



(Image Source: Garfield)

Automatic Text Simplification

Wei Xu

School of Interactive Computing
Georgia Institute of Technology

 @cocoweixu

 wei.xu@cc.gatech.edu

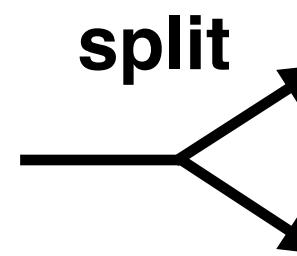
<https://cocoxu.github.io/>



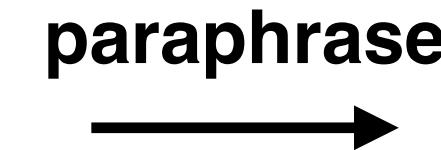
Text Simplification

Rewrite complex text into simpler language while retain its original meaning.

The ~~layers of calcified~~ plaque ~~entomb~~ the bacteria that also ~~live~~ in our mouths -- turning them into ~~small~~ fossils ~~even when we are alive~~.



And when we die, these ~~dense, calcified~~ micro-fossils remain intact, even as most of the rest of us decomposes.



The **buildup** of plaque **can trap** the bacteria that live in our mouths.

It turns them into **tiny** fossils.

Even after death, these micro-fossils **don't break down**.

Human Text Simplification

Professional editors rewrite news articles into 4 different readability levels for grade 3-12 students.

NEWSLEA

WAR & PEACE SCIENCE KIDS MONEY HEALTH

SCIENCE 1738 SHARE

MAX
1140L

960L
720L
420L

WRITE
 QUIZ

Archaeologist may have found remains of ancient Egyptian Queen Nefertiti

By Robert Gebelhoff, Washington Post.
08.17.15



The 3,330-year-old bust of Nefertiti sits in an exhibition in the Kulturforum in Berlin, Germany, March 1, 2005.
Photo: AP/Herbert Knosowski

Nefertiti — she's an ancient Egyptian queen and the source of a fantastic mystery regarding the iconic remnants of long-lost royalty.

For decades, archaeologists have speculated on the location of the queen's remains, the last royal mummy missing from the dynasty of the famous King Tutankhamun, better known as King Tut. But now, an archaeologist claims that he has found her

NEWSLEA

WAR & PEACE SCIENCE KIDS MONEY LAW HEALTH

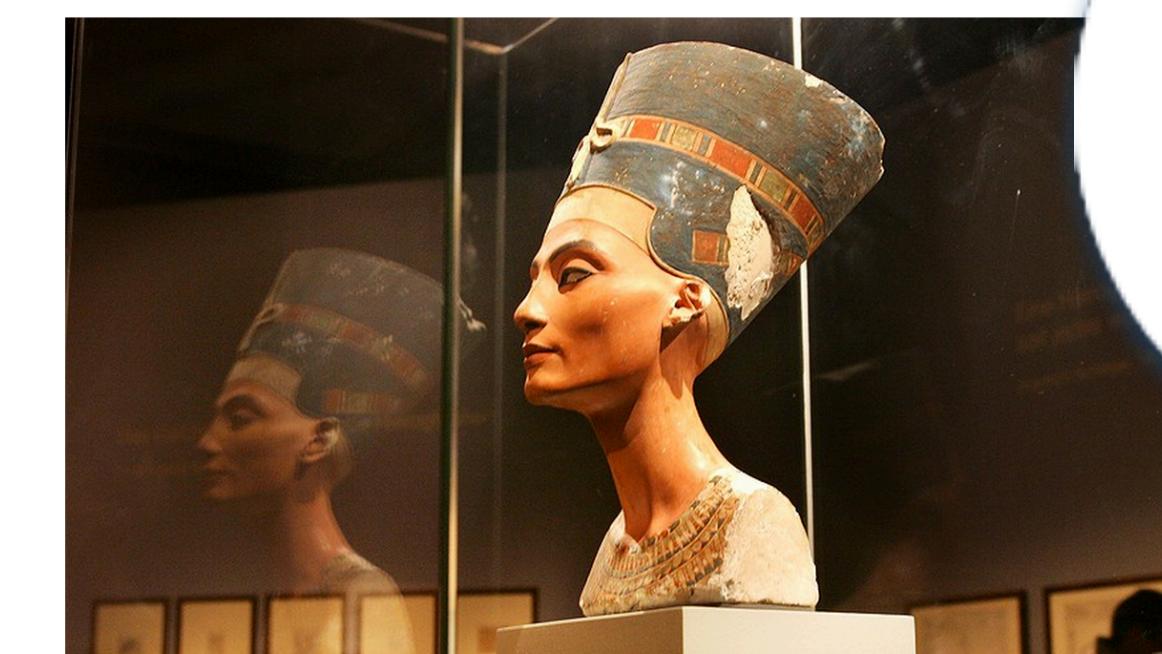
SCIENCE 1738 SHARE

1140L
960L
720L

WRITE
 QUIZ

Mystery of ancient Egypt solved? Tomb of queen may be hidden near King Tut'

By Washington Post, adapted by Newsela staff
08.17.15



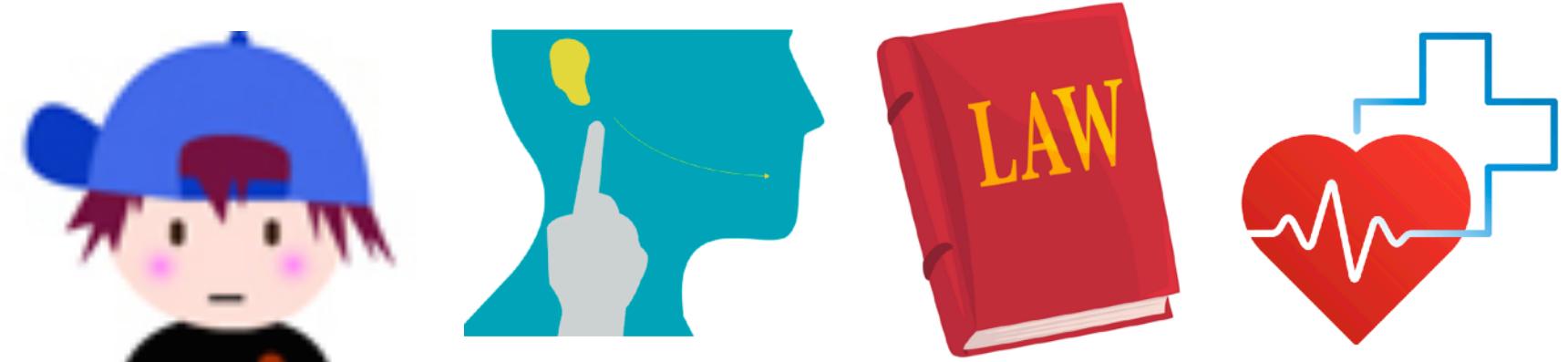
The 3,330-year-old bust of Nefertiti sits in an exhibition in the Kulturforum in Berlin, Germany, March 1, 2005.
Photo: AP/Herbert Knosowski

The ancient Egyptian Queen Nefertiti has long been at the center of a mystery. For years, archaeologists have wondered where her tomb might be hidden. Nefertiti belonged to the family line of the famous King Tutankhamun, better known as King Tut. Indeed, some believe she was Tut's mother. While the other royals in her line are

Why Text Simplification?

