

GROUNDING REPRESENTATION LEARNING FOR NLP & APPLICATIONS TO ZERO-SHOT LEARNING

Workshop on Representation Learning for Complex Data,
University Lyon 2

Friday 24th May, 2019
Éloi Zablocki

The presented research involves:



- Patrick Gallinari
- Benjamin Piwowarski
- Laure Soulier
- Patrick Bordes

1 Introduction

2 Learning Multi-Modal Word Representation Grounded in Visual Context

3 Zero-Shot Recognition with Semantic Representations and Context

Introduction

Representation learning for NLP

Encoding word semantics —Applications

- Language understanding (QA, summarization, NER, POS Tagging, sentiment analysis)
- Machine translation
- Statistical language modeling (speech recognition, dialog systems)
- Zero-Shot Learning
- ...

Representation learning for NLP

Encoding word semantics —Applications

- Language understanding (QA, summarization, NER, POS Tagging, sentiment analysis)
- Machine translation
- Statistical language modeling (speech recognition, dialog systems)
- Zero-Shot Learning
- ...

Distribution Semantic Models

- Idea: Represent the *semantic* of a word in a vector encoding the distribution of the word contexts.
- Hypothesis: Based on the *distributional hypothesis* (Harris 1954):
“*Words that occur in similar contexts should have similar meanings*”
- Property: Semantic similarity is often quantified by spatial proximity or related geometric measures

Why do we need multi-modal word embeddings?

Context

→ Language is ambiguous, biased and lacks common-sense...

Why do we need multi-modal word embeddings?

Context

→ Language is ambiguous, biased and lacks common-sense...

Example

→ When we say:

"Hand me the salt"

Why do we need multi-modal word embeddings?

Context

→ Language is ambiguous, biased and lacks common-sense...

Example

→ When we say:

"Hand me the salt"

→ We mean:

"Hand me the salt ... or rather the receptacle that contains it. It has a cylindrical shape and is about 3 inches high. This object does not float in the air and lies on the table, in other words, in direct contact with it. If I ask you this favor, it means that it is closer to you than it is to me. Besides, when you give it to me, the aperture should be on the top so that the salt is not spoiled on the table because of gravity, which makes that dropped objects fall down. When you hand it to me, I expect that you give it to by making contact with my hand, located at the end of my arm, that I am going to bring closer to you. Salt is an ionic compound that can be formed by the neutralization reaction of an acid and a base. It is usually extracted from salted water, from seas and ocean. Salt is essential for life in general, and saltiness is one of the basic human tastes. Salt is one of the oldest and most ubiquitous food seasonings. The World Health Organization recommends that adults should consume less than 2000 mg of sodium equivalent to 5 grams of salt per day. Edible salt is sold

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

→ ...while **images** are unequivocal depictions

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

→ ...while **images** are unequivocal depictions

Can language be *grounded* in the visual world?

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

→ ...while **images** are unequivocal depictions

Can language be grounded in the visual world?

Motivation

→ *Psychological studies*: meaning of words is grounded in perception

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

→ ...while **images** are unequivocal depictions

Can language be grounded in the visual world?

Motivation

- *Psychological studies*: meaning of words is grounded in perception
- *Human Reporting Bias* (van Durme 2013)

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

→ ...while **images** are unequivocal depictions

Can language be grounded in the visual world?

Motivation

→ *Psychological studies*: meaning of words is grounded in perception

→ *Human Reporting Bias* (van Durme 2013)

→ *Complementarity* between text and images

Why do we need multi-modal word embeddings?

Context

→ **Language** is ambiguous, biased and lacks common-sense...

unlikely to be mentioned	likely to be mentioned
something expected	unusual facts
trivial facts	something with value

Language \neq Real world

→ ...while **images** are unequivocal depictions

Can language be grounded in the visual world?

Motivation

→ *Psychological studies*: meaning of words is grounded in perception

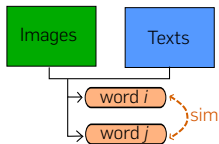
→ *Human Reporting Bias* (van Durme 2013)

→ *Complementarity* between text and images

→ *Advances in Computer Vision*

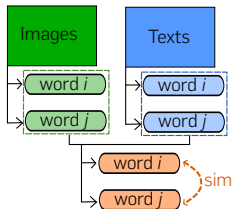
Multi-modal fusion techniques - Sequential models

Joint models

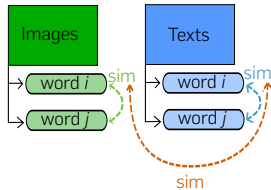


(a) *Early fusion*

Sequential models



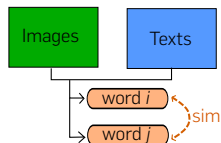
(b) *Middle fusion*



(c) *Late fusion*

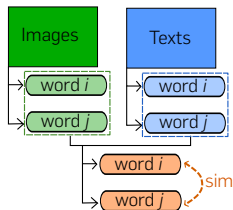
Multi-modal fusion techniques - Sequential models

Joint models

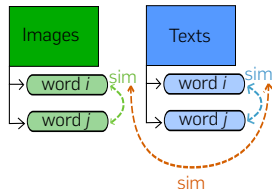


(a) *Early fusion*

Sequential models



(b) *Middle fusion*



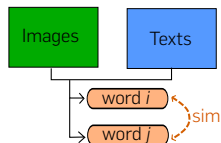
(c) *Late fusion*

Sequential models

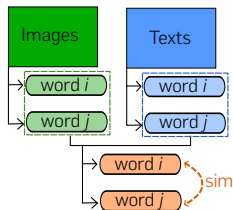
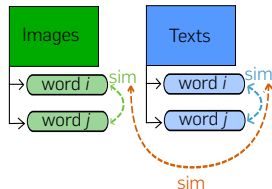
- 1 Separately construct representations for words in each modality
 - **Text**: GloVe or Word2Vec
 - **Image**: Aggregation of activations obtained from a pre-trained CNN forwarded on images

