

## CONCEPTUAL GROUNDING FOR TEXT REPRESENTATION LEARNING

Journée IA - RI

Monday 2<sup>nd</sup> December, 2019

Laure Soulier





### 1 Text grounding

#### 2 Enhancing text representation with knowledge resources

# 3 Learning Multi-Modal Word Representation Grounded in Visual Context

4 Conclusion and Perspectives

Text grounding

#### Representation learning

#### Goal: Represent the semantic of a word in a vector space



Conceptual grounding for text representation learning

#### Representation learning



#### Encoding word semantics — Applications

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



- → Human reporting bias
  - $\rightarrow~$  We are more likely to report **unusual facts** and **facts with values**

Word	Teraword	Knext	Word	Teraword	Knext
Spoke	11,577,917	372,042	Hugged	610,040	11,453
Laughed	3,904,519	179,395	Blinked	390,692	21,973
Murdered	2,843,529	16,890	Was late	368,922	31,168
Inhaled	984,613	5,617	Exhaled	168,985	4,052
Breathed	725,034	41,215	Was on time	23,997	14

Figure 1: N-gram frequencies for various verbal events and the number of times Knext learns that "A person may <x> ..."

#### Human Reporting Bias [Gordon, Van Durme, 2013]

The frequency at which objects, relations, or events occur in natural language are significantly different from their real-world frequency.



- → Human reporting bias
  - $\rightarrow~$  We are more likely to report **unusual facts** and **facts with values**...
  - $\rightarrow \ ... \ while unlikely to mention something expected or trivial facts$

#### Example

 $\rightarrow$  When we say:

"Hand me the salt"



- → Human reporting bias
  - $\rightarrow~$  We are more likely to report **unusual facts** and **facts with values**...
  - $\rightarrow \ ... \ while unlikely to mention something expected or trivial facts$

#### Example

 $\rightarrow$  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 Conceptual grouthget and the table base. It is usually extracted from



- $\rightarrow$  Word co-occurences do not capture all grammatical peculiarities
  - → Confusion between the notions of semantic similarity and conceptual association [Hill2015a]

#### Example

- [car, bike]: similar because common physical features, common function, or same category
- [car, petrol]: functional relationship, associated because they frequently occur together in space and language.



- $\rightarrow$  Word co-occurences do not capture all grammatical peculiarities
  - → Confusion between the notions of semantic similarity and conceptual association [Hill2015a]
  - → Embeddings fail to detect synonyms/antonyms [mrksic2016, Mohammad:2008]

east	expensive	British
west	pricey	American
north	cheaper	Australian
south	costly	Britain
southeast	overpriced	European
northeast	inexpensive	England

#### Critical in certain domain applications

E.g., dialog tracking for restaurant booking (expensive vs. cheap)



- $\rightarrow$  Word co-occurrences do not capture all grammatical peculiarities
  - → Confusion between the notions of semantic similarity and conceptual association [Hill2015a]
  - → Embeddings fail to detect synonyms/antonyms [mrksic2016, Mohammad:2008]
  - → Embeddings conflate the contextual evidence of different meanings of a word into a single vector [lacobacciPN15]

<b>bank</b> <sub>1</sub> <sup>n</sup> (geographical)	<b>bank</b> <sub>2</sub> <sup>n</sup> (financial)	<i>number</i> <sup>n</sup> <sub>4</sub> (phone)	$number_3^n$ (acting)	<b>hood</b> <sub>1</sub> <sup>n</sup> (gang)	$hood_{12}^n$ (convertible car)
upstream <sup><math>r</math></sup> <sub>1</sub>	commercial_bank_1^n	$calls_1^n$	appearing $_{6}^{v}$	tortures <sup>n</sup> <sub>5</sub>	$taillights_1^n$
downstream <sup>r</sup> <sub>1</sub>	financial_institution <sup><math>n</math></sup>	dialled <sup><math>v</math></sup> <sub>1</sub>	minor_roles <sup><math>n</math></sup>	vengeance <sub>1</sub> <sup>n</sup>	$grille_2^n$
$runs_6^v$	national_bank <sub>1</sub> <sup>n</sup>	$operator_{20}^n$	stage_production <sup><math>n</math></sup>	$badguy_1^n$	$bumper_2^n$
$confluence_1^n$	$trust\_company_1^n$	$telephone_network_1^n$	supporting_roles <sup><math>n</math></sup>	$brutal_1^a$	fascia <sup>n</sup> <sub>2</sub>
$river_1^n$	savings_bank <sub>1</sub> <sup><math>n</math></sup>	$telephony_1^n$	$leading_roles_1^n$	$execution_1^n$	rear_window <sub>1</sub> <sup><math>n</math></sup>
stream <sup><math>n</math></sup>	$banking_1^n$	subscriber <sub>2</sub> <sup>n</sup>	stage_shows <sub>1</sub> <sup><math>n</math></sup>	murders <sub>1</sub> <sup>n</sup>	headlights <sup>n</sup>

