

# Learning to Overcome Noise in Weak Caption Supervision for Object Detection

Mesut Erhan Unal, Keren Ye, Mingda Zhang, Christopher Thomas, Adriana Kovashka, Wei Li, Danfeng Qin, Jesse Berent.

**Abstract**—We propose the first mechanism to train object detection models from weak supervision in the form of captions at the image level. Language-based supervision for detection is appealing and inexpensive: many blogs with images and descriptive text written by human users exist. However, there is significant noise in this supervision: captions do not mention all objects that are shown, and may mention extraneous concepts. We first propose a technique to determine which image-caption pairs provide suitable signal for supervision. We further propose several complementary mechanisms to extract image-level pseudo labels for training from the caption. Finally, we train an iterative weakly-supervised object detection model from these image-level pseudo labels. We use captions from four datasets (COCO, Flickr30K, MIRFlickr1M, and Conceptual Captions) whose level of noise varies. We evaluate our approach on two object detection datasets. Weighting the labels extracted from different captions provides a boost over treating all captions equally. Further, our primary proposed technique for inferring pseudo labels for training at the image level, outperforms alternative techniques under a wide variety of settings. Both techniques generalize to datasets beyond the one they were trained on.

**Index Terms**—Language-supervised object detection, weakly-supervised object detection, vision and language

## 1 INTRODUCTION

LEARNING to localize and classify objects in images is a fundamental problem in computer vision. It has a wide range of applications, including robotics, autonomous vehicles, intelligent video surveillance, and augmented reality. Modern detectors are highly accurate [1], can run in real-time [2] and on mobile devices [3]. Despite these achievements, most modern detectors suffer from an important limitation: they are trained with expensive supervision in the form of large quantities of bounding boxes meticulously drawn by a large pool of human annotators. Due to the well-known domain shift problem [4], [5] and imperfect domain adaptation techniques, this means that when detection is to be performed in a novel domain, the expensive annotation procedure needs to be repeated.

Weakly supervised object detection (WSOD) techniques aim to alleviate the burden of collecting expensive box annotations. The classic WSOD formulation [6], [7], [8] treats an image as a bag of proposals, and learns to assign instance-level semantics to these proposals. WSOD has shown great potential for object detection, and recent methods have reached 52% mAP [9] on Pascal VOC 2012.

However, we highlight two limitations of WSOD methods. First, they depend on large-scale image-level object category labels; these require human effort that is provided in an unnatural, crowdsourced environment. Second, they make the assumption that the image-level label should be *precise*, i.e. at least one proposal instance in the image needs to be associated with the label. This assumption does not hold for real-world problems and real-world supervision.

- *The first five authors performed the work at University of Pittsburgh. The project was initiated at Google Research, Zurich, with which the last two authors are affiliated. Wei Li is now with NewsBreak Seattle.*  
E-mail: kovashka@cs.pitt.edu

Manuscript received ???.

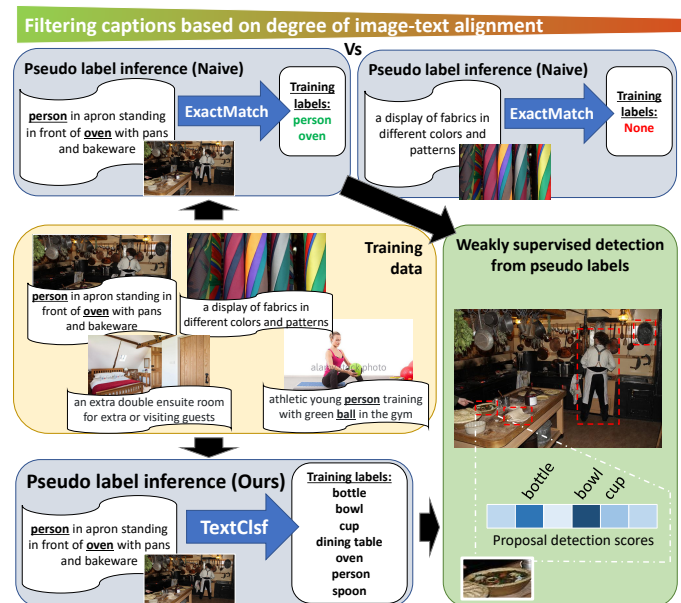


Fig. 1: We propose two mechanisms to infer pseudo training labels from captions. First (top), we determine the potential for strong object supervision signal from image-caption pairs (showing one with strong and one with weak signal). When supervision is strong, a simple training label extraction technique can be used. Second (bottom), we learn a mapping function (a text classifier) from captions to labels, which compensates for failures of exact-matching label extraction. Finally (right), we train a weakly-supervised object detection model with these pseudo image labels.

We propose mechanisms to leverage a new form of supervision for training weakly-supervised object detectors,

namely supervision in the form of natural language descriptions that web users provide when uploading their photos to social media sites such as Instagram, or their videos to video sharing platforms such as YouTube. There are tens of millions of photos uploaded to Instagram every day, and a majority have titles, tags, or descriptions. Abundant videos with subtitles or descriptive narratives are similarly available on YouTube. These annotations are “free” in that no user was *paid* to provide them; they arise out of innate needs of users to make their content available to others.<sup>1</sup>

However, existing WSOD methods cannot use such supervision. First, *natural language descriptions are unstructured*; they need to be parsed and words relevant for object recognition to be extracted, while non-object words are removed. Second, these descriptions are *both imprecise and non-exhaustive*—they might mention content that is not in the image (e.g. what event the user was attending or who they met after the photo was taken), and omit content that is in the image but is not interesting. In the bottom of Fig. 1, many large objects—e.g. dining table and bowls—were not mentioned in the human-provided description. Thus, directly feeding web data to the state-of-the-art WSOD system is infeasible, which under-utilizes the rich supervision that language on the web can provide.

To address this issue, we propose a three-part framework to build an object detector from images paired with accompanying captions (sentences). Our model bridges human-written free-form texts and visual objects, and generates accurate bounding boxes over objects in an image. Our key contributions are the first two steps, with a smaller contribution in the third step. First, we *estimate which image-caption pairs can serve as appropriate supervision* for extracting pseudo *image-level* training labels for training an object detector. In particular, we model the difference between images which are visual neighbors and those which appear with similar captions (semantic neighbors). We prioritize extracting signal from image-caption pairs where visual and semantic neighborhoods overlap, which indicates that captions closely follow the image. This enables the use of simple techniques for extracting training labels (Fig. 1 top).

Second, we devise complementary *advanced* techniques for *extracting pseudo image-level training labels from the caption*. One of our proposed strategies (Fig. 1 bottom) is to train a textual classifier to map captions to discrete object labels. Unlike the previous contribution, this classifier requires a small set of labels, and enables us to bridge the gap between what humans mention in a caption, and what truly is in an image. Alternatives include learning multimodal spaces where images and captions are projected, and using similarity in these spaces to determine which captions are similar to object words in a predefined vocabulary. This contribution and the previous have different applications: The former is fitting when no labels are available, but if they are, the latter achieves slightly stronger performance. Thus, we primarily focus on evaluating these contributions separately, as shown in Fig. 1. Both contributions generalize beyond dataset boundaries.

Third, we use the *pseudo ground truth labels at the image*

*level* (extracted in the previous step), *to train a weakly supervised object detection method*. The method we propose extracts region proposals off-the-shelf, then for each proposal and each class, learns both a class score and a detection score. These scores are then refined using an iterative approach, and combined to produce final detection results.

In our work, we first need to infer image-level pseudo labels for training from the available captions. Only then can we proceed to train a weakly-supervised detection (WSOD) algorithm, using those (potentially noisy) image-level pseudo labels. Thus, to distinguish our work from WSOD, we refer to our methods as performing **language-supervised object detection (LSOD)**.

Our paper makes four main contributions. First, we propose a new task of learning from noisy caption annotations, and set up a new benchmark. Rather than treating object categories as IDs only, we also leverage their semantics and synonyms of those object names. Second, we show the impact of multiple possible ways to map captions to image-level labels, ranging from exactly matching the captions to object category words, using learned image-text similarity scores, retrieving hand-annotated or predicted synonyms to the object categories from the captions, or training a classifier. Our proposed approach outperforms the baseline by up to 78% on noisy datasets. Third, we demonstrate the success of explicitly modeling which image-caption pairs provide strong signal for supervision, using a new metric that captures how closely the text follows the image. This alignment metric allows us to improve performance by up to 37%. Fourth, we show cross-domain results in datasets: we not only demonstrate competitive WSOD performance by training/testing on COCO captions, but also validate the benefit of our COCO-trained text classifier and alignment metric by applying it on Flickr30K, and the noisy MIR-Flickr1M and Conceptual Captions. We are not aware of other work that directly extracts labels for detection training from the latter two datasets (and only a few works pretrain for detection on Conceptual Captions). We leverage the resulting models and evaluate them on the PASCAL and COCO object detection datasets.

The remainder of the paper is organized as follows. We overview related work in Sec. 2. In Sec. 3, we discuss how to filter or weight image-caption pairs as potential signals for supervision (Sec. 3.1), different ways to reduce the gap between free-form captions and object categories (Sec. 3.2), and the backbone of our WSOD model, which combines prior work [8], [10] in a new way (Sec. 3.3). In Sec. 4, we compare to upper and lower bounds, in conjunction with state-of-the-art methods. We conclude in Sec. 5.

## 2 RELATED WORK

We formulate a new variant of weakly-supervised object detection, which we term *language-supervised*, where the supervision is even more weak but less costly than in prior work. We leverage vision-language interactions, so we also discuss work that finds alignment between image regions and text and grounds language in images. We also discuss recent work in learning visual representations from language. Finally, we describe work that investigates what kind of content humans describe in captions or models how

1. Of course, this data may be subject to license agreements limiting uses, and not all of it can truly be used for “free.”

visually concrete particular words are. We are not aware of other work that *explicitly handles noise for language-supervised object detection*, as we propose.

## 2.1 Weakly-supervised object detection

Weakly-supervised object detection (WSOD) involves localizing and categorizing objects without instance-level (bounding box) supervision. Key approaches include multiple-instance learning (MIL) where one or more regions are associated with the label of interest [8], [11], and self-training, where high-scoring proposals are treated as pseudo ground-truth [8], [12], [13]. In the multiple-instance learning (MIL) setting, proposals of an image are treated as a bag of candidate instances. If the image is labeled as containing an object, at least one of the proposals will be responsible to provide the prediction of that object. Oquab et al. [14] and Zhou et al. [15] propose a Global Average (Max) Pooling layer to learn class activation maps. Bilen et al. [6] propose Weakly Supervised Deep Detection Networks (WS-DDN) containing classification and detection data streams, where the detection stream weighs the results of the classification predictions. Tang et al. [8], [16] jointly train multiple refining models together with WSDDN, and show the final model benefits from the online iterative refinement. Diba et al. [17] and Wei et al. [7] apply a segmentation map; Wei et al. [7] further incorporate saliency. Wan et al. [18] add a min-entropy loss to reduce the randomness of the detection results. Zeng et al. [12] jointly consider bottom-up and top-down objectness from low-level measurement and CNN confidences. Ren et al. [9] aim at instance-aware self-training where they design DropBlock to zero out the most discriminative parts to avoid the part domination issue in WSOD. Earlier (pre-deep-learning) approaches include Divvala et al. [19] which rely on web search for an initial set of concepts for which to learn detection models, prune them based on model performance, and combine synonyms.

Our work is similar to these since we also represent the proposals using a MIL-weighted representation. However, prior WSOD methods require structure in the form of class labels, and these labels require dedicated human effort. Our contribution is enabling weakly-supervised detection with less costly *language supervision* which could work without explicit human annotations. In this project, we use both crowdsourced captions (from the COCO and Flickr30K datasets) and noisier ones obtained as a side product of users uploading content on the web (MIRFlickr1M, Conceptual Captions). We explicitly handle the noise in the language supervision and the misalignment between nouns (objects) that are shown but not mentioned, or mentioned but not shown. This distinguishes our work from both WSOD and self-supervised methods.

## 2.2 Vision-language tasks

Learning visual-semantic embeddings (VSE) has received tremendous interest due to its broad applications such as retrieval [20], [21], captioning [22], [23], and visual question answering [24]. VSE approaches learn a joint visual-text space, e.g. via a triplet or contrastive loss, where the distance between embedded samples reflects their semantic relationship, and cross-modal attention [25], [26]. As a side

experiment, we conducted a transformer-based pretraining involving both masked language modeling and image-text matching objectives to achieve better visual features. Still, progress in transformers is orthogonal to our primary aim as it does not consider how strong of a signal a caption provides for its co-occurring image.

There is also work to associate phrases in the caption to visually depicted objects [27] but none enable training of an independent object detector with accurate localization and classification, as we propose. In recent work, [28] predict masked words without localization, but use surrounding text at test time, unlike our models.

In Thomas and Kovashka [29], we show that image-text matching fails when the relation between an image and its corresponding (co-occurring) text is complementary rather than redundant. What this means for training object detection models from language supervision, is that the category overlap between image and co-occurring text may be low. To cope with this, our method exploits the structure of each unimodal space (image and text), and compares those structures, to compute how relevant each caption is for each image, and thus, whether the image-caption pair should be used for training object detectors.

