# Dynamic Multimodal Instance Segmentation Guided by Natural Language Queries

Edgar Margffoy-Tuay, Juan C. Pérez, Emilio Botero, and Pablo Arbeláez

Universidad de los Andes, Colombia
{ea.margffoy10, jc.perez13, e.botero10, pa.arbelaez}@uniandes.edu.co

**Abstract.** We address the problem of segmenting an object given a natural language expression that describes it. Current techniques tackle this task by either (*i*) directly or recursively merging linguistic and visual information in the channel dimension and then performing convolutions; or by (*ii*) mapping the expression to a space in which it can be thought of as a filter, whose response is directly related to the presence of the object at a given spatial coordinate in the image, so that a convolution can be applied to look for the object. We propose a novel method that integrates these two insights in order to fully exploit the recursive nature of language. Additionally, during the upsampling process, we take advantage of the intermediate information generated when downsampling the image, so that detailed segmentations can be obtained. We compare our method against the state-of-the-art approaches in four standard datasets, in which it surpasses all previous methods in six of eight of the splits for this task.

**Keywords:** Referring expressions, instance segmentation, multimodal interaction, dynamic convolutional filters, natural language processing.

## 1 Introduction

Consider the task of retrieving specific object instances from an image based on natural language descriptions, as illustrated in Fig. 1. In contrast to traditional instance segmentation, in which the goal is to label all pixels belonging to instances in the image for a set of predefined semantic classes [1][2], segmenting instances described by a natural language expression is a task that humans are able to perform without specifically focusing on a limited set of categories: we simply associate a referring expression such as "Man on the right" with what we see, as shown in Fig. 1. To learn such an association is the main goal of this paper.

In this task, the main labels to be assigned are *related to query* and *background.* Thus, the set of possible segmentation masks has few constraints, as a mask can be anything one might observe in the image, in all the ways natural language allows an object to be referred to. An algorithm to tackle this problem must then make sense of the query and relate it to what it sees and recognizes in the image, to finally output an instance segmentation map. Therefore, attempting to naively use Convolutional Neural Networks (CNNs) for this task
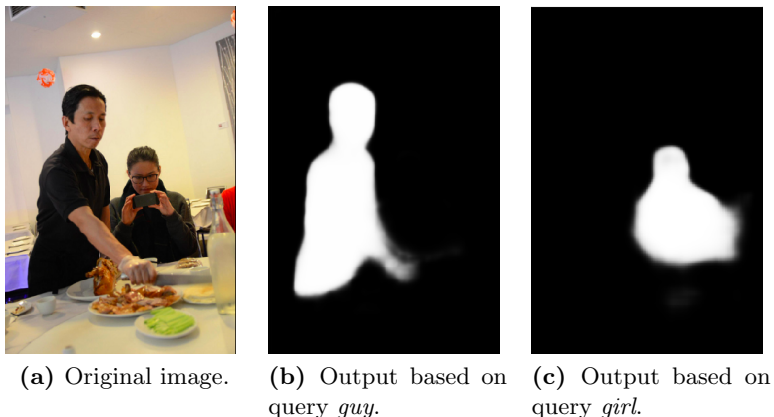
**(a)** Original image.    **(b)** Output based on query *guy*.    **(c)** Output based on query *girl*.

**Fig. 1:** Example of *segmentation based on a natural language expression*. A single mask is the output, in which the only two labels are *member of query* and *background*. Here, we show the raw output of our system, which is the pixelwise probability of belonging to the referred object instance.

falls short, since such networks do not model sequential information by nature, as is required when processing natural language. Given that the cornerstone of this task is the proper *combination* of information retrieved from multiple, dissimilar domains, we expect traditional architectures, like CNNs and Recurrent Neural Networks (RNNs), to be useful modules, but we still need to design an overall architecture that fully exploits their complementary nature.

In this paper, we introduce a modular neural network architecture that divides the task into several sub-tasks, each handling a different type of information in a specific manner. Our approach is similar to [3], [4] and [5] in that we extract visual and natural language information in an independent manner by employing networks commonly used for these types of data, *i.e.*, CNNs and RNNs, and then focus on processing this multi-domain information by means of another neural network, yielding an end-to-end trainable architecture. However, our method also introduces the usage of Simple Recurrent Units (SRUs) for efficient segmentation based on referring expressions, a Synthesis Module that processes the linguistic and visual information jointly, and an Upsampling Module that outputs highly detailed segmentation maps.

Our network, which we refer to as Dynamic Multimodal Network (DMN), is composed of several modules, as depicted in Fig. 2: *(i)* a Visual Module (VM) that produces an adequate representation of the image, *(ii)* a Language Module (LM) that outputs an appropriate representation of the meaning of the query up to a given word, *(iii)* a Synthesis Module (SM) that merges the information provided by the VM and LM at each time step and produces a single output for the whole expression and, finally, *(iv)* an Upsampling Module (UM) that incrementally upsamples the output of the SM by using the feature maps pro-
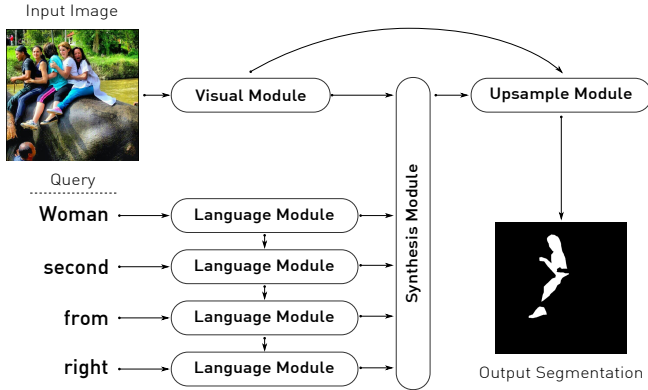
**Fig. 2:** Overview of our Dynamic Multimodal Network (DMN), involving four different modules: Visual Module (VM), Language Module (LM), Synthesis Module (SM), and Upsampling Module (UM).