It can help a lot of people!

- Children (Leonardo et al., 2018) ← research on education using Newsela data
 - Second language learners (Housel et al., 2020) ← our collaborators at RIT
 - Deaf and hard-of-hearing students (Alonzo et al., 2020) ← our collaborators at RIT
 - People with dyslexia (Rello at al., 2013)
 - People with autism spectrum disorder (González-Navarro et al., 2014)
-
- and many others ... e.g., to read medical & legal documents, etc.



Automatic Text Simplification

Now, primarily addressed by sequence-to-sequence neural network models.

Input sentence:

Since 2010, project researchers have uncovered documents in Portugal that have revealed who owned the ship



Generated Output:

Scientists have found documents in Portugal.
They have also found out who owned the ship.

Automatic Text Simplification

However, SOTA neural generation models perform mostly deletion.

Input sentence:

According to Ledford, Northrop executives said they would build substantial parts of the bomber in Palmdale, creating about 1,500 jobs.

Generated output:

Programmer-interpreter
(Dong et al., 2019)

ledford **is a big group** of bomber in palmdale.

Rerank
(Kriz et al., 2019)

ledford **is** northrop.

Reinforcement Learning
(Zhang & Lapata, 2017)

, said they would build **palmdale** parts of **the substantial** in **creating**.

Part 1 — Controllable Generation Model



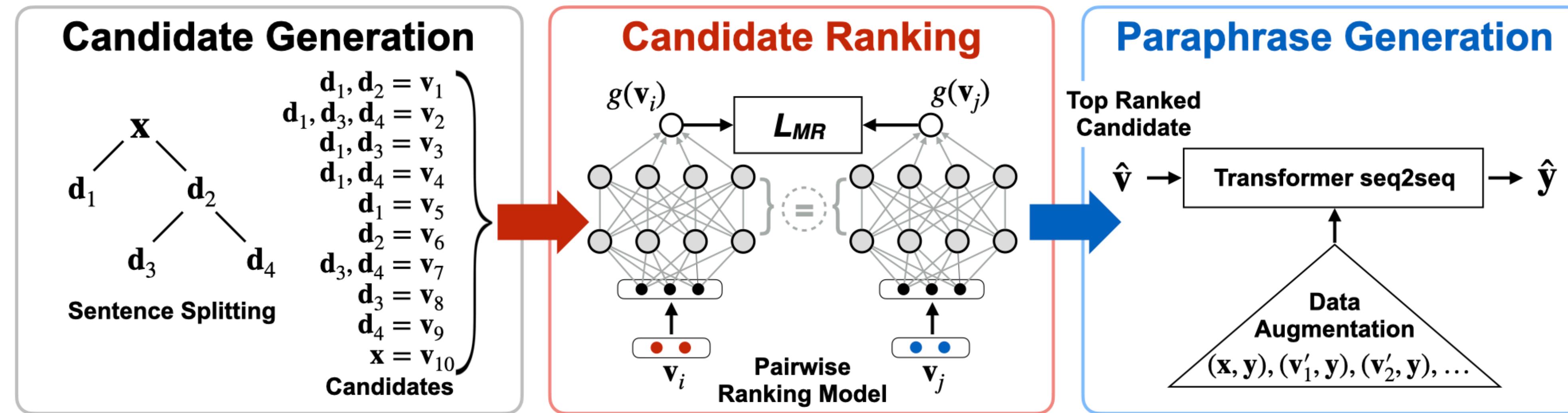
Controllable Text Simplification with Explicit Paraphrasing

Mounica Maddela, Fernando Alva-Manchego, Wei Xu (NAACL 2021)



Controllable Text Generation

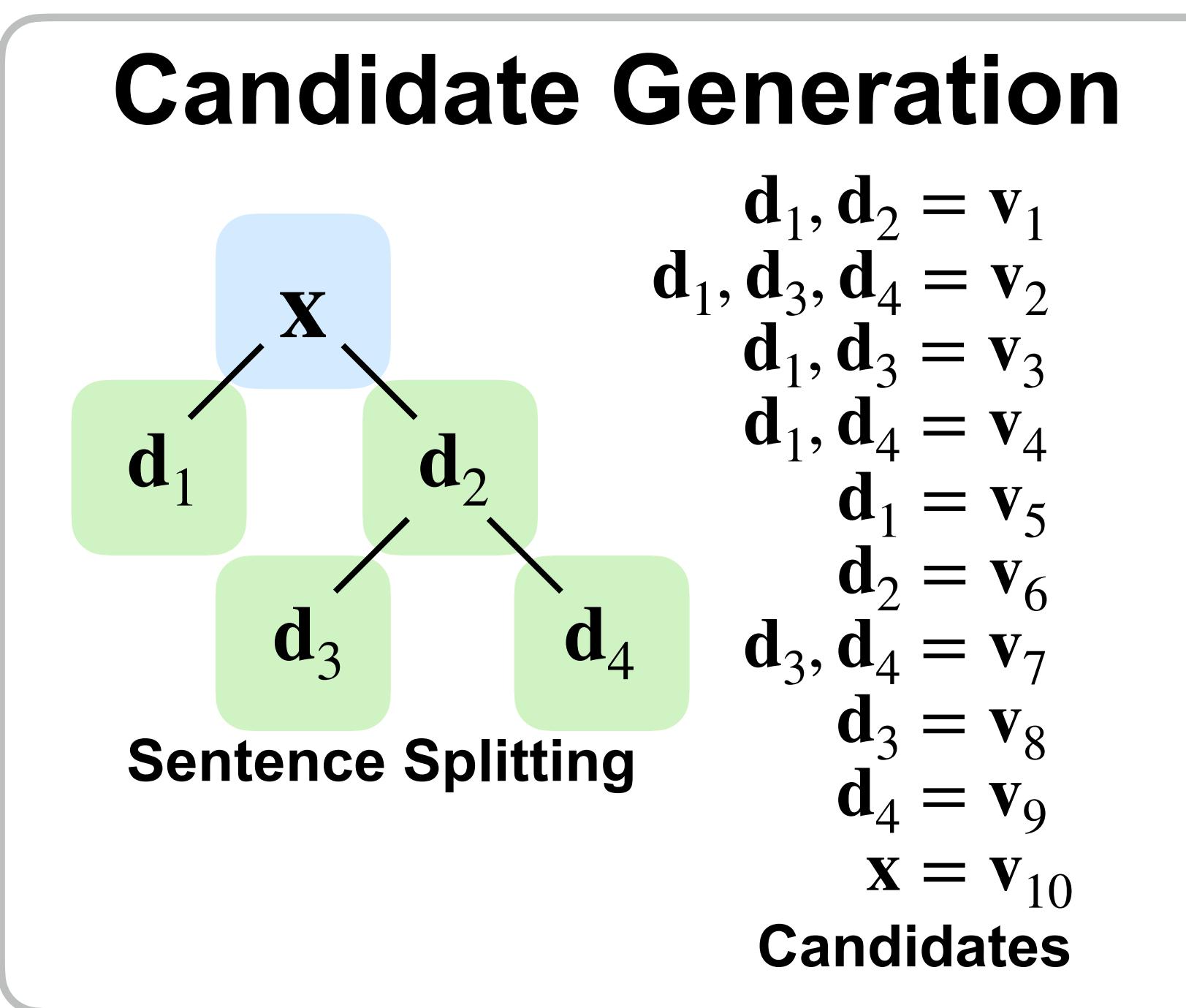
We incorporate linguistic knowledge into neural seq-to-seq models to improve controllability.