## 2.3 Learning visual representations from text

Recent work [30], [31], [32], [33], [34] aims to learn visual representations from their corresponding textual counterparts. Gomez et al. [30] predict the text LDA topic distribution from the image feature. Miech et al. [31] assume an MIL nature in videos, and use Noise Contrastive Estimation (NCE) to optimize the alignment between video clips and associated narrations. Desai and Johnson [33] harvest visual representations from training bidirectional captioning models and note the importance of predicting all caption tokens to learn a good visual representation. Radford et al. [34] optimize a classical co-attention model but learn the feature representation on a large dataset of 400M image-text pairs. However, these methods do not train standalone object detectors. Bertasius et al. [32] apply a transformer-based language model to encode the text and match the visual feature extracted by an object detection model. They only optimize the representation to classify objects, while we also care about the detection scores and learn them in the unified framework since the visual proposals we use (Selective Search) are not as accurate as detection results. Most related to our work is Chen et al. [35]. The algorithm in this work discovers and localizes new objects from documentary *videos* by associating subtitles to video tracklets. They extract keywords from the subtitles using TFIDF. However, video provides benefits we cannot leverage, e.g. numerous frames containing nearly identical object instances. Importantly, we show that only using words that actually appear in the caption (as done in [35] with TFIDF) results in suboptimal performance compared to our method. Further, many components of Chen et al.'s method, e.g. the restriction to animal classes and the reliance on tracking, limit generalizability to other vocabularies and to images.

In our preliminary work, Ye et al. [36], we show we can successfully leverage unstructured supervision (highly descriptive captions well-aligned with the visual modality)

but we do not explore any filtering or weighting of image-caption pairs. This weighting allows us to bypass the need for training a text-only classifier (which required a small amount of class labels), replacing it with techniques that require only image-caption pairs. In this work, we also include results on noisier datasets (MIRFlickr1M and Conceptual Captions) and without ImageNet pretraining.

## 2.4 Visual reporting bias and concreteness

Our results show there is a gap between what humans name in captions, and what categorical annotations they provide. [37] study a similar phenomenon they refer to as “human reporting bias”. They model the presence of an object as a latent variable, but we do the opposite—we model “what’s in the image” by observing “what’s worth saying”. Further, we use the resultant model as precise supervision to guide detection model training.

Our work also measures how abstract is the connection between an image and a co-occurring text. Prior work predicts whether image and text that co-occur have a direct or complementary relationship [38], [39], e.g. whether the relation between image and its caption is “visible”, “story”, “subjective” or “meta” [39]. Unlike our method which also implicitly measures abstractness, these methods require additional annotations, aim for a discrete rather than continuous abstractness score, and are not applied in an object detection setting. Also related is work that measures how tightly clustered the visual companions of a word are [40] but this approach only computes scores for individual words, not for the relationships within image-caption pairs. In an auxiliary task, we compute the potential of [39], [40] to predict which captions may serve as clean supervision for weakly-supervised object detection, and we show that our method is equally or more promising.

## 3 APPROACH

We train object detectors from supervision only consisting of noisy captions and corresponding images. In realistic scenarios, captions and images may contain complementary information. We hypothesize that even for crowdsourced, descriptive captions which closely follow the image (e.g. COCO), not all caption-image pairs provide equally strong supervision, as some captions will overlap with the image to a stronger degree. Fig. 1 (top) shows two images, the first with high image-text alignment, where two objects (highlighted) are both shown and mentioned. The second image contains concepts that are visually not shown or are visually ambiguous (e.g. display, fabrics), hence extracting concrete nouns (objects) is more challenging. Thus, the first step in our framework (Sec. 3.1) is to determine which image-caption pairs to use as supervision: we propose two alternative approaches, one which uses a hard cutoff over the image-text alignment score, and another which uses all image-caption pairs but gives them different weight.

After selecting image-caption pairs for training, we next extract discrete labels at the image level (Sec. 3.2). We do so through a variety of techniques, the simplest of which is looking for exact string match between nouns in the caption and object words, and the most complex being training a

classifier which takes in a caption (without a paired image) and maps this caption to a discrete set of labels (which may or may not be mentioned in the caption). Finally, given these pseudo ground-truth image-level labels, we train a variant of a prior weakly-supervised object detection technique: it first computes initial scores for each region and each object class, then refines these iteratively (Sec. 3.3).

## 3.1 Filtering captions by estimated supervision purity

We propose to filter image-caption pairs that are unlikely to be useful for training. The key idea is to estimate to what extent an image caption and the image provide overlapping (redundant) or complementary information. While complementarity is useful in general, for detection we require redundancy, i.e. the same objects being both shown and mentioned in the caption. We first describe an *example* of how to learn a joint image-text embedding; we do not require any particular technique for this part. We then compute homogeneity, i.e. how visually similar semantic neighbors of an image (images whose captions are similar in a *unimodal* word embedding space) are in the learned *joint, multimodal* embedding space. This homogeneity allows us to estimate the overlap of the image-caption pair. The computation of homogeneity follows our prior work [41], but was never used for object detection before.

### 3.1.1 Preliminaries

Let  $\mathcal{D} = \{\mathbf{I}, \mathbf{T}\}$  represent a dataset of  $n$  image-text pairs, where  $\mathbf{I} = \{x_1, x_2, \dots, x_n\}$  and  $\mathbf{T} = \{y_1, y_2, \dots, y_n\}$  are the set of images and text (captions), respectively, and  $y_i$  is text co-occurring with image  $x_i$  (the two are semantically related). To reason about the relationships of images and text, we seek a joint manifold  $\mathcal{M}$ . For images, a convolutional network  $f : \mathbf{I} \rightarrow \mathcal{M}$  is used to project images into the joint space, while a recurrent network  $g : \mathbf{T} \rightarrow \mathcal{M}$  projects text. To obtain  $\mathcal{M}$ , we can use any cross-modal retrieval method. We describe two possibilities, triplet loss [42] and polysemous embedding model (PVSE) [21].

We first consider a simple triplet loss to derive  $\mathcal{M}$ :

$$\mathcal{L}_{trip} = [\|x_i - y_i\|_2^2 - \|x_i - y_j\|_2^2 + m]_+ \quad (1)$$

where  $x_i, y_i$  appear together, while  $x_i, y_j$  do not,  $+$  denotes hinge loss, and  $m$  is a margin. Alternatively, PVSE uses a multiple-instance variant of triplet loss ( $K$  meanings hence  $K$  embeddings per sample), along with self-attention for visual and text features (not shown). The similarity of image  $x_i$  and text  $y_i$  is:

$$s(x_i, y_i) = \max_{(k_1, k_2) \in \{1, \dots, K\} \times \{1, \dots, K\}} \left\langle \frac{x_{i_{k_1}}}{\|x_{i_{k_1}}\|_2}, \frac{y_{i_{k_2}}}{\|y_{i_{k_2}}\|_2} \right\rangle \quad (2)$$

### 3.1.2 Homogeneity

To capture how well-aligned an image and its corresponding caption are, we measure visual homogeneity (similarity) of the semantic concepts that an image illustrates. In other words, are images corresponding to semantically similar texts, visually similar? To measure this homogeneity, we first discover each image-text pair’s *semantic neighbors* in text space  $\Omega(\mathbf{T})$ . Following [43], we compute neighbors in text

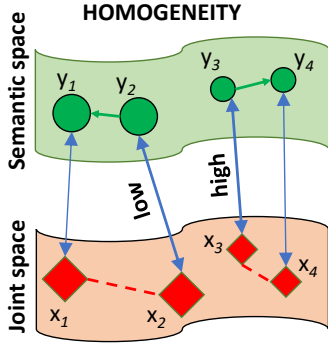


Fig. 2: Image-caption weighting: Green arrows connect neighbors in the original semantic space. Blue links show co-occurring images and text. Images whose texts are close in semantic space, which are close in the joint space (short red links), have high *homogeneity* scores ( $\gamma$ 's in Eq. 9).

space<sup>2</sup> because the text domain provides the cleanest semantic representation of the image-text pair. Let  $\Psi(\Omega(y_i)) = \{\langle x'_{i_n}, y'_{i_n} \rangle\}_{n=1}^N$  represent the semantic nearest neighbor function over  $\Omega(\mathbf{T})$ , where  $\{\langle x'_{i_n}, y'_{i_n} \rangle\}_{n=1}^N$  denotes the set of the  $N$  neighbors of  $\langle x_i, y_i \rangle$  and  $\langle x_i, y_i \rangle \notin \Psi(\Omega(y_i))$ .

We next measure the homogeneity of the semantic neighbors in both the image and text domains, using the  $f: \mathbf{I} \rightarrow \mathcal{M}$  and  $g: \mathbf{T} \rightarrow \mathcal{M}$  projections of image and text into the *joint space*. Because our formulation is equivalent for both image/text neighbors, we let  $s_i$  represent a sample from either domain but require samples  $s_i$  and  $s_j$  come from the same domain. Let  $\mathbf{s}'_i = [s'_{i_1}, s'_{i_2}, \dots, s'_{i_N}]^T$  denote the  $N \times H$  matrix of embeddings of the neighbors of  $s_i$  found via  $\Psi$ , and  $\mathbf{U} = \mathbf{s}'_i \mathbf{s}'_i{}^T$  compute the pairwise similarities between all semantic neighbors through cross-product. We compute the *homogeneity score*  $\alpha_i^{HOM}$  for  $s_i$  as follows:

$$\alpha_i^{HOM} = \frac{1}{N^2} \sum_{r=1}^N \sum_{c=1}^N \mathbf{U}_{(r,c)} \quad (3)$$

where  $r, c$  index over the rows and columns of  $\mathbf{U} = \mathbf{s}'_i \mathbf{s}'_i{}^T$ . For the different image ( $s_i = x_i$ ) and text ( $s_i = y_i$ ) domains, we compute *visual homogeneity score*  $\alpha_{(I)}^{HOM}$  and *text homogeneity score*  $\alpha_{(T)}^{HOM}$ , respectively. Both of these scores capture how aligned an image and its co-occurring text are; thus, **higher  $\alpha$  scores indicate image-text captions from which supervision signal can more reliably be extracted.** We also consider the difference of  $\alpha_{(I)}^{HOM}$  and  $\alpha_{(T)}^{HOM}$  as an indicator for supervision purity.

Note that the cost for computing homogeneity score is neglectable in that we *offline preprocessed* all examples in the training set by caching sample embeddings into a memory bank. We only find semantic neighbors once using a pre-trained Doc2Vec. Then, computing  $\alpha^{HOM}$  weights is efficient as it only requires multiplication.

### 3.1.3 Scoring captions for homogeneity

For object detection, we prefer to train from examples in which objects are both shown in the image and mentioned in the text. Sec. 3.1.2 provides a way to measure the redundancy between the image and text modalities. Here, we describe how to use this measure.

**Filtering.** We hypothesize that selecting training data and filtering out noisy image-caption pairs using  $\alpha_{(I)}^{HOM}$

2. Specifically, Doc2Vec [44] due to its appropriateness for longer texts, although BERT [45] could also be used.

and  $\alpha_{(T)}^{HOM}$  will improve the detection model training. We provide an experiment in Sec. 4.4, which selects the 30,000 image-caption pairs from COCO that have the highest homogeneity scores. Our method provides significantly better detection results than random selection (Tab. 5).

**Weighting.** Hard-cutoff filtering requires finding the right cutoff value (e.g. top-30k), and it means discarding some potentially useful data. Compared to the filtering strategy, weighting does not require a hard cutoff and is more data-efficient. It applies different weights to image-caption pairs. For image-caption pairs that are more overlapped (i.e. high homogeneity), weighting assigns large weight to the loss term in that these examples will likely be useful for training detection models. For image-caption pairs that are more complementary, weighting assigns small weights because the information may not well-aligned. In Eq. 9, we use  $\alpha_{(I)}^{HOM}$  as the heuristic weighting factor  $\gamma$ . We provide an ablation in Sec. 4.4, and Tabs. 6, 7 and 8 show the impact of using the weighting factor  $\alpha_{(I)}^{HOM}$ .

The strategy described in this section allows us to select captions with useful supervision, thus can be coupled with simple mechanisms to extract labels at the image level from captions, e.g. EXACTMATCH in the next subsection. We do not expect it to improve results when a mapping function from captions to training labels is learned with supervision.

## 3.2 Pseudo training label inference from text

After getting the image-caption pairs estimated to be well-aligned, we now proceed to extract pseudo object labels from the selected noisy captions, to benefit weakly-supervised object detection. The foundation of WSOD builds on an important assumption from MIL (Eq. 7), which suggests that *precise* image-level labels should be provided. The straightforward solution is to extract object labels from captions via lexical matching. However, it has limitations. Consider an image with three sentence descriptions:

“a *person* is riding a *bicycle* on the side of a *bridge*.”

“a *man* is crossing the street with his *bike*.”

“a *bicyclist* peddling down a busy city street.”

Only the first sentence exactly matches the categories “person” and “bicycle”. Even if we allow synonyms of “man” and “person” or “bicycle” and “bike”, only the first two precisely describe both objects, while the last one still misses the instance of “bicycle” unintentionally. When using these examples to train object detectors, the first two instances may bring positive effect, but the last one will be wastefully discarded as false negative i.e. not relevant to the categories “person” or “bicycle”. As further examples, in Fig. 1 (bottom-right), none of the captions (one shown) mention the “bowls” or “spoons” that are present, and only some mention “oven”. Finally, in Fig. 6, the caption mentions a “suit” worn by a speaker at a conference, but not the “tie”, even though one is present.

