

# From Relation Extraction to Knowledge Graphs

---

## Master Thesis Project – IC School

**Diego Antognini**

[diegoantognini@gmail.com](mailto:diegoantognini@gmail.com)

**Current PhD Student (LIA)**

**Dr. Jean-Cédric Chappelier**

**EPFL Supervisor**

**Bernard Maccari**

**Iprova Supervisor**

# Outline

---

## 1. Problem Statement

## 2. Our Method

- Relation Classification Models
- Building Knowledge Graphs

## 3. Conclusion & Future Work

# Our Objectives

1. Given a sentence, extract concepts and find the relationship among them if such one exists.  
**EEG measures brain activities**
  2. Given a corpus, build Knowledge Graphs of concepts, favoring precision over recall.
- What are concepts ?
  - What are the relationships ?
  - What are Knowledge Graphs ? Why use them ?  
⇒ Come back later !

# Concepts

---

- Short phrases made of adjectives and nouns
  - Gyroscope ✓
  - Rotational motion ✓
  - Brain electrical activity ✓
  - The new model S developed by Tesla ✗
  - Galaxy S8 of Samsung ✗

# Relations

Directed relations Iprova is interested in

Relation	Example
Cause	Those cancers were caused by radiation exposures.
Contain	My apartment has a large kitchen.
Measure	EEG measures brain activities.
Produce	A factory manufactures suits.
TypeOf	NoSQL databases such as MongoDB.
Use	Bluetooth is used in audio equipment.
Other	A misty ridge uprises from the surge.

# Relations

Directed relations Iprova is interested in

Relation	Example
Cause	Those <b>cancers</b> were caused by <b>radiation exposures</b> .
Contain	My <b>apartment</b> has a <b>large kitchen</b> .
Measure	<b>EEG</b> measures <b>brain activities</b> .
Produce	A <b>factory</b> manufactures <b>suits</b> .
TypeOf	NoSQL databases such as <b>MongoDB</b> .
Use	<b>Bluetooth</b> is used in <b>audio equipment</b> .
Other	A <b>misty ridge</b> uprises from the <b>surge</b> .

# 2 First Approaches

## 1. Relation Extraction

- Named entities (Location, Organization, Person, etc.)
- Specified relations (e.g. CoFounder, BornIn)
- Need a lot of data

## 2. Open Information Extraction

- Named entities/nominals (nouns/base noun phrases)
- No specified relations
  - ⇒ find a mapping to “ontology”

e.g. was included in ⇒  $\text{Contain}(e2, e1)$

# Relation/Open Information Extraction

**Bill Gates, Microsoft co-founder**, stepped down as **CEO** in January 2000. **Gates** was included in the **Forbes** **wealthiest** list since 1987 and **was the wealthiest** from 1995 to 2007...

It was announced that **IBM** would buy **Ciao** for an undisclosed amount. The **CEO**, **MacLorrance** has occupied the **corner office of the Hopkinton**, company

The company's storage business is also threatened by new, born-on-the Web could providers like Dropbox and Box, and ...

**RE**  
→

Co-founder(Bill Gates, Microsoft)  
Director-of (MacLorrance, Ciao)  
Employee-of (MacLorrance, Ciao)  
...

**Open IE**  
→

(Bill Gate, be, Microsoft co-founder)  
(Bill Gates, stepped down as, CEO)  
(Bill Gates, was included in, the Forbes wealthiest list)  
(Bill Gates, was, the wealthiest)  
(IBM, would buy, Ciao)  
(MacLorrance, has occupied, the corner office of the Hopkinton)  
...

<b>RE</b>		<b>Open IE</b>
<b>Input</b>	Sentences + Labeled relations	Sentences
<b>Relation</b>	Specified relations in advance	Free discovery
<b>Extractor</b>	Specified relations	Independent-relations

Image from Vo and Bagheri, 2016

# Chosen approach

---

## 3. Relation Classification

- Specified relations
- Named entities/Nominals/Concepts
  - These are given with the sentence
  - How to find concepts ?  
⇒ Using existing concept extraction system  
 $(ADJ)^* (NOUN)^+$

# **Our Method**

**Relation Classification  
Models**

**Building Knowledge  
Graphs**

# Relation Classification

## Input:

“ The [factory]<sub>e1</sub>'s products have included flower pots, Finnish rooster-whistles, pans, [trays]<sub>e2</sub>, tea pots, ash trays and air moisturisers. ”

## Output:

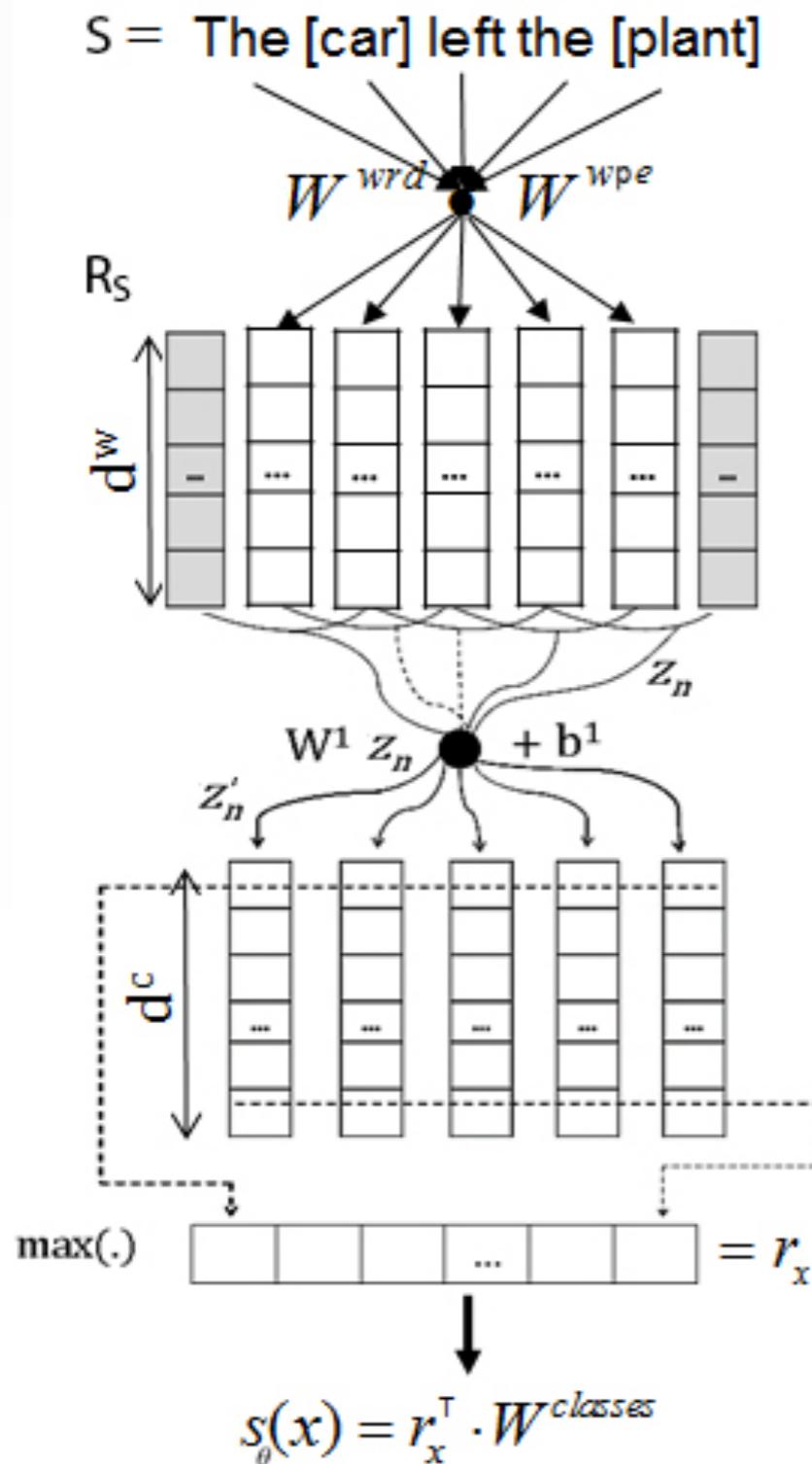
- The directed relationship among **factory** and **trays**  
⇒ **Produce(factory, trays)**

# Developed models

---

- 2 Models for **Relation Classification Task**:
  - **CR-CNN**: Convolutional neural networks
  - **BRCNN**: Recurrent & Convolutional neural networks
- These have been shown to be efficient architecture for RC

# Model 1: CR-CNN



- State of the art for 2015
- Convolutional Neural Network
- Simple features: word embeddings and relative distance
- Omit “Other” class embeddings
- Pairwise ranking loss function