duced by the VM. Our approach is a fully differentiable, end-to-end trainable neural network for segmentation based on natural language queries. Our main contributions are the following:

- The use of Simple Recurrent Units (SRUs) [6] as language *and* multi-modal processors instead of standard LSTMs [7]. We empirically show that they are efficient while providing high performance for the task at hand.
- A Synthesis Module that takes visual and linguistic information and merges them by generating "scores" for the referring expression in a visual space.
- The Synthesis Module then takes this representation as well as additional features, and exploits the spatial and sequential nature of both types of information to produce a low resolution segmentation map.
- A high resolution upsampling module that takes advantage of visual features during the upsampling procedure in order to recover fine scale details.

We validate our method by performing experiments on all standard datasets, and show that DMN outperforms all the previous methods in various splits for instance segmentation based on referring expressions, and obtains state-of-the art results. Additionally, in order to ensure reproducibility, we provide full implementation of our method and training routines, written in PyTorch[1] [8].

## 2    Related Work

The intersection of Computer Vision (CV) and Natural Language Understanding (NLU) is an active area of research that includes multiple tasks such as object detection based on natural language expressions [9][10], image captioning

---

[1] `https://github.com/BCV-Uniandes/query-objseg`

[11][12][13][14] and visual question answering (VQA) [15][16][17][18][19]. Since visual and linguistic data have properties that make them fundamentally different, *i.e.*, the former has spatial meaning and no sequentiality while the latter does not contemplate space but has a sequential nature, optimally processing both types of information is still an open question. Hence, each work in this sub-field has proposed a particular way of addressing each task.

The task studied in this paper is closest in nature to object detection based on natural language expressions, mirroring how semantic segmentation arose from object detection [20]. Indeed, in [3], object detection with NLU evolved into instance segmentation using referring expressions. We review the state-of-the-art on the task of segmentation based on natural language expressions [3][4][5], highlighting the main contributions in the fusion of multimodal information, and then compare them against our approach.

***Segmentation from Natural Language Expressions [3].*** This work processes visual and natural language information through separate neural networks: a CNN extracts visual features from the image while an LSTM scans the query. Strided convolutions and pooling operations in the CNN downsample the feature maps to a low resolution output while producing large receptive fields for neurons in the final layers. Additionally, to explicitly model spatial information, relative coordinates are concatenated at each spatial location in the feature map obtained by the CNN. Merging of visual and natural language information is done by concatenating the LSTM's output to the visual feature map at each spatial location. Convolution layers with ReLU [21] nonlinearities are applied for final classification. The loss is defined as the average over the per-pixel weighed logistic regression loss. Training has two stages: a low resolution stage, in which the ground truth mask is downsampled to have the same dimensions as the output, and a high resolution stage that trains a deconvolution layer to upsample the low resolution output to yield the final segmentation mask [3]. This seminal method does not fully exploit the sequential nature of language, as it does not make use of the learned word embeddings, it merges visual and linguistic information by concatenation, and it uses deconvolution layers for upsampling, which have been shown to introduce checkerboard artifacts in images [22].

***Recurrent Multimodal Interaction [4].*** This paper argues that segmenting the image based only on a final, memorized representation of the sentence does not fully take advantage of the sequential nature of language. Consequently, the paper proposes to perform segmentation multiple times in the pipeline. The method produces image features at every time step by generating a representation that involves visual, spatial and linguistic features. Such multimodal representation is obtained by concatenating the hidden state of the LSTM that processed the query at every spatial location of the visual representation. The segmentation mask is obtained by applying a multimodal LSTM (mLSTM) to the joint representation and then performing regular convolutions to combine the channels that were produced by the mLSTM. The mLSTM is defined as a convolutional LSTM that shares weights both across spatial location and time step, and is implemented as a $1 \times 1$ convolution that merges all these types of

information. Bilinear upsampling is performed to the network's output at test time to produce a mask with the same dimensions of the ground-truth mask. This method reduces strides of convolutional layers and uses atrous convolution in the final layers of the CNN to compensate for the downsampling. Such modification reduces the upsampling process to bilinear interpolation, but can decrease the CNN's representation capacity while also increasing the number of computations that must be performed by the mLSTM.