This observation inspires us to amplify the supervision signal that captions provide, and squeeze more information out of them. Fig. 3 (bottom) shows the approach we use to amplify the signal. This text-only model takes free-form texts as input, embeds individual words to a 300D space using GloVe [46], and projects the embedded features to a 400D latent space. We then use max-pooling to aggregate

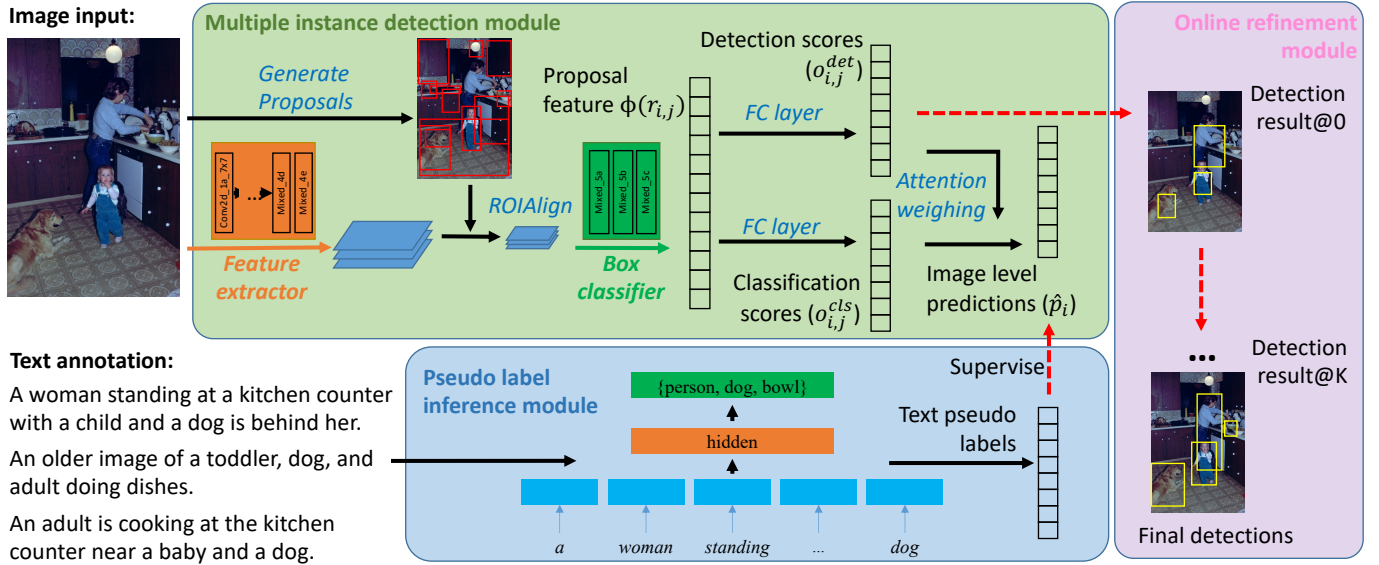


Fig. 3: Harvesting detection models from free-form text. We propose to use a pseudo training label inference module (bottom) to amplify signals in free-form texts to supervise the learning of the multiple instance detection network (top). The detection model is refined by an online refinement module (right) to produce the final detection results. Detection  $(o_{i,j}^{det} \in \mathbb{R}^C)$  and classification  $(o_{i,j}^{cls} \in \mathbb{R}^C)$  scores and image predictions  $\hat{p}_i \in \mathbb{R}^C$  refer to predictions for all classes.

the word-level representations. Then, we use this intermediate representation to predict the implied instances (e.g. 80 classes as defined in COCO, or any other categories); this prediction answers “what’s in the image” and serves as pseudo image-level training labels, to be used in Sec. 3.3.

There exists a subtle balance when using pseudo labels to train object detectors. Admittedly, our strategy increases the recall rates thus more data could be utilized. However, with the increased recall, precision will drop inevitably thus the fundamental assumption in MIL is threatened. Specifically, the *precise label* assumption makes the model very sensitive to false positive cases: when inappropriate labels are given where none of the proposals have a good response, the model gets confused, resulting in non-optimal detections. Thus, we adopt a two-steps procedure: first we look for an exact match of object labels from captions, following the intuition that *explicitly mentioned objects* should be significant and obvious enough in the image; second, when no object can be matched, we use our label inference model to predict labels as *unspoken intended objects* to guide the object detection. We show our method TEXTCLSIF outperforms several strong alternatives that also infer pseudo labels.

**Alternative strategies.** We also experiment with alternative multiple pseudo-label generation techniques when lexical matching (EXACTMATCH) fails to find a match. First, we consider a *manually constructed, hence expensive* COCO synonym vocabulary list (EXTENDVOCAB) which maps 413 words to 80 categories [47]. Another variant, GLOVE, takes advantage of GloVe word embeddings [46], assigning pseudo-labels for a sentence by looking for the category that has the smallest embedding distance to any word in the sentence. We also finetune the GloVe word embeddings on COCO using a visual-text ranking loss, and use the pseudo labels retrieved by the resultant text embedding, resulting in LEARNEDGLOVE.

**Discussion.** Our text classifier relies on both captions

and category labels. Once the bridge between captions and labels is established, it generalizes to other datasets, as we show in Tab. 1. Importantly, we only need a small fraction of labels to train this text classifier; as we show in Fig. 5, precision has a small range (between 89% and 92%) when we use between only 5% and 100% of the COCO data, while recall is stable at 62%. Thus, our text model could learn from a *single source* dataset with a *few* labels, then it could transfer the knowledge to other *target* datasets, requiring only free-form text as supervision. If no labels are available, the caption weighting strategy (Sec. 3.1) can be paired with EXACTMATCH or other alternatives to TEXTCLSIF. However, TEXTCLSIF performs slightly better overall.

### 3.3 Detection from inferred pseudo image labels

We next describe how we use the inferred pseudo labels at the image level, to train an object detection model. As shown in Fig. 3 (top), we first extract proposals with accompanying features. An image is fed into randomly initialized or pretrained (on ImageNet [48]) convolutional layers. Then, *ROIAlign* is used for cropping the proposals (at most 500 boxes per image) generated by *Selective Search* [49], resulting in fixed-sized convolutional feature maps. Finally, a box feature extractor is applied. If  $\{r_{i,j}\}_{j=1}^m$  are the  $m$  proposals of a given image  $x_i$ , this process results in proposal feature vectors  $\{\phi(r_{i,j})\}_{j=1}^m$  where each  $\phi(r_{i,j}) \in \mathbb{R}^d$ . Note that even when our model is pretrained on ImageNet, it *does not leverage* any image labels on the datasets on which we train and evaluate our detection models (PASCAL and COCO).

#### 3.3.1 Initial detection scores

We next introduce the prediction of image-level labels  $\hat{p}_{i,c}$  ( $c \in \{1, \dots, C\}$  for the  $i$ -th image, where  $C$  is the number of classes) and of detection scores. If not noted otherwise, we use  $i$  to index training examples,  $j$  to index region proposals

within an image, and  $c$  to index class labels. The method described in this section follows prior work, i.e. Bilen et al. [6], but labels used for training are potentially noisy as they come from our pseudo label inference module, Sec. 3.2.

First, we feed the proposal features  $\phi(r_{i,j})$  into two parallel fully-connected layers to compute the detection scores  $o_{i,j,c}^{\det} \in \mathbb{R}^1$  (top branch in the green MIL module in Fig. 3) and classification scores  $o_{i,j,c}^{\text{cls}} \in \mathbb{R}^1$  (bottom branch), in which both scores are related to a specific class  $c$  and the  $j$ -th proposal of image  $x_i$ :

$$o_{i,j,c}^{\text{cls}} = w_c^{\text{cls}\tau} \phi(r_{i,j}) + b_c^{\text{cls}}, \quad o_{i,j,c}^{\det} = w_c^{\det\tau} \phi(r_{i,j}) + b_c^{\det} \quad (4)$$

We convert these scores into: (1)  $p_{i,j,c}^{\text{cls}}$  the probability that object  $c$  presents in the  $j$ -th proposal; and (2)  $p_{i,j,c}^{\det}$  the probability that the  $j$ -th proposal is important for predicting image-level label  $y_{i,c}$ :

$$p_{i,j,c}^{\text{cls}} = \sigma(o_{i,j,c}^{\text{cls}}), \quad p_{i,j,c}^{\det} = \frac{\exp(o_{i,j,c}^{\det})}{\sum_{j=1}^m \exp(o_{i,j,c}^{\det})} \quad (5)$$

Finally, the aggregated image-level prediction is computed as follows, where greater values of  $\hat{p}_{i,c} \in [0, 1]$  mean higher likelihood that  $c$  is present in the image  $x_i$ :

$$\hat{p}_{i,c} = \sigma\left(\sum_{j=1}^m p_{i,j,c}^{\det} o_{i,j,c}^{\text{cls}}\right) \quad (6)$$

Assuming the label  $y_{i,c} = 1$  if and only if class  $c$  is present in the input image  $x_i$ , the multiple instance detection loss used for training the model is defined as:

$$L_{\text{mid}}(x_i, y_i) = -\sum_{c=1}^C \left[ y_{i,c} \log \hat{p}_{i,c} + (1 - y_{i,c}) \log(1 - \hat{p}_{i,c}) \right] \quad (7)$$

The weakly supervised detection score given both proposal  $r_{i,j}$  and class  $c$  is the product of  $p_{i,j,c}^{\text{cls}}$  and  $p_{i,j,c}^{\det}$  which is further refined as described in Sec. 3.3.2.

### 3.3.2 Online instance classifier refinement

The third component of our WSOD model is Online Instance Classifier Refinement (OICR), as proposed by Tang et al. [8]. Given a ground-truth class label, the top-scoring proposal, as well as proposals highly overlapping with it, are selected as references. These proposals are treated as positives for training the box classifier of this class while others are treated as negatives. The initial top-scoring proposal may only partially cover the object, so allowing highly-overlapped proposals to be treated as positives gives them a second chance to be considered as containing an object, in the subsequent model refinement. This reduces the chance of propagating incorrect predictions. In addition, sharing the convolutional features between the original and refining models makes training more robust.

Following [8], we stack multiple refining classifiers and use the output of the previous one to generate instance-level supervision to train the successor. The detection score at the 0-th iteration is computed using  $s_{i,j,c}^{(0)} = p_{i,j,c}^{\text{cls}} p_{i,j,c}^{\det}$  ( $c \in \{1, \dots, C\}$ ),  $s_{i,j,C+1}^{(0)} = 0$  (where  $C+1$  is the background class). Given the detection score  $s_{i,j,c}^{(k)}$  at the  $k$ -th iteration, we use the image-level label to get the *instance-level* supervision  $y_{i,j,c}^{(k+1)}$  at the  $(k+1)$ -th iteration. Assume that  $c'$  is a label

attached to image  $x_i$ , we first look for the top-scoring box  $r_{i,j'}$  ( $j' = \arg \max_j s_{i,j,c'}^{(k)}$ ) responsible to predict  $c'$ . We then let  $y_{i,j,c'}^{(k+1)} = 1, \forall j \in \{j' | IoU(r_{i,j}, r_{i,j'}) > \text{threshold}\}$ . When  $k > 0$ ,  $s_{i,j,c}^{(k)}$  is inferred using a  $(C+1)$ -way FC layer, as in Eq. 4. The OICR training loss is defined in Eq. 8.

$$L_{\text{oicr}}^k(x_i, y_i) = -\frac{1}{m} \sum_{j=1}^m \sum_{c=1}^{C+1} \hat{y}_{i,j,c}^{(k)} \log s_{i,j,c}^{(k)}, \quad k = 1, \dots, K \quad (8)$$

Unlike the original OICR, our WSOD module aggregates logits instead of probability scores, which in our experience stabilizes training. We also removed the reweighing of untrustworthy signals emphasized in [8] since we found it did not contribute significantly.

The final loss we optimize is Eq. 9, where  $\gamma_i$  is the *per-example* weighting factor.  $\gamma_i = 1$  for all  $i$  if we are not applying homogeneity weighting. If we use hard filtering based on  $\alpha_{(I)i}^{HOM}$ ,  $\alpha_{(T)i}^{HOM}$ , then  $\gamma_i = 1$  for samples included in training, and  $\gamma_i = 0$  for others. If using image-caption weighting,  $\gamma_i = \alpha_{(I)i}^{HOM}$ . We refine our model for 3 times ( $K = 3$ ) if not mentioned otherwise.

$$L(x_i, y_i) = \gamma_i \left( L_{\text{mid}}(x_i, y_i) + \sum_{k=1}^K L_{\text{oicr}}^k(x_i, y_i) \right) \quad (9)$$

## 3.4 Implementation details

We first obtain a joint image-caption embedding space, by training a triplet loss model [42] (with margin  $m$  set to 0.5) or PVSE [21] model (with  $K = 3$  embeddings per sample and margin 0.1, using the COCO validation set). We use this model to infer image-text joint embeddings. Then we use Eq. 3 to infer the homogeneity scores. We use ResNet-50 [50] initialized randomly or with ImageNet features for images (both types of results shown in Table 8). Images are scaled to 224x224 and augmented with random horizontal flipping. For text, we use GRU [51] with hidden state size 512, initialized with 200D word embeddings learned on the COCO dataset, using [52]'s implementation of Doc2Vec, distributed memory [44], 20 epochs with window size of 20, and ignoring words that appear less than 20 times. We use Xavier initialization [53], the Adam optimizer [54] with minibatch size of 32, learning rate 1.0e-4 (decayed by 10x after every 5 epochs of no decrease in val loss), and weight decay 1e-5. We use [55] to efficiently compute approximate nearest neighbors for  $\Psi$  and use  $N = 200$  nearest neighbors.

For the text classifier which predicts the pseudo image-level labels, we adopt a multi-layer perceptron (see Fig. 3 bottom). We first embed word tokens to 300D GloVe embeddings and project them to a 400D latent space. We use max-pooling to aggregate these word-level representations to get the fixed-length 400D caption representation. Next, we use this max-pooled intermediate representation to predict the implied instances (e.g., 80 classes in COCO). The object labels are used to supervise the text classifier learning (cross-entropy loss). We use AdaGrad optimizer [56], learning rate of 0.1, and batch size of 20.