Multi-modal fusion techniques - Sequential models

Joint models

(a) *Early fusion*

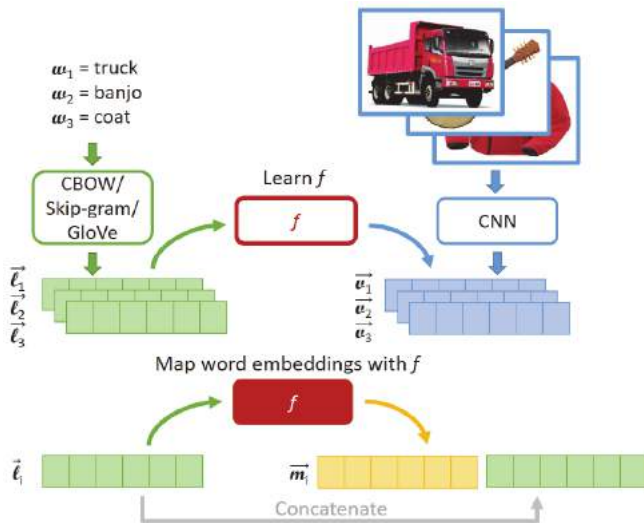
Sequential models

(b) *Middle fusion*(c) *Late fusion*

Sequential models

- 1 Separately construct representations for words in each modality
 - **Text**: GloVe or Word2Vec
 - **Image**: Aggregation of activations obtained from a pre-trained CNN forwarded on images
- 2 Combine them with different techniques
 - **Middle fusion**: form multi-modal representations (e.g. concatenation, CCA)
 - **Late fusion**: interaction in the downstream task (e.g. linear combination of similarity scores)

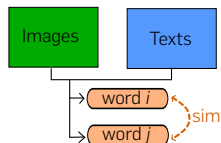
Example of sequential technique (middle fusion)



Sequential model (Collell et al. 2017)

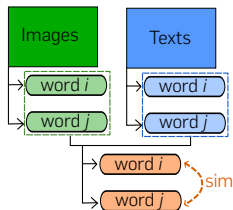
Multi-modal fusion techniques - Joint models

Joint models

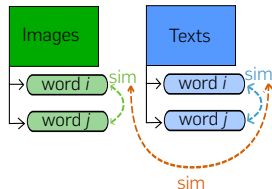


(a) *Early fusion*

Sequential models



(b) *Middle fusion*



(c) *Late fusion*

Joint models (a.k.a. early fusion)

Jointly learn multi-modal representations from multiple modalities

- With aligned text and images: bayesian models, grounded language models...
- Without aligned text and images: skip-gram extension

Multi-modal fusion techniques

Joint models without aligned texts and images

Extensions to the *skip-gram* algorithm under the **distributional hypothesis**:

“Words that occur in similar contexts should have similar meanings”

Additional hypothesis in (Hill et al. 2014): “The frequency of appearance of concrete concepts in texts correlates with the likelihood of experiencing it in the real world.”

Text corpus

Feature-Norm for concrete words

	<i>has_legs</i>	<i>can_fly</i>	<i>has_teeth</i>	<i>is_green</i>	<i>is_animal</i>	...
monkey	1	0	1	0	1	...
crocodile	1	0	1	1	1	...
plane	0	1	0	0	0	...
...

“He saw a crocodile in his backyard!”

→ Add fake sentence generated with feature-norms

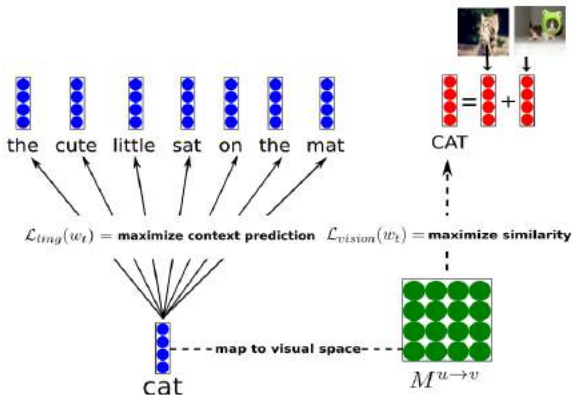
Fake sentence: “crocodile legs
crocodile teeth crocodile green”

Multi-modal fusion techniques

Joint models without aligned texts and images

Extension of (Lazaridou et al. 2015)

- Use deep visual features instead of feature norms
- Use a triplet loss: $\sum_{v^-} \max(0, \gamma - \cos(u, v) + \cos(u, v^-))$



Evaluating word representations (1/3)

→How can we **evaluate** word representations?

Word Similarity

Spearman/Pearson correlation between cosine similarity and human judgement.
SimLex, WordSim, SemSim, VisSim, MEN, ...

word 1	word 2	human judgement	model (cosine)
cat	dog	0.76	0.86
stupid	dumb	0.91	0.52
cloud	book	0.13	-0.11
attach	join	0.56	0.17
advise	baseball	0.04	0.11
...

Evaluating word representations (2/3)

Feature-Norm prediction

- Predicting feature-norms (e.g. *is red*, *can fly*) of concrete objects.
- Features are regrouped in 9 categories: encyclopedic, function, sound, tactile, taste, taxonomic, color, form_and_surface, motion (Collell et al. 2016).
 - Textual embedding better for: encyclopedic, function
 - Visual embedding better for: motion, form_and_surface, color

	has_legs	can_fly	has_teeth	is_green	is_animal	...
monkey	1	0	1	0	1	...
crocodile	1	0	1	1	1	...
plane	0	1	0	0	0	...
...			

Evaluating word representations (3/3)

Concreteness prediction

word	concreteness (human)
dog	0.97
cloud	0.72
war	0.32
hope	0.22

Analogy prediction

Predict analogies

Paris	France	London	England
man	woman	king	queen
do	doing	eat	eating
...

Several benchmarks:

https://github.com/EloiZ/embedding_evaluation

Learning Multi-Modal Word Representation Grounded in Visual Context

Research questions

Learning Multi-Modal Word Representation Grounded in Visual Context,
É. Zablocki, B. Piwowarski, L. Soulier, P. Gallinari, AAAI 2018.

Observation

- Multimodal models have shown **complementarity** of text and language ...
- ...but use **direct features** from objects and ignore **visual context**

Research questions

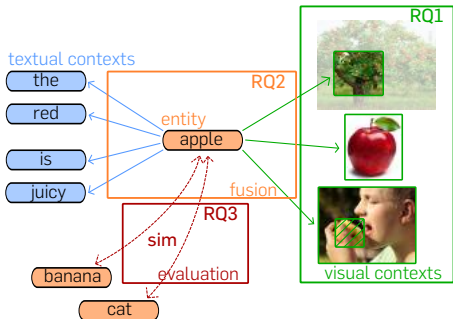
Learning Multi-Modal Word Representation Grounded in Visual Context,
É. Zablocki, B. Piwowarski, L. Soulier, P. Gallinari, AAAI 2018.

