

Language Modelling Makes Sense

BERT-based Word Sense Disambiguation
and Medical Entity Linking

Daniel Loureiro, 12th Dec 2019



Sense Embeddings

Exploiting the latest Neural Language Models (NLMs) for sense-level representation learning.



Sense Embeddings

Exploiting the latest Neural Language Models (NLMs) for sense-level representation learning.

- Beat SOTA for English Word Sense Disambiguation (WSD).
- Full WordNet in NLM-space (+100K common sense concepts).
- Concept-level analysis of NLMs.

Related Work

Related Work

[Iacobacci et al. (2016)]
[Zhong and Ng (2010)]

**Bag-of-Features
Classifiers**
(SVM)

[Luo et al. (2018b)]
[Luo et al. (2018a)]
[Vial et al. (2018)]
[Raganato et al. (2017)]

**Deep Sequence
Classifiers**
(BiLSTM)

[Peters et al. (2018)]
[Melamud et al. (2016)]
[Yuan et al. (2016)]

**Sense-level
Representations**
(k-NN)
(over NLM reprs.)

Related Work

[Iacobacci et al. (2016)]
[Zhong and Ng (2010)]

**Bag-of-Features
Classifiers**
(SVM)



[Luo et al. (2018b)]
[Luo et al. (2018a)]
[Vial et al. (2018)]
[Raganato et al. (2017)]

**Deep Sequence
Classifiers**
(BiLSTM)

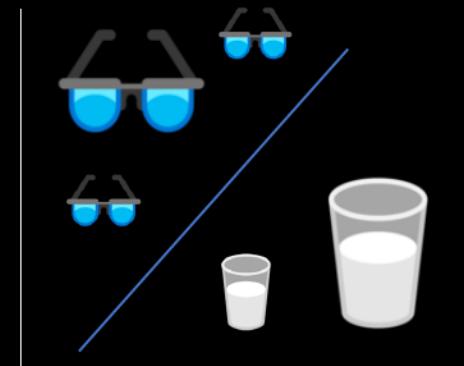
[Peters et al. (2018)]
[Melamud et al. (2016)]
[Yuan et al. (2016)]

**Sense-level
Representations**
(k-NN)
(over NLM reprs.)

Bag-of-Features Classifiers

It Makes Sense (IMS) [Zhong and Ng (2010)] :

- POS tags, surrounding words, local collocations.
- SVM for each word type in training.
- Fallback: Most Frequent Sense (MFS).
- Improved with word embedding features. [Iacobacci et al. (2016)]
- Still competitive (!)

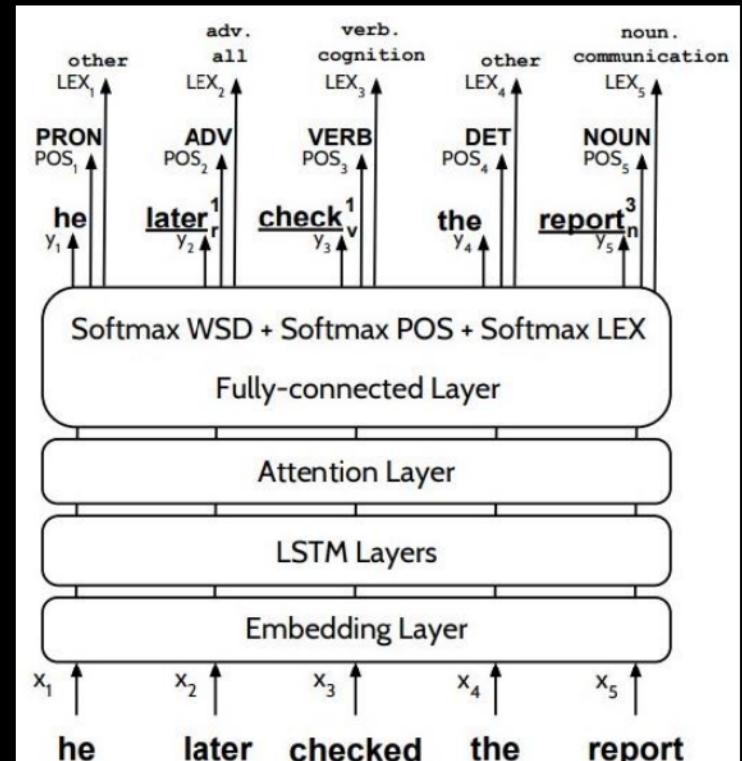


“glasses”

Deep Sequence Classifiers

Bi-directional LSTMs (BiLSTMs):

- Better with:
 - Attention (as everything else).
 - Auxiliary losses. (POS, lemmas, lexnames) [Raganato et al. (2017)]
 - Glosses, via co-attention mechanisms. [Luo et al. (2018)]
- Still must fallback on MFS.
- Not that much better than bag-of-features...

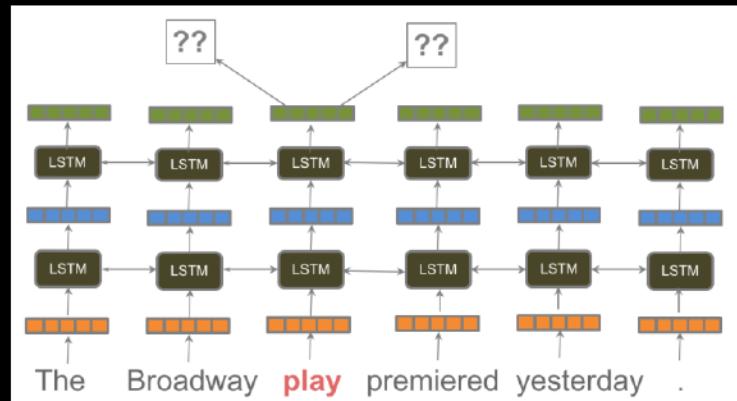


[Raganato et al. (2017)]

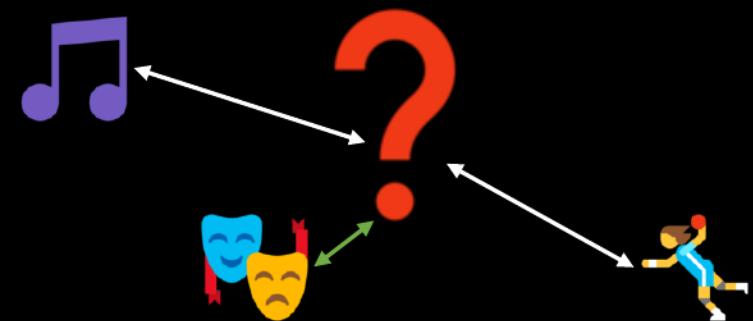
Contextual k-NN

Matching Contextual Word Embeddings:

- Produce Sense Embeddings from NLMs (averaging).
- Sense embs. can be compared with contextual embs.
- Disambiguation = Nearest Neighbour search (1-NN).
- Sense embs. limited to annotations. MFS required.
- Promising, but early attempts.



[Ruder (2018)]



Our Approach

Our Approach

- Expand the k-NN approach to full-coverage of WordNet.

Our Approach

- Expand the k-NN approach to full-coverage of WordNet.
- Matching senses becomes trivial, no MFS fallbacks needed.

Our Approach

- Expand the k-NN approach to full-coverage of WordNet.
- Matching senses becomes trivial, no MFS fallbacks needed.
- Full-set of sense embeddings in NLM-space is useful beyond WSD.

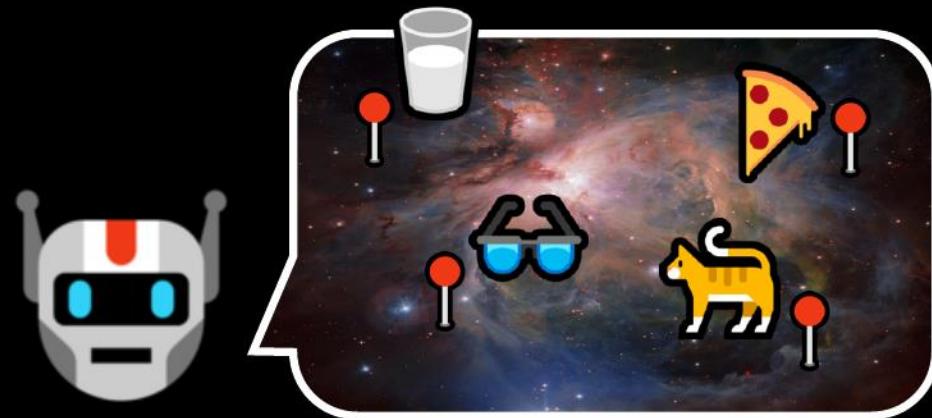
Our Approach

- Expand the k-NN approach to full-coverage of WordNet.
- Matching senses becomes trivial, no MFS fallbacks needed.
- Full-set of sense embeddings in NLM-space is useful beyond WSD.



Our Approach

- Expand the k-NN approach to full-coverage of WordNet.
- Matching senses becomes trivial, no MFS fallbacks needed.
- Full-set of sense embeddings in NLM-space is useful beyond WSD.



Challenges

Challenges

- Overcome very limited sense annotations (covers 16% senses).

Challenges

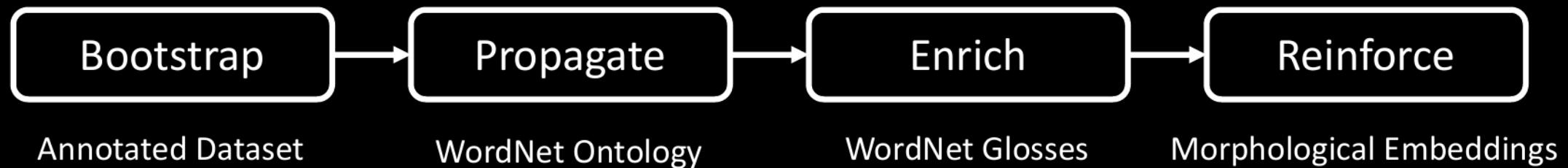
- Overcome very limited sense annotations (covers 16% senses).
- Infer missing senses correctly so that task performance improves.

Challenges

- Overcome very limited sense annotations (covers 16% senses).
- Infer missing senses correctly so that task performance improves.
- Rely only on sense embeddings, no lemma or POS features.

Challenges

- Overcome very limited sense annotations (covers 16% senses).
- Infer missing senses correctly so that task performance improves.
- Rely only on sense embeddings, no lemma or POS features.

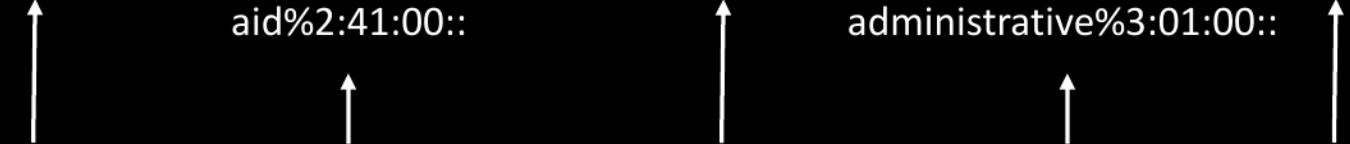


Bootstrapping Sense Embeddings

Can your insurance company aid you in reducing administrative costs ?

Would it be feasible to limit the menu in order to reduce feeding costs ?

Bootstrapping Sense Embeddings

insurance_company%1:14:00:: reduce%2:30:00:: cost%1:21:00::
aid%2:41:00::
Can your insurance company aid you in reducing administrative costs ?