To train the weakly supervised object detection model, we first use Selective Search [49] from OpenCV [57] to extract at most 500 proposals for each image. We follow the

“Selective search quality” parameter settings in [49]. We use Selective Search because it is a generic, dataset-independent proposal generation procedure, as opposed to other CNN-based alternatives which are trained end-to-end from a specific dataset in a supervised fashion. We use TensorFlow [58] as our training framework. To compute the proposal feature vectors, we use the layers (“Conv2d\_1a\_7x7” to “Mixed\_4e”) from Inception-V2 [10] to get the conv feature map, and the layers (“Mixed\_5a” to “Mixed\_5c”) from the same model to extract the proposal feature vectors after the ROIAlign operation. The Inception-V2 model is either randomly initialized, or pretrained on ImageNet [48]; the supervised detector counterpart of our model, using this architecture, was explored by [59]. To augment the training data, we resize the image randomly to one of the four scales {400, 600, 800, 1200}. We also randomly flip the image left to right at training time. At test time, we average the proposal scores from the different resolution inputs. We set the number of refinements to 3 for the OICR since it gives the best performance. For post-processing, we use non-maximum-suppression with IoU threshold of 0.4. We use the AdaGrad optimizer [56], a learning rate of 0.01, and a batch size of 2 as commonly used in WSOD methods [8], [16]. The models are usually trained for 100K iterations on Pascal VOC (roughly 40 epochs on VOC2007 and 17 epochs on VOC2012) and 500K on COCO (8.5 epochs), using a validation set to pick the best model. Our implementation is available at <https://github.com/yekeren/Cap2Det>.

## 4 EXPERIMENTS AND RESULTS

First, we present our experimental settings (Sec. 4.1). Second, we compare the accuracy of alternative strategies on *pseudo image-level training label extraction* from captions (Sec. 4.2). Third, we show that our approach achieves strong *detection* performance using supervision from captions (Sec. 4.3). By training on COCO captions, we achieve close to state-of-the-art results on weakly supervised detection on PASCAL, even though the supervision we leverage is weaker than competitor methods. Importantly, our text classifier trained on COCO generalizes to other datasets, and allows us to use Flickr30K and the noisier datasets MIRFlickr1M and Conceptual Captions which do not feature clean, descriptive captions. In all settings, our primary pseudo label inference method, EM+TEXTCLSF, outperforms the alternative techniques, including the EXACTMATCH baseline. Fourth, we show the improvements achieved by filtering and scoring noisy image-caption examples (Sec. 4.4). We conclude that the redundancy between image and text is key to train a successful weakly-supervised detection model. Finally, as a sanity check, we show our approach performs competitively to prior methods on the task of learning from *clean, ground-truth* image-level labels (Sec. 4.5).

Note that in all experiments, we focus on evaluating the impact of a single component of our model, focusing on the caption weighting and pseudo label inference. For most experiments, we determine success by comparing to an upper bound (e.g. ground-truth labels) and/or a lower bound (e.g. naive training label extraction, unweighted loss using captions equally). We show our methods’ advantages

persist when we replace Sec. 3.3 with alternative state-of-the-art techniques, e.g. Ren et al. [9].

### 4.1 Datasets and metrics

Our experiments involve the datasets COCO, PASCAL VOC, Flickr30K, MIRFlickr1M, and Conceptual Captions.

**PASCAL Visual Object Classes (VOC)** [60] is a standard image dataset for object class recognition. It focuses on a limited number of classes (20 objects). We use it to evaluate our learned object detection models (4,952 and 10,991 test examples in VOC07 and VOC12, respectively).

**Common Objects in Context (COCO)** [61] is a large-scale object detection, segmentation, and captioning dataset. We use: (1) its image-caption data to train the triplet and PVSE models necessary to infer image-caption homogeneity (Sec. 3.1); (2) its caption-label data to train our pseudo label inference module (Sec. 3.2) and test it on three other datasets; (3) its 118K training images, each paired with 5 captions, to train our detection model (Sec. 3.3) using 591,435 COCO [62] captions paired to the 118,287 *train2017* images. EXACTMATCH fails to extract any label for roughly 15K *train2017* instances. Since COCO is fully annotated with instance-level boxes, we use its evaluation server to measure performance of the resulting object detection models.

**Flickr30K and MIRFlickr1M.** We use the Flickr data to prove that our weakly-supervised object detection models can generalize to alternative datasets. Flickr30K [63] (30K images, each paired with 5 captions) contains crowdsourced captions. However, we also use 200K examples from the 1M noisy data in MIRFlickr1M [64] (subset for computational reasons). The dataset pairs images with user-generated content, which is not guaranteed to describe the content of the images in an exhaustive or precise manner. EXACTMATCH fails to extract labels on roughly 113K of those. We use captions from these datasets to train our detection models. Fig. 4 shows examples.

**Conceptual Captions** [65] contains 3.3M images annotated with captions. The captions are from the web, and illustrate the noisy web data environment. We conduct an experiment that tries to benefit from the large corpus via vision-language pretraining. We also directly train object detectors on Conceptual Captions. To keep results comparable to those when using COCO Captions, we extract a 118K subset from Conceptual Captions. EXACTMATCH fails to extract any label on roughly 54K samples in that subset. Conceptual Captions includes alt-text for images. Alt-texts are preprocessed with hypernymization to replace named entities (e.g. architect’s name) with object names (e.g. “person”). As can be seen from Fig. 4 (right), this alt-text is not descriptive of the images in the same way that COCO captions are descriptive. For example, the third caption mentions “person”, but this person likely refers to whoever made the decorative paint; *this person is not visible in the image*. Similarly, the “person enjoying tea” (fourth caption) is also not visible. Conversely, the persons in the first example are not explicitly mentioned, and neither is the car in the bottom example. This makes straightforward techniques for extracting labels (e.g. EXACTMATCH) prone to failure due to missed or incorrect labels.

**Evaluation protocols.** We follow the standard protocols used in VOC and COCO to fairly compare to other





Fig. 4: Example image-caption pairs from MIRFlickr1M and Conceptual Caption datasets. For MIRFlickr1M, captions and tags are written by the uploader and website users. For ConcCap, captions are parsed from the alt-text HTML attribute associated with web images. Both datasets did not crowdsource i.e. pay workers to label the images. Note how often “person” is mentioned on the right but not visible.

methods. On VOC, we report the mean Average Precision (AP) at IoU > 0.5 and we also report the per-class AP. On COCO, we report the same AP@IoU0.5 as in VOC, and the mean Average Precision across different IoU thresholds, i.e., mAP@IoU.5:.05:.95. We also include the full metrics in COCO such as AP of differently-sized objects. We submit our results to the VOC and COCO evaluation servers to get the VOC07 and COCO testing results.

#### 4.2 Accuracy of pseudo image label inference

We first test the generalizing power of our pseudo image-level label inference module (Sec. 3.2). Our results show that without using too much data, one can train a reasonably good text classifier inferring accurate image-level labels from paired captions.

In Fig. 5 we show the *precision* and *recall* of these label inference methods evaluated directly on the COCO image-level labels (5,000 examples of the *val2017* set). We observe that EM+EXTENDVOCAB, which uses the hand-crafted word-synonyms dictionary, provides the best recall (60.6%) among all methods but the worst precision of 81.1%. The word-embedding-based top-scoring matching methods of EM+GLOVE and EM+LEARNEDGLOVE provide precise predictions (84.5% and 84.7% respectively, which are the highest). However, our EM+TEXTCLSF achieves significantly improved precision compared to these. We would like to point out that while in Tab. 1 and 2, our method uses the full COCO training set (118,287 concatenated captions), it achieves very similar performance with even a small fraction of the data. With 5% of the data (6,000 caption-label pairs), the method achieves 89% precision (vs 92% precision with 100% of the data), both of which are much higher than any other baselines; recall is about 62% for both 5% and 100% training data. Thus, it is sufficient to use a

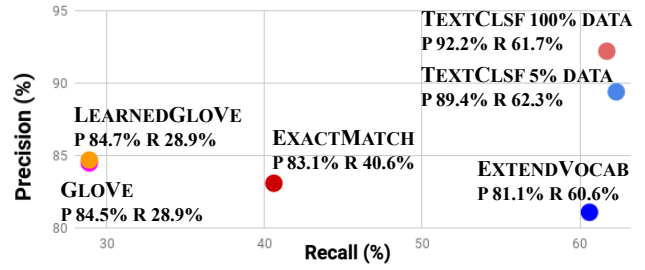


Fig. 5: Analysis of different text supervision. We compare the pseudo labels (Sec. 3.2) to COCO *val* ground-truth.

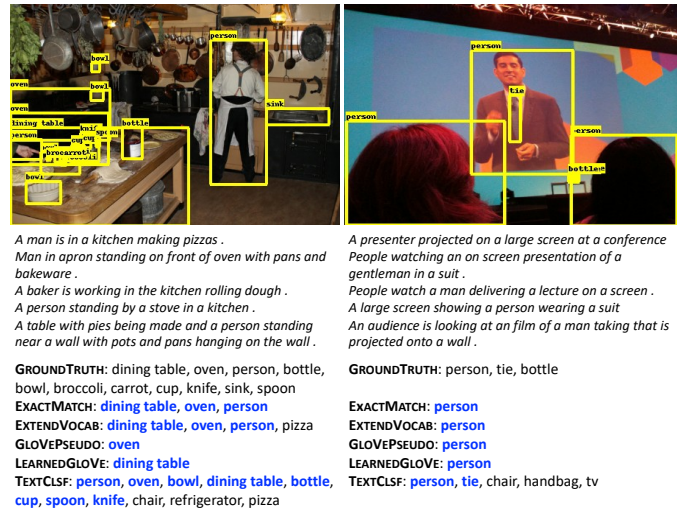


Fig. 6: Demonstration of different pseudo labels. Our method fills the gap between what is present and what is mentioned, by making inferences on the semantic level. Matches to the ground truth are shown in blue.

small portion of precise text labels to train a generalizable label inference classifier, and the knowledge can transfer to other datasets as we show in Tab. 1.

To better understand the generated labels, we show two qualitative examples in Fig. 6. The image on the right shows that our model infers “tie” from the observation of “presenter”, “conference” and “suit”, while all other methods fail to extract this object category for visual detection.

It is also interesting to measure this performance per class, as we show in Fig. 7. The lexical matching method EXACTMATCH performs similarly to EM+TEXTCLSF in terms of precision (not shown). For both EXACTMATCH and EM+TEXTCLSF, recall is very low for bottle, car, and chair, indicating these are common objects which however are *usually not mentioned in captions*. In contrast, other common objects (e.g. cat) have high recall because they are usually mentioned when present in the image. However, for classes such as boat, cow, and person, EXACTMATCH has *much lower recall rate than* EM+TEXTCLSF. We thus quantitatively explain why EM+TEXTCLSF is better than EXACTMATCH for these classes (Tab. 1 B, boat 29.9% v.s. 25.9% or 13.3% v.s. 9.6%; cow 61.2% v.s. 49.0% or 47.4% v.s. 28.0%; person 16.9% v.s. 10.4% or 10.7% v.s. 4.0%). By explicitly handling the noise in the captions (the lack of mentions of objects that do appear), we cope with the “human reporting bias” [37].

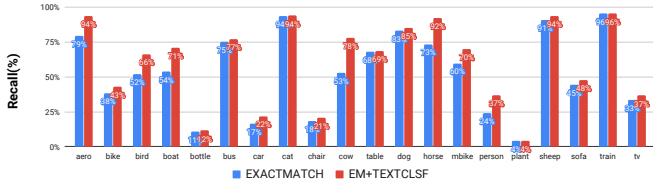


Fig. 7: Recall of PASCAL labels. We evaluate the recall of the COCO-learned text classifier, but we show only the overlapped 20 PASCAL VOC classes.

### 4.3 Comparing label inference strategies for detection

We compare our pseudo label inference module (Sec. 3.2) with alternative strategies to extract object labels from captions. All of these strategies are upper-bounded in terms of performance by using ground-truth image-level labels GT-LABEL. Note that apart from the strategy used to mine image-level labels for training, these strategies all use the same architecture and WSOD approach as our method (Sec. 3.3). We show combinations of the exact match strategy with these methods (when exact match fails), resulting in EM+GLOVE, EM+LEARNEDGLOVE, EM+EXTENDVOCAB and EM+TEXTCLSF. We examine how well these strategies leverage captions from COCO, Flickr30K, MIRFlickr1M, and Conceptual Captions for detection.

**Training with COCO captions.** Tab. 1, segments (A) and (B), show the results on PASCAL VOC 2007. At the top (A) are two upper-bound methods that train on *ground-truth* image-level labels, while methods in (B) train on labels extracted from image-level captions. EXACTMATCH (EM) performs the worst probably due to its low data utilization rate, as evidenced by the fact that all methods incorporating pseudo labels improve performance notably. Specifically, EM+GLOVE uses knowledge of the pre-trained GloVe embeddings. It alleviates the synonyms problem to a certain extent, thus it improves the mAP by 2% compared to EXACTMATCH. However, the GloVe embedding is not optimized for the specific visual-captions, resulting in noisy knowledge transformation. EM+LEARNEDGLOVE learns dataset-specific word embeddings. Its performance, as expected, is 3% better than EM+GLOVE. The strongest baseline is EM+EXTENDVOCAB, as the manually picked vocabulary list covers most frequent occurrences. However, collecting such vocabulary requires human effort, and is not a scalable and transferable strategy. Our EM+TEXTCLSF outperforms this expensive baseline, especially for categories “cat”, “cow”, “horse”, and “train”. Finally, despite the noisy supervision, our EM+TEXTCLSF almost bridges the gap to the GT-LABEL COCO upper bound in Tab. 1 (A).

For the results on COCO (Tab. 2), the gaps in performance between the different methods are smaller, but as before, our proposed EM+TEXTCLSF shows the best performance. We believe the smaller gaps are because many of the COCO objects are not described precisely via natural language, and the dataset itself is more challenging than PASCAL thus gain may be diluted by tough examples.

