# Visualizing and Understanding Neural Machine Translation

Yanzhuo Ding, Yang Liu, Huanbo Luan, Maosong Sun
Tsinghua University

Presented by: Yuchen He

source words 我 喜欢 温哥华 </s>

source word embeddings

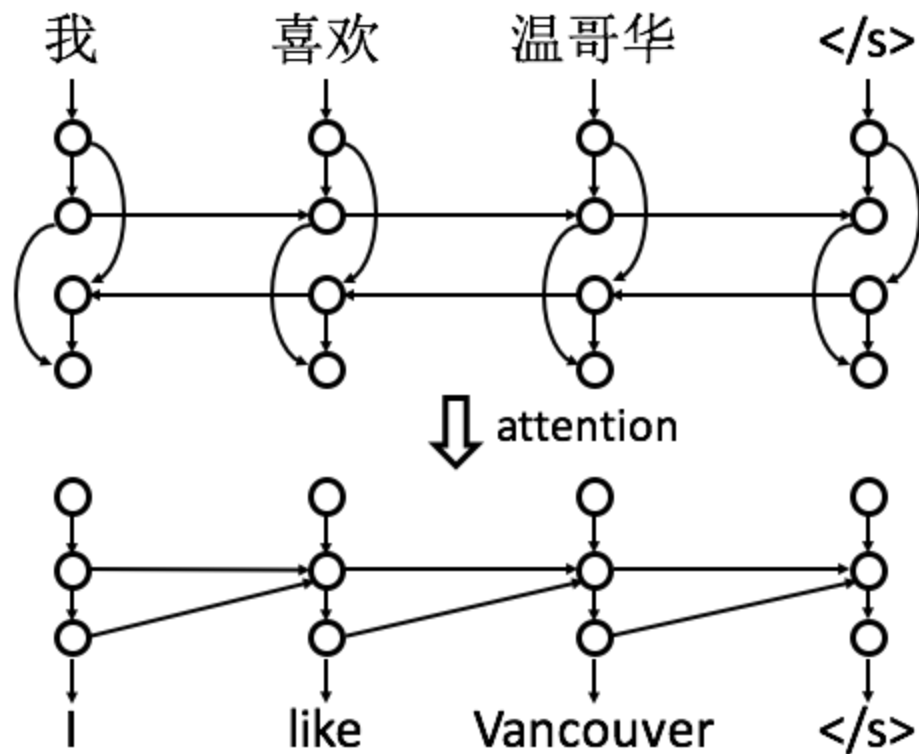source forward hidden states
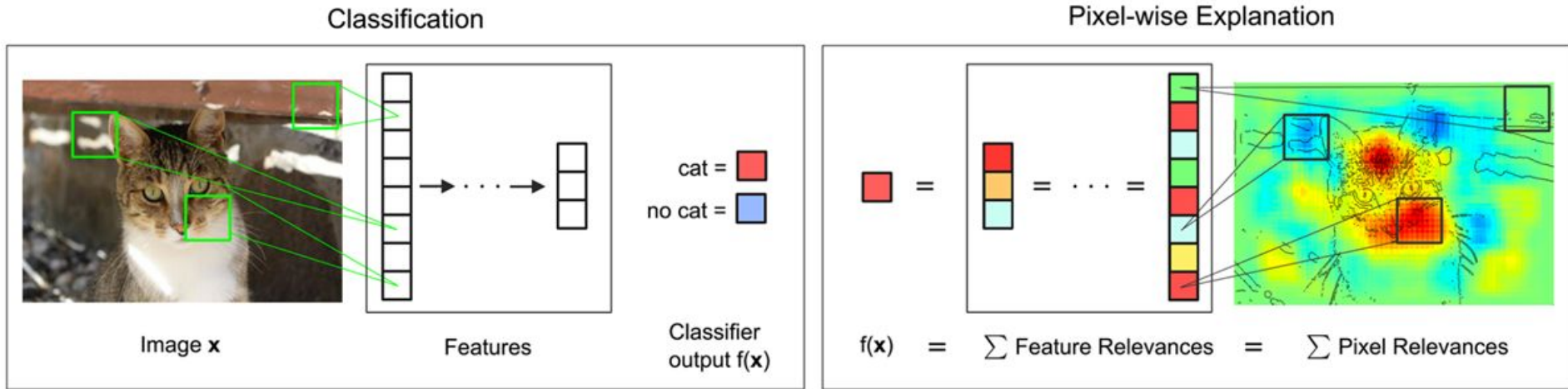
source backward hidden states

