

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Unified Visual-Semantic Embeddings: Bridging Vision and Language with Structured Meaning Representations

Anonymous CVPR submission

Paper ID 4705

Abstract

We propose Unified Visual-Semantic Embeddings (VSE) for learning a joint space for scene representation and textual semantics. It unifies the embeddings of concepts at different levels: objects, attributes, relations and full scenes. We view the sentential semantics as a combination of different semantic components such as object or relational descriptors, and align their embeddings with different regions of a scene. A contrastive learning approach is proposed for the effective learning of such fine-grained alignment from only image-caption pairs. We also present a simple yet effective approach that enforces the coverage of caption embeddings on the semantic components that appear in the sentence. We demonstrate that the Unified VSE outperforms other baselines on cross-modal retrieval tasks and the enforcement of the semantic coverage improves models' robustness in defending text-domain adversarial attacks. Moreover, such robustness empowers the use of visual cues to accurately resolve word dependencies in novel sentences.

1. Introduction

We study the problem of establishing accurate and generalizable alignments between visual concepts and textual semantics efficiently, based upon rich but few, paired but noisy, or even biased visual-textual inputs (*e.g.*, image-caption pairs). Consider the image-caption pair A shown in Fig. 1: “A white clock on the wall is above a wooden table”. The alignments are formed at multiple levels: This short sentence can be decomposed into a rich set of semantic components [3]: objects (clock, table and wall) and relations (clock above table, and clock on wall). These components are linked with different parts of the scene.

This motivates our work to introduce *Unified Visual-Semantic Embeddings* (*Unified VSE* for short) Shown in Fig. 2, Unified VSE bridges visual and textual representation in a joint embedding space that unifies the embeddings for objects (noun phrases vs. visual objects), attributes (prenominal phrases vs. visual attributes), relations (verbs or prepo-

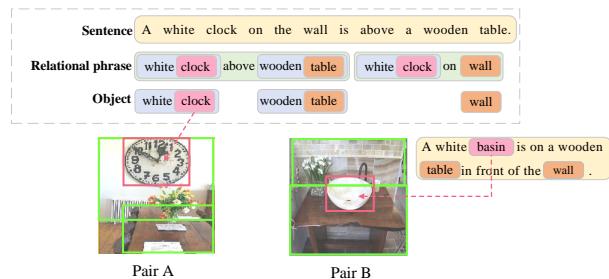


Figure 1. Two exemplar image-caption pairs. Humans are able to establish accurate and generalizable alignments between vision and language, at different levels: objects, relations and full sentences. Pair A and B form a pair of contrastive example for the concepts *clock* and *basin*.

sitional phrases vs. visual relations) and scenes (sentence vs. image).

There are two major challenges in establishing such a factorized alignment. First, the link between the textual description of an object and the corresponding image region is ambiguous: A visual scene consists of multiple objects, and thus it is unclear to the learner which object should be aligned with the description. Second, it could be problematic to directly learn a neural network that combines various semantic components in a caption and form an encoding for the full sentence, with the training objective to maximize the cross-modal retrieval performance in the training set (*e.g.*, in [50, 30, 41]). As reported by [41], because of the inevitable bias in the dataset (*e.g.*, two objects may co-occur with each other in most cases, see the table and the wall in Fig. 1 as an example), the learned sentence encoders usually pay attention to only part of the sentence. As a result, they are vulnerable to text-domain adversarial attacks: Adversarial captions constructed from original captions by adding small perturbations (*e.g.*, by changing *wall* to be *shelf*) can easily fool the model [41, 40].

We resolve the aforementioned challenges by a natural combination of two ideas: *cross-situational learning* and the enforcement of *semantic coverage* that regularizes the encoder. Cross-situational learning, or learning from contrastive examples [12], uses contrastive examples in the dataset to resolve the referential ambiguity of objects: Look-

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

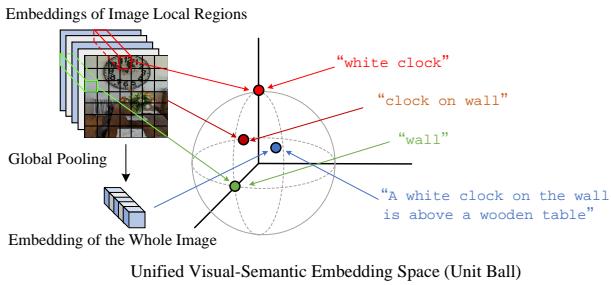


Figure 2. We build a visual-semantic embedding space, which unifies the embeddings for objects, attributes, relations and full scenes.

ing at both Pair A and B in Fig. 1, we know that *Clock* should refer to an object that occurs only in scene A but not B. Meanwhile, to alleviate the biases of datasets such as object co-occurrence, we present an effective approach that enforces the *semantic coverage*: The meaning of a caption is a composition of all semantic components in the sentence [3]. Reflectively, the embedding of a caption should have a coverage of all semantic components, while changing any of them should affect the global caption embedding.

Conceptually and empirically, Unified VSE makes the following three contributions.

First, the explicit factorization of the visual-semantic embedding space enables us to build a fine-grained correspondence between visual and textual data, which further benefits a set of downstream visual-textual tasks. We achieve this through a contrastive example mining technique that uniformly applies to different semantic components, in contrast to the sentence or image-level contrastive samples used by existing visual-semantic learning [50, 30, 11]. Unified VSE consistently outperforms pre-existing approaches on a diverse set of retrieval-based tasks.

Second, we propose a caption encoder that ensures a coverage of all semantic components appeared in the sentence. We show that this regularization helps our model to learn a robust semantic representation for captions. It effectively defends adversarial attacks on the text domain.

Furthermore, we show how our learned embeddings can provide visual cues to assist the parsing of novel sentences, including determining content word dependencies and labelling semantic roles for certain verbs. It ends up that our model can build reliable connections between vision and language using given semantic cues and in return, bootstrap the acquisition of language.

2. Related work

Visual semantic embedding. Visual semantic embedding [13] is a common technique for learning a joint representation of vision and language. The embedding space empowers a set of cross-modal tasks such as image captioning [44, 49, 8] and visual question answering [4, 48].

A fundamental technique proposed in [13] for aligning two modalities is to use the pairwise ranking to learn a dis-

tance metric from similar and dissimilar cross-modal pairs [45, 36, 23, 9, 28, 24]. As a representative, VSE++ [11] uses the online hard negative mining (OHEM) strategy [42] for data sampling and shows the performance gain. VSE-C [41], based on VSE++, enhances the robustness of the learned visual-semantic embeddings by incorporating rule-generated textual adversarial samples as hard negatives during training. In this paper, we present a contrastive learning approach based on semantic components.

There are multiple VSE approaches that also use linguistically-aware techniques for the sentence encoding and learning. Hierarchical multimodal LSTM (HM-LSTM) [34] and [47], as two examples, both leverage the constituency parsing tree. Multimodal-CNN (m-CNN) [30] and CSE [50] apply convolutional neural networks to the caption and extract the a hierarchical representation of sentences. Our model differs with them in two aspects. First, Unified VSE is built upon a factorized semantic space instead of the syntactic knowledge. Second, we employ a contrastive example mining approach that uniformly applies to different semantic components. It substantially improves the learned embeddings, while the related works use only sentence-level contrastive examples.

The learning of object-level alignment in unified VSE is also related to [19, 21, 37], where the authors incorporate pre-trained object detectors for the semantic alignment. [10] propose a selective pooling technique for the aggregation of object features. Compared with them, Unified VSE presents a more general approach that embeds concepts of different levels, while still requiring no extra supervisions.

Structured representation for vision and language. We connect visual and textual representations in a structured embedding space. The design of its structure is partially motivated by the papers on relational visual representations (scene graphs) [29, 18, 17], where a scene is represented by a set of objects and their relations. Compared with them, our model does not rely on labelled graphs during training.

Researchers have designed various types of representations [5, 33] as well as different models [26, 51] for translating natural language sentences into structured representations. In this paper, we present how the usage of such semantic parsing into visual-semantic embedding facilitates the learning of the embedding space. Moreover, we present how the learned VSE can, in return, helps the parser to resolve parsing ambiguities using visual cues.

3. Unified Visual-Semantic Embeddings

We now describe the overall architecture and training paradigm for the proposed *Unified Visual-Semantic Embeddings*. Shown in Fig. 3, given an image-caption pair, we first parse the caption into a structured meaning representation, composed by a set of semantic components: object nouns, prenominal modifiers, and relational dependencies.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

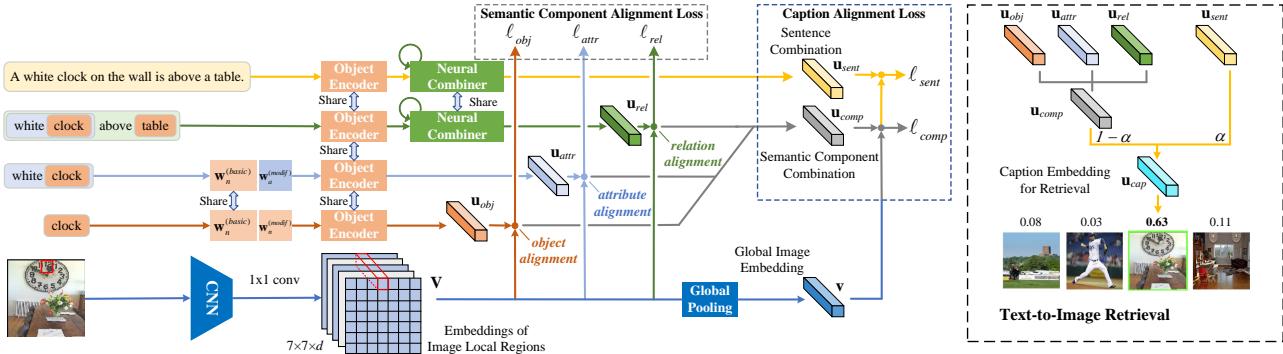


Figure 3. **Left:** the architecture of Unified VSE. The semantic component alignment is learned from contrastive examples sampled from factorized semantic space. The model also learns a caption encoder that combines the semantic components and aligns the caption with the corresponding image. **Right:** An exemplar computation graph for retrieving images from texts. The presence of \mathbf{u}_{comp} in the caption encoding enforces the coverage of all semantic components. See Sec. 3.2 for details.