Observation

- Multimodal models have shown **complementarity** of text and language ...
- ...but use **direct features** from objects and ignore **visual context**

Research Questions

- **RQ1**: What is a visual context and how can we model it?
- **RQ2**: How can we learn representations jointly from texts and images using contexts?
- **RQ3**: How can we evaluate the contribution of the visual modality to the final embeddings?



Visual skip-gram: model

Recall: Text skip-gram

$$\mathcal{L}_{\text{text}} = - \sum_{e \in \mathcal{D}} \sum_{c \in \mathcal{C}_e} [\log \sigma(t_e^\top \cdot u_c) + \sum_{c^-} \log \sigma(-t_e^\top \cdot u_{c^-})]$$

T, U embedding tables, σ sigmoid function, \mathcal{C}_e set of contexts of entity e

Visual skip-gram: model

Recall: Text skip-gram

$$\mathcal{L}_{\text{text}} = - \sum_{e \in \mathcal{D}} \sum_{c \in \mathcal{C}_e} [\log \sigma(t_e^\top \cdot u_c) + \sum_{c^-} \log \sigma(-t_e^\top \cdot u_{c^-})]$$

T, U embedding tables, σ sigmoid function, \mathcal{C}_e set of contexts of entity e

Idea: Use a skip-gram objective with visual contexts

$$\mathcal{L}_{\text{image}} = - \sum_{e \in \mathcal{D}} \sum_{c \in \mathcal{C}_e} \left[\log \sigma(t_e^\top \cdot f_\theta(c)) + \sum_{c^-} \log \sigma(-t_e^\top \cdot f_\theta(c^-)) \right]$$

- What is \mathcal{C}_e ?
- What is $c \in \mathcal{C}_e$?
- What is $f_\theta(c)$?

Visual skip-gram: instantiation



ENTITY

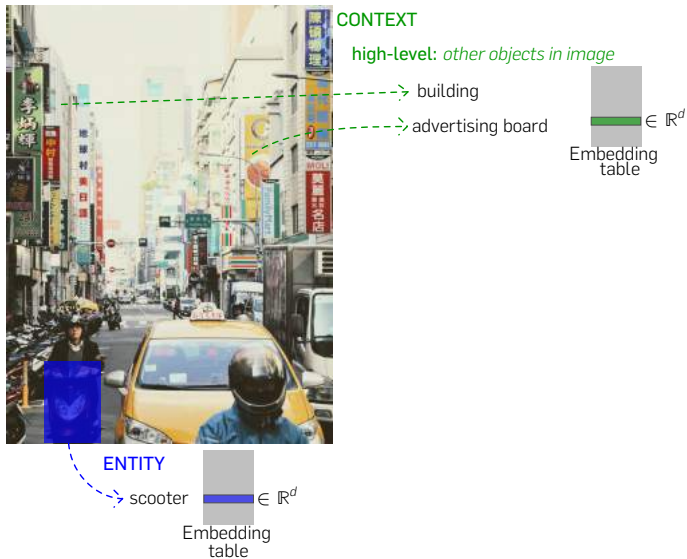
scooter



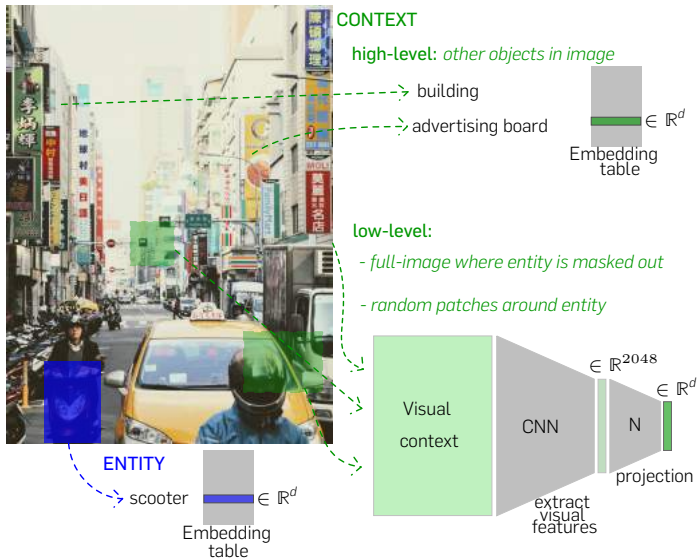
$\in \mathbb{R}^d$

Embedding
table

Visual skip-gram: instantiation



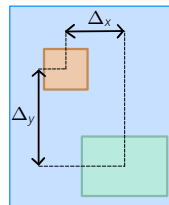
Visual skip-gram: instantiation



Visual skip-gram: spatial information

Can we use **spatial information**?

$$\mathcal{L}_{\text{image}} = - \sum_{e \in \mathcal{D}} \sum_{c \in \mathcal{C}_e} \left[\log \sigma(t_e^T \cdot f_{\theta}(c)) + \sum_{c^-} \log \sigma(-t_e^T \cdot f_{\theta}(c^-)) \right]$$



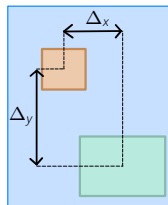
Can we enhance word representation with visual spatial information?

Form a spatial vector $s_{(e,c)}$ and integrate it in $f_{\theta}(c)$ to form $f_{\theta}^D(c, s_{(e,c)})$

Visual skip-gram: spatial information

Can we use **spatial information**?

$$\mathcal{L}_{\text{image}} = - \sum_{e \in \mathcal{D}} \sum_{c \in \mathcal{C}_e} \left[\log \sigma(t_e^T \cdot f_{\theta}(c)) + \sum_{c^-} \log \sigma(-t_e^T \cdot f_{\theta}(c^-)) \right]$$



Can we enhance word representation with visual spatial information?

