

# Phrase Grounding by Soft-Label Chain Conditional Random Field

Jiacheng Liu    Julia Hockenmaier

University of Illinois at Urbana-Champaign, Urbana, IL, USA 61801

{jl25, juliahmr}@illinois.edu

## Abstract

The phrase grounding task aims to ground each entity mention in a given caption of an image to a corresponding region in that image. Although there are clear dependencies between how different mentions of the same caption should be grounded, previous structured prediction methods that aim to capture such dependencies need to resort to approximate inference or non-differentiable losses. In this paper, we formulate phrase grounding as a sequence labeling task where we treat candidate regions as potential labels, and use neural chain Conditional Random Fields (CRFs) to model dependencies among regions for adjacent mentions. In contrast to standard sequence labeling tasks, the phrase grounding task is defined such that there may be multiple correct candidate regions. To address this multiplicity of gold labels, we define so-called Soft-Label Chain CRFs, and present an algorithm that enables convenient end-to-end training. Our method establishes a new state-of-the-art on phrase grounding on the Flickr30k Entities dataset. Analysis shows that our model benefits both from the entity dependencies captured by the CRF and from the soft-label training regime. Our code is available at [github.com/liujch1998/SoftLabelCCRF](https://github.com/liujch1998/SoftLabelCCRF)

## 1 Introduction

Given an image and a corresponding caption, the phrase grounding task aims to ground each entity mentioned by a noun phrase in the caption to a region in the image. Phrase grounding has attracted much research interest due to its application in downstream tasks including image captioning (Karpathy et al., 2014; Fang et al., 2015; Donahue et al., 2017; Xu et al., 2015), image retrieval (Chen et al., 2017a; Radenovic et al., 2016), and visual question answering (Agrawal et al., 2017; Yu et al., 2017, 2018a).

Cheerleaders at a sporting event toss a girl high up into the air .      Old man sits on rocks while working with his hands .

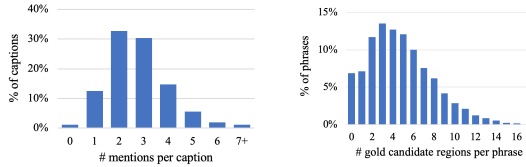


(a) Dependency between entities. The visual relationship between grounding regions for “cheerleaders” and “a girl” should agree with context “toss ... high up into the air”. (b) Gold label multiplicity. The green box is the annotated gold grounding region for entity phrase “Old man”, while the orange dash boxes are region proposals with  $\text{IoU} \geq 0.5$  with gold.

Figure 1: Example image-caption pairs from Flickr30k Entities, illustrating entity dependencies and gold label multiplicity.

Phrase grounding systems typically work by ranking a set of candidate regions (Chen et al., 2017b; Yu et al., 2018b). Region proposals are generated from the image by a vision backbone model, without conditioning on the caption. Features of the phrase to be grounded are extracted, and subsequently interact with features of candidate regions, to determine phrase-region compatibility. Candidate regions are then ranked based on this compatibility metric, and the highest-scored candidate region is selected as the predicted grounding of the phrase.

In Flickr30k Entities (Plummer et al., 2017b), each caption contains an average of 2.76 entity phrases to ground (Figure 2a; phrases with no corresponding gold regions are not counted). It therefore stands to reason that phrases in the same caption should not be grounded independently (to op-



(a) Distribution of number of entity phrases per caption. (b) Distribution of number of gold labels per entity phrase.

Figure 2: Validation set statistics for Flickr30k Entities.

timize each individual phrase-region assignment), but jointly (to optimize the global phrase-region assignment for the entire caption). Figure 1a illustrates this phenomenon. The caption contains a sequence of two entity phrases, “cheerleaders” and “a girl”, and the task is to label each phrase with a candidate region that best grounds it. Since there are several women present in the image, “a girl” has ambiguous grounding by itself, but it can be disambiguated by encouraging the visual relationship between “a girl” and “cheerleaders” to conform with context provided in the caption.

Some works are aware that dependencies between entities in the same caption play an important role in building more accurate phrase grounding systems (Wang et al., 2016; Plummer et al., 2017a; Chen et al., 2017b). The success of these structured prediction methods shows the advantage of considering entity dependencies in learning and prediction. However, these approaches capture certain relations in an *ad hoc* manner, and resort to approximate inference (Wang et al., 2016; Plummer et al., 2017a) or non-differentiable losses (Chen et al., 2017b).

To obtain models and inference algorithms that facilitate more globally consistent phrase grounding predictions, we propose to formulate phrase grounding as a sequence labeling task where we treat candidate regions as potential labels for the phrases in the input sequence. This allows us to build phrase grounding models based on Conditional Random Fields (CRFs) (Lafferty et al., 2001) that capture entity dependencies in a universal and differentiable manner. Our results indicate that systems that capture dependencies between phrases in the same caption in a principled manner outperform systems that ignore these dependencies.

A second problem lies in the use of region proposals, which distinguishes phrase grounding from other sequence labeling tasks where CRFs

are directly applicable. Following the metrics of object detection, in phrase grounding the correctness of a predicted region is judged by its overlap by Intersection-over-Union (IoU) with the gold region (Plummer et al., 2017b). To cover potential regions with high enough IoU, it is common to generate a myriad of region proposals and for these candidate regions to contain or substantially overlap with each other. As a result, there could be more than one candidate region with high IoU with the gold region, and they should all be considered as correct grounding for the phrase. This phenomenon of gold label multiplicity is illustrated in Figure 1b. We hypothesize that it is important to consider gold label multiplicity and identify all correct region proposals during training, since the model would receive contradictory training signals if some correct proposals were marked as incorrect. With region proposals generated by a Bottom-Up Attention (Anderson et al., 2018) visual backbone, in Flickr30k Entities each phrase has an average of 4.75 gold labels, and detailed statistics are presented in Figure 2b. To address this problem, we adopt the soft-label target distribution proposed by Yu et al. (2018b), and our experiments show that models trained with this regime significantly outperform those trained with one-hot target regime.