We encode different types of semantic components with type-specific encoders. A caption encoder combines the embedding of the semantic components into a caption semantic embedding. Jointly, we encode images with a convolutional neural network (CNN) into the same, *unified VSE* space. The distance between the image embedding and the sentential embedding measures the semantic similarity between the image and the caption.

We employ a multi-task learning approach for the joint learning of embeddings for semantic components (as the “basis” of the VSE space) as well as the caption encoder (as the combiner of semantic components).

3.1. Visual-Semantic Embedding: A Revisit

We begin the section with an introduction to the two-stream VSE approach. It jointly learns the embedding spaces of two modalities: vision and language, and aligns them using parallel image-text pairs (*e.g.*, image and captions from the MS-COCO dataset [27]).

Let $\mathbf{v} \in \mathbb{R}^d$ be the representation of the image and $\mathbf{u} \in \mathbb{R}^d$ be the representation of a caption matching this image, both encoded by neural modules. To archive the alignment, a bidirectional margin-based ranking loss has been widely applied [11, 50, 15]. Formally, for an image (caption) embedding \mathbf{v} (\mathbf{u}), denote the embedding of its matched caption (image) as \mathbf{u}^+ (\mathbf{v}^+). A negative (unmatched) caption (image) is sampled whose embedding is denoted as \mathbf{u}^- (\mathbf{v}^-). We define the bidirectional ranking loss ℓ_{sent} between captions and images as:

$$\begin{aligned} \ell_{sent} = & \sum_{\mathbf{u}} F_{\mathbf{v}^-} (|\delta + s(\mathbf{u}, \mathbf{v}^-) - s(\mathbf{u}, \mathbf{v}^+)|_+) \\ & + \sum_{\mathbf{v}} F_{\mathbf{u}^-} (|\delta + s(\mathbf{u}^-, \mathbf{v}) - s(\mathbf{u}^+, \mathbf{v})|_+) \end{aligned} \quad (1)$$

, where δ is a predefined margin, $|x|_+ = \max(x, 0)$ is the traditional ranking loss and $F_x(\cdot) = \max_x(\cdot)$ denotes the hard negative mining strategy [11, 42]. $s(\cdot, \cdot)$ is a similarity function between two embeddings and is usually implemented as cosine similarity [11, 41, 50].

3.2. Semantic Encodings

The encoding of a caption is made up of three steps. As an example, consider the caption shown in Fig. 3, “A white clock on the wall is above a wooden table”. 1) We extract a structured meaning representation as a collection of three types of semantic components: object (clock, wall, table), attribute-object dependencies (white clock, wooden table) and relational dependencies (clock above table, clock on wall). 2) We encode each component as well as the full sentence with type-specific encoders into the unified VSE space. 3) We represent the embedding of the caption by combining the semantic components.

Semantic parsing. We implement a semantic parser ¹ of image captions based on [39]. Given the input sentence, the parser first performs a syntactic dependency parsing. A set of rules is applied to the dependency tree and extracts object entities appeared in the sentence, adjectives that modify the object nouns, subjects/objects of the verbs and prepositional phrases. For simplicity, we consider only single-word nouns for objects and single-word adjectives for object attributes. **Encoding objects and attributes.** We use an unified object encoder ϕ for nouns and adjective-noun pairs. For each word w in the vocabulary, we initialize a basic semantic embedding $\mathbf{w}_n^{(basic)} \in \mathbb{R}^{d_{basic}}$ and a modifier semantic embedding $\mathbf{w}_n^{(modif)} \in \mathbb{R}^{d_{modif}}$.

For a single noun word w_n (*e.g.*, `clock`), we define its embedding \mathbf{w}_n as $\mathbf{w}_n^{(basic)} \oplus \mathbf{w}_n^{(modif)}$, where \oplus means the concatenation of vectors. For an (adjective, noun) pair (w_a, w_n) (*e.g.*, (white, `clock`)), its embedding $\mathbf{w}_{a,n}$ is defined as $\mathbf{w}_n^{(basic)} \oplus \mathbf{w}_a^{(modif)}$ where $\mathbf{w}_a^{(modif)}$ encodes the attribute information. In implementation, the basic semantic embedding is initialized from GloVe [35]. The modifier semantic embeddings (both $\mathbf{w}_n^{(modif)}$ and $\mathbf{w}_a^{(modif)}$) are randomly initialized and jointly learned. $\mathbf{w}_n^{(modif)}$ can be regarded as an intrinsic modifier for each nouns.

¹Code will be released upon acceptance.

324 To fuse the embeddings of basic and modifier semantics,
 325 we employ a gated fusion function:
 326

$$\begin{aligned}\phi(\mathbf{w}_n) &= \text{Norm}(\sigma(\mathbf{W}_1\mathbf{w}_n + \mathbf{b}_1)) \tanh(\mathbf{W}_2\mathbf{w}_n + \mathbf{b}_2), \\ \phi(\mathbf{w}_{a,n}) &= \text{Norm}(\sigma(\mathbf{W}_1\mathbf{w}_{a,n} + \mathbf{b}_1)) \tanh(\mathbf{W}_2\mathbf{w}_{a,n} + \mathbf{b}_2).\end{aligned}$$

330 Throughout the text, σ denotes the sigmoid function: $\sigma(x) = 1/(1 + \exp(-x))$, and Norm denotes the L2 normalization,
 331 *i.e.*, $\text{Norm}(\mathbf{w}) = \mathbf{w}/\|\mathbf{w}\|_2$. One may interpret ϕ as a GRU
 332 cell [7] taking no historical state.

333 **Encoding relations and full sentence.** Since relations and
 334 sentences are the composed based on objects, we encode
 335 them with a neural combiner ψ , which takes the embeddings
 336 of word-level semantics encoded by ϕ as input. In practice,
 337 we implement ψ as an uni-directional GRU [7], and pick the
 338 L2-normalized last state as the output.

339 To obtain a visual-semantic embedding for a relational
 340 triple (w_s, w_r, w_o) (*e.g.*, (clock, above, table)), we
 341 first extract the word embeddings for the subject, relational
 342 word and the object using ϕ . We then feed the encoded
 343 word embeddings in the same order into ψ and takes the
 344 L2-normalized last state of the GRU cell. Mathematically,
 345 $\mathbf{u}_{rel} = \psi(w_s, w_r, w_o) = \psi(\{\phi(\mathbf{w}_s), \phi(\mathbf{w}_r), \phi(\mathbf{w}_o)\})$.

346 The embedding of a sentence \mathbf{u}_{sent} is computed over the
 347 word sequence w_1, w_2, \dots, w_k of the caption:

$$\mathbf{u}_{sent} = \psi(\{\phi(\mathbf{w}_1), \phi(\mathbf{w}_2), \dots, \phi(\mathbf{w}_k)\}),$$

350 where for any word x , $\phi(\mathbf{w}_x) = \phi(\mathbf{w}_x^{(basic)} \oplus \mathbf{w}_x^{(modif)})$

351 Note that we share the weights of the encoders ψ and ϕ
 352 among the encoding processes of all semantic levels. This
 353 allows our encoders of various types of components to boot-
 354 strap the learning of each other.

355 **Combining all of the components.** A straight-forward im-
 356 plementation of the caption encoder is to directly use the
 357 sentence embedding \mathbf{u}_{sent} , as it has already combined the
 358 semantics of components in a contextually-weighted manner
 359 [25]. However, it has been revealed in [41] that such com-
 360 bination is vulnerable to adversarial attacks: Because of the
 361 biases in the dataset, the combiner ψ usually focuses on only
 362 a small set of semantic components appeared in the caption.

363 We alleviate such biases by enforcing the coverage
 364 of the semantic components appeared in the sentence.
 365 Specifically, to form the caption embedding \mathbf{u}_{cap} , the
 366 sentence embedding \mathbf{u}_{sent} is combined with an explicit
 367 bag-of-components embedding \mathbf{u}_{comp} , as illustrated in
 368 Fig. 3 (right). Mathematically, we define \mathbf{u}_{comp} as an
 369 unweighted aggregation of all components in the sentence:

$$\mathbf{u}_{comp} = \text{Norm}\left(\sum_{obj} \mathbf{u}_{obj} + \sum_{attr} \mathbf{u}_{attr} + \sum_{rel} \mathbf{u}_{rel}\right),$$

370 and encode the caption as: $\mathbf{u}_{cap} = \alpha \mathbf{u}_{sent} + (1 - \alpha) \mathbf{u}_{comp}$,
 371 where $0 \leq \alpha \leq 1$ is a scalar weight. The presence of \mathbf{u}_{comp}
 372 disallows the ignorance of any of the components in the final
 373 caption embedding \mathbf{u}_{cap} .

3.3. Image Encodings

374 We use CNN to encode the input RGB image into the
 375 unified VSE space. Specifically, we choose a ResNet-152
 376 model [14] pretrained on ImageNet [38] as the image en-
 377 coder. We apply a layer of 1×1 convolution on top of the
 378 last convolutional layer (*i.e.*, conv5_3) and obtain a conve-
 379 lutional feature map of shape $7 \times 7 \times d$ for each image. d
 380 denotes the dimension of the unified VSE space.

381 The feature map, denoted as $\mathbf{V} \in \mathbb{R}^{7 \times 7 \times d}$, can be view
 382 as the embeddings of 7×7 local regions in the image. The
 383 embedding \mathbf{v} for the whole image is defined as the aggrega-
 384 tion of the embeddings at all regions through a global spatial
 385 pooling operator.

