# Visualizing and Understanding Neural Machine Translation

Yanzhuo Ding, Yang Liu, Huanbo Luan, Maosong Sun

*Tsinghua University*

# Neural Machine Translation

- Idea: using neural networks to translate languages (Bahdanau et al., 2015)

source words         我      喜欢      温哥华      </s>

source word embeddings

source forward hidden states
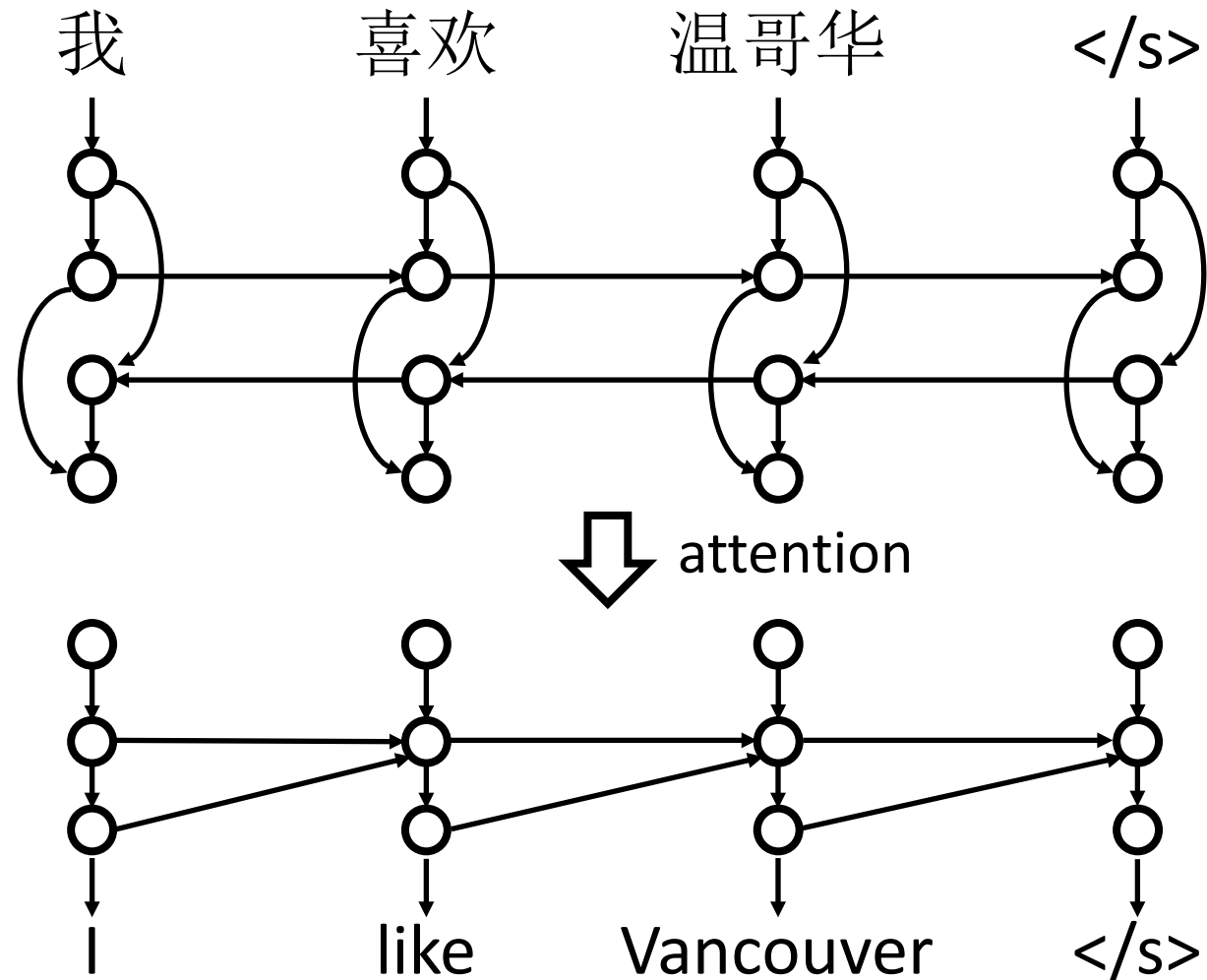
source backward hidden states

