

# Grounding language in action and perception

Simon Dobnik

31 European Summer School in Logic, Language and Information 12 August 2019, Riga Latvia



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

### Previously



- Theoretical questions related to language, action and perception
- Examples of image descriptions and challenges for their computational modelling
- Semantics of spatial language and computational models
- Experiment for vision and language
- Example of an image captioning system using DNNs



### Previously



- Theoretical questions related to language, action and perception
- Examples of image descriptions and challenges for their computational modelling
- Semantics of spatial language and computational models
- Experiment for vision and language
- Example of an image captioning system using DNNs
- TODAY: Computational systems for learning and generating grounded language



#### Outline



Words as classifiers

Words as classifiers + formal representations

Grounding a language model

Tutorial #2

The power of a language model

Tutorial #3

Summary



(日) (四) (문) (문) (문)



# Words as classifiers



▲□▶ ▲圖▶ ▲目▶ ▲目▶ 目 のへで

## Perceptual grounding





5/69

#### Words as classifiers

#### (Roy, 2002)

A scene:





centre for linguistic theory

<ロ> <四> <四> <日> <日> <日</p>

#### Training instances



- Vector of real-valued features representing the objects in the scene: r, g, b, hw\_ratio, area, x, y mm\_dimension
- Natural language descriptions
  - ► The pink square
  - The light blue square
  - The brightest green rectangle
  - The purple rectangle to the left of the pink square
  - The narrow purple rectangle below and to the right of the blue square



#### Learning



- Cluster words into classes (using probabilistic models):
  - Based on word distributions: words within a class co-occur infrequently with other words in that class
  - Their association with semantic features
  - A combination of both
- Statistical bi-gram model of classes (encodes word order constraints)



#### Generation



- For each bi-gram sequence of classes, for each class choose the most likely word given the target object:
  - ► the, the rectangle, the green rectangle, the large green rectangle, the large light green rectangle ...
- Estimate the fit of each description to the target object: the likelihood of a sequence of words to refer to the features of the object
- Contextual constraints and ambiguity of a description: ψ(Q) = fit(x<sub>target</sub>, Q) − max<sub>∀x≠target</sub> fit(x, Q)
- Combine the scores from syntactic and contextual constraints with a weighted sum



#### Generation



9/69

- For each bi-gram sequence of classes, for each class choose the most likely word given the target object:
  - ► the, the rectangle, the green rectangle, the large green rectangle, the large light green rectangle ...
- Estimate the fit of each description to the target object: the likelihood of a sequence of words to refer to the features of the object
- Contextual constraints and ambiguity of a description: ψ(Q) = fit(x<sub>target</sub>, Q) − max<sub>∀x≠target</sub> fit(x, Q)
- Combine the scores from syntactic and contextual constraints with a weighted sum
- Evaluated by 3 human judges to select the target object given a description:

< ロ > < 同 > < 回 > < 回 >

- human-generated: 89.8%
- machine-generated: 81.3%

#### Words as classifiers with robots



(Dobnik, 2009)

Difficult to write a semantic model of spatial descriptions.



(Logan and Sadler, 1996)



#### Words as classifiers with robots



#### (Dobnik, 2009)

Difficult to write a semantic model of spatial descriptions.



(Logan and Sadler, 1996)

- Difficult to pre-define a model of the world for a robot (cf. Shrdlu (Winograd, 1976))
- In SLAM (Dissanayake et al., 2001) the robot learns its environment incrementally through observations.



### Words as classifiers with robots



#### (Dobnik, 2009)

Difficult to write a semantic model of spatial descriptions.



(Logan and Sadler, 1996)

イロト イポト イヨト イヨ

- Difficult to pre-define a model of the world for a robot (cf. Shrdlu (Winograd, 1976))
- In SLAM (Dissanayake et al., 2001) the robot learns its environment incrementally through observations.
- Let's learn language by combining observations of the environment with the ways humans describe it!

#### The robot



Are robot's representations sufficient to learn spatial language?



#### The robot



11/69

- Are robot's representations sufficient to learn spatial language?
- ATRV-JR mobile robot (iRobot) primarily used for tasks such as map building, localisation and navigation and runs the following components:
  - odometry component: provides information about the robot's motion, for example (R-Heading) and (Speed);
  - SLAM component: localises the robot on a 2-dimensional map consisting of a set of points relative to some random starting point, for example (0.6234,0.2132).