3.4. Learning Paradigm

386 In this section, we present how to align vision and lan-
 387 guage into the unified space using contrastive learning on
 388 different semantic levels. The training pipeline is illustrated
 389 in Fig. 3. We start from the generation of contrastive exampls
 390 for different semantic components.

391 **Negative example sampling.** It has been discussed in [41]
 392 that to explore a large compositional space of semantics, di-
 393 rectly sampling negative captions from a human-built dataset
 394 (*e.g.*, MS-COCO captions) is not sufficient. In this paper, in-
 395 stead of manually define rules that augment the training data
 396 as in [41], we address this problem by sampling contrastive
 397 negative examples in the explicitly factorized semantic space.
 398 The generation does not require manually labelled data, and
 399 can be easily applied to any datasets. For a specific caption,
 400 we generate the following four types of contrastive negative
 401 samples.

- 402 • **Nouns.** We sample negative noun words from all nouns
 403 that do not appear in the caption. ²
- 404 • **Attribute-noun pairs.** We sample negative pairs by
 405 randomly substituting the adjective by another adjective
 406 or substituting the noun.
- 407 • **Relational triples.** We sample negative triples by ran-
 408 domly substituting the subject, or the relation, or the
 409 object. Moreover, we also sample the whole relational
 410 triples of captions in the dataset which describe other
 411 images, as the negative triples.
- 412 • **Sentences.** We sample negative sentences from the
 413 whole dataset. Meanwhile, following [13, 11], we also
 414 sample negative images from the whole dataset as con-
 415 trastive images.

416 The key motivation behind our visual-semantic alignment
 417 is that: an object appears in a local region of the image, while
 418 the aggregation of all local regions should be aligned with
 419 the full semantics of a caption.

420 **Local region-level alignment.** In detail, we propose a
 421 relevance-weighted alignment mechanism for linking textual

422 ²For the MS-COCO dataset, in all 5 captions associated with the same
 423 image. This also applies to other components.

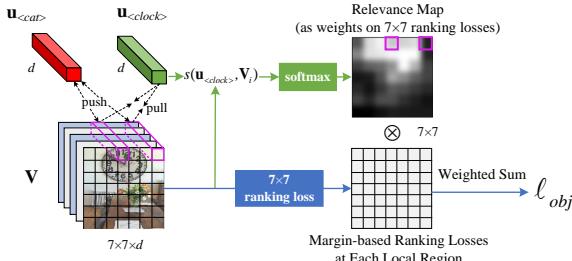


Figure 4. An illustration of our relevance-weighted alignment mechanism. The relevance map shows the similarity of each region with the object embedding $u_{<clock>}$. We weight the alignment loss with the map to reinforce the correspondence between the $u_{<clock>}$ and its matched region.

object descriptors and local image regions. As shown in Fig. 4, consider the embedding of a positive textual object descriptor u_o^+ , a negative textual object descriptor u_o^- and the set image local region embeddings V_i where $i \in 7 \times 7$ extracted from the image. We generate a relevance map $M \in \mathbb{R}^{7 \times 7}$ with $M_i, i \in 7 \times 7$ representing the relevance between u_o^+ and V_i , computed as as Eq. (2). We compute the loss for noun and (adjective, noun) pairs by:

$$M_i = \frac{\exp(s(u_o^+, V_i))}{\sum_j \exp(s(u_o^+, V_j))} \quad (2)$$

$$\ell_{obj} = \sum_{i \in 7 \times 7} (M_i \cdot |\delta + s(u_o^-, V_i) - s(u_o^+, V_i)|_+) \quad (3)$$

The intuition behind the definition is that, we explicitly try to align the embedding at each image region with u_o^+ . The losses are weighted by the matching score, thus reinforce the correspondence between u_o^+ and the matched region. This technique is related to multi-instance learning [46].

Global image-level alignment. For relational triples u_{rel} , semantic components aggregations u_{comp} and sentences u_{sent} , their semantics usually cover multiple objects. Thus, we align them with the full image embedding v via bidirectional ranking losses as Eq. (1)³. The alignment loss is denoted as ℓ_{rel} , ℓ_{comp} and ℓ_{sent} , respectively.

We want to highlight that, during training, we separately align the two type of semantic representations of the caption, *i.e.*, u_{sent} and u_{comp} , with the image. This differs from the inference-time computation of the caption. Recall that α can be viewed as a factor that balances the training objective and the enforcement of semantic coverage. This allows us to flexibly adjust α during inference.

3.5. Implementation details

We use $d = 1024$ as the dimension of the unified VSE space like [11, 41, 50]. We train the model by minimizing the alignment losses in a multi-task learning way.

$$\ell = \ell_{sent} + \eta_c \ell_{comp} + \eta_o \ell_{obj} + \eta_a \ell_{attr} + \eta_r \ell_{rel} \quad (4)$$

³Only textual negative samples are used for ℓ_{rel} .

In the first 2 epochs, we set η_c , η_o and η_a to 0.5 and η_r to 0 for learning single-object level representations. Then we turn up η_r to 1.0 to make the model learn relational semantics. To make the comparison with related works fair, we always fix the weights of the ResNet. We use the Adam [22] optimizer with learning rate at 0.001. For model details, please refer to our supplementary material.

4. Experiments

We evaluate our model on the MS-COCO [27] dataset. It contains 82,783 training images with each image annotated by 5 captions. We use the common 1K validation and test split from [19]. We also report the performance on a 5K test split for comparison with [50, 11, 43].

We begin this section with the evaluation of traditional cross-modal retrieval. Next, we validate the effectiveness of enforcing the semantic coverage of caption embeddings by comparing models on cross-modal retrieval tasks with adversarial examples. We then propose a unified text-to-image retrieval task to support the contrastive learning on various semantic components. We end this section with an application of using visual cues to facilitate the semantic parsing of novel sentences. Due to the limitation of the text length, for mode details on data processing, metrics and model implementation, we refer the readers to our supplementary material.

4.1. Overall Evaluation on Cross-Modal Retrieval.

We first show the performance of image-to-sentence and sentence-to-image retrieval tasks to evaluate learned visual-semantic embeddings. We report the R@1 (recall@1), R@5, R@10, and the median retrieval rank as in [11, 41, 50, 15]. To summarize the performance, we compute $rsum$ as the summation of R@1, R@5, and R@10.

Shown in Table 1, Unified VSE outperforms other baselines with various model architecture and training techniques [11, 50, 28, 41, 15]. This validates the effectiveness learning visual-semantic embeddings in the explicitly factorized visual-semantic embedding space. We also include the results under more challenging 5K test split. The gap between Unified VSE and other models gets further enlarged across all metrics.

4.2. Retrieval under text-domain adversarial attack

Recent works [41, 40] have raised their concerns on the robustness of the learned visual-semantic embeddings. They show that existing models are vulnerable to text-domain adversarial attacks (*i.e.*, using adversarial captions) and can be easily fooled. This is closely related to the bias in small datasets over a large, compositional semantic space [41]. To prove the robustness of the learned unified VSE, we further conduct experiments on the image-to-sentence retrieval task with text-domain adversarial attacks. Following [41], we

432	486
433	487
434	488
435	489
436	490
437	491
438	492
439	493
440	494
441	495
442	496
443	497
444	498
445	499
446	500
447	501
448	502
449	503
450	504
451	505
452	506
453	507
454	508
455	509
456	510
457	511
458	512
459	513
460	514
461	515
462	516
463	517
464	518
465	519
466	520
467	521
468	522
469	523
470	524
471	525
472	526
473	527
474	528
475	529
476	530
477	531
478	532
479	533
480	534
481	535
482	536
483	537
484	538
485	539

540	541	Task	Image-to-sentence Retrieval				Sentence-to-image Retrieval				594
			Metric	R@1	R@5	R@10	Med. r	R@1	R@5	R@10	Med. r
1K testing split (5,000 captions)											595
m-RNN [32]	41.0	73.0	83.5	2	29.0	42.2	77.0	3	345.7		596
DVSA [20]	38.4	69.9	80.5	1	27.4	60.2	74.8	3	351.2		597
MNLM [24]	43.4	75.7	85.8	-	31.0	66.7	79.9	-	382.5		598
m-CNN [30]	42.8	73.1	84.1	3	32.6	68.6	82.8	3	384.0		599
HM-LSTM[34]	43.9	-	87.8	2	36.1	-	86.7	3	-		600
Order-embedding [43]	46.7	-	88.9	2	37.9	-	85.9	2	-		601
VSE-C(open-source)[41][1]	48.0	81.0	89.2	2	39.7	72.9	83.2	2	414		602
DeepSP[45]	50.1	79.7	89.2	-	39.6	75.2	86.9	-	420.7		603
2WayNet [9]	55.8	75.2	-	-	39.7	63.3	-	-	-		604
sm-LSTM [15]	53.2	83.1	91.5	1	40.7	75.8	87.4	2	431.8		605
RRF-Net[28]	56.4	85.3	91.5	-	43.9	78.1	88.6	-	443.8		606
VSE++(open-source)[11][2]	57.7	86.0	94.0	1	42.8	77.2	87.4	2	445.1		607
CSE[50]	56.3	84.4	92.2	1	45.7	81.2	90.6	2	450.4		608
U-VSE (Ours)	62.2	88.4	94.6	1	47.9	82.3	91.3	2	466.7		609
5K testing split (25,000 captions)											610
Order-embedding [43]	23.3	-	65.0	5	18.0	-	57.6	7	-		611
VSE-C(open-source) [11][1]	22.3	51.1	65.1	5	18.7	43.8	56.7	7	257.7		612
CSE[50]	27.9	57.1	70.4	4	22.2	50.2	64.4	5	292.2		613
VSE++(open-source) [11][2]	31.7	60.9	72.7	3	22.1	49.0	62.7	6	299.1		614
U-VSE (Ours)	36.4	66.4	78.3	3	25.6	53.4	66.6	5	326.6		615

Table 1. Results of cross-modal retrieval task on MS-COCO dataset (1K and 5K testing split). All listed baselines and our models fix weights of the image encoders. For fair comparison, we do not include [10] and [16] that finetunes the image encoder or adds extra training data.