**Training with Flickr30K captions.** We also train our model on the Flickr30K [63] dataset, which contains 31,783 images and 158,915 descriptive captions. Training on

Flickr30K is more challenging: on one hand, it includes less data compared to COCO; on the other hand, we observe that the recall rate of the captions is only 48.9% with EXACTMATCH which means only half of the data can be matched to some class names. In Tab. 1 (D), we observe that due to the limited training size, the detection models trained on Flickr30K captions achieve weaker performance than those trained on COCO captions. However, given the “free” supervision, the 33.6% mAP is still very encouraging. Importantly, we observe that even though our text classifier is trained on COCO captions and labels, it generalizes well to Flickr30K captions, as evidenced by the gap between EM+TEXTCLSF and EM+EXTENDVOCAB.

**Results without pretraining on ImageNet.** For our experiments thus far, we pretrain our visual backbone on ImageNet. This is a realistic setting consistent with WSOD because while ImageNet contains clean labels at the image level, no labels are available at the box level. However, to reduce the potential interference from those labels, we also test in a setting where no image-level labels are used to learn the visual representations. The results are shown in Tab. 1 (C, E). We observe that the advantage of our EM+TEXTCLSF method remains. In particular, EM+TEXTCLSF achieves 96% (=19.30/20.15) of the GT-LABEL COCO performance, which is comparable to (even higher than) the 93% achieved when pretraining on ImageNet was used (Tab. 1 A, B). Further, EM+TEXTCLSF achieves a 5% gain (=19.30/18.46-1) over EXACTMATCH in Tab. 1 (C), vs 8% in (B). On Flickr30K, the gain of EM+TEXTCLSF over EXACTMATCH is 12% without pretraining (E), greater than 8% with pretraining (D).

**Training on noisier data: MIRFlickr1M and Conceptual Captions.** The ultimate goal of our work is to enable training of object detection models from widely-available language data, e.g. user-generated content or narratives. The results thus far meet some of the challenges of working with language dataset as supervision, but rely on relatively clean datasets. Thus, we next extend our evaluation to two noisier datasets, MIRFlickr1M and Conceptual Captions. Examples from these datasets are shown in Figs. 4 and 11. Due to the noise in these datasets (see Figs. 4 and 11), performance is significantly lower when training on captions from them, compared to the cleaner COCO and Flickr30K. However, we observe that our EM+TEXTCLSF still significantly improves upon the naive EXACTMATCH for both datasets, in Table 1 (F-H). In particular, **EM+TEXTCLSF improves upon EXACTMATCH by 78% on MIRFlickr1M (subset), and 27-50% on Conceptual Captions (without/with ImageNet pretraining, G-H)**. These are both much higher than the 8% gain of EM+TEXTCLSF over EXACTMATCH, in segment Tab. 1 (B). We note the difference in performance when using a clean (COCO) vs noisier dataset (ConcCap). EM+TEXTCLSF’s performance on the latter is reduced by 37% compared to the former (27.2 mAP in G vs 43.1 in B), using the same amount of captions.

Our text classifier inferred reasonable image-level labels on noisier captions, even though it was trained on much cleaner data (COCO). For example, on Conceptual Captions, it was able to infer ‘person’, ‘cup’, ‘chair’ and ‘dining table’ from the caption “front view of a young couple of college students drinking coffee and studying in a cafe”, and ‘person’, ‘tie’, ‘bottle’, ‘wine glass’ and ‘cup’ from the caption

Methods	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean
(A) Training on different datasets using ground-truth labels:																					
GT-LABEL VOC	68.7	49.7	53.3	27.6	14.1	64.3	58.1	76.0	23.6	59.8	50.7	57.4	48.1	63.0	15.5	18.4	49.7	55.0	48.4	67.8	48.5
GT-LABEL COCO	65.3	50.3	53.2	25.3	16.2	68.0	54.8	65.5	20.7	62.5	51.6	45.6	48.6	62.3	7.2	24.6	49.6	34.6	51.1	69.3	46.3
(B) Training on COCO dataset using captions:																					
EXACTMATCH (EM)	63.0	<b>50.3</b>	50.7	25.9	<b>14.1</b>	64.5	50.8	33.4	17.2	49.0	48.2	46.7	44.2	59.2	10.4	14.3	49.8	37.7	21.5	47.6	39.9
EM + GLOVE	<b>66.6</b>	43.7	53.3	29.4	13.6	65.3	<b>51.6</b>	33.7	15.6	50.7	46.6	45.4	47.6	<b>62.1</b>	8.0	<b>15.7</b>	48.6	46.3	30.6	36.4	40.5
EM + LEARNEDGLOVE	64.1	49.9	<b>58.6</b>	24.9	13.2	<b>66.9</b>	49.2	26.9	13.1	57.7	<b>52.8</b>	42.6	<b>53.2</b>	58.6	14.3	15.0	45.2	50.3	34.1	43.5	41.7
EM + EXTENDVOCAB	65.0	44.9	49.2	<b>30.6</b>	13.6	64.1	50.8	28.0	<b>17.8</b>	59.8	45.5	<b>56.1</b>	49.4	59.1	16.8	15.2	<b>51.1</b>	<b>57.8</b>	14.0	<b>61.8</b>	42.5
EM + TEXTCLSF	63.8	42.6	50.4	29.9	12.1	61.2	46.1	<b>41.6</b>	16.6	<b>61.2</b>	48.3	55.1	51.5	59.7	<b>16.9</b>	15.2	50.5	53.2	<b>38.2</b>	48.2	<b>43.1</b>
(C) Training on COCO dataset without ImageNet pretraining:																					
GT-LABEL COCO	40.7	17.2	5.7	0.8	0.0	31.9	21.4	30.6	0.3	31.0	26.0	30.5	23.9	43.8	13.8	7.0	25.5	2.7	15.3	34.9	20.15
EXACTMATCH (EM)	43.8	13.9	6.0	0.8	0.0	<b>39.6</b>	31.0	<b>34.2</b>	0.3	20.6	<b>28.8</b>	25.6	25.2	35.9	4.9	<b>3.3</b>	27.3	3.0	9.9	15.1	18.46
EM + GLOVE	44.6	13.4	<b>9.1</b>	1.4	0.1	36.3	29.7	30.7	0.2	17.2	21.2	<b>29.3</b>	25.1	36.9	5.9	2.6	6.7	2.7	8.7	16.8	16.93
EM + EXTENDVOCAB	<b>44.9</b>	12.8	8.3	1.1	0.0	35.5	<b>31.6</b>	27.0	<b>0.4</b>	<b>25.7</b>	27.1	27.1	22.0	34.8	<b>10.3</b>	0.2	<b>29.6</b>	<b>4.1</b>	<b>11.7</b>	<b>26.8</b>	19.05
EM + TEXTCLSF	44.5	<b>18.4</b>	6.7	<b>1.5</b>	<b>0.3</b>	35.6	27.6	32.2	0.2	<b>25.7</b>	28.0	28.3	<b>26.0</b>	<b>38.8</b>	6.2	0.1	28.4	2.7	9	25.8	<b>19.30</b>
(D) Training on Flickr30K dataset using captions:																					
EXACTMATCH (EM)	<b>46.6</b>	<b>42.9</b>	42.0	9.6	7.7	31.6	44.8	53.2	13.1	28.0	39.1	43.2	31.9	<b>52.5</b>	4.0	<b>5.1</b>	38.0	28.7	<b>15.8</b>	41.1	31.0
EM + EXTENDVOCAB	37.8	37.6	35.5	11.0	<b>10.3</b>	18.0	47.9	51.3	<b>17.7</b>	25.5	37.0	47.9	<b>35.2</b>	46.1	<b>15.2</b>	0.8	27.8	35.6	5.8	42.0	29.3
EM + TEXTCLSF	24.1	38.8	<b>44.5</b>	<b>13.3</b>	6.2	<b>38.9</b>	<b>49.9</b>	<b>60.4</b>	12.4	<b>47.4</b>	<b>39.2</b>	<b>59.3</b>	34.8	48.1	10.7	0.3	<b>42.4</b>	<b>39.4</b>	14.1	<b>47.3</b>	<b>33.6</b>
(E) Training on Flickr30K dataset using captions, without ImageNet pretraining:																					
EXACTMATCH (EM)	1.4	14.5	2.4	<b>0.7</b>	0.0	<b>9.4</b>	19.2	3.5	0.2	<b>3.5</b>	9.2	<b>11.0</b>	10.0	23.0	1.6	0.0	0.4	2.1	2.8	2.8	5.9
EM + EXTENDVOCAB	0.9	15.0	<b>2.9</b>	0.6	0.0	1.0	11.2	4.0	0.2	0.3	13.0	10.8	<b>10.9</b>	21.6	2.7	<b>0.1</b>	0.6	<b>3.9</b>	0.8	<b>6.6</b>	5.4
EM + TEXTCLSF	<b>1.8</b>	<b>21.0</b>	<b>1.9</b>	<b>0.7</b>	0.0	5.4	<b>20.6</b>	<b>5.7</b>	0.2	2.1	<b>12.4</b>	<b>11.0</b>	10.4	<b>24.7</b>	<b>3.4</b>	<b>0.1</b>	<b>1.0</b>	1.4	<b>3.7</b>	<b>4.1</b>	<b>6.6</b>
(F) Training on MIRFlickr1M (200k subset) using captions, without ImageNet pretraining:																					
EXACTMATCH (EM)	7.6	7.5	8.5	0.5	0.1	13.0	18.9	15.2	0.4	2.9	2.6	6.6	3.7	21.9	3.3	0.2	2.4	<b>2.3</b>	1.4	16.2	6.8
EM + TEXTCLSF	<b>31.7</b>	<b>10.5</b>	<b>8.9</b>	<b>0.8</b>	<b>0.3</b>	<b>15.5</b>	<b>28.9</b>	<b>28.1</b>	<b>0.6</b>	<b>6.3</b>	<b>6.5</b>	<b>20.5</b>	<b>12.9</b>	<b>26.6</b>	<b>15.1</b>	<b>0.6</b>	<b>4.1</b>	1.0	<b>3.5</b>	<b>19.1</b>	<b>12.1</b>
(G) Training on Conceptual Captions (118k subset):																					
EXACTMATCH (EM)	26.1	21.2	17.0	10.0	7.7	19.1	<b>31.1</b>	15.2	<b>7.0</b>	22.3	<b>35.9</b>	32.6	13.1	25.5	3.5	<b>4.6</b>	34.1	12.0	11.8	12.2	18.1
EM + TEXTCLSF	<b>60.9</b>	<b>36.5</b>	<b>35.7</b>	<b>21.4</b>	<b>8.9</b>	<b>28.3</b>	20.1	<b>26.2</b>	4.4	<b>34.8</b>	17.5	<b>41.7</b>	<b>22.8</b>	<b>51.2</b>	<b>11.3</b>	0.3	<b>42.4</b>	<b>29.3</b>	<b>21.3</b>	<b>28.4</b>	<b>27.2</b>
(H) Training on Conceptual Captions (118k subset), without ImageNet pretraining:																					
EXACTMATCH (EM)	1.1	1.6	<b>0.7</b>	0.3	0.0	3.1	<b>10.4</b>	<b>14.2</b>	0.2	0.3	1.2	7.0	0.9	2.7	1.8	0.1	<b>0.5</b>	<b>0.7</b>	0.6	<b>0.2</b>	2.38
EM + TEXTCLSF	<b>1.4</b>	1.6	0.5	0.3	0.0	<b>3.9</b>	7.0	13.5	0.2	<b>0.4</b>	<b>3.4</b>	<b>8.8</b>	<b>1.7</b>	<b>8.9</b>	<b>2.9</b>	0.1	0.3	0.6	<b>4.1</b>	0.1	<b>3.02</b>

TABLE 1: Average precision (in %) on the VOC 2007 test set (learning from COCO, Flickr30K, MIRFlickr1M, and Conceptual Captions captions). We evaluate on only the overlapping 20 VOC objects. Best method per column (except GT methods) in **bold**. Our proposed EM+TEXTCLSF achieves 93-96% of GT-LABEL COCO, and exceeds EXACTMATCH (EM) by 8% with COCO captions and 27-78% with noisy captions (MIRFlickr1M and Conceptual Captions).

Methods	Avg. Precision, IoU			Avg. Precision, Area			Avg. Recall, #Dets			Avg. Recall, Area		
	0.5:0.95	0.5	0.75	S	M	L	1	10	100	S	M	L
GT-LABEL	10.6	23.4	8.7	3.2	12.1	18.1	13.6	20.9	21.4	4.5	23.1	39.3
EXACTMATCH (EM)	8.9	19.7	7.1	2.3	10.1	16.3	<b>12.6</b>	<b>19.3</b>	<b>19.8</b>	3.4	20.3	37.4
EM + GLOVE	8.6	19.0	6.9	2.2	10.0	16.0	12.2	18.7	18.9	2.9	19.0	37.6
EM + LEARNEDGLOVE	8.9	19.7	7.2	2.5	10.4	<b>16.6</b>	12.3	19.1	19.6	<b>3.5</b>	20.0	37.7
EM + EXTENDVOCAB	8.8	19.4	7.1	2.3	10.5	16.1	12.1	19.0	19.5	3.4	20.3	37.5
EM + TEXTCLSF	<b>9.1</b>	<b>20.2</b>	<b>7.3</b>	<b>2.6</b>	<b>10.8</b>	<b>16.6</b>	12.5	<b>19.3</b>	<b>19.8</b>	<b>3.5</b>	<b>20.6</b>	<b>37.8</b>

TABLE 2: COCO test-dev results (learning from COCO captions), measured by COCO eval server. Best method in **bold**. Our EM+TEXTCLSF achieves 86% of the GT-LABEL performance, and improves upon EXACTMATCH (EM) by 2.5%.