< ロ > < 同 > < 回 > < 回 >

#### From sensors to spatial geometry



#### Human



Film: 3D laser point cloud generation



Robot



- A robot an a human describer situated in a room.
- The robot is guided manually by another person.
- Human describers (4) freely generate descriptions from the perspective of the robot (speech recognition).
- ► All linguistic and non-linguistic observations are logged.



## Linguistic descriptions

UNIVERSITY OF GOTHENBURG

Descriptions of robotic motion

- You're going forward slowly.
- Now you're turning right.

#### Descriptions of relations between objects

- The chair is to the left of you.
- The table is further away than the chair.



#### Supervised off-line learning

- Automatic extraction of useful information from datasets
- Matching observations in time ... noise
- A lot of perceptual data for few descriptions





### Supervised off-line learning

- Automatic extraction of useful information from datasets
- Matching observations in time ... noise
- A lot of perceptual data for few descriptions
- Some extracted features:

. . .

R-Head	ling	${\sf Speed}$	Ver	b	
0.001		0.234	mov	ving	
0.535		0.122	turr	ning	
0.123		-0.364	reve	ersing	
LO_x	LO_y	/ REI	=0_x	REFO_y	Relation
0.632	0.53	6 0.0	01	0.321	to_the_right_of
0.212	0.44	7 0.34	46	0.342	to_the_right_of
0.573	0.73	1 0.50	64	0.632	near

< ロ > < 同 > < 回 > < 回 >

15/69

## Applying the classifiers

UNIVERSITY OF GOTHENBURG

Two interactive systems:

pDescriber: generates descriptions (of objects, of robot's motion)



# Applying the classifiers

UNIVERSITY OF GOTHENBURG

16/69

Two interactive systems:

- pDescriber: generates descriptions (of objects, of robot's motion)
- pDialogue: generates motion and answers user's questions
  - M Motion requests: Go forward slowly. Go forward right fast.
  - A Locating objects: Where is the table? The table is to the left of the chair? Where are you? I'm behind the sofa.
  - B Confirming object description: Is the table to the left of the chair? Yes, the table is to the left of the chair. No, the table is near the chair.
  - C Finding objects: What is to the left of the chair? The pillars, the tyres and the wall are to the left of the chair.
  - D Referencing an object: What is the chair to the left of? The chair is to the left of the table, the desk and the wall.

< ロ > < 同 > < 回 > < 回 >

#### Was learning successful?

- Performance of the classifiers vs subjective opinion
- One context vs different contexts
- New room, 5 subjects (pDescriber), 13 subjects (pDialogue)
  - pDescriber: Is this a good description? Yes/No.
  - ▶ pDialogue: How natural is the answer? 1 to 5.

Question type	Accuracy (%)	Classifier	Accuracy (%)
pDescriber	59.28	relation	69.12
A	43.51	relation	69.12
В	54.17	relation	69.12
С	54.70	lo_x	48.80
		lo_y	72.80
D	56.92	refo_x	65.60
		refo_y	82.24
Mean	52.33		67.71
			CLASP Inguistic theory and studies in prot
		4 🗖	



# Words as classifiers + formal representations



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

# Connecting grounded and formal representations

- (Harnad, 1990): language in a domain of symbolic computation
- Not true but there are benefits of this view from the computational perspective
- Top-down filtering of bottom-up induced knowledge
- "Deeper" cognitively inspired representations vs representations induced from patterns
- Can examine and interpret the beliefs obtained
- Possibilities of using techniques for logical inference with machine learning methods
- The need for an over-arching framework for perception and language

イロト イポト イヨト イヨ

19/69

(Dobnik and Kelleher, 2017)

# NL, formal language and grounding

# Parsing natural language to a robot control language and grounding it in action and perception (Matuszek et al., 2012b)

go to	the	second	junction	and	go left
S/ NP	NP/NP	NP/N	N	S\S/S	S
(move-to forward)	[null]	(do-n-times 2 x)	(until (junction current-loc) y)	(do-seq g f)	(turn-left)
			S\ S		
		(do-n-times 2 (	(do-seq g turn-le ft)		
(do-n-times 2 (			NP until (junction current-loc) y))		
(do-n-times	2 (until (	S junction current-lo	oc) (move-to forward)))		
(do-seq	(do-n-tim	es 2 (until (junctio	S on current-loc) (move-to forward)	)) (turn-left))	