		Object attack				Attribute attack				Relation attack				total sum
		Metric	R@1	R@5	R@10	rsum	R@1	R@5	R@10	rsum	R@1	R@5	R@10	rsum
VSE++	32.3	69.6	81.4	183.3	19.8	59.4	76	155.2	26.1	66.8	78.7	171.6	510.1	617
VSE-C	41.1	76	85.6	202.7	26.7	61	74.3	162	35.5	71.7	81.5	188.1	552.8	618
U-VSE ($\mathbf{u}_{sent} + \mathbf{u}_{comp}$)	45.5	77.6	86.5	219.8	34.6	70.9	81.7	187.2	37.9	75.2	84.8	197.9	604.9	619
U-VSE (\mathbf{u}_{sent})	40.7	76.4	85.5	202.6	30	70.5	80.6	181.1	32.6	72.6	83.5	188.7	572.4	620
U-VSE ($\mathbf{u}_{sent} + \mathbf{u}_{obj}$)	42.9	77.2	85.6	216.4	30.1	69	79.8	178.9	34	71.2	83.6	188.8	584.1	621
U-VSE ($\mathbf{u}_{sent} + \mathbf{u}_{attr}$)	40.1	73.9	83.3	209.4	37.4	72	81.9	191.3	30.5	70	81.9	182.4	583.1	622
U-VSE ($\mathbf{u}_{sent} + \mathbf{u}_{rel}$)	45.4	77.1	85.5	218.2	29.2	68.1	78.5	175.8	42.8	77.5	85.6	205.9	599.9	623

Table 2. Results on image-to-sentence retrieval task with text-domain adversarial attacks. For each caption, we generate 5 adversarial fake captions which do not match the images. Thus, the models need to retrieve 5 positive captions from 30,000 candidate captions.

first design several types of adversarial captions by adding perturbations to existing captions.

1. **Object attack:** Randomly replace / append by an irrelevant one in the original caption.
2. **Attribute attack:** Randomly replace / add an irrelevant attribute modifier for one object in the original caption.
3. **Relational attack:** 1) Randomly replace the subject/relation/object word by an irrelevant one. 2) Randomly select an entity as a subject/object and add an irrelevant relational word and object/subject.

We include VSE++ and VSE-C as the baselines and show the results in Table 2 where different columns represent different types of attacks. We also visualize some examples for a qualitative view of the performance in Fig. 5.

VSE++ performs worst as it is only optimized for the retrieval performance on the dataset. Its sentence encoder is insensitive to a small perturbation in the text. VSE-C explicitly generates the adversarial captions based on human-designed rules as hard negative examples during training,

which makes it relatively robust to those adversarial attacks. Unified VSE shows strong robustness across all types of adversarial attacks and outperforms all baselines.

It is worth noting that VSE-C shows inferior performances in the normal retrieval tasks without adversarial captions (see Table 1), even compared with VSE++. Considering that VSE-C shares the exactly the same model architecture as VSE++, we can conclude that directly adding adversarial captions during training, although improves models' robustness, may sacrifice the performance on other tasks. In contrast, the ability of Unified VSE to defend adversarial texts comes almost for free: we present zero adversarial captions during training. Unified VSE builds fine-grained semantic alignments via the contrastive learning of semantic components. It use the explicit aggregation of the components \mathbf{u}_{comp} to alleviate the dataset biases.

Ablation study: semantic components. We now delve into the effectiveness of different semantic components by choosing different combinations of components for the caption embedding. Shown in Table 2, we use different subsets of

CVPR 2019 Submission #4705. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 5. Examples showing the top-5 image-to-text retrieval results with the presence of adversarial samples. We highlight the positive captions in blue. Words in red indicates the irrelevant words in the adversarial captions. The score in the front of each sentence is the similarity score of the caption and image computed by different model.

Task	obj	attr	rel	obj (det)	sum
VSE++	29.95	26.64	27.54	50.57	134.70
VSE-C	27.48	28.76	26.55	46.20	128.99
U-VSE _{all}	39.23	34.11	39.15	58.12	170.61
U-VSE _{obj}	39.71	33.37	34.38	56.84	164.3
U-VSE _{attr}	31.31	37.51	34.73	52.26	155.81
U-VSE _{rel}	37.55	32.7	39.57	59.12	168.94

Table 3. The mAP performance on the unified text-to-image retrieval task. Please refer to the text for details.

the semantic components to form the bag-of-component embeddings \mathbf{u}_{comp} . For example, in U-VSE_{obj}, only object nouns are selected and aggregated as \mathbf{u}_{comp} .

The results demonstrate the effectiveness of the enforcement of semantic coverage: even if the semantic components have got fine-grained alignment with visual concepts, directly using \mathbf{u}_{sent} as the caption encoding still degenerates the robustness against adversarial examples. Consistent with the intuition, enforcing of coverage of a certain type of components (*e.g.*, objects) helps the model to defend the adversarial attacks of the same type (*e.g.*, defending adversarial attacks of nouns). Combining all components leads to the best performance.

Choice of the combination factor: α . We study the choice of α by conducting experiments on both normal retrieval tasks and the adversarial one. Fig 4.2 shows the R@1 performance under the normal/adversarial retrieval scenario w.r.t. different choices of α . We observe that the \mathbf{u}_{comp} term contributes little on the normal cross-modal retrieval tasks but largely on tasks with adversarial attacks. Recall that α can be viewed as a factor that balances the training objective and the enforcement of semantic coverage. By choosing α from a reasonable region (*e.g.*, from 0.6 to 0.8), our model can effectively defend adversarial attacks, with no sacrifice for the overall performance.

4.3. Unified Text-to-Image Retrieval

We extend the word-to-scene retrieval used by [41] into a general *unified text-to-image retrieval* task. In this task, models receive queries of different semantic levels, including single words (*e.g.*, “Clock.”), noun phrases (*e.g.*, “White clock.”), relational phrases (*e.g.*, “Clocks on wall”) and full



[1] (0.592) two passenger jets are on the tarmac by the airport.
[2] (0.590) two passenger jets **swing** on the tarmac by the airport.

VSE++ [3] (0.579) two passenger jets **contain** on the tarmac by the airport.

[4] (0.573) an airplane parked on the tarmac **stick** an airport.

[5] (0.571) two passenger jets **underneath** on the tarmac by the airport.

[1] (0.515) a large white and blue plane is **hit** the airport.

[2] (0.513) two passenger jets are on the tarmac by the airport.

VSE-C [3] (0.513) two passenger jets **contain** on the tarmac by the airport.

[4] (0.508) a large white and blue plane is **at** the airport.

[5] (0.504) two passenger jets **swing** on the tarmac by the airport.

[1] (0.603) a large white and blue plane is **at** the airport.

[2] (0.562) a large white and blue plane is **hit** the airport.

U-VSE [3] (0.544) two passenger jets **underneath** on the tarmac by the airport.

[4] (0.539) two passenger jets are on the tarmac by the airport.

[5] (0.538) an airplane parked on the tarmac at an airport

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

sentences. For all baselines, the texts of different types as treated as full sentences. The result is presented in Table 3.

We generate positive image-text pairs by randomly choosing an image and a semantic component from 5 matched captions with the chosen image. It is worth mention that the semantic components extracted from captions may not cover all visual concepts in the corresponding image, which makes the annotation noisy. To address this, we also leverage the MS-COCO detection annotations to facilitate the evaluation (see $obj(det)$ column). We treat the labels for detection bounding boxes as the annotation of objects in the scene.

Ablation study: contrastive learning of components. We evaluate the effectiveness of using contrastive samples for different semantic components. Shown in Table 3, U-VSE_{obj} denotes the model trained with only contrastive samples of noun components. The same notation applies to other models. The U-VSE trained with a certain type of contrastive examples (*e.g.*, U-VSE_{obj} with contrastive nouns) consistently improves the retrieval performance of the same type of queries (*e.g.*, retrieving images from a single noun). U-VSE trained with all kinds of contrastive samples performs best in overall and shows a significant gap w.r.t. other baselines.

Visualization of the semantic alignment. We visualize the semantic-relevance map on an image w.r.t. a given query \mathbf{u}_q for a qualitative evaluation of the alignment performance of various semantic components. The map M_i is computed as the similarity between each image region \mathbf{v}_i and \mathbf{u}_q , in a similar way as Eq. (2). Shown as Fig. 6, this visualization helps to verify that our model successfully aligns different semantic components with the corresponding image regions.

4.4. Semantic Parsing with Visual Cues

As a side application, we show how the learned unified VSE space can provide the visual cues to help the semantic parsing of sentences. Fig. 7 shows the general idea. When parsing a sentence, ambiguity may occur, *e.g.*, the subject of the relational word *eat* may be *sweater* or *burger*. It is not easy for a textual parser to decide which one is correct because of the innate syntactic ambiguity. However, we can use the image which is depicted by this sentence to assist the parsing by. This is related to previous works on using image segmentation models to facilitate the sentence parsing [6].

	Query: black dog	Query: white dog	Query: player swing bat																																																																		
Retrieved Image																																																																					
Relevance Map																																																																					
Grounded Area																																																																					
Matching Score	0.257	0.255	0.247	0.211	0.205	0.302	0.286	0.255	0.251	0.247	0.490	0.406	0.404	0.393	0.359	0.810	0.811	0.812	0.813	0.814	0.815	0.816	0.817	0.818	0.819	0.820	0.821	0.822	0.823	0.824	0.825	0.826	0.827	0.828	0.829	0.830	0.831	0.832	0.833	0.834	0.835	0.836	0.837	0.838	0.839	0.840	0.841	0.842	0.843	0.844	0.845	0.846	0.847	0.848	0.849	0.850	0.851	0.852	0.853	0.854	0.855	0.856	0.857	0.858	0.859	0.860	0.861	0.862	0.863

Figure 6. The relevance maps and grounded areas obtained from the retrieved images w.r.t. three queries. The temperature of the softmax for visualizing the relevance map is $\tau = 0.1$. Pixels in white indicates a higher matching score. Note that the third image of the query “black dog” contains two dogs, while our model successfully locates the black one (on the left). It also succeeded in finding the white dog in the first image of “white dog”. Moreover, for the query “player swing bat”, although there are many players in the image, our model only attend to the man swinging the bat.