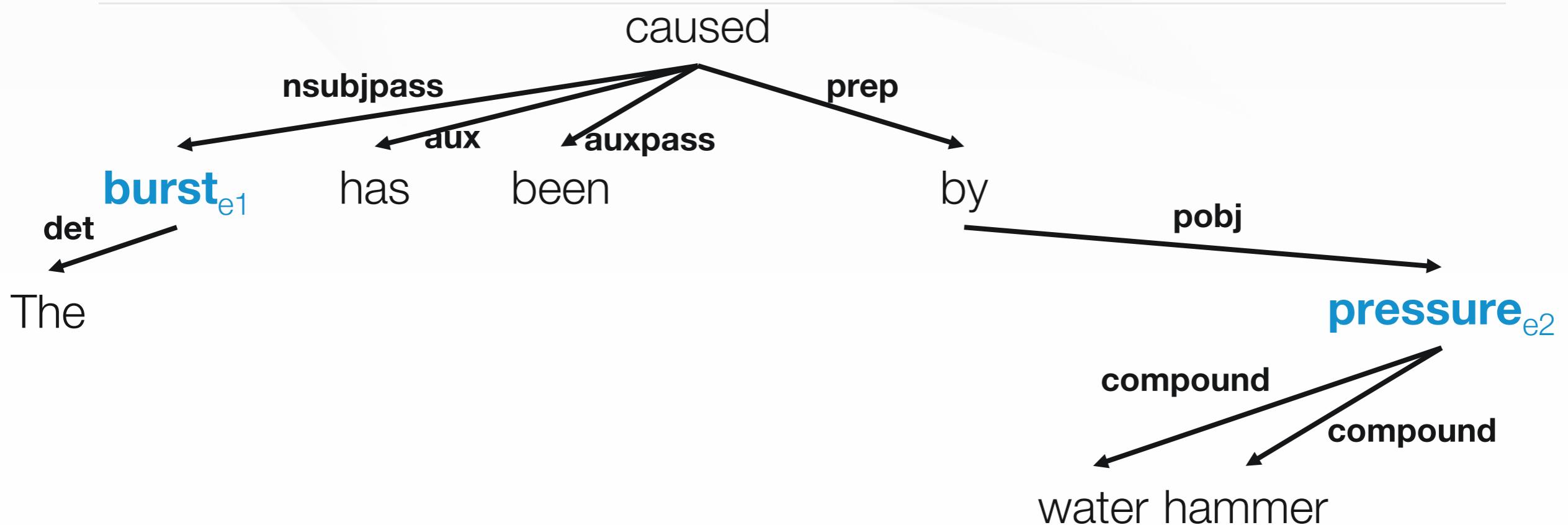
$$L = \log(1 + \exp(\gamma(m^+ - s_\theta(x)_{y^+}))) + \log(1 + \exp(\gamma(m^- + s_\theta(x)_{c^-})))$$

# Model 2: BR-CNN

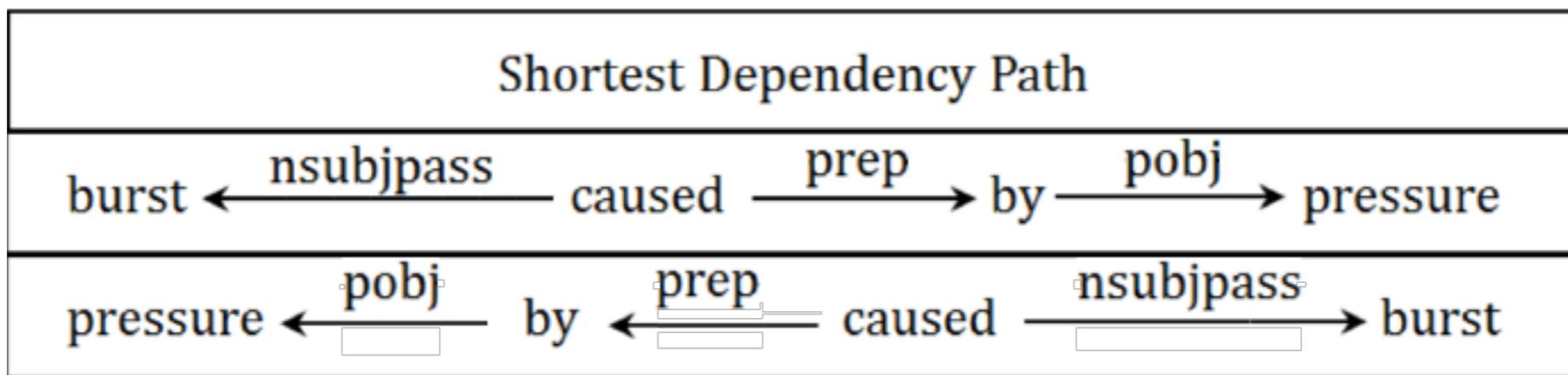
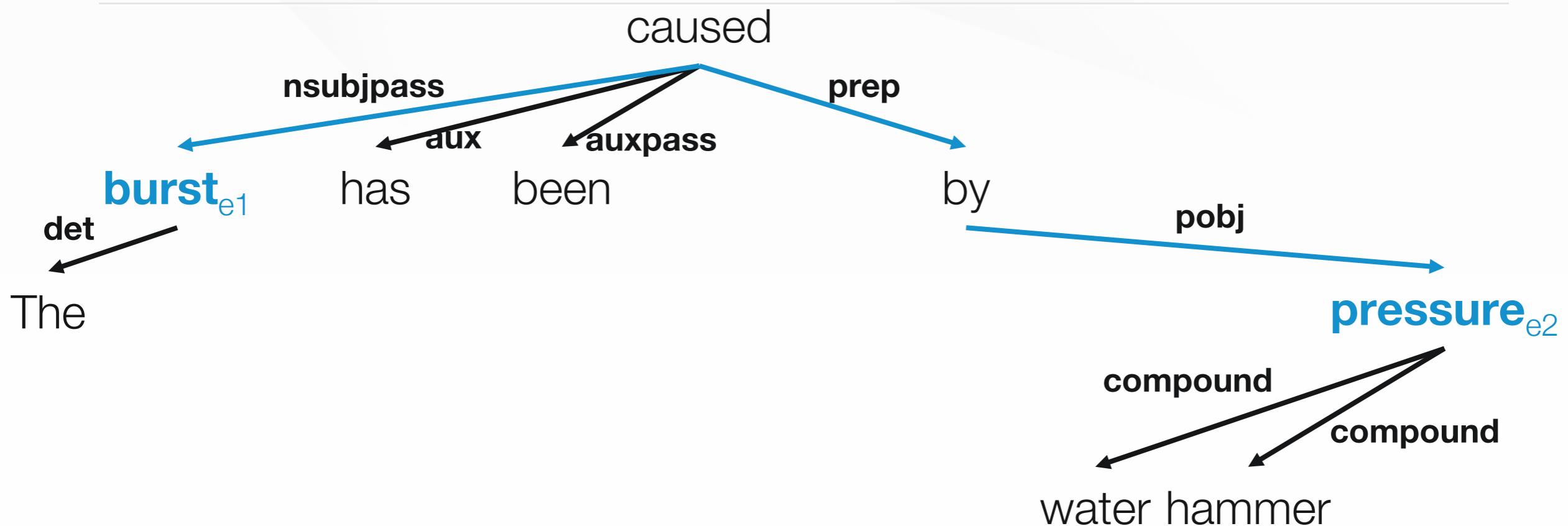
- State of the art since 2016
  - Bi-Recurrent Convolutional Neural Network
  - Use shortest dependency path
- BRCNN<sub>1</sub>** • Word embeddings, dependency tag embeddings
- BRCNN<sub>2</sub>** + POS tags, NER tags and WordNet hypernyms
- Training objective: cross entropy

$$J(x_i) = - \sum_{i=1}^{2K+1} \vec{t}_i \log \vec{y}_i - \sum_{i=1}^{2K+1} \overleftarrow{t}_i \log \overleftarrow{y}_i - \sum_{i=1}^{K+1} t_i \log y_i$$