Methods	Avg. Precision, IoU	
	0.5:0.95	0.5
GT-LABEL COCO	11.2	22.8
EXACT MATCH (EM)	10.0	21.1
EM + TEXTCLSF	10.4	22.0

TABLE 3: Evaluation of our pseudo label inference, using Ren et al. [9] as our WSOD method, on COCO2017-val.

“group of business people raising a toast with champagne at office”. On MIRFlickr1M, it inferred ‘tie’ from the caption “For G. aaron kilt wedding whitby dunsleyhall”, and ‘knife’, ‘bowl’, ‘broccoli’, ‘carrot’ from the caption “Roasted Veggies! cauliflower asparagus limes roasted food cooking kitchen

vegetables dinner supertime eating healthy explore”. Thus, EM+TEXTCLSF can infer labels that are implied but not stated, and leverage them for training.

**Contribution of pseudo label inference using alternative WSOD method.** For all of our experiments thus far, we have used OICR [8] as our WSOD method. In other words, once pseudo labels for the training set are inferred at the image level by one of the method alternatives, these labels are used to train a WSOD model, using OICR. Here, we experiment with replacing OICR with a more recent WSOD technique, namely Ren et al. [9]. We compare the performance of two label inference techniques, EXACTMATCH and EM+TEXTCLSF, against using ground-truth labels at the image level (GT-LABEL). We show the

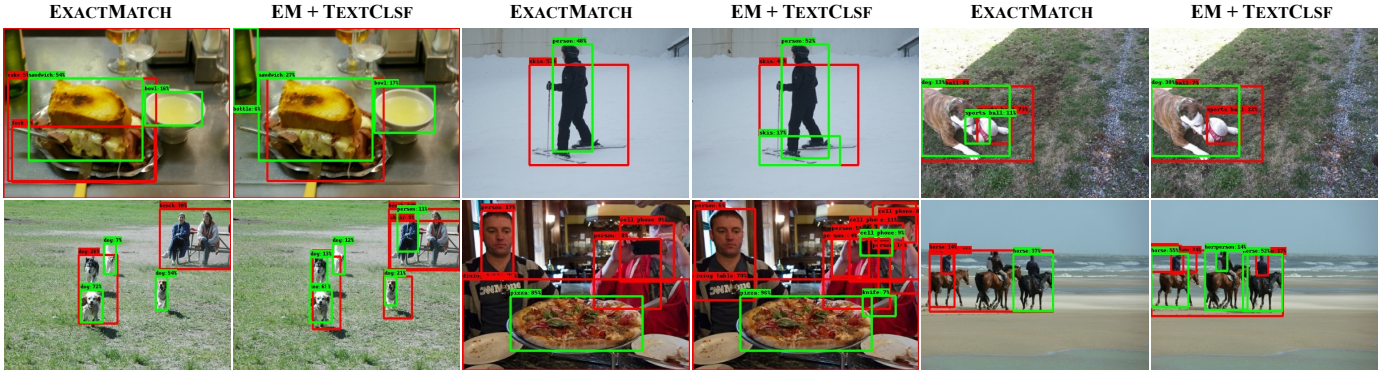


Fig. 8: Visualization of our Cap2Det model results on COCO *val* set. We show boxes with confidence scores > 5%. Green boxes denote correct detection results ( $IoU > 0.5$ ) while red boxes indicate incorrect ones. Best viewed with 300% zoom-in.

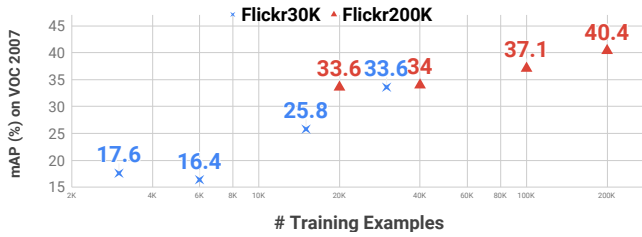


Fig. 9: Data vs. Performance. Our text classifier learned on COCO generalized well on Flickr30K and the noisier Flickr200K data (subset of MIRFlickr1M) formed by user-generated content tags.

results in Tab. 3, where all methods rely on Ren et al.’s WSOD method. We see that our proposed EM+TEXTCLSF achieves 96% of the GT performance (using IoU 0.5), and a 4% boost over EXACTMATCH. These results are even better than Tab. 2 where our method achieved 86% of the ground-truth performance, and 2.5% gain over EXACTMATCH.

**Data vs. performance.** We show the potential of our model using Flickr30K and MIRFlickr1M [64]. For the latter, we concatenate the title and all user-generated content tags to form caption annotation. We then use our text classifier learned on COCO to rule out examples unlikely to mention our target classes. This filtering results in a dataset with around 20% of the original data, and we refer to it as Flickr200K. We use 10%, 20%, 50%, 100% data from both datasets, and report average precision on VOC 2007. We see from Fig. 9 that as training data increases, mAP increases accordingly. To estimate model potential, we fit a square root function to the rightmost four points in the figure and use it to estimate 54.4 mAP at 1 million samples.

**Qualitative results on COCO.** We provide qualitative examples on the COCO *val* set, in Fig. 6 and Fig. 8. We compare EXACTMATCH and our EM+TEXTCLSF side-by-side in Fig. 8. Our proposed method EM+TEXTCLSF provides better detection results than the baseline EXACTMATCH. Thus, we conclude that it has squeezed more useful and precise information than the EXACTMATCH baseline.

#### 4.4 Impact of filtering noisy captions

Not all nouns in captions are object words, and not all object words are mentioned in the caption: for example, an image with caption “guests are sitting at a table during a wedding” will show object “table”, no object “wedding”, and additional objects (e.g. plates). The purpose of Sec. 3.1 is to determine which image-caption pairs contain significant alignment that would allow us to extract quality pseudo training labels at the image level. As discussed in Sec. 1, this is an alternative to our EM+TEXTCLSF (Sec. 3.2). The motivation for this alternative is that it does not require labels for training the text classifier. Thus, **our goal is for homogeneity scoring to improve our detection results with EXACTMATCH, EM+EXTENDVOCAB or EM+GLOVE, not with EM+TEXTCLSF.** Results with weighting on top of EM+TEXTCLSF are still shown, but marked with gray shading, in Tabs. 5, 6, 7, and 8.

**Ranking captions by potential purity of objects mentions.** We test to what extent homogeneity can be used to estimate if a caption contains clean labels. We use both captions and labels in COCO to compute a ground-truth *ranking of images* by the overlap between caption and object label words. We use the EXACTMATCH and EXTENDEDMATCH methods to compute the overlap, and rank images by precision (fraction of caption nouns that are also object labels). We then use our method to also compute an approximate ranking of images by the  $\alpha^{HOM}$  scores, calculated individually for image and paired captions (Eqs. 3), as well as their difference. Finally, we compute the correlation between ground-truth rankings and our methods’ rankings, using Kendall’s  $\tau$  and Spearman’s  $\rho$ . If the correlation is high, our method is a good indicator of how well-aligned a caption and object labels are, thus how likely it is that weakly-supervised detection will succeed if we extract labels from this particular caption.

We also compute approximate ranks using two baselines. HESSEL [40] computes visual concreteness for a word using the purity of images co-occurring with this word. We rank images by the average concreteness of nouns in their paired captions. ALIKHANI [39] collects annotations for the type of relation between an image and its caption, including “visible” (the most direct relation) and other less direct ones (e.g. “story”). We train a classifier using image and captions

GT Ranking	Pred Ranking	Image		Label	
		$\rho$	$\tau$	$\rho$	$\tau$
EXACTMATCH	HESSEL	<b>0.100</b>	<b>0.067</b>	0.293	0.182
	ALIKHANI	-0.081	-0.054	0.348	0.229
	HOM-IMAGE	0.088	0.058	<b>0.382</b>	<b>0.230</b>
	HOM-TEXT	0.058	0.038	0.326	0.212
	HOM-DIFF	-0.003	-0.003	0.003	-0.001
EXTENDMATCH	HESSEL	0.043	0.029	0.229	0.142
	ALIKHANI	-0.009	-0.007	0.376	<b>0.251</b>
	HOM-IMAGE	0.130	0.086	<b>0.388</b>	0.225
	HOM-TEXT	<b>0.179</b>	<b>0.120</b>	0.341	0.223
	HOM-DIFF	0.177	0.118	0.053	0.037

TABLE 4: Ranking images and object categories by how well captions overlap with the true image-level labels. Higher correlations are better. The best number per GT ranking is **bolded**. HOM-IMAGE is the best performer overall.

and the annotations from [39]. It outputs the probability of a caption being visible, and we sort images by their captions’ average visibility.

We also *rank object categories*: by their concreteness scores for HESSEL, by average visibility of images containing this category for ALIKHANI, and by average  $\alpha^{HOM}$  values of images containing this category for our methods. In ground-truth rankings EXACTMATCH and EXTENDEDMATCH, accuracy of category prediction (i.e. fraction of images where category  $c$  is correctly predicted) is used.

We show the results in Tab. 4. HOM-IMAGE is the best performer overall, outperforming the alternatives in 3 cases (vs 2 for HOM-TEXT, 2 for HESSEL, and 1 for ALIKHANI). We observe that the image ranking obtained by HESSEL is more correlated with the ground-truth ranking acquired by EXACTMATCH. However, it fails to sustain its performance when ground-truth image ranking is EXTENDMATCH, probably due to not being able to capture visual concreteness for synonyms of object categories. In this setting, HOM-TEXT is the best. In the label ranking task, HOM-IMAGE is the most correlated ranking overall.

**Detection results using image-caption filtering.** We use a limited 30,000 image-caption subset from COCO *train2017* split for training, assuming a setting of restricted computation resources and training time. We keep the most useful examples while removing the others. We use the metric of *homogeneity* (Sec. 3.1.2) to measure the image-caption relevance. The higher this metric, the better alignment between the both modalities, and more likely the captions describe the visual objects in detail. HOM-IMAGE and HOM-TEXT use  $\alpha_{(I)}^{HOM}$  and  $\alpha_{(T)}^{HOM}$  to filter examples, respectively. We use random sampling of 30K examples as a baseline.

Tab. 5 shows the results. We see the performances of EXTENDVOCAB and GLOVE are improved using the filtered training data. If we use a random selection of 30K examples, the performances are 36.7% and 38.6% respectively. Using *image homogeneity score* (HOM-IMAGE,  $\alpha_{(I)}^{HOM}$ ) for filtering improved these methods by 9% (40.1% v.s. 36.7%) and 7% (41.3% v.s. 38.6%). The TEXTCLSF column is shaded in gray because we do not expect a boost in performance from filtering. TEXTCLSF had already explained the gap between the image and text thus is not sensitive to the improved filtered training data. However, this text classifier requires a small number of ground-truth labels. In contrast, HOM-IMAGE and HOM-TEXT with GLOVE achieve competitive

	Label inference		
Im-cap scoring	EXTENDVOCAB	GLOVE	TEXTCLSF
Random (30K)	36.7	38.6	40.4
HOM-IMAGE (30K)	40.1	<b>41.3</b>	<b>40.5</b>
HOM-TEXT (30K)	<b>40.4</b>	40.8	39.9

TABLE 5: Comparing the filtering strategies with the random sampling baseline, using AP (in %) on VOC 2007 test. Both filtering mechanisms improve results under for EXTENDVOCAB and GLOVE pseudo label inference. Gray cells indicate we do not intend or expect improvement compared to no filtering.

	Label inference		
Im-cap scoring	EXTENDVOCAB	GLOVE	TEXTCLSF
No weighting (118K)	42.5	40.5	<b>43.1</b>
HOM-IMAGE Weighting (118K)	<b>43.5</b>	<b>42.6</b>	42.2

TABLE 6: Evaluating the caption weighting strategy, using AP (in %) on VOC 2007 test.

results to TEXTCLSF with Random, but do not require any labels. Thus, homogeneity can be used to determine which captions provide strong supervision for object detection, without the need for *any* ground-truth labels. We omit HESSEL, ALIKHANI and HOM-DIFF in Tab. 5 because their overall performance in Tab. 4 does not exceed HOM-IMAGE, while ALIKHANI also requires image-caption relation labels.

**Results using image-caption weighting.** One weakness of the filtering approach is that it requires a hard cutoff of the dataset examples. In comparison, weighting applies a soft “cutoff” to the data. It never drops data, thus is data-efficient. We use  $\alpha_{(I)}^{HOM}$  as the per-example weighting factor  $\gamma$  to weigh the pseudo image-level labels extracted from different image-caption pairs in Eq. 9.

Comparing Tab. 6 to Tab. 5, we see that even with a good filtering strategy, e.g., HOM-IMAGE Filtering (30K), using 30K “clean” training examples is still inferior than training on the full COCO dataset (No weighting 118K), for two of the three label inference methods (columns). However, if we apply the HOM-IMAGE *weighting* on the loss, the performance (HOM-IMAGE Weighting 118K) is improved (43.5% v.s. 42.5%, 42.6% v.s. 40.5%, i.e. 2-5% gain), except for TEXTCLSF which requires annotations. **This is an important finding with important ramifications for multimodal learning.** Approaches to learn visual representations have benefited greatly from widely available web/video data, and **our method suggests how the useful signal and noise in such data can be distinguished to boost the quality of the learned representations**, without requiring annotations. We note that EXTENDVOCAB with weighting is comparable (only 1% better) to TEXTCLSF with no weighting.