INPUT x : The exhibition, which opened Oct. 8 and runs through Jan. 3, features 27 self-portraits. **REFERENCE y** : The show started Oct. 8. It ends Jan. 3.
 d_1 : The exhibition features 27 self-portraits. **d_2** : The exhibition opened Oct. 8 and runs through Jan. 3.
 d_3 : The exhibition opened Oct. 8. **d_4** : The exhibition runs through Jan. 3. $\hat{v} = v_7$: The exhibition opened Oct. 8. The exhibition runs through Jan. 3.

Step 1 —

We use a rule-based method (Niklaus et al., 2019) + a seq2seq model for splitting and deletion.



Input sentence:

The exhibition, which opened Oct. 8 and runs through Jan. 3, features 27 self-portraits.

Split sentences:

The exhibition features 27 portraits.

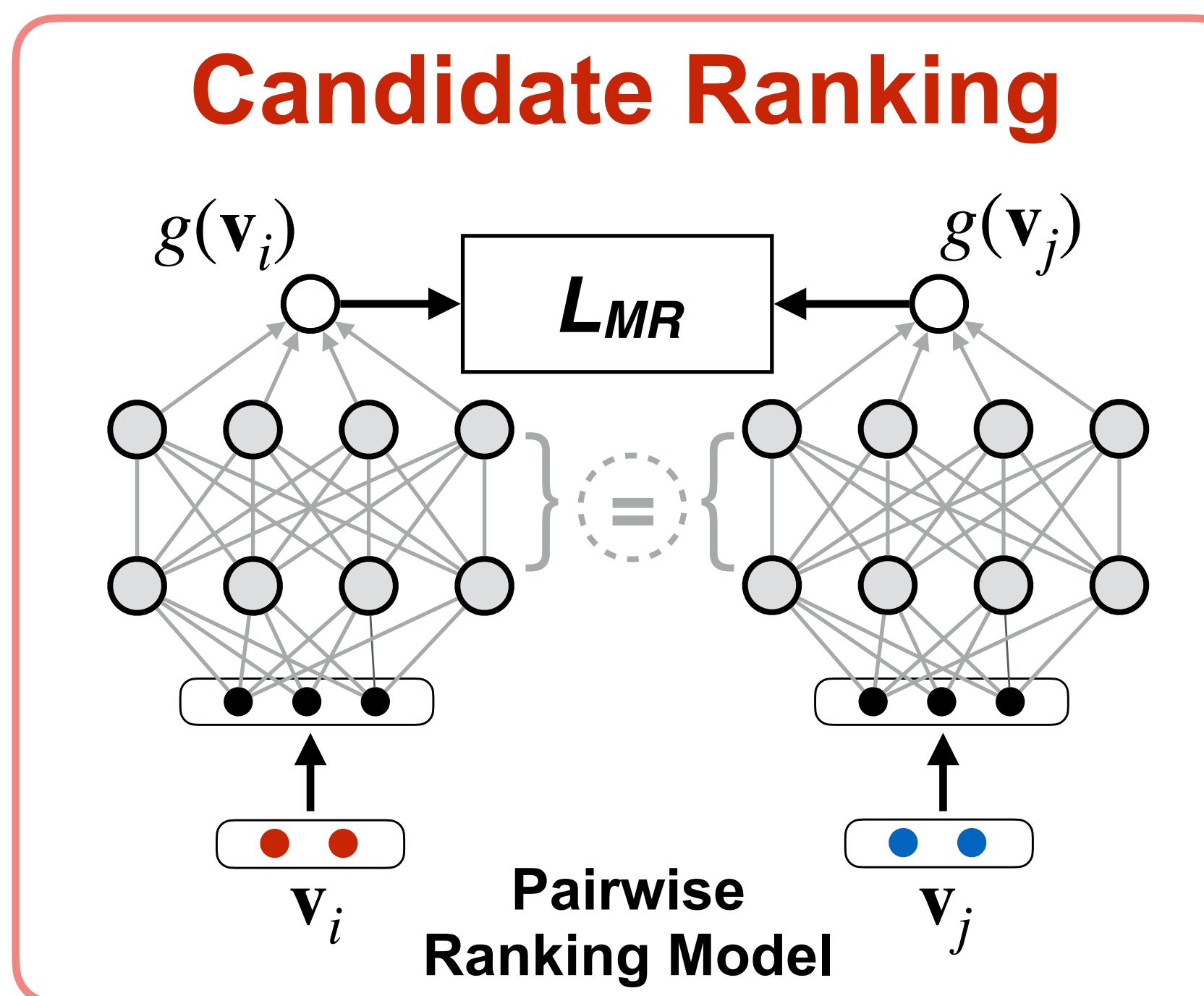
The exhibition opened Oct. 8 and runs through Jan. 3.

The exhibition opened Oct. 8.

The exhibition runs through Jan. 3.

Step 2 —

Then, we rank all the intermediate outputs (after splitting & deletion).



Candidates:

The exhibition opened Oct. 8. The exhibition runs through Jan. 3.

The exhibition opened Oct. 8 and runs through Jan. 3.

The exhibition features 27 portraits. The exhibition opened Oct. 8.

The exhibition features 27 portraits. The exhibition opened Oct. 8 and runs through Jan. 3.

The exhibition features 27 portraits.