source hidden states

attention

source contexts

target hidden states

target word embeddings

target words      I      like      Vancouver      </s>

# Challenge

- It is hard to visualize and understand the internal workings

source words

我　　　喜欢　　　温哥华　　　</s>

source word embeddings

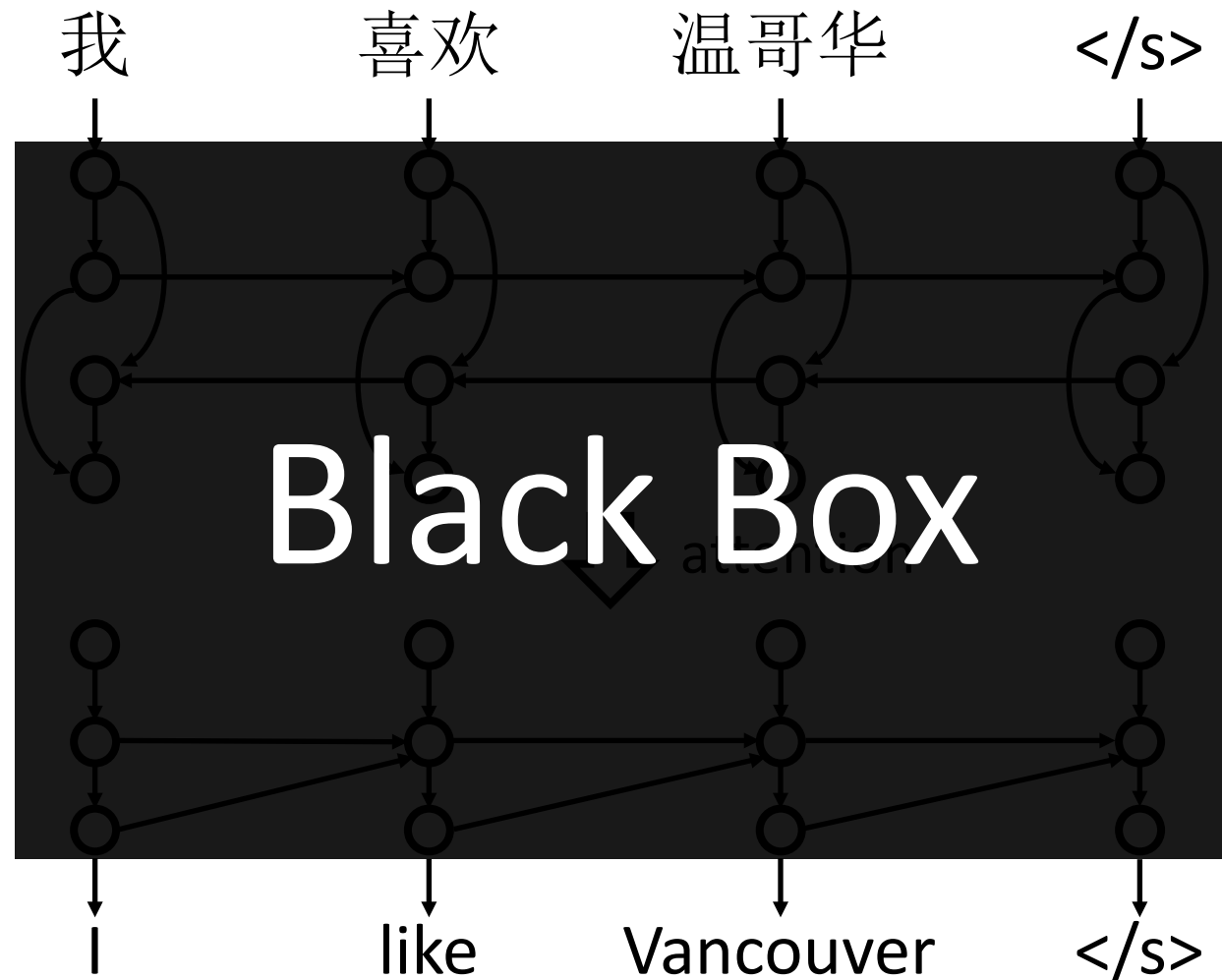source forward hidden states

source backward hidden states

source hidden states

**Black Box**

source contexts

target hidden states

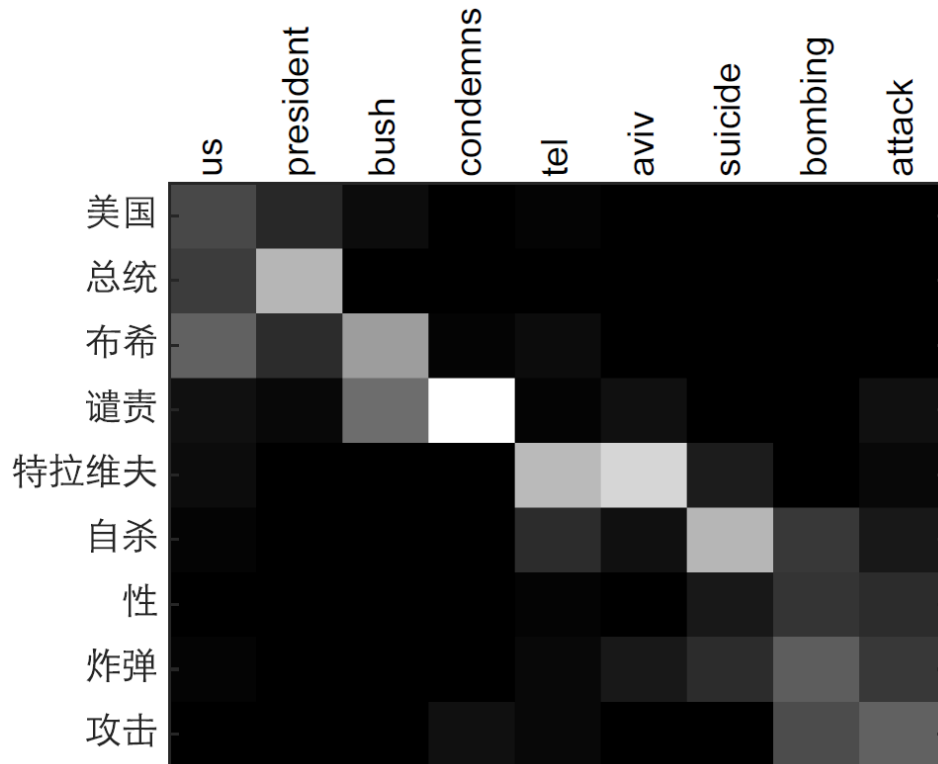target word embeddings

target words

I　　　like　　Vancouver　　</s>
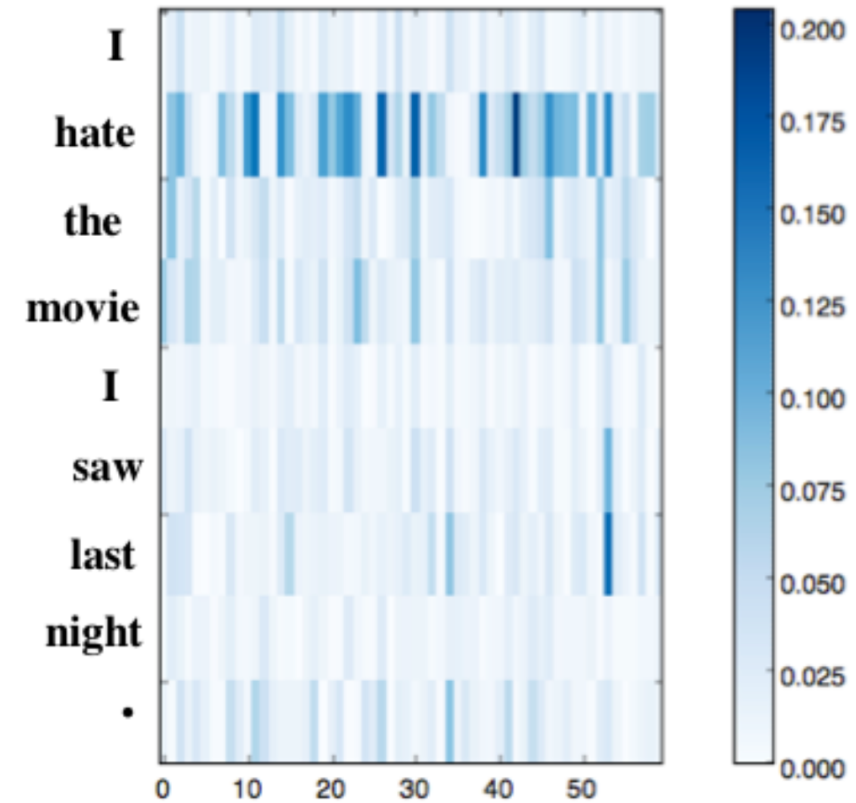
# Related Work

## attention mechanism
(Bahdanau et al., 2015)



*restricted to the connection between input and output*
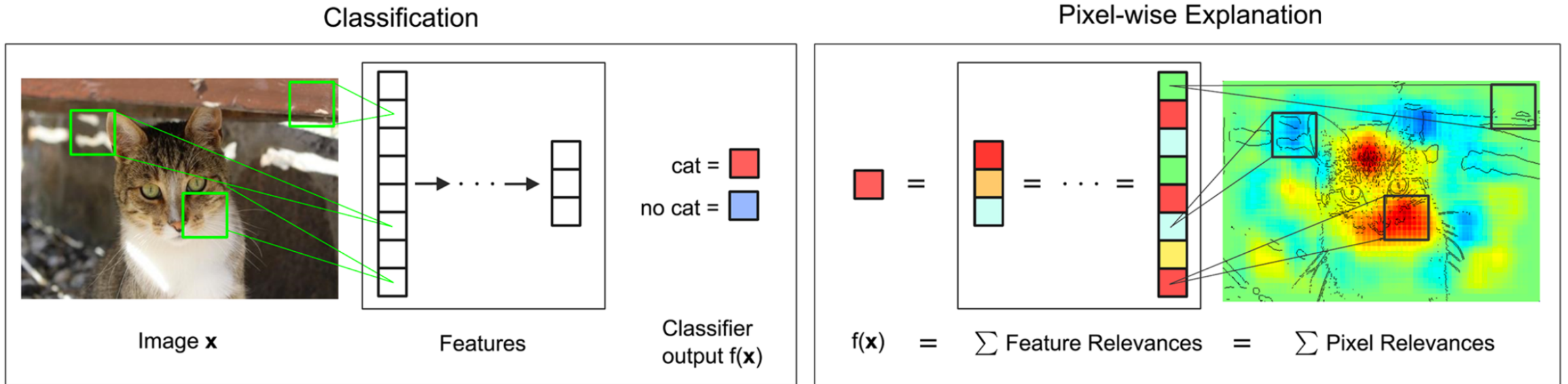
## first-derivative saliency
(Li et al., 2016)



*require neural activations to be differentiable*

# Related Work

## layer-wise relevance propagation (LRP)
(Bach et al., 2015)



*calculating the relevance between two arbitrary neurons*
*without requiring differentiability*