# Type Theory with Records

UNIVERSITY OF GOTHENBURG

Type Theory with Records (Cooper, in prep; Dobnik et al., 2013; Larsson, 2015; Cooper et al., 2015; Dobnik and Cooper, 2017)

- meaning relative to agent
- judgements
- of situations, of speech events (and of neural events)
- meaning representations as record types (and a few basic types)
- types of perceptual readings to types of dialogue game-boards

types are cognitive and intensional

Types of objects



(Dobnik and Cooper, 2017)

#### Perceptual domain

► [[34,24,48],[56,78,114]...]: PointMap PointMap ⊑ list(list(Real))

**Conceptual domain** 



Types of objects

(Dobnik and Cooper, 2017)

#### Perceptual domain

- [[34,24,48],[56,78,114]...]: PointMap  $PointMap \square$  list(list(Real))
- Object detection function  $(\textit{Pointmap} \rightarrow \mathsf{set}(\left[\begin{array}{cc} \mathsf{reg} & : & \textit{Pointmap} \\ \mathsf{pfun} & : & (\textit{Ind} \rightarrow \textit{Tvpe}) \end{array}\right]))$ pfun =  $\lambda x$ : Ind.chair(x)

#### **Conceptual domain**





# Types of objects

(Dobnik and Cooper, 2017)

#### Perceptual domain

- [[34,24,48],[56,78,114]...]: PointMap  $PointMap \square$  list(list(Real))
- Object detection function

 $pfun = \lambda x$ : Ind. chair(x)

Individuation function

 $\lambda r: \begin{vmatrix} \operatorname{reg:} Pointmap \\ pfun: (Ind \rightarrow Type) \end{vmatrix}$ 

 $(Pointmap \rightarrow set(\begin{bmatrix} reg : Pointmap \\ pfun : (Ind \rightarrow Type) \end{bmatrix}))$ 

$$\int_{-\infty}^{\infty} \left[ \begin{array}{ccc} a & : & Ind \\ loc & : & location(a, r.reg) \\ c & : & r.pfun(a) \end{array} \right]$$

< ロ > < 同 > < 回 > < 回 >

Conceptual domain

22 / 69







23 / 69

# (Matsson, Dobnik, and Larsson, 2019): pyTTR for Visual Question Answering (VQA)



A B > A B >

https://github.com/arildm/imagettr

#### Question answering



#### A polar question as a subtype check

Scene

loc <sub>0</sub> :		w : Int				
	cy : Int					
	•	cx : Int	w =	422		
		h : Int	cy =	355		
			cx =	435		
		L	h =	242		
<b>y</b> <sub>4</sub>	:	Indao				
cp1	:	$\langle \lambda a : Ind . person(a), [x_9] \rangle$				
cr3	2	$\langle \lambda a : Ind , \lambda b : Ind , right(a, b), [y_4, x_{11}] \rangle$				
X11	2	Indaz				
cp2	1	$\langle \lambda a : Ind . dog(a), [x_{11}] \rangle$				
cr <sub>6</sub>	1	$(\lambda a : Ind . \lambda b :$	Ind . a	above(a	$(b), [x_9, y_4])$	
X9	:	Inda				
			w	Int		
cli		$\langle \lambda k: Ind , \lambda l:$	cy :	: Int	Investment of the test	
			cx :	: Int	$10cation(k, t), [x_9, toc_1]$	
			h	Int		
cr5	2	$\langle \lambda a : Ind . \lambda b :$	Ind . r	ight(a,	b), $[x_9, x_{11}]$	
cro	1	$\langle \lambda a : Ind . \lambda b :$	Ind . b	below(a	$(b), [y_4, x_9]$	
		and the second sec		1.1.1		

#### Question

 $\begin{array}{c} \begin{tabular}{ccc} x & : & Ind \\ y & : & Ind \\ c_0 & : & kite(x) \\ c_1 & : & person(y) \\ c_2 & : & above(x,y) \end{tabular} \end{tabular} \end{array}$ 





# Grounding a language model



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

Grounding bottom-up without formal representations



#### (Ghanimifard and Dobnik, 2017)

- Compositionality as a strength of formal grammars
- A probabilistic/neural language model learns the associations between words in a sequence
- Is bottom-up grounding compositional?
  - Composed strings of words
  - Composed perceptual representations
- Grounded language model  $Pr(w_{1:T}|c) = \prod_{t=1}^{T} Pr(w_t|w_{1:t}, c)$ where a word sequence  $w_{1:T}$  is a description of an image with

the visual features c.



