



# Word Sense Disambiguation

Didier SCHWAB  
Didier.Schwab@imag.fr

# What is Word Sense Disambiguation ?

- Natural languages are ambiguous:

*The mouse ate some cheese*

# What is Word Sense Disambiguation ?

- Natural languages are ambiguous:

*The **mouse** ate some cheese*

# What is Word Sense Disambiguation ?

- Natural languages are ambiguous:

*The **mouse** ate some cheese*



# What is Word Sense Disambiguation ?

- Natural languages are ambiguous:

*The mouse ate some cheese*



# What is Word Sense Disambiguation ?

- Natural languages are ambiguous:

*The mouse ate some cheese*



# What is Word Sense Disambiguation ?

- Most words have several possible meanings
- => Very few have a single meaning
- Monosemic : '*neuroleptic*', '*daucus carota*',
- Polysemic : '*mouse*', '*rabbit*', '*carot*'
- In English : the 121 most frequent nouns
  - On average 1 word out of five in actual texts
  - ~7.8 meanings per word (in Princeton WordNet)
- What is (often) really easy task for a human is difficult for a computer
- Finding a better sense for a word in a text is called  
**Word Sense Disambiguation**

# What is Word Sense Disambiguation ?

- Aim of WSD: selecting a sense for each word in a text from an inventory (set) of predefined possibilities
- A word sense is the meaning of a word in a given context
- Inventories are produced from dictionaries, raw texts, ...
- How to represent word senses ?
- How to fetch the meanings of a word ?



# Sets of Word Senses

- How to fetch the meanings of a word ?
  - With respect to a dictionary, a lexical base...
    - **mouse#1** : any of numerous small rodents...
    - **mouse#2** : a hand-operated electronic device...
  - With respect to the translation in a second language
    - **mouse#1** : tikus
    - **mouse#2** : tetikus

# Sets of Word Senses

- How to fetch the meanings of a word ?
  - With respect to the context where it occurs...
    - **mouse#1** : „The cat hurt the mouse“ ; “The mouse is eating the cheese“ ; ...
    - **mouse#2** : „The mouse is linked to the computer.“ ; „My mouse is broken.“ ; ...
  - With respect to relations it shares in a semantic network
    - **mouse#1** : hypernyms (kind-of) : '*rodent*', '*mammal*',... ; related-to : '*mousy*', '*mousey*'
    - **mouse#2** : hypernyms : '*electronic device*' ; related-to : '*to mouse*'
  - Others
  - Combinations



# Sense Tagging

# Sense Tagging

- Given a pre-defined inventory of word senses
- Given a text
- Tag each ambiguous word occurrence with the most likely word sense
- Example :
- 'The cat is eating the mouse'

# Sense Tagging

*'The cat  
is eating  
the mouse'*

# Sense Tagging

*'The cat  
is eating  
the mouse'*

Word Sense  
Disambiguator

# Sense Tagging

**cat#1** : feline  
**cat#2** : caterpillar

**mouse#1** : rodent  
**mouse#2** : device

*'The cat  
is eating  
the mouse'*

input

Word Sense  
Disambiguator

# Sense Tagging

**cat#1** : feline  
**cat#2** : caterpillar

**mouse#1** : rodent  
**mouse#2** : device

*'The cat  
is eating  
the mouse'*

Word Sense  
Disambiguator

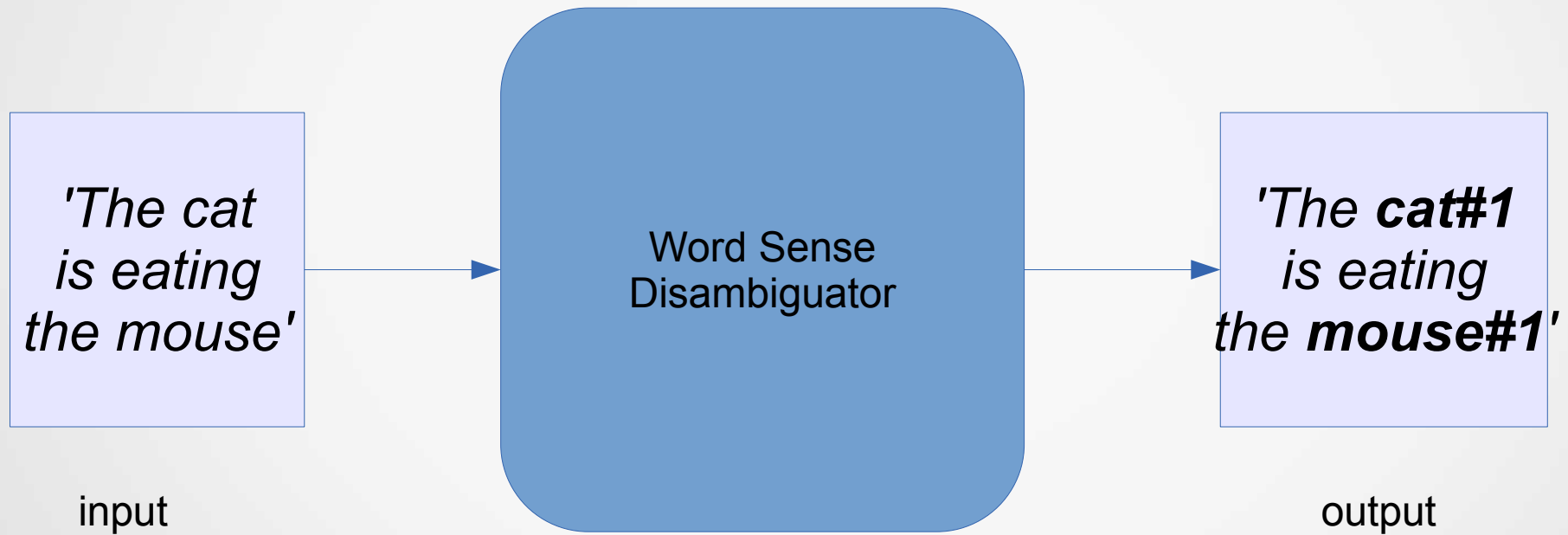
*'The **cat#1**  
is eating  
the **mouse#1**'*

input

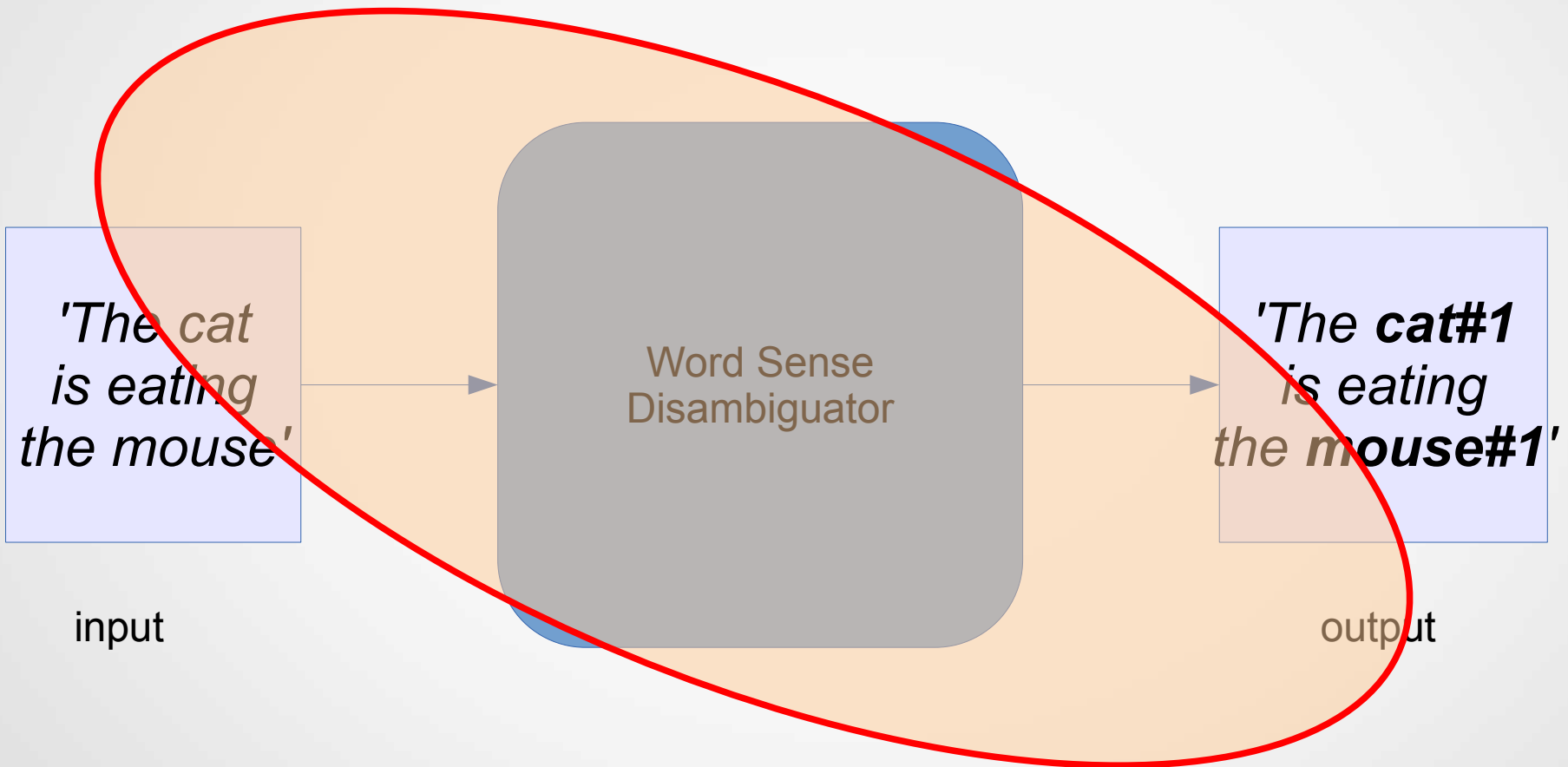
output



# Sense Tagging



# Sense Tagging





# Practical Applications

# WSD for machine translation

- Which translation of "mouse" ?



tetikus



tikus

- Which translation of "bank" in French?