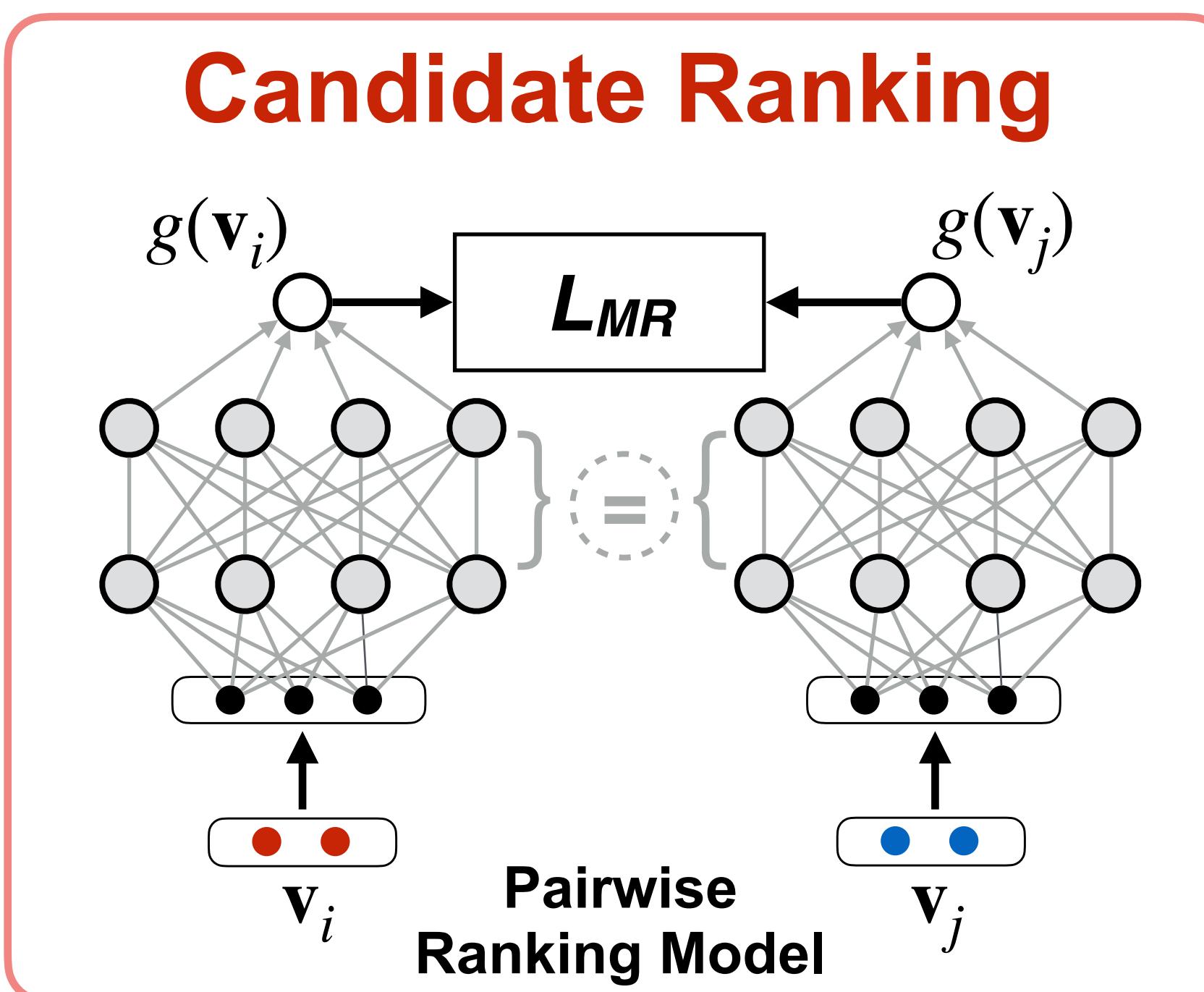
... (and more)

Human reference:

The show started Oct. 8. It ends. Jan 3.

Step 2 —

During training, we access each candidate using BERTScore (Zhang et al. 2019) with length penalty.



Scoring function:

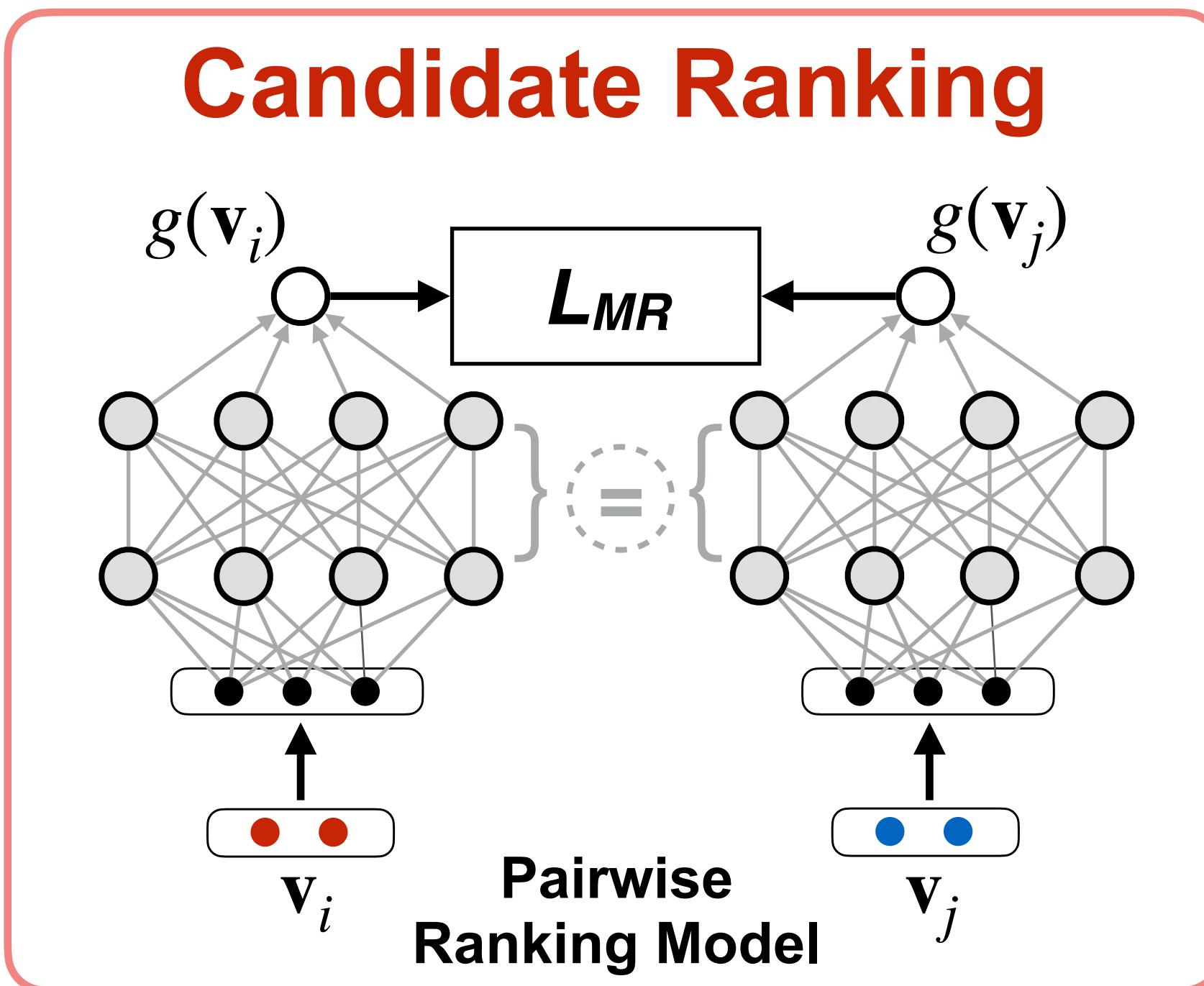
target compression ratio

$$g^*(\mathbf{v}_i, \mathbf{y}) = e^{-\lambda \|\phi_{\mathbf{v}_i} - \phi_{\mathbf{y}}\|} \times BERTScore(\mathbf{v}_i, \mathbf{y})$$

candidate reference

Step 2 —

During training, we access each candidate using BERTScore (Zhang et al. 2019) with length penalty.