To combine the benefits brought by structured prediction from CRFs and by soft-label training regime, we define Soft-Label Chain CRFs, a variation of standard chain CRFs that allows us to work with gold label multiplicity. We adapt learning and inference algorithms from chain CRFs and develop an end-to-end training algorithm for our proposed model.

We evaluate the effectiveness of Soft-Label Chain CRF on phrase grounding by conducting experiments on the Flickr30k Entities dataset (Plummer et al., 2017b) and comparing grounding accuracy with strong baseline models, as well as with existing structured prediction methods and current state-of-the-art models. Experimental results show that our Soft-Label Chain CRF model outperforms its hard-label CRF counterpart by 2.43%, a vanilla non-CRF soft-label model by 0.40%, and the previous best results by about 1.4%, demonstrating that both of our contributions, modeling phrase grounding as a sequence labeling task, and training with soft label targets, matter for this task.

## 2 Related Work

**Phrase Grounding.** The phrase grounding task was first postulated by Karpathy and Fei-Fei (2017) and Plummer et al. (2017b), both of which moved from the holistic image captioning to the finer-grained task of matching regions with phrases in the caption. Datasets for this task include Flickr30k Entities (Plummer et al., 2017b), RefCOCO (Yu et al., 2016), and Visual Genome (Krishna et al., 2017). The general framework of proposal-generation-ranking has become adopted by most approaches to phrase grounding, and research in this area has focused on improving specific components of this framework. Our work can be viewed as an improvement to the training and prediction aspects.

**Structured Prediction in Phrase Grounding.** We summarize some works that consider entity dependencies by structured prediction. Structured Matching (Wang et al., 2016) formulates phrase grounding as a bipartite matching process between phrases and candidate regions, and encourages the spatial relationship between two grounding regions to conform to an extracted *partial coreference* relation between their corresponding phrases. The resulting discrete optimization problem is then relaxed into a linear program to enable end-to-end training. Phrase-Region CCA (Plummer et al., 2017a) mines *frequent patterns* of semantically related paired phrases and trains a separate model for each pattern. The addition of this pairwise score makes the optimization a quadratic programming problem that requires approximate inference. QRC Net (Chen et al., 2017b) assumes that phrases in a caption refer to distinct entities, and thus predicted grounding regions are penalized for *spatial overlapping*. However, overlapping regions can be penalized only after prediction, so this loss is not differentiable, and one has to resort to reinforcement learning. In these works, partial coreference extraction, frequent patterns mining and spatial overlap penalties are *ad hoc* entity dependency capturing, while we aim to universally encompass the spectrum of such dependencies.

**Soft-Label Training Regime.** Conventionally, region proposal ranking is done by predicting a probability distribution over all candidate regions for grounding a given entity phrase, which is learned to match a target distribution. Chen et al. (2017b) and Rohrbach et al. (2016) define the tar-

get distribution as a one-hot vector which only gives credit to the candidate region with highest IoU with the gold region, and cross-entropy loss is used as training objective. Under this hard-label training regime, the model is trained to pick only the best candidate region while rejecting all the inferior-than-best candidate regions, which is intuitively not a good behavior. Yu et al. (2018b) proposes a soft-label target distribution which gives weighted credit to all good candidate regions (i.e. those with above-threshold IoU with the gold region), and uses Kullback-Leibler (KL) divergence loss as training objective.

**Conditional Random Fields.** CRFs (Lafferty et al., 2001) are discriminative probabilistic models that have been found useful in sequence labeling tasks by capturing label dependencies (Ma and Hovy, 2016; Lample et al., 2016). We summarize some works relevant to CRFs learned in soft-label or multi-label settings. Multi-CRFs (Dredze et al., 2009) learn CRFs with noisy annotated data, where annotators may disagree on the label for input tokens. The assumption is that there is always only one gold label for each token, so the model favors single label while conforming to the prior distribution of labels set by annotators. To work with soft-label targets, it employs a mode-seeking, exclusive KL divergence definition, which does not imply moment-matching, a desired property of CRFs (and in general, exponential family models) that we show in Section 3.1 and 3.2 for the mean-seeking, inclusive KL divergence definition in our model. Rodrigues et al. (2014) models the latent reliability of individual annotators, and use this information to guide the selection of trustworthy annotation sources and estimation of real gold labels. Note that both works always assume one gold label per input token, where the ambiguity comes from unreliability of annotations, while our work focuses on cases where there may be multiple gold labels per input token by the nature of the task.

## 3 Soft-Label Chain CRF

CRFs model the probability of a label sequence  $\mathbf{y} = y^{1:T}$  conditioned on an input sequence  $\mathbf{x} = x^{1:T}$  in terms of a score function  $s(\mathbf{x}, \mathbf{y})$ :

$$p(\mathbf{y}|\mathbf{x}) = \frac{\exp s(\mathbf{y}, \mathbf{x})}{\sum_{\mathbf{y}'} \exp s(\mathbf{y}', \mathbf{x})}$$

For a given training example  $\{(\mathbf{x}, \mathbf{y})\}$ , the negative log-likelihood loss (i.e. cross-entropy loss

w.r.t. a one-hot target distribution that gives credit to the gold label only) is

$$L = -\log p(\mathbf{y}|\mathbf{x}) = -s(\mathbf{y}, \mathbf{x}) + \log Z(\mathbf{x})$$

where  $Z(\mathbf{x}) = \sum_{\mathbf{y}'} \exp s(\mathbf{y}', \mathbf{x})$ . The gradient of this loss w.r.t. score function is

$$\frac{\partial L}{\partial s(\mathbf{y}', \mathbf{x})} = -\mathbb{I}(\mathbf{y}' = \mathbf{y}) + p(\mathbf{y}'|\mathbf{x})$$

which is known as moment-matching. This allows us to train CRFs with gradient methods and conveniently connect to backpropagation when the score function is modeled by a neural architecture.