#### Text grounding





#### Contributions



#### General objectives

#### Incorporating common-sense and word knowledge in text representations

- Exploring the potential of grounded on NLP & IR tasks, not to compete with all representation learning baselines
- → Knowledge-empowered text representations
  - $\rightarrow$  Incorporating word senses through concepts in knowledge resources
  - $\rightarrow~$  Leveraging word association through concept relations

Collaboration with IRIT: Gia-Hung Nguyen, Lynda Tamine, Nathalie Souf - ESWC 2018 & ACM TOIS 2019

- → Visual-grounded word representations
  - $\rightarrow$  Grounding words in images to improve word description
  - $\rightarrow~$  Leveraging the visual context to incorporate word functionalities
  - → Extension to sentence representations

#### MĽA

Work with Eloi Zablocki, Patrick Bordes, Benjamin Piwowarski, Patrick Gallinari (CHIST-ERA MUSTER project) - AAAI 2018 & EMNLP 2019

Conceptual grounding for text representation learning

Enhancing text representation with knowledge resources

#### **Motivations**





pêche (fruit vs pêcher)

#### **Motivations**





#### Sémantique relationnelle

- Niveau représentation : expansion de représentations lexicales (Navigli and Velardi 2003; Pal et al. 2014 ; Xiong and Callan 2015a)
- Niveau appariement : inférence logique sur les graphes de concepts (Koopman et al. 2016; Xiong and Callan, 2015b)

#### \* Sémantique distributionnelle

- Niveau représentation : représentation distribuée LSA (Deerwester et al., 1990), PLSA (Hofmann, 1990), WordEmbedding (Mikolov et al., 2013; Pennington et al., 2014)
- Niveau appariement : réseau de neurones siamois, convolutif (Huang et al., 2013 ; Guo et al., 2016 ; Mitra et al., 2017)

#### Offline vs. online representation learning



#### Offline learning

## Retrofitting text embeddings [Faruqui2015, mrksic2016]



Figure 1: Word graph with edges between related words showing the observed (grey) and the inferred (white) word vector representations.

$$\mathcal{L} = \sum_{i=1}^{n} \left[ lpha_i \left\| q_i - \hat{q}_i 
ight\|^2 + \sum_{(i,j) \in \mathcal{E}} eta_{ij} \left\| q_i - q_j 
ight\|^2 
ight]$$

## Online learning

Joint learning of word, concept, text embeddings [Liu2016b, DBLP:journals/corr/ManciniCIN16]



$$\mathcal{L} = -\log\left(\rho\left(w_{t}|W^{t}, S^{t}\right)\right) - \sum_{s \in S_{t}}\log\left(\rho\left(s|W^{t}, S^{t}\right)\right)$$

#### Conceptual grounding for text representation learning

#### Previous work: Offline models

- → Learning word representations
- $\rightarrow~$  Updating word vectors with constraints in the knowledge base

#### Retrofitting Word Vectors to Semantic Lexicons »

[Faruqui et al., 2015]





#### Previous work: Offline models



- → Learning word representations
- ightarrow Updating word vectors with constraints in the knowledge base

 « Counter-fitting Word Vectors to Linguistic constraints » [Mrksic et al., 2016]