# Our Work

source words                我          喜欢        温哥华        </s>

# Our Work

source words                我         喜欢        温哥华        &lt;/s&gt;

source word embeddings     O         O        O         O

# Our Work

source words       我      喜欢      温哥华      </s>

source word embeddings

source forward hidden states

# Our Work

1.00

source words 　　　　　　　　我　　　　　喜欢　　　　温哥华　　　　</s>

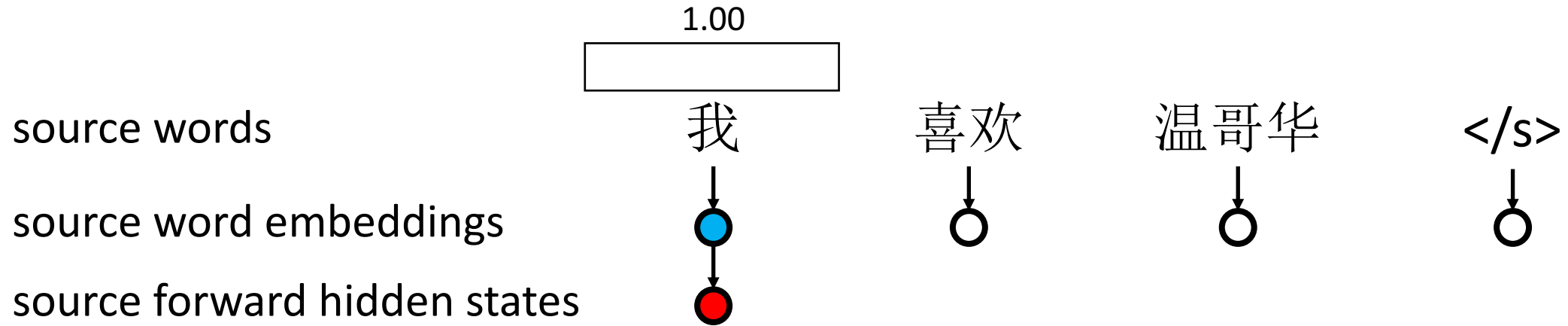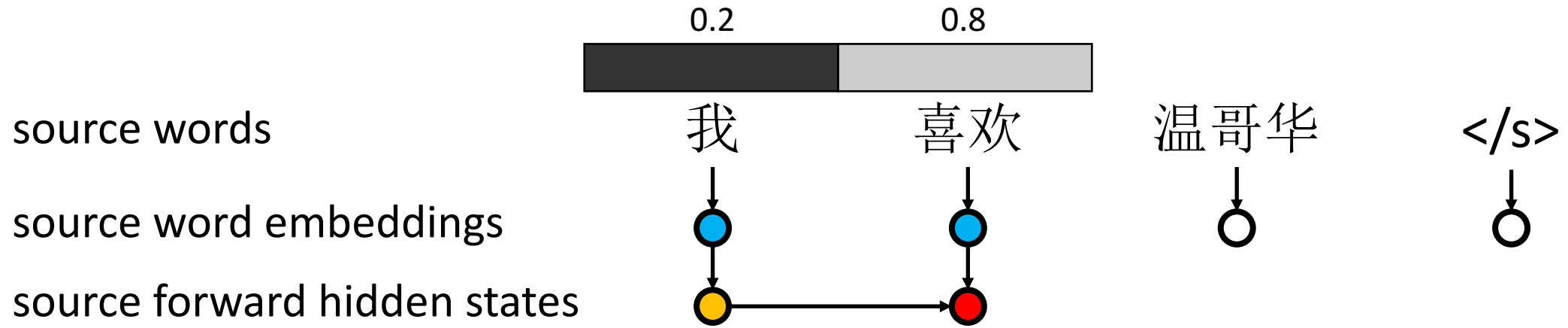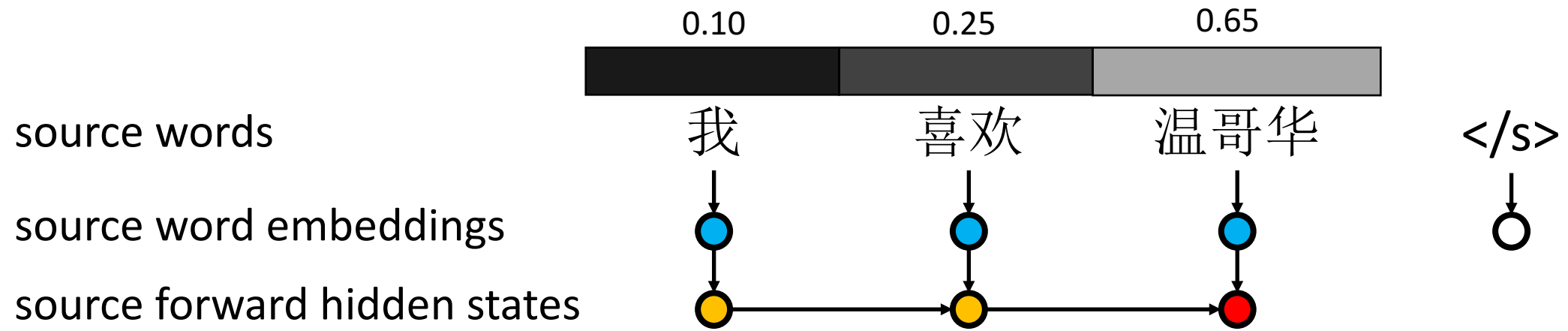source word embeddings