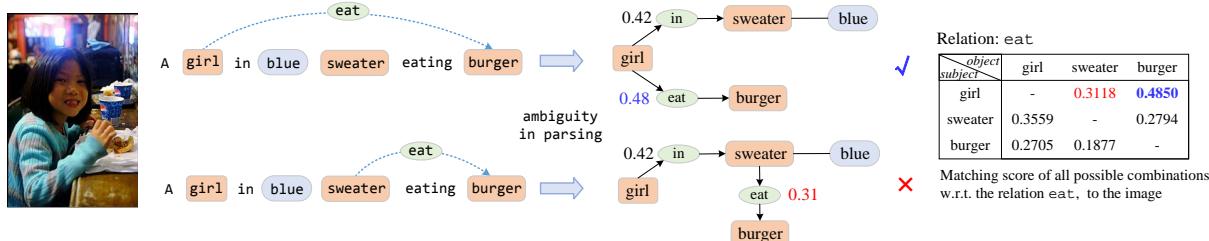


Figure 7. Example showing our model can leverage image to assist semantic parsing when there is ambiguity in the sentence. We can infer that the matching score of “girl eat burger” is much higher than “sweater eat burger”, which can help to eliminate the ambiguity. Note that the other components in the scene graph are also correctly inferred by our model.

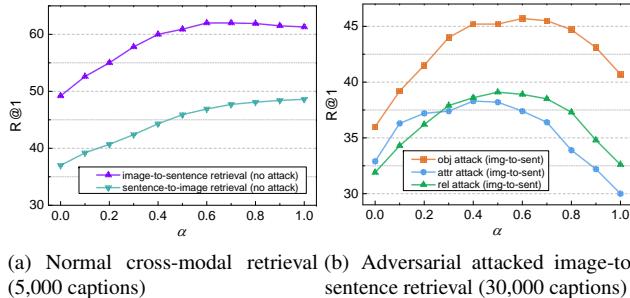


Figure 8. The performance of U-VSE on cross-modal retrieval tasks with different combination weight α . Our model can effectively defend adversarial attacks, with no sacrifice for the performance on other tasks by choosing a reasonable α (thus we set $\alpha = 0.75$ in all other experiments).

This motivates us to design two tasks, 1) recovering the dependency between attributes and entities, and 2) recovering the relational triples. In detail, we first extract the entities, attributes and relational words from the raw sentence without knowing their dependencies. For each possible combination of certain semantic component, our model computes its embedding in the unified joint space. *E.g.*, in Fig. 7, there are in total $3 \times (3 - 1) = 6$ possible dependencies for eat. We choose the combination with the highest matching score with the image to decide the subject/object dependencies of the relation eat. We use parsed semantic components as the ground-truth and report the accuracy, defined as the fraction of the number of correct dependency resolution and the total number of attributes/relations.

Task	attributed object	relational phrase
Random	37.41	31.90
VSE++	41.12	43.31
VSE-C	43.44	41.08
U-VSE	64.48	60.76

Table 4. The accuracy of different models on recovering word dependencies with visual cues. In the “Random” baseline, we randomly assign the word dependencies.

Table 4 reports the results on assisting semantic parsing with visual cues, compared with other baselines. Fig. 7 shows a real case in which we successfully resolve the textual ambiguity.

5. Conclusion

We present a *unified visual-semantic embedding* approach that learns a joint representation space of vision and language in a factorized manner: Different levels of textual semantic components such as objects and relations get aligned with regions of images. A contrastive learning approach for semantic components is proposed for the efficient learning of the fine-grained alignment. We also introduce the enforcement of semantic coverage: each caption embedding should have a coverage of all semantic components in the sentence. Unified VSE shows superiority on multiple cross-modal retrieval tasks and can effectively defend text-domain adversarial attacks. We hope the proposed approach can empower machines that learn vision and language jointly, efficiently and robustly.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] VSE-C open-sourced code. <https://github.com/ExplorerFreda/VSE-C>. 6
- [2] VSE++ open-sourced code. <https://github.com/fartashf/vsepp>. 6
- [3] O. Abend, T. Kwiatkowski, N. J. Smith, S. Goldwater, and M. Steedman. Bootstrapping language acquisition. *Cognition*, 164:116–143, 2017. 1, 2
- [4] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 2
- [5] L. Banarescu, C. Bonial, S. Cai, M. Georgescu, K. Griffitt, U. Hermjakob, K. Knight, P. Koehn, M. Palmer, and N. Schneider. Abstract Meaning Representation for Sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186, 2013. 2
- [6] G. Christie, A. Laddha, A. Agrawal, S. Antol, Y. Goyal, K. Kochersberger, and D. Batra. Resolving language and vision ambiguities together: Joint segmentation & prepositional attachment resolution in captioned scenes. *arXiv preprint arXiv:1604.02125*, 2016. 7
- [7] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. In *NIPS 2014 Workshop on Deep Learning, December 2014*, 2014. 4
- [8] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-Term Recurrent Convolutional Networks for Visual Recognition and Description. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2625–2634, 2015. 2
- [9] A. Eisenschitz and L. Wolf. Linking Image and Text with 2-way Nets. In *IEEE conference on computer vision and pattern recognition (CVPR)*, pages 1855–1865, 2017. 2, 6
- [10] M. Engilberge, L. Chevallier, P. Pérez, and M. Cord. Finding Beans in Burgers: Deep Semantic-Visual Embedding with Localization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3984–3993, 2018. 2, 6
- [11] F. Faghri, D. J. Fleet, J. R. Kiros, and S. Fidler. VSE++: Improving Visual-Semantic Embeddings with Hard Negatives. 2018. 2, 3, 4, 5, 6, 11
- [12] A. Fazly, A. Alishahi, and S. Stevenson. A probabilistic computational model of cross-situational word learning. *Cognitive Science*, 34(6):1017–1063, 2010. 1
- [13] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al. Devise: A Deep Visual-Semantic Embedding Model. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2121–2129, 2013. 2, 4
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4
- [15] Y. Huang, W. Wang, and L. Wang. Instance-Aware Image and Sentence Matching with Selective Multimodal LSTM. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7254–7262, 2017. 3, 5, 6

- [16] Y. Huang, Q. Wu, and L. Wang. Learning Semantic Concepts and Order for Image and Sentence Matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6163–6171, 2017. 6
- [17] J. Johnson, A. Gupta, and L. Fei-Fei. Image Generation from Scene Graphs. 2018. 2
- [18] J. Johnson, R. Krishna, M. Stark, L. Li, D. A. Shamma, M. S. Bernstein, and F. Li. Image Retrieval using Scene Graphs. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3668–3678, 2015. 2
- [19] A. Karpathy and L. Fei-Fei. Deep Visual-Semantic Alignments for Generating Image Descriptions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3128–3137, 2015. 2, 5
- [20] A. Karpathy and L. Fei-Fei. Deep Visual-Semantic Alignments for Generating Image Descriptions. In *IEEE conference on computer vision and pattern recognition (CVPR)*, pages 3128–3137, 2015. 6
- [21] A. Karpathy, A. Joulin, and L. F. Fei-Fei. Deep Fragment Embeddings for Bidirectional Image Sentence Mapping. In *Conference on Neural Information Processing Systems (NIPS)*, pages 1889–1897, 2014. 2
- [22] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*, 2017. 5
- [23] R. Kiros, R. Salakhutdinov, and R. Zemel. Multimodal Neural Language Models. In *International Conference on Machine Learning*, pages 595–603, 2014. 2
- [24] R. Kiros, R. Salakhutdinov, and R. S. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *arXiv preprint arXiv:1411.2539*, 2014. 2, 6
- [25] O. Levy, K. Lee, N. Fitzgerald, and L. Zettlemoyer. Long Short-Term Memory as a Dynamically Computed Element-wise Weighted Sum. *arXiv preprint arXiv:1805.03716*, 2018. 4
- [26] P. Liang, M. I. Jordan, and D. Klein. Learning Dependency-Based Compositional Semantics. *Computational Linguistics*, 39(2):389–446, 2013. 2
- [27] 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 (ECCV)*, pages 740–755, 2014. 3, 5
- [28] Y. Liu, Y. Guo, E. M. Bakker, and M. S. Lew. Learning a Recurrent Residual Fusion Network for Multimodal Matching. In *IEEE international conference on computer vision (ICCV)*, pages 4127–4136, 2017. 2, 5, 6
- [29] C. Lu, R. Krishna, M. Bernstein, and L. Fei-Fei. Visual relationship detection with language priors. In *European Conference on Computer Vision (ECCV)*, pages 852–869. Springer, 2016. 2
- [30] L. Ma, Z. Lu, L. Shang, and H. Li. Multimodal Convolutional Neural Networks for Matching Image and Sentence. In *IEEE international conference on computer vision (ICCV)*, pages 2623–2631, 2015. 1, 2, 6
- [31] L. v. d. Maaten and G. Hinton. Visualizing Data Using t-SNE. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008. 12