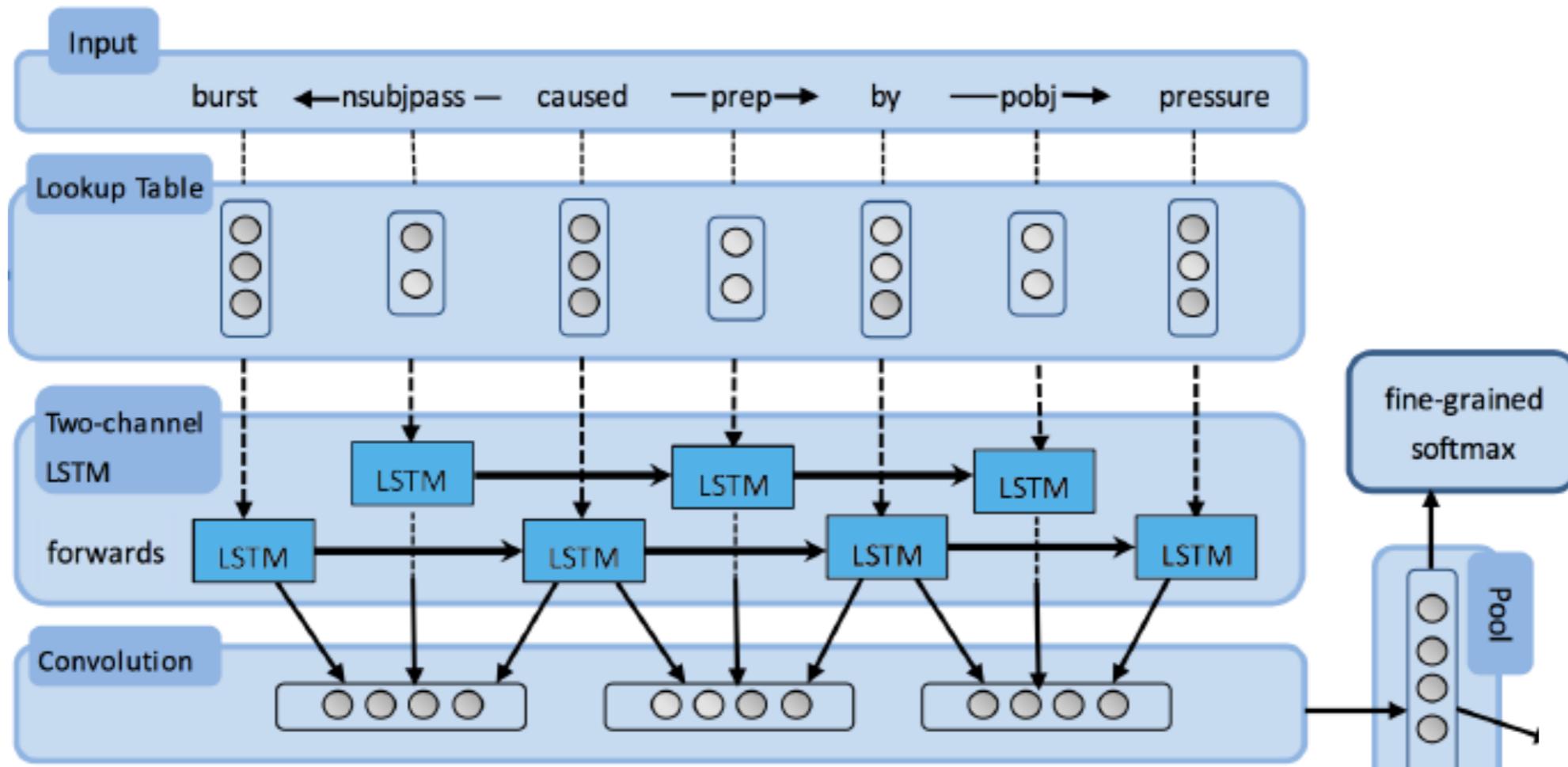
# Shortest Dependency Path



# Shortest Dependency Path



# Model 2: BR-CNN



# Experiments

- Run on the 3 datasets 5 times (mean + stdev)
- Compute macro  $F_1$ -Score excluding class Other
- Tune on validation set and evaluate on test set
- Comparison with 3 baselines
  - **UTD**: Support Vector Machine with lexical features
  - **SPTree**: bi-Recurrent Neural Networks
  - **DRNN**: deep Recurrent Neural Networks

# Datasets

---

## 1. SemEval-2010 Task 8

- Established benchmark for Relation Classification
- Most of the sentences are either short or average
- **2x9 relations + 1 Other  $\Rightarrow$  19 relations**

## 2. KBP37

- Named entities
- Longer sentences
- **2x18 relations + 1 Other  $\Rightarrow$  37 relations**

# Datasets

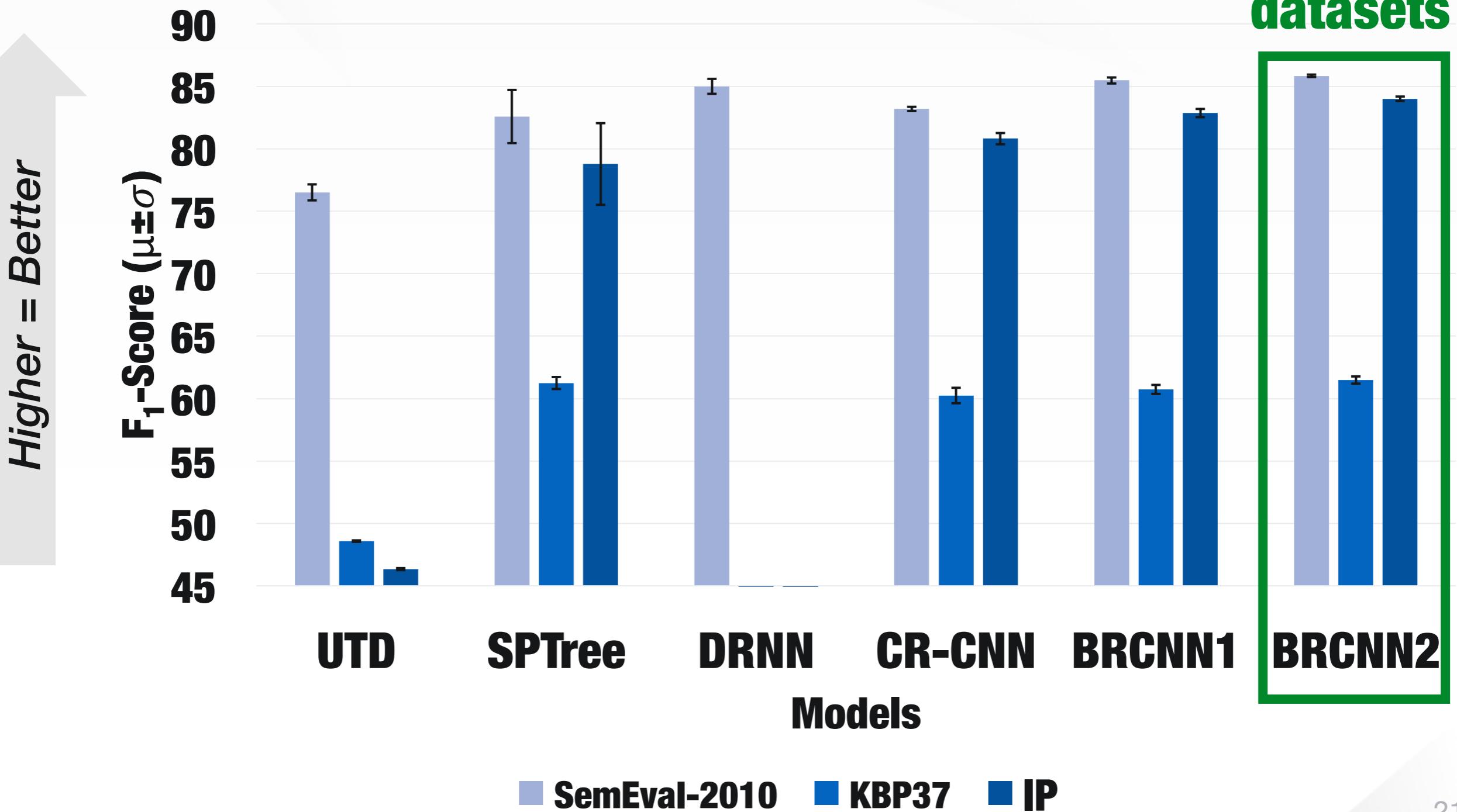
---

## 3. IP

- Training set partially based on SemEval-2007/2010
- Manually gathered sentences from various websites & manual searches on the Internet
- Most of the sentences are either short or average
- Relations of interest for Iprova  
**2x6 relations + 1 Other ⇒ 13 relations**

# Results

Best on all  
datasets



# Improvements

Gum disease rates were highest in [males]<sub>e<sub>1</sub></sub>, Mexican Americans, adults with less than a high school education, adults below the [poverty line]<sub>e<sub>2</sub></sub> and current smokers.

