# (Open) Information Extraction: Where are we going?

Claudio Delli Bovi
July 18th, 2016

SAPIENZA
UNIVERSITÀ DI ROMA

Linguistic Computing Laboratory

A2

ALLEN INSTITUTE
for ARTIFICIAL INTELLIGENCE

# About me

**dellibovi@di.uniroma1.it**

`http://wwwusers.di.uniroma1.it/~dellibovi`

bn:17381128n

Second-year PhD student

LCL group @ Sapienza

Advisor:  prof. Roberto Navigli

Focus (so far):  Disambiguation, (Open) Information Extraction

# Outline

**BabelNet and friends**: some background
Research work @ LCL Sapienza

**DefIE**: OIE from textual definitions
Delli Bovi, Telesca, Navigli:  **TACL 2015**

**KBUnify**: KB disambiguation and unification
Delli Bovi, Espinosa-Anke, Navigli: **EMNLP 2015**

# Outline

**BabelNet and friends**: some background
Research work @ LCL Sapienza

**DefIE**: OIE from textual definitions
Delli Bovi, Telesca, Navigli:  **TACL 2015**

**KBUnify**: KB disambiguation and unification
Delli Bovi, Espinosa-Anke, Navigli: **EMNLP 2015**

# Linguistic Computing Laboratory (LCL) @ Sapienza University of Rome

- Part of the **Computer Science Department** of Sapienza, focused on **Natural Language Processing**

- Some projects we have been involved in:
  - **MultiJEDI** (**1.3M €**): ERC Starting Grant
  - **LIDER** (**1.5 M €**): EU CSA
  - **Google Focused Research Award** (**300k $**)

# Multijedi_

## Multilingual joint word sense disambiguation

# Project

MultiJEDI is a 5-year ERC Starting Grant (2011-2016) headed by Prof. Roberto Navigli at the Linguistic Computing Laboratory of the Sapienza University of Rome. The project has two main objectives: creating large-scale lexical resources for dozens of languages, and enabling multilingual text understanding. The project has received funding from the European Union's specific programme 'Ideas' implementing the seventh framework programme (FP7-IDEAS-ERC) under grant agreement no. 259234.
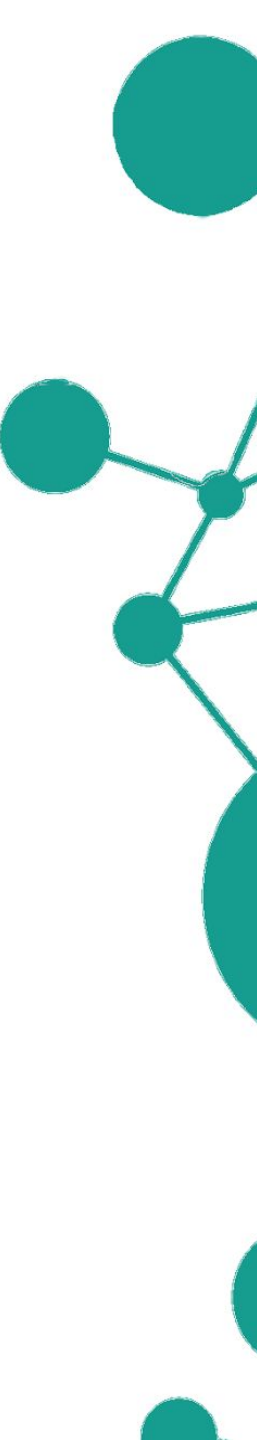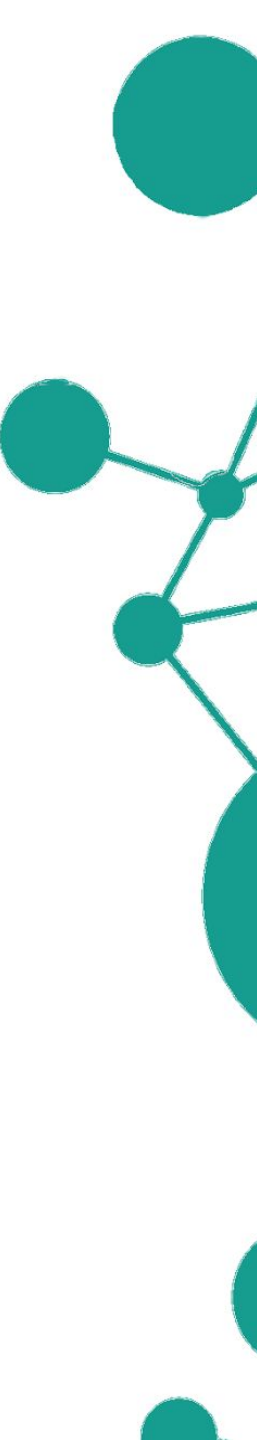
# BabelNet

- To the best of our knowledge, the largest **multilingual encyclopedic dictionary** and **semantic network** (almost **14M** entries in **271** languages and **380M** semantic connections)

# BabelNet

- To the best of our knowledge, the largest **multilingual encyclopedic dictionary** and **semantic network** (almost **14M** entries in **271** languages and **380M** semantic connections)

