# PMBots

From Raw Parallel Documents to
the Documents with several annotation layers
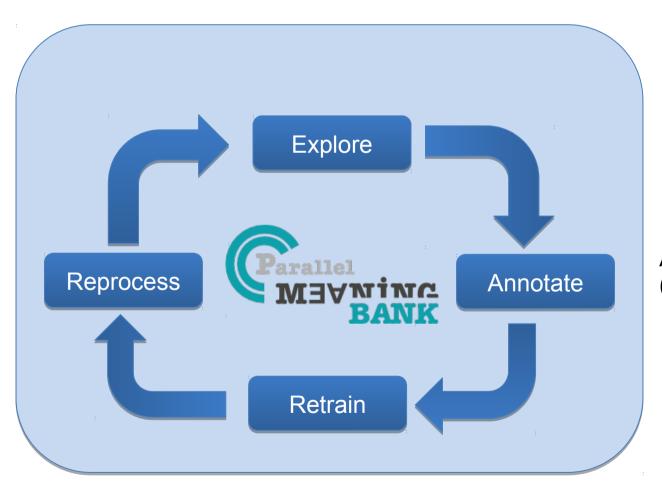
285K documents (863K sentences, 11.4m words)

5 m words

1,1m words

3,9m words

1,4m words

# Annotated English Document

| My | old | ex | is | a | thief | ! |
|---|---|---|---|---|---|---|
| HAS | IST | CON | NOW | DIS | CON | NIL |
| speaker | old | ex | be | a | thief | ! |
| NP/(N/PP) | N/N | N/PP | (S[dcl]\NP)/NP | NP/N | N | S[dcl]\S[dcl] |

>B     <

**old ex**       **is a thief**

N/PP        S[dcl]\NP

>

**My old ex**

NP

<

**My old ex is a thief**

S[dcl]

<

**My old ex is a thief !**

S[dcl]

# Usage of the PMB

- Train and evaluate several systems for the languages
  - Tokenizer
  - Semantic tagger
  - Coreference resolution system
  - CCG parser
  - Semantic parser
- Studying compositional semantics in a data-driven fashion
- Cross-lingual projection of annotations:

  From an annotated English document

  to the (DE/NL/IT) document

# Annotation with Bootstrapping



Explore

Annotate

Retrain

Reprocess

Parallel MEANING BANK

Adding BoWs
(Bits of Wisdoms)

# Adding Bulk of BoWs

## 667 search results for |CD|gold

The search finds tokens in all accepted documents of the PMB with the help of the several filters. Filtering is done based on the predefined list of values. An exception is a search according to a token/lemma, where a search term is identical (case-insensitive) to the token or the lemma of a word token. If you want to filter by semantic class, try the semantic lexicon. To search for multi-word expressions, such as New York, which are seen as a single token, use the tilde, e.g. 'New~York'.

By pressing the 'Edit'-button on the right, you can edit semantic tags of the found tokens.

Page: previous 1 2 3 4 5 6 7 next

Save changes  Canc

| document | category | pos | roles | token/lemma | sense | semtag | snippet | subcorpus | part |
|----------|----------|-----|-------|-------------|-------|--------|---------|-----------|------|
| **FILTER BY:** | [all] ▼ | CD ▼ | [all] ▼ | | [all] ▼ | [all] ▼ | | [all] ▼ | gold ▼ |
| **SET ALL TO:** | | | | | QUC ▼ | | | | |
| 00/0713 | N/N | CD | | 5 | O / o ▼ | QUC / quc ▼ | ø **Five** men have been killed by a serial killer in ø Mumbai . | RTE | 00 |
| 00/1164 | N/N | CD | | 5 | O / o ▼ | QUC / quc ▼ | ø **Five** rings | INTERSECT | 00 |
| 00/1177 | N/N | CD | | 15 | O / o ▼ | QUC / quc ▼ | In ø 1937 , the company moved to its current location in ø Reading , about **15** miles from ø central Cincinnati . | INTERSECT | 00 |
| 00/1177 | N | CD | | 1937 | O | YOC / quc ▼ | In ø **1937** , the company moved to its current location in ø Reading , about 15 miles from ø central Cincinnati . | INTERSECT | 00 |
| 00/1348 | N | CD | | 6 | O / o ▼ | NTH / quc ▼ | ø Article **6** . ø Marriage and ø family enjoy the special protection of the state . | INTERSECT | 00 |
| 00/1357 | N | CD | | 7 | O / o ▼ | NTH / quc ▼ | ø Article **7** . All separate alliances and all treaties of a political nature between ø Cantons are prohibited . | INTERSECT | 00 |
| 00/1543 | N/N | CD | | 2 | O / o ▼ | QUC / quc ▼ | There are ø only **two** possibilities . | Tatoeba | 00 |
| 00/1555 | NP/PP | CD | | one | O / o ▼ | DIS / quc ▼ | The bridge collapsed when **one** of the cables broke . | Tatoeba | 00 |
| 00/1861 | N/N | CD | | 1 | O / o ▼ | QUC / quc ▼ | The box is leaning to ø **one** side . | Tatoeba | 00 |
| 00/1936 | N/N | CD | | 2 | O / o ▼ | QUC / quc ▼ | My watch may be ø one or **two** minutes fast . | Tatoeba | 00 |
| 00/1936 | N/N | CD | | 1 | O / o ▼ | QUC / quc ▼ | My watch may be ø **one** or two minutes fast . | Tatoeba | 00 |
| 00/1978 | N/N | CD | | 3 | O / o ▼ | QUC / quc ▼ | They got married ø **three** months later . | Tatoeba | 00 |
| 00/1981 | N | CD | | 2 | O / o ▼ | QUC / quc ▼ | It takes ø **two** to tango . | Tatoeba | 00 |
| 00/2145 | N/N | CD | | 6 | O / o ▼ | QUC / quc ▼ | The admission costs ø **six** euros but on ø Sundays it 's free . | Tatoeba | 00 |