### 3.1 Soft-Label CRF

In the standard CRF above, each input  $x^t$  corresponds to a single gold label  $y^t$ . To account for gold label multiplicity in training stage, we replace the sequence of gold labels  $\mathbf{y}$  with a sequence of distributions  $\mathbf{q} = q^{1:T}$  where  $q^t \in \mathbb{R}^K$  is the gold label distribution over all  $K$  possible labels for input  $x^t$ . Note that this distribution should not be interpreted as the confidence of each label being correct; rather, it should be understood as a probabilistic gold label model: if we randomly choose a gold label, how likely is each label to be selected. With independence assumption, the gold probability of an arbitrary label sequence  $\mathbf{y}$  is

$$q(\mathbf{y}|\mathbf{x}) = \prod_t q(y^t|\mathbf{x}) = \prod_t q(y^t|x^t) = \prod_t q_{y^t}^t$$

It is easy to see that  $q(\mathbf{y}|\mathbf{x})$  is a distribution:

$$\sum_{\mathbf{y}} q(\mathbf{y}|\mathbf{x}) = \sum_{\mathbf{y}} \prod_t q_{y^t}^t = \prod_t \sum_{y^t} q_{y^t}^t = 1$$

And our goal is to learn this target distribution.

Since this target distribution is no longer degenerate, we use Kullback-Leibler (KL) divergence to measure the discrepancy between the model and the target distribution. Our training objective is the KL divergence loss (in mean-seeking, inclusive form):

$$L = \sum_{\mathbf{y}} \left\{ q(\mathbf{y}|\mathbf{x}) \log \frac{q(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x})} \right\}$$

which also gives gradients that demonstrate moment-matching:

$$\frac{\partial L}{\partial s(\mathbf{y}', \mathbf{x})} = -q(\mathbf{y}'|\mathbf{x}) + p(\mathbf{y}'|\mathbf{x})$$

Note that if we had defined the KL divergence loss in its mode-seeking, exclusive form  $\sum_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) \log \frac{p(\mathbf{y}|\mathbf{x})}{q(\mathbf{y}|\mathbf{x})}$ , we would have lost this desired moment-matching property.

### 3.2 Factorization of Soft-Label Chain CRF

Learning CRFs of general graphs requires inference in unit of cliques, which is usually computationally intractable. By restricting to local, pairwise potentials, we reduce the model to a first-order linear chain CRF, whose scoring function factorizes as

$$\begin{aligned} s(\mathbf{y}, \mathbf{x}) &= \sum_t s(y^t, y^{t-1}, \mathbf{x}) \\ &= \sum_t \left\{ \tau(y^t, y^{t-1}, \mathbf{x}) + \varepsilon(y^t, \mathbf{x}) \right\} \end{aligned}$$

where  $\tau(\cdot, \cdot, \cdot)$  is the transition score between labels at  $t-1$  and  $t$  that captures the dependency between labels for adjacent input tokens, and  $\varepsilon(\cdot, \cdot)$  is the emission score between label and input at  $t$ .

Combining this factorization with soft-label targets gives the formal definition of Soft-Label Chain CRF. The loss can be written as

$$\begin{aligned} L &= \sum_{\mathbf{y}} \left\{ q(\mathbf{y}|\mathbf{x}) \log \frac{q(\mathbf{y}|\mathbf{x})}{p(\mathbf{y}|\mathbf{x})} \right\} \\ &= \sum_{\mathbf{y}} \left\{ q(\mathbf{y}|\mathbf{x}) [\log q(\mathbf{y}|\mathbf{x}) - s(\mathbf{y}, \mathbf{x}) + \log Z(\mathbf{x})] \right\} \\ &\quad \text{(Expand } p(\mathbf{y}|\mathbf{x}) \text{ by CRF modeling)} \\ &= \sum_{\mathbf{y}} \left\{ q(\mathbf{y}|\mathbf{x}) [\log q(\mathbf{y}|\mathbf{x}) - s(\mathbf{y}, \mathbf{x})] \right\} + \log Z(\mathbf{x}) \\ &\quad \text{(Marginalize } q(\mathbf{y}|\mathbf{x})) \\ &= \sum_t \sum_{y^t} \left\{ q(y^t|\mathbf{x}) \log q(y^t|\mathbf{x}) \right\} \\ &\quad - \sum_{\mathbf{y}} \left\{ q(\mathbf{y}|\mathbf{x}) s(\mathbf{y}, \mathbf{x}) \right\} + \log Z(\mathbf{x}) \\ &\quad \text{(Independence of } q(y^t|\mathbf{x}) \text{ across } t) \\ &= \sum_t \sum_{y^t} \left\{ q(y^t|\mathbf{x}) \log q(y^t|\mathbf{x}) \right\} \\ &\quad - \sum_{\mathbf{y}} \left\{ q(\mathbf{y}|\mathbf{x}) \sum_t s(y^t, y^{t-1}, \mathbf{x}) \right\} + \log Z(\mathbf{x}) \\ &\quad \text{(Factorization of } s(\mathbf{y}, \mathbf{x})) \\ &= \sum_t \sum_{y^t} \left\{ q(y^t|\mathbf{x}) \log q(y^t|\mathbf{x}) \right\} \\ &\quad - \sum_t \sum_{y^t, y^{t-1}} \left\{ s(y^t, y^{t-1}, \mathbf{x}) \sum_{\mathbf{y}|y^t, y^{t-1}} q(\mathbf{y}|\mathbf{x}) \right\} \\ &\quad + \log Z(\mathbf{x}) \\ &\quad \text{(Reorganize sums by } s(y^t, y^{t-1}, \mathbf{x})) \\ &= \sum_t \sum_{y^t} \left\{ q(y^t|\mathbf{x}) \log q(y^t|\mathbf{x}) \right\} \end{aligned}$$