source forward hidden states

🔴　　　targeted vector of neurons

🔵　　　relevant vector of neurons

🟠　　　intermediate vector of neurons

⚪　　　irrelevant vector of neurons

1.0　　　relevance
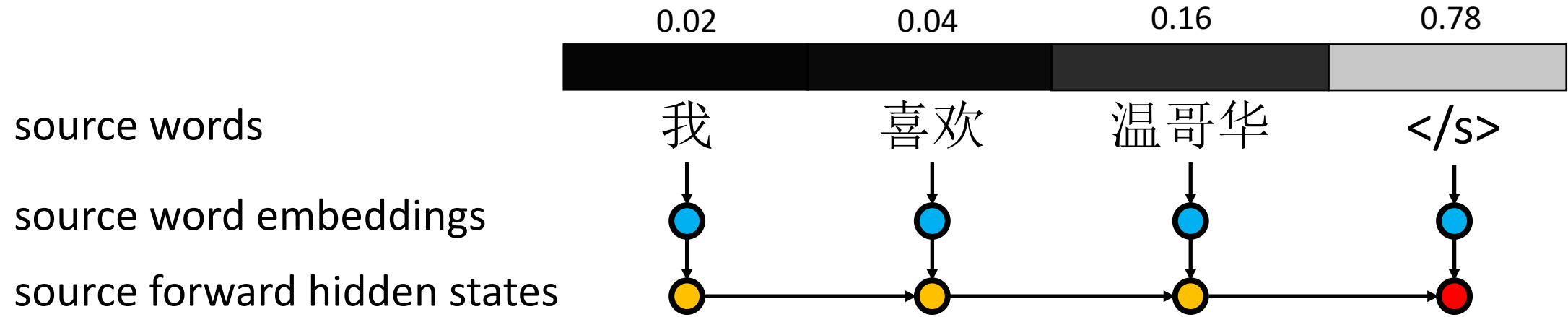
# Our Work

0.2            0.8

source words                    我            喜欢            温哥华            </s>

source word embeddings

source forward hidden states

# Our Work

# Our Work

| 0.02 | 0.04 | 0.16 | 0.78 |

**source words**  我  喜欢  温哥华  </s>

**source word embeddings**

**source forward hidden states**

# Our Work

1.00

source words       我       喜欢       温哥华       \</s\>
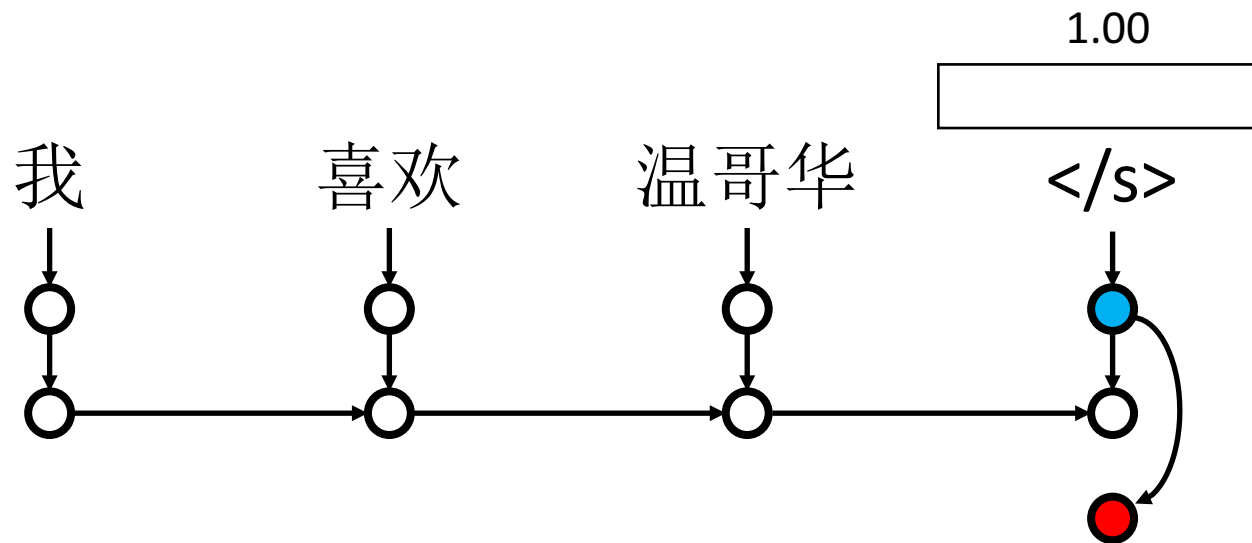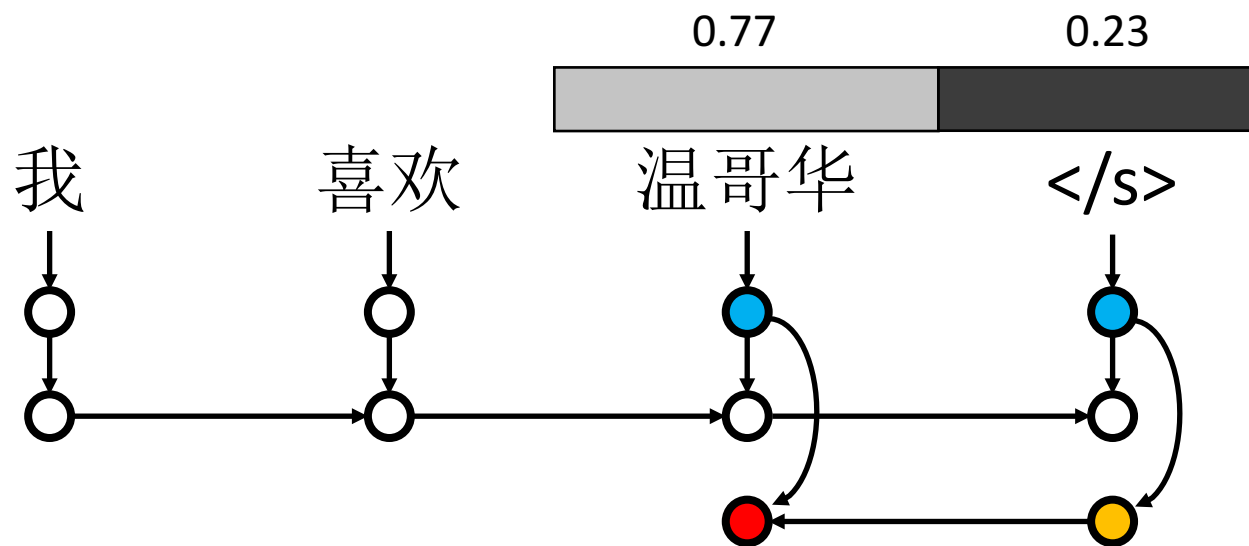
source word embeddings

source forward hidden states

source backward hidden states

# Our Work

0.77                0.23

source words           我        喜欢       温哥华       </s>

source word embeddings

source forward hidden states

source backward hidden states

# Our Work

# Our Work

# Our Work

# Our Work

# Our Work

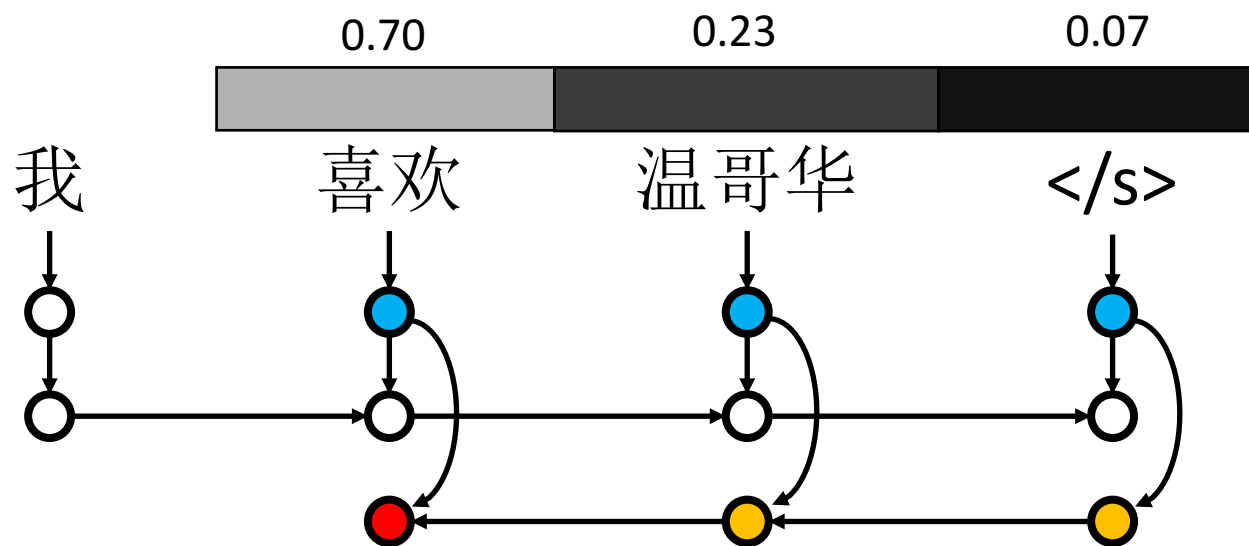| 0.06 | 0.12 | 0.80 | 0.02 |

source words　　　　　　我　　　　　喜欢　　　　温哥华　　　　</s>
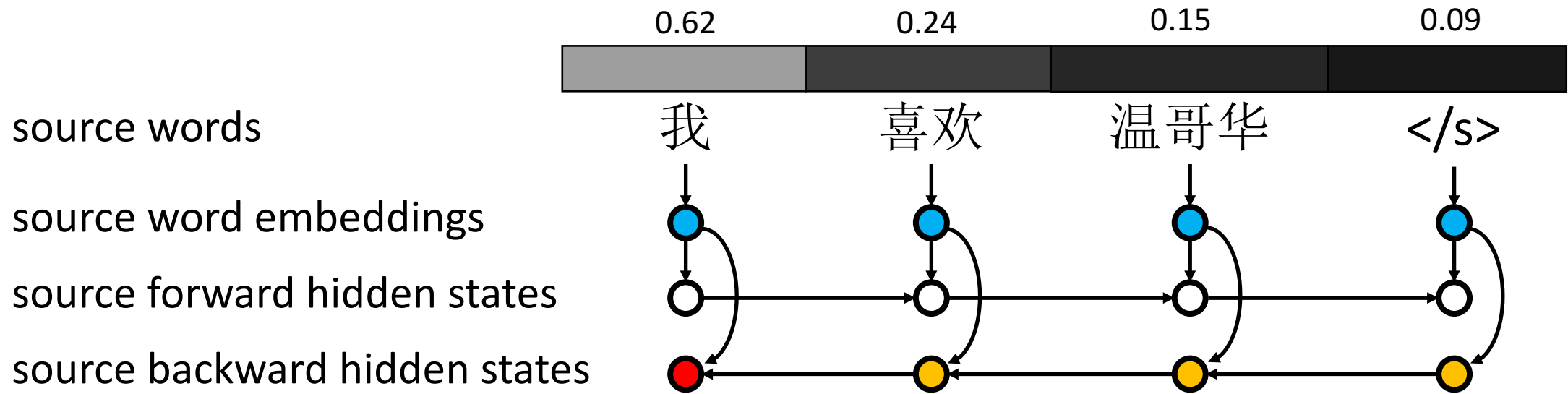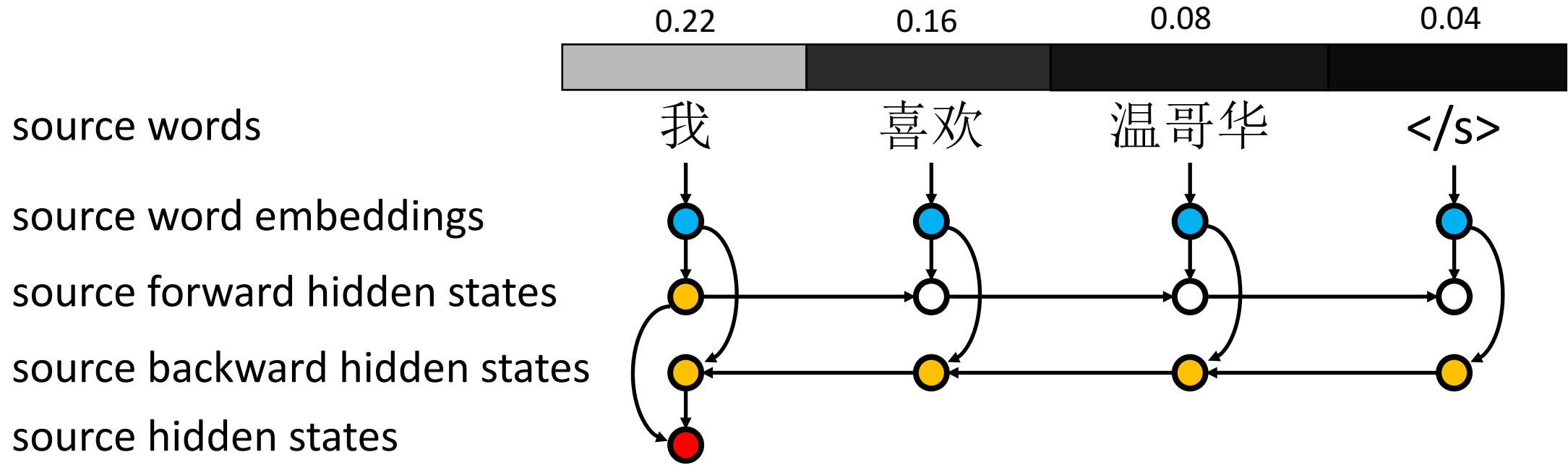
source word embeddings