Attraction des synonymes



$$\mathrm{SA}(V') = \sum_{(u,w)\in S} \tau \left( d(\mathbf{v}'_u, \mathbf{v}'_w) - \gamma \right)$$

Abrogation des antonymes

$$\operatorname{AR}(V') = \sum_{(u,w)\in A} \tau \left( \delta - d(\mathbf{v}'_u, \mathbf{v}'_w) \right)$$

#### Previous work: Online models



- $\rightarrow$  Joint learning of word and concept representations
- → revisited PV-DM

#### Improving Lexical Embeddings with Semantic Knowledge »

[Yu and Dredze, 2014]



#### Previous work: Online models



- → Joint learning of word and concept representations
- → revisited PV-DM

#### « RC-Net: a general framework for incorporating knowledge into word representations » [Xu et al., 2014]



#### Previous work: Online models



- $\rightarrow\,$  Joint learning of word and concept representations
- → revisited PV-DM
- « Embedding Words and Senses Together via Joint Knowledge-Enhanced training » [Mancini et al., 2017]



#### Previous work



Knowledge-enhanced representation learning					
Joi	nt learning of e	embeddings (On	line learning)		
Embeddings learning	word embedding	concept embedding	document embedding	with relational constraints	
[Liu et al., 2016]	x			X	
[Yu and Dredze, 2014]	x			X	
[Jauhar et al., 2015]	x			X	
[Liu et al., 2018]	х	Х		X	
[Mancini et al., 2016]	х	х			
[Cheng et al., 2015]	х	х			
[Yamada et al., 2016]	Х	Х			
Our model	x	X	х	X	
R	Retrofitting of embeddings (Offline learning)				
[Faruqui et al., 2014]	Х			X	
[Glavas et Vulic, 2018]	Х			X	
[Mrksic et al., 2016]	Х			X	
[Jauhar et al., 2015]	X			X	

#### Online representation learning with knowledge resources



- $\rightarrow$  Learning documents embeddings using PV-DM model
- → Associating words with concepts/synsets



$$\Psi(D) = \frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} \frac{1}{|d_{wc}|} \sum_{w_i \in d_{wc}} [\log p(w_i | w_{i \pm W}, c_{i \pm W}, d) + \log p(c_i | w_{i \pm W}, c_{i \pm W}, d) - \frac{\gamma}{|d_{wc}|} ||\hat{d}_{wc}||^2]$$
(1)

#### Online representation learning with knowledge resources





$$\sum_{d \in D} \sum_{w_t \in d} \left[ \log P(w_t | w_{t \pm k}, c_{t \pm k}, d) + \log P(c_t | w_{t \pm k}, c_{t \pm k}, d) - \frac{\gamma}{|d|} \|\vec{d}\|^2 \right] \\ + \alpha_W \sum_{(w_i w_j) \in \mathcal{R}_W} \operatorname{sim}\left(\overrightarrow{w_i}, \overrightarrow{w_j}\right) \\ + \alpha_C \sum_{(c_k, c_l) \in \mathcal{R}_W} \operatorname{sim}\left(\overrightarrow{c_k}, \overrightarrow{c_l}\right)$$



Data sats			
Data sets	ROBUST	OHSUMED	TREC Med
Туре	site de nouvelles	articles médicaux	rapport de visites médicales
# documents	~528 000	~348 000	~17 000
# queries	250	63	35
Sample of query	"Best Retirement Country"	"adult respiratory distress syndrome"	"patients with gastroesophageal reflux disease"

Resources	DBpedia	MeSH
# extracted concepts	250 000	18 000
used relations	"gold:hypernym"	"IS-A"



	TF - IDF	représentation classique en RI	
les de ence	AWE	moyenne des représentations de mots du document	
/lodè référ	AWE <sub>R</sub>	moyenne des représentations de mots ajustées par une relation [Faruqui et al. 2015]	
2	D2V	représentations de document apprises par ParagraphVector [Le and Mikolov., 2014]	

hes	Tâches TALN	Similarité des phrases/ documents	Classification des phrases
Tâc	Tâches RI	Réordonnancement des documents	Expansion de la requête