- Initially created as an integration of **Wikipedia** and **WordNet**, now BabelNet is a merger of many different resources (Wiktionary, Wikidata, OmegaWiki, VerbNet, ImageNet, …)
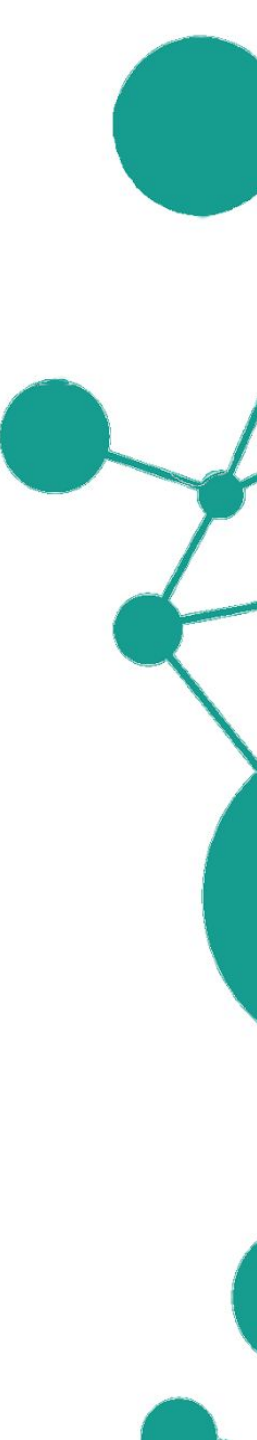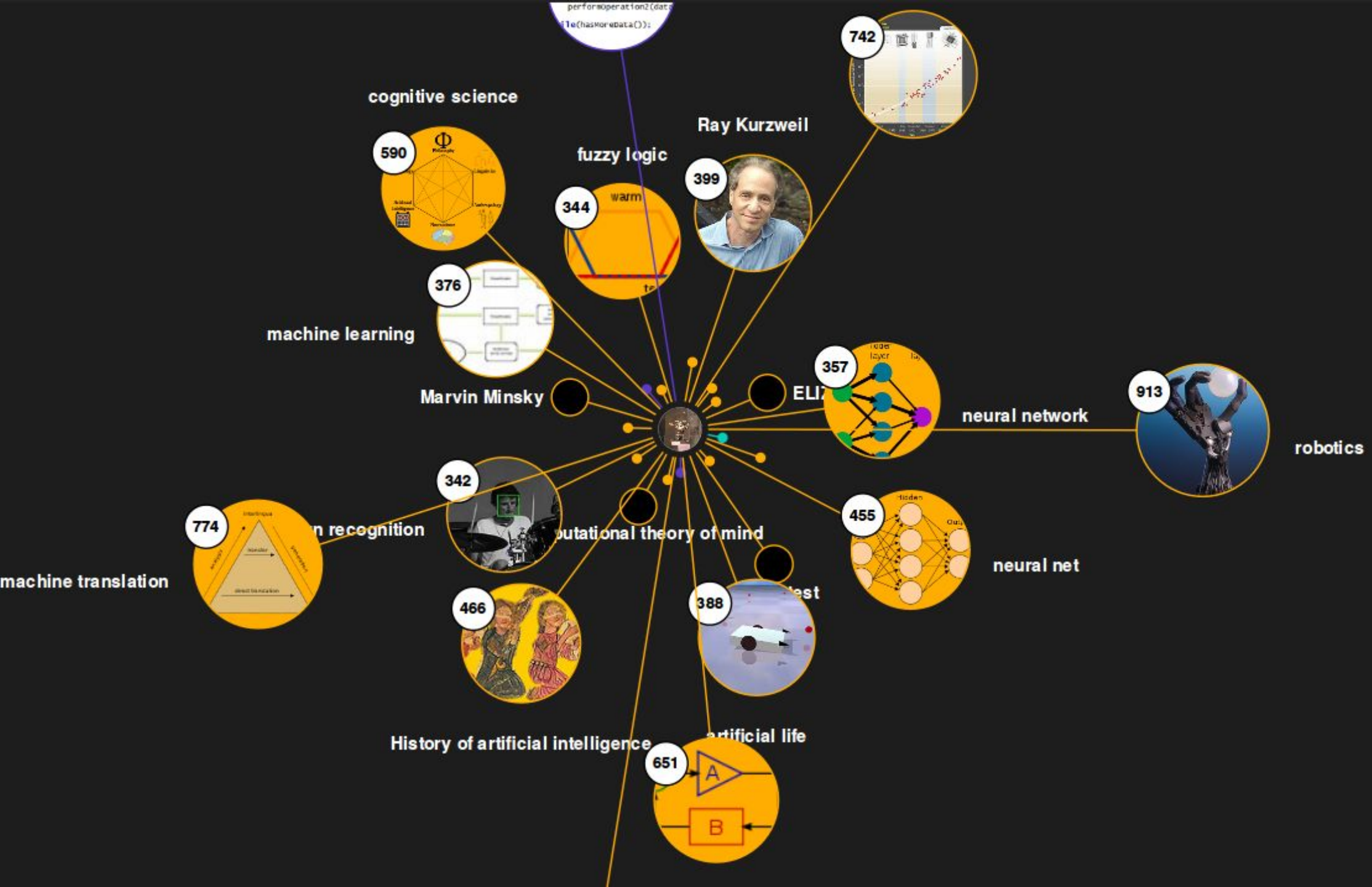
# BabelNet

- The integration is performed via an **automatic linking algorithm** and by filling in lexical gaps with the aid of **Machine Translation**
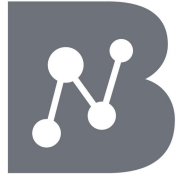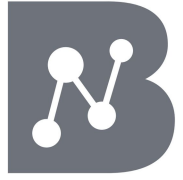
# BabelNet

- The integration is performed via an **automatic linking algorithm** and by filling in lexical gaps with the aid of **Machine Translation**

- BabelNet is composed of **Babel Synsets**, concepts or entities **lexicalized** ("WordNet-style") in many languages and featuring:

  - **is-a** relations
  - **domain** and **categories**
  - **images** and **definitions**
  - **translations**

performOperation2(dat

1e(hasMoreData());

cognitive science

**742**

**Ray Kurzweil**

**590**

fuzzy logic

**399**

warm

**344**

**376**

machine learning

Marvin Minsky

**357**

ELIZ

neural network

**913**

robotics

342

n recognition

utational theory of mind

**455**

neural net

**774**

machine translation

**466**

**388**

est

History of artificial intelligence

artificial life

**651**

A
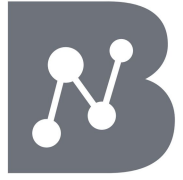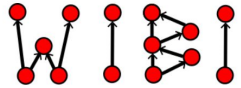
B

# BabelNet and friends

# BabelNet and friends

## Babelfy
A graph-based algorithm for multilingual joint **Word Sense Disambiguation** and **Entity Linking**, based on BabelNet

# BabelNet and friends

## Babelfy
A graph-based algorithm for multilingual joint **Word Sense Disambiguation** and **Entity Linking**, based on BabelNet

## The Wikipedia Bitaxonomy
An iterative algorithm for the automatic creation of a "**bitaxonomy**" for Wikipedia pages and categories

**... and much more!**

# BabelNet and my research

- BabelNet (especially in its early stages) was conceived as a **lexico-semantic resource** more than an actual **knowledge base**:
  - semantic connections are mostly **lexical relations** from WordNet or unspecified "**relatedness edges**" derived from Wikipedia hyperlinks
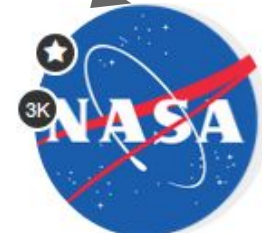


semantically related

**Atom Heart Mother**

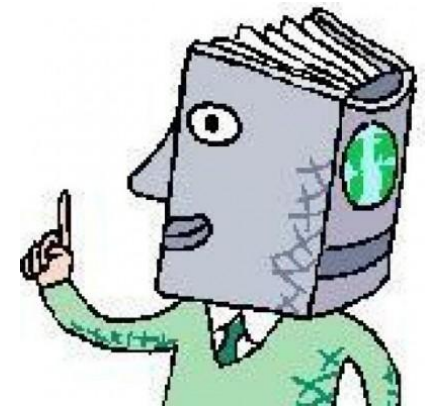**Pink Floyd**

semantically related

**Neil Armstrong**

**NASA**

# BabelNet and my research

- BabelNet (especially in its early stages) was conceived as a **lexico-semantic resource** more than an actual **knowledge base**:
  - semantic connections are mostly **lexical relations** from WordNet or unspecified "**relatedness edges**" derived from Wikipedia hyperlinks

- Construct from BabelNet a proper knowledge base with **labeled relations** (X is album by Y, X worked at Y, … )
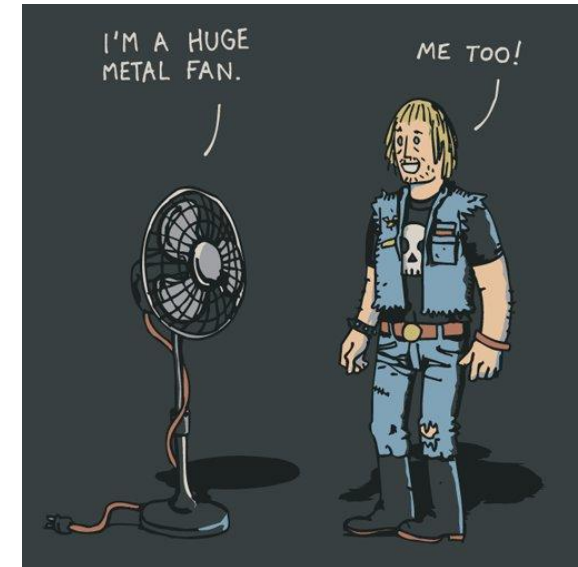
- Use **Open Information Extraction**!