source forward hidden states

source backward hidden states

source hidden states

# Our Work

source words 我 喜欢 温哥华 </s>

0.01  0.02  0.07  0.90

source word embeddings

source forward hidden states

source backward hidden states

source hidden states

# Our Work

0.58    0.32    0.07    0.03

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

# Our Work

# Our Work

# Our Work

0.18 0.29 0.25 0.03

source words    我    喜欢    温哥华    </s>

source word embeddings

source forward hidden states

source backward hidden states

source hidden states
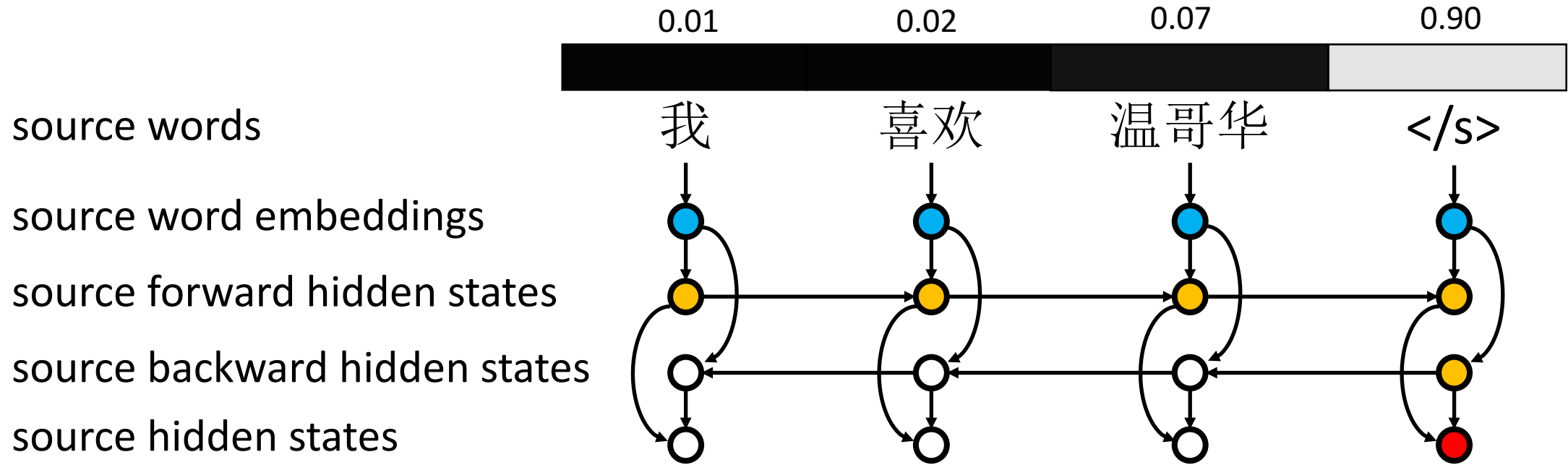
attention

source contexts

target hidden states

target word embeddings

target words

0.23

# Our Work

| 0.14 | 0.33 | 0.24 | 0.04 |

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

0.25

# Our Work



| 0.11 | 0.45 | 0.11 | 0.02 |

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

0.31

# Our Work

0.09      0.15      0.51      0.04

source words     我     喜欢     温哥华     &lt;/s&gt;

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

0.21

# Our Work

0.05      0.16      0.35      0.02

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

0.11      0.31

# Our Work



source words 我 喜欢 温哥华 </s>

0.04   0.10   0.42   0.02

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

0.05   0.37

# Our Work

0.05          0.09          0.11          0.52

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
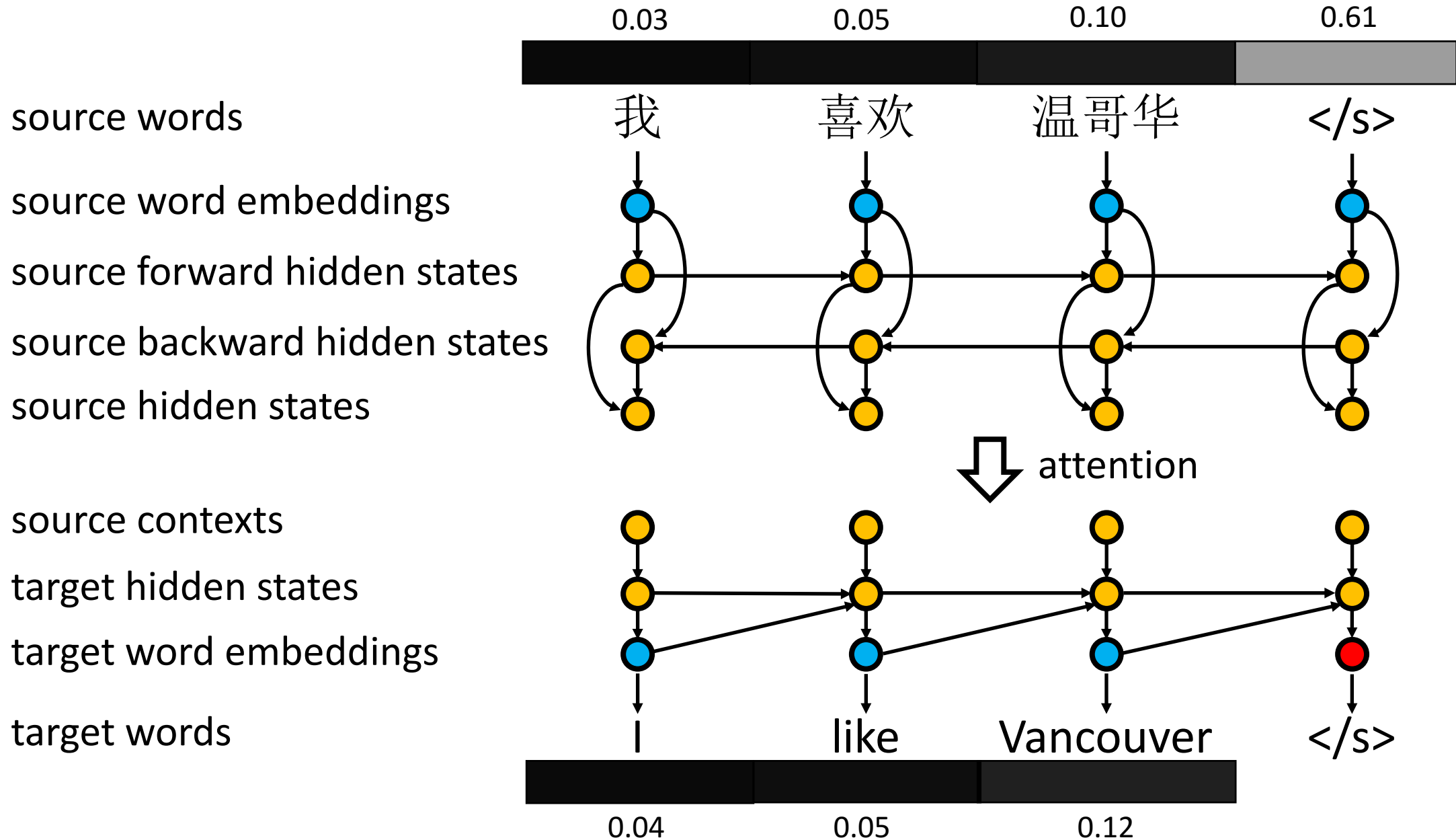
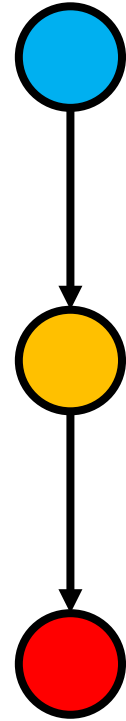target words          I          like          Vancouver
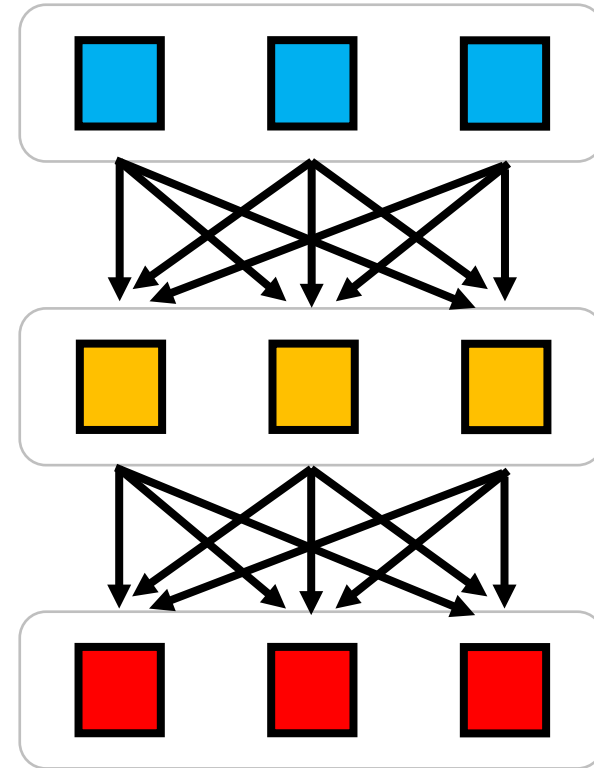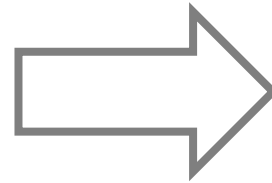
0.06          0.17