Loss function:

$$L_{MR} = \frac{1}{m} \sum_{k=1}^m \frac{1}{n_k^2} \sum_{i=1}^{n_k} \sum_{j=1, i \neq j}^{n_k} \max(0, 1 - l_{ij}^k d_{ij}^k)$$

$$d_{ij}^k = g(\mathbf{v}_i^k) - g(\mathbf{v}_j^k)$$

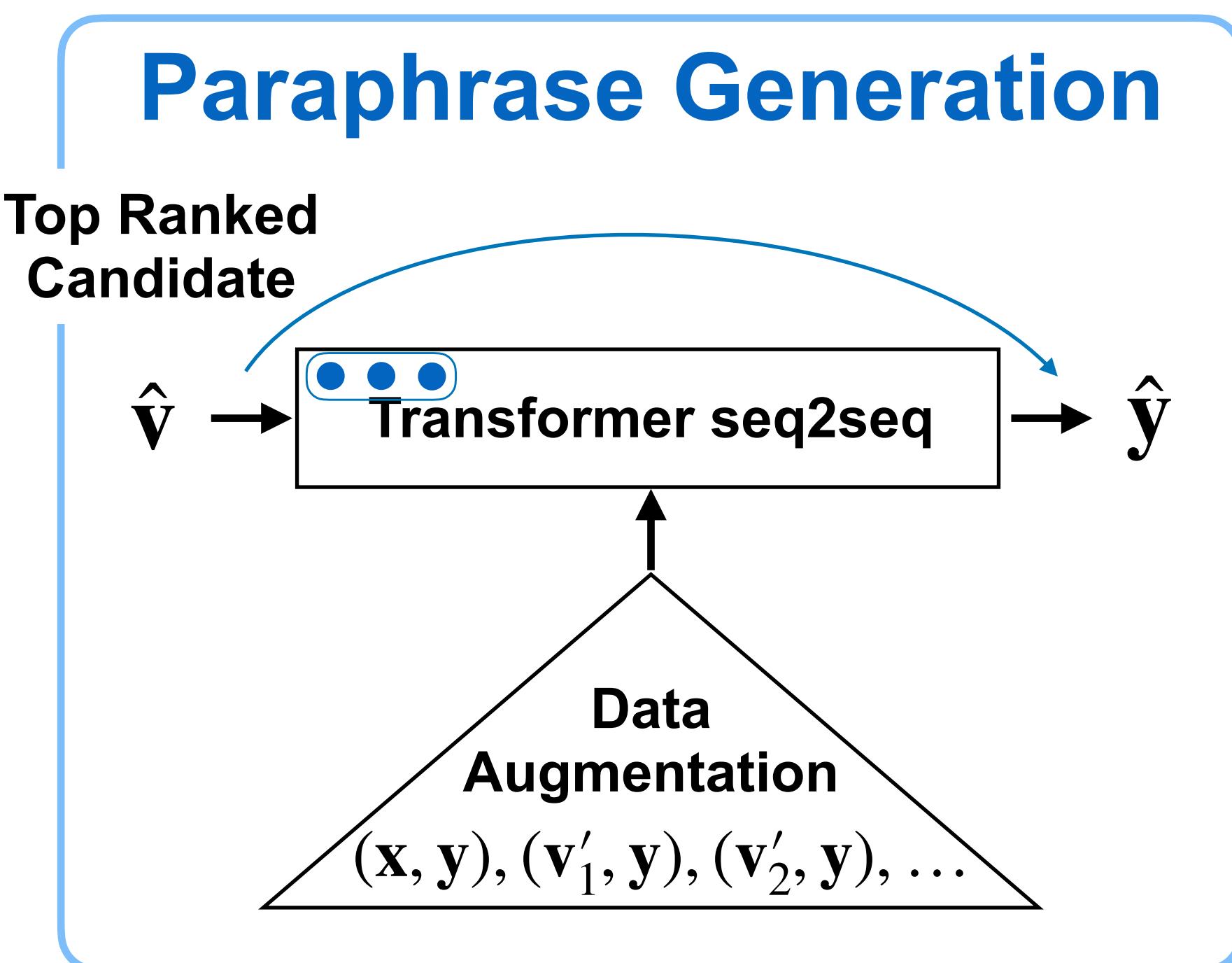
$$l_{ij}^k = \text{sign} \left(g^*(\mathbf{v}_i^k, \mathbf{y}^k) - g^*(\mathbf{v}_j^k, \mathbf{y}^k) \right)$$

Length-penalized BERTScore

Features: number of words in v_i and x , compression ratio of v_i with respect to x , Jaccard similarity between v_i and x , the rules applied on x to obtain v_i , and the number of rule applications.

Step 3 —

Finally, we have a paraphrase generation model trained with augmented training data.
(some selected candidates, in addition to the original input, are paired with the human reference)



Additional control over the degree of paraphrasing:

- A copy-control token as soft constraint.
- An auxiliary task (whether a word should be copied) using a **monolingual word aligner** to derive noisy training labels.

Controllable Text Generation

We can control the degree of sentence splitting, deletion, and paraphrasing.

Input: Experts say China's air pollution exacts a tremendous toll on human health.

Reference: China's air pollution is very unhealthy.

Our Model
($cp = 0.6$)

experts say china's air pollution **is a big problem for** human health.

Our Model
($cp = 0.7$)

experts say china's air pollution **can cause a lot of damage on** human health.

Our Model
($cp = 0.8$)

experts say china's air pollution **is a huge** toll on human health.

Hybrid-NG

experts say **government's** air pollution exacts a tremendous toll on human health.

LSTM

experts say china's air pollution exacts a tremendous toll on human health.