source hidden states

attention

source contexts

target hidden states

target word embeddings

target words I like Vancouver </s>

# Layer-wise relevance propagation (LRP)



Classification | Pixel-wise Explanation

Image **x** | Features | Classifier output f(**x**)

cat = 🟥
no cat = 🟦

f(**x**) = ∑ Feature Relevances = ∑ Pixel Relevances

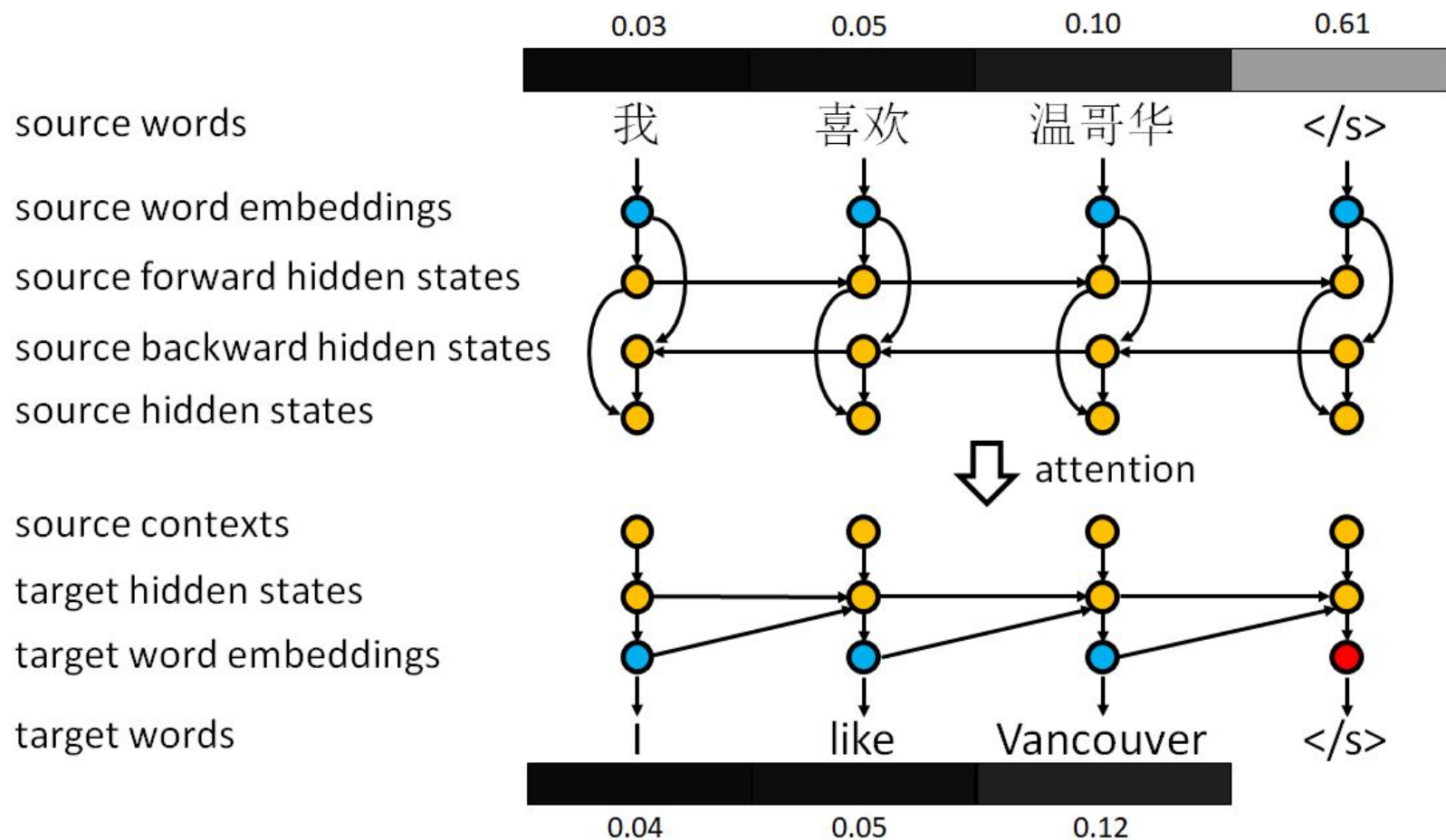Can calculate the relevance between two arbitrary neurons

Measures/visualizes how much each pixel is related to the final classification

# Goal

- To quantify and visualize the relevance between a neural network layer and contextual word vectors(source & target word embeddings)