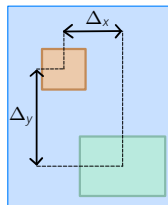
Form a spatial vector $s_{(e,c)}$ and integrate it in $f_{\theta}(c)$ to form $f_{\theta}^p(c, s_{(e,c)})$

$s_{(e,c)}$	spatial information in a 4-d vector, two modelings → Low-level: $s_{(e,c)} = (\Delta_x, \Delta_y, \Delta_{\text{width}}, \Delta_{\text{height}})$ → High-level: $s_{(e,c)} = (\mathbb{1}_{\text{above}}, \mathbb{1}_{\text{below}}, \mathbb{1}_{\text{beside}}, \mathbb{1}_{\text{bigger}})$
-------------	--

Visual skip-gram: spatial information

Can we use **spatial information**?

$$\mathcal{L}_{\text{image}} = - \sum_{e \in \mathcal{D}} \sum_{c \in \mathcal{C}_e} \left[\log \sigma(t_e^T \cdot f_\theta(c)) + \sum_{c^-} \log \sigma(-t_e^T \cdot f_\theta(c^-)) \right]$$



Can we enhance word representation with visual spatial information?

Form a spatial vector $s_{(e,c)}$ and integrate it in $f_\theta(c)$ to form $f_\theta^p(c, s_{(e,c)})$

$s_{(e,c)}$	spatial information in a 4-d vector, two modelings → Low-level: $s_{(e,c)} = (\Delta_x, \Delta_y, \Delta_{\text{width}}, \Delta_{\text{height}})$ → High-level: $s_{(e,c)} = (\mathbb{1}_{\text{above}}, \mathbb{1}_{\text{below}}, \mathbb{1}_{\text{beside}}, \mathbb{1}_{\text{bigger}})$
f_θ^p	integration of spatial feature, two integration methods → Linear: $f_\theta^p(c, s_{(e,c)}) = M \cdot (u_c \oplus s_{(e,c)})$, where $M \in \mathbb{R}^{d \times (d+4)}$ → Bilinear: $f_\theta^p(c, s_{(e,c)}) = s_{(e,c)} M u_c$, where $M \in \mathbb{R}^{4 \times d \times d}$

Multi-modal skip-gram

Model

$$\mathcal{L}(T, U, \theta) = \mathcal{L}_{\text{text}}(T, U) + \alpha \mathcal{L}_{\text{image}}(T, \theta), \text{ where } \alpha \in [0, 1]$$

- T **shared** multi-modal word embeddings
- U **textual** context parameters
- θ **visual** context parameters
- T, U and θ learned with SGD ; α found with cross-validation

Grounding words in visual context: Experiments

Evaluation

→ **Word Similarity**: correlation between cosine similarity and human judgement. e.g. $sim('cat', 'dog') = 0.8$; $sim('cloud', 'book') = 0.1$;

Grounding words in visual context: Experiments

Evaluation

- **Word Similarity:** correlation between cosine similarity and human judgement. e.g. $sim('cat', 'dog') = 0.8$; $sim('cloud', 'book') = 0.1$;
- **Feature-Norm Prediction:** predict objects' attributes from embedding with a linear SVM. e.g. $has_legs('cat') = True$, $is_red('dog') = False$

Grounding words in visual context: Experiments

Evaluation

- **Word Similarity:** correlation between cosine similarity and human judgement. e.g. $sim('cat', 'dog') = 0.8$; $sim('cloud', 'book') = 0.1$;
- **Feature-Norm Prediction:** predict objects' attributes from embedding with a linear SVM. e.g. $has_legs('cat') = True$, $is_red('dog') = False$
- **Concreteness Prediction:** predict words' concreteness with a linear SVM. e.g. $conc('dog')=0.9$, $conc('life')=0.1$

Grounding words in visual context: Experiments

Evaluation

- **Word Similarity:** correlation between cosine similarity and human judgement. e.g. $sim('cat', 'dog') = 0.8$; $sim('cloud', 'book') = 0.1$;
- **Feature-Norm Prediction:** predict objects' attributes from embedding with a linear SVM. e.g. $has_legs('cat') = True$, $is_red('dog') = False$
- **Concreteness Prediction:** predict words' concreteness with a linear SVM. e.g. $conc('dog')=0.9$, $conc('life')=0.1$

Data

- Wikipedia (4.5 million articles)
- Visual Genome (108k images, 30 objets per image)

Multi-modal skip-gram: results

Results

- **Multi-modal embeddings** > text-only embeddings
 - 9% average improvement on all evaluation benchmarks

Multi-modal skip-gram: results

Results

- **Multi-modal embeddings** > text-only embeddings
 - 9% average improvement on all evaluation benchmarks
- **Visual context** (objects surroundings) > Visual features from object
 - 3.2% average improvement on word similarity tasks

Multi-modal skip-gram: results

Results

- **Multi-modal embeddings** > text-only embeddings
 - 9% average improvement on all evaluation benchmarks
- **Visual context** (objects surroundings) > Visual features from object
 - 3.2% average improvement on word similarity tasks
- Visual context is **complementary** to visual features from objects
 - ensemble model performs 6% better

Multi-modal skip-gram: results

Results

- **Multi-modal embeddings** > text-only embeddings
 - 9% average improvement on all evaluation benchmarks
- **Visual context** (objects surroundings) > Visual features from object
 - 3.2% average improvement on word similarity tasks
- Visual context is **complementary** to visual features from objects
 - ensemble model performs 6% better
- **Spatial information** helps to improve word representations.
 - and all spatial modelings and integration perform equally

Multi-modal skip-gram: results

Results

- **Multi-modal embeddings** > text-only embeddings
 - 9% average improvement on all evaluation benchmarks
- **Visual context** (objects surroundings) > Visual features from object
 - 3.2% average improvement on word similarity tasks
- Visual context is **complementary** to visual features from objects
 - ensemble model performs 6% better
- **Spatial information** helps to improve word representations.
 - and all spatial modelings and integration perform equally
- **High-level context** > low-level context
 - 1% average improvement on all tasks

Multi-modal skip-gram: results

Results

- **Multi-modal embeddings** > text-only embeddings
 - 9% average improvement on all evaluation benchmarks
- **Visual context** (objects surroundings) > Visual features from object
 - 3.2% average improvement on word similarity tasks
- Visual context is **complementary** to visual features from objects
 - ensemble model performs 6% better
- **Spatial information** helps to improve word representations.
 - and all spatial modelings and integration perform equally
- **High-level context** > low-level context
 - 1% average improvement on all tasks
- **Joint model** > Sequential CCA
 - 5% average improvement

Conclusion and Perspectives

Conclusion

- Images can help NLP ...
- ...for richer semantics to captured in word representations ...
- with knowledge about visual contexts and spatial organization.