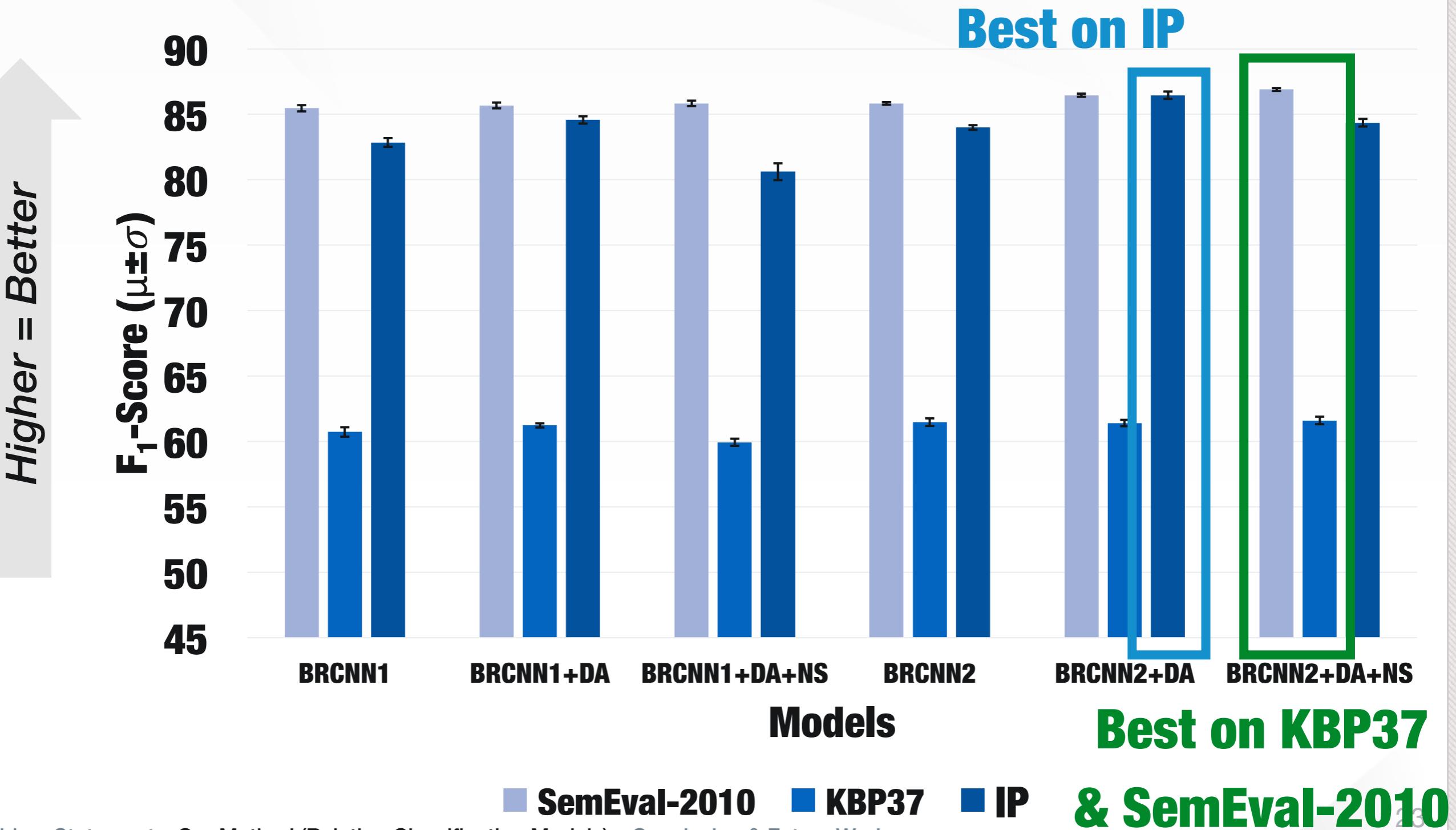
- **Data augmentation**
  - Replace some words with neighbors in Word2Vec space

Gum infection rates were highest in [males]<sub>e<sub>1</sub></sub>, Peruvian Americans, adults with less than a junior-senior school education, adults below the [poverty line]<sub>e<sub>2</sub></sub> and former nonsmokers.

- **Negative Sampling**
  - Assign tags  $\square_{e_1}$ ,  $\square_{e_2}$  to other words

Gum disease rates were highest in males, Mexican Americans, adults with less than a high school education, [adults]<sub>e<sub>1</sub></sub> below the poverty line and current [smokers]<sub>e<sub>2</sub></sub>.

# Results



# Recap

---

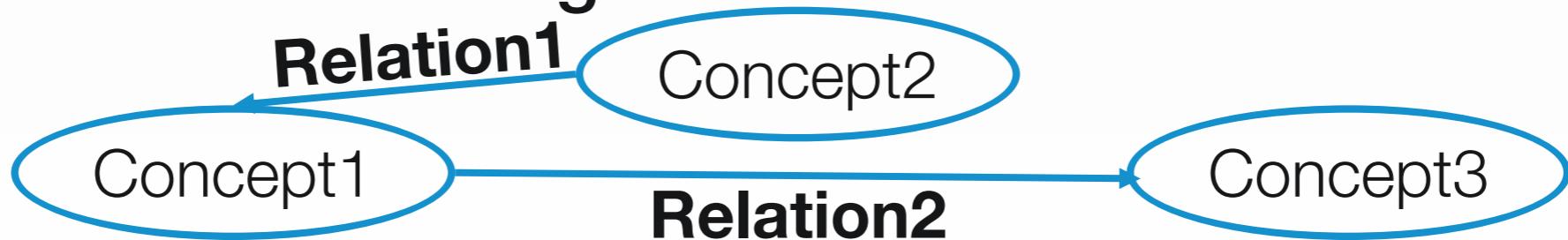
- Incorporating linguistic information in network's architecture is still important and beneficial
- Data Augmentation & Negative Sampling techniques help to strengthen classifiers
- BRCNN<sub>2</sub>+DA+NS outperforms all models on Sem & KBP
- BRCNN<sub>2</sub>+DA outperforms all models on IP dataset  
⇒ **Will be used to build Knowledge Graphs**

# **Our Method**

**Relation Classification  
Models  
Building Knowledge  
Graphs**

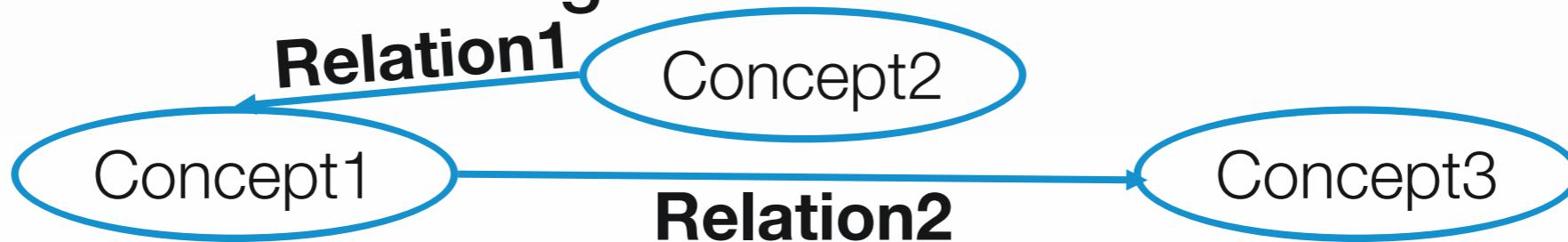
# What, Why, How ?

- **What:** structured representation of semantic knowledge and relations among nodes



# What, Why, How ?

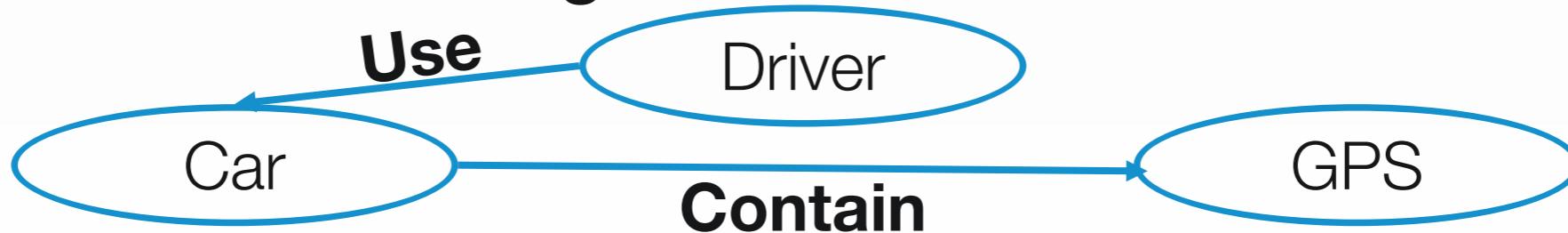
- **What:** structured representation of semantic knowledge and relations among nodes



- **Why:** model domains of interest, infer new relations, basis for a Question-Answering system, etc.

# What, Why, How ?

- **What:** structured representation of semantic knowledge and relations among nodes



- **Why:** model domains of interest, infer new relations, basis for a Question-Answering system, etc.

# What, Why, How ?

- **What:** structured representation of semantic knowledge and relations among nodes



- **Why:** model domains of interest, **infer new relations**, basis for a Question-Answering system, etc.