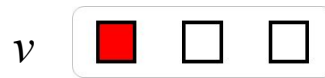Offers more insights in interpreting how target words are generated

# Relevance vector
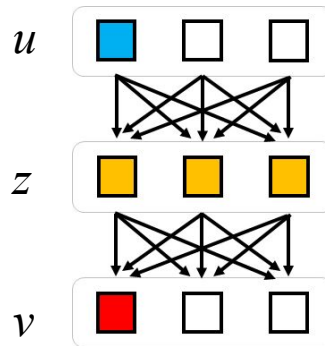
# Calculating Neuron-Level Relevance

Base case: (relevance of $v$ to itself)

$$r_{v \leftarrow v} = v$$

for any neuron $v$

$v$ 

Recursive case: (relevance of $u$ to $v$)

$$r_{u \leftarrow v} = \sum_{z \in \text{OUT}(u)} w_{u \rightarrow z} r_{z \leftarrow v}$$
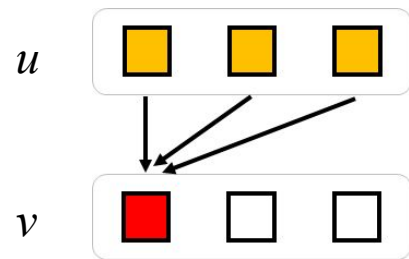
for any neurons $u, v$



$\text{OUT}(u)$ comprises all $u$'s ***directly connected descendant*** neurons in the network.

# Calculating Weight Ratios

$$w_{u \to v} = \frac{\mathbf{W}_{u,v}u}{\sum_{u' \in \text{IN}(v)} \mathbf{W}_{u',v}u'}$$
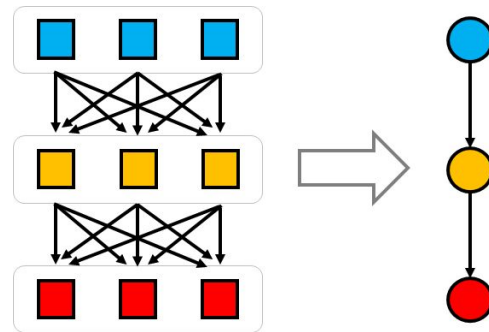
for any neurons $u$, $v$



$\mathbf{W}_{u,v}u$ is the weight of $u$ to $v$ in the existing neural network

$\text{IN}(u)$ comprises all $u$'s **directly connected ancestor** neurons in the network.