Would it be feasible to limit the menu in order to reduce feeding costs ?

feasible%5:00:00:possible:00 menu%1:10:00:: feeding%1:04:01::
limit%2:30:00:: reduce%2:30:00:: cost%1:21:00::

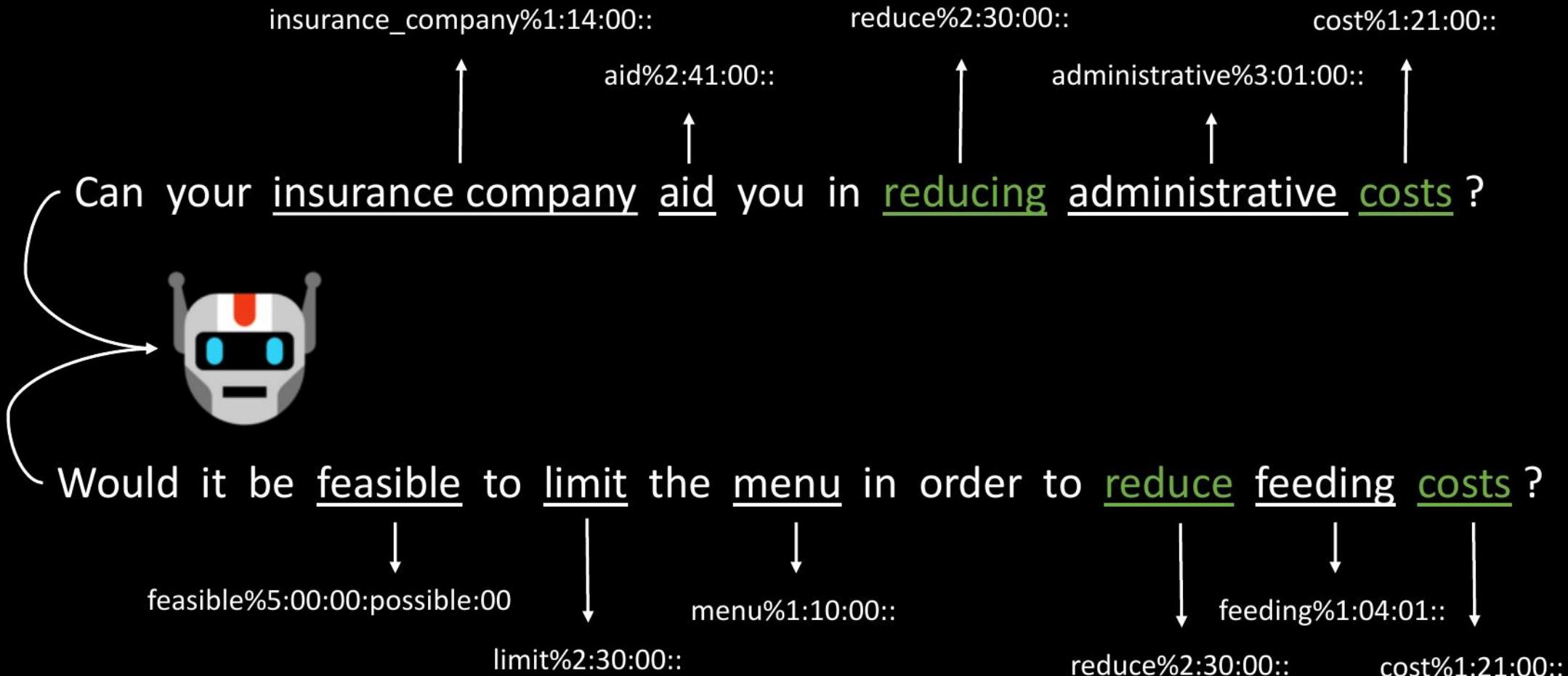

Bootstrapping Sense Embeddings

insurance_company%1:14:00:: reduce%2:30:00:: cost%1:21:00::
aid%2:41:00:: administrative%3:01:00::
Can your insurance company aid you in reducing administrative costs ?

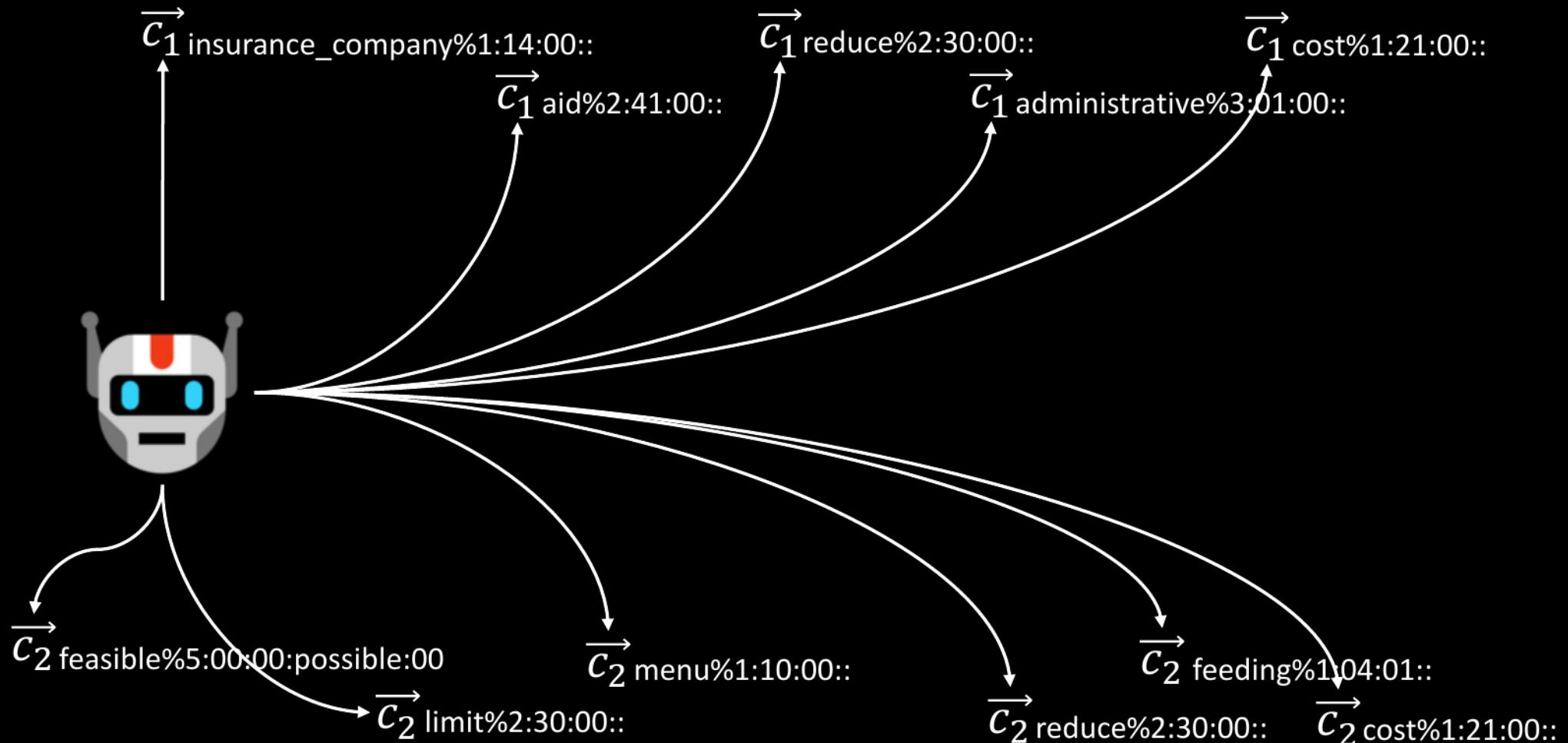
Would it be feasible to limit the menu in order to reduce feeding costs ?



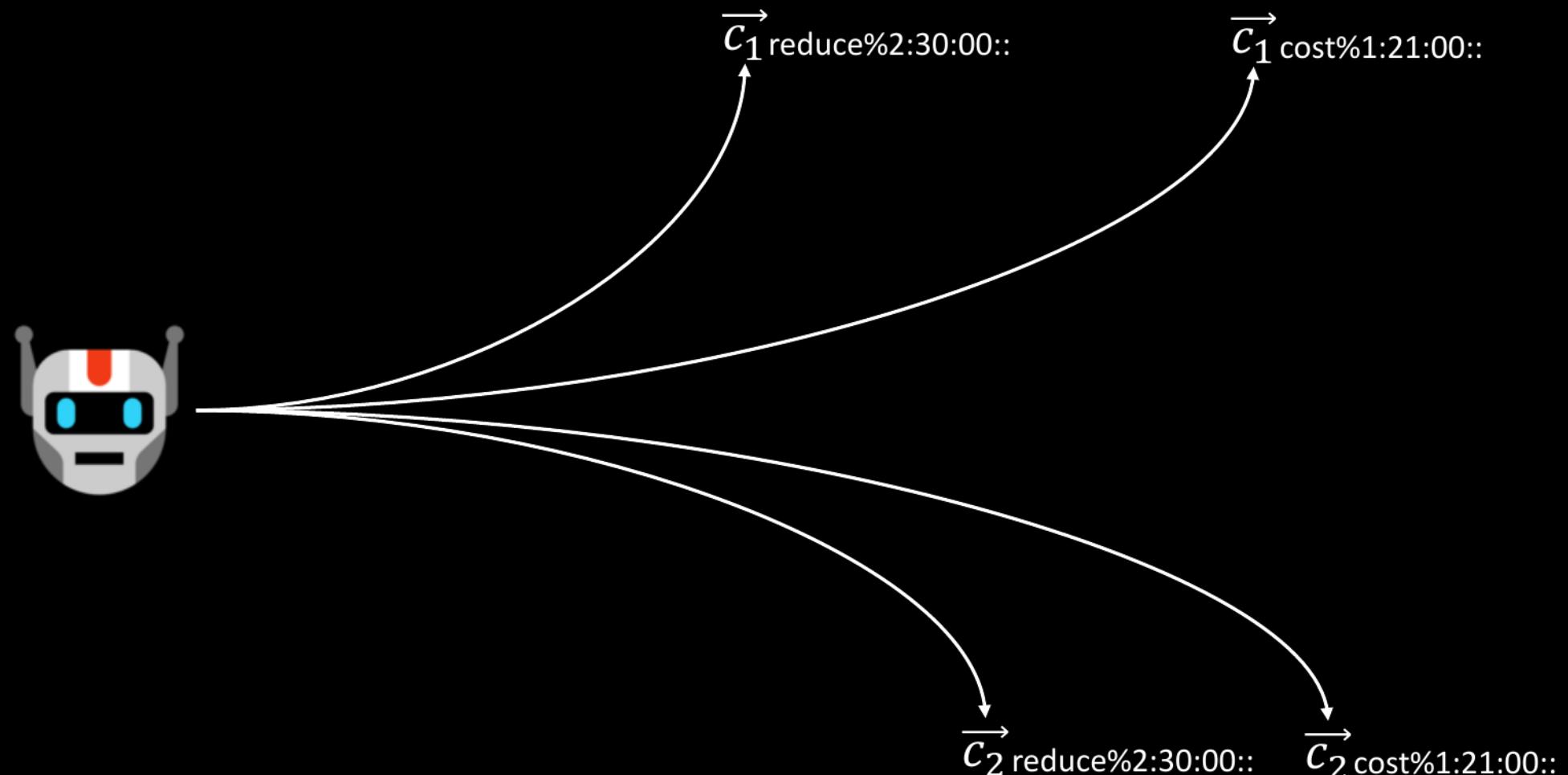
Bootstrapping Sense Embeddings



Bootstrapping Sense Embeddings



Bootstrapping Sense Embeddings



Bootstrapping Sense Embeddings

$$\vec{v}_{\text{reduce}} = \frac{\vec{c}_1_{\text{reduce}} + \vec{c}_2_{\text{reduce}} + \dots + \vec{c}_n_{\text{reduce}}}{n}$$

$$\vec{v}_{\text{cost}} = \frac{\vec{c}_1_{\text{cost}} + \vec{c}_2_{\text{cost}} + \dots + \vec{c}_n_{\text{cost}}}{n}$$

Bootstrapping Sense Embeddings

$$\vec{v}_{\text{reduce}\%2:30:00::} = \frac{\vec{c}_1_{\text{reduce}\%2:30:00::} + \vec{c}_2_{\text{reduce}\%2:30:00::} + \dots + \vec{c}_n_{\text{reduce}\%2:30:00::}}{n}$$

$$\vec{v}_{\text{cost}\%1:21:00::} = \frac{\vec{c}_1_{\text{cost}\%1:21:00::} + \vec{c}_2_{\text{cost}\%1:21:00::} + \dots + \vec{c}_n_{\text{cost}\%1:21:00::}}{n}$$