**Results using large corpus pretraining.** We conduct a study on the impact of vision-language pretraining, which has gained popularity for both VQA and cross-modal retrieval [25], [26]. However, no previous study had shown where such pre-trained models can also improve fully- or weakly-supervised object detection (except Zareian et al. [66] which use bounding boxes for some object categories). Our pre-training pipeline is shown in Fig. 10. We evenly distribute anchor boxes (different in scales and locations) across the image and treat each anchored visual region as

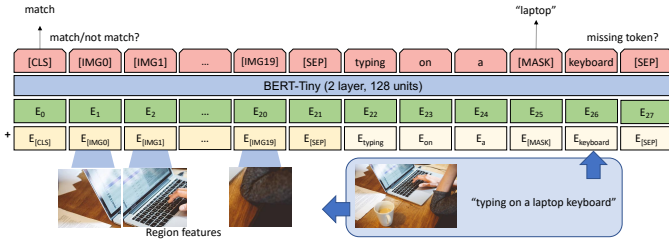


Fig. 10: Vision-language pretraining for weakly-supervised object detection. The inputs to the BERT-Tiny model are positional embeddings ( $E_0, E_1, \dots$ ) and token embeddings (proposal features ( $E_{[IMG0]}, E_{[IMG1]}, \dots$ ) and word embeddings ( $E_{typing}, E_{on}, \dots$ )).

Label inference	EXTENDVOCAB	GLOVE	TEXTCLS F
Im-cap scoring			
No weighting, no pretraining	42.5	40.5	43.1
+Weighting	43.5	42.6	42.2
+Pretraining	43.9	41.5	43.6
+Pretraining, +weighting	42.5	42.0	42.6

TABLE 7: The impact of large corpus pretraining, measuring AP (in %) on VOC 2007 test. Both pretraining and weighting help, but weighting uses no external data.

a visual token. We then concatenate the visual tokens to the caption tokens, forming a sequence: “[CLS] [v1] [v2] ... [SEP] [t1] [t2] ... [SEP]”, where the “[v]” and “[t]” are the visual and textual tokens, respectively. “[SEP]” is a special token to separate sequences, and “[CLS]” is the special classification token. We use the FastRCNN features and word embeddings to represent visual and text token embeddings, respectively. Then, we contextualize the sequence features using two self-attention layers (each with 128 hidden units). Finally, we add a linear classification layer on top of the “[CLS]” representation to predict a 0/1 value denoting if the image implies the caption’s semantic meaning (matching). Besides the image-caption matching modeling, we also process the masked language modeling optimization. We randomly (with a probability of 15%) replace a text token with “[MASK]” and require the model to reconstruct the token, given the visual and text contexts. The visual model is the same as Sec. 3.4, while we use BERT-Tiny to initialize the text token embeddings and the self-attention layers. We use the Adam optimizer [54] with a batch size of 5, a learning rate of  $1e-5$ , a weight decay of  $1e-8$ , and we trim the maximum gradient norm to 1.0. Based on the above pre-training setting, we trained for 600K steps, roughly 1 epoch on the Conceptual Captions dataset. The training costs around 44 hours on 5 GeForce GTX1080Ti GPUs, using Tensorflow distributed training [58]. After pretraining the model using Conceptual Captions, we process the weakly supervised object detection training using the COCO images and texts and evaluate on the VOC07.

Tab. 7 shows the results. We see that pretraining improved the baseline by 3% (43.9% v.s. 42.5%), 2% (41.5% v.s. 40.5%), 1% (43.6% v.s. 43.1%). In the only pseudo label inference setting that requires neither a hand-defined vocabulary of synonyms, nor object labels, namely **GLOVE**, **our weighting achieves stronger results than pretraining**. For EXTENDVOCAB, pretraining and weighting achieve

Label inference	EXACTMATCH	TEXTCLS F
Im-cap scoring		
(A) ImageNet pretraining in detection:		
No weighting	18.1	27.2
HOM-IMAGE Weighting (Triplet w/ IN PT)	24.7	18.8
HOM-IMAGE Weighting (Triplet)	24.0	25.5
(B) No ImageNet pretraining in detection:		
No weighting	2.38	3.02
HOM-IMAGE Weighting (Triplet)	2.99	2.75

TABLE 8: Weighting on Conceptual Captions, using AP (in %) on VOC 2007 test set. Triplet models are trained on the COCO dataset. Weighting improves results across dataset boundaries. Gray cells indicate we do not intend or expect improvement compared to no weighting.

comparable results. Pretraining can be seen as a state-of-the-art method akin to Zareian et al. [66], as a way to use vision-language data, which could be alternative to our EM+TEXTCLS F. However, pretraining uses the external and large 3.3M Conceptual Captions dataset, while our EM+TEXTCLS F only uses the 118K COCO captions. Given inconclusive gains from pretraining over weighting (sometimes worse, sometimes better), and its large cost, it is not warranted in our setting. Applying both pretraining and weighting did not further boost results.

#### Generalization of weighting on Conceptual Captions.

The results described thus far in this subsection all apply filtering or weighting over COCO captions. The weighting model was trained on COCO, so in this part, we test the generalization of the weighting model, by applying it on the Conceptual Captions subset described in Sec. 4.1. Note that when applying weights to the captions, the gradient magnitude is reduced. Thus, rather than retuning the learning rate, we apply scaling to the weights such that the sum of weights over all samples remains the same as in the unweighted version. We include the result in Tab. 8. Part (A) uses ImageNet pretraining in extracting visual features for detection (Sec. 3.3), while part (B) does not. The second and third rows in (A) differ in the use of an ImageNet-trained backbone for the image-text model (trained with triplet loss). As before, our focus is on weighting improving the performance of the EXACTMATCH method. We observe that our weighting model does generalize to Conceptual Captions. On EXACTMATCH, we achieve a 33-37% boost in performance when using weighting (compared to no weighting) in (A), and 26% in (B). Thus, on Conceptual Captions, the impact of weighting is even more significant than in Tab. 6. TEXTCLS F is 10% better than EXACTMATCH with weighting ( $=27.2/24.7-1$ ) but requires some labels.

**Simplified computation of weights.** The results in Tab. 8 replace the use of PVSE [21] to compute an image-text joint embedding, with a much simpler model trained with triplet loss and no attention. Thus, the gains from weighting that we achieve are not due to the complexity of PVSE.

**Image-caption pairs with high/low scores.** In Fig. 11, we show examples of image-caption pairs from each of COCO, Conceptual Captions, and MIRFlickr1M, that achieved high or low homogeneity scores. Among COCO samples, images with low homogeneity scores are usually more complex than ones with high scores, and mentioned concepts may be abstract (“symbols”, “stop and go”, “display”, “fabrics”).

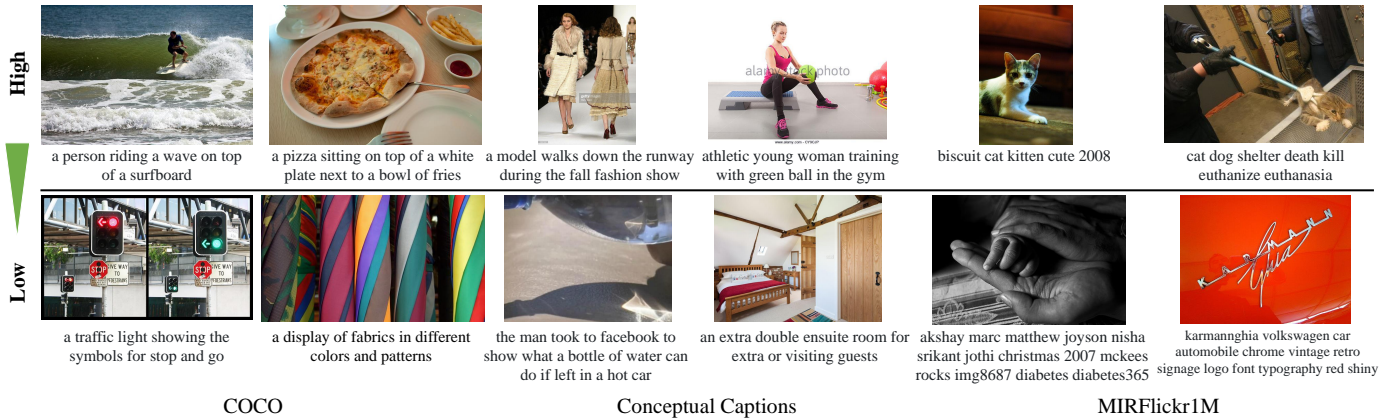


Fig. 11: Image-caption pairs with high homogeneity scores on the top, and low scores on the bottom.

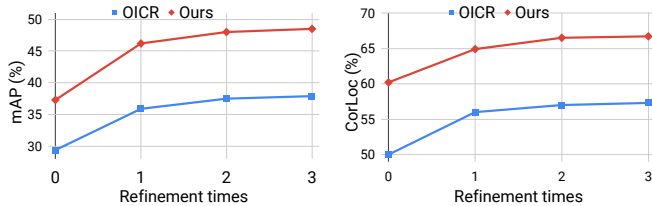


Fig. 12: Analysis of our basic network and OICR components on VOC 2007. Comparison of the performance of our model and OICR VGG\_M after iterative refinement.

On Conceptual Captions, high-scoring pairs describe the content in a literal fashion, and many objects are mentioned. In low-scoring pairs, mentioned objects are not present (e.g. “man”) or present objects are not mentioned (e.g. the bed). On MIRFlickr1M, the high-scoring images mention object categories, while the low-scoring ones are significantly more abstract or non-object-like. Thus, we qualitatively showed that the homogeneity scores measure the relevance and redundancy between the image and text modalities, in terms of ability to extract object labels. Homogeneity helps to rule out less useful examples to better train a detector.

#### 4.5 Verifying WSOD with ground-truth image labels

We finally show the performance of our method in the classic WSOD setting where *clean image-level labels for training are available* and do not need to be inferred. These results validate the method component in Sec. 3.3. Our goal is not to exceed the very latest WSOD models, but to perform on par with recent ones. Note that all methods tested in this section, including ours, use *ground-truth image-level labels*, but differ in terms of architectures and WSOD techniques. We refer to our method here as a WSOD variant, to distinguish it from the main focus of our work on using language supervision (which is not utilized here). Also note that the multi-scale training and test time augmentation mentioned in our Sec. 3.4 are widely adopted in WSOD. We verified that all baseline methods in this section use them. Thus, our comparisons to the SOTA methods in this section are fair.

**Results on PASCAL VOC.** For each image, we extract object categories from all the ground-truth bounding boxes, and only keep these *image-level* labels for training, discarding box information. For VOC 2007 and 2012, we train on

5,011 and 11,540 *trainval* images respectively and evaluate on 4,952 and 10,991 *test* images. We report mean Average Precision (mAP) at IoU > 0.5, and compare against multiple strong WSOD baselines, in Tab. 9. The WSOD variant of our model performs on par with or better than the baselines on both VOC 2007 and 2012.

**Effects of the basic network and OICR.** The performance gain of our model comes from two aspects: (1) a more advanced detection model backbone architecture and (2) the online instance classifier refinement (OICR). Fig. 12 shows the performance of the WSOD variant of our method and that of Tang [8] (OICR VGG\_M), both refining for 0, 1, 2, 3 times. With no (0) refinement, our basic network architecture outperforms the VGG\_M backbone of Tang *et al.* by 27% in mAP. But the basic architecture improvement is not sufficient to achieve top results. If we use OICR to refine the models 1, 2, or 3 times, we gain 24%, 29%, and 30% respectively while Tang achieve smaller improvement (22%, 28%, and 29% gains).

**Results on COCO.** We train the WSOD variant of our model on the 118,287 *train2017* images, using the image-level ground truth labels. We report mAP at IoU=.50:.05:.95 and mAP@0.5, on the 20,288 *test-dev2017* images. We compare to a representative fully-supervised detection model [1]; “Faster Inception-V2” [59] which is our WSOD variant’s supervised detection counterpart (using bounding-box annotations), and a recent WSOD model, PCL-OB-G Ens + FRCNN [16]. As demonstrated in Tab. 10, our model outperforms the WSOD model by 15% in terms of mAP, but as expected, the gap between with the supervised method is still large due to the disparate supervision strength.

## 5 CONCLUSIONS

We showed how we can successfully leverage naturally arising, weak supervision in the form of captions. We amplify the signal that captions provide by learning to bridge the gap between what human annotators mention, and what is present in the image. We also learn how to weigh the contribution of different captions as supervision, based on the expected alignment between the image and caption. In the future, we will extend our method to incorporate raw supervision in the form of spoken descriptions in video.