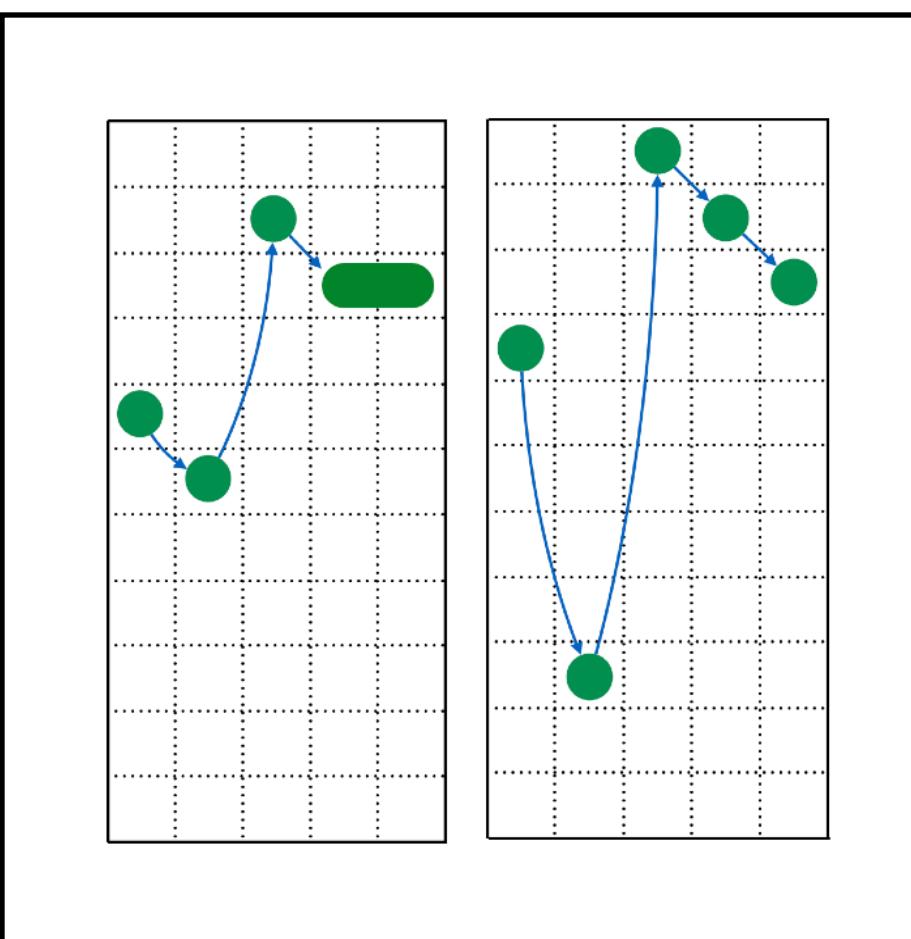
Transformer

experts say china's air pollution exacts a tremendous **effect** on human health.

EditNTS

experts say china's air pollution **can cause** human health.

Part 2 — Monolingual Word Alignment



Neural semi-Markov CRF for Monolingual Word Alignment

Wuwei Lan*, Chao Jiang*, Wei Xu (ACL 2021)



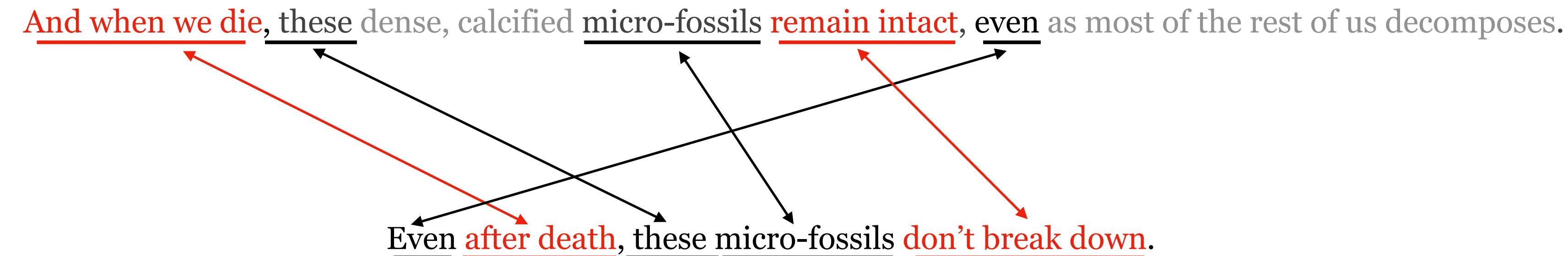
Span Monolingual Word Alignment

Can support not only text-to-text generation tasks, but also natural language understanding tasks.