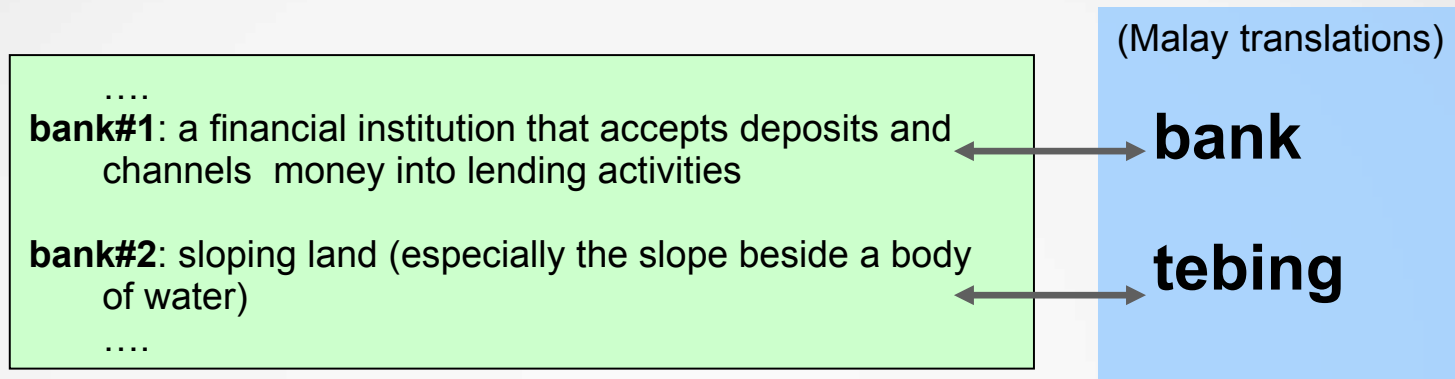
Bank → Berge



Bank → Banque



# WSD for machine translation



...withdraw money from the **bank**...

sense-tag  
(WSD)

...withdraw money from the **bank#1**...

select  
translation  
word

**Malay output**

...mengeluarkan wang dari **bank**...

# WSD for Information Retrieval



*mouse*



*mouse*



*house*



# WSD for Information Retrieval

Query :

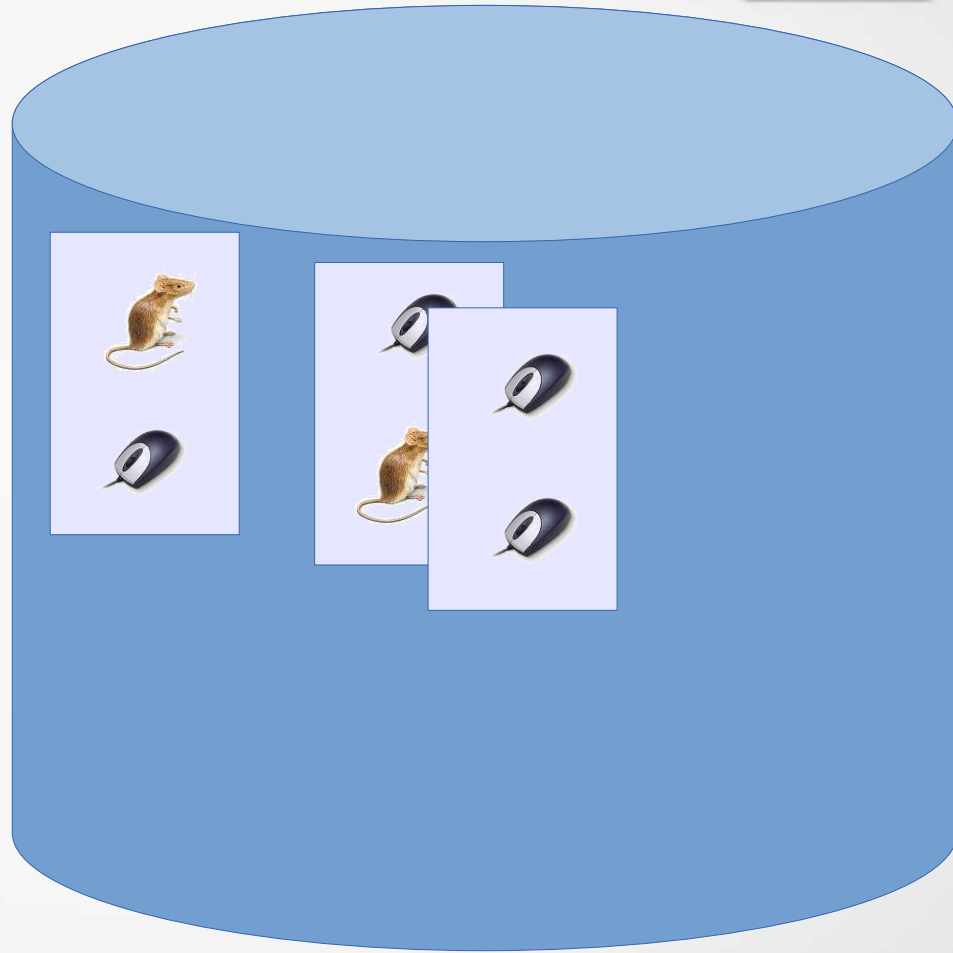
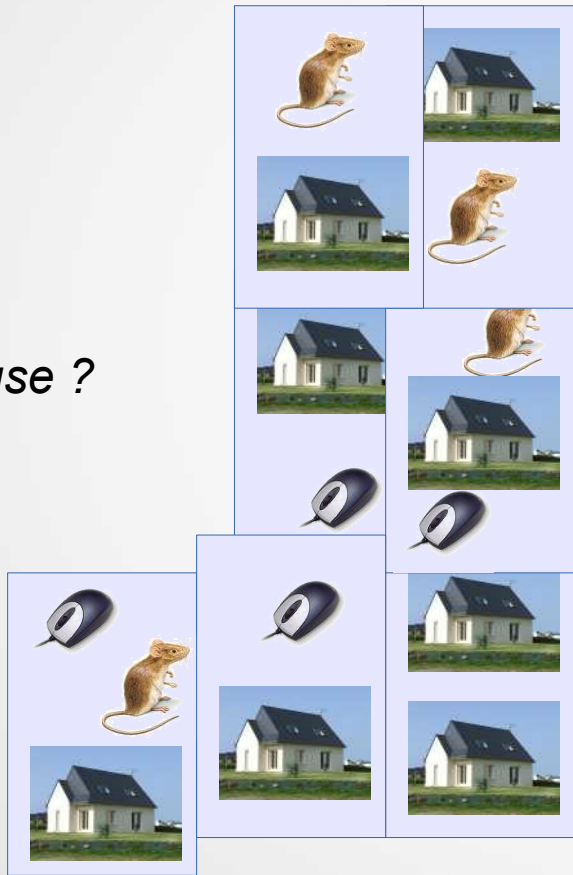
*house ?*



# WSD for Information Retrieval

Query :

*house ?*





# WSD for Information Retrieval

Query :

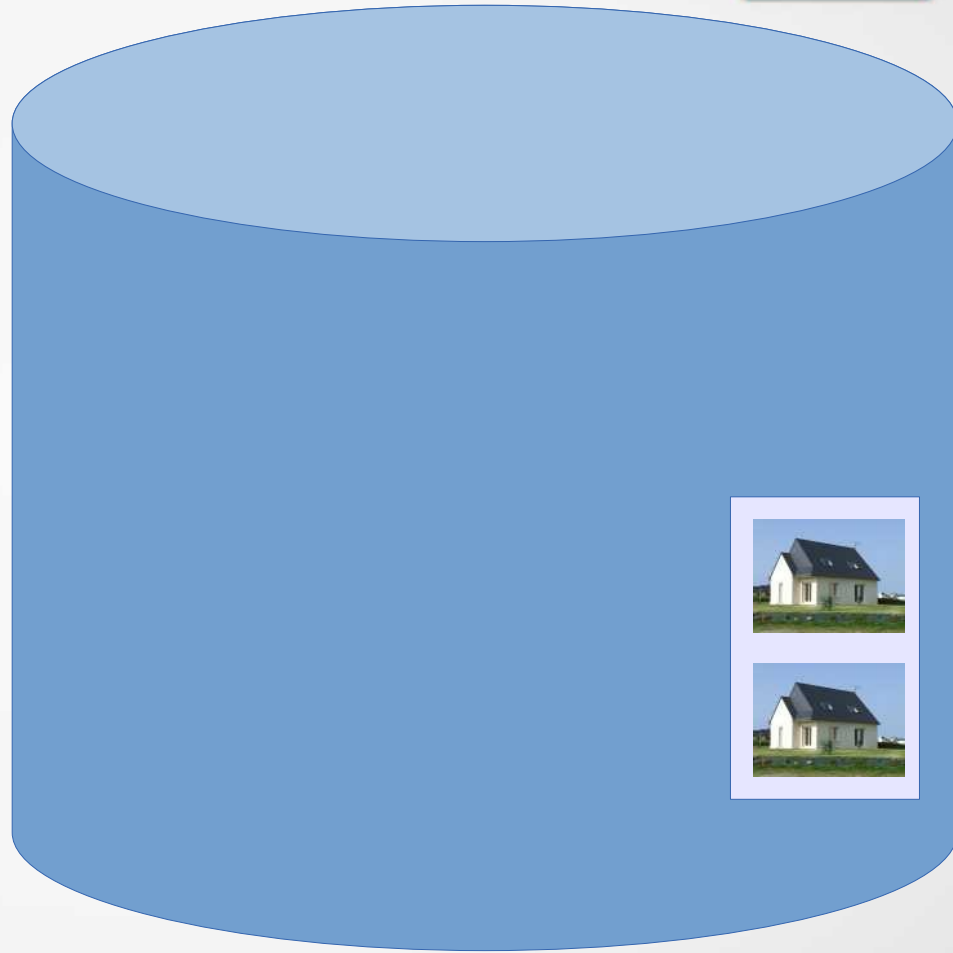
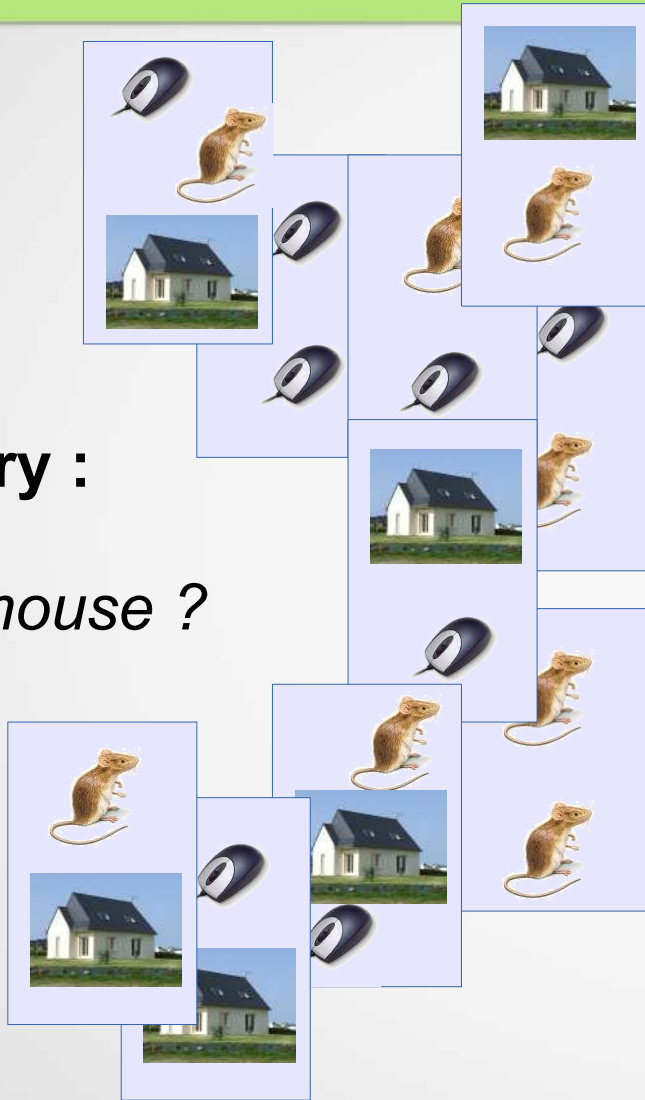
*mouse ?*