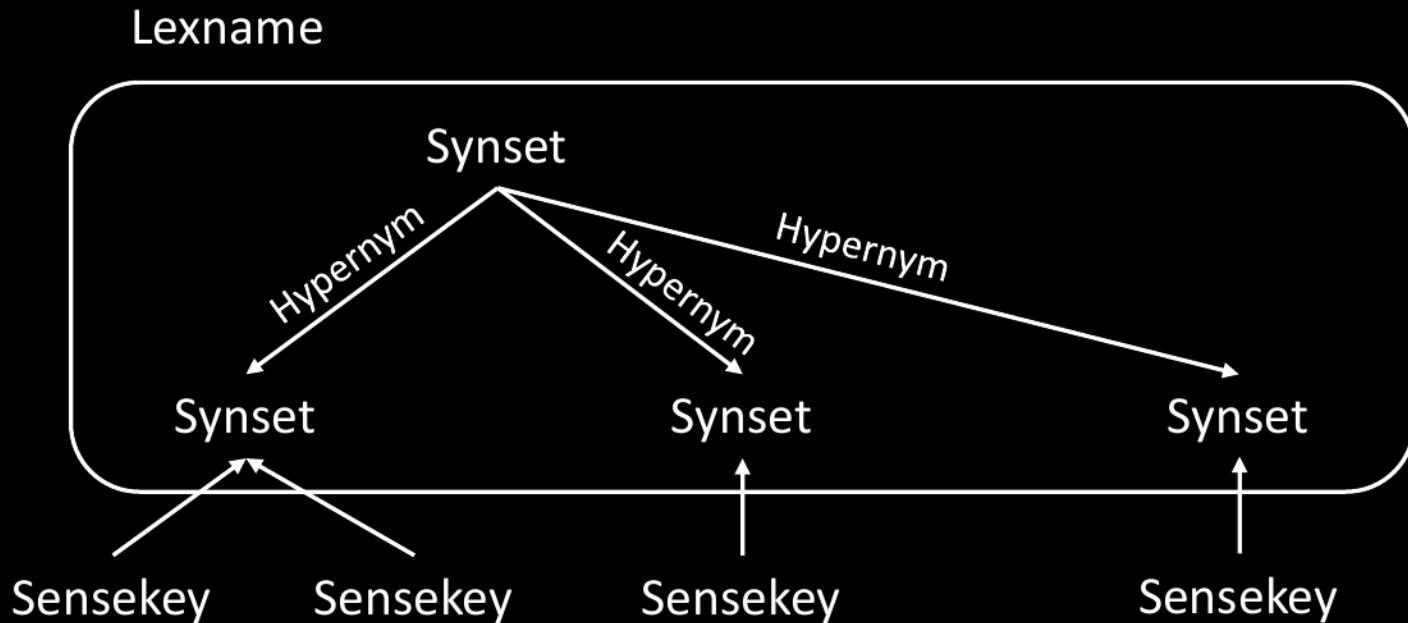
Outcome: 33,360 sense embeddings (16% coverage)

Propagating Sense Embeddings

WordNet's units, synsets, represent concepts at different levels.

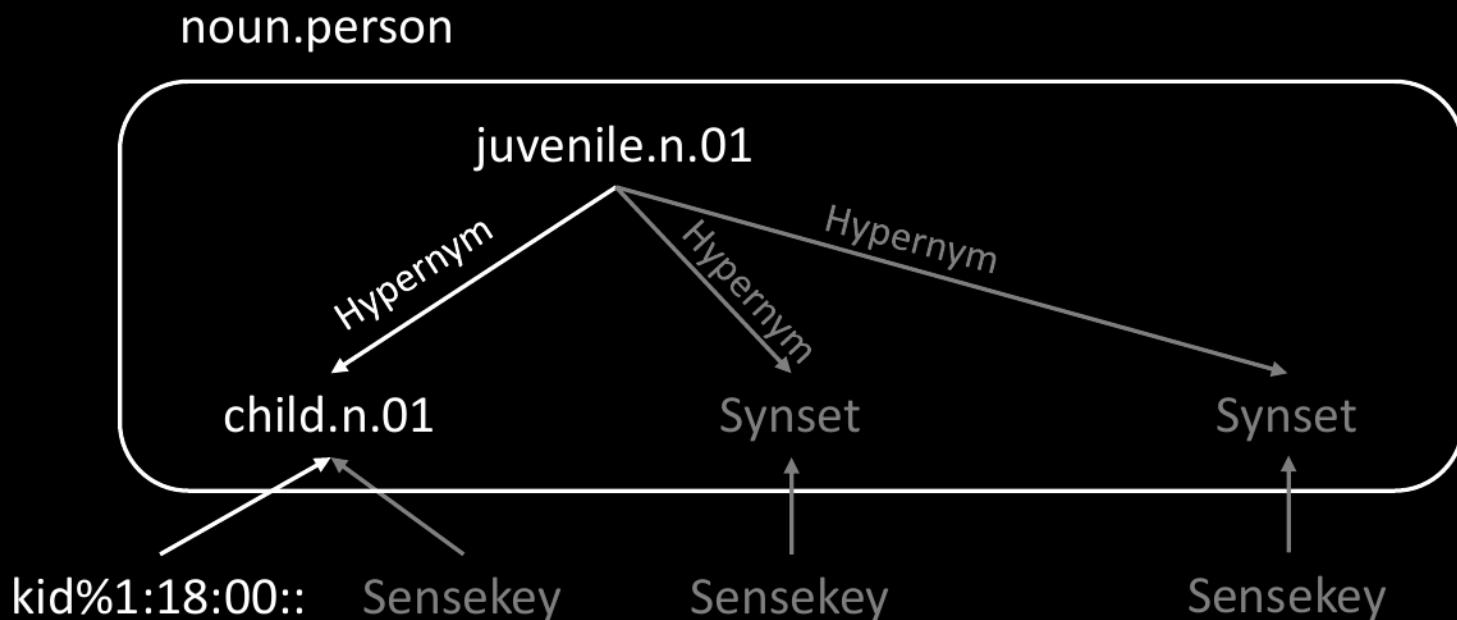
Propagating Sense Embeddings

WordNet's units, synsets, represent concepts at different levels.



Propagating Sense Embeddings

WordNet's units, synsets, represent concepts at different levels.



Propagating Sense Embeddings

burger%1:13:00::

hotdog%1:18:00::

hamburger%1:13:01::

sandwich%1:13:00::

wrap%1:13:00::

potato_chip%1:13:00::

Propagating Sense Embeddings

burger%1:13:00::

hotdog%1:18:00::

hamburger%1:13:01::

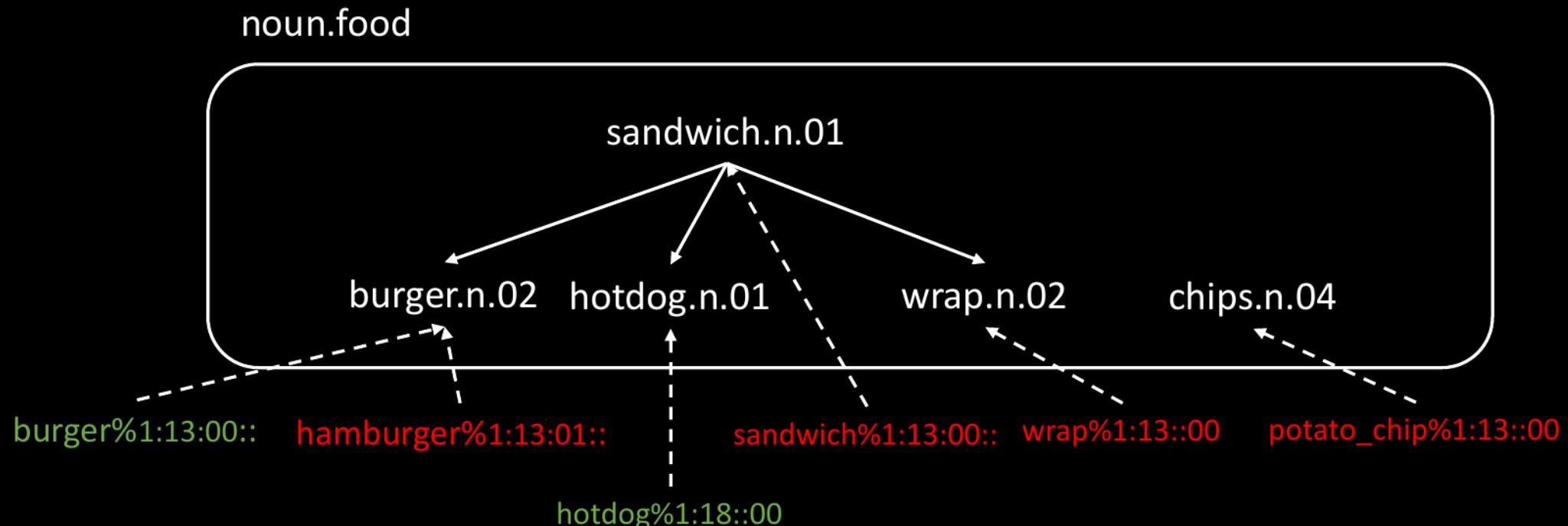
sandwich%1:13:00::

wrap%1:13:00::

potato_chip%1:13:00::