# Artificial dataset



 Average acceptability/probability scores over locations (Logan and Sadler, 1996)

• 
$$freq(w_{1:T}, c) = n \times Pr(w_{1:T}, c)$$

- Artificial composition
  - Simple language with connectives  $g_{\wedge}$  :  $(v_i, v_i) \rightarrow [v_i, \text{``and''}, v_i]$  $g_{\vee}$  :  $(v_i, v_i) \rightarrow ["either", v_i, "or", v_i]$  $g_{\neg}$  :  $v \rightarrow [$ "not", v] Language with "distractor" words  $g_1$  :  $(v*) \rightarrow [v*]$  $g_2$  :  $(v*) \rightarrow [``it'', ``is'', v*]$  $g_3$  :  $(v*) \rightarrow [``it", ``is", v*, ``the'', ``box'']$  $g_{4}$  :  $(v*) \rightarrow ["the", "ball", "is", v*, "the", "box"]$  $g_5$  :  $(v*) \rightarrow ["the", "object", "is", v*, "the", "box"]$ Locations
### Generated templates







28 / 69

## Grounded neural language model



UNIVERSITY OF GOTHENBURG

ASP Contre for Inguistic theory and studies in probability

29 / 69



The probability of the generated sequence for a particular location is a judgement score

$$\begin{array}{lll} \mathcal{T}_{w_{1:\mathcal{T}}} &=& \{\mathit{Score}_{w_{1:\mathcal{T}},c}\}_{c\in L} \\ \hat{\mathcal{T}}_{w_{1:\mathcal{T}}} &=& \{\mathit{Pr}(w_{1:\mathcal{T}}|c)\}_{c\in L} \\ \rho(\mathcal{T}_{w_{1:\mathcal{T}}}, \hat{\mathcal{T}}_{w_{1:\mathcal{T}}}) & & \text{Spearman's rank correlation coefficient} \end{array}$$



## Evaluation: composition



#### Performance on different datasets

Descriptions	-0%	+Distractors	-10%	-20%	-30%	-80%
One word	0.92	0.91				
Phrases	0.83	0.84				
AND-phrases $*$	0.87	0.85	0.84	0.80	0.78	0.53
OR-phrases*	0.79	0.80	0.74	0.73	0.69	0.38
NEG-phrases*	0.72	0.82				

\*Contains single words and their negations



## Some examples of new compositions



-50%



centre for Inguistic theory and studies in probability

A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A

32 / 69

## Evaluation: decomposition



#### Trained on all descriptions with some removed

Removed	-0%	-10%	-20%	-30%	-40%
–AND-phrases*	0.83	0.86	0.80	0.77	0.81
$-NEG\operatorname{-phrases}^*$	0.83	0.83	0.64	0.59	0.43
$-OR ext{-}phrases^*$	0.83	0.73	0.78	0.68	0.69
-One word	0.92	0.90	0.90	0.84	0.87

\* also excludes single words and their negation



## Evaluation: decomposition



#### Trained on all descriptions with some removed

Removed	-0%	-10%	-20%	-30%	-40%
–AND-phrases*	0.83	0.86	0.80	0.77	0.81
$-NEG\operatorname{-phrases}^*$	0.83	0.83	0.64	0.59	0.43
$-OR ext{-}phrases^*$	0.83	0.73	0.78	0.68	0.69
-One word	0.92	0.90	0.90	0.84	0.87

\* also excludes single words and their negation







## Tutorial #2



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

## Tutorial #2: Learning to compose



by Mehdi Ghanimifard

In this tutorial we look at the code that was used for the paper. You can:

- adjust the descriptions and the spatial templates used to generate the artificial dataset
- remove individual words or composed phrases from the training dataset
- train and evaluate the models on them

Code on Github





# The power of a language model



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

Interacting objects



#### Coventry et al. (2005)





## Degrees of contributions

UNIVERSITY OF GOTHENBURG

The umbrella is over/above the man.