---

**Algorithm 1** Modified forward algorithm to compute the KL divergence loss for Soft-Label Chain CRFs

---

**procedure** SOFTLABELCHAINCRFLOSS( $\mathbf{q}, \varepsilon(y^t, \mathbf{x}), \tau(y^t, y^{t-1}, \mathbf{x})$ )

**for all** label  $y^0$  **do**

$$\alpha_{y^0}^0 \leftarrow 0$$

$$g_{y^0}^0 \leftarrow 0$$

**for**  $t = 1 \dots T$  **do**

**for all** label  $y^t$  **do**

$$\alpha_{y^t}^t \leftarrow \sum_{y^{t-1}} \left\{ \alpha_{y^{t-1}}^{t-1} \exp [\tau(y^t, y^{t-1}, \mathbf{x}) + \varepsilon(y^t, \mathbf{x})] \right\}$$

$$g_{y^t}^t \leftarrow \sum_{y^{t-1}} \left\{ [g_{y^{t-1}}^{t-1} + \tau(y^t, y^{t-1}, \mathbf{x})] q_{y^{t-1}}^{t-1} + [\varepsilon(y^t, \mathbf{x}) - \log q_{y^t}^t] \right\}$$

$$Z \leftarrow \sum_{y^T} \alpha_{y^T}^T$$

$$G \leftarrow \sum_{y^T} g_{y^T}^T q_{y^T}^T$$

$$L \leftarrow -G + \log Z$$

**return**  $L$

---

$$-\sum_t \sum_{y^t, y^{t-1}} \left\{ q(y^t, y^{t-1} | \mathbf{x}) s(y^t, y^{t-1}, \mathbf{x}) \right\} \\ + \log Z(\mathbf{x})$$

which gives moment-matching gradients

$$\frac{\partial L}{\partial s(y^t, y^{t-1}, \mathbf{x})} = -q(y^t, y^{t-1} | \mathbf{x}) + p(y^t, y^{t-1} | \mathbf{x})$$

$$\frac{\partial L}{\partial \tau(y^t, y^{t-1}, \mathbf{x})} = -q(y^t, y^{t-1} | \mathbf{x}) + p(y^t, y^{t-1} | \mathbf{x})$$

$$\frac{\partial L}{\partial \varepsilon(y^t, \mathbf{x})} = -q(y^t | \mathbf{x}) + p(y^t | \mathbf{x})$$

where

$$q(y^t | \mathbf{x}) = q_{y^t}^t \\ q(y^t, y^{t-1} | \mathbf{x}) = q_{y^t}^t q_{y^{t-1}}^{t-1}$$

are the probability of local label(s) marginalized over all possible non-local labels. Smoothing inference  $p(y^t | \mathbf{x})$  and  $p(y^t, y^{t-1} | \mathbf{x})$  can be computed with forward-backward algorithm.

### 3.3 As an Extension of Soft-Label Model

Note that if we omit all transition terms in Soft-Label Chain CRF, the loss reduces to

$$L' = \sum_t \sum_{y^t} \left\{ q(y^t | \mathbf{x}) [-\varepsilon(y^t, \mathbf{x}) + \log q(y^t | \mathbf{x})] \right\} \\ + \log Z(\mathbf{x}) \\ = \sum_t \sum_{y^t} \left\{ q(y^t | \mathbf{x}) \log \frac{q(y^t | \mathbf{x})}{p(y^t | \mathbf{x})} \right\}$$

which is a total factorization over time. This is as if each label is predicted independently using a soft-label training regime, which is exactly

the KL divergence loss proposed by Yu et al. (2018b). Therefore, our Soft-Label Chain CRF can be viewed as an extension of this soft-label discriminative model.

### 3.4 Modified Forward Algorithm

For chain CRFs, computing the loss only requires forward algorithm, while computing the gradients requires a full forward-backward algorithm. It can be proved that backpropagation on the loss gives the same result as running forward-backward. This is a commonly used trick in modern deep learning frameworks to eliminate the need of implementing the backward pass. Algorithm 1 presents a modified forward algorithm that computes the loss for Soft-Label Chain CRF. In Section 1 and 2 of the Supplementary Materials, we prove the correctness of this algorithm, and that its backpropagation is also equivalent to forward-backward.

## 4 Phrase Grounding as Sequence Labeling

### 4.1 Task Formulation

We formulate phrase grounding as a sequence labeling task. Given an image  $I$ , a caption sentence  $[c^1 \dots c^L]$  where  $c^l$  is a word token, and a set of non-overlapping noun phrase spans  $[p^1 \dots p^T]$  where  $p^t = (s^t, e^t)$  denotes that the  $t$ 'th phrase covers tokens  $c^{s^t}$  to  $c^{e^t}$  (inclusive), we generate a set of region proposals  $\{r_1 \dots r_K\}$ , label each phrase with a candidate region, and refine the region by performing a bounding box regression.