# Putting things together

Sum up $r_{u_n \leftarrow v_m}$ and get vector-level relevance $R_{\mathbf{u} \leftarrow \mathbf{v}}$
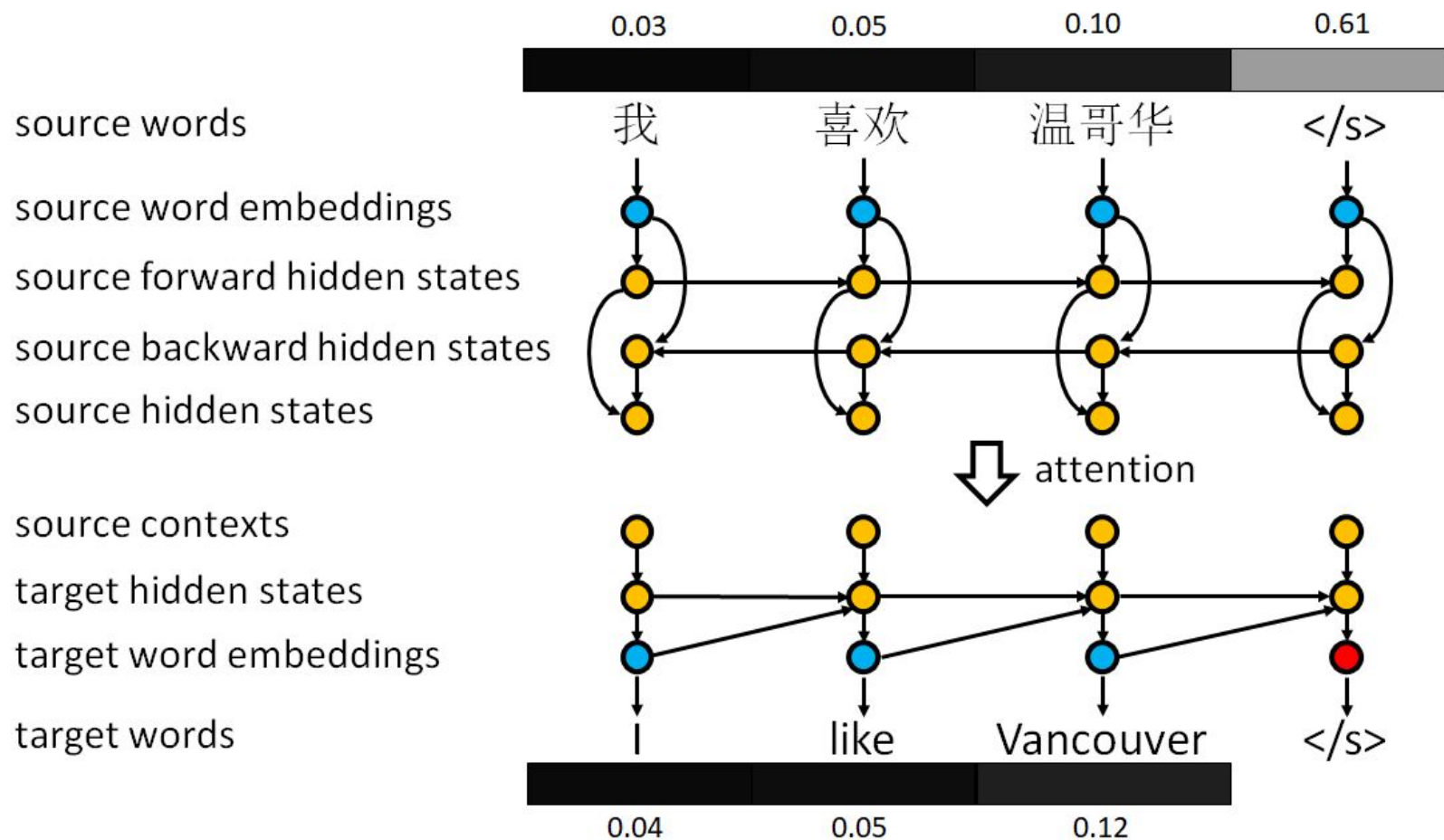
$$R_{\mathbf{u} \leftarrow \mathbf{v}} = \sum_{m=1}^{M} \sum_{n=1}^{N} r_{u_n \leftarrow v_m}$$



Generate and normalize relevance vector $R_{\mathbf{v}}$ as a sequence of $R_{\mathbf{u} \leftarrow \mathbf{v}}$ for all related contextual word vectors

$$R_{\mathbf{v}} = \{R_{\mathbf{u}_1 \leftarrow \mathbf{v}}, \ldots, R_{\mathbf{u}_{|\mathcal{C}(\mathbf{v})|} \leftarrow \mathbf{v}}\}$$

# Relevance vector



source words | 我 | 喜欢 | 温哥华 | </s>

0.03     0.05     0.10     0.61

source word embeddings

source forward hidden states

source backward hidden states

source hidden states

attention

source contexts

target hidden states

target word embeddings

target words    I    like    Vancouver    </s>
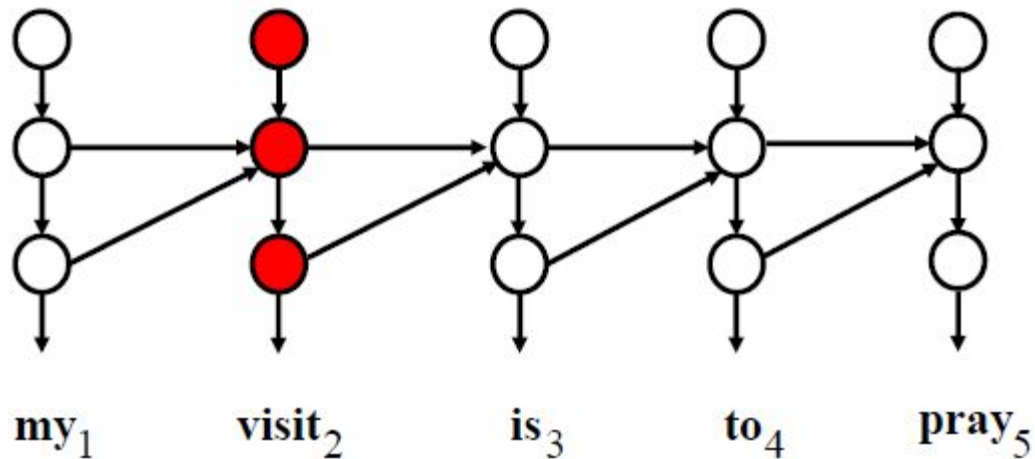
0.04     0.05     0.12