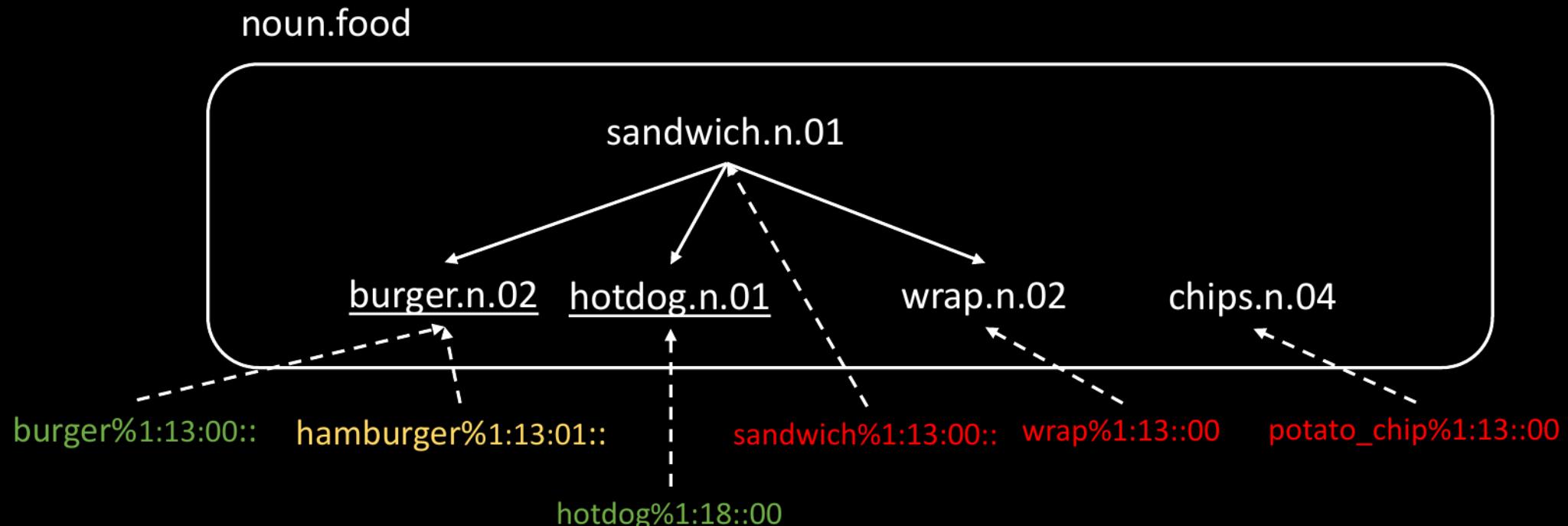
Propagating Sense Embeddings

Retrieve Synsets, Relations and Categories



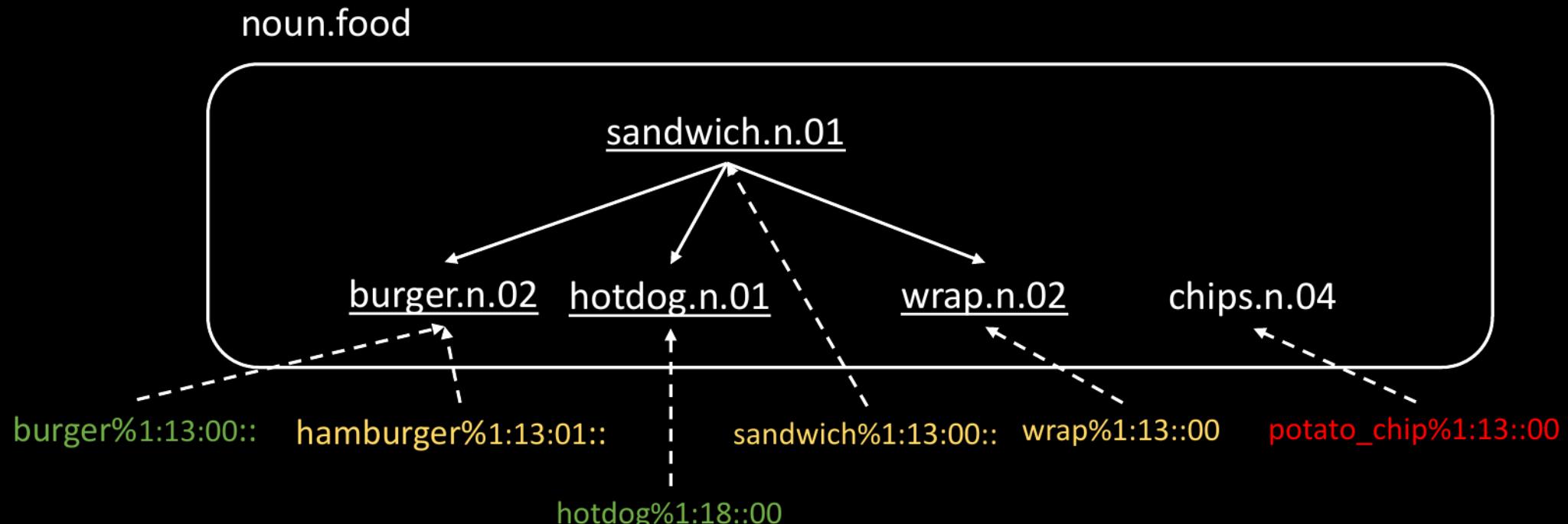
Propagating Sense Embeddings

1st stage: Synset Embeddings



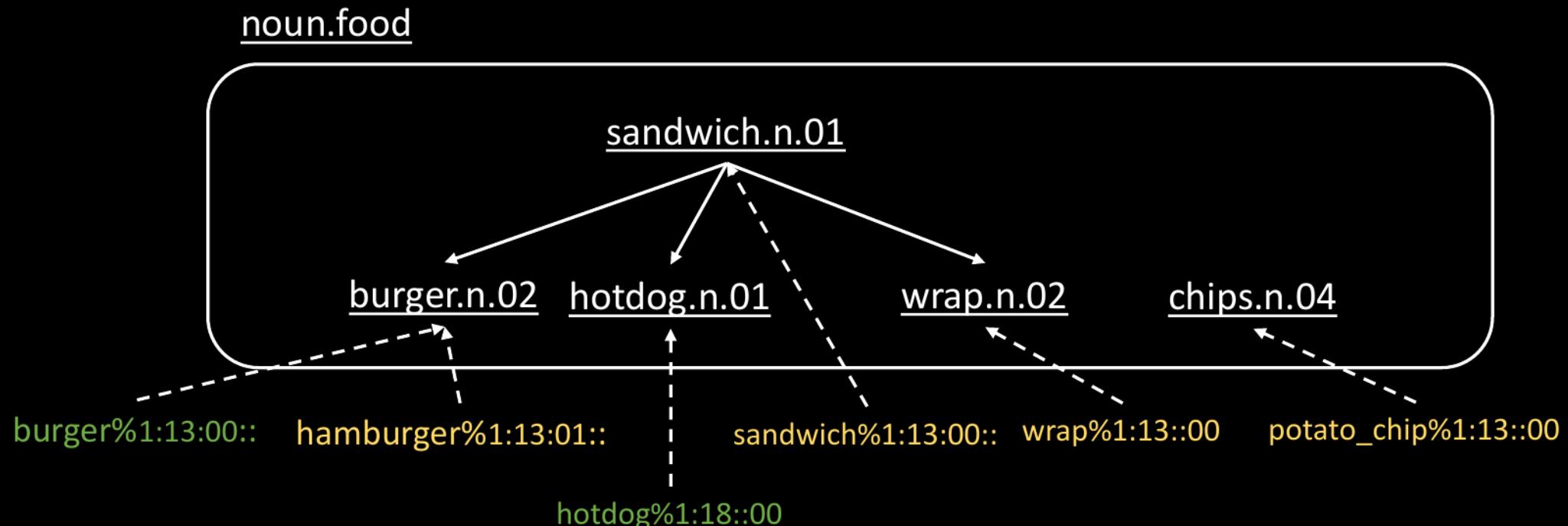
Propagating Sense Embeddings

2nd Stage: Hypernym Embeddings (ind. Synsets)



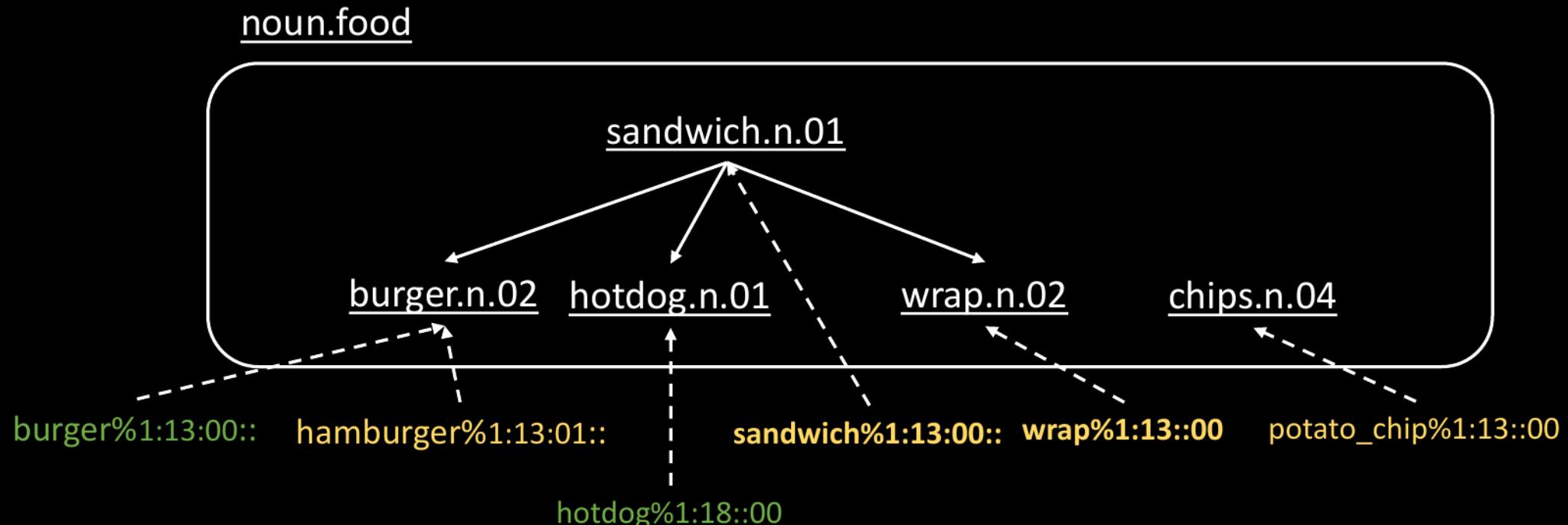
Propagating Sense Embeddings

3rd Stage: Lexname Embeddings



Propagating Sense Embeddings

But  !=  ...



Enriching Sense Embeddings

Leverage Synset Definitions and Lemmas for Differentiation

Enriching Sense Embeddings

Leverage Synset Definitions and Lemmas for Differentiation



sandwich:%1:13:00:: (sandwich.n.01)

Definition: two (or more) slices of bread with a filling between them

Lemmas: sandwich



wrap:%1:13:00:: (wrap.n.02)

Definition: a sandwich in which the filling is rolled up in a soft tortilla

Lemmas: wrap, tortilla

Enriching Sense Embeddings

Compose a new context



sandwich:%1:13:00:: (sandwich.n.01)

sandwich - two (or more) slices of bread with a filling between them



wrap:%1:13:00:: (wrap.n.02)

wrap, tortilla - a sandwich in which the filling is rolled up in a soft tortilla

Enriching Sense Embeddings

Make the context specific to sensekey (repeat lemma)



sandwich:%1:13:00::

sandwich - sandwich - two (or more) slices of bread with a filling between them



wrap%1:13:00::

wrap - wrap, tortilla - a sandwich in which the filling is rolled up in a soft tortilla

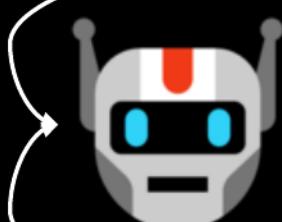
Enriching Sense Embeddings

Make the context specific to sensekey (repeat lemma)



sandwich:%1:13:00::

sandwich - sandwich - two (or more) slices of bread with a filling between them

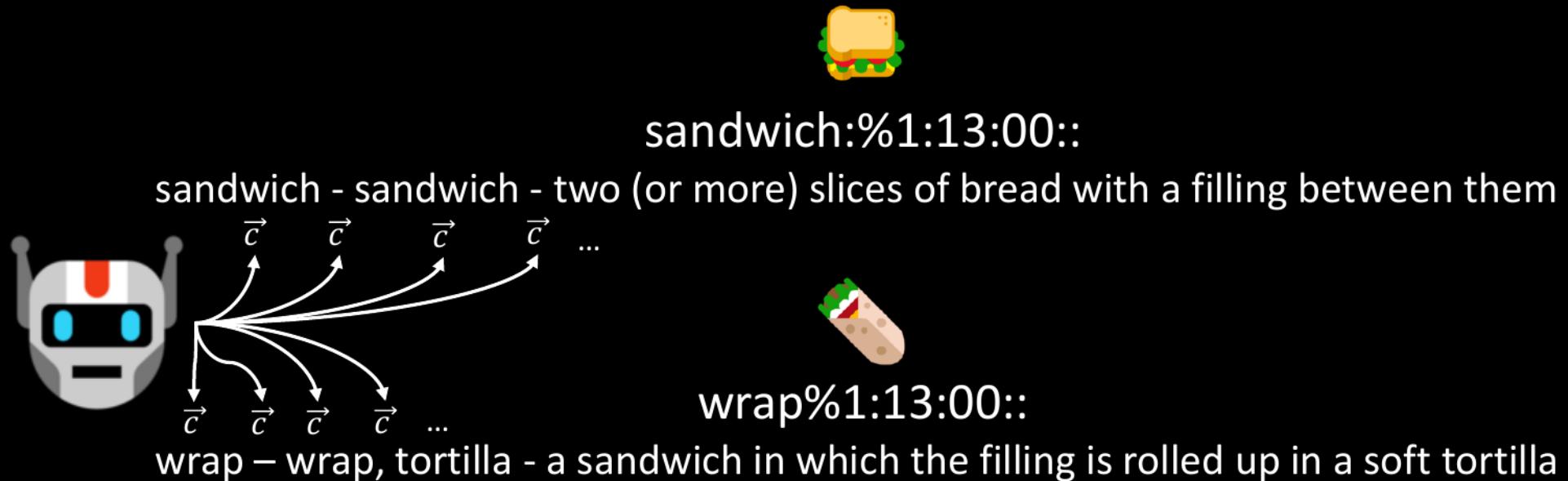


wrap%1:13:00::

wrap - wrap, tortilla - a sandwich in which the filling is rolled up in a soft tortilla