Rephrase

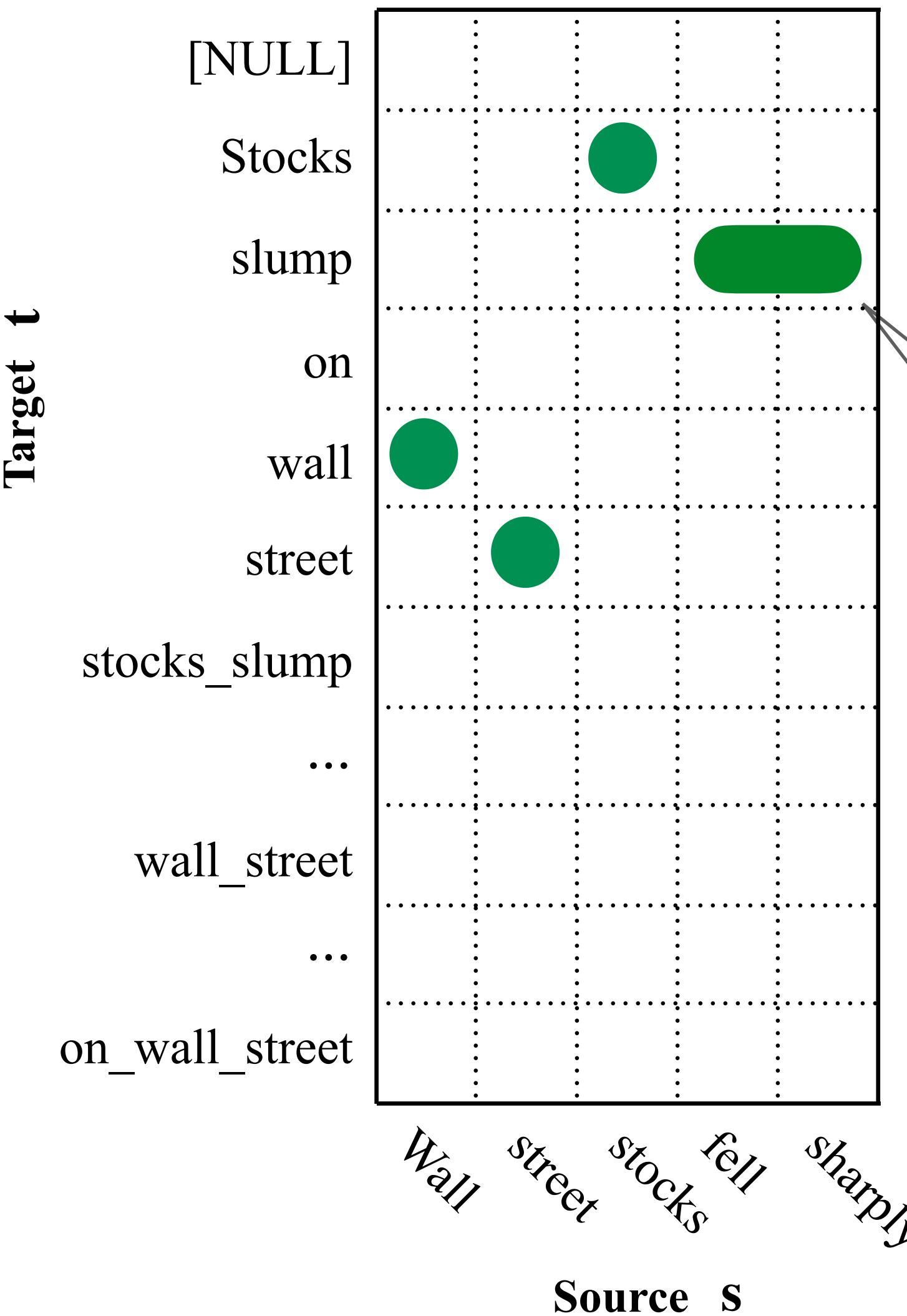
Keep

Delete



Semi-CRF Word Alignment Model

Span Interaction Matrix



Span representation based on SpanBERT (Joshi et al. 2020)

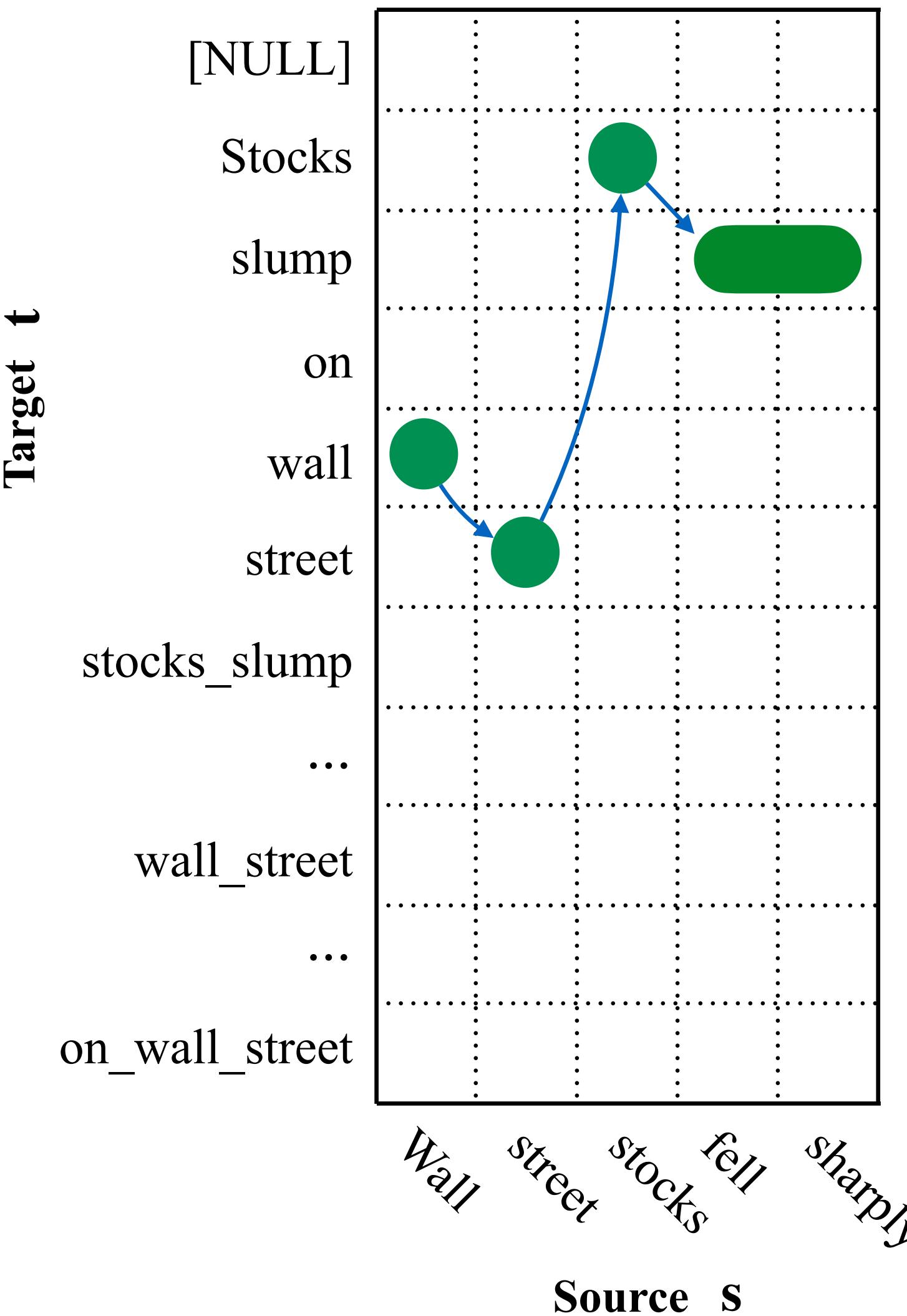
$$h_i^s = (e_{start(i)}; e_{end(i)}; attn_i)$$

$$score(s_i, t_j) = \text{FFNN}(h_i^s; h_j^t; |h_i^s - h_j^t|; h_i^s \circ h_j^t)$$

2-layer FFNN to capture semantic similarity between (s_i, t_j)

Semi-CRF Word Alignment Model

Alignment Label Transition



semi-Markov Conditional Random Fields for span alignment

$$\Psi(\mathbf{a}, \mathbf{s}, \mathbf{t}) = \sum_i score(s_i, t_{a_i}) + T(a_{i-1}, a_i) + cost(\mathbf{a}, \mathbf{a}^*)$$