# (Open) Information Extraction

**OIE is great, but…**

**Sparsity**:  many relation phrases express the same relationship (e.g. synonyms, paraphrases)

**Ambiguity**:  arguments (and relation phrases) are ambiguous!

# Outline

**BabelNet and friends**: some background
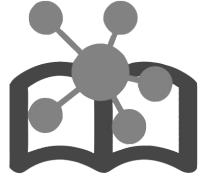Research work @ LCL Sapienza

**DefIE**: OIE from textual definitions
Delli Bovi, Telesca, Navigli: **TACL 2015**

**KBUnify**: KB disambiguation and unification
Delli Bovi, Espinosa-Anke, Navigli: **EMNLP 2015**

# DefIE: OIE from textual definitions

## The idea:

instead of targeting massive and noisy corpora (like the web) and then trying to find a smart way to cope with the noise

target smaller but "denser" (and virtually noise-free) corpora of **definitional knowledge**.
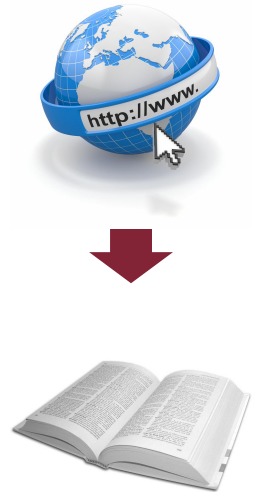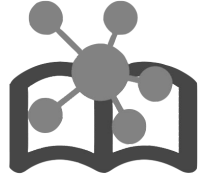
# DefIE: OIE from textual definitions

## The idea:

instead of targeting massive and noisy corpora (like the web) and then trying to find a smart way to cope with the noise

target smaller but "denser" (and virtually noise-free) corpora of **definitional knowledge**.
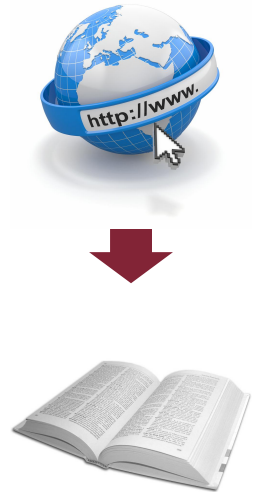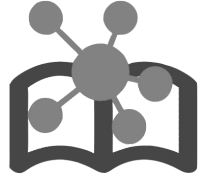
Apply OIE techniques to extract as much information as possible!

# DefIE: OIE from textual definitions

## The tools:

– An underlying **inventory/knowledge base** (to which arguments and relation patterns will be connected)

– A **WSD/EL system** (to disambiguate concepts and entity mentions across the input text)

– A **syntactic parser** (to construct meaningful relation patterns and avoid sparsity)

# DefIE: OIE from textual definitions

## The tools:

- An underlying **inventory/knowledge base** (to which arguments and relation patterns will be connected)

- A **WSD/EL system** (to disambiguate concepts and entity mentions across the input text)

- A **syntactic parser** (to construct meaningful relation patterns and avoid sparsity)

http://babelnet.org

**BabelNet**

**14 million** entries

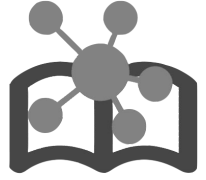both **lexicographic** and **encyclopedic** knowledge

# DefIE: OIE from textual definitions

## The tools:

- An underlying **inventory/knowledge base** (to which arguments and relation patterns will be connected)

- A **WSD/EL system** (to disambiguate concepts and entity mentions across the input text)

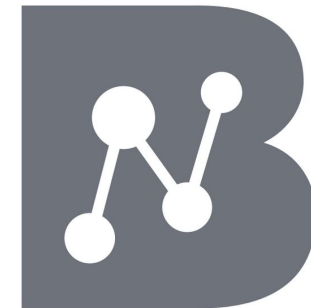- A **syntactic parser** (to construct meaningful relation patterns and avoid sparsity)

http://babelfy.org

**Babelfy**

unified graph-based approach to **EL** and **WSD**

unsupervised, based on **BabelNet**

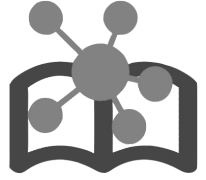# DefIE: OIE from textual definitions

## The tools:

– An underlying **inventory/knowledge base** (to which arguments and relation patterns will be connected)

– A **WSD/EL system** (to disambiguate concepts and entity mentions across the input text)

– A **syntactic parser** (to construct meaningful relation patterns and avoid sparsity)

**http://svn.ask.it.usyd.edu.au/ trac/candc**

**C&C tools**

log-linear parser and supertagger based on **CCG**

(theoretically) suited to **long-distance dependencies**

# DefIE: How it works

## 1. Extracting relation instances

*"Atom Heart Mother is the fifth
album by English band Pink Floyd."*

Textual definition  $d$

# DefIE: How it works

## 1. Extracting relation instances

Parsing

Dependency graph $G_d$

"Atom Heart Mother is the fifth album by English band Pink Floyd."