Perspectives

- Evaluation on downstream tasks
- Leverage other perceptual information sources (audio, KB, ...)
- Grounded sentences representations (videos)

Changing perspectives

Two complementary approaches

- Images can help for learning NLP representations

Changing perspectives

Two complementary approaches

→ Images can help for learning NLP representations

Conversely,

→ Using NLP representations can help Computer Vision

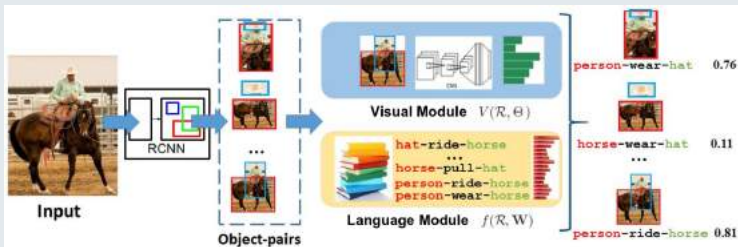
Changing perspectives

Two complementary approaches

- Images can help for learning NLP representations
- Conversely,
 - Using NLP representations can help Computer Vision

NLP can assist computer vision, e.g. when no/few supervision available

- Visual relationship detection



Visual Relationship Detection with Language Priors (from [Lu et al. 2016])

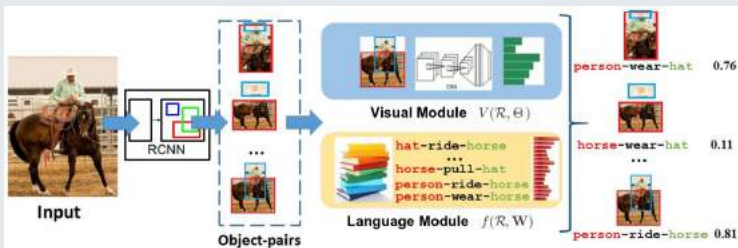
Changing perspectives

Two complementary approaches

- Images can help for learning NLP representations
- Conversely,
 - Using NLP representations can help Computer Vision

NLP can assist computer vision, e.g. when no/few supervision available

- Visual relationship detection



Visual Relationship Detection with Language Priors (from [Lu et al. 2016])

- Zero-shot learning

Zero-Shot Recognition with Semantic Representations and Context

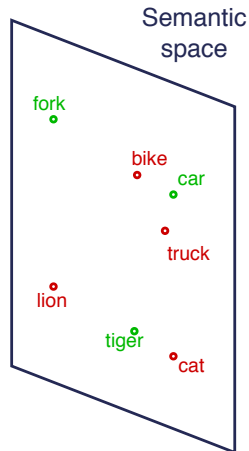
Zero-Shot Learning —Task & Hypothesis

Two domains ...

- **Source domain \mathcal{S}**
Classes with training instances
- **Target domain \mathcal{T}**
Zero-shot classes

...in a common space

- Attributes
- Semantic representation



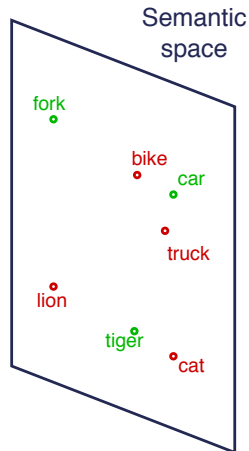
Zero-Shot Learning —Task & Hypothesis

Two domains ...

- **Source domain \mathcal{S}**
Classes with training instances
- **Target domain \mathcal{T}**
Zero-shot classes

...in a common space

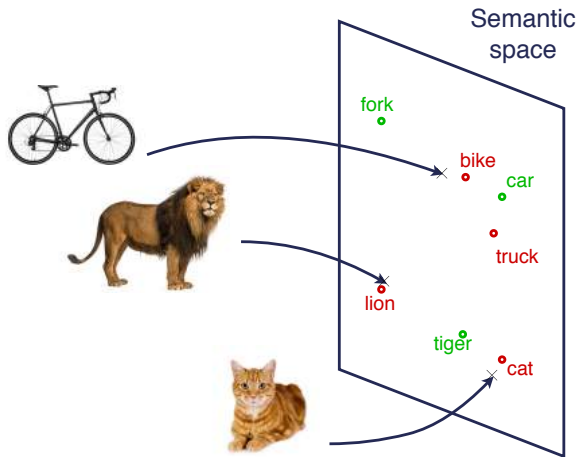
- Attributes
- Semantic representation



Maximize likelihood of $P(i, \mathcal{V})$ → null probability to classes of the target domain

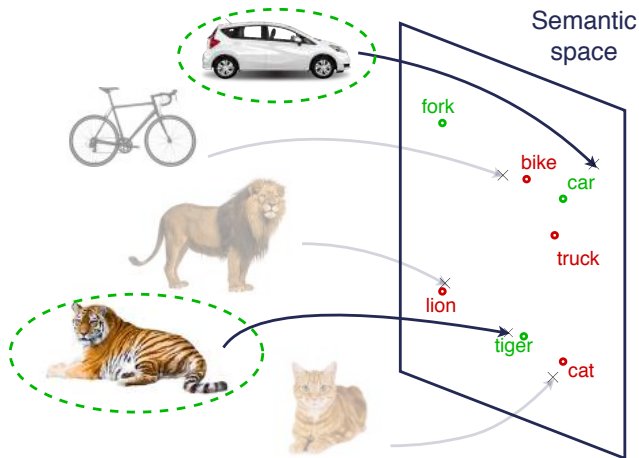
Idea: project images to the semantic space (Socher et al. 2013)

Zero-Shot Learning — Training



$$\mathcal{L}_V = \mathbb{E}_{i, \mathcal{V} \sim P^*} \mathbb{E}_{j \sim \mathcal{U}} [\gamma - f(\mathcal{V})^T w_i + f(\mathcal{V})^T w_j]_+$$

Zero-Shot Learning — Inference



$$i^* = \arg \min_{i \in \mathcal{T}} \cos(w_i, f(\mathcal{V}))$$

Motivation

→ Can language help visual understanding?

Motivation

→ Can language help visual understanding?

	Vision	Language
Region of Interest	Classification	Zero-Shot learning

Motivation

- Can language help visual understanding?
- Can visual context help visual understanding?

	Vision	Language
Region of Interest	Classification	Zero-Shot learning
Visual context	Context-aware classification	

Motivation

- Can language help visual understanding?
- Can visual context help visual understanding?

	Vision	Language
Region of Interest	Classification	Zero-Shot learning
Visual context	Context-aware classification	????