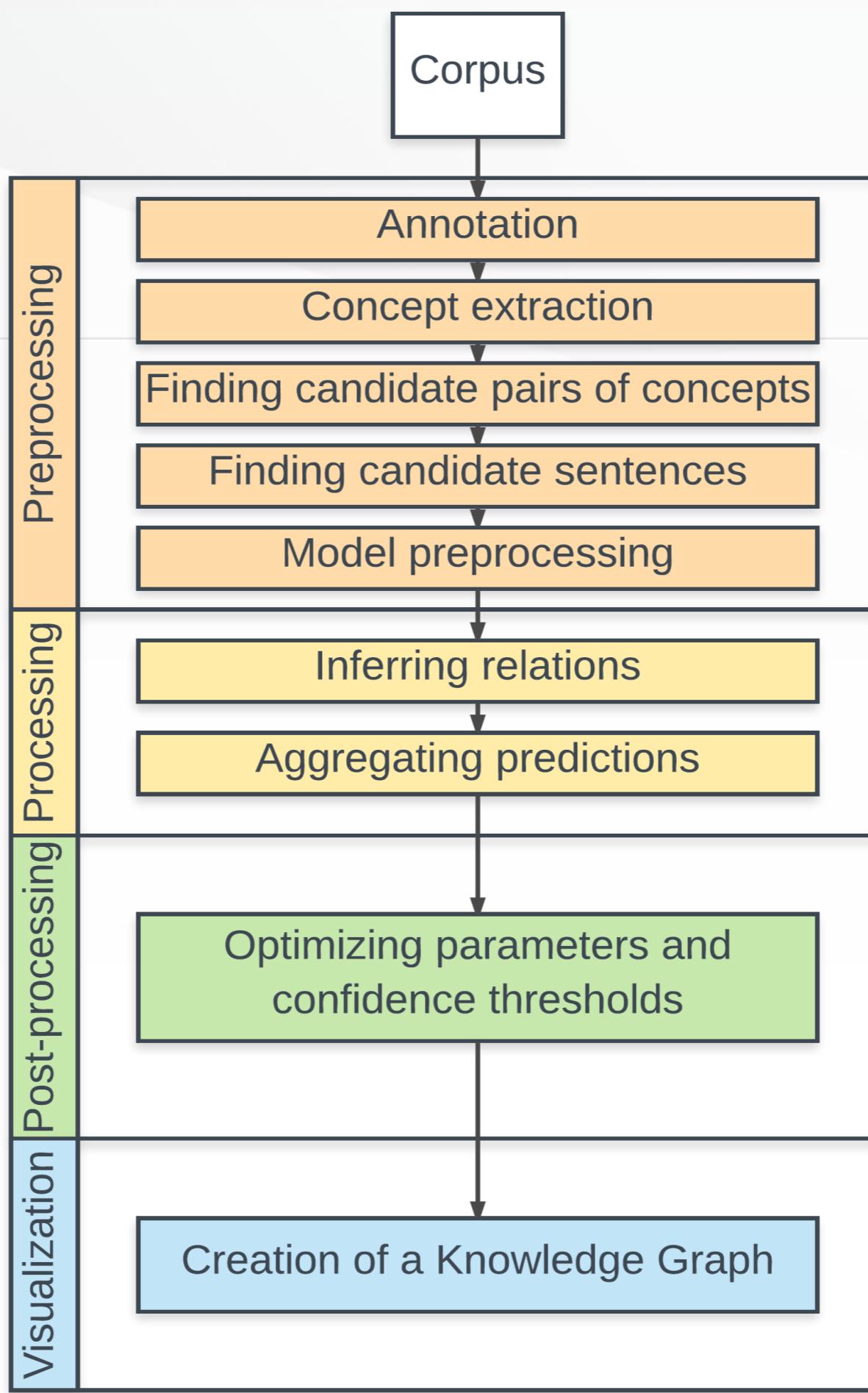
# What, Why, How ?

- **What:** structured representation of semantic knowledge and relations among nodes

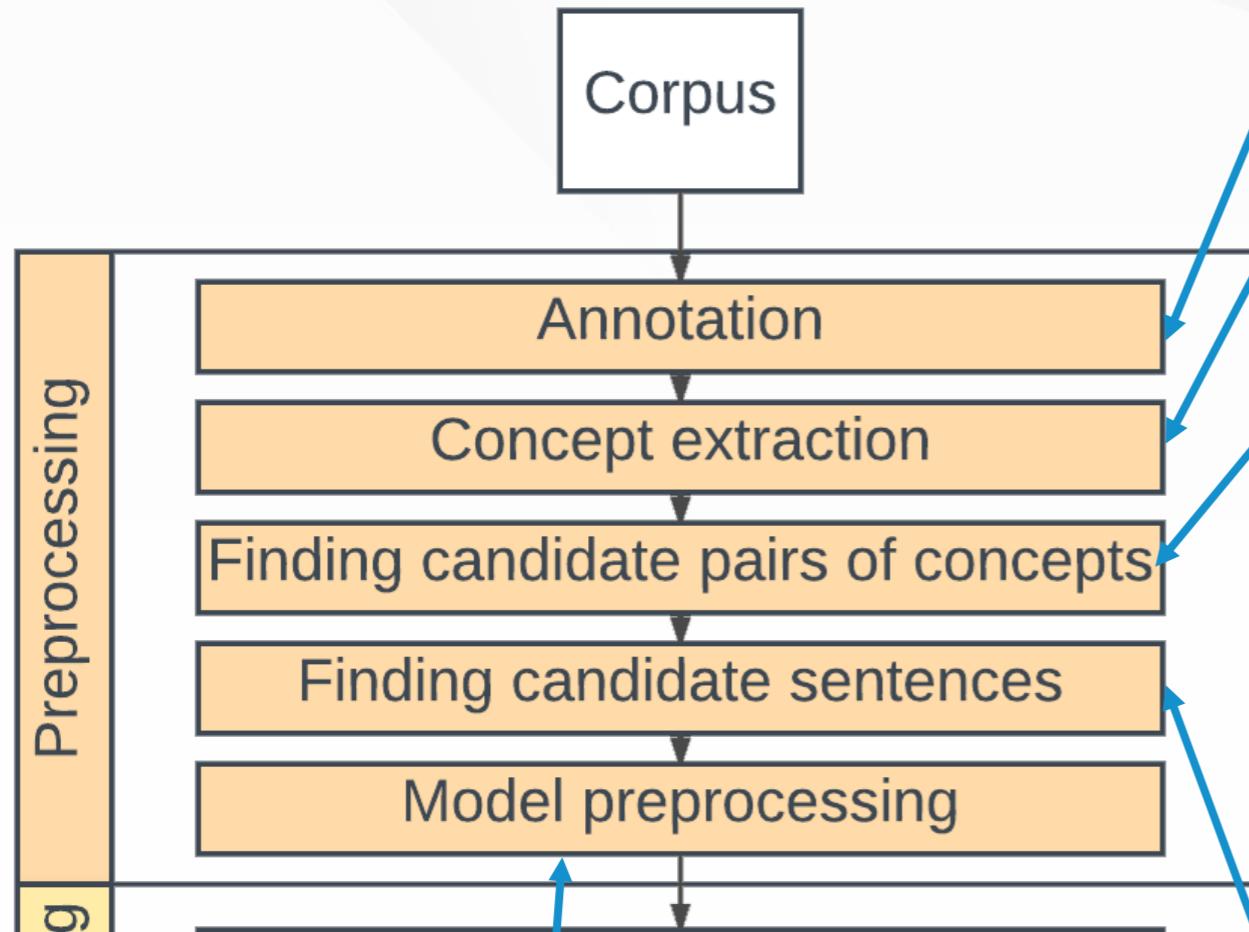


- **Why:** model domains of interest, infer new relations, basis for a Question-Answering system, etc.
- **How:**
  - Extracting pairs of concepts from large corpora
  - Infer relations with best model on IP dataset:**BRCNN<sub>2</sub>+DA**

# Pipeline



# Preprocessing



Preprocessing of  
**BRCNN<sub>2</sub>+DA** model

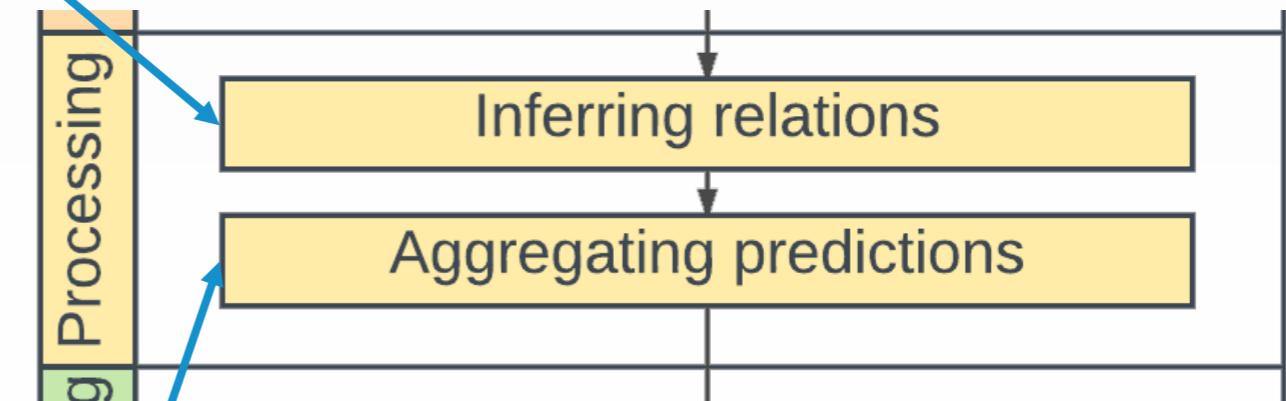
Stanford's tools

Iprova's concept extractor

- At least 2 sentences
- Concepts not too far away
- Containing both concepts  
A concept might be part of  
bigger concept e.g.  
**diabetes** ∈ type 2 diabetes  
**schedule** ∈ rotating schedule

# Processing

Using model **BRCNN2+DA**



What if **R1(Concept1, Concept2)** & **R2(Concept1, Concept2)** ?  
⇒ **Aggregate probability distribution vectors by class-label**