Disambiguation

bn:02070902n

bn:03292767n

**Atom Heart Mother** is the fifth **album** by **English** **band** **Pink Floyd**

bn:00002488n

bn:00102248a

bn:00008280n

Sense mappings $S_d$

# DefIE: How it works

## 1. Extracting relation instances

Syntactic-Semantic Graph $S_d^{sem}$



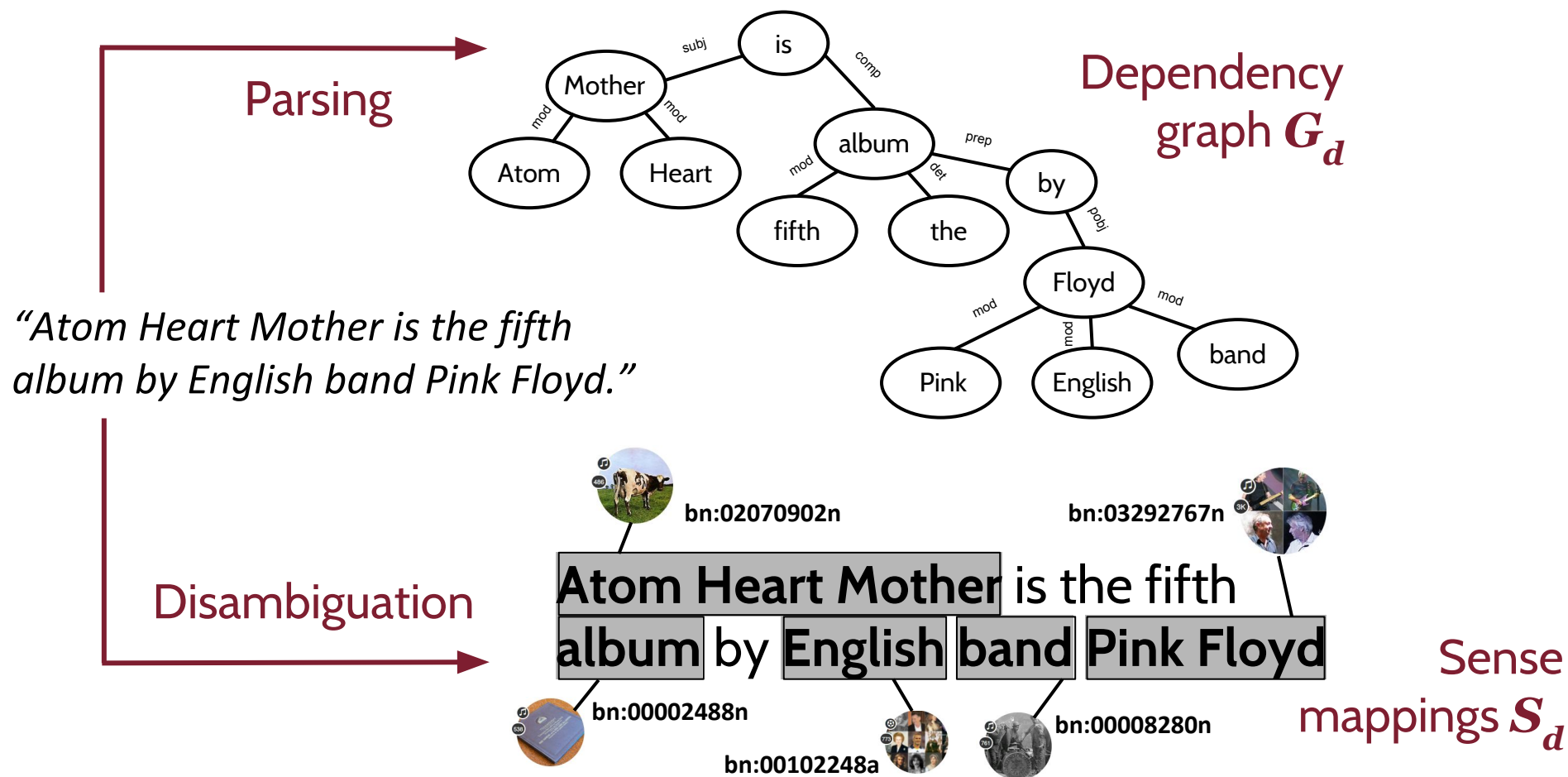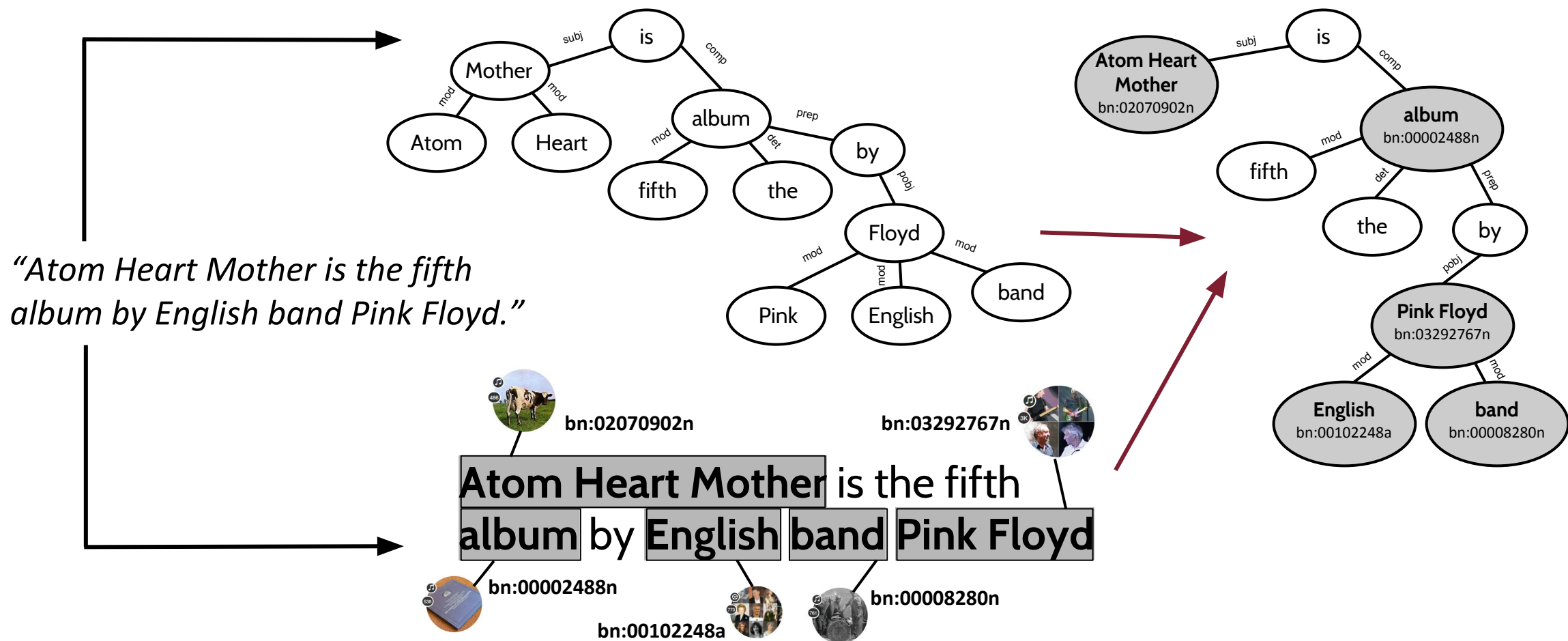*"Atom Heart Mother is the fifth album by English band Pink Floyd."*

# DefIE: How it works

## 1. Extracting relation instances

# DefIE: How it works

## 1. Extracting relation instances

**Extraction 2**

$$X \rightarrow \text{is} \rightarrow Y$$

$X$ = **Atom Heart Mother** bn:02070902n

$Y$ = **album** bn:00002488n

**Extraction 1**

$$X \rightarrow \text{is} \rightarrow \textbf{album} \rightarrow \text{by} \rightarrow Y$$
bn:00002488n

$X$ = **Atom Heart Mother** bn:02070902n

$Y$ = **Pink Floyd** bn:03292767n

## 1. Extracting relation instances

$$R_1: \quad X \rightarrow \texttt{is} \rightarrow Y$$

⟨Atom Heart Mother, album⟩
⟨Pink Floyd, band⟩
⋮
⟨Seattle, city⟩

$$R_2: \quad X \rightarrow \texttt{is} \rightarrow \textbf{album} \rightarrow \texttt{by} \rightarrow Y$$

bn:00002488n

⟨Atom Heart Mother, Pink Floyd⟩
⟨Mutter, Rammstein⟩
⋮
⟨Can't Get Enough, Barry White⟩

# DefIE: How it works

## 1. Extracting relation instances

R₁:  $X \rightarrow is \rightarrow Y$

Domain      Range

⟨Atom Heart Mother, album⟩
⟨Pink Floyd, band⟩
⋮
⟨Seattle, city⟩

R₂:  $X \rightarrow is \rightarrow$ **album** $\rightarrow by \rightarrow Y$
bn:00002488n

⟨Atom Heart Mother, Pink Floyd⟩
⟨Mutter, Rammstein⟩
⋮
⟨Can't Get Enough, Barry White⟩

# DefIE: How it works

## 2. Relation typing and scoring

# DefIE: How it works

## 2. Relation typing and scoring

For each relation R:

Substitute each domain and range argument with its **hypernym h** (using the BabelNet taxonomy) and generate a **probability distribution over semantic types** for the two sets

Compute the **entropy** of R as $\quad H_R = -\sum_{i=1}^{n} p(h_i) \, log_2 \, p(h_i)$

# **DefIE: How it works**

## 2. Relation typing and scoring

For each relation R:

Compute the **score** of R as

Total number of extracted instances for R

$$score(R) = \frac{|S_R|}{(H_R + 1)\, length(r)}$$

Domain and range entropy of R

Length of the relation pattern of R

# DefIE: How it works

## 2. Relation typing and scoring

| Pattern | Score | Entropy |
|---|---|---|
| X *directed by* Y | 4 025.80 | 1.74 |
| X *known for* Y | 2 590.70 | 3.65 |
| X *is* $\mathtt{election\ district}_{bn}^{1}$ *of* Y | 110.49 | 0.83 |
| X *is* $\mathtt{composer}_{bn}^{1}$ *from* Y | 39.92 | 2.08 |
| X *is* $\mathtt{street}_{bn}^{1}$ *named after* Y | 1.91 | 2.24 |
| X *is* $\mathtt{village}_{bn}^{2}$ *founded in 1912 in* Y | 0.91 | 0.18 |

# DefIE: How it works

## 3. Relation taxonomization

## 3. Relation taxonomization



Hypernym generalization

# DefIE: How it works

## 3. Relation taxonomization



Hypernym generalization

Substring generalization

# DefIE: Setup

**Dataset:**
whole set of English textual definitions in BabelNet 2.5

**4 357 327** items from **5** different sources (Wikipedia, WordNet, Wikidata, Wiktionary, OmegaWiki)

BabelNet

**EN  Atom Heart Mother** ◄)) · **Lulubelle III** ◄)) · **The Cow Album** ◄))

Atom Heart Mother is the fifth studio album by the English progressive rock band Pink Floyd. ◄))

⊖ *Fewer definitions*

W  1970 album by Pink Floyd. ◄))

▐▌▌  Album by Pink Floyd ◄))

**EN  Syd Barrett** ◄)) · **Syd Barett** ◄)) · **Syd barratt** ◄)) · **Barrett, Syd** ◄)) · **Bi5** ◄))

Roger Keith "Syd" Barrett was an English musician, composer, singer, songwriter and painter. ◄))

⊖ *Fewer definitions*

W  The late British singer and musician, formerly of Pink Floyd ◄))

◄))  Syd Barrett, born Roger Keith Barrett, was an English singer, songwriter, guitarist and artist. ◄))

# DefIE: Results

| | DefIE | NELL | PATTY | ReVerb | WiSeNet |
|---|---|---|---|---|---|
| **# Relations** | 255 881 | 298 | **1 631 531** | 664 746 | 245 935 |
| **Avg. extractions** | 81.68 | **7 013.03** | 9.68 | 22.16 | 9.24 |
| **# Extractions** | **20 352 903** | 2 089 883 | 15 802 946 | 14 728 268 | 2 271 807 |
| **# Entities** | 2 398 982 | 1 996 021 | 1 087 907 | **3 327 425** | 1 636 307 |
| **# Edges in the taxonomy** | 44 412 | - | 20 339 | - | - |

# DefIE: Results

**Other evaluations:**

- **Precision** and **coverage** of relations

- **Novelty** of information

- Quality of relation **taxonomization**

- Quality of **entity linking/disambiguation**

- **Impact** of definition sources

...

# DefIE: Future work

**Where from here?**

- Relation **clustering** (as in PATTY and WiSeNet)

- **Multilinguality**

- Relational **learning** and KB completion

- Harvest definitions from the **web**

- Adapt to "**general**" text

...

# Outline

**BabelNet and friends**: some background
Research work @ LCL Sapienza

**DefIE**: OIE from textual definitions
Delli Bovi, Telesca, Navigli:  **TACL 2015**

**KBUnify**: KB disambiguation and unification
Delli Bovi, Espinosa-Anke, Navigli: **EMNLP 2015**

# KB-Unify: Knowledge base unification via sense embeddings and disambiguation

## The idea:

PATTY
WiseNet
...

**Open Information
Extraction system**

NELL
ReVerb
...

**Linked Resources**

⟨ W **Armstrong**, *has worked at*, W **NASA** ⟩

**Unlinked Resources**

⟨ *Armstrong*, *works for*, *NASA* ⟩

# KB-Unify: Knowledge base unification via sense embeddings and disambiguation

## The idea:

PATTY
WiseNet
...

**Open Information Extraction system**

NELL
ReVerb
...

**Linked Resources**

**Unlinked Resources**

⟨ Armstrong , $r_{work}$ , NASA ⟩

$r_{work}$ = { has worked at ,

works for ,

employed at,

... }

# KB-Unify: Knowledge base unification via sense embeddings and disambiguation

## The tools:

- A **WSD/EL system** (to disambiguate unlinked resources)

- A unified **sense inventory S** (to make the various resources "speak to each other")

- A unified **vector space $V_S$** (to associate a vector with each item of **S**)

# KB-Unify: Knowledge base unification via sense embeddings and disambiguation

## The tools:

- A **WSD/EL system** (to disambiguate unlinked resources)

**Babelfy**

- A unified **sense inventory S** (to make the various resources "speak to each other")

**Babelnet**

- A unified **vector space** $V_S$ (to associate a vector with each item of **S**)

# KB-Unify: Knowledge base unification via sense embeddings and disambiguation

## The tools:

- A **WSD/EL system** (to disambiguate unlinked resources)

- A unified **sense inventory S** (to make the various resources "speak to each other")

- A unified **vector space V$_S$** (to associate a vector with each item of **S**)

**SensEmbed**
(Iacobacci et al., 2015)

Sense-based embedding model

Popular word2vec architecture (**skip-gram**) trained on a **sense-annotated corpus**

# KB-Unify: How it works

## A bird's-eye view

Linked Resources $K_D$

$K_{D1}$  $K_{D2}$

Unlinked Resources $K_U$

$K_{U1}$  $K_{U2}$

# KB-Unify: How it works

## A bird's-eye view



use **BabelNet mappings** to redefine each linked resource

**disambiguate** each unlinked resource using BabelNet as sense inventory (more on this later!)

# KB-Unify: How it works

👉 **Disambiguation**

# KB-Unify: How it works

## 👉 Disambiguation

Two basic intuitions:

1. Among all triples in target knowledge base, some of them (even if ambiguous) will be **easier to disambiguate**

   e.g.         ⟨ **Armstrong , works for ,  NASA** ⟩

# KB-Unify: How it works

## 👉 Disambiguation

Two basic intuitions:

1. Among all triples in target knowledge base, some of them (even if ambiguous) will be **easier to disambiguate**

   e.g. 〈 **Armstrong , works for , NASA** 〉

2. In general, the disambiguation strategy should vary according to the **degree of specificity** of each relation

# KB-Unify: How it works

## 👉 Disambiguation

Group the set of unlinked triples by relation

For each relation **r**:

- Extract and disambiguate a subset of high-confidence **seed argument pairs** for **r** ;

- Estimate the **specificity** of **r** by looking at the distribution of its disambiguated seeds in the vector space **V$_s$** ;

- Disambiguate the remaining argument pairs of **r** with Babelfy either **triple-by-triple** (if **r** is general) or **all at once** (if **r** is specific).