- 972 [32] J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille.
973 Deep Captioning with Multimodal Recurrent Neural Net-
974 works (m-RNN). In *International Conference on Learning
975 Representations (ICLR)*, 2015. 6
- 976 [33] R. Montague. Universal Grammar. *Theoria*, 36(3):373–398,
977 1970. 2
- 978 [34] Z. Niu, M. Zhou, L. Wang, X. Gao, and G. Hua. Hierarchical
979 Multimodal LSTM for Dense Visual-Semantic Embedding.
980 In *IEEE international conference on computer vision (ICCV)*,
981 pages 1899–1907, 2017. 2, 6
- 982 [35] J. Pennington, R. Socher, and C. Manning. GloVe: Global
983 Vectors for Word Representation. In *Proceedings of the
984 2014 Conference on Empirical Methods in Natural Language
985 Processing (EMNLP)*, pages 1532–1543, 2014. 3
- 986 [36] Z. Ren, H. Jin, Z. Lin, C. Fang, and A. Yuille. Joint Image-
987 Text Representation by Gaussian Visual-Semantic Embed-
988 ding. In *ACM Multimedia (ACM-MM)*, pages 207–211, 2016.
989 2
- 990 [37] Z. Ren, H. Jin, Z. Lin, C. Fang, and A. Yuille. Multiple
991 Instance Visual-Semantic Embedding. In *British Machine
992 Vision Conference (BMVC)*, 2017. 2
- 993 [38] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma,
994 Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg,
995 and L. Fei-Fei. ImageNet Large Scale Visual Recognition
996 Challenge. *International Journal of Computer Vision (IJCV)*,
997 115(3):211–252, 2015. 4
- 998 [39] S. Schuster, R. Krishna, A. Chang, L. Fei-Fei, and C. D. Man-
999 ning. Generating semantically precise scene graphs from tex-
1000 tual descriptions for improved image retrieval. In *Workshop
1001 on Vision and Language (VL15)*, Lisbon, Portugal, September
1002 2015. Association for Computational Linguistics. 3
- 1003 [40] R. Shekhar, S. Pezzelle, Y. Klimovich, A. Herbelot, M. Nabi,
1004 E. Sangineto, and R. Bernardi. FOIL it! Find One Mis-
1005 match between Image and Language Caption. *arXiv preprint
arXiv:1705.01359*, 2017. 1, 5
- 1006 [41] H. Shi, J. Mao, T. Xiao, Y. Jiang, and J. Sun. Learning
1007 Visually-Grounded Semantics from Contrastive Adversarial
1008 Samples. In *Proceedings of the 27th International Conference
1009 on Computational Linguistics (COLING)*, pages 3715–3727,
1010 2018. 1, 2, 3, 4, 5, 6, 7, 11, 12
- 1011 [42] A. Shrivastava, A. Gupta, and R. Girshick. Training Region-
1012 Based Object Detectors with Online Hard Example Mining.
1013 In *IEEE Conference on Computer Vision and Pattern Recog-
1014 nition (CVPR)*, pages 5253–5262, 2016. 2, 3
- 1015 [43] I. Vendrov, R. Kiros, S. Fidler, and R. Urtasun. Order-
1016 Embeddings of Images and Language. In *International Con-
1017 ference on Learning Representations (ICLR)*, 2016. 5, 6
- 1018 [44] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and
1019 Tell: A Neural Image Caption Generator. In *IEEE Conference
1020 on Computer Vision and Pattern Recognition (CVPR)*, pages
3156–3164, 2015. 2
- 1021 [45] L. Wang, Y. Li, and S. Lazebnik. Learning Deep Structure-
1022 Preserving Image-Text Embeddings. In *IEEE conference
1023 on computer vision and pattern recognition (CVPR)*, pages
5005–5013, 2016. 2, 6
- 1024 [46] J. Wu, Y. Yu, C. Huang, and K. Yu. Deep multiple instance
1025 learning for image classification and auto-annotation. In *The
1026 IEEE Conference on Computer Vision and Pattern Recog-
1027 nition (CVPR)*, June 2015. 5
- 1028 [47] F. Xiao, L. Sigal, and Y. J. Lee. Weakly-Supervised Vi-
1029 sual Grounding of Phrases with Linguistic Structures. In
1030 *IEEE Conference on Computer Vision and Pattern Recog-
1031 nition (CVPR)*, pages 5945–5954, 2017. 2
- 1032 [48] H. Xu and K. Saenko. Ask, Attend and Answer: Explor-
1033 ing Question-Guided Spatial Attention for Visual Question
1034 Answering. In *European Conference on Computer Vision
1035 (ECCV)*, pages 451–466. Springer, 2016. 2
- 1036 [49] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov,
1037 R. Zemel, and Y. Bengio. Show, Attend and Tell: Neural
1038 Image Caption Generation with Visual Attention. In *Inter-
1039 national Conference on Machine Learning (ICML)*, pages
2048–2057, 2015. 2
- 1040 [50] Q. You, Z. Zhang, and J. Luo. End-to-End Convolutional
1041 Semantic Embeddings. In *IEEE Conference on Computer
1042 Vision and Pattern Recognition (CVPR)*, pages 5735–5744,
1043 2018. 1, 2, 3, 5, 6
- 1044 [51] L. S. Zettlemoyer and M. Collins. Learning to Map Sentences
1045 to Logical Form: Structured Classification with Probabilistic
1046 Categorical Grammars. In *Proceedings of the Twenty-First
1047 Conference on Uncertainty in Artificial Intelligence (UAI)*,
1048 pages 658–666, 2005. 2

1080 This supplementary material is organized as follows.
 1081 First, in Appendix 6, we provide more details for the imple-
 1082 mentation of our model and the training method. Second,
 1083 in Appendix 7, we provide the experiment setups, metrics,
 1084 baseline implementations, qualitative examples and analysis
 1085 for each experiment we discussed in the main text. We end
 1086 this section with the visualization of the learned unified VSE
 1087 space of different semantic levels.
 1088

6. Implementation Details

6.1. Generating Negative samples

To generate negative samples in sentence level, we follow the sampling paradigm introduced by [11]: We sample negative examples from all other captions/images in the dataset in a training batch. Note that as [11] shown, the batch size will largely affect the models' performance. For a fair comparison, we set the batch size as 128 which is the same as [11, 41]. In the rest of this section, we discuss in detail how we sample negative semantic components.

As for nouns, we sample 16 negative nouns from a fixed set of nouns: This noun set is extracted from nouns with frequency more than 100 (in total 1,205 nouns extracted in MS-COCO dataset).

As for attribute-noun pairs, we randomly sample 8 other attributes in a fixed attribute set and replace the original attribute in the pair, as negative examples. The attribute set is composed by the frequently appeared attributes in the MS-COCO dataset. In detail, we extract in total 37 attributes, *i.e.*, white, black, red, green, brown, yellow, orange, pink, gray/grey, purple, young, wooden, old, snowy, grassy, cloudy, colorful, sunny, beautiful, bright, sandy, fresh, morden, cute, dry, dirty, clean, polar, crowded, silver, plastic, concrete, rocky, wooded, messy, square. We also randomly replace nouns in the pairs to generate another set of negative attribute-noun pairs. For each attribute-noun pair, we randomly draw 16 negative examples.

We separately compute the ranking loss corresponding to two types of negatives, denoted as $\ell_{attr_{negnoun}}$ and $\ell_{attr_{negattr}}$. Both of them are computed by a uni-directional ranking loss with negative examples drawn in text-domain. OHEM strategy is not applied on them. The final loss is the sum of them, *i.e.*, $\ell_{attr} = \ell_{attr_{negnoun}} + \ell_{attr_{negattr}}$.

Here we add a small note for the reproducibility. In cases with multiple modifiers on the nouns (*e.g.*, old black dog), for simplicity, in our implementation, we always extract the first modifier of each noun phrases as its attribute (old dog in this case).

As for relational triples, we randomly sample 4 relational words and 2 negatives subjects (nouns) and 2 negative objects (nouns) to replace the corresponding parts in the triple,

as negative examples. In total, we have 8 negative triples for each relational triple. The choice of this small number of negative examples is attributed to the trade-off between the computational efficiency and stability of training. Empirically, we find that increasing the number of negative triples does not bring much improvement to the performance.

We also sample negative relational triples from other captions within the training batch. In detail, we sample 1 negative rational triple for each other caption within the batch. This results in at most $128 - 1 = 127$ negative examples for each relational triple (“at most” means some captions may not contain relational phrases). Similar as attribute-noun pairs, we individually compute the ranking loss on each type of negatives and sum them together as the ℓ_{rel} . The losses are computed by uni-directional ranking loss without OHEM.

As for negative bag-of-components, we sample negative ones in a similar manner as we do for sentences: We draw them from the bag-of-components in other captions within the training batch. We also draw other images from batch as negative images. The loss ℓ_{comp} is computed by bi-directional ranking loss with OHEM strategy.

6.2. Hyperparameters

We set the dimension d_{basic} of basic semantic embeddings as 300. The embeddings are initialized by GloVe word embeddings pre-trained on the Common Crawl dataset: <http://nlp.stanford.edu/data/glove.840B.300d.zip>. The dimension d_{modif} of modifier semantic embeddings is set to 100. The embeddings are randomly initialized. During training, we fix the basic semantic embeddings of words $w^{(basic)}$. The learning rate of the Adam optimizer is fixed to 0.001 at first 6 epochs and is exponentially decayed by 2 for each next epoch until it reaches 1e-5.

7. Experiment Details

7.1. Cross-modal Retrieval

Visualizations. We show a set of examples of the image-to-sentence retrieval in Fig. 9 and sentence-to-image retrieval in Fig. 10.

7.2. Retrieval under text-domain adversarial attack

Experiment setup. We use the 1K test split (including 5,000 captions) for generating adversarial attacks. For each caption, we generate five adversarial captions under one type of attack setting. The detailed settings of the three types of adversarial attack are listed below.

1. **Object attack:** We randomly replace / append by an irrelevant noun for both 50% probability. The replacing/appending place is randomly selected in nouns of

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

1188 the caption. For the case of appending extra noun, the
 1189 word and is also added before the appended noun, *e.g.*,
 1190 A dog eats meat → A dog eats meat **and**
 1191 **table**. The irrelevant nouns are drawn from the set
 1192 containing nouns with high concreteness (manually ex-
 1193 tracted).

