecture in the method proposed. Table 4 shows the quantitative results of adding components in the baseline model. Namely, we evaluate the effect of using self-attention and the multi-modal reasoning (MMR) module. We successively add to the attended global features ( $G_{fa}$ ), local features ( $V_f$ ), textual features ( $T_f$ ) and the bounding boxes ( $bboxes$ ) of both used in the Positional Encoder. In order to assess the effectiveness of the multi-modal reasoning graph module, we compare a model that uses the Faster R-CNN ROIs without the usage of the MMR. It is observed that solely by using the Faster R-CNN features, an important boost is achieved. One of the biggest improvements is reached by the usage of scene text, which enforces the idea that textual information is essential to successfully discriminate between visually similar classes. By the incorporation of scene text, an improvement of 9.7% is gained in Con-Text and 2.5% in the Drink-Bottle datasets. Nonetheless, the improvement is accentuated by the usage of the MMR module, which produces as output richer local and textual features coming from the graph nodes. Finally by adding the positional encoder module into the MMR, another increase in the results is achieved. This encourages us to think that the MMR module learns relationships coming from semantic and spatial information. Insights into the at-

| Features  | Con-Text     | Drink Bottle |
|---|--------------|--------------|
| $G_f$   | 62.11        | 65.15        |
| $G_{fa}$  | 63.78        | 66.62        |
| <b>without MMR</b>                              |              |              |
| $G_{fa} + V_f$                                  | 70.48        | 73.21        |
| $G_{fa} + V_f + T_f$                            | 78.72        | 76.43        |
| $G_{fa} + V_f + T_f + bboxes$                   | 80.12        | 77.51        |
| <b>with MMR</b>                                 |              |              |
| $G_{fa} + V_f$                                  | 72.88        | 74.96        |
| $G_{fa} + V_f + T_f$                            | 82.51        | 77.46        |
| $V_f + T_f + bboxes$                            | 84.33        | 75.42        |
| <b><math>G_{fa} + V_f + T_f + bboxes</math></b> | <b>85.81</b> | <b>79.87</b> |

Table 4. Quantitative results of the different components that form the proposed model.  $G_f$ : Global features,  $G_{fa}$ :  $G_f$  + Self-Attention,  $V_f$ : Local Features,  $T_f$ : Text Features,  $bboxes$ : Bounding Box information used by the Positional Encoder, MMR: Multi-modal Reasoning. Results are shown in terms of the mAP(%).

| Projection         | Fusion        | Con-Text     | Drink Bottle |
|--------------------|---------------|--------------|--------------|
| Attention          | MLB [27]      | 80.83        | 78.26        |
| Attention          | Block [5]     | 80.82        | 78.42        |
| Attention          | Concat        | 81.09        | 78.45        |
| GRU                | MLB [27]      | 83.12        | 78.21        |
| GRU                | Block [5]     | 83.8         | 78.74        |
| GRU                | Concat        | 83.93        | 78.89        |
| Avg Pooling        | MLB [27]      | 84.23        | 78.56        |
| Avg Pooling        | Block [5]     | 85.11        | 79.15        |
| <b>Avg Pooling</b> | <b>Concat</b> | <b>85.81</b> | <b>79.87</b> |

Table 5. Results obtained by employing different Projection and Fusion strategies on all the modules of our pipeline. Results are shown in terms of the mAP(%).

tention masks learned and the reasoning coming from the MMR by using visual and textual regions can be found in the Supplementary Material section.

Furthermore, we explore in our work several projection and fusion methods which are shown in Table 5. In our experiments, Projection refers to the strategy used to reduce the dimensionality of the output tensor coming from the MMR as  $V_G$  to obtain a single vector  $V_{Gf}$ . Late Fusion showcases the method employed to combine the features coming from  $V_{Gf}$  and  $G_{fa}$ . Due to several works showing performance gains by the usage of attention [56, 54] and Recurrent Neural Networks [30, 9] as reasoning modules, we explored those alternatives, however, no improvements were found. In the same manner, as it is presented by [38], we explored two additional fusion mechanisms, MLB [27] and Block [5] but no gains were obtained compared to feature concatenation.

#### 4.6. Qualitative Results

Qualitative results of the fine-grained image classification task are shown in Figure 3. By reviewing the samples obtained, we can note that our model is capable of learning a semantic space which combines successfully visual



Figure 3. Classification predictions. The top-3 probabilities of a class are shown as well as the Ground Truth label performed on the test set. Without recognizing textual instances some images are extremely hard to classify even for humans. Text in blue and red is used to show correct and incorrect predictions respectively. Best viewed in color.

and textual signals coming from a single image. Classified samples such as “Pizzeria”, “Tea House” and “Diner” often contain similar semantic classes ranked on second and third positions. Images belonging to the Drink Bottle dataset on the second row, are correctly classified even though text instances belong to specific brands, thus showing generalization capability of our method. The seventh image on the first row is wrongly classified as “Theatre” due to OCR recognition errors and a lack of strong enough visual cues. The remaining wrongly classified images are very challenging and contain some degree of ambiguity even for humans.