# WSD for Information Retrieval

Query :

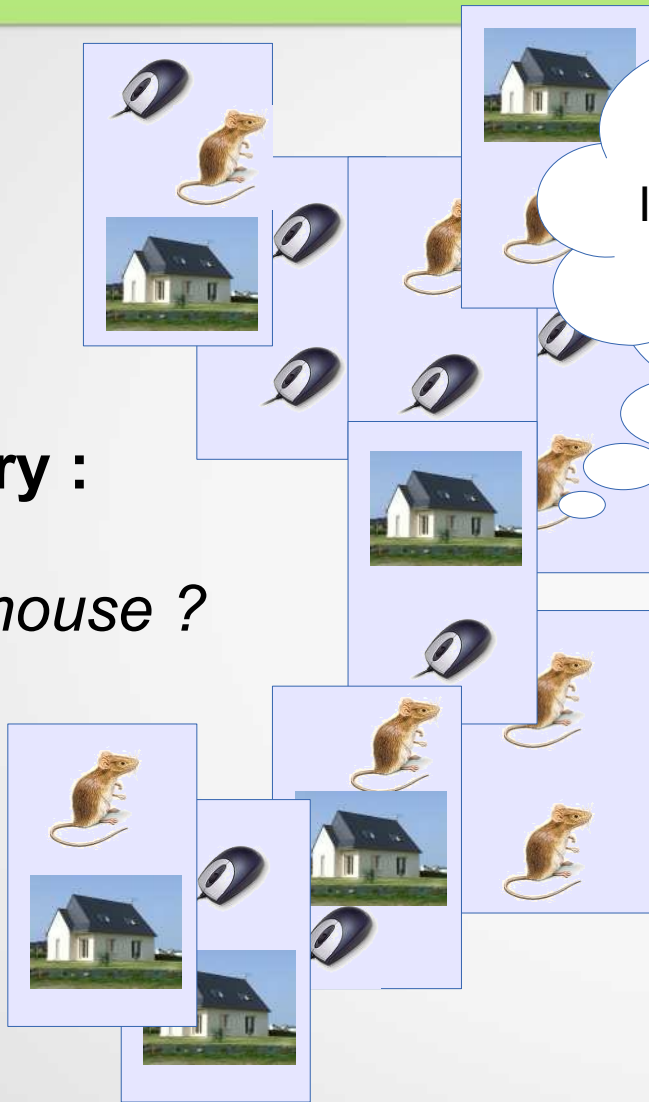
*mouse ?*



# WSD for Information Retrieval

Query :

*mouse ?*



Too much text,  
I just want information  
about rodents



# WSD for Information Retrieval

**Query :**

*mouse*  
*rodent ?*



# WSD for Information Retrieval

Query :

*mouse*  
*rodent ?*



# WSD for Question Answering

- Systems that automatically answer questions posed by humans in a natural language
- Examples :
  - Where is the Eiffel Tower ?
  - What time is it ?
  - When did George Bush enter in White House ?

# WSD for Question Answering

*When did George Bush enter in White House ?*

# WSD for Question Answering

*When did George Bush enter in office?*

**Which George Bush ?**





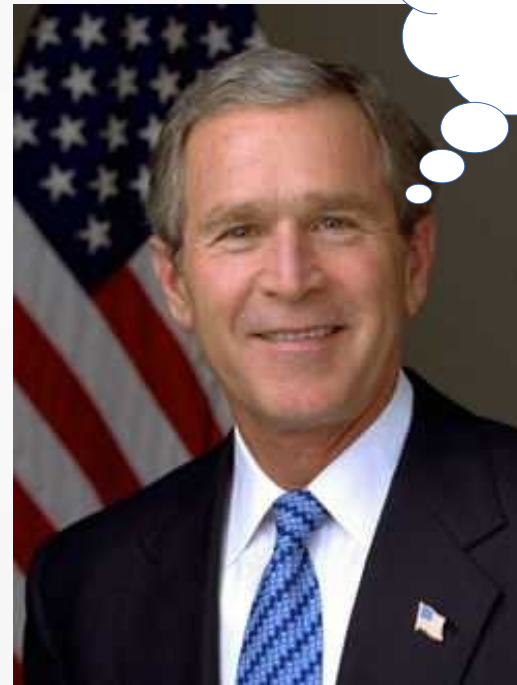
# WSD for Question Answering

When did George W. Bush enter in White House ?

W. Bush ?

1989

2001

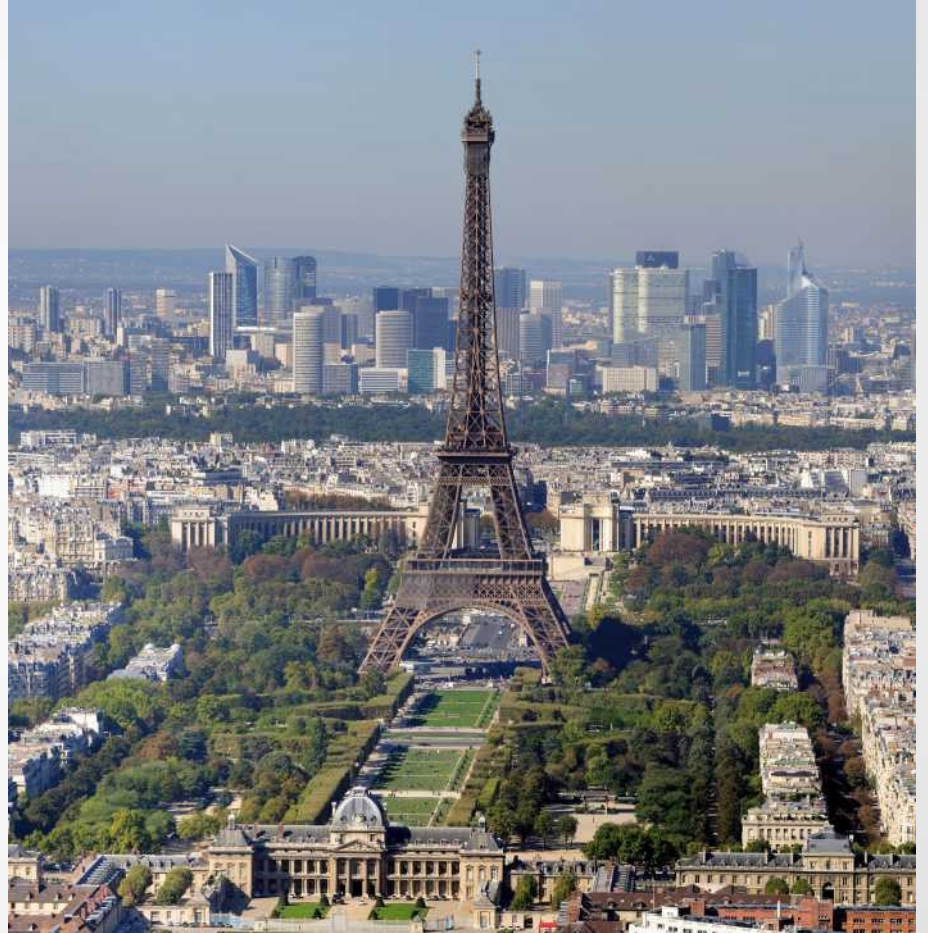


# Knowledge Acquisition

*The liberation of Paris was in 1944*



Kentucky, USA



France

# Knowledge Acquisition

*Mozart est mort à Vienne*



Austria



France

# WSD for speech synthesis

- Artificial production of human speech from written text
- Integrated in some operating systems
- Useful for:
  - Blind people
  - Mutes
  - System interaction through phones

# WSD for speech synthesis

French : fils (yarn)



[fi|]



[fis]