# KB-Unify: How it works

**Identifying seed argument pairs**

# KB-Unify: How it works

## Identifying seed argument pairs



$$\mathbf{v}_d = \{ v_d^1 \quad v_d^2 \quad v_d^3 \}$$

$$\langle \boxed{\text{Armstrong}},$$

works for ,

$$\boxed{\text{NASA}} \rangle$$

$$\mathbf{v}_g = \{ v_g^1 \}$$

SENSEMBED

$$\langle v_d^*, v_g^* \rangle = \operatorname*{argmax}_{v_d \in \mathbf{v}_d,\, v_g \in \mathbf{v}_g} \frac{v_d \cdot v_g}{\|v_d\|\,\|v_g\|}$$

# KB-Unify: How it works

**Identifying seed argument pairs**

$$\mathbf{v}_d = \{ v_d^1 \quad v_d^2 \quad v_d^3 \quad \}$$

$\langle$ **Armstrong**,

works for ,

**NASA** $\rangle$

$$\mathbf{v}_g = \{ v_g^1 \quad \}$$

**SENSEMBED**

$$\langle v_d^*, v_g^* \rangle = \zeta_{\text{dis}}$$

Seed Disambiguation Confidence

# KB-Unify: How it works

### Ranking relations by specificity

$$\mu_k = \frac{1}{|\mathbf{v}_k|} \sum_{v \in \mathbf{v}_k} \frac{v}{\|v\|} \ , \ k \in \{D, G\}$$

**Domain/Range
Centroids**

$$\sigma_k^2 = \frac{1}{|\mathbf{v}_k|} \sum_{v \in \mathbf{v}_k} \left(1 - \cos\left(v, \mu_k\right)\right)^2$$

**Domain/Range
Variances**

# KB-Unify: How it works

## 🏆 Ranking relations by specificity

$$\mu_k = \frac{1}{|\mathbf{v}_k|} \sum_{v \in \mathbf{v}_k} \frac{v}{\|v\|} \, , \; k \in \{D, G\}$$

**Domain/Range Centroids**

$$\sigma_k^2 = \frac{1}{|\mathbf{v}_k|} \sum_{v \in \mathbf{v}_k} \left(1 - \cos\left(v, \mu_k\right)\right)^2$$

**Domain/Range Variances**

$$Gen(r) = \frac{\sigma_D^2 + \sigma_G^2}{2}$$

Specificity threshold:

$$\boldsymbol{\delta}_{\text{spec}}$$



High $Gen(r)$

$\mathbf{v}_D$    $\mathbf{v}_G$

Low $Gen(r)$

$\mathbf{v}_D$    $\mathbf{v}_G$

# KB-Unify: How it works

## Disambiguation with Relation Context



unlinked triples — disambiguated seeds — specificity ranking — $\delta_{spec}$ — general → triple-by-triple disambiguation — Babelfy — specific → relation-by-relation disambiguation

# KB-Unify: How it works

## A bird's-eye view



represent each relation in the **unified vector space $V_S$** and compare them pairwise

Linked Resources $K_D$

$K_{D1}$ $K_{D2}$

Inter-Resource Linking

Unlinked Resources $K_U$

$K_{U1}$ $K_{U2}$

Disambiguation

Redefined Resources $K^S$

$K_{D1}^S$ $K_{D2}^S$ $K_{U1}^S$ $K_{U2}^S$

Relation Alignment

Unified Resource KB*

# KB-Unify: How it works

**Relation alignment**

# KB-Unify: How it works

## Relation alignment
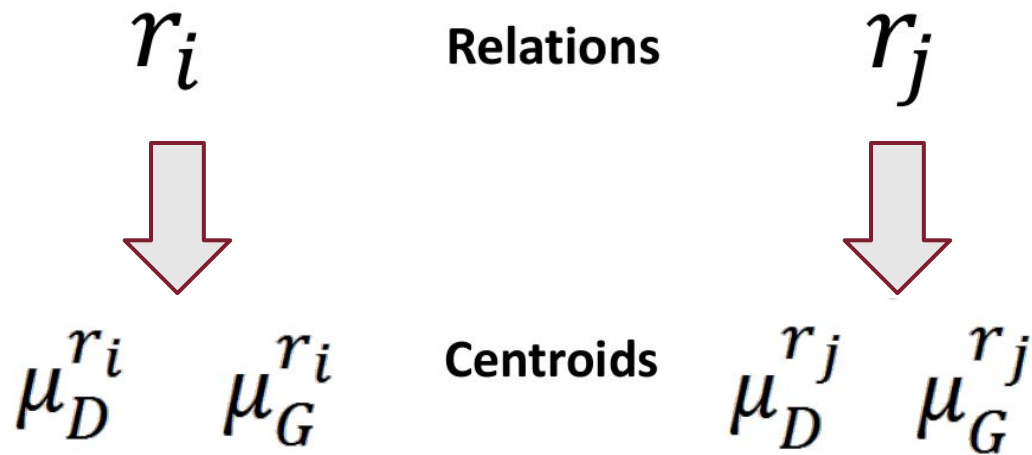
For each relation pair $\langle r_i, r_j \rangle$:

$r_i$      **Relations**      $r_j$

$\mu_D^{r_i}$   $\mu_G^{r_i}$     **Centroids**     $\mu_D^{r_j}$   $\mu_G^{r_j}$

# KB-Unify: How it works

## Relation alignment

For each relation pair $\langle r_i , r_j \rangle$:

$r_i$ **Relations** $r_j$

$\mu_D^{r_i}$ $\mu_G^{r_i}$ **Centroids** $\mu_D^{r_j}$ $\mu_G^{r_j}$

Compare domain and range centroids pairwise:

$$s_k = \frac{\mu_k^{r_i} \cdot \mu_k^{r_j}}{\|\mu_k^{r_i}\| \, \|\mu_k^{r_j}\|}$$

$k \in \{D, G\}$
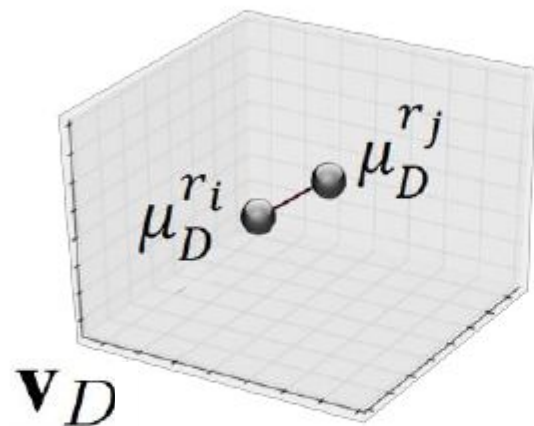
**Relation Centroid Similarity**
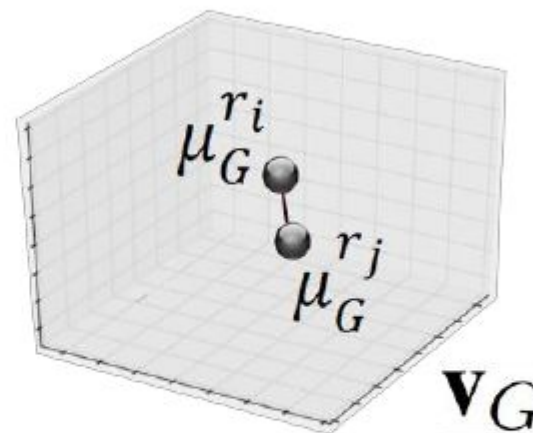
# KB-Unify: How it works

## ⇅ Relation alignment

Fix a similarity threshold $\delta_{align}$:



**Domain Centroids**

**Range Centroids**

$\frac{1}{2}(s_D + s_G) \geq \delta_{align}$? Align $r_i$ and $r_j$ and merge them in the same cluster

# KB-Unify: How it works

## Relation alignment

Fix a similarity threshold $\delta_{align}$:



**Domain Centroids** — $\mu_D^{r_i}$, $\mu_D^{r_j}$, $\mathbf{v}_D$

**Range Centroids** — $\mu_G^{r_i}$, $\mu_G^{r_j}$, $\mathbf{v}_G$

$\frac{1}{2}(s_D + s_G) < \delta_{align}$? Leave $r_i$ and $r_j$ in separate clusters

# KB-Unify: Experiments

## Evaluation

Experimental setup:

**Linked Resources $K_D$ :**

{ **PATTY** , **WISENET** }

| | |
|---|---|
| 1,631,531 relations | 245,935 relations |
| 15,802,946 triples | 2,271,807 triples |