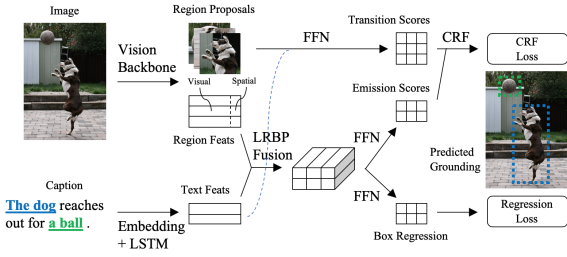


Figure 3: Our model for phrase grounding as a sequence labeling task. The  $K \times K$  transition score matrix is derived from the features of  $K$  region proposals. The  $T \times K$  emission score matrix is derived from a joint representation of phrase-region pairs, which is fused from features of region proposals and  $T$  entity phrases. Bounding box regression is applied to the sequence of regions predicted by the CRF. **Cyan dashed line:** contextualized transition score prediction (Section 4.2).

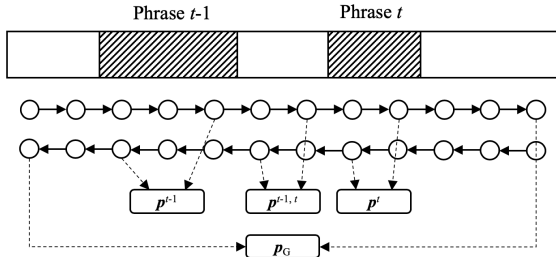


Figure 4: Text feature extraction for phrases in a caption. Shaded regions are entity phrase spans; circles represent LSTM cells. For phrase  $t$  hidden states at its span boundaries are concatenated to form its text features  $p^t$ , which is used in fusion with region features. For the contextualized transition score between phrases  $t-1$  and  $t$ , hidden states at the boundaries of the context between them are concatenated into a context feature vector  $p^{t-1,t}$ , which can be further extended by phrase features  $p^{t-1}$  and  $p^t$  as well as global text features  $p_G$ .

## 4.2 Model Specification

Figure 3 outlines our phrase grounding model.  $K$  region proposals and their visual and spatial features are extracted from an object detection vision backbone. We feed the token embeddings of the caption into a bi-directional LSTM (Hochreiter and Schmidhuber, 1997), and then concatenate the forward hidden state at the ending boundary of the phrase with the backward hidden state at the starting boundary of the phrase (see Figure 4). This phrase representation captures context both preceding and following the phrase in the caption.

$$\begin{aligned} \overrightarrow{(h^{1:L}, h^{1:L})} &= \text{BiLSTM}(\text{Embed}([c^1 \dots c^L])) \\ p^t &= [\overrightarrow{h^{e^t}} || \overleftarrow{h^{s^t}}] \end{aligned}$$

We use low-rank bilinear pooling (LRBP) (Kim et al., 2017) to fuse text and region features. Compared to simple concatenation, LRBP supports pairwise interaction between bimodal feature channels while keeping a reasonable computation overhead. Given a text feature vector  $p^t \in \mathbb{R}^{d_{\text{text}}}$  and a region feature vector  $r_k \in \mathbb{R}^{d_{\text{vis}}}$ , LRBP fuses them into a joint representation  $f_k^t \in \mathbb{R}^{d_{\text{joint}}}$ :

$$f_k^t = P^\top (U^\top p^t \circ V^\top r_k) + b$$

where  $U \in \mathbb{R}^{d_{\text{text}} \times r}$ ,  $V \in \mathbb{R}^{d_{\text{vis}} \times r}$ , pooling matrix  $P \in \mathbb{R}^{r \times d_{\text{joint}}}$ , bias  $b \in \mathbb{R}^{d_{\text{joint}}}$ , and  $\circ$  is the Hadamard (i.e. element-wise) product.

As discussed in Section 3.2, the CRF score function consists of emission score and transition score. The emission score  $\varepsilon(r_k, p^t)$  models the compatibility between each phrase and each candidate region. We feed the joint representation to a single-layer feed-forward neural network:

$$\varepsilon(r_k, p^t) = \text{FFN}(f_k^t)$$

The transition score  $\tau(r_k, r_{k'}, p^{1:T})$  is modeled by a two-layer feed-forward neural network with ReLU activation for the hidden layer:

$$\tau(r_k, r_{k'}, p^{1:T}) = \text{FFN}(\sigma(\text{FFN}([r_k || r_{k'}])))$$

To condition the transition scores on local and global context from the caption, we can extend the input  $[r_k || r_{k'}]$  with the following text features: context in between the two phrases (feature vector  $p^{t-1,t}$ ), context from phrase features  $p^{t-1}$  and  $p^t$ , and global context  $p_G$ .

One important difference between the standard use of CRFs for sequence labeling and our task is that our "labels" do not correspond to a fixed set of classes that can be predicted for any input, but are as specific to the particular input example as the sequences to be labeled themselves. Hence, our transition and emission scores do not depend on the (arbitrary) indices of regions to be ground, but on their visual and spatial features (as well as on their corresponding linguistic contexts). Finally, although our approach could in principle be extended to higher-order CRFs, we restrict our attention here to first-order CRFs for computational efficiency. As a consequence, our models can only capture dependencies between string-adjacent phrases.

