# Natural Language Visual Grounding with Keyword-Aware Attention Network

## Jinpeng Mi

Universität Hamburg
Fakultät für Mathematik, Informatik und Naturwissenschaften
Fachbereich Informatik

**Technische Aspekte Multimodaler Systeme**

8. Januar 2019

# Gliederung

# Natural Language Visual Grounding

▶ task: given a referring expression, localize the referred object or area in an image



referring expression: a glass of water on the table

Universität Hamburg

- ▶ applications: visual understanding systems, dialogue systems, natural language based interaction with intelligence agents, e.g., robots

- ▶ main difficulties:

- - how to learn the correlation between natural language referring expression and visual domain (image region)
- - how to locate the target object (the spatial relationship between objects)

Universität Hamburg

visual grounding is re-formulated three sub-problems:

▶ which words to focus on in a referring expression

▶ where to look in an image

▶ which object to locate

Universität Hamburg

# public datasets

▶ RefCOCO: 19994 images, 142210 expressions

RefCOCO

woman on right in white shirt
woman on right
right woman

▶ RefCOCO+: 19992 images, 141564 expressions

▶ RefCOCOg: 25799 images, 95010 expressions (no test set)

RefCOCOg Val



1 a young boy in a blue shirt
2 a woman in a white shirt
  and black shorts
3 a woman in a white shirt

# Attention Mechanism

- ▶ inspired by how the human visual cortex employs visual attention mechanism to focus on informative regions in visual scenes
- ▶ first proposed in machine translation[1], image captioning[2]
- ▶ type: hard attention and soft attention

[1]Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate, ICLR 2014.
[2]Xu, K., Ba, J., Kiros, ..., Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention, ICML 2015.
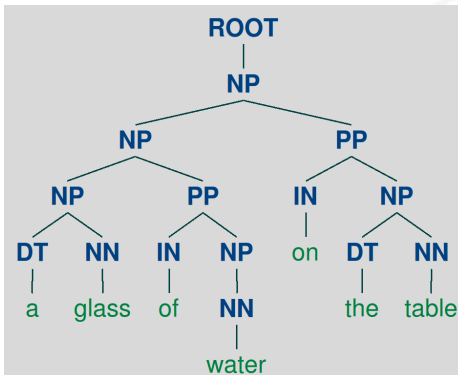
Universität Hamburg

# Architecture

- ▶ which words to focus on
- ▶ where to look in an image
- ▶ which object to locate

Universität Hamburg

# which words to focus on

▶ Syntactic Parsing

# which words to focus on

- ▶ referring expression filtering

filer insignificant words: determiner, coordinating conjunction, "to", interjection, modal words, linking verb, etc.

- ▶ examples

raw: young man with blond hair wearing a white shirt and dark tie in a ballroom
filtered: young man with blond hair wearing white shirt dark tie in ballroom


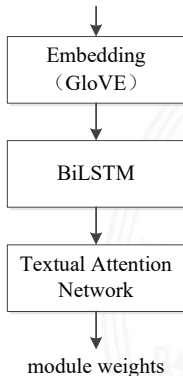raw:a person standing behind a snowboarder with a blue jacket and black pants
filtered:person standing behind snowboarder with blue jacket black pants

Universität Hamburg

# which words to focus on

▶ acquire different weights for different words

glass of water on table

$\downarrow$

```
Embedding
（GloVE）
```

$\downarrow$

```
BiLSTM
```

$\downarrow$

```
Textual Attention
Network
```

$\downarrow$

module weights

UH
Universität Hamburg

# which words to focus on

▶ deep representation of a referring expression

$$e_t = embedding(w_t), t \in [1, T] \tag{1}$$

$$\overrightarrow{h_t} = BiLSTM(e_t, \overrightarrow{h_{t-1}}) \tag{2}$$

$$\overleftarrow{h_t} = BiLSTM(e_t, \overleftarrow{h_{t-1}}) \tag{3}$$

$$h_t = [\overrightarrow{h_t}, \overleftarrow{h_t}] \tag{4}$$

where $T$ is the length of a filtered referring expression.

# which words to focus on

▶ Textual Attention Network

$$u_t = tanh(W_w h_t + b_w) \tag{5}$$

$$\alpha_t = \frac{exp(u_T^t \beta_w)}{\sum_t^T exp(u_T^t \beta_w)} \tag{6}$$

$$r_t = FC(\alpha_t \odot h_t) \tag{7}$$

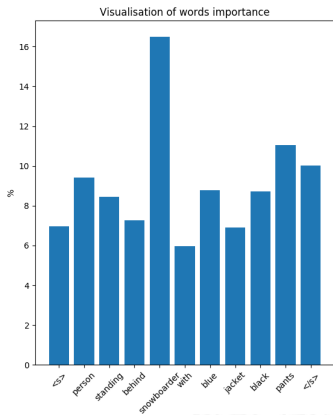where $W_w$, $b_w$ and $\beta_w$ are trainable vectors, $r_t$ is calculated weights, $\odot$ denotes element-wise production.

* Yang Z, Yang D, Dyer C, et al. Hierarchical attention networks for document classification. Proceedings of NAACL-HLT 2016.

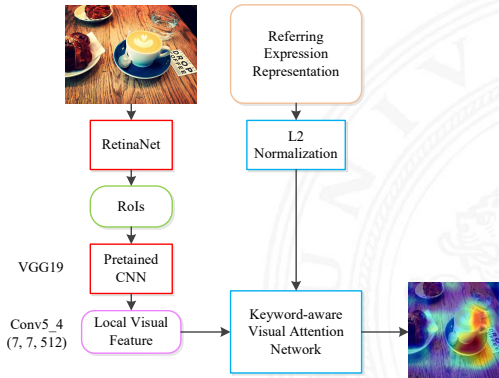# which words to focus on

► result



Visualisation of words importance

Universität Hamburg

# which words to focus on



Visualisation of words importance

Universität Hamburg

# where to look in an image

► Keyword-aware Visual Attention Network

# where to look in an image

$$v' = Conv(v) \tag{8}$$

$$s_r = f(W_s r_t + b_s) \tag{9}$$

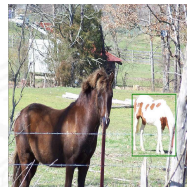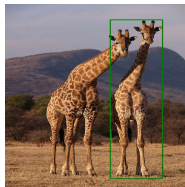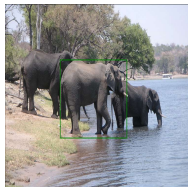$$M_{atten} = softmax(s_r \odot v') \tag{10}$$

where $v'$ denotes projected feature map, $f$ is non-linear function, $W_s$ and $b_s$ are trainable vectors, $M_{atten}$ is generated attention map, $\odot$ denotes element-wise production.
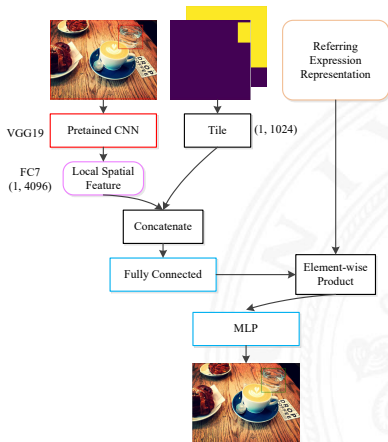
Universität Hamburg

# where to look in an image

Universität Hamburg

# where to look in an image

Universität Hamburg

# which object to locate

Universität Hamburg

# ToDo

- ▶ debug and train
- ▶ adjust parameters
- ▶ improve architecture
- ▶ grasping experiments on PR2

Universität Hamburg

# Thank you for your attention!