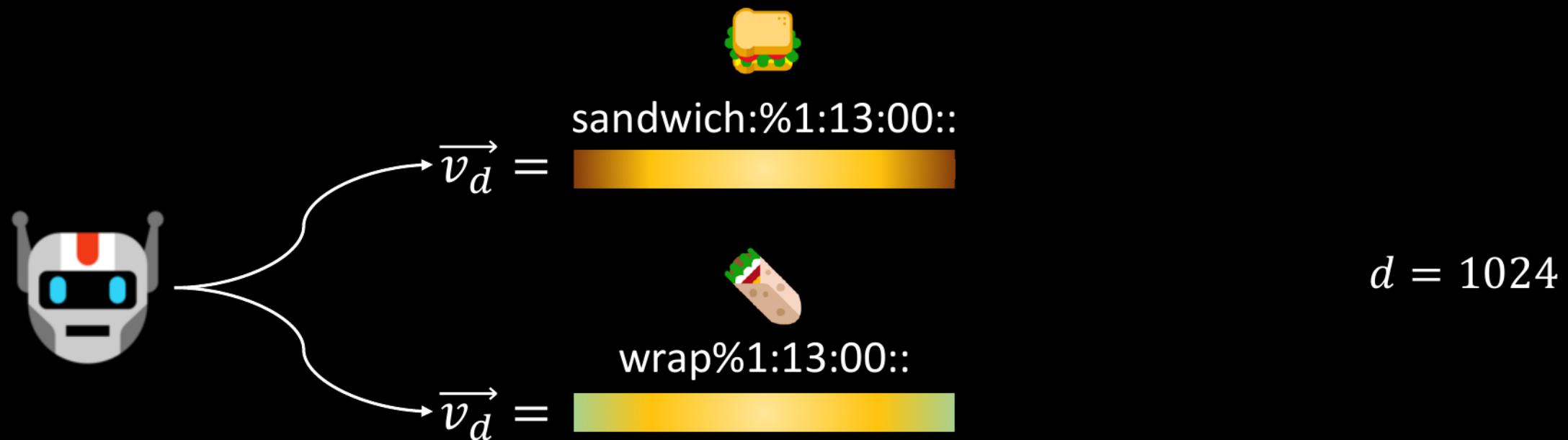
Enriching Sense Embeddings

Obtain contextual embeddings for every token



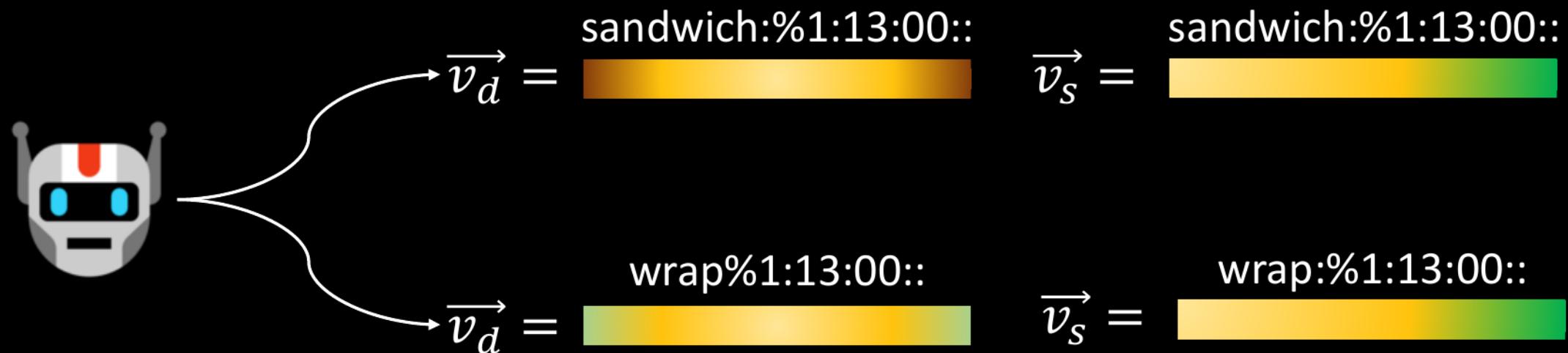
Enriching Sense Embeddings

Sentence Embedding from avg. of Contextual Embeddings



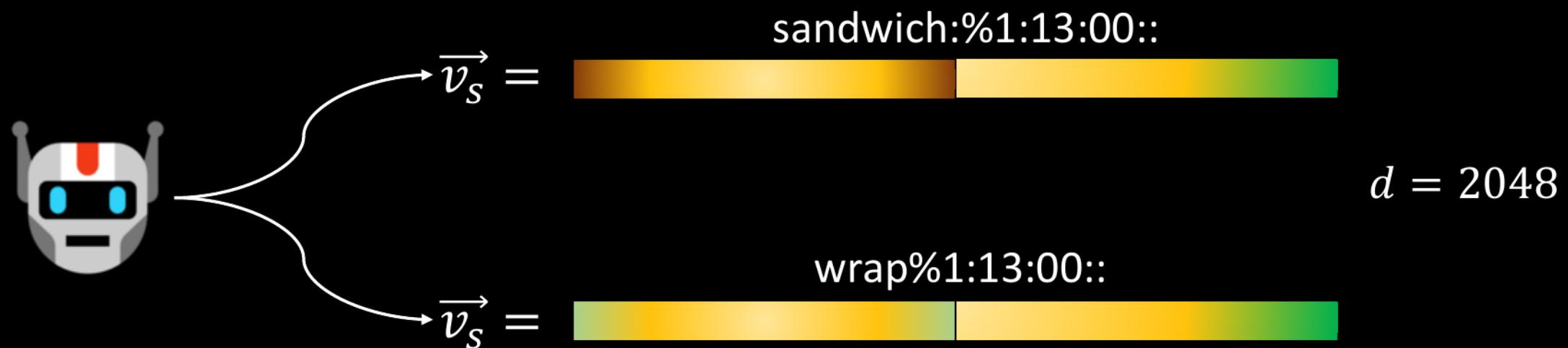
Enriching Sense Embeddings

Merge Sentence Embedding with previous Sense Embedding



Enriching Sense Embeddings

Merge Sentence Embedding with previous Sense Embedding

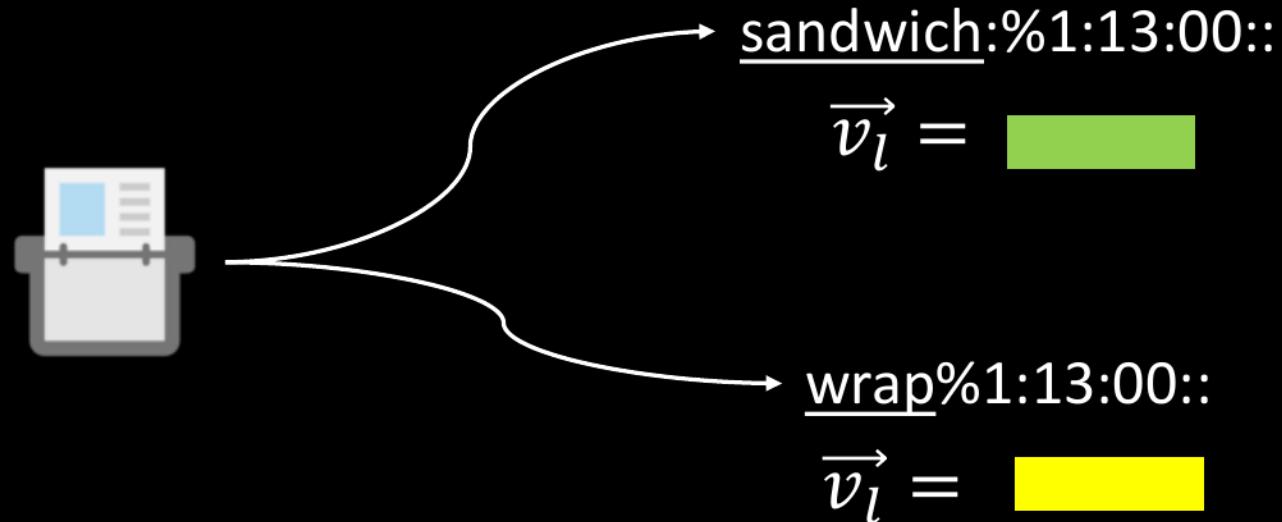


Reinforcing Sense Embeddings

Contextual Embeddings aren't good at preserving morphological relatedness

Reinforcing Sense Embeddings

Retrieve char-ngram embeddings (static) for lemmas



Reinforcing Sense Embeddings

Merge with previous sense embeddings

sandwich:%1:13:00::



wrap%1:13:00::



Reinforcing Sense Embeddings

Merge with previous sense embeddings

sandwich:%1:13:00::

$$\vec{v}_s = \text{[color bar]} \quad d = 2348$$


wrap%1:13:00::

$$\vec{v}_s = \text{[color bar]}$$


Matching Sense Embeddings

The glasses are in the cupboard.

Matching Sense Embeddings



The glasses are in the cupboard.

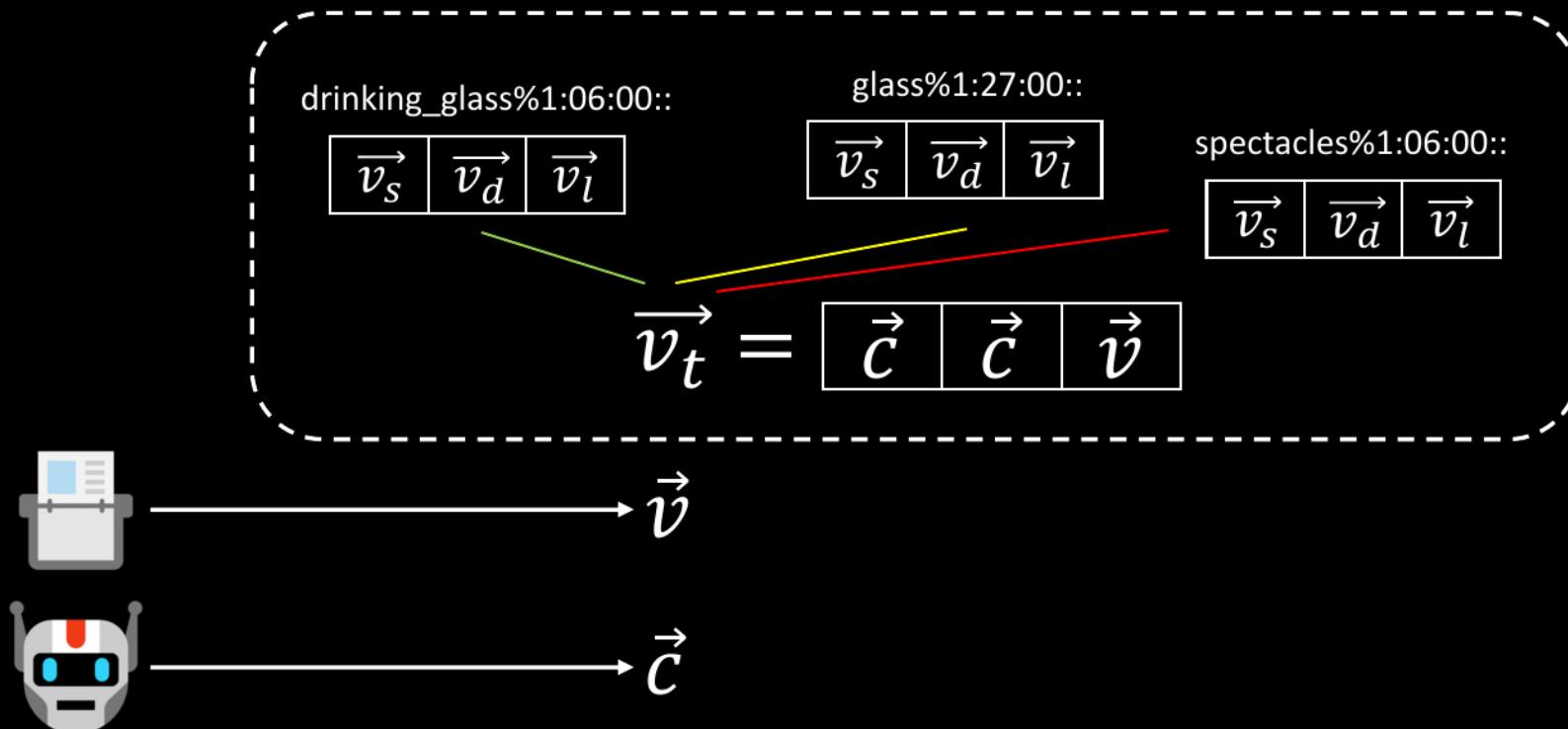
Matching Sense Embeddings

$$\vec{v}_t = \boxed{\vec{c} \mid \vec{c} \mid \vec{v}}$$



The glasses are in the cupboard.