### 4.3 Training Objectives

For each image-caption instance, the loss is a linear combination of the labeling and bounding box regression loss:

$$L = L_{\text{label}} + \gamma L_{\text{reg}}$$

$L_{\text{label}}$  is the CRF loss defined in Section 3.2.  $L_{\text{reg}}$  (Ren et al., 2017) is defined as

$$L_{\text{reg}} = (\beta, \hat{\beta}) = \sum_{i \in \{x, y, w, h\}} \text{SmoothL1}(\hat{\beta}_i - \beta_i)$$

with the ground truth regression parameterization

$$\beta = \left[ \frac{x - x_a}{w_a}, \frac{y - y_a}{h_a}, \log \frac{w}{w_a}, \log \frac{h}{h_a} \right]$$

and

$$\text{SmoothL1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

## 5 Experiments

### 5.1 Experiment Setup

**Dataset.** We train and evaluate our models on the Flickr30k Entities dataset (Plummer et al., 2017b), which contains 31,783 images, each accompanied by 5 captions. In keeping with previous work on this dataset, we assume that entity phrase boundaries are given, so inferring which phrases to ground is not part of our task. Following Plummer et al. (2017b), we merge all regions that are ground to the same phrase into one larger bounding box, and split the dataset into 29,783 training images, 1k validation images and 1k test images.

We do not apply our method to RefCOCO (Yu et al., 2016) or Visual Genome (Krishna et al., 2017) because they consist of independently grounded entity phrases without any entity dependencies that CRFs could leverage.

**Implementation details.** For text feature extraction, we use the 1024-d contextualized word embeddings from the last layer of ELMo (Peters et al., 2018), followed by a bi-directional LSTM (Hochreiter and Schmidhuber, 1997) encoder with hidden dimension  $d_{\text{hidden}} = 512$  for each direction, so that the text feature vector has dimension  $d_{\text{text}} = 1024$ . We use the Bottom-Up Attention model (Anderson et al., 2018) to generate region proposals and extract visual features, as in the state-of-the-art BAN (Kim et al., 2018) and DDPN (Yu et al., 2018b) models.  $K = 100$  region

proposals are generated for each image. Each candidate region with coordinates  $(x_{\min}, y_{\min})$ ,  $(x_{\max}, y_{\max})$  is represented by a  $d_{\text{vis}} = 2053$  feature vector that consists of 2048-d visual features concatenated with 5-d spatial features  $[x_{\min}/W, y_{\min}/H, x_{\max}/W, y_{\max}/H, wh/WH]$ . The low-rank bilinear pooling (LRBP) layer used for text-region bimodal feature fusion has rank  $r = 1024$  and output dimension  $d_{\text{joint}} = 1024$ . We train with a mini-batch size of 16 image-caption instances. Each instance contains all entity phrases to be grounded in the caption. Weights are initialized with Xavier (Glorot and Bengio, 2010). We apply a dropout rate of  $p = 0.2$  after the word embedding layer, LSTM layer, and LRBP fusion layer. The loss weighting parameter  $\gamma$  is 10.0. All gradients are clipped by  $\infty$ -norm of 10.0 to prevent gradient explosion. We do not fine-tune ELMo or the Bottom-Up Attention model. All models are trained for 50k iterations using Adam (Kingma and Ba, 2015) with learning rate  $5e - 5$  and  $\beta_1 = 0.9, \beta_2 = 0.98$ . Model snapshots are taken every 5k iterations and the model with the highest validation set accuracy is selected.

**Metrics.** We predict one grounded region for each entity phrase. Following Plummer et al. (2017b), a prediction is deemed accurate if it has at least 0.5 IoU overlap with the gold region. We report the percentage of accurately grounded phrases.

### 5.2 Quantitative Results

We compare our Soft-Label Chain CRF model against three baselines: a Hard-Label non-CRF model, a Hard-Label CRF, and a Soft-Label non-CRF model. The non-CRF models ground each phrase independently with a loglinear model. The Hard-Label models are trained with a standard one-hot training regime. The Soft-Label models use the soft-label training regime described above. The Soft-Label non-CRF model corresponds to the reduced form of the Soft-Label Chain CRF in Section 3.3.

Table 1 shows the performance of previous structured prediction models, current state-of-the-art models, our baseline models and the Soft-Label Chain CRF model. For a fair comparison with BAN (Kim et al., 2018), we also report result of the hard-label baseline with GloVe (Pennington et al., 2014) embeddings, while we obtain 0.33% higher result with ELMo. Training a non-CRF model on soft-label target distributions

Method	Vision Backbone	Grounding Accuracy (%)
<i>Compared methods</i>		
Structured Matching (Wang et al., 2016)	Fast R-CNN (Girshick, 2015)	42.08
Phrase-Region CCA (Plummer et al., 2017a)	Fast R-CNN (Girshick, 2015)	55.85
QRC Net (Chen et al., 2017b)	Fast R-CNN (Girshick, 2015)	65.14
BAN (Kim et al., 2018)	Bottom-Up Attention (Anderson et al., 2018)	69.69
DDPN (Yu et al., 2018b)	Bottom-Up Attention (Anderson et al., 2018)	73.3
<i>Our methods</i>		
Hard-Label (GloVe (Pennington et al., 2014))	Bottom-Up Attention (Anderson et al., 2018)	71.88
Hard-Label (HL)	Bottom-Up Attention (Anderson et al., 2018)	72.21
Soft-Label (SL)	Bottom-Up Attention (Anderson et al., 2018)	74.29
Hard-Label Chain CRF (HL-CCRF)	Bottom-Up Attention (Anderson et al., 2018)	72.26
Soft-Label Chain CRF (SL-CCRF)	Bottom-Up Attention (Anderson et al., 2018)	<b>74.69</b>

Table 1: Performance of different phrase grounding methods on Flickr30k Entities (test set). Our CRF models has transition scores conditioned on features of context in between the two phrases (“M” in Table 2). Our methods, unless explicitly specified, uses ELMo (Peters et al., 2018) as word embeddings.