# Application

Help debug attention-based NMT systems

- Word omission
- Word repetition
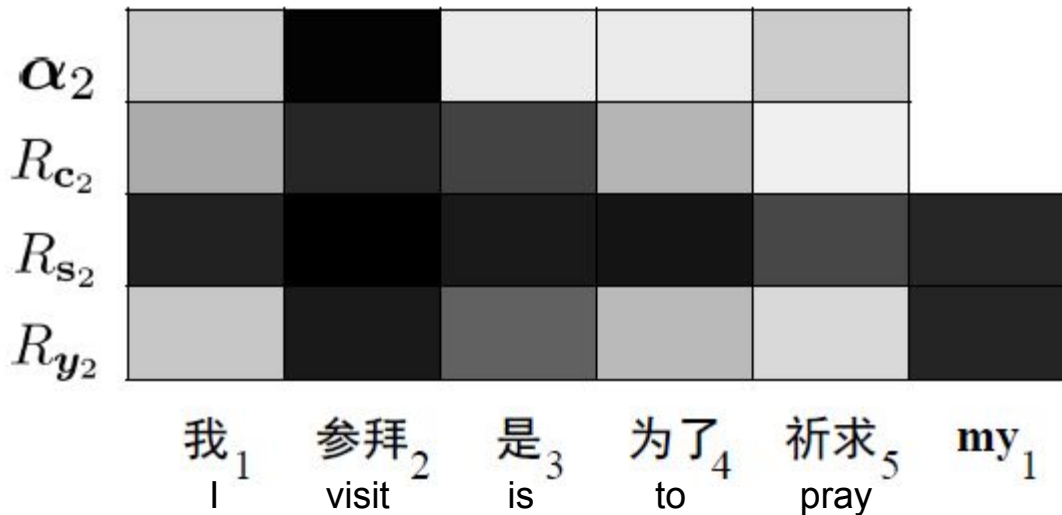- Unrelated words
- Negation reversion

"Relevance matrix"

attention weights $\alpha_2$
source context vector $R_{\mathbf{c}_2}$
target hidden state $R_{\mathbf{s}_2}$
target word embedding $R_{\mathbf{y}_2}$

$\mathbf{my}_1$    $\mathbf{visit}_2$    $\mathbf{is}_3$    $\mathbf{to}_4$    $\mathbf{pray}_5$

我$_1$ 参拜$_2$ 是$_3$ 为了$_4$ 祈求$_5$ $\mathbf{my}_1$
I    visit    is$_3$    to$_4$    pray

"Relevance matrix"