Coventry et al. (2001)



## Extracting knowledge about object interaction



- Encoded in the language model, cf. the success of distributional semantics
- Use the predictions as a filter in description generation
- Predict the bias of a spatial relation to functional or geometric knowledge:
  - A functional spatial relation is more selective of their target and landmark objects
  - A geometric relation will occur with any kind of objects.

(Dobnik and Kelleher, 2013, 2014)



## Corpora of image descriptions



40 / 69



a yellow building with white columns in the background; two palm trees in front of the house; cars parked in front of the house; a woman and a child are walking over the square;

A B > A
 A
 B > A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

## Choosing a relation



FG	Prep	$-2 log \lambda$	$H_2$ vs. $H_1$
people*square	on	655.66*	$2.37 \times 10^{142}$
people*square	in	133.63*	$1.04 \  imes 10^{29}$
people*square	at	1.81	2.47
people*umbrella	with	16.06*	3076.878
boy*umbrella	under	12.16*	436.788
table*umbrella	under	9.39*	109.447
child*umbrella	under	8.35*	65.006
sculpture*umbrella	with	6.88*	31.25
woman*umbrella	with	6.83*	30.428
woman*umbrella	under	6.78*	29.592
girl*umbrella	with	4.59*	9.921
man*umbrella	with	2.29	3.15
child*umbrella	with	1.53	2.153

\*: *p* < 0.05

CL

## (Normalised) entropy and object variation



11	D	ГСТ	TIL			GOTHENBURG
#	Preposition	FG-Types	Tokens	Norm FG ent		
1	on_left_side_of	5	31	0.35448		
2	underneath	31	74	0.65535		
3	in	7584	34846	0.6714		
4	onto	49	86	0.79109		
5	down	83	142	0.81099		
6	over	440	736	0.83106		
7	at	1393	2726	0.83148		
8	on_top_of	61	87	0.83409		
9	against	50	68	0.85171		
10	on	4897	10085	0.852		
11	on_side_of	46	63	0.87644		
15	on_back_of	9	11	0.89489		
16	through	179	245	0.89738		
17	in_front_of	1278	1938	0.90998		
22	under	167	220	0.92096		
23	above	145	190	0.9228		
26	below	13	14	0.96248		
					LASP	centre for linguistic theory and studies in probabili
						ີ 🗉 🔊 🖉

42 / 69

## Neural language models and perplexity





unbalanced plain perplexity average in 10 folds

(Dobnik, Ghanimifard, and Kelleher, 2018)



## Perplexities by NLMs as vectors



centre for linguistic theory and studies in probability

44 / 69

A B > A B >



(Ghanimifard and Dobnik, 2019)

## When language model takes over



45/69

#### Demo from yesterday





< ロ > < 同 > < 回 > < 回 >

#### (Ghanimifard and Dobnik, 2018b)

On the linguistic bias of vision and language datasets (Agrawal et al., 2017)



## Tutorial #3



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

Tutorial #3: bottom-up grounding of vision and lange  $\pi$ 

UNIVERSITY OF GOTHENBURG

#### by Mehdi Ghanimifard

In this tutorial we look at the extension of the code from the previous tutorial that replaces simple spatial locations with visual features that are also trained from the data. It will examine

- how the visual features are trained, represented and used in the model
- how objects are detected
- how image descriptions are generated
- Iimitations and further work in generating image descriptions

Code on Github





# Summary



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで





- Grounding words classifiers
- Integration of grounded words with a language model
  - probabilistic
  - rule-based
- End-to-end grounding with a grounded neural language model
  - Compositionality?
  - Information encoded in language alone
- End-to-end grounding of visual features and a language model