1194
 1195 **Attribute attack:** If a caption contains attribute-noun
 1196 pairs. We randomly select one pair and replace the
 1197 attribute by a negative one. If a caption does not
 1198 contain any attributes, we randomly choose one noun in
 1199 the caption and append an attribute on it. The negative
 1200 attribute is generated from the attribute set excluding
 1201 the attributes (and its similar attributes) in the caption.
 1202 The similar attribute group is defined as the following.
 1203 {white, snowy, polar}, {red, pink},
 1204 {blue, cloudy}, {green, grassy},
 1205 {brown, sandy, yellow, orange},
 1206 {rocky, concrete}.

1207
 1208 **Relational attack:** For those captions containing rela-
 1209 tional phrases, we randomly select one relation triple
 1210 and with equivalent probability to choose one in the
 1211 triple to be replaced by an irrelevant one. *e.g.*, A dog
 1212 eats meat → A dog **plays** meat. For those
 1213 captions which do not have any relational phrases, we
 1214 first randomly select one noun in the caption and re-
 1215 gard it as a subject/object with 50% / 50% probabili-
 1216 ty. Then we draw a relational word and an irrele-
 1217 vant noun as the object/subject to form a new fake rela-
 1218 tion. *e.g.*, A dog is sleeping → A dog **in**
 1219 **sky** is sleeping.

1220 **Baselines.** We train the VSE-C according to the setting
 1221 in [41] with the officially open-sourced code. In the orig-
 1222 inal VSE-C paper, The VSE-C is trained by generating ei-
 1223 ther noun-typed/numeral-typed/relation-typed or all of these
 1224 three types of adversarial samples. We use the setting of
 1225 training under all types of adversarial samples as a compara-
 1226 ble competitor in this evaluation.

1227 **Visualizations.** We show a set of examples of image-to-text
 1228 retrieval under text-domain adversarial attack in Fig. 11.

7.3. Unified text-to-image retrieval

1229 **Experiment setup.** We use the 1K test split as the retrieval
 1230 set. The queries are generated from frequent semantic com-
 1231 ponents extracted by the semantic parser from the training
 1232 set. We regard a query as a valid one if at least 3 images (5
 1233 for noun-level retrieval) in the test set contain the query. For
 1234 the obj(det) queries, we directly use the class names of the
 1235 MS-COCO object detection / segmentation annotations.

1236 **Baselines.** For VSE++ and VSE-C, as they do not have an
 1237 object-level encoder. For any query, we always regard it as a
 1238 short sentence and feed it into the sentence encoder to get

1239 the embedding of the query text. For U-VSE, as it has the
 1240 object-level encoder which means a noun/attribute-noun pair
 1241 can be either encoded by the object encoder ϕ or by neural
 1242 combiner ψ by regarding the query as a short sentence. We
 1243 select the encoder having higher performance on a validation
 1244 set and report the results.

1245 **Visualizations.** We show a set of retrieved image by queries
 1246 of various types in Fig. 12.

7.4. Semantic Parsing

1247 **Experiment setup.** We also use the 1K test split for this
 1248 experiment. For each caption, we first extract nouns,
 1249 adjectives and relational words. We call adjective and relational
 1250 words as content words. The model should recover the de-
 1251 pendencies linked with them. We exclude some relational
 1252 words whose lexical meanings are usually ambiguous, such
 1253 as *include*, *to*, *of*, *etc*.

1254 Given a content word (either an adjective or a relational
 1255 word), we generate all possible dependencies among nouns
 1256 in the sentence to form candidate dependencies. Each candi-
 1257 date dependency, which is either an adjective-noun pair or a
 1258 subject-relation-object triple, will get a matching score w.r.t.
 1259 the image (the *visual cue*). We select the dependency that
 1260 has the highest score as the recovered dependency w.r.t. the
 1261 chosen content word.

1262 **Metrics.** We report the accuracy of the recovered semantic
 1263 dependencies. In detail, for an attribute-noun dependency,
 1264 the model gets a correct count if the dependency having
 1265 the highest matching score is identical to the ground-truth.
 1266 For the dependency of a relation, the model gets 0.5 correct
 1267 counts if the subject/object of the answer is the same as the
 1268 ground-truth. If both of them are the same as the ground-
 1269 truth, the model gets 1 correct count. The reported accuracy
 1270 computed as the fraction between total correct counts and
 1271 the total number of dependencies.

1272 **Visualizations and failure case study.** Shown in Fig. 13,
 1273 we visualize some successful and failure cases in seman-
 1274 tic parsing with visual cues. Error source analysis is also
 1275 provided.

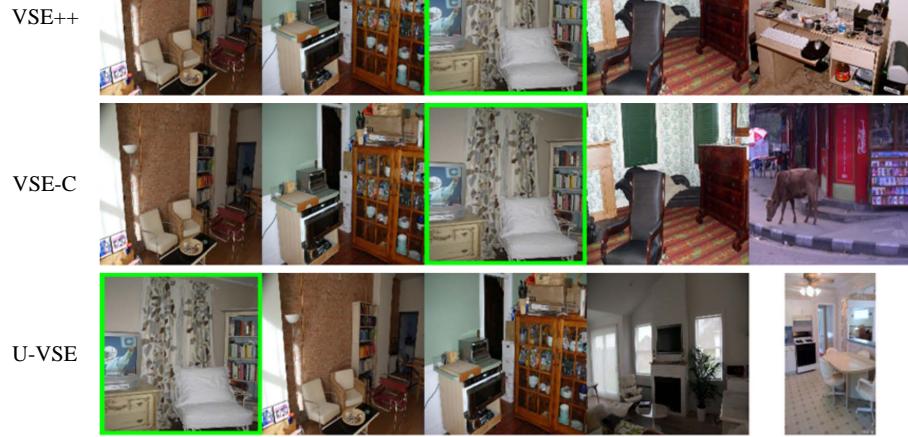
7.5. Embedding Visualization

1276 We visualize the semantic space of different semantic
 1277 levels by t-SNE [31]. The result can be found in Fig. 7.5.
 1278 Through the joint learning of vision and language, our unified
 1279 VSE space successfully recovers the similarities between
 1280 semantic components at various levels.

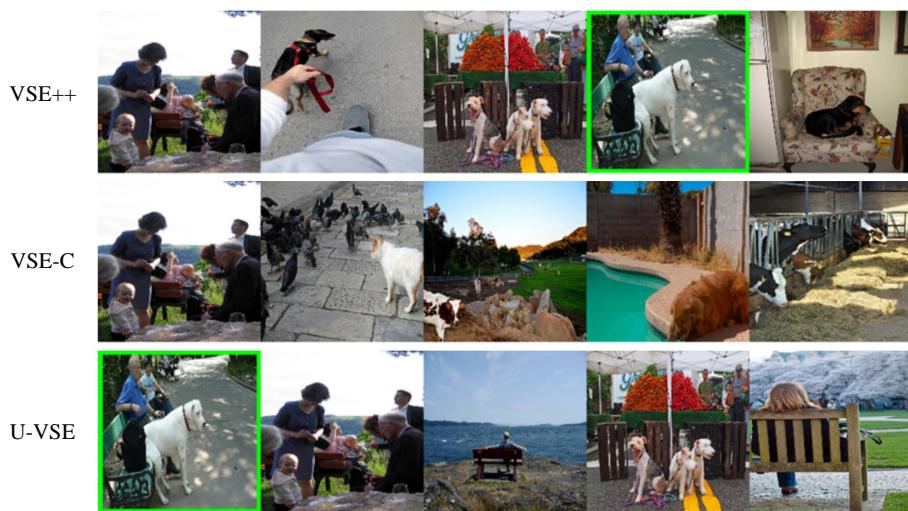
CVPR 2019 Submission #4705. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