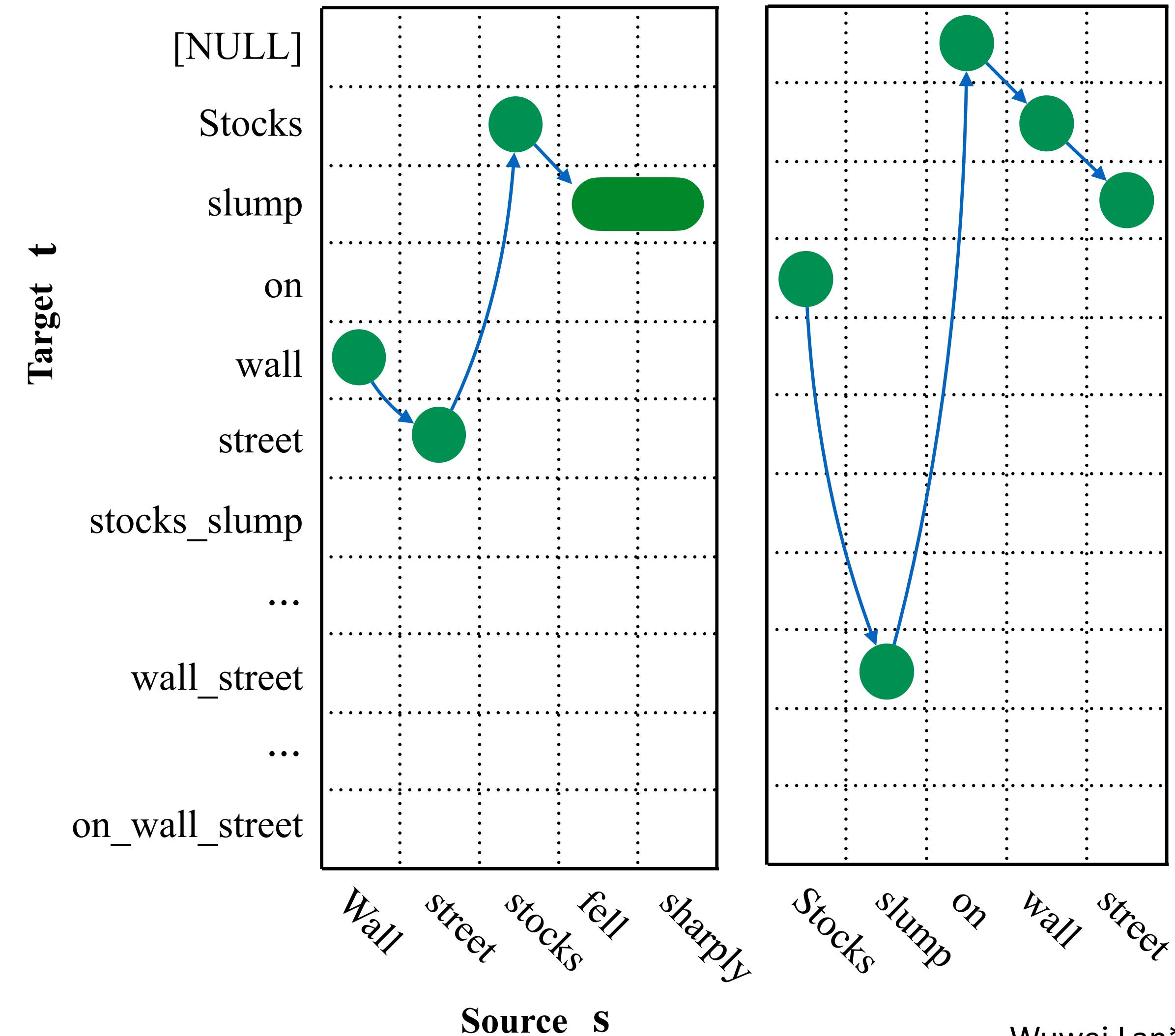
Negative Log-likelihood Loss Hamming Loss

$$P(\mathbf{a} | \mathbf{s}, \mathbf{t}) = \frac{\exp (\Psi(\mathbf{a}, \mathbf{s}, \mathbf{t}))}{\sum_{\mathbf{a} \in A} \exp (\Psi(\mathbf{a}, \mathbf{s}, \mathbf{t}))}$$

all possible alignments over variable length spans

Semi-CRF Word Alignment Model

Bi-directional Training / Decoding



Training objective:

$$\sum_{s,t,a} -\log P(a_{s2t} | s, t) - \log P(a_{t2s} | t, s)$$

Source-to-target

Target-to-source

Decoding:

Viterbi-like Algorithm + Intersect + Expand

Experiments on MultiMWA Benchmark

We annotate a Multi-Genre Monolingual Word Alignment dataset that covers four different text genres.

	In-domain	Out-of-domain		
		MTReference	Newsela	arXiv
JacanaToken (Yao et al. 2013a)	76.2	79.8	95.8	95.8
JacanaPhrase (Yao et al. 2013b)	75.8	79.4	93.7	94.9
PipelineAligner (Sultan et al. 2014)	74.8	80.3	96.5	97.1
Our Neural CRF aligner	90.8	86.6	95.7	97.0
Our Neural semi-CRF aligner	92.4	87.2	97.3	97.4

🚀 16.2 F1

🚀 6.9 F1

🚀 0.8 F1

🚀 0.3 F1



Our Work — Automatic Text Simplification

- **Controllable Generation Model**

- Neural semi-Markov CRF for Monolingual Word Alignment (Lan*, Jiang* & Xu, ACL 2021)
 - Also useful for semantics and natural language understanding.
- Controllable Text Simplification with Explicit Paraphrasing (Maddela, Alva-Manchego & Xu, NAACL 2021)
 - How to incorporate linguistic rules with neural networks?

- **High-quality Training Data**

- Neural CRF Model for Sentence Alignment in Text Simplification (Jiang, Maddela, Lan, Zhong & Xu, ACL 2020)
 - Performance gains from better data are huge!
- BiSECT: Learning to Split and Rephrase Sentences with Bitexts (Kim*, Maddela*, Kriz, Xu, Callison-Burch, EMNLP 2021)
- Discourse Level Factors for Sentence Deletion in Text Simplification (Zhong, Jiang, Xu & Li, AAAI 2020)
- A Neural Readability Ranking Model and A Word-Complexity Lexicon for Lexical Simplification (Maddela & Xu, EMNLP 2018)
- Optimizing Statistical Machine Translation for Text Simplification (Xu et al., TACL 2016)
- Problems in Current Text Simplification Research: New Data Can Help (Xu et al., TACL 2015)



SimplePPDB++

A database of 14.1 million paraphrase rules with improved complexity ranking scores.

Paraphrase Rule	Score
→ <i>self-supporting</i>	0.93
<i>self-reliant</i> → <i>self-sufficient</i>	0.48
→ <i>self-sustainable</i>	-0.60
→ <i>possible</i>	0.94
<i>viable</i> → <i>realistic</i>	0.15
→ <i>plausible</i>	-0.91
→ <i>in-depth review</i>	0.89
<i>detailed assessment</i> → <i>careful examination</i>	0.28
→ <i>comprehensive evaluation</i>	-0.87

complex

All our code/data/models are open-source! You can find them at: <https://cocoxu.github.io/>

NLP X Research Lab

We design machine learning algorithms to help computer to understand and generate human languages.
 X = Machine Learning, Education, Accessibility, Social Media, ...



Natural Language Processing

- Text generation and evaluation
- Semantics and structured prediction
- Misinformation, social media
- Offensive/biased language
- Interactive NLP systems



Wei Xu
Assistant Professor

Mounica Maddela
PhD student

Chao Jiang
PhD student

Yao Dou
PhD student

Jagriti Sikka
MS student

Angana Borah
MS student

Machine Learning

- Sequence-to-sequence models
- Controllability of neural networks
- Cross-lingual transfer learning
- Learning from noisy data
- Data augmentation



David Heineman
Undergrad

Ema Goh
Undergrad

Jonathan Zheng
Undergrad

Michael Ryan
Undergrad

Alex Kelso
Undergrad

Dylan Small
Undergrad

Srushti Nandu
Intern