## References I



- Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Aniruddha Kembhavi. 2017. Don't just assume; look and answer: Overcoming priors for visual question answering. *arXiv*, arXiv:1712.00377 [cs.CV]:1–15.
- Robin Cooper. in prep. Type theory and language: from perception to linguistic communication. Draft at https: //sites.google.com/site/typetheorywithrecords/drafts.
- Robin Cooper, Simon Dobnik, Shalom Lappin, and Staffan Larsson.
   2015. Probabilistic type theory and natural language semantics.
   *Linguistic Issues in Language Technology LiLT*, 10(4):1–43.
- Kenny R. Coventry, Angelo Cangelosi, Rohanna Rajapakse, Alison Bacon, Stephen Newstead, Dan Joyce, and Lynn V. Richards. 2005. Spatial prepositions and vague quantifiers: Implementing the functional geometric framework. In Christian Freksa, Markus Knauff, Bernd Krieg-Brückner, Bernhard Nebel, and Thomas Barkowsky, editors, *Spatial Cognition IV. Reasoning, Action, Interaction*, volume 3343 of *Lecture Notes in Computer Science*, pages 98–110. Springer Berlin Heidelberg.

< ロ > < 同 > < 回 > < 回 >

## References II



- Kenny R. Coventry, Mercè Prat-Sala, and Lynn Richards. 2001. The interplay between geometry and function in the apprehension of Over, Under, Above and Below. *Journal of Memory and Language*, 44(3):376–398.
- Robert Dale and Ehud Reiter. 1995. Computational interpretations of the Gricean maxims in the generation of referring expressions. *Cognitive science*, 19(2):233–263.
- M. W. M. G Dissanayake, P. M. Newman, H. F. Durrant-Whyte, S. Clark, and M. Csorba. 2001. A solution to the simultaneous localization and map building (SLAM) problem. *IEEE Transactions on Robotic and Automation*, 17(3):229–241.
- Simon Dobnik. 2009. *Teaching mobile robots to use spatial words*. Ph.D. thesis, University of Oxford: Faculty of Linguistics, Philology and Phonetics and The Queen's College, Oxford, United Kingdom.
- Simon Dobnik and Robin Cooper. 2017. Interfacing language, spatial perception and cognition in Type Theory with Records. *Journal of Language Modelling*, 5(2):273–301.

## References III



< ロ > < 同 > < 回 > < 回 >

- Simon Dobnik, Robin Cooper, and Staffan Larsson. 2013. Modelling language, action, and perception in type theory with records. In Denys Duchier and Yannick Parmentier, editors, Constraint Solving and Language Processing - 7th International Workshop on Constraint Solving and Language Processing, CSLP 2012, Orleans, France, September 13-14, 2012. Revised Selected Papers, number 8114 in Publications on Logic, Language and Information (FoLLI). Springer, Berlin, Heidelberg.
- Simon Dobnik, Mehdi Ghanimifard, and John D. Kelleher. 2018. Exploring the functional and geometric bias of spatial relations using neural language models. In Proceedings of the First International Workshop on Spatial Language Understanding (SpLU 2018) at NAACL-HLT 2018, pages 1–11, New Orleans, Louisiana, USA. Association for Computational Linguistics.

## References IV



Simon Dobnik and John D. Kelleher. 2013. Towards an automatic identification of functional and geometric spatial prepositions. In Proceedings of PRE-CogSsci 2013: Production of referring expressions – bridging the gap between cognitive and computational approaches to reference, pages 1–6, Berlin, Germany.

- Simon Dobnik and John D. Kelleher. 2014. Exploration of functional semantics of prepositions from corpora of descriptions of visual scenes. In Proceedings of the Third V&L Net Workshop on Vision and Language, pages 33–37, Dublin, Ireland. Dublin City University and the Association for Computational Linguistics.
- Simon Dobnik and John D. Kelleher. 2017. Modular mechanistic networks: On bridging mechanistic and phenomenological models with deep neural networks in natural language processing. In *Proceedings of the Conference on Logic and Machine Learning in Natural Language* (*LaML 2017*), *Gothenburg*, 12–13 June 2017, volume 1 of *CLASP Papers in Computational Linguistics*, pages 1–11, Gothenburg, Sweden. Department of Philosophy, Linguistics and Theory of Science

A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

## References V



(FLOV), University of Gothenburg, CLASP, Centre for Language and Studies in Probability.