# Speech recognition (~WSD)

- Artificial production of text from human speech
- Homophones: Two words that sound the same but have different meanings



night

[nɪt]



knight

# Speech recognition (~WSD)



ancre

[ancre]



encre



# Evaluating Word Sense Disambiguation Performance



# Evaluation of WSD Systems

- *In vivo* evaluation
  - WSD systems evaluated through their contributions to the overall performance of a particular NLP application
  - The most natural way to evaluate
  - But the harder to set up
- *In vitro* evaluation
  - WSD task defined independently of any particular application
  - Systems evaluated using specially constructed benchmarks

# *In Vitro* Evaluation

- A benchmark : a sense-annotated corpus
- The same corpus without annotations

# Evaluation of WSD Systems

- *In vivo* evaluation (extrinsic)
  - WSD systems evaluated through their contributions to the overall performance of a particular NLP application
  - The most natural way to evaluate
  - But the most difficult to set up
- *In vitro* evaluation (intrinsic)
  - WSD task defined independently from any particular application
  - Systems evaluated using specifically constructed benchmarks

# *In Vitro* Evaluation

- A benchmark (gold-standard):reference sense-annotated corpus
- The same corpus without annotations

d001 d001.s001.t001 editorial%1:10:00:: !! lemma=editorial#n  
d001 d001.s001.t002 ill%3:00:01:: !! lemma=Ill#a  
d001 d001.s001.t003 homeless%1:14:00:: !! lemma=Homeless#n  
d001 d001.s001.t004 refer%2:42:00:: !! lemma=refer#v  
d001 d001.s001.t005 research%1:09:00:: !! lemma=research#n  
d001 d001.s001.t006 six%5:00:00:cardinal:00 !! lemma=six#a  
d001 d001.s001.t007 colleague%1:18:01:: !! lemma=colleague#n  
d001 d001.s001.t008 report%2:32:13:: !! lemma=report#v

# *In Vitro* Evaluation

- A benchmark (gold-standard):reference sense-annotated corpus
- The same corpus without annotations

d001 d001.s001.t001 editorial%1:10:00:: !! lemma=editorial#n

d001 d001.s001.t002 ill%3:00:01:: !! lemma=Ill#a

d001 d001.s001.t003 homeless%1:14:00:: !! lemma=Homeless#n

d001 d001.s001.t004 refer%2:42:00:: !! lemma=refer#v

d001 d001.s001.t005 research%1:09:00:: !! lemma=research#n

d001 d001.s001.t006 six%5:00:00:cardinal:00 !! lemma=six#a

d001 d001.s001.t007 colleague%1:18:01:: !! lemma=colleague#n

d001 d001.s001.t008 report%2:32:13:: !! lemma=report#v

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- The same corpus without annotations

First document

d001 d001.s001.t001 editorial%1:10:00:: !! lemma=editorial#n  
d001 d001.s001.t002 ill%3:00:01:: !! lemma=Ill#a  
d001 d001.s001.t003 homeless%1:14:00:: !! lemma=Homeless#n  
d001 d001.s001.t004 refer%2:42:00:: !! lemma=refer#v  
d001 d001.s001.t005 research%1:09:00:: !! lemma=research#n  
d001 d001.s001.t006 six%5:00:00:cardinal:00 !! lemma=six#a  
d001 d001.s001.t007 colleague%1:18:01:: !! lemma=colleague#n  
d001 d001.s001.t008 report%2:32:13:: !! lemma=report#v

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- The same corpus without annotations

First document

d001 d001.s001.t001 editorial%1:10:00:: !! lemma=editorial#n  
d001 d001.s001.t002 ill%3:00:01:: !! lemma=Ill#a  
d001 d001.s001.t003 homeless%1:14:00:: !! lemma=Homeless#n  
d001 d001.s001.t004 refer%2:42:00:: !! lemma=refer#v  
d001 d001.s001.t005 research%1:09:00:: !! lemma=research#n  
d001 d001.s001.t006 six%5:00:00:cardinal:00 !! lemma=six#a  
d001 d001.s001.t007 colleague%1:00:00: !! lemma=colleague#n  
d001 d001.s001.t008 report%1:00:00: !! lemma=report#v

5th term  
of the first sentence  
of the first document

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- The same corpus without annotations

First document

Solution  
(best sense)

d001 d001.s001.t001 editorial%1:10:00:: !! lemma=editorial#n  
d001 d001.s001.t002 ill%3:00:01:: !! lemma=Ill#a  
d001 d001.s001.t003 homeless%1:14:00:: !! lemma=Homeless#n  
d001 d001.s001.t004 refer%2:42:00:: !! lemma=refer#v  
d001 d001.s001.t005 research%1:09:00:: !! lemma=research#n  
d001 d001.s001.t006 six%5:00:00:cardinal:00 !! lemma=six#a  
d001 d001.s001.t007 colleague%1:00:00: !! lemma=colleague#n  
d001 d001.s001.t008 report%1:00:00: !! lemma=report#v

5th term  
of the first sentence  
of the first document

lemma



# *In Vitro* Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

```
...
```

# *In Vitro* Evaluation

- A benchmark : a reference sense-annotated corpus
- The same corpus without annotations

# *In Vitro* Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

**<text id="d001">**

<sentence id="d001.s001">

Your Oct. 6

<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>

``The

<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>

<instance id="d001.s001.t003" lemma="Homeless"

pos="n">Homeless</instance>

...

# In Vitro Evaluation

- A benchmark corpus
- The same corpus annotations

First text

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

```
...
```

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

Your Oct. 6

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

``The

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

...

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus annotations**

First sentence  
of the first text

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

Your Oct. 6

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

"The

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

...

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

**Your Oct. 6**

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

**The**

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

...

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same sense-annotations**

- Raw Texts Unevaluated parts

Your Oct. 6 editorial "Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

Your Oct. 6

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

``The

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

...



# *In Vitro* Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The Ill Homeless" referred to research by us and six of our colleagues that was reported in the Sept. 8 issue of the Journal of the American Medical Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

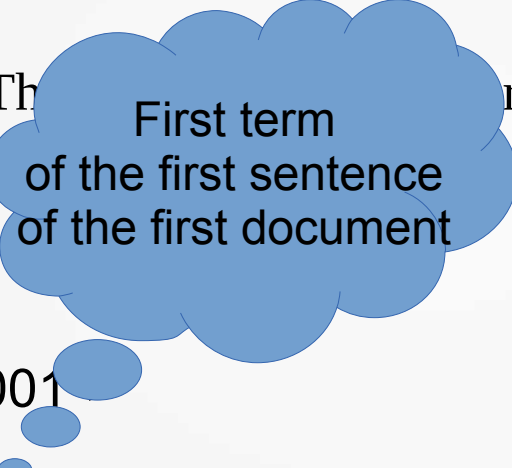
```
...
```

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The ... rred to research by us and six of our colleagues that was reported in the Journal of the American Medical Association .



First term  
of the first sentence  
of the first document

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

```
...
```

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The ... rred ... us and six of our colleagues that was reported ... of the ... American Medical Association .

First term  
of the first sentence  
of the first document

lemma

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001"
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

```
...
```

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The ... rred ... us and six of our colleagues that was reported ... of the ... American Medical Association .

First term  
of the first sentence  
of the first document

lemma

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001">
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

```
...
```

Part  
of speech

# In Vitro Evaluation

- A benchmark : a sense-annotated corpus
- **The same corpus without sense-annotations**

- Raw Texts

Your Oct. 6 editorial "The ... rred ... us and ... ur  
colleagues that was repo ... of the ... An ...  
Association .

- Texts

```
<text id="d001">
```

```
<sentence id="d001.s001"
```

```
Your Oct. 6
```

```
<instance id="d001.s001.t001" lemma="editorial" pos="n">editorial</instance>
```

```
``The
```

```
<instance id="d001.s001.t002" lemma="Ill" pos="a">Ill</instance>
```

```
<instance id="d001.s001.t003" lemma="Homeless"
```

```
pos="n">Homeless</instance>
```

```
...
```

First term  
of the first sentence  
of the first document

lemma

word

Part  
of speech

## *In Vitro* Evaluation : metrics

$$\textit{precision} = \frac{\textit{words correctly tagged}}{\textit{tagged words}}$$

$$\textit{recall} = \frac{\textit{words correctly tagged}}{\textit{words}}$$

$$\textit{F - measure} = \frac{2 \times \textit{precision} \times \textit{recall}}{\textit{precision} + \textit{recall}}$$

If all words are tagged

$$P = R \rightarrow \textit{F - measure} = \frac{2 \times P \times P}{P + P} = \frac{2 \times P^2}{2 \times P} = P$$

## *In Vitro* Evaluation : metrics

$$\textit{precision} = \frac{\textit{words correctly tagged}}{\textit{tagged words}}$$

$$\textit{recall} = \frac{\textit{words correctly tagged}}{\textit{words}}$$

$$\textit{F - measure} = \frac{2 \times \textit{precision} \times \textit{recall}}{\textit{precision} + \textit{recall}}$$

If all words are tagged

$$P = R = \textit{F - measure}$$

# *In Vitro* Evaluation : example

- Example :
  - 100 words to tag
  - The system tags 75 words
  - 50 are correctly tagged
  - Precision :  $50/75 = 66\%$
  - Recall :  $50/100 = 50\%$
  - F-measure  $\approx 56.9\%$



# Bounds of performance

- Evaluating performance of an algorithm relative to the difficulty of the benchmark
- Lower bound (baseline)
  - random assignment: average score obtained when a random sense is chosen for each words in the text

$$\text{random baseline} = \frac{1}{n} \sum_{i=1}^n \frac{1}{|\text{senses}(w_i)|}$$

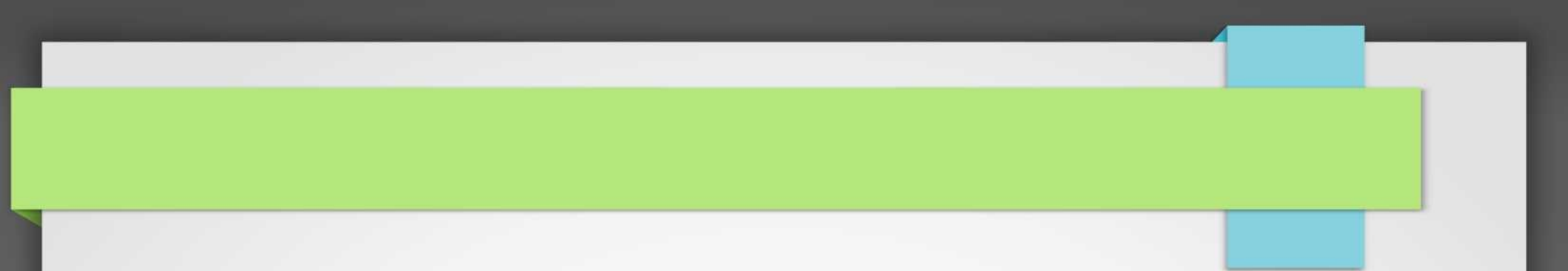
- most frequent sense: score when the most frequent sense in the language is chosen for each word in the text
- Upper bound
  - Highest performance reasonably attainable
  - Average human interannotator agreement : Around 90%

# Example: Semeval 2007 task 7

- All-words task: sense labelling task over all parts-of-speech (nouns, verbs, adjectives, adverbs)
- 2269 words over 5 texts: journalism, book review, travel, computer science, biography
- Disambiguated reference tagged with WordNet senses  
Evaluation in terms of Precision, Recall, F1 score
- Currently the most recent general English All-words disambiguation task available.