*Tracking by Natural Language Specification [5].* In this paper, the main task is object tracking in video sequences. A typical user interaction in tracking consists in providing the bounding box of the object of interest in the first frame. However, this type of interaction has the issue that, for the duration of the video, the appearance and location of objects may change, rendering the initial bounding box useless in some cases. The main idea is to provide an alternative to this approach, by noting that *(i)* the semantic meaning of the object being tracked does not vary for the duration of the video as much as the appearance, and *(ii)* this semantic meaning may be better defined by a linguistic expression. This approach is substantially different from [4] and [3]: visual and linguistic information is never merged *per se*, but rather the linguistic information is mapped to a space in which it can be interpreted as having visual meaning. The visual input is thus processed by a modified VGG [23] to yield a feature map. An LSTM scans the linguistic input, and a single layer perceptron is applied to the LSTM's last hidden state to generate a vector that can be interpreted as being a *filter for a 2D convolution* that is to be performed on the feature map. The *dynamic convolutional visual filter*, generated based on the expression, is computed to produce a strong response to the elements being referred to the expression, and a weak response to those *not* being referred to. This response is interpreted as a "score" for the referring expression, so that a segmentation can be produced. This method proposes a new paradigm for combining information from the visual and linguistic domains, but assumes a non linear combination of the last hidden state is sufficient for modeling a filter that responds to the query.

*Our approach.* The approach of [3] merges multi-domain information by concatenation of linguistic information, subsequent $1 \times 1$ convolutions for segmentation and a deconvolution layer to perform upsampling. The method in [4] follows the same logic as [3] but introduces *recursion* into the approach, exploiting the linguistic information further; however, the upsampling module is an interpolation that produces rather coarse results, to which the authors apply a post-processing DenseCRF, making the architecture not trainable end-to-end. Finally, [5] has a different approach, in which linguistic information is never merged with feature maps, but is rather transformed so that it can detect the locations in the image to which the referring expression has strong response; nonetheless, like [3], it does not fully exploit linguistic information in a sequential manner. Moreover, all these methods fail to utilize information acquired in the downsampling process in the upsampling process.
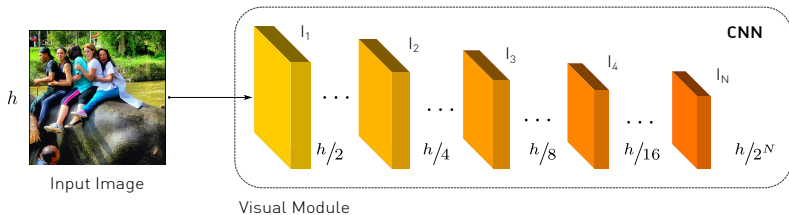
**Fig. 3:** The Visual Module outputs feature maps at $N$ different scales with the aim of using them in the segmentation process and in the upsampling.

Our approach takes advantage of the previous insights, and consists of a modularized network that exploits both the possibility of segmentation based on combinations of multi-domain information, and the feasibility of producing filters that respond to objects being referred to by processing the linguistic information. Following the spirit of [24][25][26], we use skip connections between the downsampling process and the upsampling module to output finely-defined segmentations. We employ the concatenation strategy of [3] but include richer visual and language features. Furthermore, we make use of dynamic filter computation, like [5], but in a sequential manner. Lastly, we introduce the use of a more efficient alternative to LSTMs in this domain, namely SRUs. We demonstrate empirically that SRUs can be used for modeling language *and* multimodal information for this task, and that they can be up to $3\times$ faster than LSTMs, allowing us to train more expressive models.

## 3   Dynamic Multimodal Network

### 3.1   Overall Architecture

Fig. 2 illustrates our overall architecture. Given an input consisting of an image $I$, and a query composed of $T$ words, $\{w_t\}_{t=1}^T$, the Visual Module (VM) takes $I$ as input and produces feature maps at $N$ different scales: $\{I_n\}_{n=1}^N$. The Language Module (LM) processes $\{w_t\}_{t=1}^T$ and yields a set of features $\{r_t\}_{t=1}^T$ and a set of dynamic filters $\{\{f_{k,t}\}_{k=1}^K\}_{t=1}^T$. Given the VM's last output, $I_N$, $\{r_t\}_{t=1}^T$, and $\{\{f_{k,t}\}_{k=1}^K\}_{t=1}^T$, the Synthesis Module (SM) processes this information and produces a single feature map for the entire referring expression. This output, along with the feature maps given by the VM, is processed by the Upsampling Module (UM), that outputs a heatmap with a single channel, to which a sigmoid activation function is applied in order to produce the final prediction.

### 3.2   Visual Module (VM)

Fig. 3 depicts the Visual Module. We extract deep visual features from the image using as backbone a Dual Path Network 92 (DPN92) [27], which has shown competitive performance in various tasks, and is efficient in parameter usage. The VM can be written as a function returning a tuple:
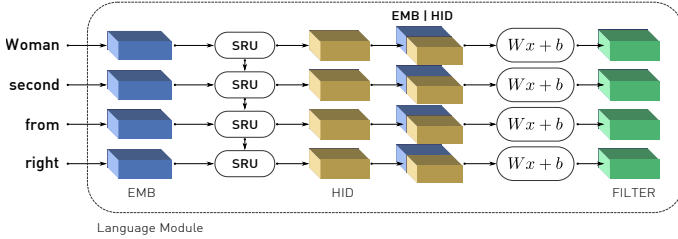
**Fig. 4:** The Language Module uses an SRU, instead of the traditional LSTM, to output enriched features of the query and dynamic filters based on such features.

$$\{I_n\}_{n=1}^{N} = \text{VM}(I) \tag{1}$$

Where $I$ is the original image, and $I_n, n \in \{1, \dots, N\}$ are the downsampled feature maps of dimensions equal to $\frac{1}{2^n}$ of the dimensions of $I$. In the experiments, we use $N = 5$, which considers all convolutional layers in the visual encoder. Note that, since our architecture is fully convolutional, we are not restricted to a fixed image size.