Matching Sense Embeddings



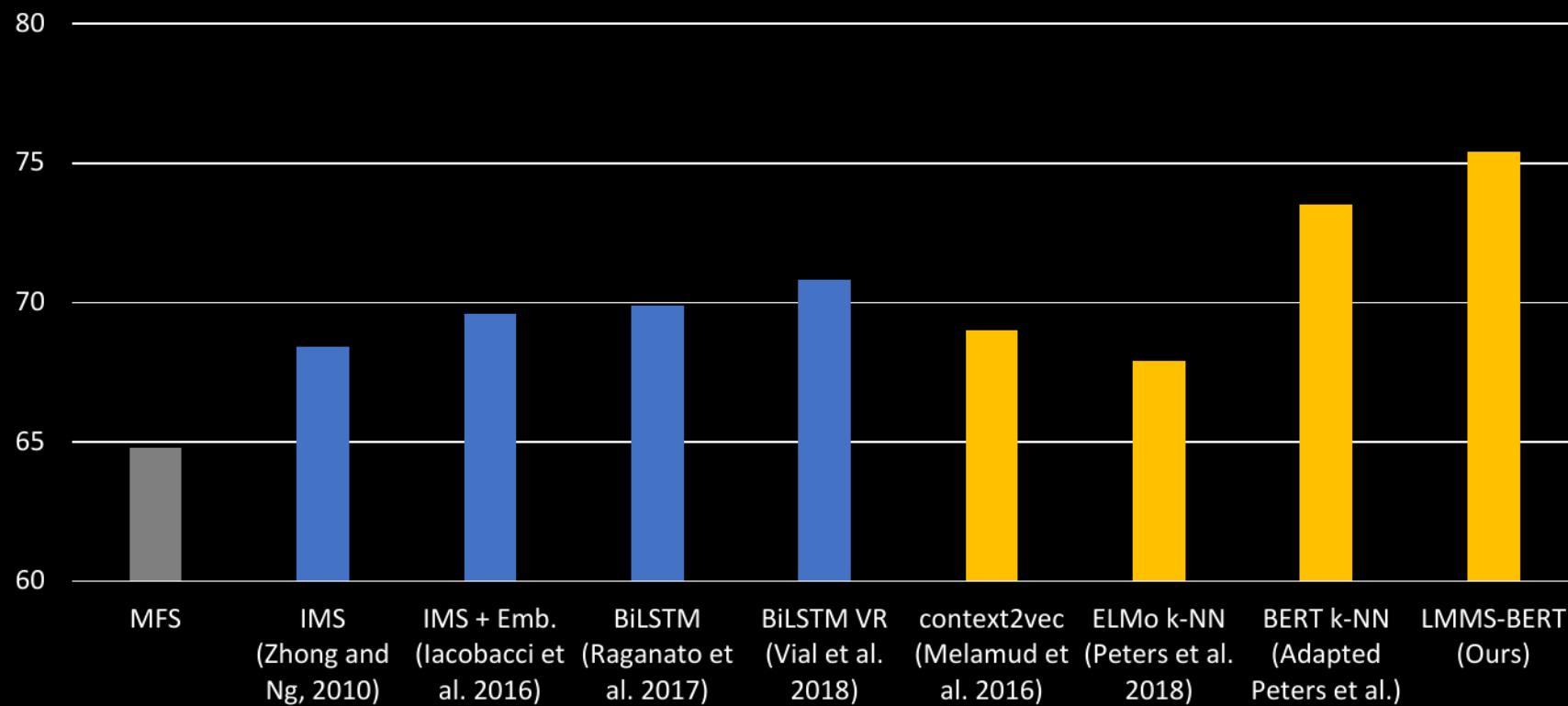
The glasses are in the cupboard.

WSD Results

WSD Results

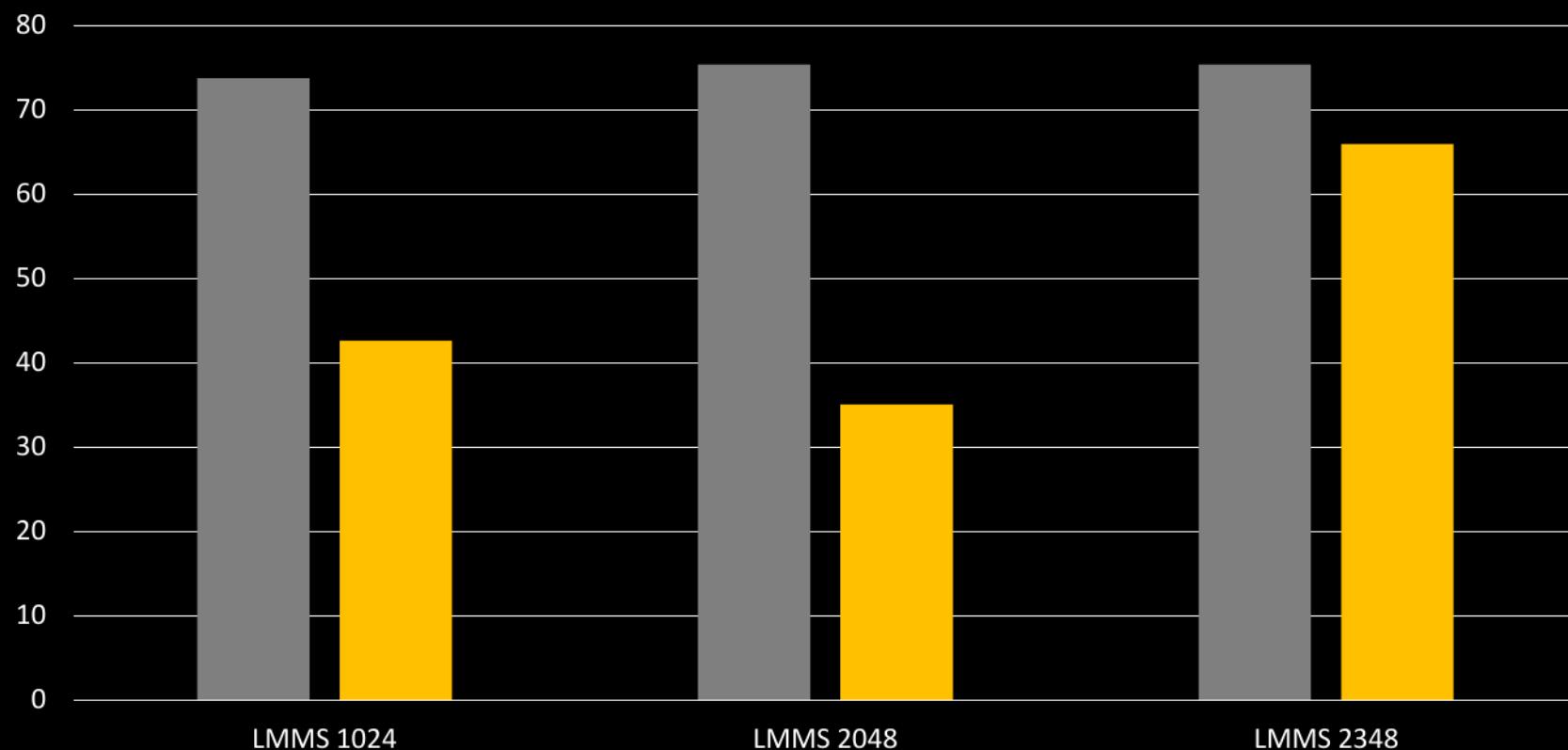
Standard English WSD Evaluation

F1 on ALL set of the WSD Evaluation Framework (Raganato et al. 2017)

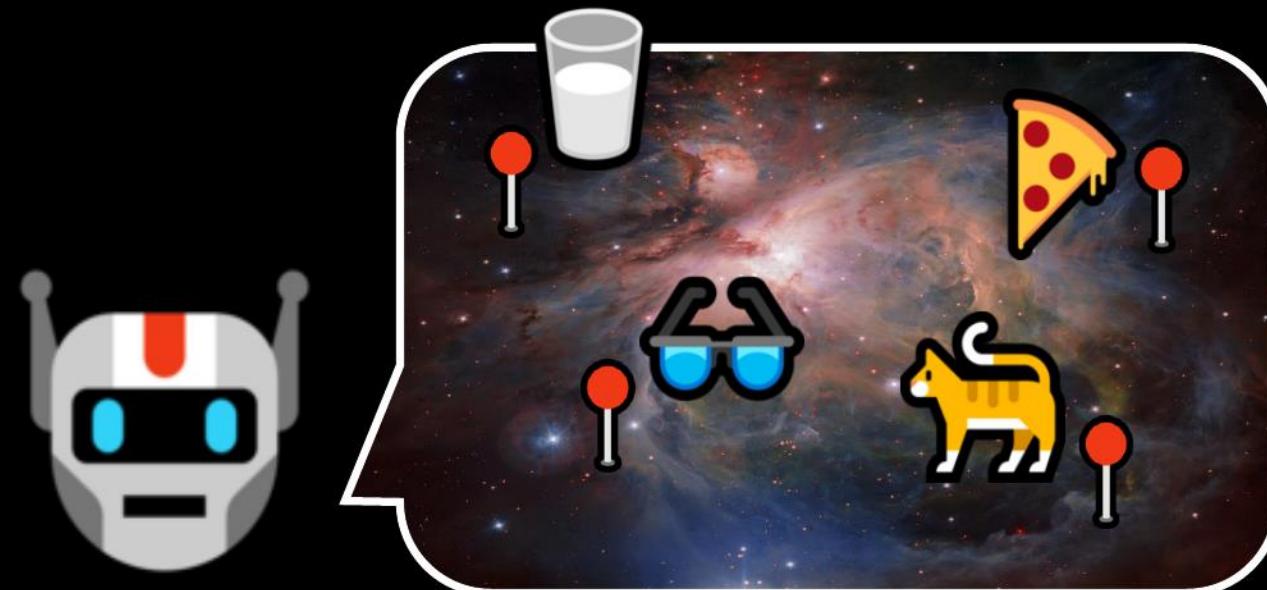


WSD Results

Uninformed Sense Matching (matching +200K)
Same standard but without filtering candidates by lemmas or POS



Applying Sense Embeddings



World Knowledge in NLMs



What's BERT thinking about when he reads?

World Knowledge in NLMs

[E1] played [E2] in [E3]

Marlon*	Brando*	played	Corleone*	in	Godfather*
<i>person¹_n</i>	<i>person¹_n</i>	<i>act³_v</i>	<i>syndicate¹_n</i>	<i>movie¹_n</i>	<i>location¹_n</i>
<i>womanizer¹_n</i>	<i>group¹_n</i>	<i>make⁴²_v</i>	<i>mafia¹_n</i>	<i>telefilm¹_n</i>	<i>here¹_n</i>
<i>bustle¹_n</i>	<i>location¹_n</i>	<i>emote¹_v</i>	<i>person¹_n</i>	<i>final_cut¹_n</i>	<i>there¹_n</i>

act³_v: play a role or part; *make⁴²_v*: represent fictiously, as in a play, or pretend to be or act like; *emote¹_v*: give expression or emotion to, in a stage or movie role.