attention weights
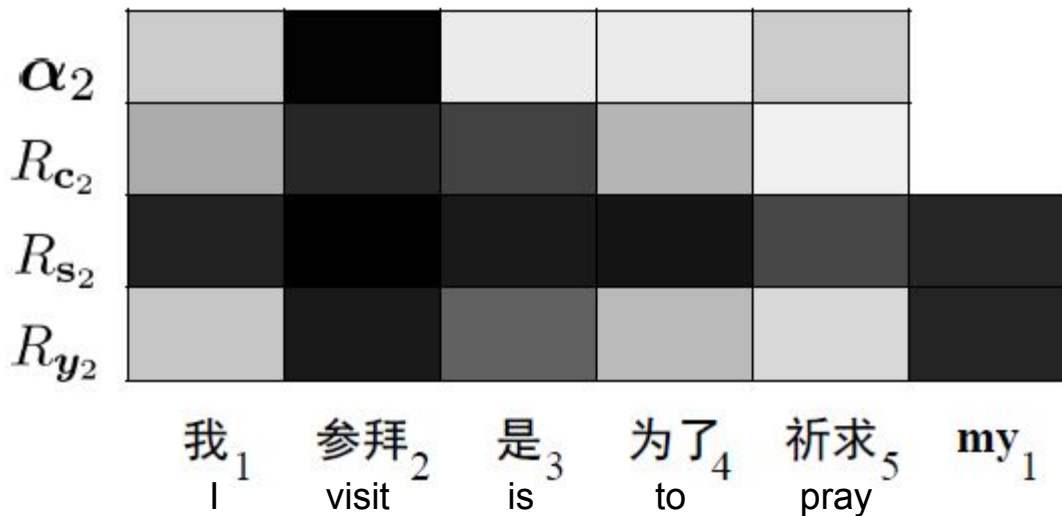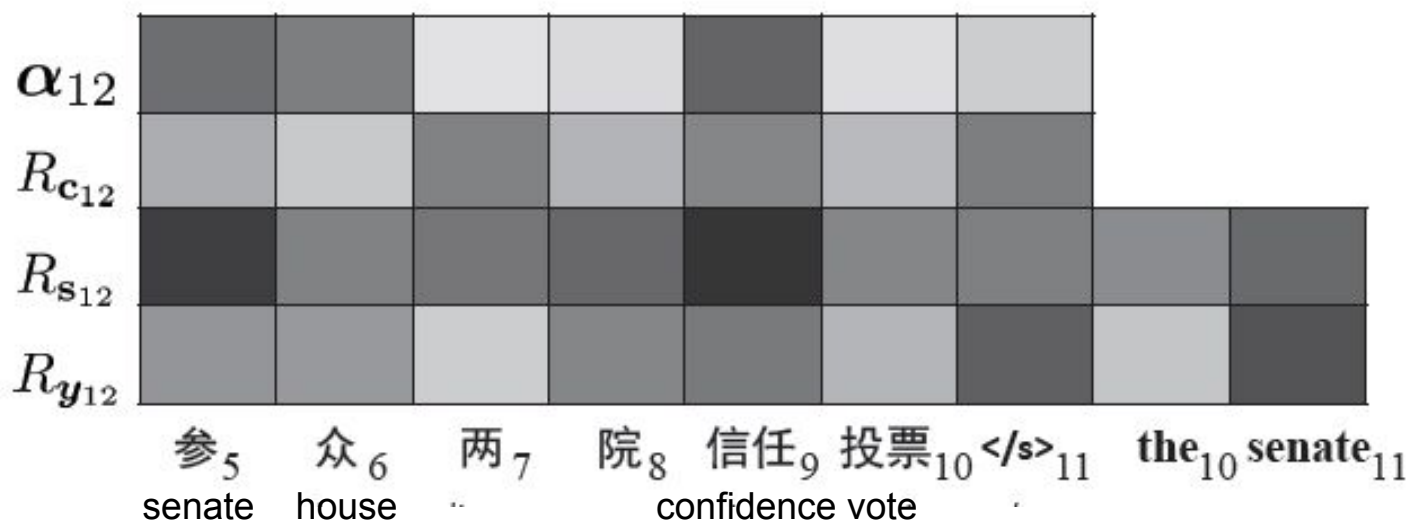
source context vector

target hidden state

target word embedding

# Word omission

| Input | 巴基斯坦总统穆沙拉夫赢得参**众**两院信任投票 |
|-------|------|
| Reference | Pakistani president Musharraf wins votes of confidence in senate <span style="color:red">and house</span> |
| Output | Pakistani president win over democratic vote of confidence in senate <span style="color:red">(missing words)</span> |

vote$_6$    of$_7$    confidence$_8$    in$_9$    the$_{10}$    senate$_{11}$    &lt;/s&gt;$_{12}$

$\boldsymbol{\alpha}_{12}$

$R_{\mathbf{c}_{12}}$

$R_{\mathbf{s}_{12}}$

$R_{\boldsymbol{y}_{12}}$

参$_5$   众$_6$   两$_7$   院$_8$   信任$_9$   投票$_{10}$   &lt;/s&gt;$_{11}$   the$_{10}$   senate$_{11}$

senate    house           confidence vote

vote$_6$    of$_7$    confidence$_8$    in$_9$    the$_{10}$    senate$_{11}$    </s>$_{12}$

$\boldsymbol{\alpha}_{12}$
$R_{\mathbf{c}_{12}}$
$R_{\mathbf{s}_{12}}$
$R_{\boldsymbol{y}_{12}}$

参$_5$    众$_6$    两$_7$    院$_8$    信任$_9$    投票$_{10}$    </s>$_{11}$    the$_{10}$    senate$_{11}$

senate    house        confidence vote

# Word repetition

| Input | 美国人历史上有讲诚信的传统，有犯错认错的传统 |
|---|---|
| Reference | In history, Americans have the tradition of honesty and would not hesitate to admit their mistakes |
| Output | In the history of the history of the history of the Americans, there is a tradition of faith in the history of mistakes |