Pair	Conf. R1	Conf. R2	Conf. R3
(c1,c2)	0.5	0.3	0.2
(c1,c2)	0.8	0.2	0.0
(c1,c2)	0.2	0.1	0.7

Median



Pair	Conf. R1	Conf. R2	Conf. R3
(c1,c2)	0.5	0.2	0.2

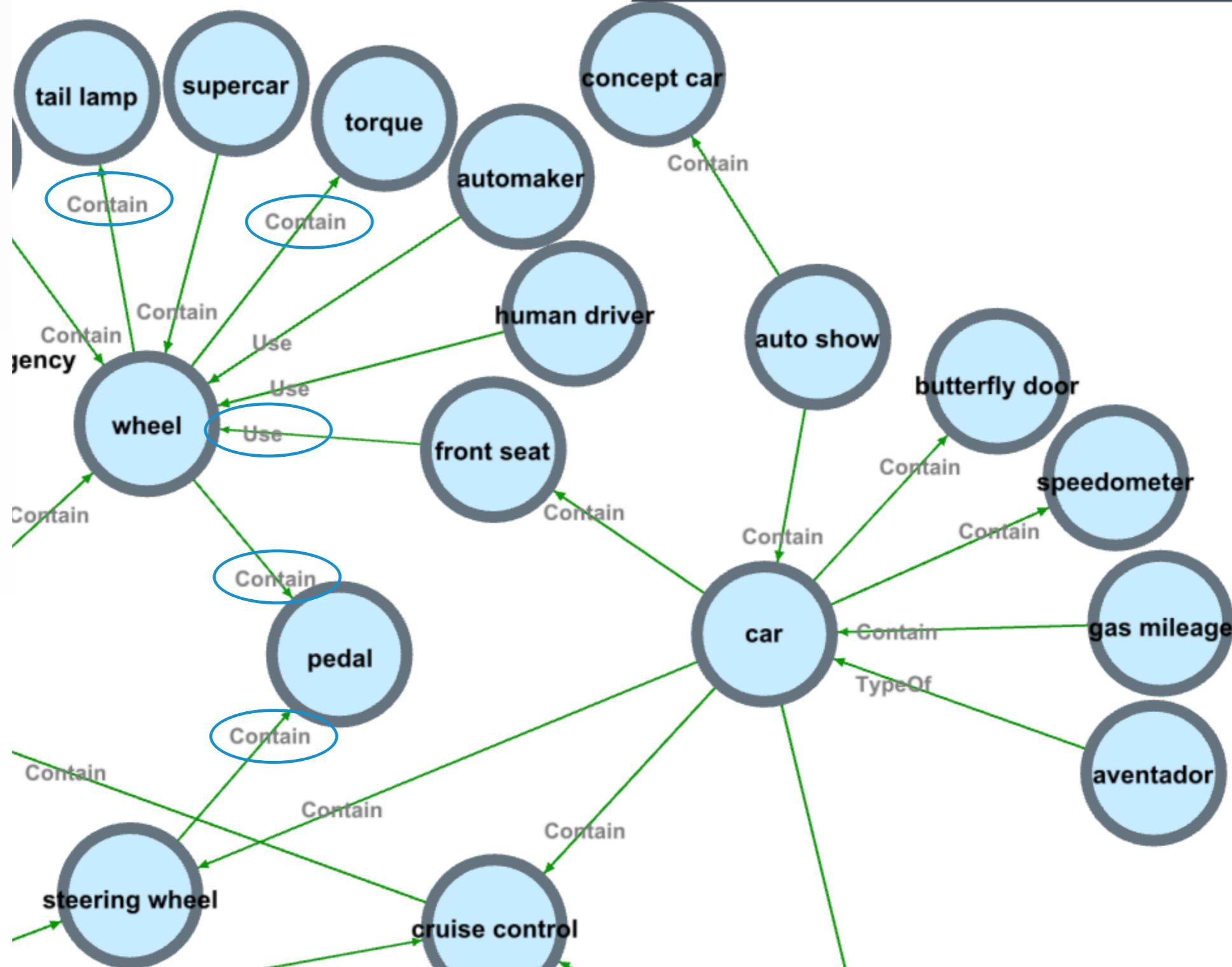
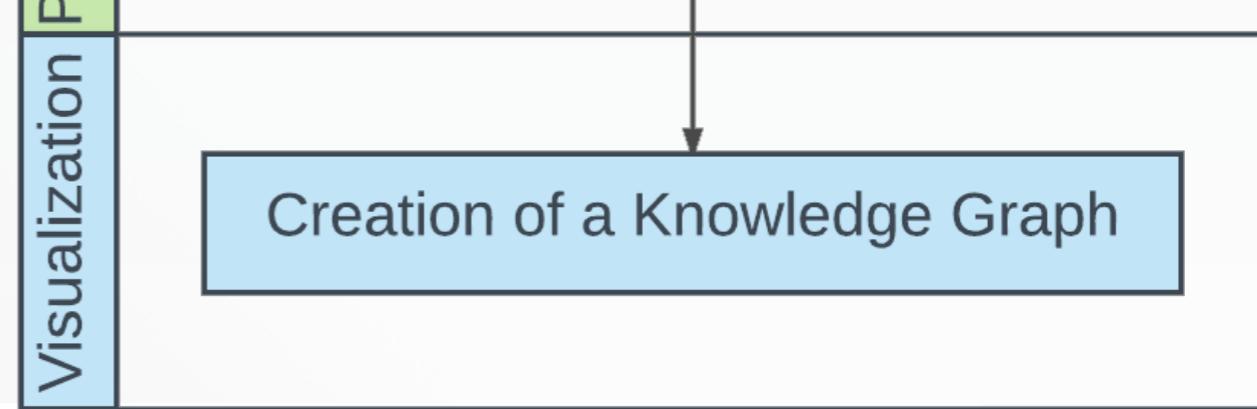
# Post-Processing

- **Goal:** filter out noise
- **Parameters:** free to setup during visualization by Iprova
- **Confidence Thresholds**  
Confidence threshold for each class
  - $\geq$  threshold  $\Rightarrow$  keep relation
  - $<$  threshold  $\Rightarrow$  Other

Post-processing

Optimizing parameters and confidence thresholds

# Visualization



# Qualtitative Evaluation

---

- Build representative Knowledge Graphs from 3 corpora with and without confidence thresholds (CT)
- **Manually assess** quality of the predictions on 2% of each KG
- Classify each sample in one of the four classes:
  1. **Makes sense**, e.g. Contain(car, wheels)
  2. **Reversed direction**, e.g. Contain(wheels, car)
  3. **Might make sense**, e.g. Use(racing, drivers)
  4. **Nonsense**, e.g. TypeOf(neck, tail)