Serena*	Williams	played	Kerber*	in	Wimbledon*
<i>person¹_n</i>	<i>professional_tennis¹_n</i>	<i>play¹_v</i>	<i>person¹_n</i>	<i>win¹_v</i>	<i>tournament¹_n</i>
<i>therefore¹_r</i>	<i>tennis¹_n</i>	<i>line_up⁶_v</i>	<i>group¹_n</i>	<i>romp³_v</i>	<i>world_cup¹_n</i>
<i>reef¹_n</i>	<i>singles¹_n</i>	<i>curl⁵_v</i>	<i>take_orders²_v</i>	<i>carry³⁸_v</i>	<i>elimination_tournament¹_n</i>

play¹_v: participate in games or sport; *line_up⁶_v*: take one's position before a kick-off; *curl⁵_v*: play the Scottish game of curling.

David	Bowie*	played	Warszawa*	in	Tokyo
<i>person¹_n</i>	<i>person¹_n</i>	<i>play¹⁴_v</i>	<i>poland¹_n</i>	<i>originate_in¹_n</i>	<i>tokyo¹_n</i>
<i>amati²_n</i>	<i>folk_song¹_n</i>	<i>play⁶_v</i>	<i>location¹_n</i>	<i>in¹_r</i>	<i>japan¹_n</i>
<i>guarnerius³_n</i>	<i>fado¹_n</i>	<i>riff²_v</i>	<i>here¹_n</i>	<i>take_the_field²_v</i>	<i>japanese¹_a</i>

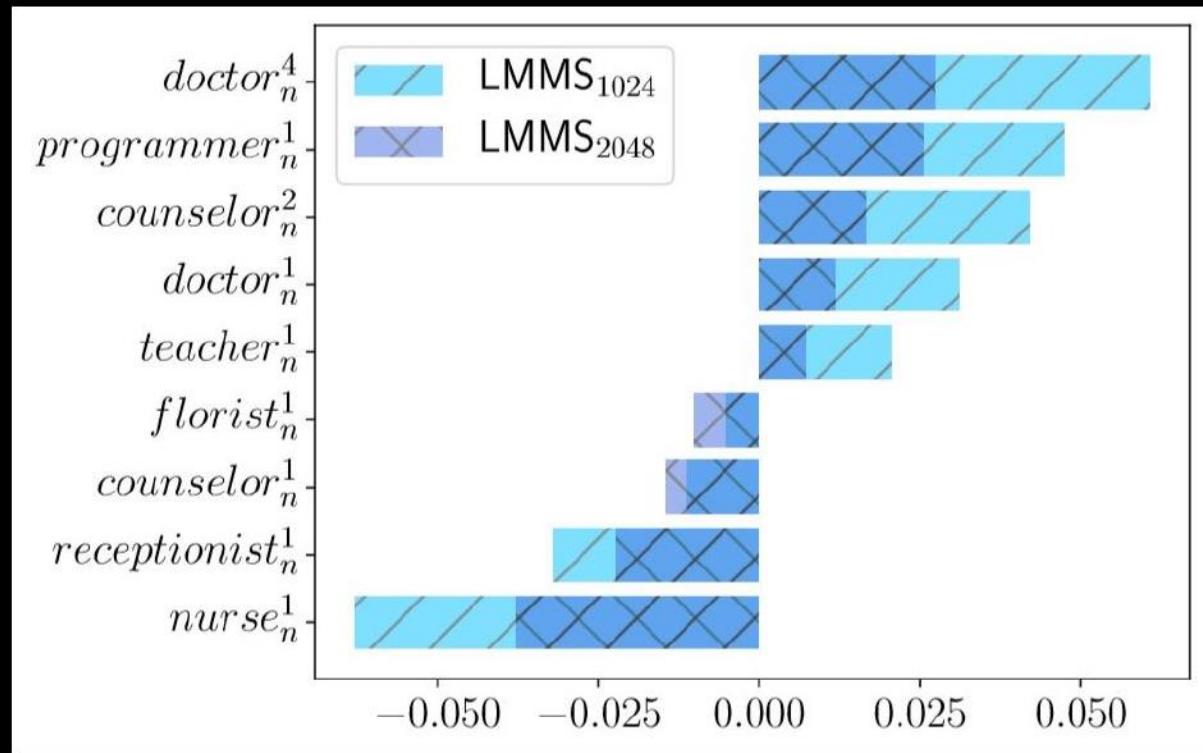
play¹⁴_v: perform on a certain location; *play⁶_v*: replay (as a melody); *riff²_v*: play riffs.

Checking for Biases in NLMs



Putting BERT on the spot

Checking for Biases in NLMs

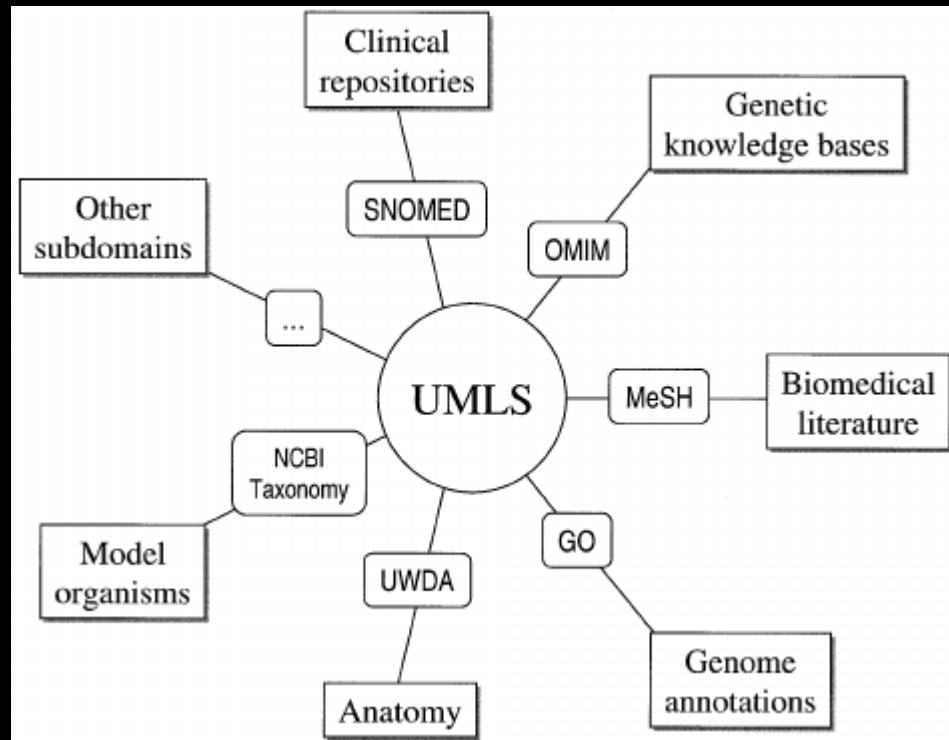


$$bias(s) = sim(\vec{v}_{man_n^1}, \vec{v}_s) - sim(\vec{v}_{woman_n^1}, \vec{v}_s)$$

MedLinker: WSD Hard Mode



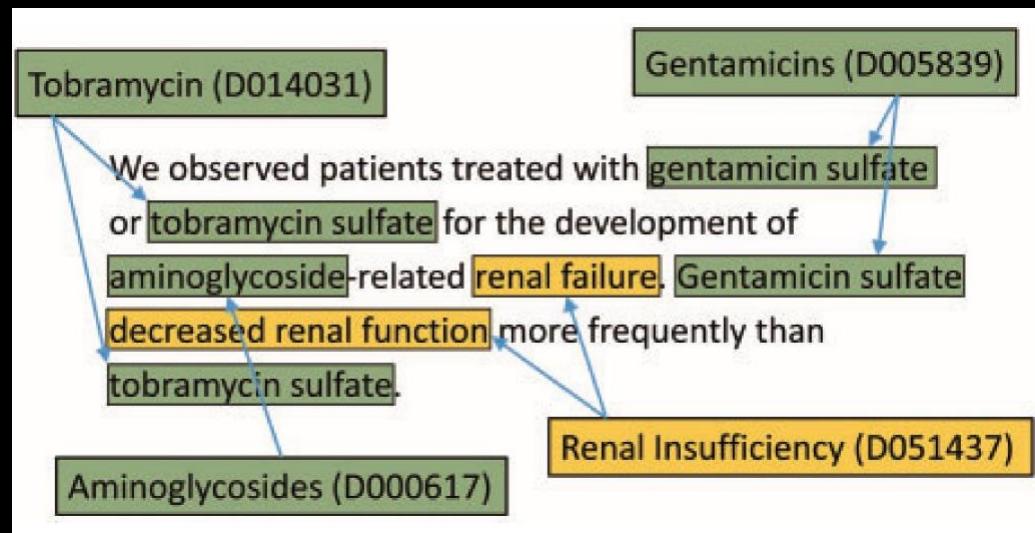
MedLinker: WSD Hard Mode



UMLS Ontology:

- +2M Medical Domain Concepts
- From professions to genes
- 10x size of WordNet

MedLinker: WSD Hard Mode



MedMentions (Mohan and Li, 2019)

- 200k annotations
- 18k concepts (entities)
- 50% train/test overlap
- 1% coverage