### 3.3  Language Module (LM)

Fig. 4 shows a diagram of the Language Module. Given an expression consisting of $T$ words $\{w_t\}_{t=1}^{T}$, each word is represented by an embedding (WE), $e_t = WE(w_t)$ ($EMB$ in Fig. 4), and the sentence is scanned by an RNN to produce a hidden state $h_t$ for each word ($HID$ in Fig. 4). Instead of using LSTMs as recurrent cells, we employ SRUs [6], which allow the LM to process the natural language queries more efficiently than when using LSTMs. The SRU is defined by:

$$\tilde{x}_t = Wx_t \tag{2}$$
$$f'_t = \sigma\left(W_f x_t + b_f\right) \tag{3}$$
$$r_t = \sigma\left(W_r x_t + b_r\right) \tag{4}$$
$$c_t = f'_t \odot c_{t-1} + (1 - f'_t) \odot \tilde{x}_t \tag{5}$$
$$h_t = r_t \odot g(c_t) + (1 - r_t) \odot x_t \tag{6}$$

Where $\odot$ is the element-wise multiplication. The function $g(\cdot)$ can be selected based on the task; here we choose $g(\cdot)$ to be the sigmoid function. For further details regarding the SRU definition and implementation, please refer to [6].

We concatenate the hidden state $h_t$ with the word embedding $e_t$ to produce the final language output: $r_t = [e_t, h_t]$. This procedure yields an enriched language representation of the concept of the sentence up to word $t$. Moreover, we compute a set of dynamic filters $f_{k,t}$ based on $r_t$, defined by:

$$f_{k,t} = \sigma\left(W_{f_k} r_t + b_{f_k}\right), k = 1, \dots, K \tag{7}$$
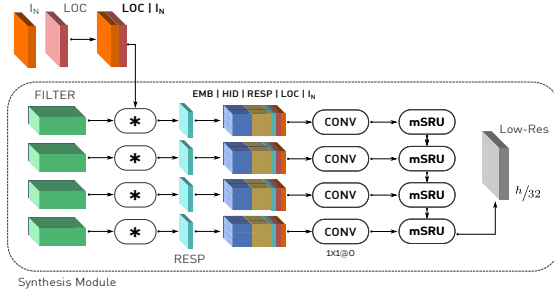
**Fig. 5:** The Synthesis Module takes into account the response to dynamic filters, language features, spatial coordinates representation, and visual features in a recurrent manner to output a single response map.

thus, we define the LM formally as:

$$\left(\{r_t\}_{t=1}^T, \{\{f_{k,t}\}_{k=1}^K\}_{t=1}^T\right) = \text{LM}\left(\{w_t\}_{t=1}^T\right) \tag{8}$$

### 3.4 Synthesis Module (SM)

Fig. 5 illustrates the Synthesis Module. The SM is the core of our architecture, as it is responsible for merging multimodal information. We first concatenate $I_N$ and a representation of the spatial coordinates ($LOC$ in Fig. 5), following the implementation of [3], and convolve this result with each of the filters computed by the LM to generate a response map ($RESP$ in Fig. 5) consisting of $K$ channels: $F_t = \{f_{k,t} * I_N\}_{k=1}^K$. Next, we concatenate $I_N$, $LOC$, and $F_t$ along the channel dimension to obtain a representation $I'$, to which $r_t$ is concatenated *at each spatial location*, as to have all the multimodal information in a single tensor. Finally, we apply a $1 \times 1$ convolutional layer that merges all the multimodal information, providing tan output corresponding to each time step $t$, denoted by $M_t$. Formally, $M_t$ is defined by:

$$M_t = \text{Conv}_{1\times1}([I_N, F_t, LOC, r_t]) \tag{9}$$

Next, in the pursuit of performing a recurrent operation that takes into account the sequentiality of the set and also the information of each of the channels in $M_t$, we propose the use of a multimodal SRU (mSRU), which we define as a $1 \times 1$ convolution, similar to [4] but using SRUs. We apply the mSRU to the whole set $\{M_t\}_{t=1}^T$, so that all the information in each $M_t$, including the sequentiality of the set, is used in the segmentation process. The final hidden states are gathered to produce a 3D tensor that is interpreted as a feature map. This tensor, which we denominate $R_N$, due to its size being $\frac{1}{2^N}$ of the image's original size, has the same dimensions as $M_t$ and has as many channels as there are entries in the hidden state of the mSRU. We define the SM as a function returning $R_N$:
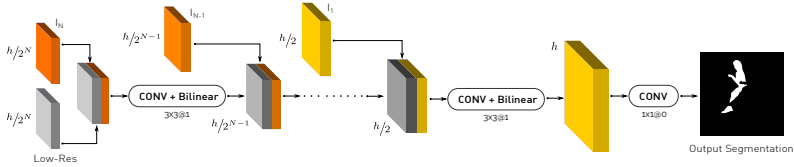
**Fig. 6:** The Upsampling Module makes use of all the feature maps that were generated in the feature extraction process to provide more detailed segmentations.

$$R_N = \text{SM}\left(\{M_t\}_{t=1}^T\right) = \text{mSRU}\left(\{M_t\}_{t=1}^T\right), \tag{10}$$

where $M_t$ is reshaped appropriately to make sense of the sequential nature of the information at each time step.