1296				1350
1297				1351
1298				1352
1299				1353
1300				1354
1301				1355
1302		[1] (0.476) A few people that are playing tennis on a court. [2] (0.472) Two people playing a match of tennis on a court. VSE++ [3] (0.460) Several men playing with a soccer ball in a park. [4] (0.457) Two young men playing game of soccer. [5] (0.455) There is a man running on a field with a soccer ball.		1356
1303		[1] (0.428) There are two soccer teams playing a game on the field. [2] (0.401) A few people that are playing tennis on a court. VSE-C [3] (0.395) Several boys on a field playing with a frisbee. [4] (0.381) A group of people playing soccer in a field. [5] (0.374) There are people playing a game of tennis.		1357
1304		[1] (0.456) A man carrying a soccer ball down a field. [2] (0.454) There is a man running on a field with a soccer ball.		1358
1305		U-VSE [3] (0.440) A man that is on a soccer field with a ball. [4] (0.429) A man kicking a soccer ball while standing on a field. [5] (0.411) The soccer player is bringing back the ball into play.		1359
1306				1360
1307				1361
1308				1362
1309				1363
1310				1364
1311				1365
1312		[1] (0.520) A bowl with something in it with a banana next to it. [2] (0.500) A banana sits by two oranges, a bowl and a white plate on a white tray. VSE++ [3] (0.498) The banana is laying next to an almost empty bowl. [4] (0.492) A banana and a nearly empty bowl of food resting on top of a table. [5] (0.467) A white tray with a banana and two tangerines and a plate and bowl.		1366
1313		[1] (0.465) The banana is laying next to an almost empty bowl. [2] (0.440) A banana and a nearly empty bowl of food resting on top of a table. VSE-C [3] (0.423) A white tray with a banana and two tangerines and a plate and bowl. [4] (0.414) A bowl with something in it with a banana next to it. [5] (0.360) A bowl filled with leftover food sitting next to a banana.		1367
1314		[1] (0.551) A banana and two oranges sit on a tray next to a bowl and a plate. [2] (0.519) A bowl with something in it with a banana next to it.		1368
1315		U-VSE [3] (0.506) A banana sits by two oranges, a bowl and a white plate on a white tray. [4] (0.502) A white tray with a banana and two tangerines and a plate and bowl. [5] (0.498) The banana is laying next to an almost empty bowl.		1369
1316				1370
1317				1371
1318				1372
1319				1373
1320				1374
1321				1375
1322				1376
1323		[1] (0.373) A couple of horses standing in a field. [2] (0.353) Two giraffes standing in front of each other. VSE++ [3] (0.346) A big heard of cows walking down a road in a row with green tags on their ears. [4] (0.345) Sheep that have been sheared standing in a pen. [5] (0.344) Mythical character with white horse standing on grooved surface.		1377
1324		[1] (0.325) Ten porcelain pieces with floral patterns painted on them. [2] (0.310) Two horses have feathers on their head.		1378
1325		VSE-C [3] (0.304) Two giraffes standing in front of each other. [4] (0.303) Horses standing in shallow water in a wooded area. [5] (0.298) Two dogs lay next to each other on a brown couch.		1379
1326		[1] (0.363) Three different horse figurines are placed beside each other. [2] (0.360) A couple of white horses standing in front of a building.		1380
1327		U-VSE [3] (0.345) Three plastic horse figurines standing next to each other on a shelf. [4] (0.344) Two horses with red feathers on top of their heads. [5] (0.339) Three model horses on a table in front of a pegboard backdrop.		1381
1328				1382
1329				1383
1330				1384
1331				1385
1332				1386
1333				1387
1334		[1] (0.410) A basketball player holds a basketball for a picture. [2] (0.395) A woman standing in the dark holding up a cell phone.		1388
1335		VSE++ [3] (0.383) A person with a basketball stands in front of a goal. [4] (0.371) A young woman is posing for camera. [5] (0.352) A woman standing next to another woman in a building.		1389
1336		[1] (0.322) A woman hugging a girl who is holding a suitcase. [2] (0.297) A young woman is posing for a camera.		1390
1337		VSE-C [3] (0.296) A young man in green jersey is holding a ball. [4] (0.290) A woman with her arms around a girl who is holding a suitcase. [5] (0.288) A woman standing in the dark holding up a cell phone.		1391
1338		[1] (0.363) A basketball player holds a basketball for a picture. [2] (0.360) A uniformed boy is holding a basketball with his back to the hoop.		1392
1339		U-VSE [3] (0.345) A person with a basketball stands in front of a goal. [4] (0.344) Two basketball players reach up for the hoop. [5] (0.339) Two basketball players jump to the hoop to block another from scoring.		1393
1340				1394
1341				1395
1342				1396
1343				1397
1344	Figure 9. Examples showing the top-5 image-to-text retrieval results. We highlight the positive captions in blue. The score in the front of each sentence is the similarity score of the caption and image computed by different model.			1398
1345				1399
1346				1400
1347				1401
1348				1402
1349				1403

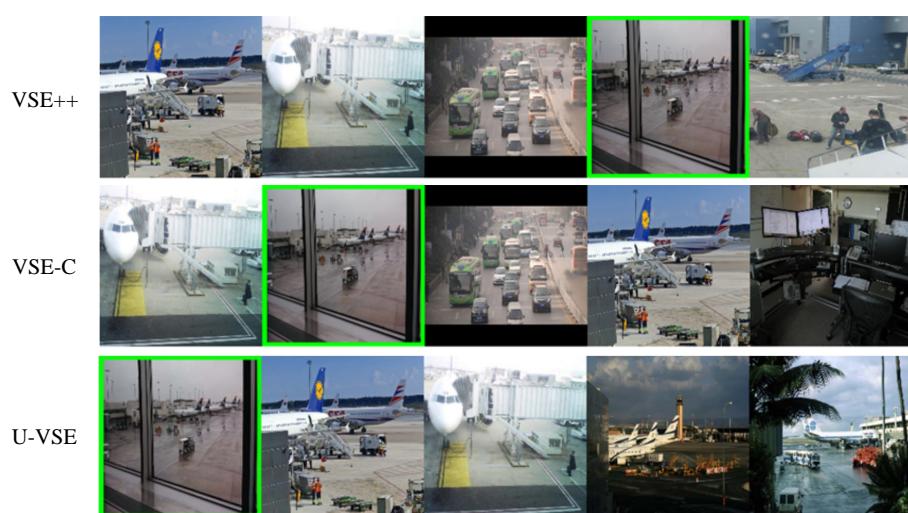
CVPR 2019 Submission #4705. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420

(a) Query: A white chair, books and shelves and a TV on in this room.

1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438

(b) Query: A couple of people sitting on a bench next to a dog.

1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

(c) Query: Window view from the inside of airplanes, baggage carrier and tarmac.

Figure 10. Examples showing the top-5 sentence-to-image retrieval results. We highlight the correct images in green box.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

1512		1566
1513		1567
1514		1568
1515		1569
1516		1570
1517		1571
1518		1572
1519		1573
1520		1574
1521		1575
1522		1576
1523		1577
1524		1578
1525		1579
1526		1580
1527		1581
1528		1582
1529		1583
1530		1584
1531		1585
1532		1586
1533		1587
1534		1588
1535		1589
1536		1590
1537		1591
1538		1592
1539		1593
1540		1594
1541		1595
1542		1596
1543		1597
1544		1598
1545		1599
1546		1600
1547		1601
1548		1602
1549		1603
1550		1604
1551		1605
1552		1606
1553		1607
1554		1608
1555		1609
1556		1610
1557		1611
1558		1612
1559	Figure 11. Examples showing the top-5 image-to-sentence retrieval results with the presence of adversarial samples. We highlight the positive captions in blue. Captions with red words are adversarial samples generated from the original captions. Words in red indicates the irrelevant words in the adversarial captions.	1613
1560		1614
1561		1615
1562		1616
1563		1617
1564		1618
1565		1619

CVPR 2019 Submission #4705. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 12. The top-20 retrieved image in the 1K test split set by queries different types: attribute-object pairs and relational triples.

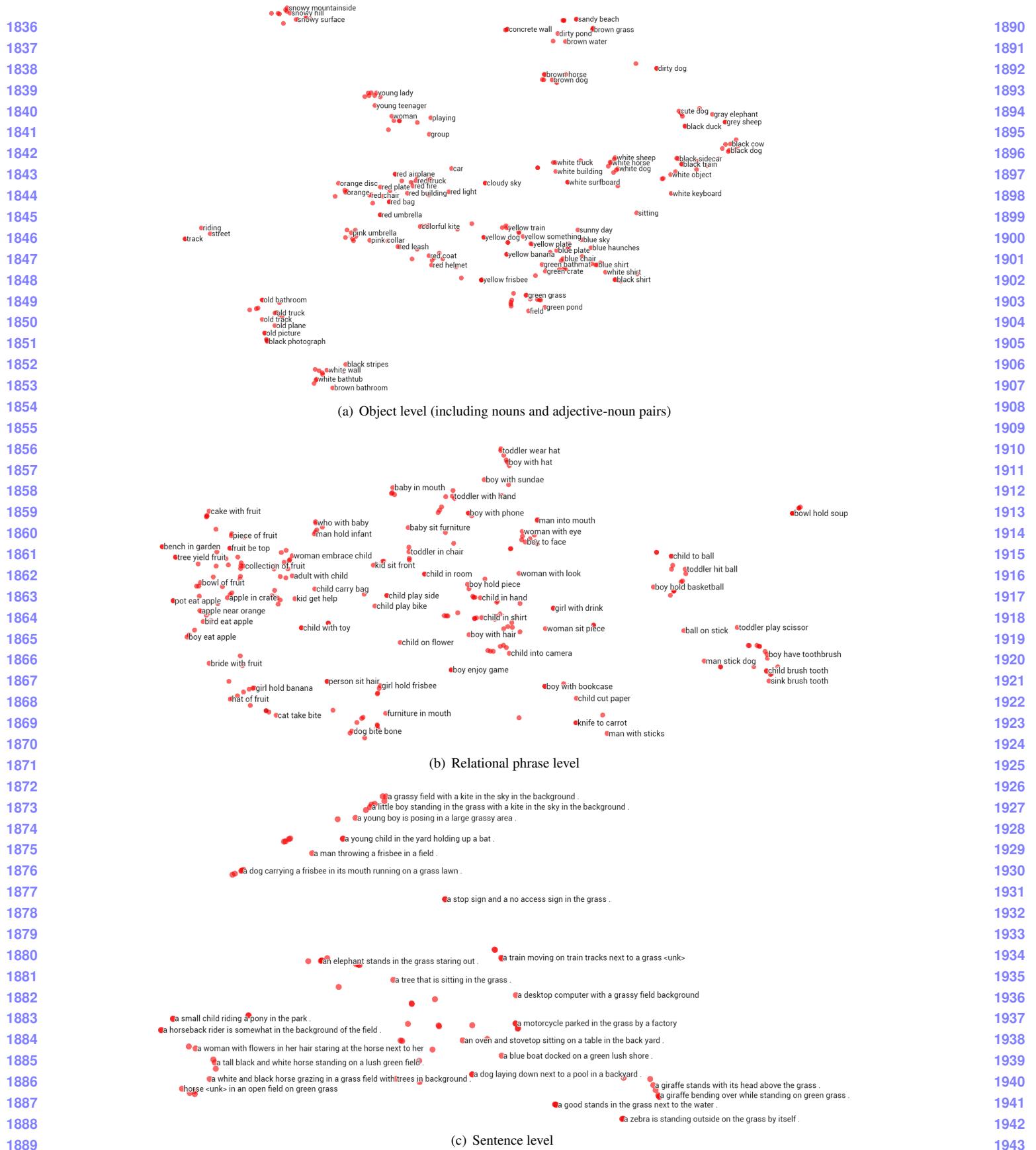


Figure 14. The visualization of the semantic embedding space of different semantic levels. The unified VSE space successfully recovers the similarities between semantic components at various levels.