Model	Transition Context	Accuracy (%)
SL-CCRF	–	74.28
SL-CCRF	M	<b>74.69</b>
SL-CCRF	M+LR	74.45
SL-CCRF	M+LR+G	74.48

Table 2: Performance of Soft-Label Chain CRF models by conditioning transition scores on different sets of context features. –: input to transition score prediction is  $[r_k || r_k]$ . M: input extended by features of context  $p^{t-1,t}$  in between the two phrases. M+LR: input further extended by features of LHS phrase  $p^{t-1}$  and RHS phrase  $p^t$ . M+LR+G: input further extended by features of global context  $p_G$ .

Decoding Algorithm	HL-CCRF	SL-CCRF
Viterbi (MAP)	72.26	74.69
Smoothing	<b>72.30</b>	<b>74.73</b>

Table 3: Decoding algorithms’ impact on performance.

improves accuracy by a further 2.08%. On top of that, Soft-Label Chain CRF improves accuracy by another 0.40%, which shows the effectiveness of treating phrase grounding as a sequence labeling task and using CRFs to capture entity dependencies. We also observe that the Hard-Label Chain CRF outperforms the hard-label baseline by a mere margin of 0.05%, so our conjecture is that using chain CRFs works well only with a suitable choice of training regime. Soft-Label Chain CRF gives an overall improvement of 2.48% over the hard-label baseline; it significantly outperforms previous structured prediction models including Structured Matching (Wang et al., 2016), Phrase-Region CCA (Plummer et al., 2017a) and QRC Net (Chen et al., 2017b), and surpasses the state-

of-the-art BAN (Kim et al., 2018) and DDPN (Yu et al., 2018b) models by a margin of 5.00% and about 1.4%, respectively.

We conduct an ablation study to find the most appropriate combination of context features for the transition scores in the SL-CCRF model. Table 2 shows that we obtain the best results by including the context in between the two phrases, which gives an improvement of 0.41%. We did not see any benefit from adding further text features from the left and right side of the phrases, or from the entire caption.

Besides the Viterbi decoding algorithm used in prediction in CRFs, we also experiment with a smoothing decoding algorithm. While Viterbi finds the MAP label sequence conditioned on the input sequence  $\arg \max_{\mathbf{y}} p(\mathbf{y}|\mathbf{x})$ , smoothing decoding finds the best label for each input  $x^t$ :  $\arg \max_{y^t} p(y^t|\mathbf{x})$ . This makes sense in some scenarios where we want to refine the predicted grounding of one entity by referring to the context instead of attempting to ground all entities mentioned in the description. Table 3 shows that in both Hard-Label Chain CRF and Soft-Label Chain CRF, smoothing decoding gives a prediction accuracy 0.04% higher than Viterbi decoding.

Without bounding box regression, the Soft-Label Chain CRF model has an accuracy of 69.85%, a 4.84% reduction compared to the setting with bounding box regression.

### 5.3 Qualitative Results

We visualize some phrase grounding results in the validation set of Flickr30k Entities in Figure 5. In (a), our CRF model avoids the error in grounding “a lounge chair” by constraining its relative posi-



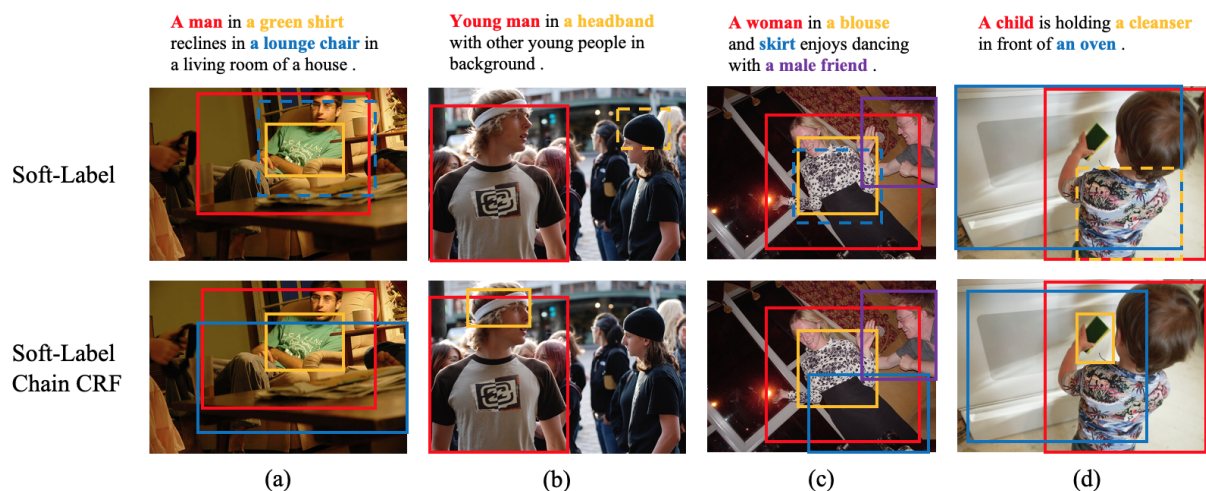


Figure 5: Selected visualization of phrase grounding results in the validation set of Flickr30k Entities. Solid boxes are correct predicted groundings, while dashed boxes are incorrect predicted groundings. Gold regions are not shown. Each entity phrase and its predicted grounding are marked with same color. Best viewed in color.

tion to “a man”. In (b), although it may not have learned to distinguish “headband” and “hat”, the CRF constrains the spatial position of “headband” to agree with the ownership dependency provided in context. In (c), it avoids the error in grounding “skirt” by spatially discriminating it from “a blouse”. In (d), it avoids the error in grounding “a cleanser” by constraining its relative size w.r.t. “a child”. These examples indicate that the CRF model may avoid grounding errors made by non-CRF models by leveraging entity dependencies, including relative position, spatial overlapping, and relative size.