· Similarité de phrases



22/40



	TF - IDF	représentation classique en RI	
les de ence	AWE	moyenne des représentations de mots du document	
1odèl référe	AWE <sub>R</sub>	moyenne des représentations de mots ajustées par une relation [Faruqui et al. 2015]	
2 -	D2V	représentations de document apprises par ParagraphVector [Le and Mikolov., 2014]	

hes	Tâches TALN	Similarité des phrases/ documents	Classification des phrases
Tâc	Tâches RI	Réordonnancement des documents	Expansion de la requête

Ré-ordonnancement des documents

$$RSV(q, d) = \alpha \cdot IRScore(q, d) + (1 - \alpha) \cdot NeuralScore_{KB}(q, d)$$

Représentations distribuées de la requête et du document → Même niveau de granularité



	TF - IDF	représentation classique en RI	
les de ence	AWE	moyenne des représentations de mots du document	
Aodè référ	AWE <sub>R</sub>	moyenne des représentations de mots ajustées par une relation [Faruqui et al. 201	
<	D2V	représentations de document apprises par ParagraphVector [Le and Mikolov., 2014]	

hes	Tâches TALN	Similarité des phrases/ documents	Classification des phrases
Tâc	Tâches RI	Réordonnancement des documents	Expansion de la requête

• Expansion de la requête

m : mot ou concept candidat

$$p(m|q*) = \alpha p_{mle}(m|q) + (1-\alpha)p_{emb}(m|q)$$

$$p_{emb}(m|q) = \frac{\sigma(\vec{e}_m, \vec{q})}{\sum_{m' \in V} \sigma(\vec{e}_{m'}, \vec{q})} \xrightarrow[\text{Représentations districtions districtions}]{Représentations districtions}$$

Représentations distribuées de la requête et des éléments candidats → Niveaux de granularité différents

#### Results



#### Sentence relatedness / classification tasks

→ Improvements for out-of-domain datasets

Training dataset			Ohsumed				TREC Med					
Eval. benchmark	SUBJ	MPQA	TREC	MRPC	SUBJ	MPQA	ATREC	MRPC	SUBJ	MPQA	ATREC	MRPC
TF-IDF	72.13	68.45	79.98	69.12	33.13	25.35	31.48	30.32	22.55	21.99	21.48	19.75
AWE	73.10	68.04	79.52	68.05	32.50	25.74	32.12	29.35	21.92	21.87	20.51	19.25
$AWE_{KB}$	75.71	69.08	81.91	68.75	35.61	26.63	34.18	31.20	22.63	22.45	21.24	20.60
D2V	73.52	69.35	79.30	70.81	31.56	25.01	32.07	28.23	21.55	21.58	20.68	18.56
SD2V <sub>off</sub>	76.15	72.11	79.50	$71.90^{\Delta}_{\bullet}$	32.65	25.40	32.15°	27.88°	22.00	21.18	20.92	18.21 <sup>•</sup>
$SD2V_{on}$	75.44	70.89 <b>•</b>	79.56	$72.04^{\bullet}$	32.99	25.53	32.57 <b>•</b>	28.80°	21.81	21.38	21.00	$18.15^{\bullet}$

#### Results



	Document re-ranking							Query expansion							
	Robust		Ohsumed		TREC Med		Robust		Ohsumed		TREC Med				
	MAP	%Chg	MAP	%Chg	MAP	%Chg	MAP	%Chg	MAP	%Chg	MAP	%Chg			
BM25	0.251		0.2147		0.312		0.251		0.2147		0.312				
AWE	0.250	-0.40%	0.201	-2.24%	0.349 <sup>•</sup>	+11.83%	0.250	-0.40 %	0.252°	+17.51%	0.289 <sup>•</sup>	-7.08%			
$AWE_{KB}$	0.251	+0.00%	0.301°	+40.20%	0.350°	+12.24%	0.251	0.00%	0.254 <sup>•</sup>	+18.30%	0.2901°	-7.02%			
D2V	0.2505	-0.20%	0.300°	+39.78%	0.356°	+14.07%	0.2511	+0.04%	0.255°	+19.19%	0.291 <sup>•</sup>	-6.67%			
SD2Voff	0.251	0.00%	0.3018	+40.57%	0.3591	$^{\Delta}+15.10\%$	0.2464	-1.83%	0.258 <sup>•</sup>	+20.17%	0.3205	<sup>4</sup> +2.72%			
SD2Von	0.2507	-0.12%	0.302°	+40.66%	0.3554°	+13.91%	0.2443	-2.67%	0.2599	<sup>∆</sup> +21.05%	0.2889°	-7.40%			