### 3.5   Upsampling Module (UM)

Finally, the Upsampling Module is shown in Fig. 6. Inspired by skip connections [24][25] [28], we construct an upsampling architecture that takes into account the feature maps $\{I_n\}_{n=1}^N$ at all stages in order to recover fine-scale details. At each stage we concatenate $R_n$ with $I_n$, perform $3 \times 3$ convolution over this result, and then scale the size by a factor of 2 via bilinear interpolation to generate $R_{n-1}$. We apply this process $\log_2(N)$ times, to produce an output mask of the same size of the input $R_1$. We apply $1 \times 1$ convolution over $R_1$ to generate a single channel and, finally, a sigmoid layer to obtain scores between 0 and 1.

## 4   Experimental Setup

### 4.1   Datasets

We conduct experiments on the four standard datasets for this task: ReferIt, UNC, UNC+ [29], and GRef [30]. UNC, UNC+ and GRef are based on MS COCO [1]. The type of objects that appear in the referring expressions, length of the expressions, and relative size of the referred objects are the main differences between the datasets. The high variability of those characteristics across the datasets evidences the challenge of constructing models for this task that are capable of generalization.

**ReferIt** [29] is a crowd-sourced database that contains images and referring expressions to objects in those images. Currently it has 130,525 expressions, referring to 96,654 distinct objects, in 19,894 photographs of natural scenes.

**UNC** [31], was collected interactively in the ReferIt game, with images that were selected to contain two or more objects of the same object category [31], which means that an expression making reference to a determined type of object will need to be further analysed to determine *which* object the query is referring to, since ambiguity arises when only guided by semantic instance class cues. It consists of 142,209 referring expressions for 50,000 objects in 19,994 images.

**UNC+** [31], is similar to UNC but has an additional restriction regarding words describing location: expressions must be based only on appearance rather than location. Such restriction implies that the expression will depend on the perspective of the scene and the semantic class of the object.

**GRef** [30], was collected on Amazon's Mechanical Turk and contains 85,474 referring expressions for 54,822 objects in 26,711 images selected to contain between two and four objects of the same class, and thus, it presents similar challenges to those of UNC.

## 4.2   Performance Metrics

We use the standard metrics from the literature to allow for direct comparison with respect to the state-of-the-art. We perform experiments with the proposed method on the four standard datasets described above by training on the training set and evaluating the performance in each of the validation or test sets. We evaluate results by using two standard metrics: *(i)* mean Intersection over Union (mIoU), defined as the total intersection area between the output and the Ground Truth (GT) mask, divided by the total union area between the output and the GT mask, added over all the images in the evaluation set, and *(ii) Precision@X*, or *Pr@X*, ($X \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$), defined as the percentage of images with IoU higher than $X$. We report mIoU in the validation and test splits of each dataset, when available, using optimal thresholds from the training or validation splits, respectively.

## 4.3   Implementation Details

All the models are defined and trained with DPN92 [27] as the backbone, which outputs 2688 channels in the last layer. We use $N = 5$ scales in the VM. We use the following hyperparameters, which we optimized on the UNC+ *val* set: WE size of 1000, 2-layered SRU with hidden state size of 1000, $K = 10$ filters, 1000 $1 \times 1$ convolution filters in the SM, 3-layered mSRU's hidden state size of 1000. The increased number of layers presented here in both the SRU and the mSRU, with respect to usual number of layers in LSTMs, are in response to the increased need of layers for an SRU to work as expected, according to [6]. We train our method in two stages: at low resolution (*i.e.*, without using the UM) and then finetune the UM to obtain high resolution segmentation maps.

Training is done with Adam optimizer [32] with an initial learning rate of $1 \times 10^{-5}$, a scheduler that waits 2 epochs for loss stagnation to reduce the learning rate by a factor of 10, and batch size of 1 image-query pair.

# 5   Results

## 5.1   Control Experiments

We assess the relative importance of our modules in the final result by performing ablation experiments. The control experiments were trained until convergence on

the UNC dataset. Accordingly, we compare them to a version of our full method trained for a similar amount of time. Table 1 presents the results.

The "Only VM" experiment in row 1 consists on training only the VM and upsampling the low resolution output with bilinear interpolation, without using the query. At test time, the VM processes the image and the resulting segmentation map is upsampled using the UM and compared to the GT mask. Results show how this method performs poorly in comparison to our full approach, which confirms our hypothesis that naively using a CNN falls short for the task addressed in this paper. However, it is interesting that performance is rather high for a method that does not use linguistic information at all. This result reveals that many of the objects annotated in this dataset are salient, and so the network is able to learn to segment salient objects without help from the query.

The experiment in row 2 consists of defining $r_t = h_t$, instead of using the concatenation of $h_t$ and $e_t$, which affects both the LM (when computing the dynamic filters) and the SM. Results show that using the learned embeddings provides a small gain in the full method, particularly for stricter overlap thresholds.