Research Question

Can **language** be leveraged along with **visual context** to help visual understanding?

→ Context-aware Zero-Shot Learning

Zero-Shot Learning — Adding context ?

- “A **tiger** is most recognizable for its **dark vertical stripes** on **reddish-orange** fur with a **lighter underside**”
- “I saw two **tigers** in the **savanna**, one was sleeping under a **tree**, the other one was hunting an **antelope**”
- “We went to the **zoo**! An **animal keeper** was feeding a **tiger** with a **dead boar**. The **tiger** was shredding the poor animal with its **huge teeth** -_-”



Overview

Context-Aware Zero-Shot Learning for Object Recognition,
É. Zablocki, P. Bordes, B. Piwowarski, L. Soulier, P. Gallinari, ICML 2019.

Objective

Simultaneously use:

- \mathcal{V} : The visual area containing an object of interest
- \mathcal{C} : The visual context surrounding the object of interest

Hypothesis

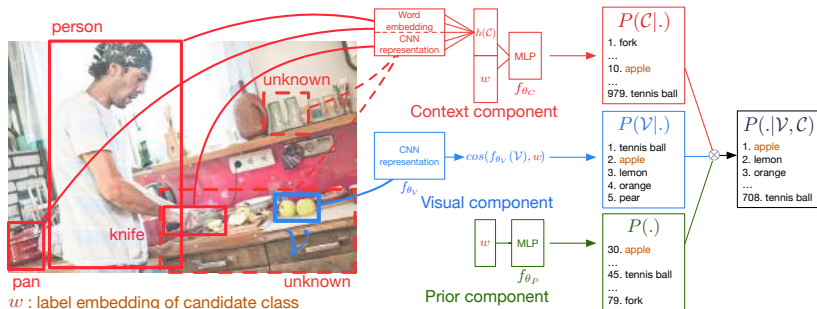
\mathcal{V} and \mathcal{C} conditionally independent given class i :

$$(\mathcal{V} \perp\!\!\!\perp \mathcal{C}) \mid i$$

Model

$$P(i \mid \mathcal{V}, \mathcal{C}) \propto P(\mathcal{V} \mid i)P(\mathcal{C} \mid i)P(i)$$

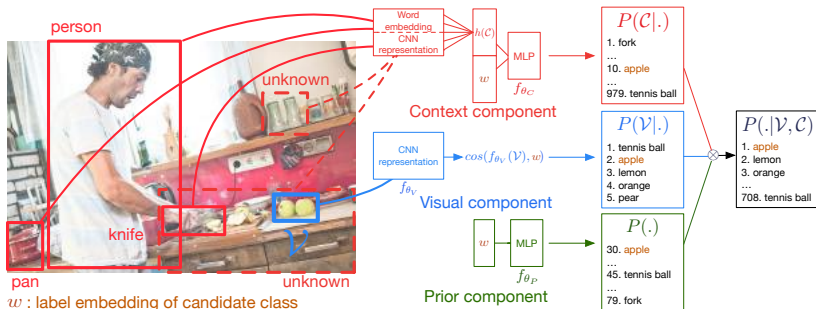
Overview



Model

$$P(i | \mathcal{V}, C) \propto P(\mathcal{V} | i)P(C | i)P(i)$$

Model's components



For a class i , and its semantic representation w_i :

Visual component

$$\log P(V | i; \theta_V) \propto \cos(f_{\theta_V}(V), w_i) := \log \tilde{P}_{visual}$$

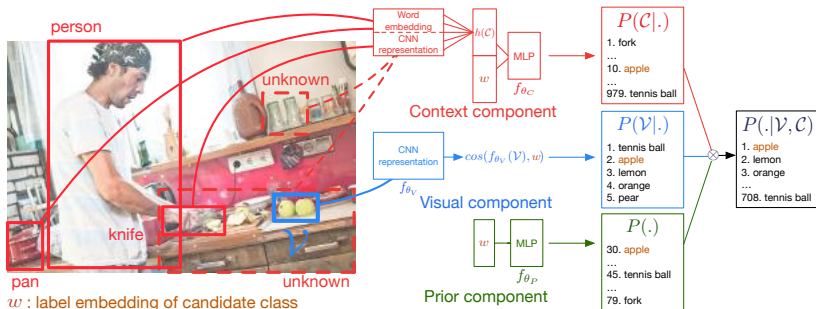
$$f_{\theta_V}(V) = W_V \cdot \text{CNN}(V) + b_V$$

Prior component

$$\log P(i; \theta_P) \propto f_{\theta_P}(w_i) := \log \tilde{P}_{prior}$$

$$f_{\theta_P}(w_i) = \text{MLP}_{\theta_P}(w_i)$$

Model's components

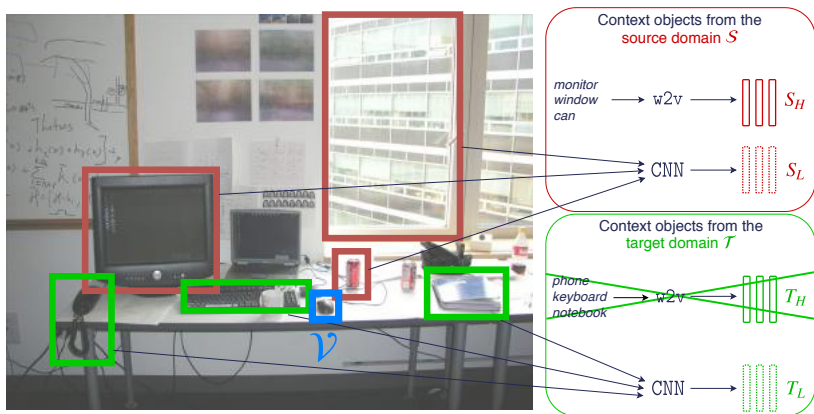


Context component

$$\log P(C | i; \theta_C) \propto f_{\theta_C}(h(C) \oplus w_i) := \log \tilde{P}_{context}$$

$$f_{\theta_C} = \text{MLP}$$

Modeling contextual information



$$h(C_{S_H \cap T_L}) = \underbrace{\begin{matrix} S_H & T_L \\ \text{[red boxes]} & \text{[green boxes]} \end{matrix}}_{\text{Average}} \rightarrow \text{[vertical bar]}$$