#### 4.7. Fine-Grained Image Retrieval

As an additional experiment that highlights the capabilities of the proposed model, we show the results obtained in Table 6 by performing query-by-example (QbE) image retrieval. In QbE, a system must return images in the form of a ranked list that belongs to the same class as the image used as a query. To provide comparable results and following the work from [4, 38], we use the vector of class probabilities as the image descriptor without using a specific metric-learning method. This vector is used to retrieve the nearest samples computed by the usage of the cosine similarity as a distance metric.

In our experiments, the query, as well as the database is formed by unseen samples at training time. The results demonstrate that a very significant boost of 10.98% and 2.48% in Con-Text and Drink-Bottle is achieved respectively. The lower gain in the Drink-Bottle dataset directly depends on the harder to recognize text instances, as well as the low image quality of several samples that directly affects the model performance.

| Method    | Con-Text     | Drink Bottle |
|-----------|--------------|--------------|
| Bai[4]    | 62.87        | 60.80        |
| Mafra[38] | 64.52        | 62.91        |
| Proposed  | <b>75.50</b> | <b>65.39</b> |

Table 6. Retrieval results on the evaluated datasets. The retrieval scores are depicted in terms of the mAP(%).

Qualitative results that show the robustness of the model, as well as experiments addressing the importance of text can be found in the Supplementary Material section.

## 5. Conclusions

In this paper, we have presented a simple end-to-end model that employs a Multi-Modal Reasoning graph to encounter semantic and positional relationships between text and salient visual regions. The learned space is composed of enriched features obtained from nodes in a graph, module that acts as an appropriate reasoning scheme. Exhaustive experiments in two datasets and two different tasks validate the robustness of the presented model which achieves state-of-the-art results by a significant margin over previous methods. Moreover, our end-to-end pipeline does not require pre-computed handcrafted features or a collection of ensemble models as earlier works. In the future, we expect to explore the effectiveness of this approach in other vision and language-related tasks.

## References

- [1] Jon Almazán, Albert Gordo, Alicia Fornés, and Ernest Valveny. Word spotting and recognition with embedded at-



- tributes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(12):2552–2566, 2014.
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.
- [3] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [4] Xiang Bai, Mingkun Yang, Pengyuan Lyu, Yongchao Xu, and Jiebo Luo. Integrating scene text and visual appearance for fine-grained image classification. *IEEE Access*, 6:66322–66335, 2018.
- [5] Hedi Ben-Younes, Rémi Cadene, Nicolas Thome, and Matthieu Cord. Block: Bilinear superdiagonal fusion for visual question answering and visual relationship detection. *arXiv preprint arXiv:1902.00038*, 2019.
- [6] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- [7] Fedor Borisjuk, Albert Gordo, and Viswanath Sivakumar. Rosetta: Large scale system for text detection and recognition in images. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 71–79, 2018.
- [8] Michal Buřta, Yash Patel, and Jiri Matas. E2e-mlt-an unconstrained end-to-end method for multi-language scene text. In *Asian Conference on Computer Vision*, pages 127–143. Springer, 2018.
- [9] Charles Chen, Ruiyi Zhang, Eunye Koh, Sungchul Kim, Scott Cohen, Tong Yu, Ryan Rossi, and Razvan Bunescu. Figure captioning with reasoning and sequence-level training. *arXiv preprint arXiv:1906.02850*, 2019.
- [10] Xiaoxue Chen, Lianwen Jin, Yuanzhi Zhu, Canjie Luo, and Tianwei Wang. Text recognition in the wild: A survey. *arXiv preprint arXiv:2005.03492*, 2020.
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009.
- [12] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016.
- [13] Difei Gao, Ke Li, Ruiping Wang, Shiguang Shan, and Xilin Chen. Multi-modal graph neural network for joint reasoning on vision and scene text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12746–12756, 2020.
- [14] Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. Compact bilinear pooling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 317–326, 2016.
- [15] ZongYuan Ge, Chris McCool, Conrad Sanderson, Peng Wang, Lingqiao Liu, Ian Reid, and Peter Corke. Exploiting temporal information for DCNN-based fine-grained object classification. In *International Conference on Digital Image Computing: Techniques and Applications*, 2016.
- [16] Lluís Gomez, Andres Mafla, Marçal Rusinol, and Dimosthenis Karatzas. Single shot scene text retrieval. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [17] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, pages 369–376. ACM, 2006.
- [18] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in natural images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2315–2324, 2016.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [20] Tong He, Zhi Tian, Weilin Huang, Chunhua Shen, Yu Qiao, and Changming Sun. An end-to-end textspotter with explicit alignment and attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5020–5029, 2018.
- [21] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision*, 116(1):1–20, 2016.
- [22] Sezer Karaoglu, Ran Tao, Theo Gevers, and Arnold WM Smeulders. Words matter: Scene text for image classification and retrieval. *IEEE Transactions on Multimedia*, 19(5):1063–1076, 2017.
- [23] Sezer Karaoglu, Jan C van Gemert, and Theo Gevers. Context: text detection using background connectivity for fine-grained object classification. In *Proceedings of the 21st ACM international conference on Multimedia*, pages 757–760. ACM, 2013.
- [24] Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. ICDAR 2015 competition on robust reading. In *Proc. of the IEEE International Conference on Document Analysis and Recognition*, pages 1156–1160, 2015.
- [25] Vahid Kazemi and Ali Elqursh. Show, ask, attend, and answer: A strong baseline for visual question answering. *arXiv preprint arXiv:1704.03162*, 2017.
- [26] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for fine-grained image categorization: Stanford dogs. In *Proc. CVPR Workshop on Fine-Grained Visual Categorization (FGVC)*, volume 2, page 1, 2011.
- [27] Jin-Hwa Kim, Kyoung-Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang. Hadamard