# Example: semeval 2007 task 7

- Coarse-grained evaluation : close senses are counted as equivalent (e.g. snow/precipitation and snow/cover)
- Two ways to use this benchmark
  - *A Posteriori*
    - Input: fine-grained (WordNet Senses)
    - Random baseline: 61,27%
    - First sense baseline: 78,89%
  - *A priori*
    - Input: coarse-grained
    - Random baseline: 52,57%
    - First sense baseline: 78,89%



# General Overview of Word Sense Disambiguation Systems

# Word Sense Disambiguation Process

- Composed of 3 steps
  - Build/select raw lexical material(s)
  - Build an elaborate resource
  - Use that resource to lexically disambiguate a text

# Build/Select of Raw Lexical Material(s)

- One or more of several types of materials can be used:
  - Dictionaries, encyclopedias, lexical databases
  - Unannotated corpora, Sense-annotated corpora
- Among existing material, some:
  - Are generated/built automatically
  - Require significant human effort and supervision

# Build an elaborate resource

- Computational representation of an inventory of possible word senses
- Two ways of obtaining inventories of word senses:
  - Induction from word contexts
    - When only non-annotated corpora are available
  - Human experts
    - e.g. Dictionaries, Structured Lexical Resources
- Many underlying computational representations:
  - Semantic Networks (graphs)
  - Bags of words & n-gram models
  - Vector spaces

# Use the resource to disambiguate

- The Word Sense Disambiguation algorithm
  - More or less complex
  - SVMs, Naive Bayes, Deep Neural Network, etc.
  - PageRank, Ant Colony algorithms, genetic algorithms, etc.
- Several common parameters are involved:
  - Context : window, phrase, sentence, text,...
  - Depth in a graph



# Resources

- In WSD, we consider two kinds of resources
  - Knowledge
    - Machine readable dictionaries
    - Lexical Databases
    - Encyclopedias
  - Corpus
    - Non-sense-annotated corpus
    - Sense-annotated corpus

# Resources : knowledge

- Machine readable dictionaries
  - Longman, Oxford Advanced Learner's dictionary,...
  - Until the 1990's for English
- Lexical Databases
  - WordNet from the 1990's [Miller]
  - BabelNet [Navigli, 2012]
- Encyclopedias
  - Wikipedia from 2007 [Mihalcea, 2007]

# Resources: non-sense-annotated corpora

- A set of texts
- Covers one or more domains
- One or more languages
- Up to dozens of millions of words
- Can be lemmatized and tagged with part of speech information
- Various sources :
  - Newspapers, books, encyclopedias, Web,...

# Resources: sense-annotated corpora

- SemCor [Miller et al., 1993]
- Subset of the Brown Corpus (1961)
  - 700,000 words
  - 30,000 words manually tagged with Wordnet synsets
  - 352 texts
    - For 186 texts, nouns, verbs, adjectives, and adverbs tagged : 192,639 words
    - For 166, only verbs are tagged : 41,497 words

# Resources: sense-annotated corpora

- The Defense Science Organisation corpus [Ng & Lee, 1996]
  - Non-freely available sense- annotated English corpus
  - 192800 word occurrences manually tagged with WordNet synsets
  - Annotations cover
    - 121 nouns (113,000 occurrences)
    - 70 verbs (79,800 occurrences)
  - The most frequent, as ambiguous possible.
  - Coverage corresponding to 20% of verb and noun occurrences in English texts

# Resources: Sense-annotated corpora

- Corpora from evaluation campaigns
  - Most of them in English
  - But also on Japanese, Spanish, Chinese
  - Uncommonly beyond 5000 tagged words
- Other languages:
  - Dutch SemCor [Vossen et al., 2012]
    - 250,000 manually tagged words
  - Basque SemCor [Agirre, 2006]

# Sense-annotated corpora : limitations

- Really difficult task compared to other annotation tasks
- Penn Treebank [Taylor et al., 2003]
  - Part of speech tagged corpus
  - Only 45 possible tags
  - 3000 annotations per hour
- WordNet synset-annotated corpus
  - 117,000 possible tags
  - Example for the Defense Science Organisation corpus
    - 191 different nouns, 1800 possible tags
    - 1 man-year for 192000 word occurrences 150-250 annotations per hour

# Sense-annotated corpora : limitations

- Have to be repeated for
  - each sense inventory;
  - each language;
  - each domain.
  - ...
  - With updated corpus (new senses, new words...).

Ex : mouse in SemCor based on the Brown Corpus (1961)



# Mitigating the limitations

- Improving annotation speeds
  - [Mihalcea & Chklovski, 2003] WSD algorithm on corpus -  
> Then human verification
  - Not much improvement
- Usage of new kinds of sense-annotated corpora
  - E.g. Wikipedia and its internal links [Mihalcea, 2007]
  - A page can be considered as a sense
- More languages
  - BabelCor

# UFSAC: Unification of Sense Annotated Corpora and Tools [Vial et al., 2018]

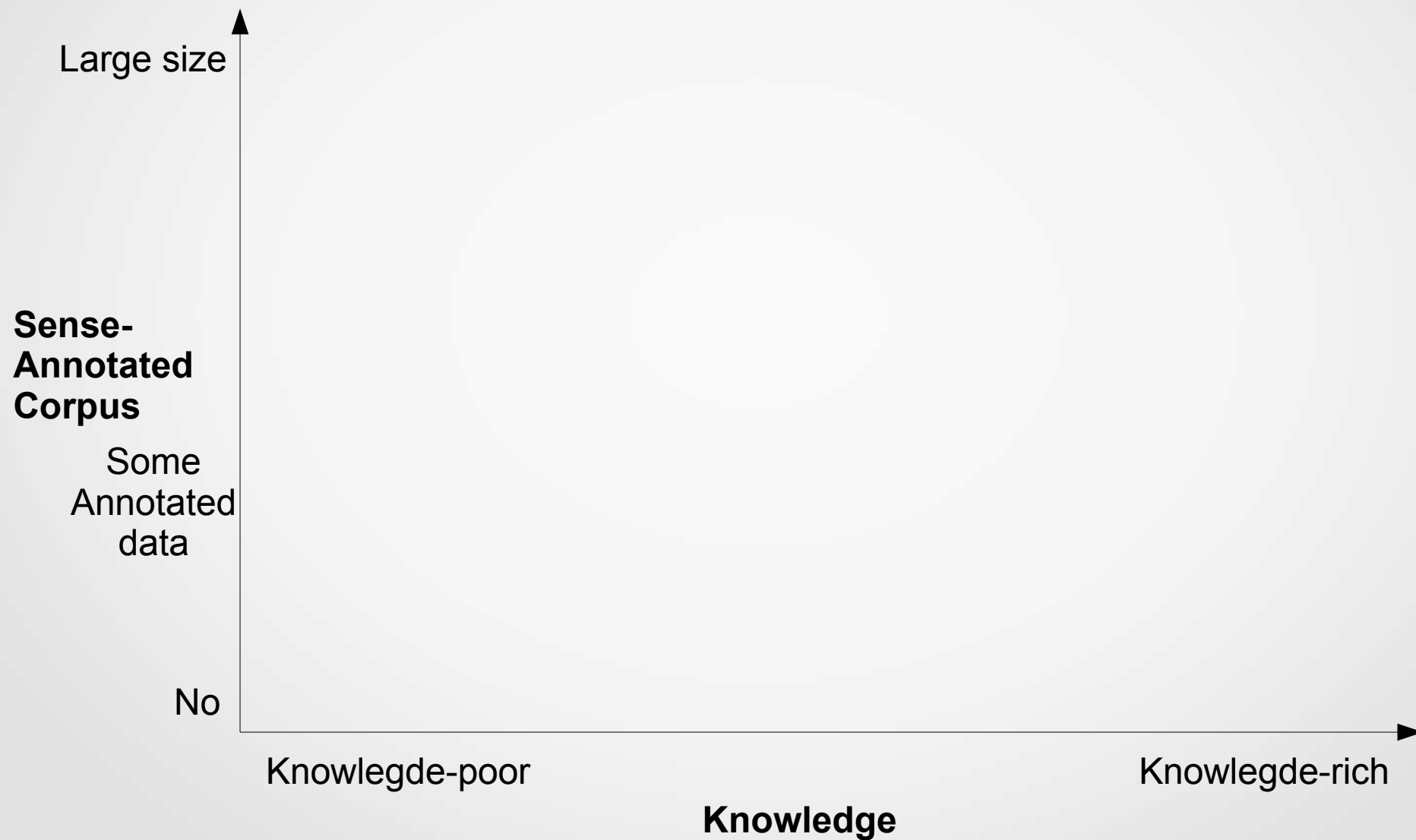
- In English, there are a dozen of manually annotated sense annotated corpora, but their file formats are very different from one another.
- Unification of these corpora in a format
  - easy to use
  - Easy to understand
- Facilitate
  - the creation of new WSD systems
  - the evaluation of existing ones