# Our Work

# Our Work

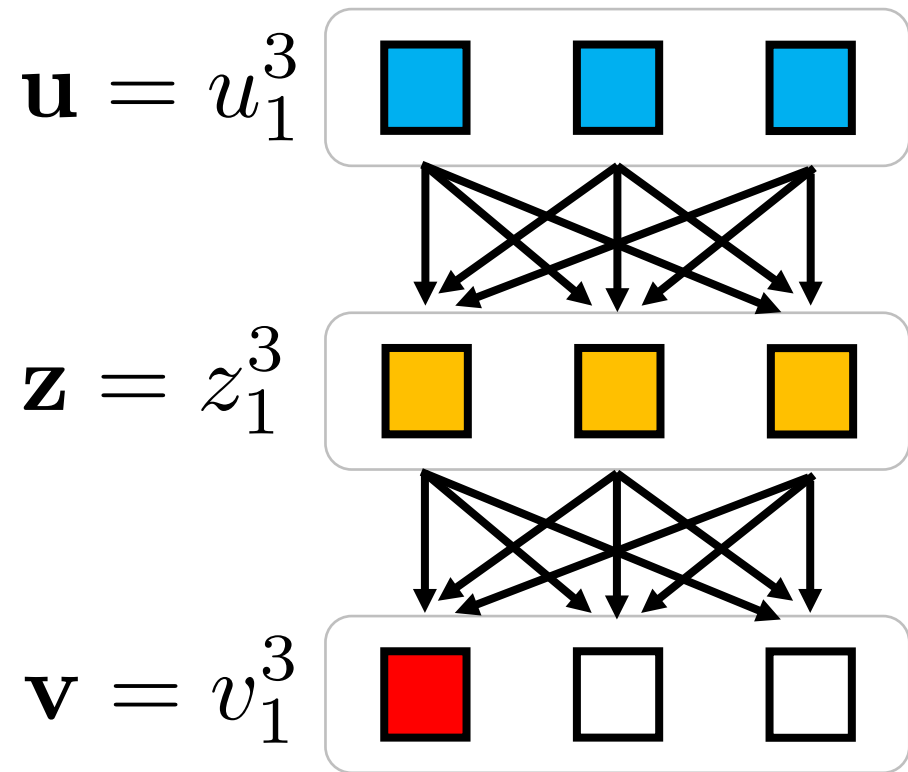# Vector- and Neuron-Level Relevance



vector-level relevance     neural-level relevance

# Neuron-Level Relevance

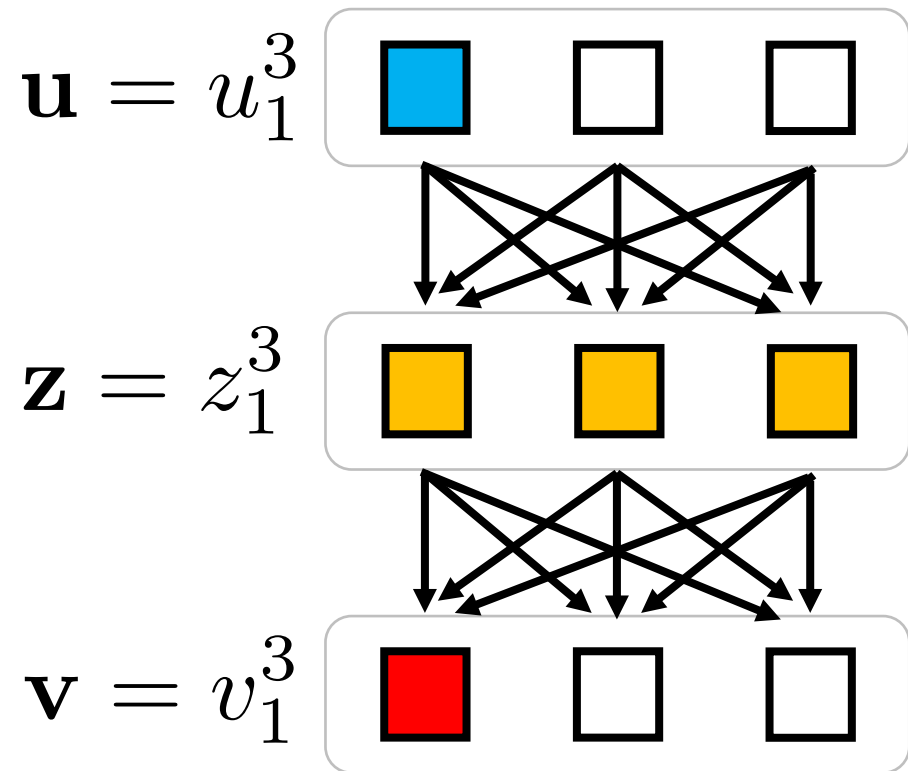- Idea: decompose the activation of the targeted neuron among relevant neurons

$$\mathbf{u} = u_1^3$$



$$\mathbf{z} = z_1^3$$

$$\mathbf{v} = v_1^3$$

$$v_m = \sum_{\mathbf{u} \in \mathcal{C}(v_m)} \sum_{n=1}^{N} r_{u_n \leftarrow v_m}$$