Next, in row 3 we assess the importance of skip connections in the UM, which is a measure of the usefulness of features extracted in the downsampling process for the upsampling module. The large drop in performance with respect to the full method shows that the skip connections allow the network to exploit finer details that are otherwise lost, showing how the upsampling strategy can benefit from performing convolutions followed by bilinear interpolations instead of deconvolutions, as done in [3].

We next study the effects of removing features from $M_t$. In rows 4 and 5 we remove the set of responses to the dynamic filters $F$, as well as the concatenation of $r_t$ in the SM, respectively. We observe that the dynamic filters generate useful scores for the natural language queries in the visual space, and that reusing features from the LM in the SM does not help the network significantly.

Our results show that the key components of our network have significant impact in overall performance. High performance is not achieved by either using only linguistic information ($r_t$) or the response to filters ($F$): both must be properly combined. Additionally, the UM allows the network to properly exploit features from the downsampling stage and perform detailed segmentation.

Table 1: Precision@$X$ and mIoU for ablation study in the UNC testA split.

| Method | Pr@0.5 | Pr@0.6 | Pr@0.7 | Pr@0.8 | Pr@0.9 | mIoU |
|---|---|---|---|---|---|---|
| Only VM | 15.26 | 6.36 | 2.96 | 0.91 | 0.14 | 30.92 |
| Only $h_t$ in LM and SM | 65.38 | **57.99** | **47.07** | 27.38 | 4.63 | 54.80 |
| No skip connections in UM | 56.58 | 42.77 | 26.32 | 9.22 | 1.07 | 49.26 |
| No dynamic filters | 57.53 | 48.70 | 38.27 | 20.64 | 3.00 | 50.34 |
| No concatenation of $r_t$ | 64.52 | 56.69 | 45.16 | 25.56 | 4.38 | 54.69 |
| DMN | **65.83** | 57.82 | 46.80 | **27.64** | **5.12** | **54.83** |

**Table 2:** Comparison with the state-of-the-art in mIoU performance across the different datasets. Blank entries where authors do not report performance.

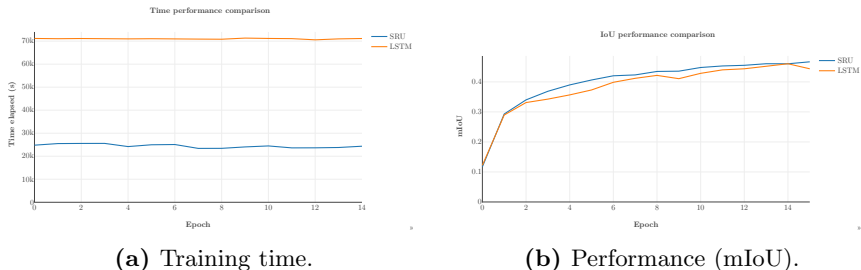| Method | Referit | GRef | UNC | | | UNC+ | | |
|---|---|---|---|---|---|---|---|---|
| | test | val | val | testA | testB | val | testA | testB |
| [3] | 48.03 | 28.14 | - | - | - | - | - | - |
| [33] | 49.91 | 34.06 | - | - | - | - | - | - |
| [5] | 54.30 | - | - | - | - | - | - | - |
| [4] | **58.73** | 34.52 | 45.18 | 45.69 | **45.57** | 29.86 | 30.48 | 29.50 |
| DMN | 52.81 | **36.76** | **49.78** | **54.83** | 45.13 | **38.88** | **44.22** | **32.29** |



(a) Training time.



(b) Performance (mIoU).

**Fig. 7:** SRUs *vs.* LSTMs comparison on UNC *testA* (at low resolution).

## 5.2  Comparison with the State-of-the-Art

Next, we proceed to compare our full method with the state-of-the-art, for which we evaluate on all the datasets described above. Table 2 compares the mIoU of our method with the state-of-the-art in this task [3][4][5]. The results show that our method outperforms all other methods in six out of eight splits of the datasets. By including enriched linguistic features at several stages of the process, and by combining them in different ways, our network learns appropriate associations between queries and the instances they refer to.

Interestingly, the performance gain in the *testB* splits of UNC and UNC+ is not as large as in *testA*. One possible reason for the smaller performance gain across splits is their difference: visual inspection of results shows how *testA* splits are biased towards queries related to segmenting persons. The *testB* splits, however, contain more varied queries and objects, which is why the increase in mIoU is not as marked. This behavior can also be observed for the method proposed by [4], as shown in the second-to-last line of Table 2.