Positive examples sampled from true **data distribution** P^*

Prior component

Negative sampled in the **uniform distribution** U

$$\mathcal{L}_P = \mathbb{E}_{i \sim P^*} \mathbb{E}_{j \sim U} [\gamma_P - f_{\theta_P}(w_i) + f_{\theta_P}(w_j)]_+$$

Hence,

$$P^*(i) > P^*(j) \Rightarrow f_{\theta_P}(i) > f_{\theta_P}(j)$$

Learning

Positive examples sampled from true **data distribution** P^*

Prior component

Negative sampled in the **uniform distribution** U

$$\mathcal{L}_P = \mathbb{E}_{i \sim P^*} \mathbb{E}_{j \sim U} [\gamma_P - f_{\theta_P}(w_i) + f_{\theta_P}(w_j)]_+$$

Hence,

$$P^*(i) > P^*(j) \Rightarrow f_{\theta_P}(i) > f_{\theta_P}(j)$$

Context and Visual component

Negative sampled in the true **data distribution** P^*

$$\mathcal{L}_V = \mathbb{E}_{i, \mathcal{V} \sim P^*} \mathbb{E}_{j \sim P^*} [\gamma_V - f_{\theta_V}(\mathcal{V})^\top w_i + f_{\theta_V}(\mathcal{V})^\top w_j]_+$$

$$\mathcal{L}_C = \mathbb{E}_{i, \mathcal{C} \sim P^*} \mathbb{E}_{j \sim P^*} [\gamma_C - f_{\theta_C}(\mathcal{C}, w_i) + f_{\theta_C}(\mathcal{C}, w_j)]_+$$

Hence,

$$P^*(\mathcal{V} | i) > P^*(\mathcal{V} | j) \Rightarrow f_{\theta_V}(\mathcal{V})^\top w_i > f_{\theta_V}(\mathcal{V})^\top w_j$$

$$P^*(\mathcal{C} | i) > P^*(\mathcal{C} | j) \Rightarrow f_{\theta_C}(\mathcal{C}, w_i) > f_{\theta_C}(\mathcal{C}, w_j)$$

Inference

$$\log P(\mathcal{C} | i) \propto \log \tilde{P}_{context}$$

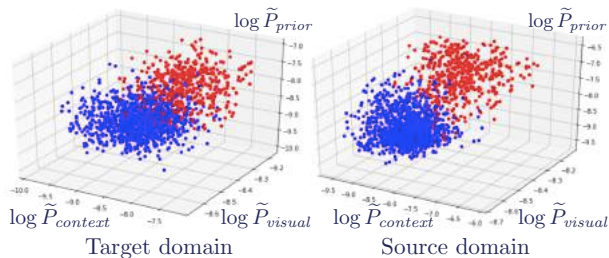
$$\log P(\mathcal{V} | i) \propto \log \tilde{P}_{visual}$$

$$\log P(i) \propto \log \tilde{P}_{prior}$$

Hypothesis

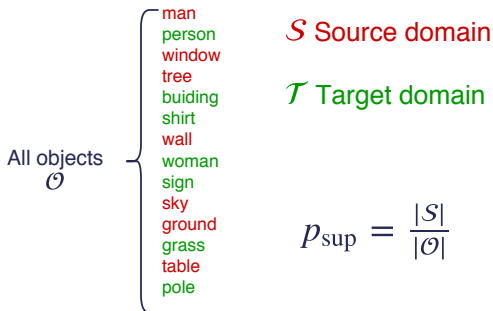
The normalization constants are scalars, independent of the class i , hence:

$$P(i | \mathcal{V}, \mathcal{C}) = \underbrace{(\tilde{P}_{context})^{\alpha_C}}_{P(\mathcal{C} | i)} \cdot \underbrace{(\tilde{P}_{visual})^{\alpha_V}}_{P(\mathcal{V} | i)} \cdot \underbrace{(\tilde{P}_{prior})^{\alpha_P}}_{P(i)}$$

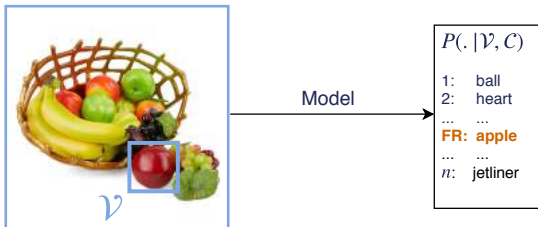


Visual genome

- 108k images
 - 3.8M instances, 105K unique objects, 31 instances per image
- Filtered (> 10 occurrences, need a w2v embedding)
 - 3.4M instances, 4.8K unique objects
- Highly imbalanced:
 - 10% of most frequent classes → 84% of instances



Evaluation



Metric

For a collection of M test images:

$$\text{MFR} = \frac{1}{M} \sum_{i=1}^M \frac{2 * (\text{FR}_i - 1)}{n - 1}$$

- Lower the better
- Random: MFR = 100%
- Perfect model: MFR = 0%

Importance of context

Model	p_{sup} Probability	Target domain \mathcal{T}			Source domain \mathcal{S}		
		10%	50%	90%	10%	50%	90%
<i>Random</i>	\mathcal{U}	100	100	100	100	100	100
$M(\emptyset)$	$P(\cdot)$	38.6	23.7	13.8	12.0	10.6	11.2
$M(\mathcal{V})$	$P(\mathcal{V} \cdot)P(\cdot)$	20.5	10.7	6.0	1.5	2.6	3.6
$M(\mathcal{C})$	$P(\mathcal{C} \cdot)P(\cdot)$	28.7	14.4	9.1	4.2	4.3	4.4
$M(\mathcal{C}, \mathcal{V})$	$P(\mathcal{C} \cdot)P(\mathcal{V} \cdot)P(\cdot)$	18.1	9.0	5.2	1.1	1.9	2.4
$\delta_{\mathcal{C}}$ (%)		11.6	16.4	12.1	23.7	27.3	31.5

Table 1: MFR scores (%). $\delta_{\mathcal{C}}$ is the relative improvement of $M(\mathcal{C}, \mathcal{V})$ over $M(\mathcal{V})$.