Methods	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mean
VOC 2007 results:																					
OICR VGG16 [8]	58.0	62.4	31.1	19.4	13.0	<b>65.1</b>	62.2	28.4	24.8	44.7	30.6	25.3	37.8	65.5	15.7	<b>24.1</b>	41.7	46.9	<b>64.3</b>	62.6	41.2
PCL-OB-G VGG16 [16]	54.4	<b>69.0</b>	39.3	19.2	<b>15.7</b>	62.9	<b>64.4</b>	30.0	<b>25.1</b>	52.5	44.4	19.6	39.3	<b>67.7</b>	<b>17.8</b>	22.9	46.6	<b>57.5</b>	58.6	63.0	43.5
TS <sup>2</sup> C [7]	59.3	57.5	43.7	27.3	13.5	63.9	61.7	59.9	24.1	46.9	36.7	45.6	39.9	62.6	10.3	23.6	41.7	52.4	58.7	56.6	44.3
OICR Ens.+FRCNN [8]	65.5	67.2	47.2	21.6	22.1	68.0	68.5	35.9	5.7	63.1	49.5	30.3	64.7	66.1	13.0	25.6	50.0	57.1	60.2	59.0	47.0
PCL-OB-G Ens.+FRCNN [16]	63.2	69.9	47.9	22.6	27.3	71.0	69.1	49.6	12.0	60.1	51.5	37.3	63.3	63.9	15.8	23.6	48.8	55.3	61.2	62.1	48.8
<b>Ours</b>	<b>68.7</b>	49.7	<b>53.3</b>	<b>27.6</b>	14.1	64.3	58.1	<b>76.0</b>	23.6	<b>59.8</b>	<b>50.7</b>	<b>57.4</b>	<b>48.1</b>	63.0	15.5	18.4	<b>49.7</b>	55.0	48.4	<b>67.8</b>	<b>48.5</b>
VOC 2012 results:																					
OICR VGG16 [8]	67.7	61.2	41.5	25.6	22.2	54.6	49.7	25.4	19.9	47.0	18.1	26.0	38.9	67.7	2.0	22.6	41.1	34.3	37.9	55.3	37.9
PCL-OB-G VGG16 [16]	58.2	<b>66.0</b>	41.8	24.8	<b>27.2</b>	55.7	<b>55.2</b>	28.5	16.6	<b>51.0</b>	17.5	28.6	<b>49.7</b>	<b>70.5</b>	7.1	<b>25.7</b>	47.5	36.6	44.1	<b>59.2</b>	40.6
TS <sup>2</sup> C [7]	67.4	57.0	37.7	23.7	15.2	<b>56.9</b>	49.1	64.8	15.1	39.4	19.3	<b>48.4</b>	44.5	67.2	2.1	23.3	35.1	40.2	<b>46.6</b>	45.8	40.0
OICR Ens.+FRCNN [8]	71.4	69.4	55.1	29.8	28.1	55.0	57.9	24.4	17.2	59.1	21.8	26.6	57.8	71.3	1.0	23.1	52.7	37.5	33.5	56.6	42.5
PCL-OB-G Ens.+FRCNN [16]	69.0	71.3	56.1	30.3	27.3	55.2	57.6	30.1	8.6	56.6	18.4	43.9	64.6	71.8	7.5	23.0	46.0	44.1	42.6	58.8	44.2
<b>Ours</b>	<b>74.2</b>	49.8	<b>56.0</b>	<b>32.5</b>	22.0	55.1	49.8	<b>73.4</b>	<b>20.4</b>	47.8	<b>32.0</b>	39.7	48.0	62.6	<b>8.6</b>	23.7	<b>52.1</b>	<b>52.5</b>	42.9	59.1	<b>45.1</b>

TABLE 9: Average precision (in %) on the Pascal VOC test set using *ground-truth* image-level labels. The top shows VOC 2007 and the bottom shows VOC 2012 results. The best single model is in **bold**, and best ensemble in *italics*.

Methods	Avg. Precision, IoU	
	0.5:0.95	0.5
Faster RCNN [1]	21.9	42.7
Faster Inception-V2 [59]	28.0	-
PCL-OB-G VGG16 [16]	8.5	19.4
PCL-OB-G Ens.+FRCNN [16]	9.2	19.6
<b>Ours</b>	<b>10.6</b>	<b>23.4</b>

TABLE 10: COCO detection using *ground-truth* image labels, with supervised models at the top, best WSOD in **bold**.

## ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation under Grant Numbers 1566270, 1718262 and 2046853. It was also supported by Google Faculty Research Awards and an NVIDIA hardware grant.

## REFERENCES

- [1] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [2] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [3] K. Ye, A. Kovashka, M. Sandler, M. Zhu, A. Howard, and M. Forni, "Spotpatch: Parameter-efficient transfer learning for mobile object detection," in *Proceedings of the Asian Conference on Computer Vision (ACCV)*, November 2020.
- [4] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Domain adaptive faster r-cnn for object detection in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [5] C. Thomas and A. Kovashka, "Artistic object recognition by unsupervised style adaptation," in *Asian Conference on Computer Vision (ACCV)*. Springer, 2018, pp. 460–476.
- [6] H. Bilen and A. Vedaldi, "Weakly supervised deep detection networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [7] Y. Wei, Z. Shen, B. Cheng, H. Shi, J. Xiong, J. Feng, and T. Huang, "Ts2c: Tight box mining with surrounding segmentation context for weakly supervised object detection," in *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
- [8] P. Tang, X. Wang, X. Bai, and W. Liu, "Multiple instance detection network with online instance classifier refinement," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] Z. Ren, Z. Yu, X. Yang, M.-Y. Liu, Y. J. Lee, A. G. Schwing, and J. Kautz, "Instance-aware, context-focused, and memory-efficient weakly supervised object detection," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [10] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [11] F. Wan, C. Liu, W. Ke, X. Ji, J. Jiao, and Q. Ye, "C-mil: Continuation multiple instance learning for weakly supervised object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [12] Z. Zeng, B. Liu, J. Fu, H. Chao, and L. Zhang, "Wsd2: Learning bottom-up and top-down objectness distillation for weakly-supervised object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [13] Z. Ren, Z. Yu, X. Yang, M.-Y. Liu, Y. J. Lee, A. G. Schwing, and J. Kautz, "Instance-aware, context-focused, and memory-efficient weakly supervised object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [14] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Is object localization for free? - weakly-supervised learning with convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [15] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [16] P. Tang, X. Wang, S. Bai, W. Shen, X. Bai, W. Liu, and A. L. Yuille, "Pcl: Proposal cluster learning for weakly supervised object detection," *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [17] A. Diba, V. Sharma, A. Pazandeh, H. Pirsiavash, and L. Van Gool, "Weakly supervised cascaded convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [18] F. Wan, P. Wei, J. Jiao, Z. Han, and Q. Ye, "Min-entropy latent model for weakly supervised object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [19] S. K. Divvala, A. Farhadi, and C. Guestrin, "Learning everything about anything: Webly-supervised visual concept learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3270–3277.
- [20] M. Carvalho, R. Cadène, D. Picard, L. Soulier, N. Thome, and M. Cord, "Cross-modal retrieval in the cooking context: Learning semantic text-image embeddings," in *International ACM SIGIR Conference on Research & Development in Information Retrieval*, 2018, pp. 35–44.
- [21] Y. Song and M. Soleymani, "Polysemous visual-semantic embedding for cross-modal retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1979–1988.



- [22] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128–3137.
- [23] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*, 2015, pp. 2048–2057.
- [24] H. Xu and K. Saenko, "Ask, attend and answer: Exploring question-guided spatial attention for visual question answering," in *European Conference on Computer Vision*. Springer, 2016, pp. 451–466.
- [25] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Advances in Neural Information Processing Systems*, 2019, pp. 13–23.
- [26] H. Tan and M. Bansal, "LXMERT: Learning cross-modality encoder representations from transformers," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Nov. 2019, pp. 5100–5111.
- [27] A. Rohrbach, M. Rohrbach, R. Hu, T. Darrell, and B. Schiele, "Grounding of textual phrases in images by reconstruction," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [28] D. Surís, D. Epstein, H. Ji, S.-F. Chang, and C. Vondrick, "Learning to learn words from visual scenes," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [29] C. Thomas and A. Kovashka, "Preserving semantic neighborhoods for robust cross-modal retrieval," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [30] L. Gomez, Y. Patel, M. Rusiñol, D. Karatzas, and C. Jawahar, "Self-supervised learning of visual features through embedding images into text topic spaces," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4230–4239.
- [31] A. Miech, J.-B. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman, "End-to-end learning of visual representations from uncurated instructional videos," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [32] G. Bertasius and L. Torresani, "Cobe: Contextualized object embeddings from narrated instructional video," in *Advances in Neural Information Processing Systems*, 2020.
- [33] K. Desai and J. Johnson, "Virtex: Learning visual representations from textual annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [34] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [35] K. Chen, H. Song, C. Change Loy, and D. Lin, "Discover and learn new objects from documentaries," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [36] K. Ye, M. Zhang, A. Kovashka, W. Li, D. Qin, and J. Bernt, "Cap2det: Learning to amplify weak caption supervision for object detection," in *International Conference on Computer Vision (ICCV)*, 2019.
- [37] I. Misra, C. Lawrence Zitnick, M. Mitchell, and R. Girshick, "Seeing through the human reporting bias: Visual classifiers from noisy human-centric labels," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [38] M. Zhang, R. Hwa, and A. Kovashka, "Equal but not the same: Understanding the implicit relationship between persuasive images and text," in *British Machine Vision Conference*, 2018.
- [39] M. Alikhani, P. Sharma, S. Li, R. Soricut, and M. Stone, "Clue: Cross-modal coherence modeling for caption generation," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020.
- [40] J. Hessel, D. Mimno, and L. Lee, "Quantifying the visual concreteness of words and topics in multimodal datasets," in *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2018.
- [41] C. Thomas and A. Kovashka, "Emphasizing complementary samples for non-literal cross-modal retrieval," in *Multimodal Learning and Applications (CVPR Workshop)*, 2022.
- [42] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [43] C. Thomas and A. Kovashka, "Preserving semantic neighborhoods for robust cross-modal retrieval," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [44] Q. Le and T. Mikolov, "Distributed representations of sentences and documents," in *International conference on machine learning*, 2014, pp. 1188–1196.
- [45] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2019.
- [46] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [47] J. Lu, J. Yang, D. Batra, and D. Parikh, "Neural baby talk," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [48] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2009.
- [49] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders, "Selective search for object recognition," *International Journal of Computer Vision*, vol. 104, no. 2, pp. 154–171, 2013.
- [50] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [51] K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," in *Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation (SSST-8)*, 2014, 2014.
- [52] R. Rehůřek and P. Sojka, "Software Framework for Topic Modelling with Large Corpora," in *Proceedings of the LREC Workshop on New Challenges for NLP Frameworks*. Valletta, Malta: ELRA, May 2010, pp. 45–50, <http://is.muni.cz/publication/884893/en>.
- [53] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010, pp. 249–256.
- [54] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [55] Y. A. Malkov and D. A. Yashunin, "Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs," *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2016.
- [56] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of machine learning research*, vol. 12, no. 7, 2011.
- [57] G. Bradschi, "The opencv library," *Dr. Dobb's Journal of Software Tools*, 2000.
- [58] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, "Tensorflow: A system for large-scale machine learning," in *Proceedings of the USENIX Conference on Operating Systems Design and Implementation*, 2016, pp. 265–283.
- [59] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy, "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [60] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [61] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

- [62] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick, "Microsoft coco captions: Data collection and evaluation server," *arXiv preprint arXiv:1504.00325*, 2015.
- [63] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *Transactions of the Association for Computational Linguistics*, vol. 2, pp. 67–78, 2014.
- [64] M. J. Huiskes, B. Thomee, and M. S. Lew, "New trends and ideas in visual concept detection: the mir flickr retrieval evaluation initiative," in *Proceedings of the international conference on Multimedia information retrieval*. ACM, 2010, pp. 527–536.
- [65] P. Sharma, N. Ding, S. Goodman, and R. Soricut, "Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018, pp. 2556–2565.
- [66] A. Zareian, K. D. Rosa, D. H. Hu, and S.-F. Chang, "Open-vocabulary object detection using captions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 393–14 402.



**Mingda Zhang** obtained his Ph.D. in Computer Science from the University of Pittsburgh in December 2021. He is now a software engineer at Google Cloud AI and Industry Solutions, New York. His research focus is on the intersection of computer vision and natural language processing. He has completed research internships in Google AI, Seattle. Before coming to Pitt, he obtained his B.Sc. in Chemical Biology in 2013 from Peking University, China.



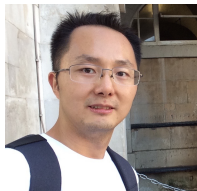
**Adriana Kovashka** is an Associate Professor in Computer Science at the University of Pittsburgh. She got her PhD in August 2014 at the University of Texas at Austin. Her work has been published in CVPR, ICCV, ECCV, TPAMI, IJCV, NeurIPS, AAAI and ACL, and has been funded by NSF, Google, Amazon and Adobe. She is the recipient of a NSF CAREER award in 2021. She served as Area Chair for CVPR 2018-2021, and will serve as ICCV 2025 Program Co-Chair.



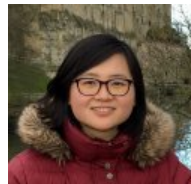
**Mesut Erhan Unal** is a fourth-year Ph.D. student in Computer Science at the University of Pittsburgh (Pitt). The goal of his research is to better understand, interpret and manipulate visual data using natural language supervision and contextual information. He holds a B.Sc. degree in Computer Science and Engineering from Hacettepe University, Turkey. Before starting his Ph.D. at Pitt, he worked for Jotform as a software engineer for two years.



**Wei Li** is a research scientist at NewsBreak Seattle. His current research focus is on video understanding, editing and synthesis. His work has been published at CVPR, ICCV, PAMI and etc. He got his master degree from CUHK and bachelor degree from Tsinghua University.



**Keren Ye** obtained his Ph.D. in Computer Science from the University of Pittsburgh in August 2021. He is now a Senior Applied Research Scientist at Cruise, San Francisco. His interests lie in object detection, multi-modal learning, and knowledge representation. Before studying at Pitt, he worked as a software engineer at Baidu Inc. for 5 years. He got both of his Bachelors and Masters degrees (2004-2011) from Beihang University, China.



**Danfeng Qin** is a software engineer at Google AI. Her current work focuses on label efficient learning with web data. She got her PhD degree from the Computer Vision Lab in ETH Zurich.



**Christopher Thomas** is an Assistant Professor in Computer Science at Virginia Tech. Previously he was a postdoctoral researcher at Columbia University, mentored by Prof. Shih-Fu Chang and working as part of the DARPA SemaFor program. He received his Ph.D. degree in Computer Science from the University of Pittsburgh in August 2020. His research interests include semantic and pragmatic image understanding, structured, weak, and unsupervised learning, vision and language, and image generation.



**Jesse Berent** is a research scientist and tech lead manager at Google AI Zurich. His group focuses on image and video analysis, digital ink recognition and machine learning. Prior to joining Google in 2009, Jesse was a post-doctoral researcher at Imperial College London focusing on multi-view image analysis. He obtained his PhD in 2008 in Communications and Signal Processing at Imperial College London and his Masters in Microengineering from the EPFL, Switzerland.