- product for low-rank bilinear pooling. *arXiv preprint arXiv:1610.04325*, 2016.
- [28] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [29] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017.
- [30] Kunpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. Visual semantic reasoning for image-text matching. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4654–4662, 2019.
- [31] Xiangyang Li and Shuqiang Jiang. Know more say less: Image captioning based on scene graphs. *IEEE Transactions on Multimedia*, 21(8):2117–2130, 2019.
- [32] Minghui Liao, Baoguang Shi, and Xiang Bai. Textboxes++: A single-shot oriented scene text detector. *IEEE Transactions on Image Processing*, 27(8):3676–3690, 2018.
- [33] Minghui Liao, Baoguang Shi, Xiang Bai, Xinggang Wang, and Wenyu Liu. Textboxes: A fast text detector with a single deep neural network. In *AAAI*, pages 4161–4167, 2017.
- [34] Chunxiao Liu, Zhendong Mao, Tianzhu Zhang, Hongtao Xie, Bin Wang, and Yongdong Zhang. Graph structured network for image-text matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10921–10930, 2020.
- [35] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019.
- [36] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [37] Xuebo Liu, Ding Liang, Shi Yan, Dagui Chen, Yu Qiao, and Junjie Yan. Fots: Fast oriented text spotting with a unified network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5676–5685, 2018.
- [38] Andres Mafla, Sounak Dey, Ali Furkan Biten, Lluís Gomez, and Dimosthenis Karatzas. Fine-grained image classification and retrieval by combining visual and locally pooled textual features. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 2950–2959, 2020.
- [39] Andrés Mafla, Rubèn Tito, Sounak Dey, Lluís Gómez, Marçal Rusiñol, Ernest Valveny, and Dimosthenis Karatzas. Real-time lexicon-free scene text retrieval. *Pattern Recognition*, page 107656, 2020.
- [40] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- [41] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [42] Yair Movshovitz-Attias, Qian Yu, Martin C Stumpe, Vinay Shet, Sacha Arnoud, and Liron Yatziv. Ontological supervision for fine grained classification of street view storefronts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1693–1702, 2015.
- [43] Medhini Narasimhan, Svetlana Lazebnik, and Alexander Schwing. Out of the box: Reasoning with graph convolution nets for factual visual question answering. In *Advances in neural information processing systems*, pages 2654–2665, 2018.
- [44] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [45] Florent Perronnin and Christopher Dance. Fisher kernels on visual vocabularies for image categorization. In *2007 IEEE conference on computer vision and pattern recognition*, pages 1–8. IEEE, 2007.
- [46] Joseph Redmon and Ali Farhadi. YOLO9000: better, faster, stronger. *arXiv preprint arXiv:1612.08242*, 2016.
- [47] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [48] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*, pages 593–607. Springer, 2018.
- [49] Ajeet Kumar Singh, Anand Mishra, Shashank Shekhar, and Anirban Chakraborty. From strings to things: Knowledge-enabled vqa model that can read and reason. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4602–4612, 2019.
- [50] Sebastian Sudholt, Neha Gurjar, and Gernot A Fink. Learning deep representations for word spotting under weak supervision. *arXiv preprint arXiv:1712.00250*, 2017.
- [51] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [52] Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. Coco-text: Dataset and benchmark for text detection and recognition in natural images. *arXiv preprint arXiv:1601.07140*, 2016.
- [53] Xiu-Shen Wei, Jianxin Wu, and Quan Cui. Deep learning for fine-grained image analysis: A survey. *arXiv preprint arXiv:1907.03069*, 2019.
- [54] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015.
- [55] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10685–10694, 2019.

- [56] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image captioning with semantic attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4651–4659, 2016.
- [57] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6720–6731, 2019.
- [58] Michael R Zhang, James Lucas, Geoffrey Hinton, and Jimmy Ba. Lookahead optimizer: k steps forward, 1 step back. *arXiv preprint arXiv:1907.08610*, 2019.