the$_2$    history$_3$    of$_4$    the$_5$    history$_6$

$\boldsymbol{\alpha}_6$

$R_{\mathbf{c}_6}$

$R_{\mathbf{s}_6}$

$R_{\boldsymbol{y}_6}$

美国人$_1$   历史$_2$   上$_3$   有$_4$   of$_4$   the$_5$

Americans   history        have

$\boldsymbol{\alpha}_6$

$R_{\mathbf{c}_6}$

$R_{\mathbf{s}_6}$

$R_{\boldsymbol{y}_6}$

the$_2$  history$_3$  of$_4$  the$_5$  history$_6$

美国人$_1$  历史$_2$  上$_3$  有$_4$  of$_4$  the$_5$

Americans  history  have

# Unrelated words

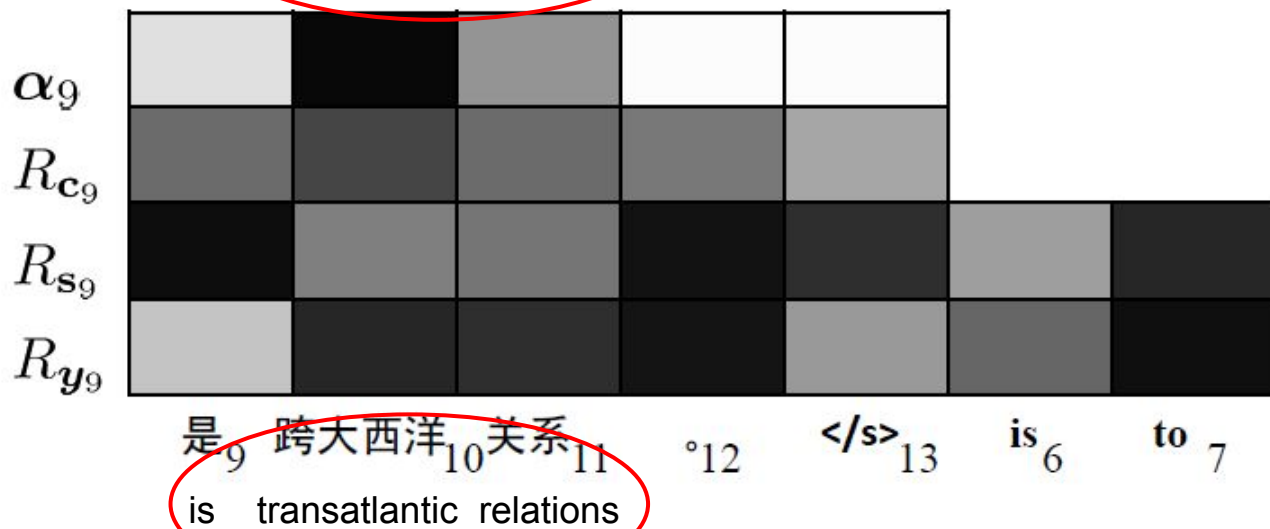| | |
|---|---|
| Input | 此次会议的一个重要议题是<span style="color:red">跨大西洋关系</span> |
| Reference | One of the top agendas of the meeting is to <span style="color:red">discuss the transatlantic relations</span> |
| Output | A key topic of the meeting is to <span style="color:red">forge ahead</span> |

$\mathbf{is}_7$  $\mathbf{to}_8$  $\mathbf{forge}_9$  $\mathbf{ahead}_{10}$  $\cdot_{11}$  $\mathbf{</s>}_{12}$

$\boldsymbol{\alpha}_9$

$R_{\mathbf{c}_9}$

$R_{\mathbf{s}_9}$

$R_{\boldsymbol{y}_9}$

是$_9$ 跨大西洋$_{10}$关系$_{11}$  。$_{12}$  $\mathbf{</s>}_{13}$  $\mathbf{is}_6$  $\mathbf{to}_7$