## 5.3  Efficiency Comparison: SRU *vs.* LSTM

In order to assess the efficiency and the performance of SRUs when compared to the more commonly used LSTMs, both as language *and* multi-modal processors, we conduct an experiment in which we replaced the SRUs with LSTMs in our final system, both in the LM and the SM, we trained on the UNC dataset,
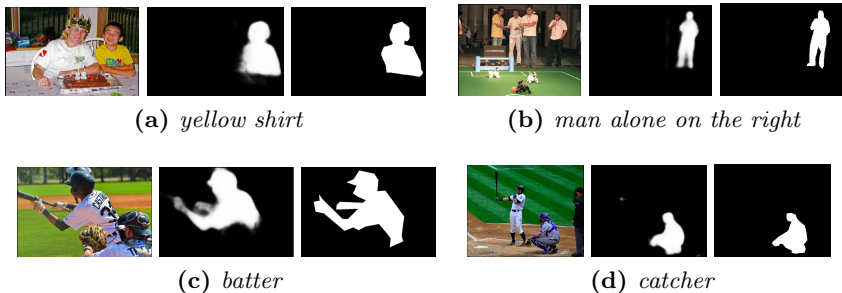
**(a)** *yellow shirt*        **(b)** *man alone on the right*

**(c)** *batter*        **(d)** *catcher*

**Fig. 8:** Qualitative examples of the output of the network. From left to right in each subfigure: original image, heatmap produced by our method, and ground-truth mask. Each caption is the query that produced the output.

and we measured performance on the *testA* split. In terms of model complexity, when using SRUs, the LM and the SM have 9M and 10M trainable parameters, respectively. When switching to LSTMs, the number of parameters increases to 24M and 24.2M, respectively, multiplying training time by a factor of three, as shown in Fig. 7*a*. Regarding accuracy, Fig. 7*b*. shows that both systems perform similarly, with a small advantage for SRUs. Therefore, when compared to LSTMs, SRUs allow us to design architectures that are more compact, train significantly faster, and generalize better.

### 5.4   Qualitative Results

Fig. 8 shows qualitative results in which the network performed well. These examples demonstrate DMN's flexibility for segmenting based on different information about a particular class or instance: attributes, location or role. We emphasize that understanding a role is not trivial, as it is related to the object's context and appearance. Additionally, a semantic difficulty that our network seems to overcome is that the role coexists with the object's class: an instance can be "batter" and be "person". Notice in Fig. 8 that thanks to the upsampling module, our network segments fine details such as legs, heads and hands. In Fig. 8*a* the query refers to the kid by one of his *attributes*: the color of his shirt; in Fig. 8*b* the man is defined by his *location* and the fact that he is alone (although that could be dropped, as there is no ambiguity); in Figs. 8*c* and *d* the reference is based on the person's *role*.

Typical failure cases are depicted in Fig. 9. In Fig. 9*a* the network segments (arguably) the incorrect person, since the correct segmentation was the person at the border of the image whose face is partially shown. Several failure cases we found had exactly the same issue: ambiguity in the expression that could confuse even a human. Fig. 9*b* shows an example of strong failure, in which a weak segmentation is produced. The model appears to have only focused on the word "right". We attribute this failure to the network's inability to make

**(a)** *person on left*

**(b)** *person sitting on the right with a hat that has a white stripe*

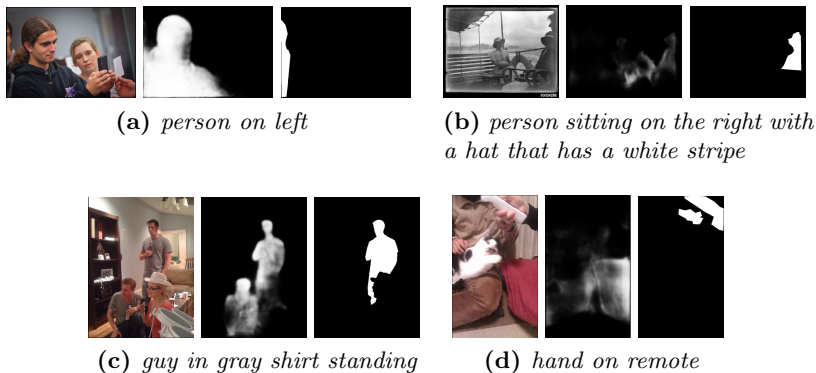**(c)** *guy in gray shirt standing*

**(d)** *hand on remote*

**Fig. 9:** Negative examples of the output of the network. From left to right in each Subfigure: original image, heatmap produced by our method, and ground-truth mask. Each caption is the query that produced the output.

sense of such a relatively long sentence, which, while unambiguously defining an object, is a confusing way of referring to it. Fig. 9*c* is an interesting example of the network's confusion. While the woman is not segmented, two subjects that share several attributes (*guy*, *gray* and *shirt*) are confused and are *both* segmented. However, the network does not manage to use the word "standing" to resolve the ambiguity. Finally, in Fig. 9*d* a failure is observed, where nothing related to the query is segmented. The mask that is produced only reflects a weak attempt of segmenting the red object, while ignoring the upper part of the image, in which both the hand *and* the remote were present.

## 6    Conclusions

We propose Dynamic Multimodal Network, a novel approach for segmentation of instances based on natural language expressions. DMN integrates insights from previous works into a modularized network, in which each module has the responsibility of handling information from a specific domain. Our Synthesis Module combines the outputs from previous modules and handles this multi-modal information to produce features that can be used by the Upsampling Module. Thanks to the incremental use of feature maps obtained in the encoding part of the network, the Upsampling Module delivers great detail in the final segmentations. Our method outperforms the state-of-the-art methods in six of the eight standard dataset splits for this task.