Table 7. Comparison of online/online learning approaches on IR evaluation tasks: re-ranking and query expansion. Metric: MAP. %*Chg* denotes the effectiveness improvement of models w.r.t. BM25.  $\Delta$ : significance test of offline vs. online models. •: significance test of BM25 vs. our model.

#### IR tasks

- $\rightarrow\,$  Improvements vs. BM25 are important in the medical field (between 15% and 40%)
- → Higher growth rates on re-ranking task (40.6% vs. 21.16% for query expansion)

→ Better performance compared to neural models (+3%) Conceptual grounding for text representation learning

#### Results: Validating relational constraints



#### Relational constraints

- $\rightarrow$  C1: word-word relations.
- → C2: concept-concept relations.

Training dataset	Robust	Ohsumed	TREC Med
SD2V <sub>on</sub>	0.14	0.02	0.02
SD2VReg <sub>on</sub>	0.18	0.14	0.14

Table 1: Validation of the relational constraints C1 and C2 on pivotal words. Metric: P@10.

#### Results



Training dataset	Robust				Ohsumed				TREC Med			
Eval. benchmark	SUBJ	MPQA	ATREC	MRPC	SUBJ	MPQA	ATREC	MRPC	SUBJ	MPQA	ATREC	MRPC
SD2V <sub>off</sub>	76.15	72.11	79.50	71.90	32.65	25.40	32.15	27.88	22.00	21.18	20.92	18.21
SD2VReg <sub>off</sub>	77.02	72.78	80.89	•73.37•	33.83	26.34	33.12	28.77	23.02	21.99	22.48	18.93
SD2VIns <sub>off</sub>	76.41	72.00	80.20	73.77 <b>•</b>	34.37	•27.42	•34.18	•29.31•	22.94	20.94	21.32	$19.44^{\bullet}$
SD2Von	75.44	70.89	79.56	72.04	32.99	25.53	32.57	28.80	21.81	21.38	21.00	18.15
$SD2VReg_{on}$	75.69	71.09	80.01	73.62 <b>•</b>	34.33°	27.28	33.21	29.44	22.78	22.74	•22.03	19.61 <sup>•</sup>
SD2VIns <sub>on</sub>	76.78•	71.13	79.43	72.27	33.31	26.08	32.76	28.96	22.40	22.05	20.34	19.28°

Table 14. Comparison of the relation constraint strategies on the on classification tasks (SentEval benchmark). Metric: Accuracy (%). %*Chg* denotes the effectiveness improvement of models w.r.t. BM25.  $\Delta$ : significance test of regularization vs. training instance. •: significance test w.r.t model scenario without relation.

- → Constraining the learning with relational knowledge is effective in both NLP and IR tasks.
- → The learning leverages from both word-word relations and concept-concept relations

#### Limitations and Perspectives



- → (Main) Pending issues and Perspectives
  - → Robustness of the models: significant performance variation depending on multiple factors (knowledge resource, task, annotation quality, etc.)
  - $\rightarrow$  Exploiting contextual representation learning models (Elmo, BERT, ...)
  - $\rightarrow$  Considering the relation types in the learning objective

Learning Multi-Modal Word Representation Grounded in Visual Context

#### Multi-modal fusion techniques





 Images
 Texts

 Word7
 word7

 word7
 word7

 word7
 word7

 word7
 sim

 word7
 sim

 word7
 sim

 (b) Middle fusion
 (c) Late fusion

Sequential models

- → Joint models: aligned/non aligned text and images (skip-gram extension, grounded models)
- → Sequential models: combination of text representations (e.g., GloVe or Word2Vec) and image representations (pre-trained CNN):
  - → Middle fusion: form multi-modal representations (e.g. concatenation, CCA)
  - → Late fusion: interaction in the downstream task (e.g. linear combination of scores)