- Raquel Fernández. 2013. Rethinking overspecification in terms of incremental processing. In *Proceedings of the PRE-CogSci 2013* Workshop on the Production of Referring Expressions.
- Mehdi Ghanimifard and Simon Dobnik. 2017. Learning to compose spatial relations with grounded neural language models. In *Proceedings of IWCS 2017: 12th International Conference on Computational Semantics*, pages 1–12, Montpellier, France. Association for Computational Linguistics.
- Mehdi Ghanimifard and Simon Dobnik. 2018a. Knowing when to look for what and where: Evaluating generation of spatial descriptions with adaptive attention. In *Computer Vision – ECCV 2018 Workshops. ECCV 2018*, volume 11132 of *Lecture Notes in Computer Science (LNCS)*, pages 1–9, Proceedings of the Workshop on Shortcomings in Vision and Language (SiVL), ECCV 2018, Munich, Germany. Springer, Cham.

< ロト < 同ト < ヨト < ヨ)

## References VI



< ロ > < 同 > < 回 > < 回 >

Mehdi Ghanimifard and Simon Dobnik. 2018b. Visual grounding of spatial relations in recurrent neural language models. In Proceedings of the 3rd Workshop on Models and Representations in Spatial Cognition (MRSC-3) at 11th International Conference on Spatial Cognition 2018, pages 1–7, Tübingen, Germany. https://dobnik.net/simon/events/mrsc-3/#programme.

Mehdi Ghanimifard and Simon Dobnik. 2019. What a neural language model tells us about spatial relations. In Proceedings of the Combined Workshop on Spatial Language Understanding (SpLU) and Grounded Communication for Robotics (RoboNLP) at NAACL-HLT 2019, pages 1–13, Minneapolis, Minnesota, USA. North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019), Association for Computational Linguistics.

Stevan Harnad. 1990. The symbol grounding problem. *Physica D*, 42(1–3):335–346.

## References VII



- Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137.
- John D. Kelleher. 2010. How to preposition a robot: A case-study in symbol grounding. Presentation at the University College Dublin Research Seminars.
- John D. Kelleher and Simon Dobnik. 2017. What is not where: the challenge of integrating spatial representations into deep learning architectures. In *Proceedings of the Conference on Logic and Machine Learning in Natural Language (LaML 2017), Gothenburg, 12–13 June,* volume 1 of *CLASP Papers in Computational Linguistics,* pages 41–52, Gothenburg, Sweden. Department of Philosophy, Linguistics and Theory of Science (FLOV), University of Gothenburg, CLASP, Centre for Language and Studies in Probability.
- Emiel Krahmer and Kees van Deemter. 2011. Computational generation of referring expressions: A survey. Computational Linguistics, 38(1):173–218.

イロト イポト イヨト イヨ

## References VIII



- Staffan Larsson. 2015. Formal semantics for perceptual classification Journal of Logic and Computation, 25(2):335–369.
- Gordon D. Logan and Daniel D. Sadler. 1996. A computational analysis of the apprehension of spatial relations. In Paul Bloom, Mary A. Peterson, Lynn Nadel, and Merrill F. Garrett, editors, *Language and Space*, pages 493–530. MIT Press, Cambridge, MA.
- David G Lowe. 1999. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pages 1150–1157. IEEE.
- Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2016. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. *arXiv*, arXiv:1612.01887 [cs.CV]:1–10.
- Arild Matsson, Simon Dobnik, and Staffan Larsson. 2019. ImageTTR: Grounding type theory with records in image classification for visual question answering. In *Proceedings of the IWCS 2019 Workshop on Computing Semantics with Types, Frames and Related Structures,* pages 55–64, Gothenburg, Sweden. Association for Computational Linguistics.

## References IX



Cynthia Matuszek, Nicholas FitzGerald, Luke Zettlemoyer, Liefeng Bo, and Dieter Fox. 2012a. A joint model of language and perception for grounded attribute learning. In *Proceedings of the 29th International Conference on Machine Learning (ICML 2012)*, Edinburgh, Scotland.

- Cynthia Matuszek, Evan Herbst, Luke Zettlemoyer, and Dieter Fox. 2012b. Learning to parse natural language commands to a robot control system. In *Proceedings of the 13th International Symposium on Experimental Robotics (ISER)*.
- Margaret Mitchell, Kees van Deemter, and Ehud Reiter. 2013. Generating expressions that refer to visible objects. In *HLT-NAACL*, pages 1174–1184.
- Deb Roy. 2002. Learning visually-grounded words and syntax for a scene description task. *Computer speech and language*, 16(3):353–385.
- Andrea Vedaldi. 2016. Convolutional networks for computer vision applications. IV&L summer school on vision and language, Malta. http://www.robots.ox.ac.uk/~vedaldi/assets/teach/vedaldi16deepcv.pdf.