# References

1. Lin, T.Y., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: European Conference on Computer Vision (ECCV). (2014)

2. Hariharan, B., Arbeláez, P., Girshick, R., Malik, J.: Simultaneous detection and segmentation. In Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T., eds.: Computer Vision – ECCV 2014, Cham, Springer International Publishing (2014) 297–312

3. Hu, R., Rohrbach, M., Darrell, T.: Segmentation from natural language expressions. European Conference on Computer Vision (ECCV) (2016)

4. Liu, C., Lin, Z.L., Shen, X., Yang, J., Lu, X., Yuille, A.L.: Recurrent multimodal interaction for referring image segmentation. International Conference on Computer Vision (ICCV) (2017) 1280–1289

5. Li, Z., Tao, R., Gavves, E., Snoek, C.G.M., Smeulders, A.W.M.: Tracking by natural language specification. Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

6. Lei, T., Zhang, Y., Artzi, Y.: Training RNNs as fast as CNNs. CoRR **abs/1709.02755** (2017)

7. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8) (1997) 1735–1780

8. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: Conference on Neural Information Processing Systems (NIPS). (2017)

9. Hu, R., Xu, H., Rohrbach, M., Feng, J., Saenko, K., Darrell, T.: Natural language object retrieval. Conference on Computer Vision and Pattern Recognition (CVPR) (2016) 4555–4564

10. Guadarrama, S., Rodner, E., Saenko, K., Darrell, T.: Understanding object descriptions in robotics by open-vocabulary object retrieval and detection. The International Journal of Robotics Research **35**(1-3) (2016) 265–280

11. Hendricks, L.A., Venugopalan, S., Rohrbach, M., Mooney, R.J., Saenko, K., Darrell, T.: Deep compositional captioning: Describing novel object categories without paired training data. Conference on Computer Vision and Pattern Recognition (CVPR) (2016) 1–10

12. Gan, Z., Gan, C., He, X., Pu, Y., Tran, K., Gao, J., Carin, L., Deng, L.: Semantic compositional networks for visual captioning. Conference on Computer Vision and Pattern Recognition (CVPR) (2017) 1141–1150

13. Johnson, J., Karpathy, A., Fei-Fei, L., Li, C., Li, Y.W., fei Li, F.: Densecap: Fully convolutional localization networks for dense captioning. Conference on Computer Vision and Pattern Recognition (CVPR) (2016) 4565–4574

14. Yang, Z., Yuan, Y., Wu, Y., Cohen, W.W., Salakhutdinov, R.: Review networks for caption generation. In: Conference on Neural Information Processing Systems (NIPS). (2016)

15. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the V in VQA matter: Elevating the role of image understanding in visual question answering. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017) 6325–6334

16. Li, C., Groth, O., Bernstein, M.S., Fei-Fei, L., Li, Y.W., Li, F.F.: Visual7w: Grounded question answering in images. Conference on Computer Vision and Pattern Recognition (CVPR) (2016) 4995–5004

17. Krishna, R., Li, C., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Shamma, D.A., Bernstein, M.S., Fei-Fei, L.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. International Journal of Computer Vision (IJCV) **123** (2016) 32–73
18. Agrawal, A., Lu, J., Antol, S., Mitchell, M., Zitnick, C.L., Parikh, D., Batra, D.: VQA: Visual question answering. International Conference on Computer Vision (ICCV) (2015) 2425–2433
19. Teney, D., Liu, L., van den Hengel, A.: Graph-structured representations for visual question answering. Conference on Computer Vision and Pattern Recognition (CVPR) (2017) 3233–3241
20. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International Journal of Computer Vision **88**(2) (June 2010) 303–338
21. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: Proceedings of the 27th international conference on machine learning (ICML-10). (2010) 807–814
22. Odena, A., Dumoulin, V., Olah, C.: Deconvolution and checkerboard artifacts. Distill (2016)
23. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. International Conference on Learning Representations (ICLR) (2015)
24. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015. (2015) 234–241
25. Shelhamer, E., Long, J., Darrell, T.: Fully convolutional networks for semantic segmentation. Conference on Computer Vision and Pattern Recognition (CVPR) (2015) 3431–3440
26. Hariharan, B., Arbeláez, P., Girshick, R., Malik, J.: Object instance segmentation and fine-grained localization using hypercolumns. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **39**(4) (2017) 627–639
27. Chen, Y., Li, J., Xiao, H., Jin, X., Yan, S., Feng, J.: Dual path networks. In: Conference on Neural Information Processing Systems (NIPS). (2017)
28. Huang, G., Liu, Z., Weinberger, K.Q.: Densely connected convolutional networks. Conference on Computer Vision and Pattern Recognition (CVPR) (2017) 2261–2269
29. Kazemzadeh, S., Ordonez, V., Matten, M., Berg, T.L.: Referit game: Referring to objects in photographs of natural scenes. In: EMNLP. (2014)
30. Mao, J., Huang, J., Toshev, A., Camburu, O., Yuille, A.L., Murphy, K.: Generation and comprehension of unambiguous object descriptions. Conference on Computer Vision and Pattern Recognition (CVPR) (2016) 11–20
31. Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. In: European Conference on Computer Vision (ECCV). (2016)
32. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. International Conference on Learning Representations (ICLR) (2015)
33. Hu, R., Rohrbach, M., Venugopalan, S., Darrell, T.: Utilizing large scale vision and text datasets for image segmentation from referring expressions. European Conference on Computer Vision (ECCV) (2016)