**Unlinked Resources $K_U$ :**

{ **NELL** , **REVERB** }

| | |
|---|---|
| 298 relations | 1,299,844 relations |
| 2,245,050 triples | 14,728,268 triples |

# KB-Unify: Experiments

## Disambiguation



Seed Precision:

$\zeta_{dis}$
- 0.5 - 0.7
- 0.7 - 0.9
- 0.9 - 1.0

PATTY: 0.98 0.98 1
WiSeNet: 0.96 0.96 0.97
NELL: 0.95 0.99 1
ReVerb: 0.93 0.94 0.95

Coverage:

$\delta_{spec}$
- 0.8
- 0.5
- 0.3

PATTY: 0.62 0.52 0.41
WiSeNet: 0.6 0.54 0.53
NELL: 0.77 0.51 0.45
ReVerb: 0.41 0.25 0.13

# KB-Unify: Experiments

## Specificity ranking

For each ranked relation compute $Gen(r)$ against the **average argument similarity** $\bar{s}$:

# KB-Unify: Experiments

## Specificity ranking

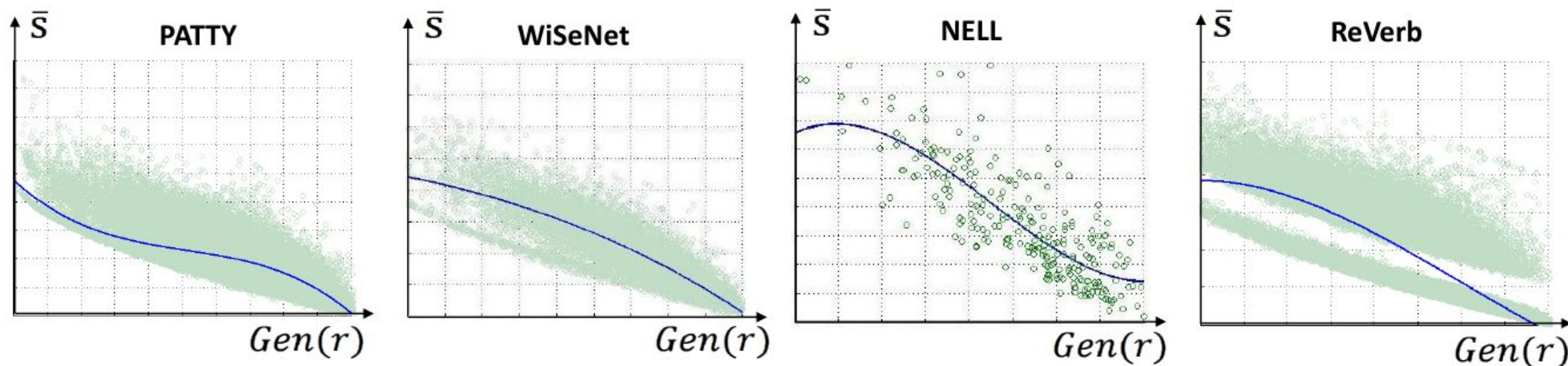For each ranked relation compute $Gen(r)$ against the **average argument similarity** $\bar{s}$:

# KB-Unify: Experiments

## Specificity ranking

For each ranked relation compute $Gen(r)$ against the **average argument similarity** $\bar{s}$:
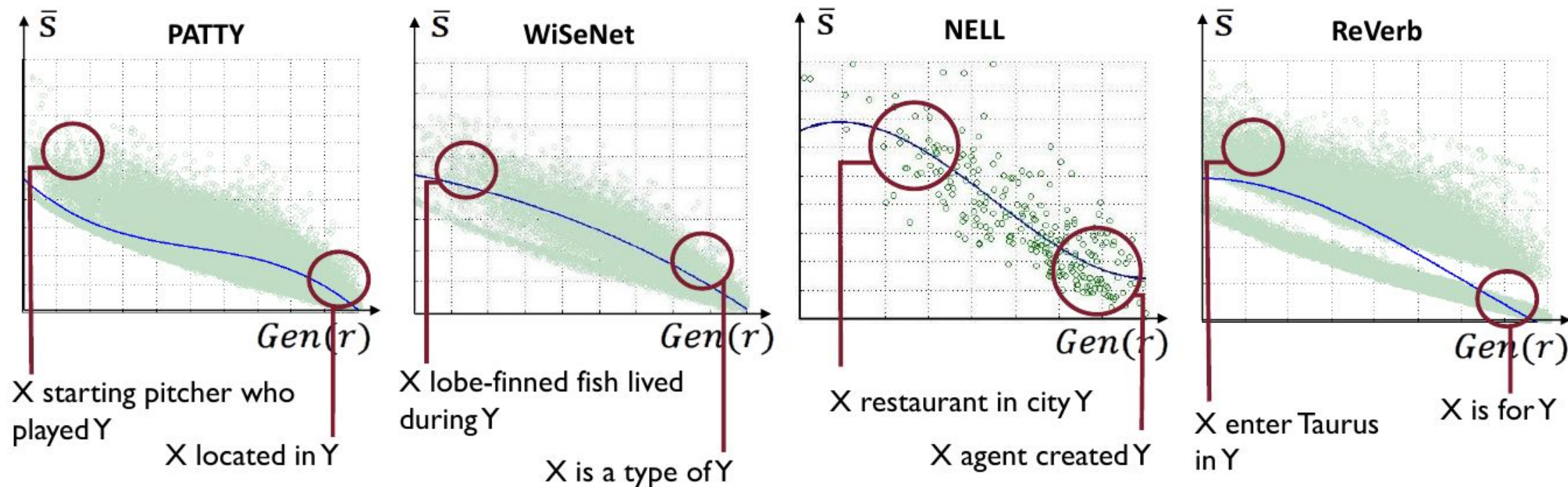


PATTY

X starting pitcher who played Y

X located in Y

WiSeNet

X lobe-finned fish lived during Y

X is a type of Y

NELL

X restaurant in city Y

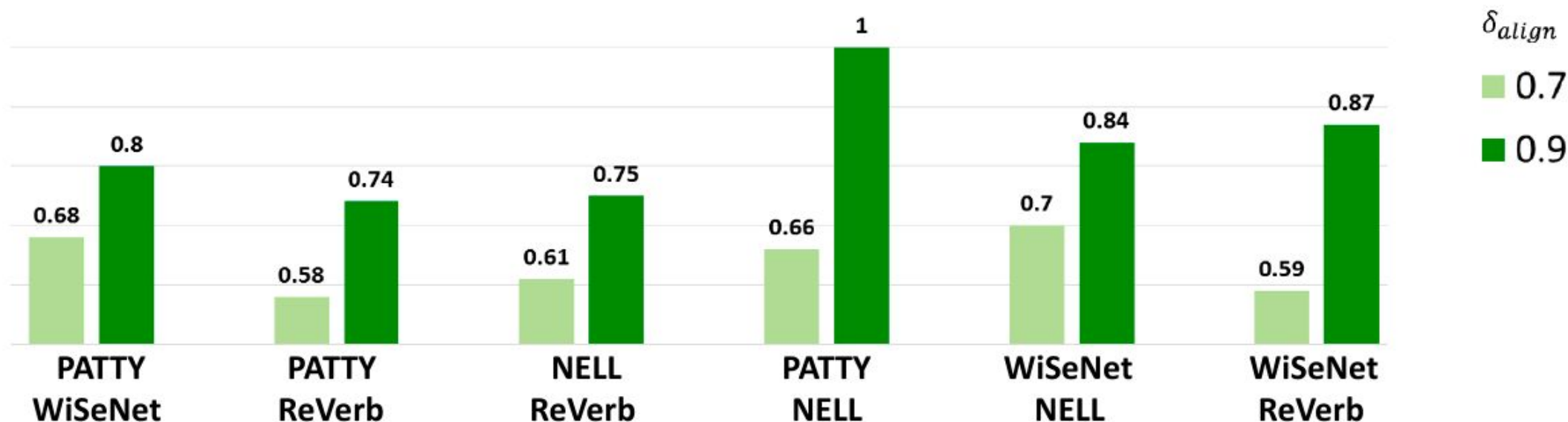X agent created Y

ReVerb

X enter Taurus in Y

X is for Y

# KB-Unify: Experiments

## Cross-resource relation alignment

Samples of **150 candidate alignments** for different alignment thresholds $\delta_{align}$ manually evaluated (in terms of **paraphrasing**) by two human judges