- Terry Winograd. 1976. *Understanding Natural Language*. Edinburgh University Press.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. *arXiv*, arXiv:1502.03044 [cs.LG]:1–22.



## Referring to what matters





### From (Fernández, 2013)



60 / 69

## Generating referring expressions (GRE)



The Incremental Algorithm, Dale and Reiter (1995)

- Each target object is associated with certain properties, e.g. Colour, Type and Position
- Each target object is assigned a set of distractor objects having some property in common
- An unambiguous referring expression matches the target object but none of the distractors
- Properties are assigned a preference order (Colour, Type, Position) based on their salience in that domain
- Add properties in this order to the description plan iff a property has a discriminatory power (reduces the set of distractor objects) at that point in the order: "the red"
- Stop if the description is uniquely identifies the target: "the red chair"

イロト イポト イヨト イヨ

## Not just visual properties

UNIVERSITY OF GOTHENBURG

< ロト < 同ト < ヨト < ヨ)

Other factors may influence the amount of information speakers include in a referring description:

- Complexity of the domain: the number of properties
- Cardinality of the target: plural targets are more likely to be over-specified
- Cross-linguistic differences
- Dialogue context and intent

Further reading on GRE: (Krahmer and van Deemter, 2011)
# Rule-based generation (Mitchell et al., 2013) I



63 / 69



TUNA corpus

GRE3D3 corpus

Represented as properties

tg	color:yellow	size:(63,63)	type:ball	loc:right-hand
lm	color:red	size:(345,345)	type:cube	loc:right-hand
obj3	color:yellow	size:(70,70)	type:cube	loc:left-hand

# Rule-based generation (Mitchell et al., 2013) II



- Process through this graph of attributes, calculating the likelihood of generating a property based on its prior α<sub>att</sub> and a description length penalty γ
   f(A ∪ {x}) = γα<sub>att</sub>
- Add the property if f > n where n is random  $0 \le n \le 1$
- $\blacktriangleright$  Scan through objects; if there are more objects like the referent object, generate properties that distinguish them constrained by  $\gamma$

A B > A B >



64 / 69

## Engineering visual features

- Scale Invariant Image Transform (SIFT) features (Lowe, 1999)
- Creating visual words







65 / 69

### Learning visual features with CNNs





A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classif cation with deep convolutional neural networks. In Proc. NIPS, 2012.

Image from Vedaldi (2016)

CLASP ortho for y inglatic theory and studies in probability イロトィボット イヨト モンモンタへ 66 / 69

# Generation of image descriptions with deep learning



(Xu et al., 2015)



GOTHENBUR

## Attention



#### (Xu et al., 2015; Lu et al., 2016)

- Align visual features with words
- Combine the image features V and the hidden state h<sub>t</sub> of the LSTM through a single layer followed by a softmax z<sub>t</sub> = w<sub>h</sub><sup>T</sup> tanh(W<sub>v</sub>V + (W<sub>g</sub>h<sub>t</sub>)1<sup>T</sup>) α<sub>t</sub> = softmax(z<sub>t</sub>)

• 
$$c_t = \sum_{i=1}^k \alpha_{ti} v_i$$

- c<sub>t</sub>: attended visual features at time t
- ▶ i: a region of k regions of an image
- $v_i$ : visual representation of a region *i*
- $\alpha_{ti}$ : the attentional weight to the region *i*
- $c_t$  and  $h_t$  are combined to predict the next word  $y_{t+1}$
- ► Adaptive attention (Lu et al., 2016):

$$\hat{c}_t = \beta_t s_t + (1 - \beta_t) c_t$$

 $s_t$  is obtained from the memory state of the language model.

イロト イポト イヨト イヨ

## Attention and different description types

Attention as spatial templates?



Kelleher and Dobnik (2017); Ghanimifard and Dobnik (2018a)



- Visual attention is high in general: higher with objects than relations
- Spatial relations depend more on the language model
- Spatial relations are attended in less focused way: not geometric relation

69 / 69

< 🗇 🕨