# Corpora

---

## 1. Common Crawl

- based on ScienceDaily and Phys.org
- ~20 millions sentences

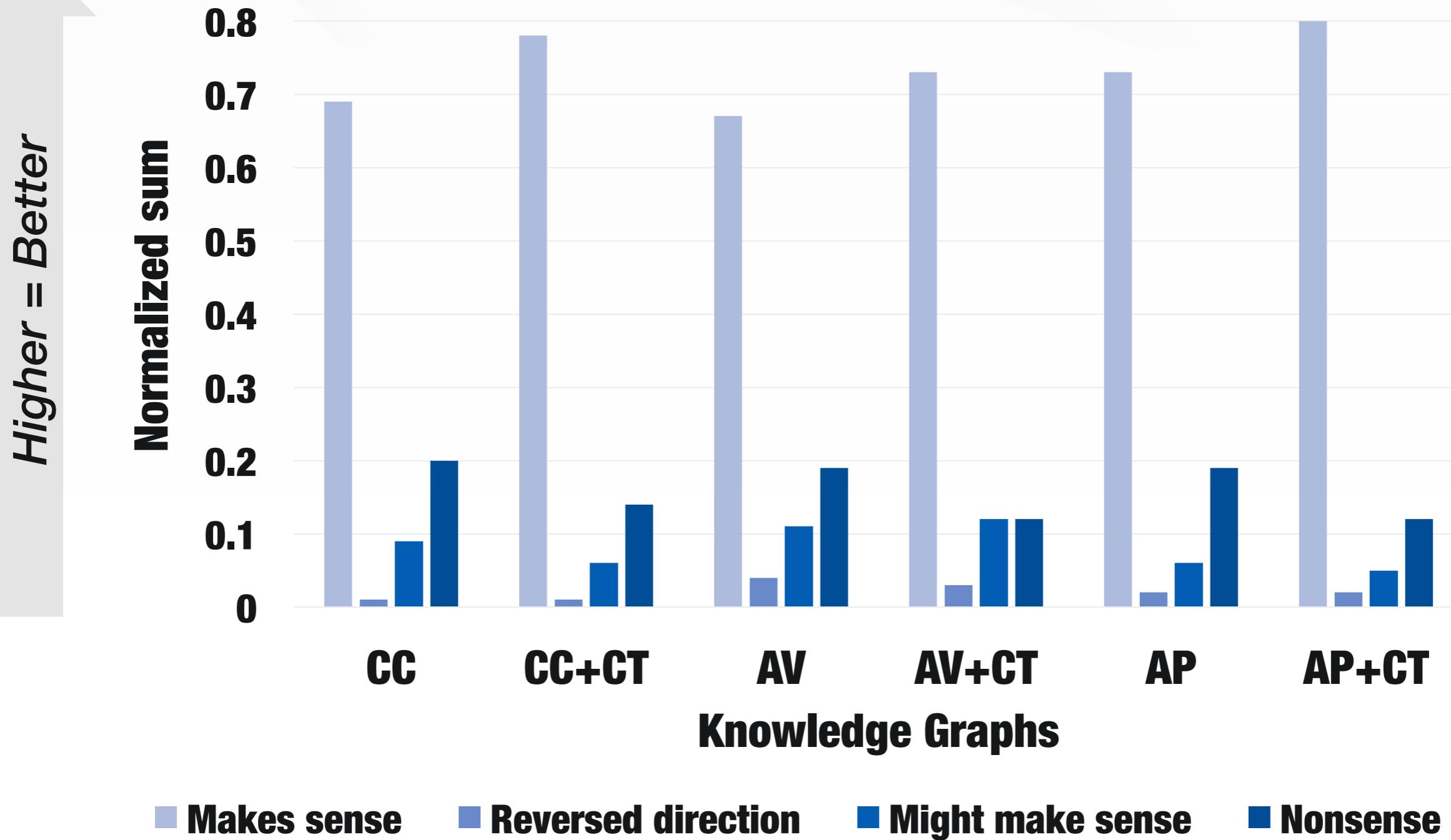
## 2. Autonomous Vehicles documents

- ~  $\frac{1}{2}$  million sentences

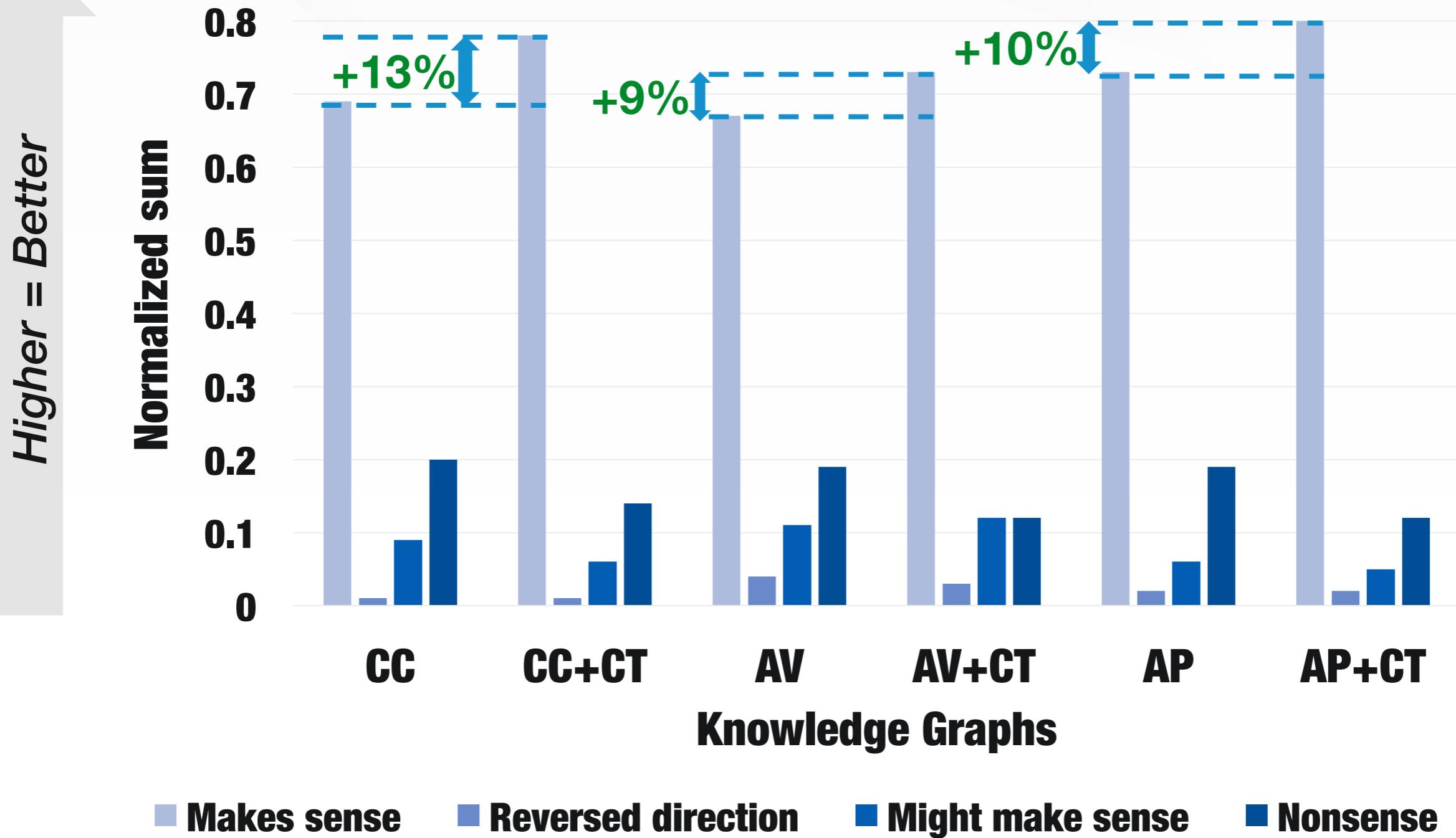
## 3. Air Purifier documents

- ~1 million sentences

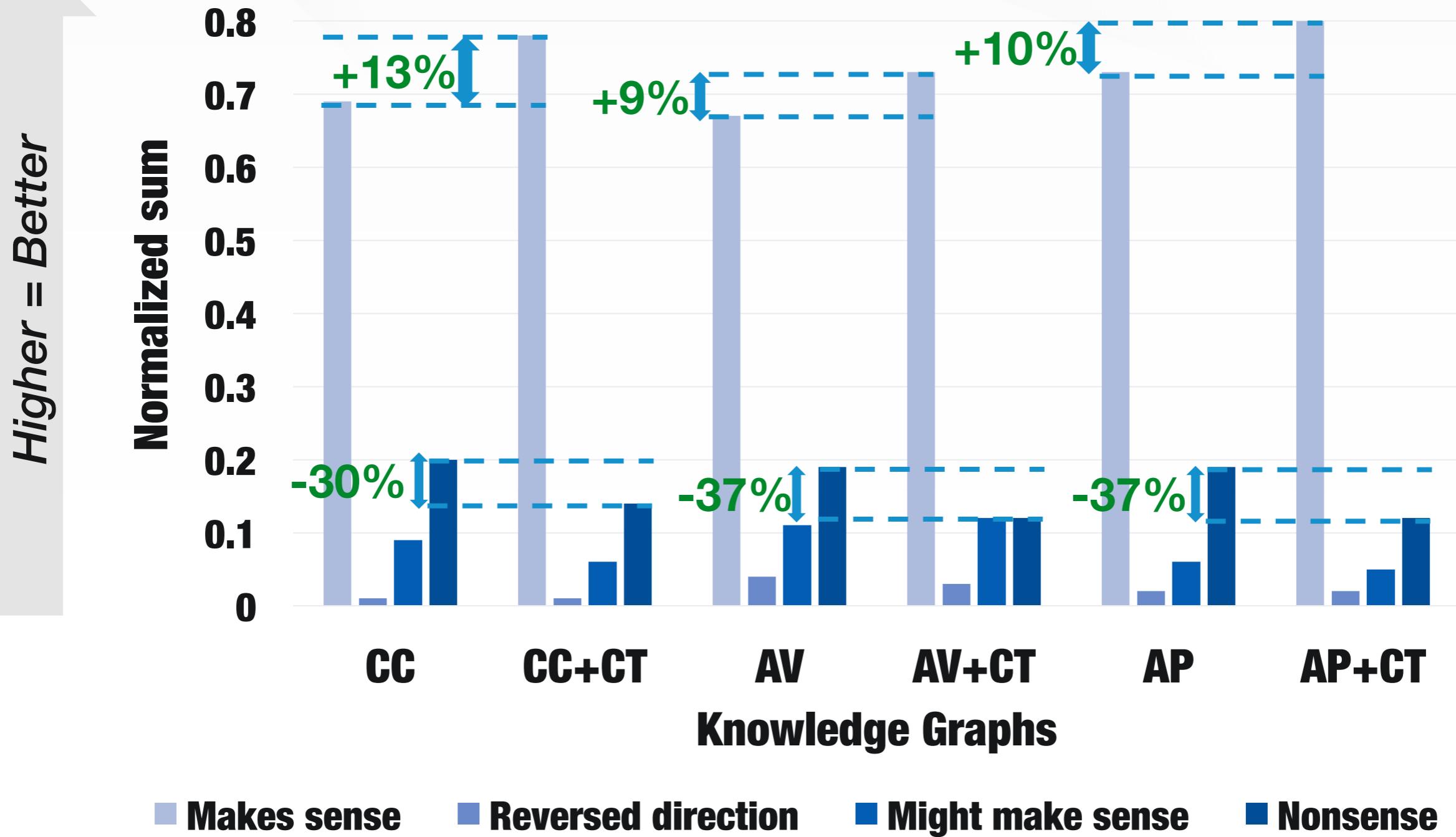
# Results



# Results

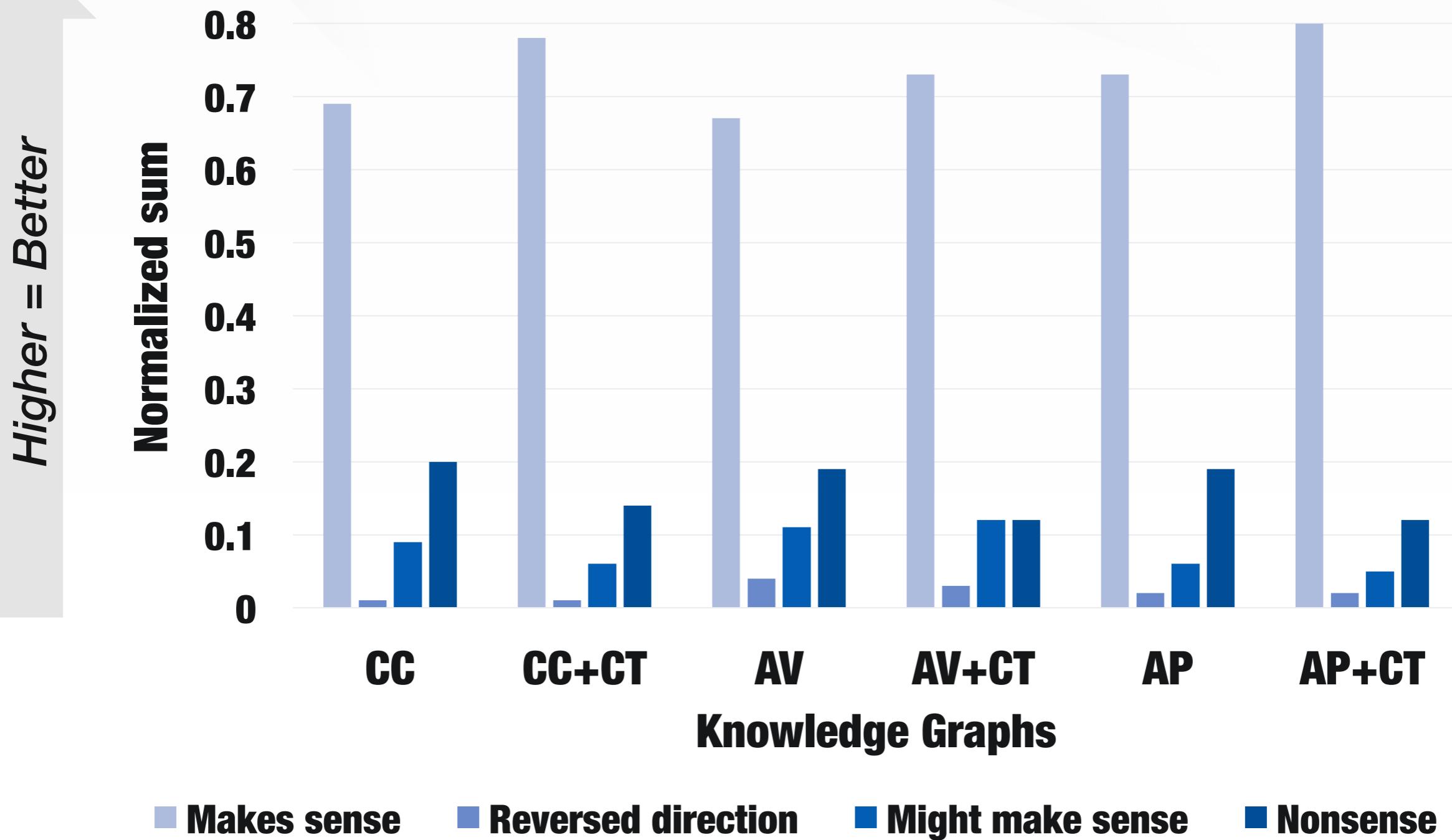


# Results



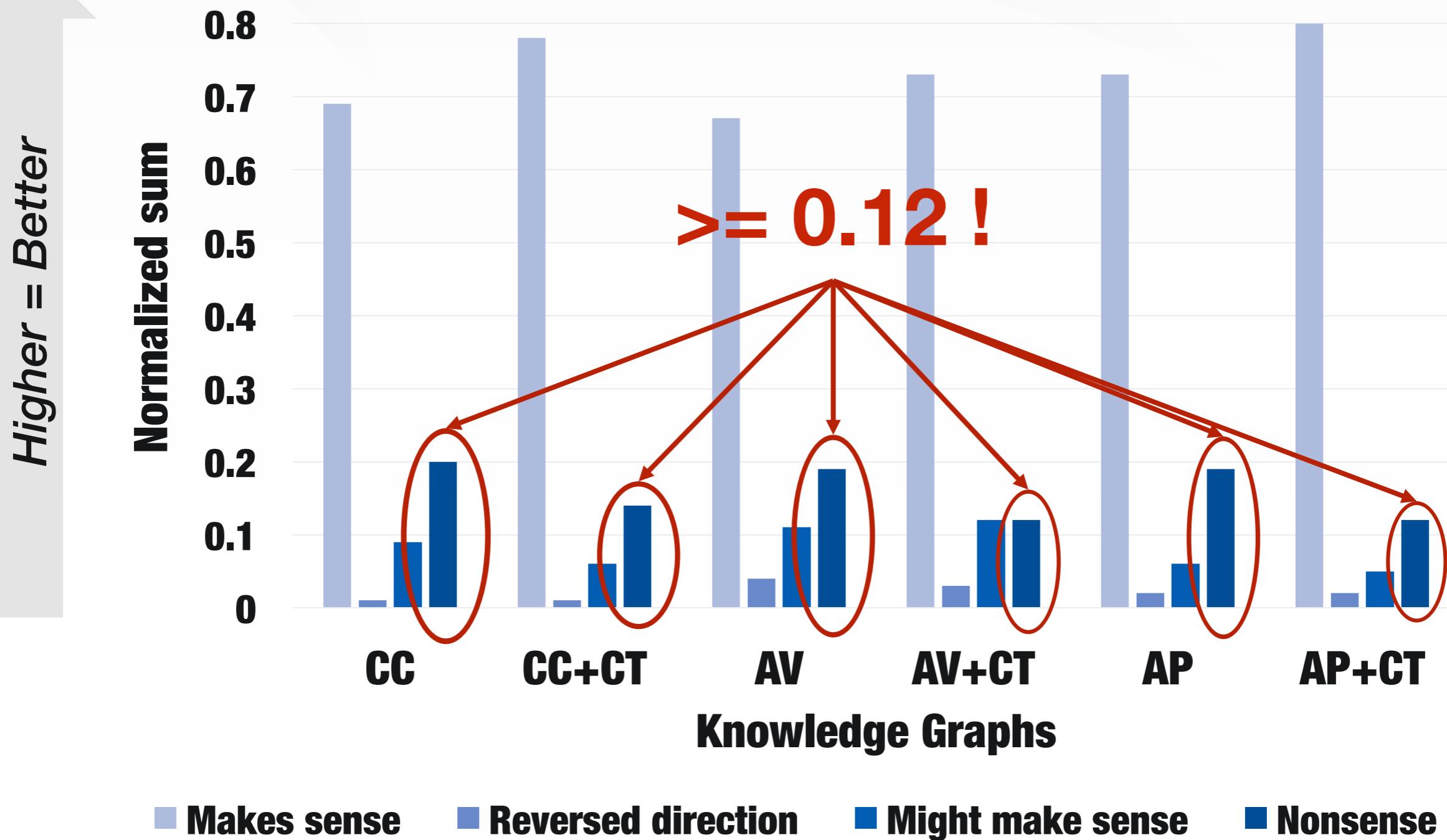
# Results

**High precision**



# Results

↑ High precision



# Limitations

---

- **Potential overlaps among relations**  
e.g. Use(laptop, processor) & Contain(laptop, processor)
- **Delimitation of the relations**  
e.g. Contain(mouse, genes)
- **Mixture of semantic meanings**  
e.g. Contain(mouse, brain) & TypeOf(mouse, device)
- **Hypothetic relations**  
e.g. Contain(artery, clot)

# **Conclusion & Future Work**

# Conclusion

---

- State of the art model for Relation Classification task
  - Linguistic information features and +/- sampling help !
- Create a dataset fitting Iprova's needs and build KGs
  - Precision not high enough yet
  - Has some limitations
  - Can be used as an help for humans
- This kind of Knowledge Graphs doesn't exist  
⇒ **We provide a tool to model domains of interest**

# Future Work

---

- Inferring new relations by using prior knowledge from KG
- Training pair-words embeddings
- Use pairwise ranking loss function
- Better filtering for Knowledge Graphs
- Improve concept extraction system

# Questions ?