For example

$$v_1 = \sum_{n=1}^{3} r_{u_n \leftarrow v_1}$$

# Calculating Neuron-Level Relevance

- Recursive calculation in a backward propagation



$$r_{u \leftarrow v} = \sum_{z \in \mathrm{OUT}(u)} w_{u \rightarrow z} r_{z \leftarrow v}$$
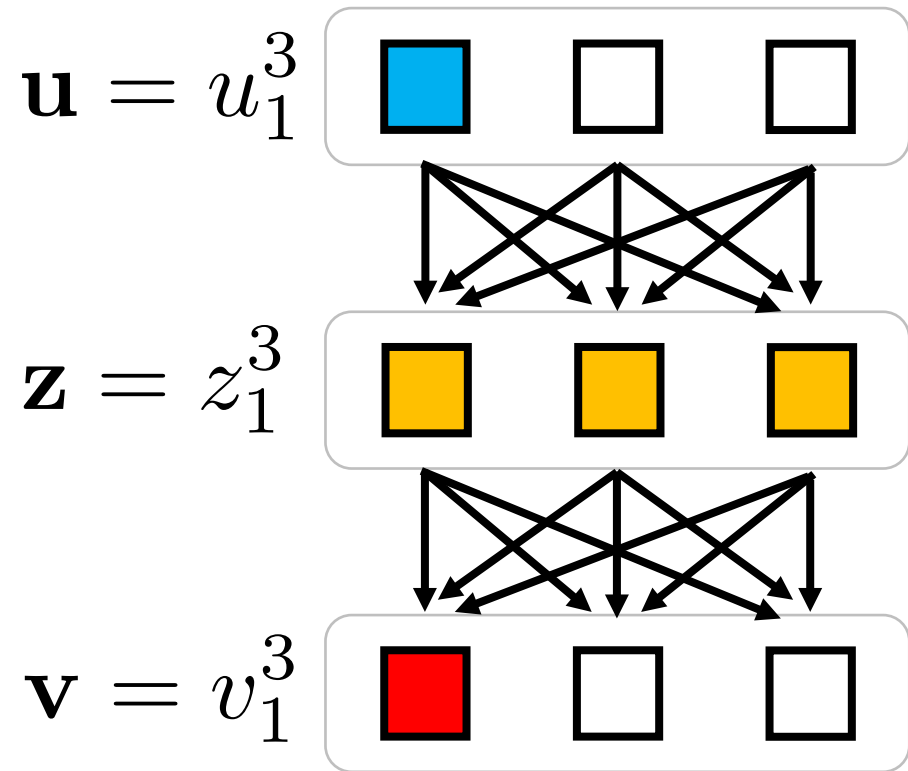
For example

$$r_{u_1 \leftarrow v_1} = \sum_{k=1}^{3} w_{u_1 \rightarrow z_k} r_{z_k \leftarrow v_1}$$

$$r_{z_k \leftarrow w_1} = w_{z_k \rightarrow v_1} v_1$$

# Calculating Weight Ratios
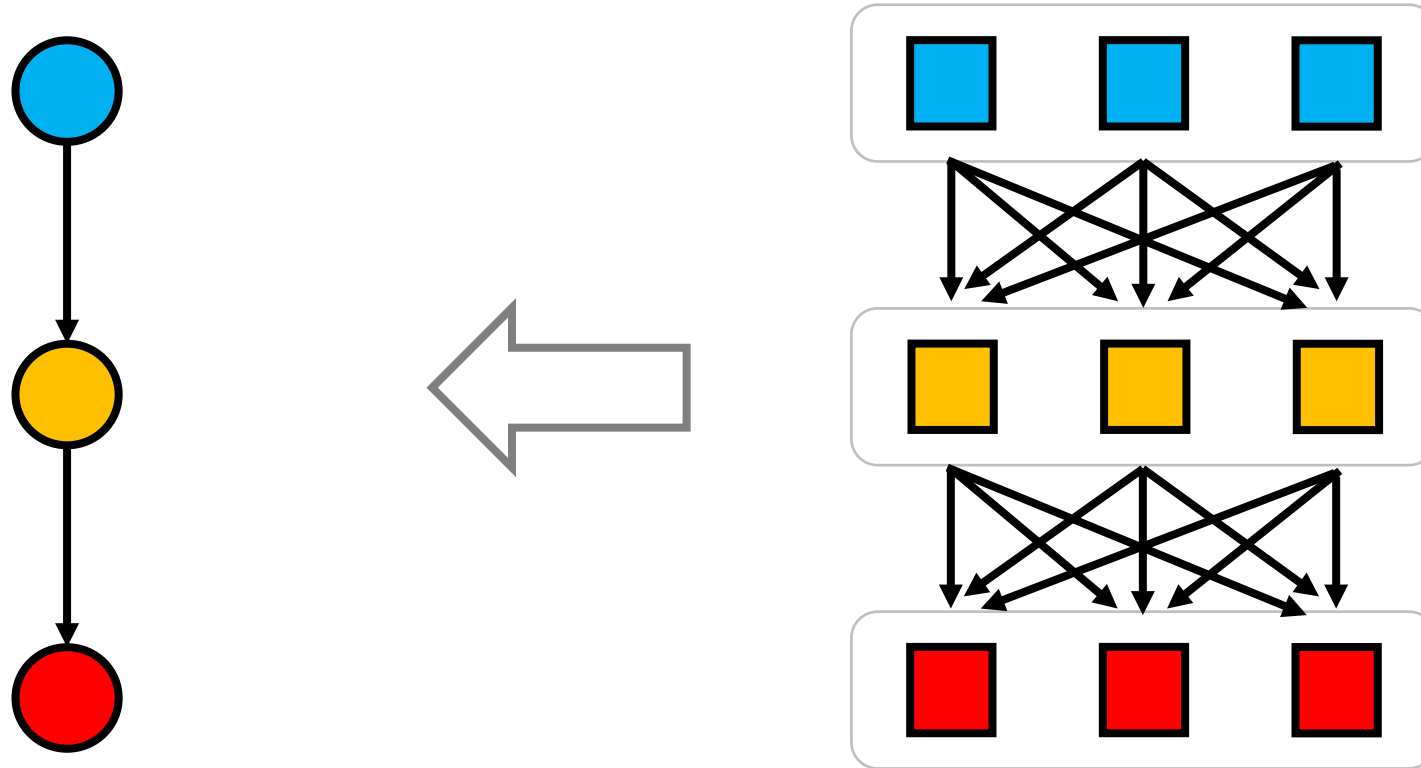
- Recursive calculation in a forward propagation



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

$\mathbf{u} = u_1^3$

$\mathbf{z} = z_1^3$

For example

$\mathbf{v} = v_1^3$

$$w_{u_1 \to z_1} = \frac{\mathbf{W}_{1,1}^{(1)} u_1}{\mathbf{W}_{1,1}^{(1)} u_1 + \mathbf{W}_{2,1}^{(1)} u_2 + \mathbf{W}_{3,1}^{(1)} u_3}$$

# Calculating Vector-Level Relevance



$$R_{\mathbf{u} \leftarrow \mathbf{v}} = \sum_{m=1}^{M} \sum_{n=1}^{N} r_{u_n \leftarrow v_m}$$
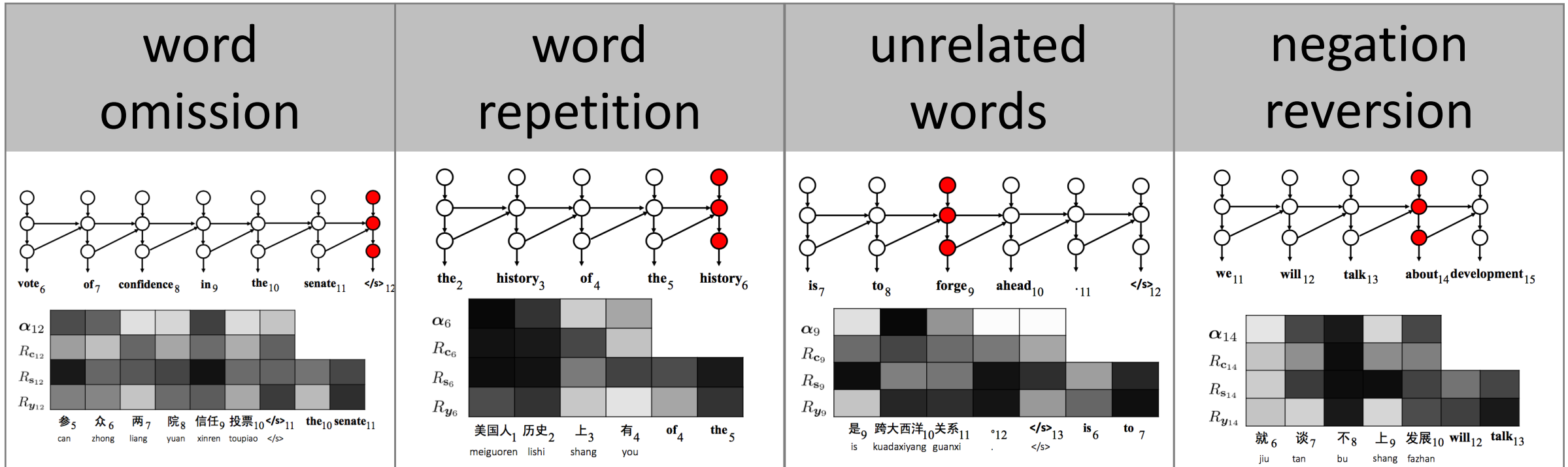
# Algorithm