# Phrase Search

## Phrase search

Query: `pos=CD pos=NNS of pos=NN`

All queries are interpreted as phrase queries, with tokens separated by whitespace. For each token, the default is to match lemma OR token, other constraints can be specified using the equals sign, multiple constraints for the same token can be combined using pipe symbols. POS tags allow `*` as a wildcard.

For example, use `pos=MD not` to search for a modal followed by *not*, `can|pos=NN` for *can* tagged as a singular common noun, `pos=VB*` `without` for finding any verb form followed by *without*, and so on.

| | |
|---|---|
| 00/0027 | out of the huge riches extracted from their tribal lands - where the bulk of Nigeria 's 2.3 **million barrels of petroleum** are pumped daily . |
| 00/0294 | government and rebels of the Lord 's Resistance Army are nearing a peace deal to end more than **20 years of conflict** . Officials close to the negotiations say an accord signed Friday provides for the disarmament and demobilization of the |
| 00/0633 | in Darfur . The government is still rejecting any U.N. deployment to the region , where more than **three years of violence** has killed nearly 200,000 people , and displaced some two million others . |
| 01/0025 | from patients . The official Xinhua news agency says investigators from China 's Health Ministry uncovered more than **200 cases of hospital** staff members buying and selling medicine for personal profit . The health workers are said to have received $ |
| 01/0121 | television station . Reporters Without Borders and the Burma Media Association welcomed the release of the reporters after **two days of detention** . They say the journalists are in good shape - despite being shaken by the incident . The two |
| 01/0141 | off the coast of Indonesia 's Sumatra Island since early January . The probe of the area where **two pieces of ocean** floor collided shows a ridge of mud hundreds of meters thick where the seawater was forced up to form |
| 01/0168 | ) . In addition to those killed , 728 claims were filed for employees who missed more than **four days of work** . The inspector 's report says " Iraq 's unsettled security environment continues to present grave risks for contractors |
| 01/0368 | engine says it will install enough solar grids at its complex in Mountain View , California to generate **1.6 megawatts of electricity** . That 's enough energy to light up about 1,000 California homes . The company hopes the sun will |
| 01/0379 | peace prize laureate has been vocal about pursuing democratic reforms since her release November 13 from more than **seven years of house** arrest . But she has also been careful not to verbally challenge Burma 's ruling generals . Nambiar also |
| 01/0451 | sides say they hope to reach agreement at the Cairo summit on a formal truce ending more than **four years of violence** . Meanwhile , Palestinian gunmen interrupted a two-week old de~facto truce Thursday with an attack on an Israeli military |
| 01/0490 | Kashmir . The report describes the human rights situation as poor in Afghanistan , a country recovering from **20 years of war** . It notes violence against women and minorities , as well as restrictions on personal freedoms . Similarly |

Despite this adding BOWs manually is associated with a lot of human annotation efforts

# Idea of PMBots

- PMBots (prudent monitoring bots) are bots that monitor the annotated documents and with high reliability identify (plausible) annotation errors.

- PMBots use context sensitive rules:

    [sem=QUC; cat=N/N]   [ ]$_X$   [tok=of]   [cat=N]

    ===> [ sem=UOM,CON ]$_X$

- PMBots can be prudent (verify) or confident (rewrite)