<https://github.com/getalp/UFSAC>

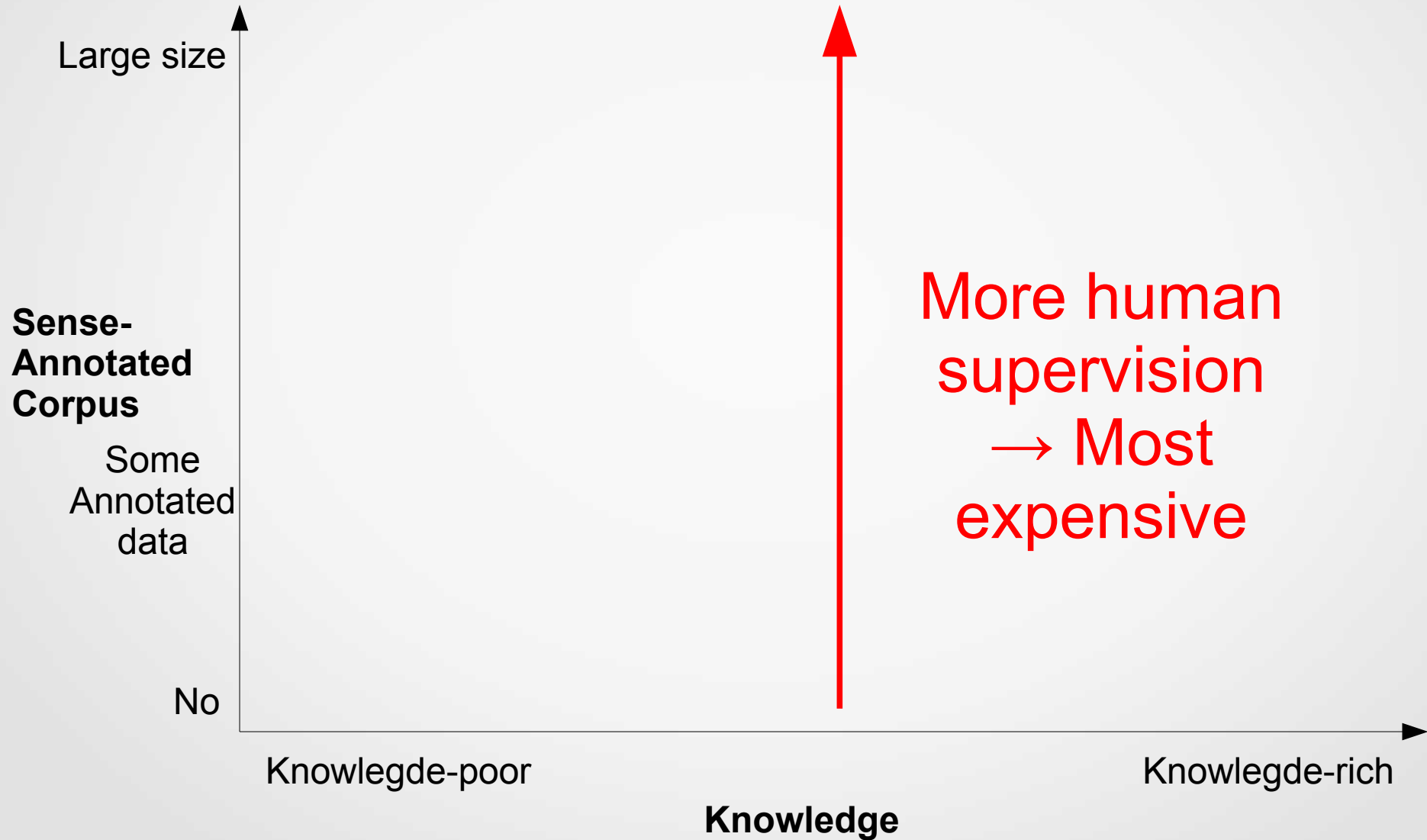
# UFSAC: Unification of Sense Annotated Corpora and Tools [Vial et al., 2018]

Corpus	Sentences	Words		Annotated parts of speech			
		Total	Annotated	Nouns	Verbs	Adj.	Adv.
SemCor [7]	37176	778587	229517	87581	89037	33751	19148
DSO [11]	178119	5317184	176915	105925	70990	0	0
WordNet GlossTag [6]	117659	1634691	496776	232319	62211	84233	19445
MASC [4]	34217	596333	114950	49263	40325	25016	0
OMSTI [14]	820557	35843024	920794	476944	253644	190206	0
Ontonotes [3]	21938	435340	52263	9220	43042	0	0
Senseval 2 [2]	238	5589	2301	1061	541	422	277
Senseval 3 task 1 [13]	300	5511	1957	886	723	336	12
SemEval 2007 task 07 [10]	245	5637	2261	1108	591	356	206
SemEval 2007 task 17 [12]	120	3395	455	159	296	0	0
SemEval 2013 task 12 [9]	306	8142	1644	1644	0	0	0
SemEval 2015 task 13 [8]	138	2638	1053	554	251	166	82