## 6 Conclusion

In this paper, we formulate phrase grounding as a sequence labeling task and propose the Soft-Label Chain CRF model that successfully combines the benefits brought by global structured prediction and soft-label training regime that addresses the gold label multiplicity problem. Experimental results show that we achieve an overall improvement of 2.48% on grounding accuracy compared to a strong baseline, and that our model outperforms previous methods on phrase grounding.

## 7 Acknowledgements

This material is based upon work supported by the National Science Foundation under Grants No. 1405883 and 1563727. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not

necessarily reflect the views of the National Science Foundation.

## References

- Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Devi Parikh, and Dhruv Batra. 2017. [VQA: visual question answering - www.visualqa.org](http://www.visualqa.org). *International Journal of Computer Vision*, 123(1):4–31.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. [Bottom-up and top-down attention for image captioning and visual question answering](#). In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 6077–6086.
- Kan Chen, Trung Bui, Chen Fang, Zhaowen Wang, and Ram Nevatia. 2017a. [AMC: attention guided multimodal correlation learning for image search](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6203–6211.
- Kan Chen, Rama Kovvuri, and Ram Nevatia. 2017b. [Query-guided regression network with context policy for phrase grounding](#). In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 824–832.
- Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. 2017. [Long-term recurrent convolutional networks for visual recognition and description](#). *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(4):677–691.
- Mark Dredze, Partha Pratim Talukdar, and Koby Crammer. 2009. Sequence learning from data with mul-

- multiple labels. In *ECMLPKDD 2009 workshop on learning from multi-label data*, page 39.
- Hao Fang, Saurabh Gupta, Forrest N. Iandola, Rupesh Kumar Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C. Platt, C. Lawrence Zitnick, and Geoffrey Zweig. 2015. [From captions to visual concepts and back](#). In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1473–1482.
- Ross B. Girshick. 2015. [Fast R-CNN](#). In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 1440–1448.
- Xavier Glorot and Yoshua Bengio. 2010. [Understanding the difficulty of training deep feedforward neural networks](#). In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, pages 249–256.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural Computation*, 9(8):1735–1780.
- Andrej Karpathy and Li Fei-Fei. 2017. [Deep visual-semantic alignments for generating image descriptions](#). *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(4):664–676.
- Andrej Karpathy, Armand Joulin, and Fei-Fei Li. 2014. [Deep fragment embeddings for bidirectional image sentence mapping](#). In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 1889–1897.
- Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. 2018. [Bilinear attention networks](#). In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada.*, pages 1571–1581.
- Jin-Hwa Kim, Kyoung Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang. 2017. [Hadamard product for low-rank bilinear pooling](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. [Visual genome: Connecting language and vision using crowdsourced dense image annotations](#). *International Journal of Computer Vision*, 123(1):32–73.
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. [Conditional random fields: Probabilistic models for segmenting and labeling sequence data](#). In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML 2001), Williams College, Williamstown, MA, USA, June 28 - July 1, 2001*, pages 282–289.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. [Neural architectures for named entity recognition](#). In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 260–270.
- Xuezhe Ma and Eduard H. Hovy. 2016. [End-to-end sequence labeling via bi-directional lstm-cnns-crf](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [Glove: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1532–1543.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 2227–2237.
- Bryan A. Plummer, Arun Mallya, Christopher M. Cervantes, Julia Hockenmaier, and Svetlana Lazebnik. 2017a. [Phrase localization and visual relationship detection with comprehensive image-language cues](#). In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1946–1955.
- Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2017b. [Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models](#). *International Journal of Computer Vision*, 123(1):74–93.
- Filip Radenovic, Giorgos Tolias, and Ondrej Chum. 2016. [CNN image retrieval learns from bow: Unsupervised fine-tuning with hard examples](#). In *Com*

*puter Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 3–20.

Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. 2017. [Faster R-CNN: towards real-time object detection with region proposal networks](#). *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6):1137–1149.

Filipe Rodrigues, Francisco C. Pereira, and Bernardete Ribeiro. 2014. [Sequence labeling with multiple annotators](#). *Machine Learning*, 95(2):165–181.

Anna Rohrbach, Marcus Rohrbach, Ronghang Hu, Trevor Darrell, and Bernt Schiele. 2016. [Grounding of textual phrases in images by reconstruction](#). In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 817–834.

Mingzhe Wang, Mahmoud Azab, Noriyuki Kojima, Rada Mihalcea, and Jia Deng. 2016. [Structured matching for phrase localization](#). In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII*, pages 696–711.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. 2015. [Show, attend and tell: Neural image caption generation with visual attention](#). In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 2048–2057.

Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. 2016. [Modeling context in referring expressions](#). In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*, pages 69–85.

Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. 2017. [Multi-modal factorized bilinear pooling with co-attention learning for visual question answering](#). In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1839–1848.

Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. 2018a. [Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering](#). *IEEE Trans. Neural Netw. Learning Syst.*, 29(12):5947–5959.

Zhou Yu, Jun Yu, Chenchao Xiang, Zhou Zhao, Qi Tian, and Dacheng Tao. 2018b. [Rethinking diversified and discriminative proposal generation for visual grounding](#). In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 1114–1120.