Results

- Context is useful: $M(\mathcal{C})^{\text{MFR}} < \text{Random}$
- Visual frequency can be learned from textual semantics: $M(\emptyset)$ generalizes
- Context is complementary to visual appearance: $M(\mathcal{C}, \mathcal{V})^{\text{MFR}} < M(\mathcal{V})$

Modeling contextual information

Context objects in \mathcal{T}		Context objects in \mathcal{S}	\mathcal{T}	\mathcal{S}
		No context	10.72	2.64
		Full image	9.19	2.13
\emptyset		image	9.01	2.05
image		\emptyset	9.00	2.13
\emptyset		label	8.96	1.92
image		image	8.60	1.93
image		label	8.52	1.86
image		image & label	8.31	1.79

Increasing supervision helps:



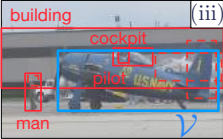
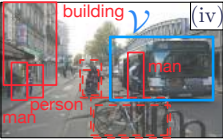
- 1 When bounding boxes are available...
- 2 When context objects are labeled...
- 3 When more context objects are used...
- 4 When low-level info. is used complementarily

Best model:


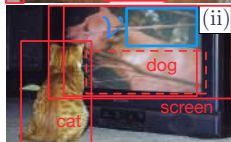
→ 22% relative improvement in the target domain \mathcal{T}

(32% in \mathcal{S})

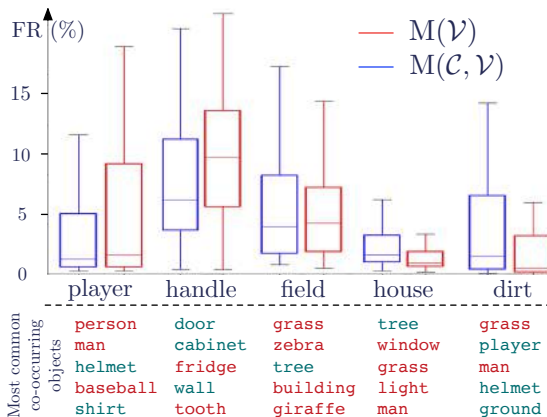
Qualitative results

 <p>(i)</p> <p>vase water stems grass</p>	$P(C .)$ <ol style="list-style-type: none"> lilies flower garden carnations orchids 	$P(V .)$ <ol style="list-style-type: none"> needle fingertip kitten ... 223. flower 	$P(.)$ <ol style="list-style-type: none"> tree woman car ... 8. flower 	$P(. V,C)$ <ol style="list-style-type: none"> flower tip hair ... 29. needle
 <p>(ii)</p> <p>clouds lake</p>	$P(C .)$ <ol style="list-style-type: none"> water river sea ... 9. boat 	$P(V .)$ <ol style="list-style-type: none"> filters connector indicator ... 1757. boat 	$P(.)$ <ol style="list-style-type: none"> tree woman car ... 25. boat 	$P(. V,C)$ <ol style="list-style-type: none"> boat ... 31. water ... 376. filters
 <p>(iii)</p> <p>building cockpit pilot man</p>	$P(C .)$ <ol style="list-style-type: none"> jetliner engines air airplane rotor 	$P(V .)$ <ol style="list-style-type: none"> flight KLM jetliner ... 11. airplane 	$P(.)$ <ol style="list-style-type: none"> tree ... 256. airplane ... 1717. jetliner 	$P(. V,C)$ <ol style="list-style-type: none"> airplane flight runway air hand
 <p>(iv)</p> <p>building man person man</p>	$P(C .)$ <ol style="list-style-type: none"> men woman widow ... 752. bus 	$P(V .)$ <ol style="list-style-type: none"> bus business parking shops downtown 	$P(.)$ <ol style="list-style-type: none"> tree woman car ... 13. bus 	$P(. V,C)$ <ol style="list-style-type: none"> bus car building woman vehicle

Negative results

 <p>(i)</p>	$P(C .)$ <ol style="list-style-type: none"> 1. freezer 2. ovens 3. heater 4. ... 981. books 	$P(V .)$ <ol style="list-style-type: none"> 1. shelving 2. books 3. bookshelf 4. cartons 5. papers 	$P(. V)$ <ol style="list-style-type: none"> 1. books 2. boxes 3. wall 4. shelving 5. table 	$P(. V,C)$ <ol style="list-style-type: none"> 1. table 2. room 3. boxes ... 40. books
 <p>(ii)</p>	$P(C .)$ <ol style="list-style-type: none"> 1. specks 2. whiskers 3. paws 4. ... 1242. leaves 	$P(V .)$ <ol style="list-style-type: none"> 1. sole 2. mold 3. scrape 4. leaves 5. branch 	$P(. V)$ <ol style="list-style-type: none"> 1. tree 2. canopy 3. hand 4. wall 5. leaves 	$P(. V,C)$ <ol style="list-style-type: none"> 1. hand 2. wall 3. nose ... 74. leaves

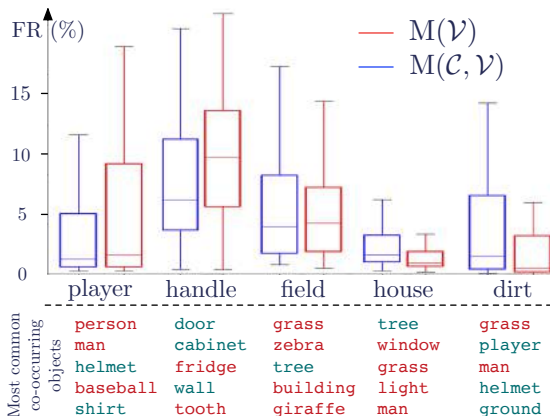
Qualitative results



Context can help ...

- refining predictions
- disambiguating generic shapes

Qualitative results



Context can help ...

- refining predictions
- disambiguating generic shapes

...but not always

- Objects with lots of possible context
- $\rho(\#(\text{context}), \delta_c) = -0.33$

Conclusion and Perspectives

Conclusion

- Language can help to bridge the gap of missing supervision in computer vision
- Visual contexts and priors are contained within word representation

Conclusion and Perspectives

Conclusion

- Language can help to bridge the gap of missing supervision in computer vision
- Visual contexts and priors are contained within word representation

Perspectives

- Learn spatial arrangement from language data
- Zero-shot object detection

Conclusion

Images can help NLP

- to learn grounded representations
- model the semantic gap between language and images?
- useful for downstream tasks? (open-domain QA, machine translation, ...)

Language can help Computer Vision

- using a semantic space (build with NLP techniques)
- evaluating NLP models (captioning, VQA)
- low/few supervision (e.g. zero-shot learning)

Questions?

—Thanks for your attention —

Questions ?