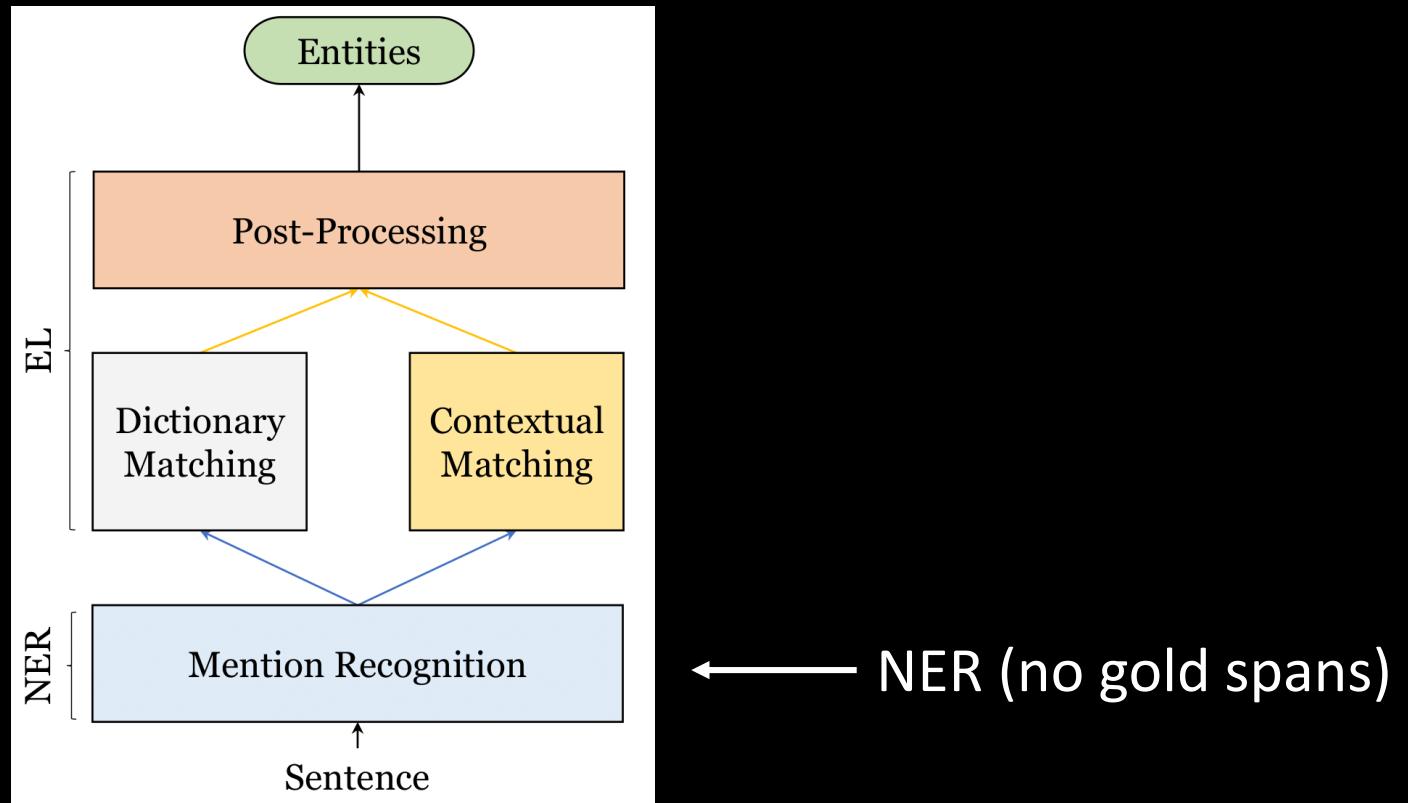
MedLinker: WSD Hard Mode



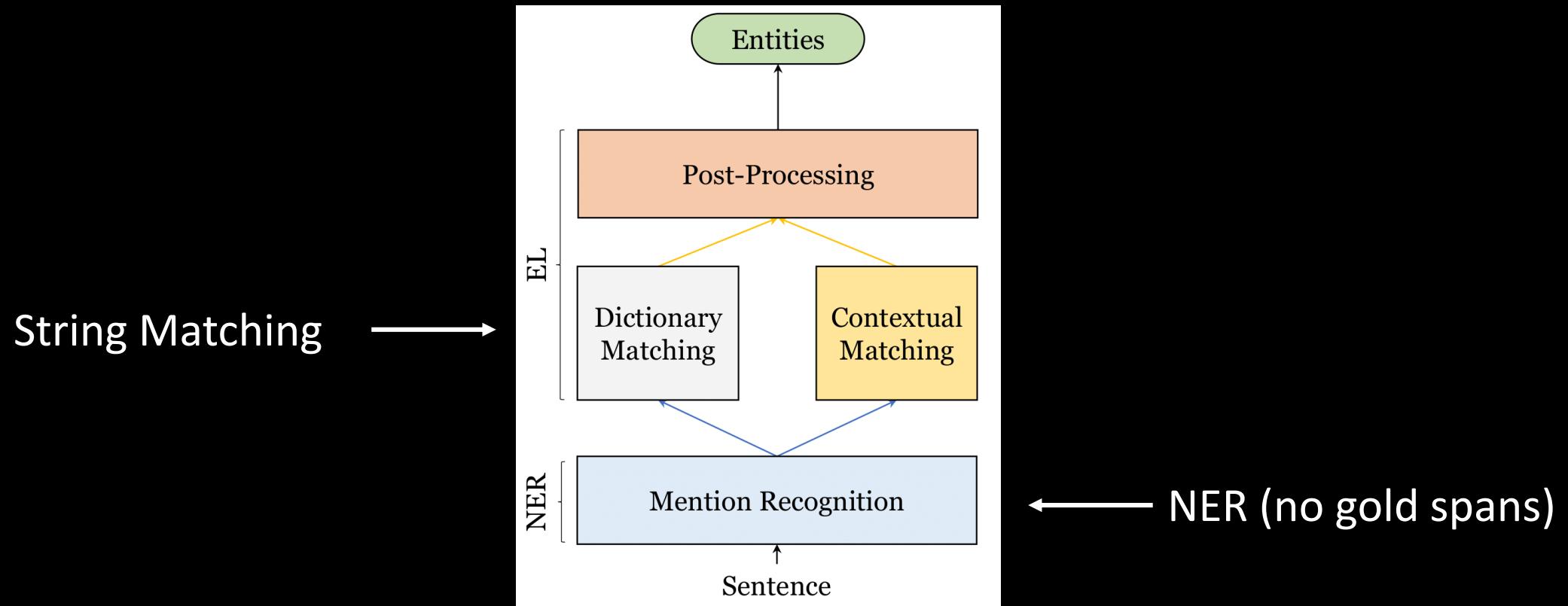
Specialized BERT models

- Trained on PubMed, etc.
- NCBI BERT
- BioBert
- SciBERT

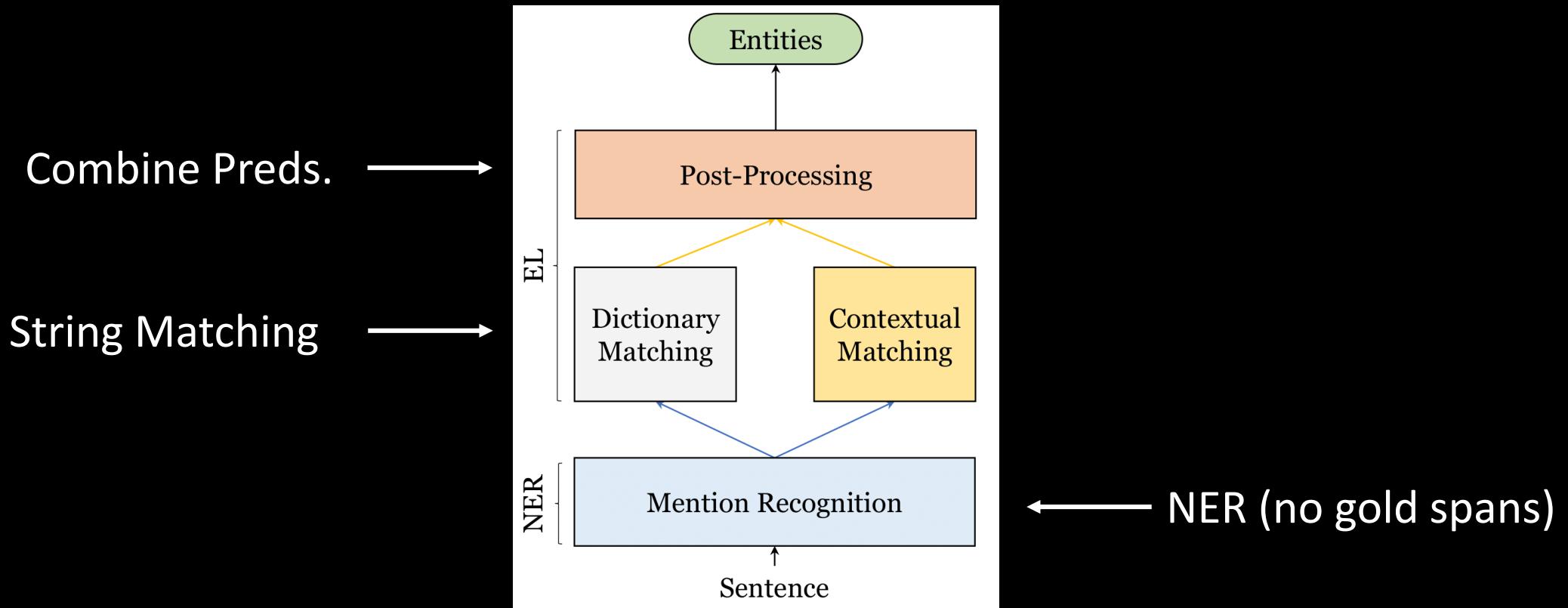
MedLinker: Solution



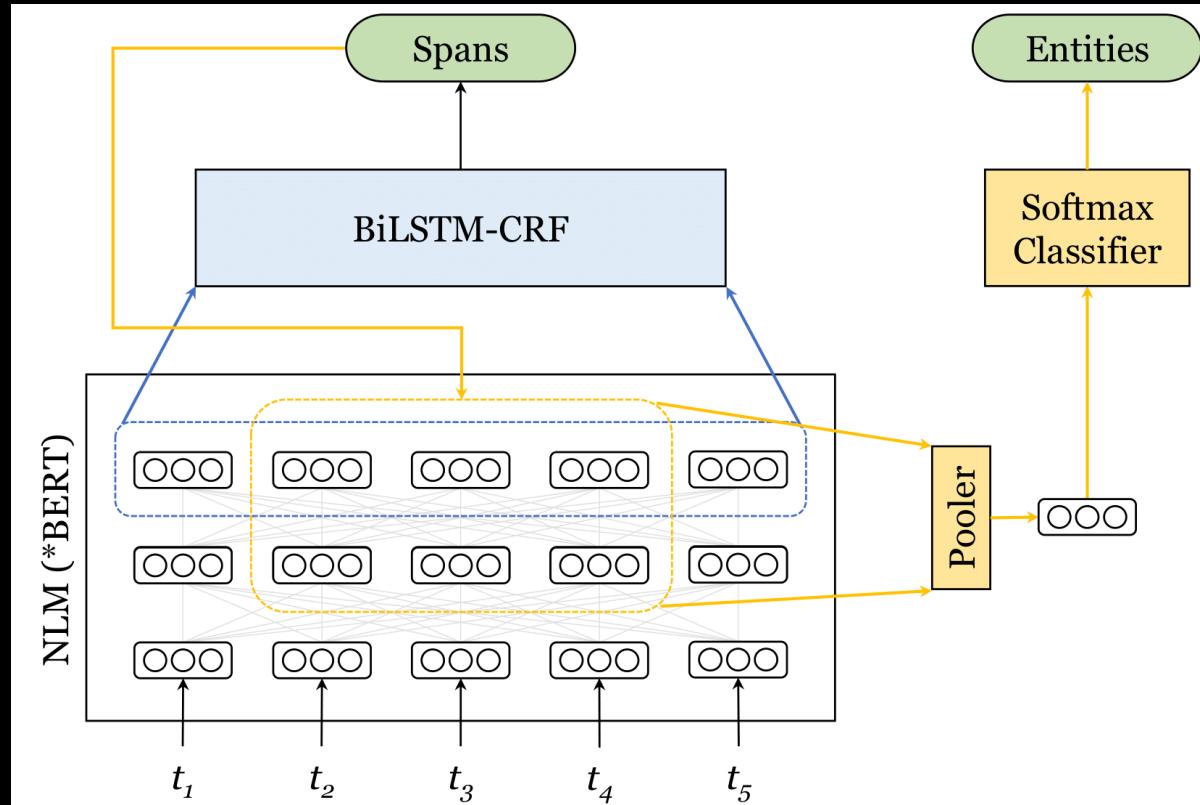
MedLinker: Solution



MedLinker: Solution



MedLinker: Solution



MedLinker: Solution

$$scoreSTR(s, e) = \max_{a \in map(e)} \cos(\hat{s}, \hat{a})$$

← All aliases(e), for all entities

$$scoreCLF(s, e) = P(e = j | \vec{s}) = \frac{\exp^{f(\vec{s})_j}}{\sum_{i=1}^{|E|} \exp^{f(\vec{s})_i}}$$

← Entities in train set

$$scoreSTR_CLF(s, e) = \max(scoreSTR(s, e), scoreCLF(s, e))$$

MedLinker: Results

Mention Recognition			
Model	P	R	F1
Exact Match	51.32	32.96	40.14
NCBI BERT (Uncased)	69.44	69.38	69.41
BioBERT 1.1 (Cased)	70.00	70.43	70.21
SciBERT (SciVocab)			
- Uncased	69.42	71.81	70.59
- Cased	69.16	71.30	70.22

MedLinker: Results

Concept (CUI) Linking			
Model	P	R	F1
Exact Match	47.12	31.11	37.48
TaggerOne [1]	47.10	43.60	45.30
MedLinker			
- scoreSTR	33.03	47.34	38.91
- score1NN	33.61	55.16	41.77
- scoreCLF	32.21	52.66	39.97
- scoreSTR_1NN	40.46	59.69	48.23
- scoreSTR_CLF	40.70	59.59	48.37
- scoreSTR_CLF (t=0.70)	48.43	50.07	49.24

Semantic Type (STY) Linking			
Model	P	R	F1
Exact Match	49.04	31.97	38.71
Fraser et al. 2019 [8]			
- BioBERT	61	66	63
- BioBERT BERT-base	63	65	64
MedLinker			
- scoreSTR	48.31	56.81	52.22
- score1NN	46.62	62.67	53.47
- scoreCLF	58.62	64.63	61.48
- scoreSTR_1NN	53.06	65.94	58.80
- scoreSTR_CLF	59.23	67.81	63.23
- scoreSTR_CLF (t=0.45)	63.13	63.69	63.41

Conclusion

- Recent NLMs allow for unprecedented gains in WSD.
- The same principles show similar improvements in Entity Linking tasks.
- 1NN (and w/softmax) methods are straightforward enough to assume progress should follow with better NLMs.
- Sense/entity embeddings can be useful for probing NLMs for world knowledge and downstream tasks.

Thanks



Code and Sense Embeddings:
github.com/danlou/LMMS



dloureiro@fc.up.pt



@danielbloureiro