is    transatlantic  relations

$\boldsymbol{\alpha}_9$

$R_{\mathbf{c}_9}$

$R_{\mathbf{s}_9}$

$R_{\boldsymbol{y}_9}$

is$_7$    to$_8$    forge$_9$    ahead$_{10}$    ·$_{11}$    </s>$_{12}$

是$_9$  跨大西洋$_{10}$ 关系$_{11}$  ·$_{12}$  </s>$_{13}$  is$_6$  to$_7$

is    transatlantic  relations

# Negation reversion

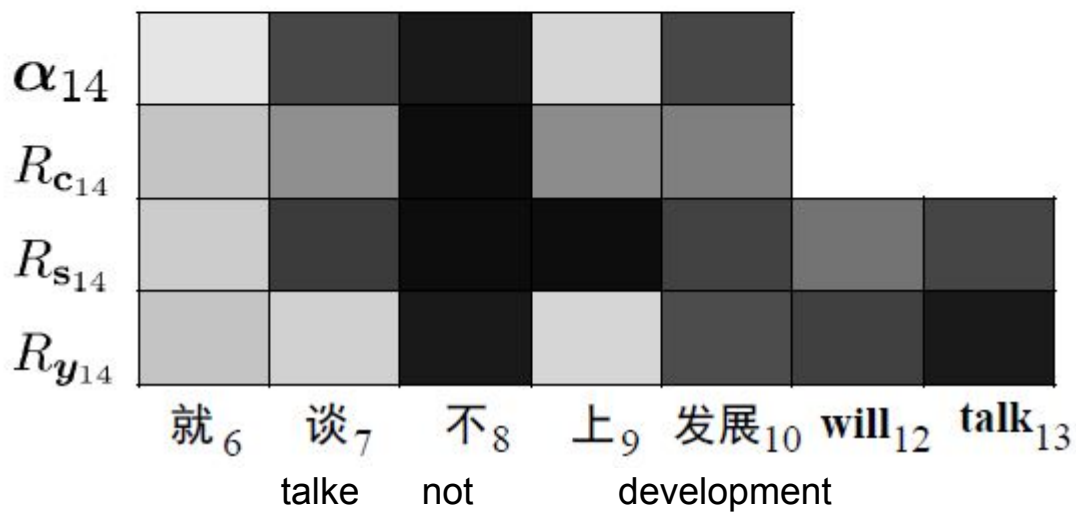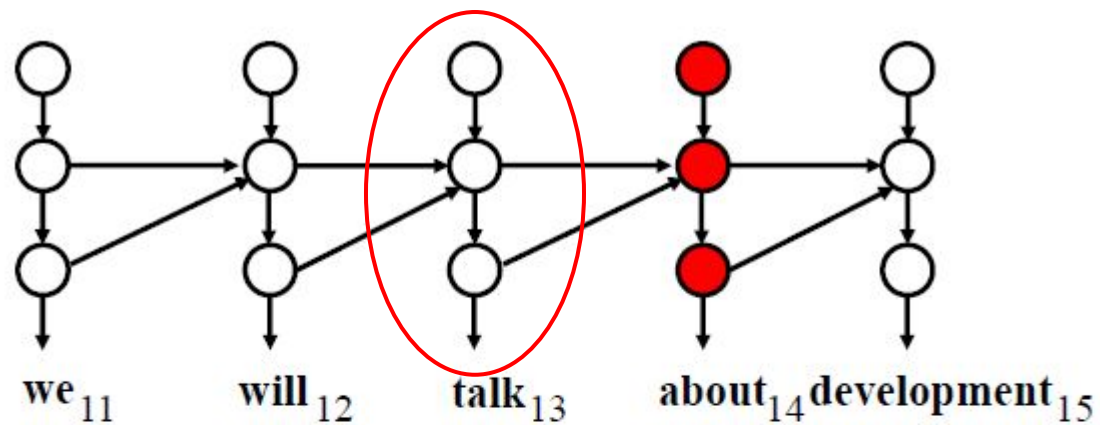| | |
|---|---|
| Input | 不解决生存问题，就谈<span style="color:red">不</span>上发展，更谈不上可持续发展 |
| Reference | Without solving the issue of subsistence, there will be no development to speak of , let alone sustainable development |
| Output | If we do not solve the problem of living , we will talk about development and still less can we talk about sustainable development |

| | 就$_6$ | 谈$_7$ | 不$_8$ | 上$_9$ | 发展$_{10}$ | will$_{12}$ | talk$_{13}$ |
|---|---|---|---|---|---|---|---|
| $\boldsymbol{\alpha}_{14}$ | | | | | | | |
| $R_{\mathbf{c}_{14}}$ | | | | | | | |
| $R_{\mathbf{s}_{14}}$ | | | | | | | |
| $R_{\boldsymbol{y}_{14}}$ | | | | | | | |

talke     not     development

$\boldsymbol{\alpha}_{14}$

$R_{\mathbf{c}_{14}}$

$R_{\mathbf{s}_{14}}$

$R_{\boldsymbol{y}_{14}}$

we$_{11}$    will$_{12}$    talk$_{13}$    about$_{14}$ development$_{15}$

就$_6$    谈$_7$    不$_8$    上$_9$    发展$_{10}$    will$_{12}$    talk$_{13}$

talke    not    development

# Thank you