Conceptual grounding for text representation learning

#### Example of sequential technique (middle fusion)





Sequential model (Collell et al. 2017)

#### Observation

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

#### Research questions



#### **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

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

*T*, *U* embedding tables,  $\sigma$  sigmoid function,  $C_e$  set of contexts of entity e

#### Visual skip-gram: model



#### Recall: Text skip-gram

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

*T*, *U* embedding tables,  $\sigma$  sigmoid function,  $C_e$  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 the context  $c \in C_e$ ?
- → What is the context modeling function  $f_{\theta}(c)$ ?

#### Visual skip-gram: instantiation





#### Visual skip-gram: instantiation





#### Visual skip-gram: instantiation





#### 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
- $\rightarrow \theta$  visual context parameters
- $\rightarrow~$  T, U and  $\theta$  learned with SGD ;  $\alpha$  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;
- → 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

#### Baselines

- → Text only
- → Text + direct visual features from objects
- → Text + visual contexts (sequential CCA)

#### Data

- → Wikipedia (4.5 million articles)
- → Visual Genome (108k images)

#### Multi-modal skip-gram: results



#### Results

- → Multi-modal embeddings > text-only embeddings
  - $\rightarrow$  9% average improvement on all evaluation benchmarks
- → Visual context (objects surroundings) > Visual features from object
  - $\rightarrow$  3.2% average improvement on word similarity tasks
- → Visual context is **complementary** to visual features from objects
  - → ensemble model performs 6% better
- → High-level context > low-level context
  - $\rightarrow$  1% average improvement on all tasks
- → Joint model > Sequential CCA
  - → 5% average improvement

## Ongoing work: grounded sentences representations



## We distinguish two sources of information

- → Cluster information: implicit knowledge that sentences associated with the same video refer to the same underlying reality
- → Perceptual information: high-level information extracted from a video using a pre-trained CNN

To preserve textual semantics and to avoid an over-constrained textual space: grounded space which partially transfer the structure of visual space in the textual one.



## Conclusion and Perspectives

### **Conclusion and Perspectives**



## External resources can help NLP/IR tasks

- → ...for richer semantics/common sense captured in text representations
- → Investigating the the semantic gap/reporting bias between language and images/knowledge resources
- $\rightarrow$  Evaluating on real tasks (open-domain QA, machine translation, ...)

#### Language can help building external resources

- → Knowledge base completion
- → Commonsense mining

#### Language can help Computer Vision

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

#### Open question



#### Language Models as Knowledge Bases?

Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, Sebastian Riedel

(Submitted on 3 Sep 2019 (v1), last revised 4 Sep 2019 (this version, v2))

Recent progress in pretraining language models on large textual corpora led to a surge of improvements for downstream NLP tasks. W these models may also be storing relational knowledge present in the training data, and may be able to answer queries structured as " Language models have many advantages over structured knowledge bases: they require no schema engineering, allow practitioners to relations, are easy to extend to more data, and require no human supervision to train. We present an in-depth analysis of the relational (without fine-tuning) in a wide range of state-of-the-art pretrained language models. We find that (i) without fine-tuning, BERT contai competitive with traditional NLP methods that have some access to oracle knowledge, (ii) BERT also does remarkably well on open-don supervised baseline, and (iii) certain types of factual knowledge are learned much more readily than others by standard language model surprisingly strong ability of these models to recall factual knowledge without any fine-tuning demonstrates their potential as unsuper The code to reproduce our analysis is available at this https URL.

Comments: accepted at EMNLP 2019 Subjects: Computation and Language (cs.CL) Cite as: arXiv:1909.01066 [cs.CL] (or arXiv:1909.0106602 [cs.CL] for this version)





## —Thanks for your attention —

## Questions ?

Conceptual grounding for text representation learning