- Specify targeted vector of neurons

- Calculate weight ratios in a forward propagation

- Calculate relevance in a backward propagation

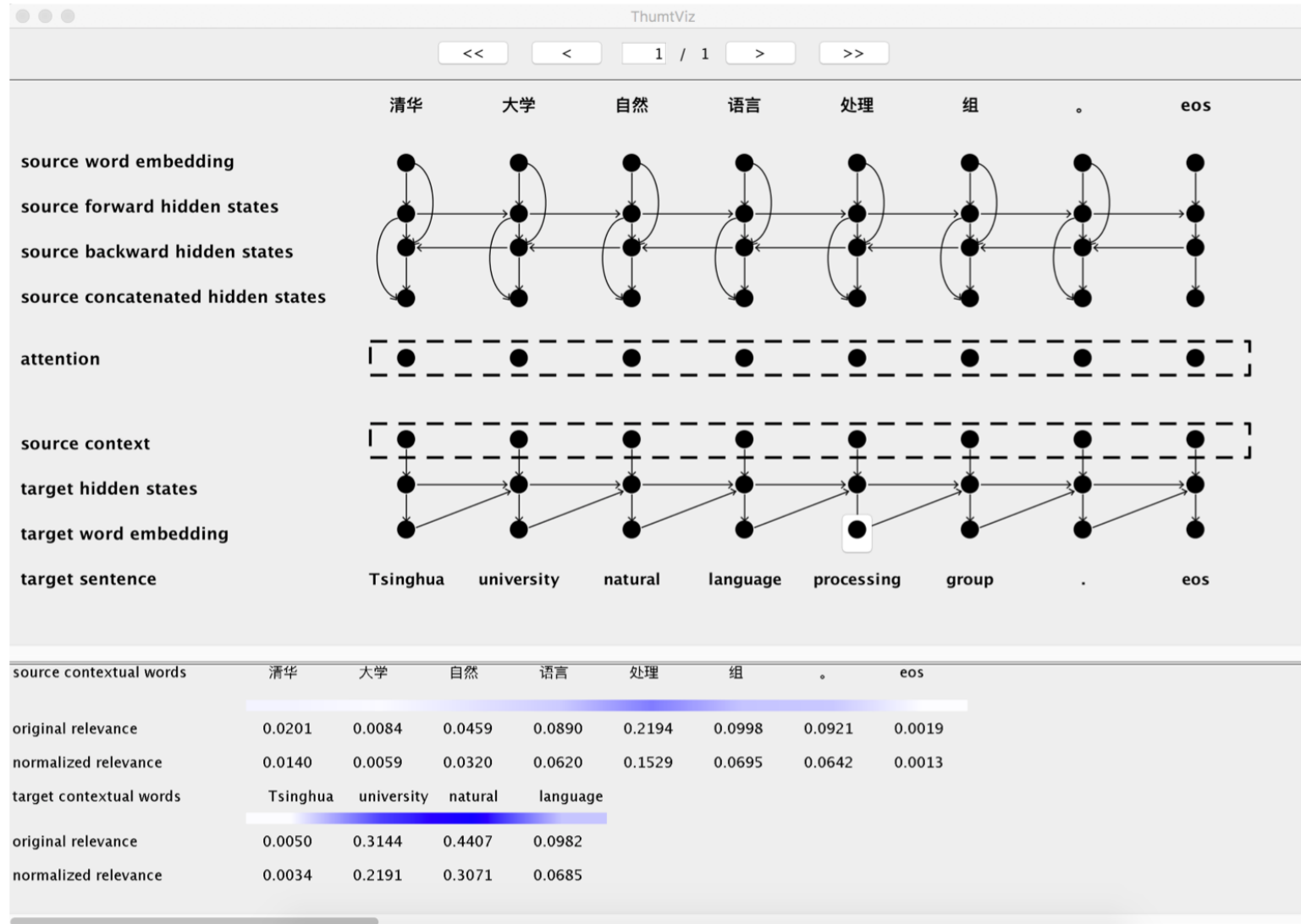*layer-wise propagation for neural machine translation*

# Application

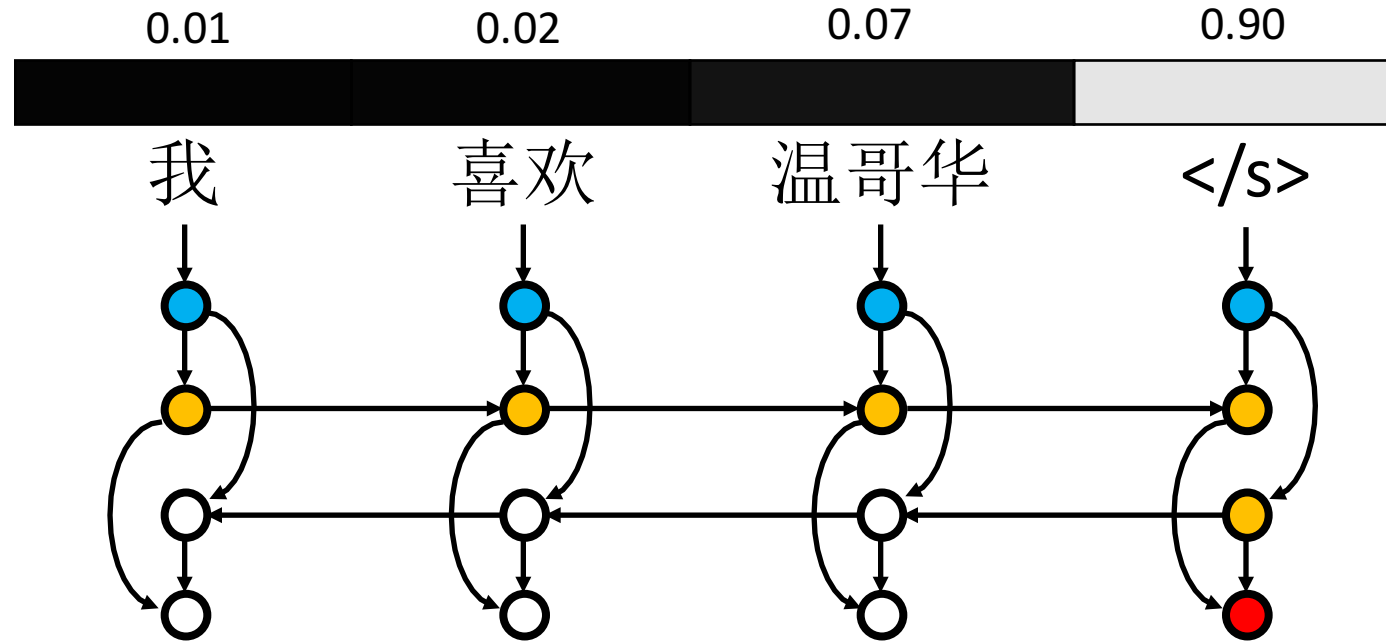- Help to debug attention-based NMT systems



| word omission | word repetition | unrelated words | negation reversion |
|---|---|---|---|

analyzing major translation error types by visualizing relevance step by step

# Open-Source Toolkit



http://thumt.thunlp.org

# Conclusion



- It is challenging to interpret how neural networks work
- We leverage layer-wise relevance propagation to visualize NMT
- Our approach can be applied to networks in other NLP tasks

# Thanks

http://thumt.thunlp.org