# KB-Unify: Experiments

## Cross-resource relation alignment

Some examples:

| PATTY-WISENET | | $\zeta_{align}$ |
|---|---|---|
| portrayed | 's character | 0.84 |
| debuted in | first appeared in | 0.86 |

| NELL-PATTY | | $\zeta_{align}$ |
|---|---|---|
| worksfor | was hired by | 0.72 |
| riveremptiesintoriver | tributary of | 0.89 |

| PATTY-REVERB | | $\zeta_{align}$ |
|---|---|---|
| language in | is spoken in | 0.81 |
| mostly known for | plays the role of | 0.70 |

| NELL-WISENET | | $\zeta_{align}$ |
|---|---|---|
| animaleatfood | feeds on | 0.72 |
| teamhomestadium | play their home games at | 0.88 |

| NELL-REVERB | | $\zeta_{align}$ |
|---|---|---|
| bookwriter | is a novel by | 0.88 |
| personleadscity | is the mayor of | 0.60 |

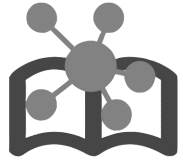| REVERB-WISENET | | $\zeta_{align}$ |
|---|---|---|
| has a selection of | offers | 0.82 |
| had grown up in | was born and raised in | 0.85 |

# KB-Unify: Future work

**Where from here?**

- Less "naïve" **relation alignment** procedure

- **Iterative** algorithm for disambiguation and alignment (EM-style)

- Unify OIE-based KBs with **hand-curated resources** (Wikidata, DBpedia, etc.)
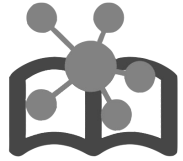
...

# Wrap up and Conclusion

# Wrap up and Conclusion

**DefIE:** A full-fledged OIE pipeline targeted to textual definitions, with explicit semantic characterization of both arguments and relation patterns

# Wrap up and Conclusion

**DefIE:** A full-fledged OIE pipeline targeted to textual definitions, with explicit semantic characterization of both arguments and relation patterns

**KB-Unify:** An approach to knowledge base disambiguation and unification based on a shared sense inventory and a sense-based vector space model

# Wrap up and Conclusion

**Take-home message(s):**

Web-scale OIE is absolutely great, but...

# Wrap up and Conclusion

**Take-home message(s):**

Web-scale OIE is absolutely great, but…

1. **Definitional knowledge is important**: sometimes it is worth just stepping back and analyze from where valuable information is extracted (**quality vs. quantity**)
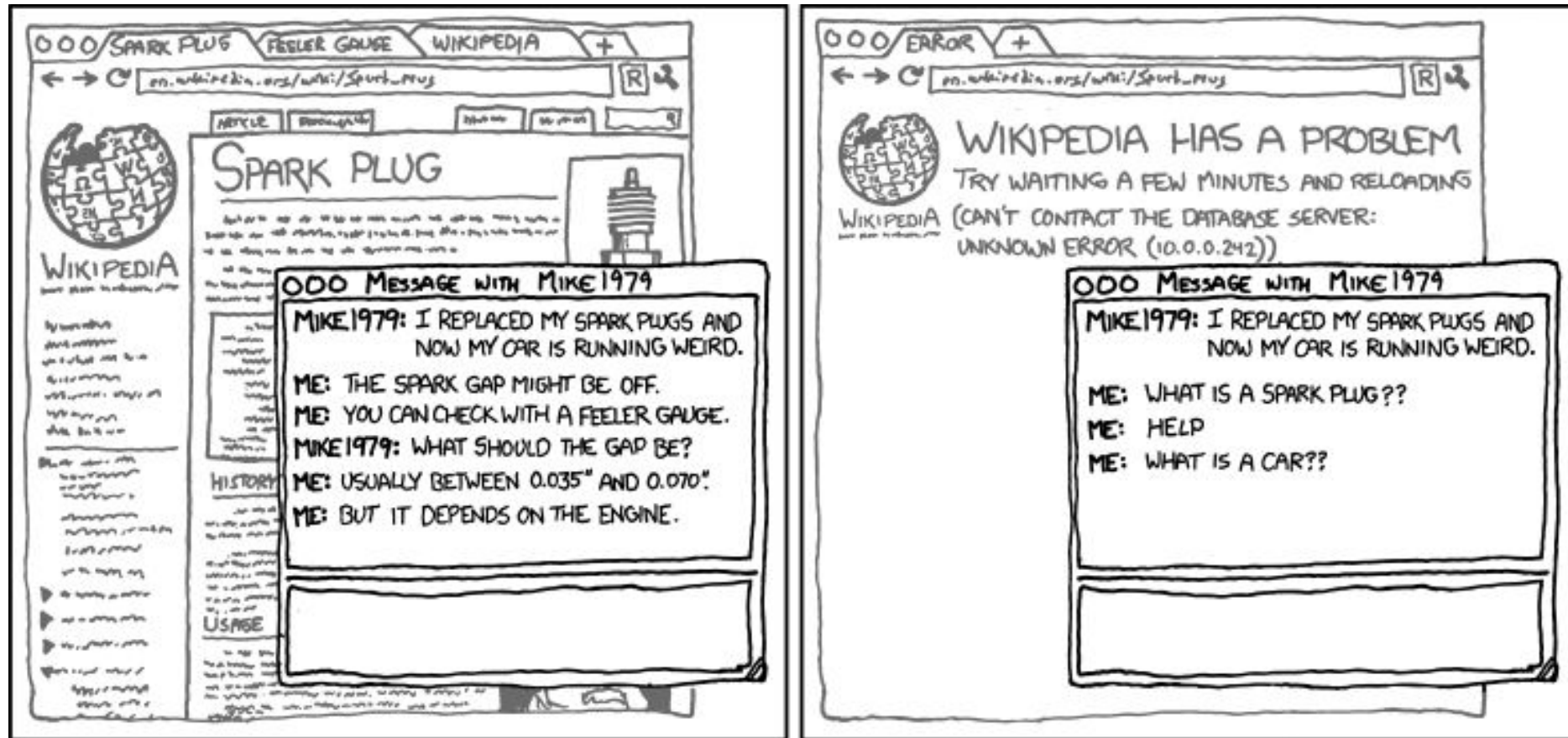
# Wrap up and Conclusion

**Take-home message(s):**

Web-scale OIE is absolutely great, but...

1. **Definitional knowledge is important**: sometimes it is worth just stepping back and analyze from where valuable information is extracted (**quality vs. quantity**)

2. **Making sense of the output is important:** semantic analysis can be used to let different OIE outputs "speak to each other" and benefit from mutual enrichment

# Thank you!



xkcd, "Extended Mind"