# Analysis of resources for WSD



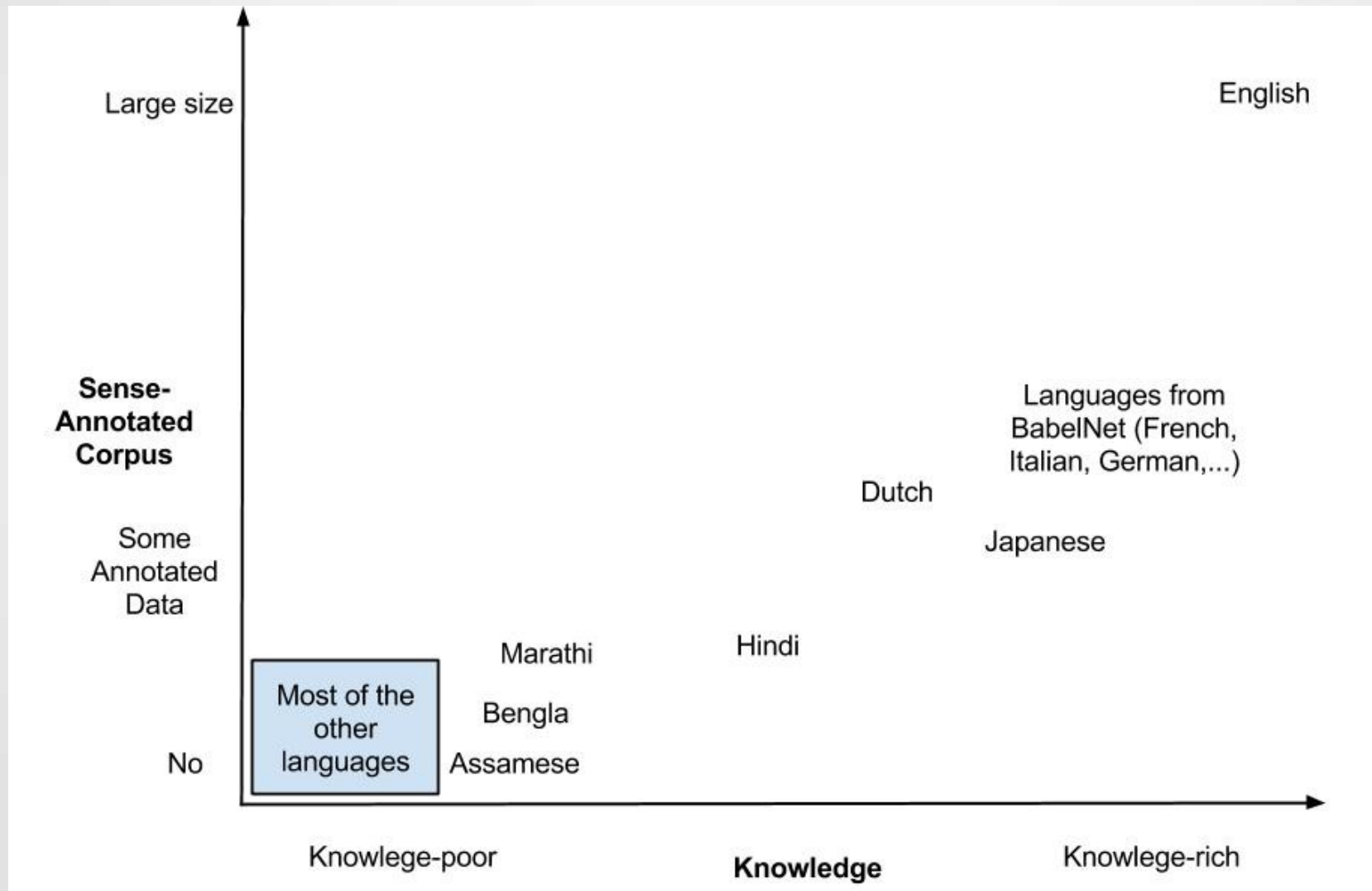
# Analysis of resources for WSD



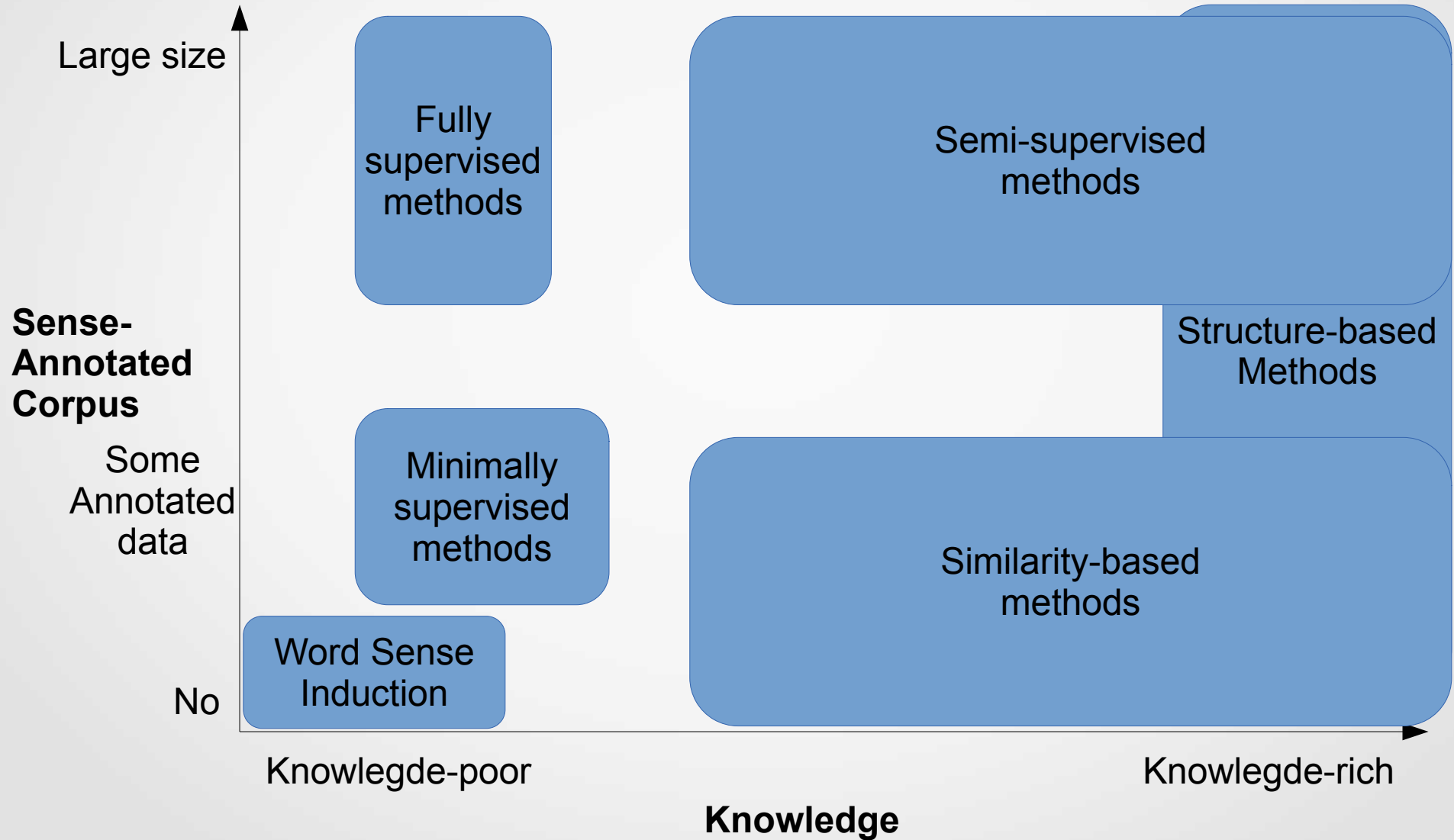
# Analyse of resources for WSD



# Languages Resources Available

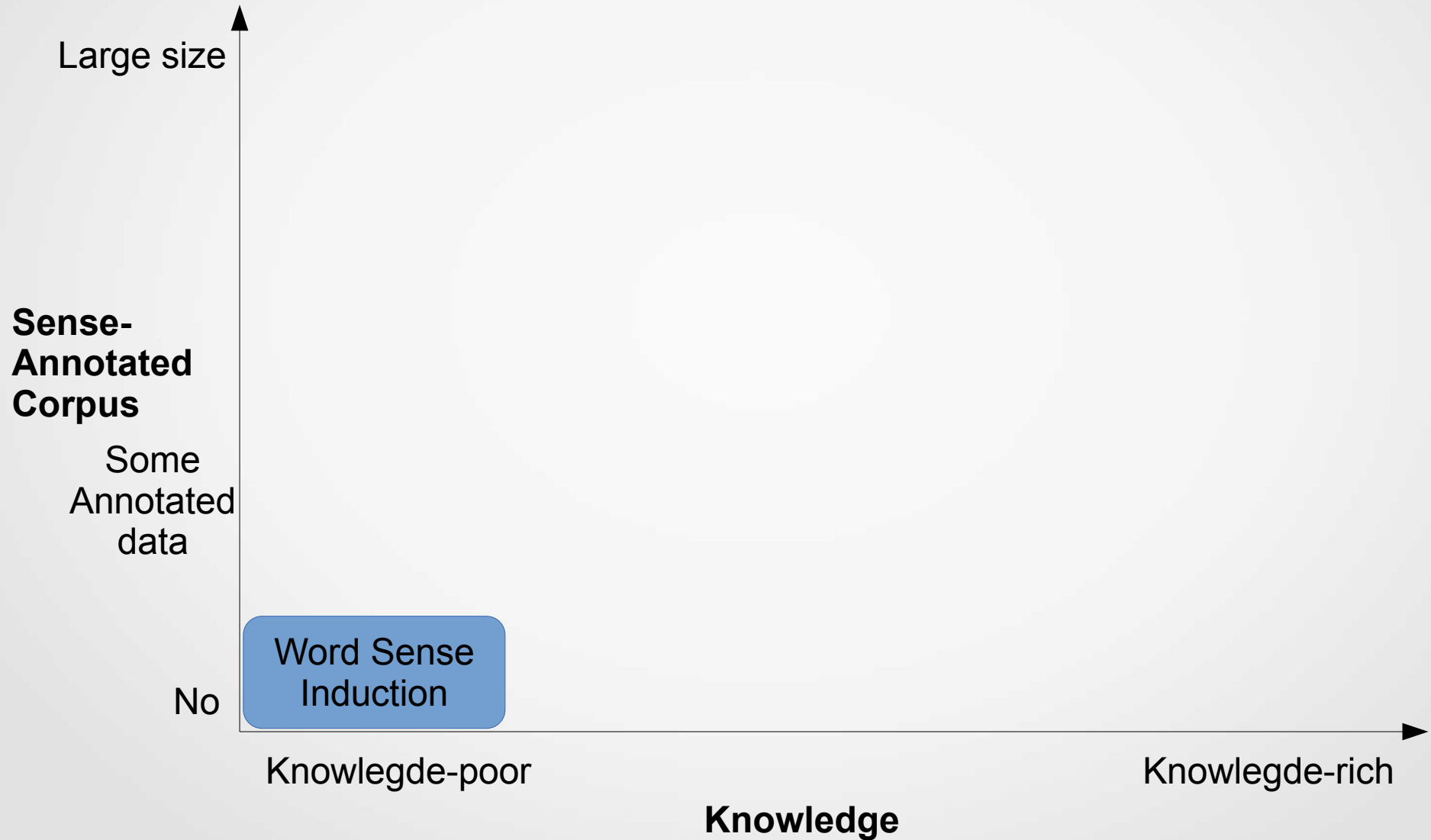


# WSD Approaches





# WSD Approaches



# Word Sense Disambiguation Process

- Composed of 3 steps
  - Build/select of raw lexical material(s)
  - Build an elaborate resource
  - Use that resource to lexically disambiguate a text

# Word Sense induction (WSI)

- Word Sense induction (or discrimination)
- Build/select raw lexical material(s)
  - Only raw (no sense annotations) corpora
- Build an elaborate resource
  - Induce word senses from contexts
- Use that resource to lexically disambiguate a text
  - Open

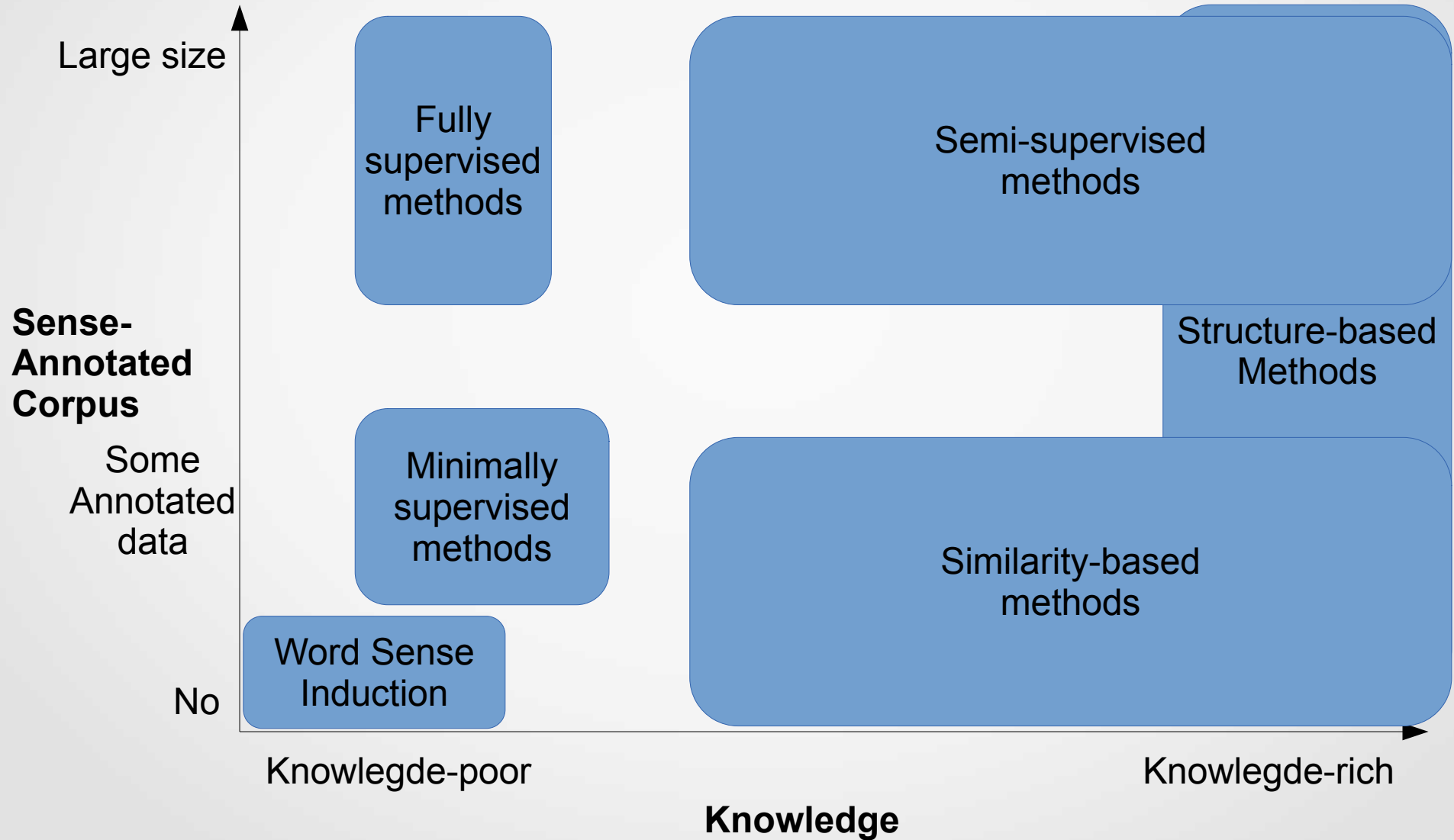
# WSI : Build an elaborate resource

- Use only raw corpora
- Induce word senses from contexts
- Harris' (1954) Distributional semantics principle -
  - Hypothesis : the meaning of a word comes from its context
- Example:
  - „The mouse is eating cheese“, „The cat is hunting a mouse“
  - „The mouse is linked to the computer“, „my mouse is broken“

# WSI : Build an elaborate resource

- Induce word senses from input text by clustering word occurrences
- Computational representation:
  - Vectors, Bag of words
- Clustering algorithms : Kmean,...
- Graphs: each node is a word and edges are coocurences, senses are given by identification of hubs (clusters)

# WSD Approaches



# Useful heuristics

- Based on observations
- One sense per discourse [Gale *et al.*, 1991]
- One sense per collocation [Yarowsky, 1993]

# One sense per discourse [Gale et al., 1991]

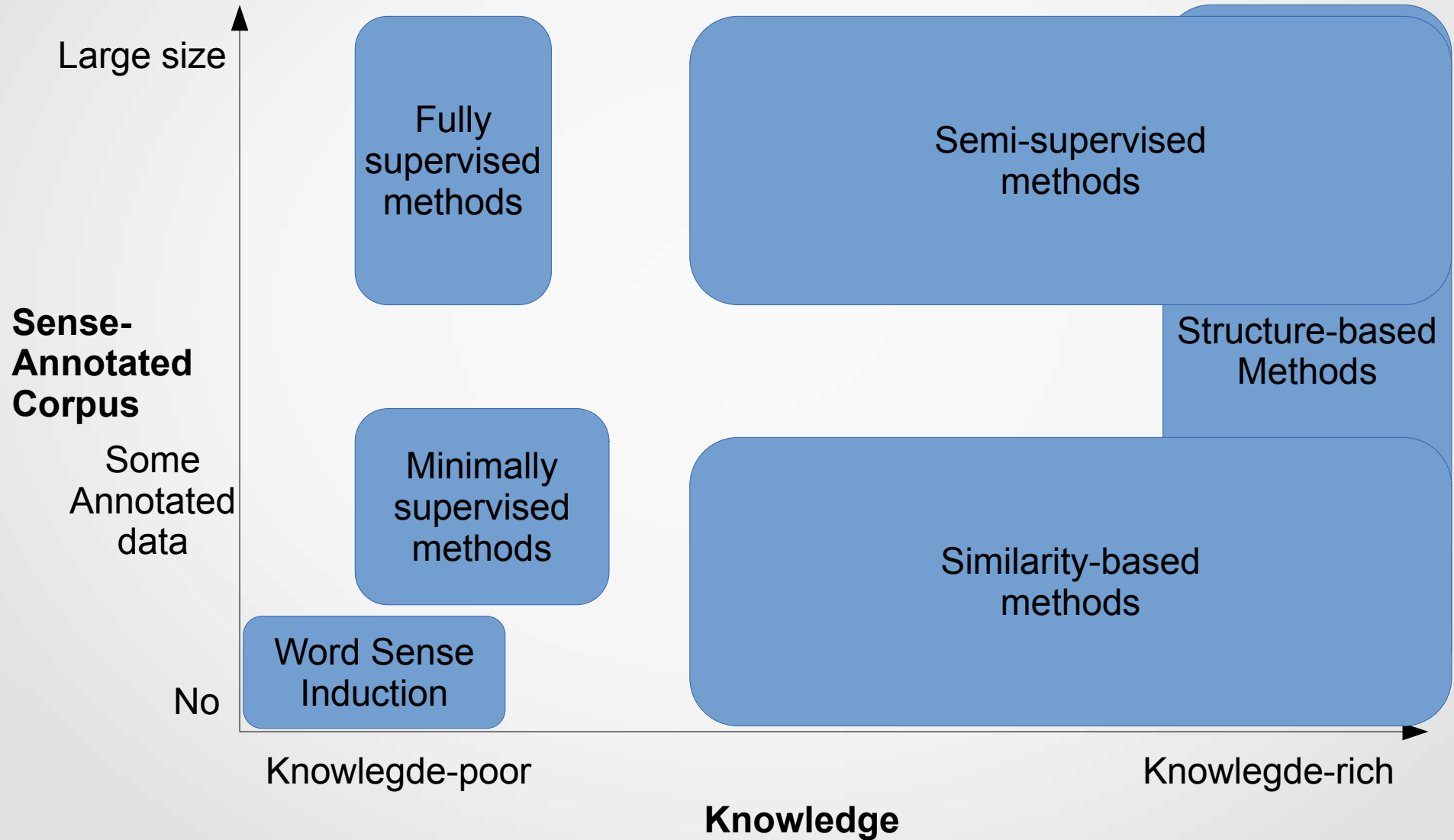
- Random sample of 108 nouns
- 300 articles studied
- 3 judges
- Only 6 articles judged to contain multiple senses of one of the test words
- Tendency to share senses in the same discourse extremely strong: 98%



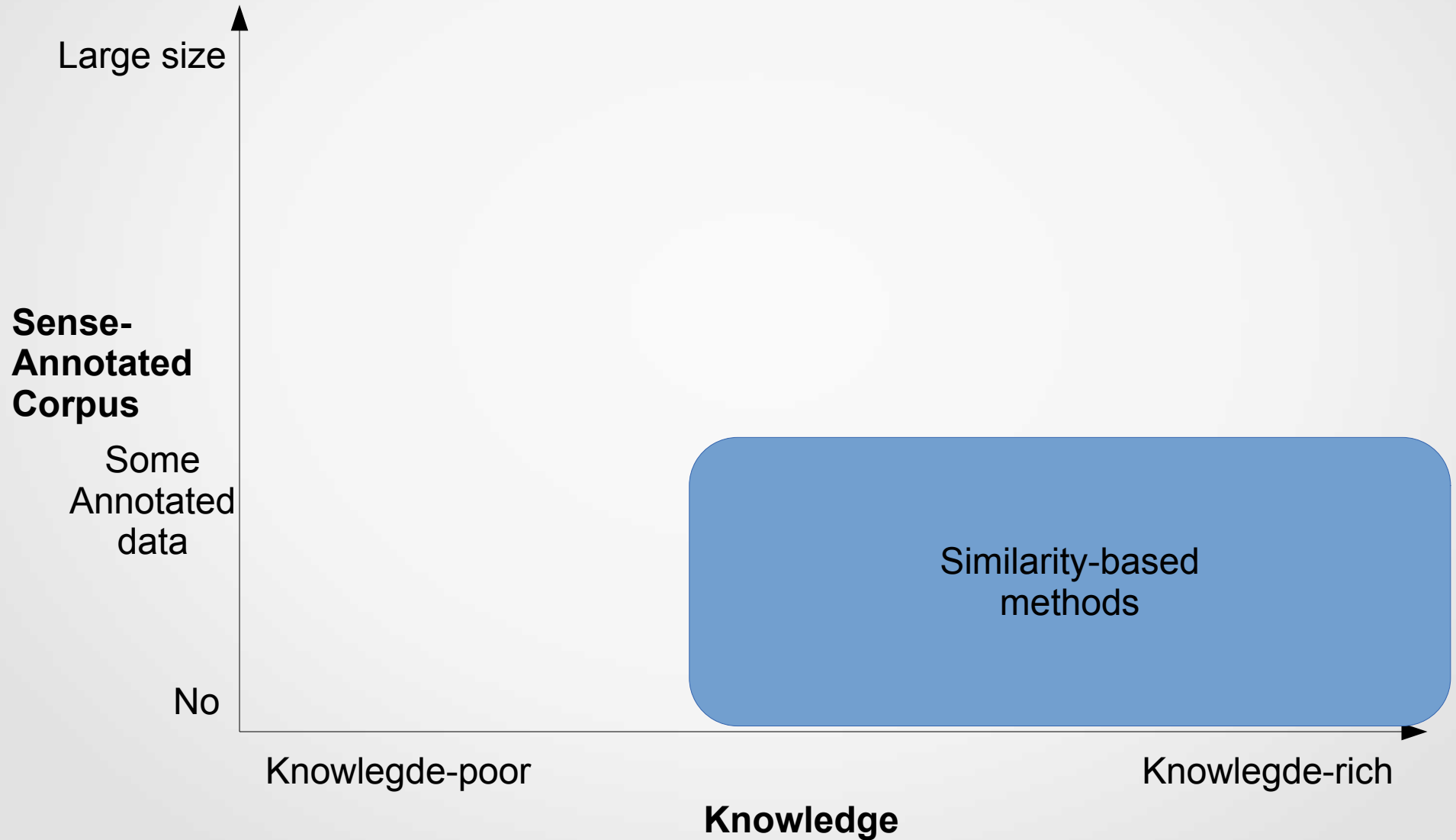
# One Sense per Collocation [Yarowsky, 1993]

- Collocation : sequence of words or terms that co-occur more often than would be expected by chance
- Types of collocations:
  - adjective+noun : *peur bleue, strong fever*
  - noun+noun (such as collective nouns): *meute de loups, douzaine d'œufs, wolf pack, dozen egg*
  - verb+noun: *prendre une gifle, prendre l'escalier, chair a meeting, conduct an experiment*
- 90% to 99% for an average of 95% share senses in texts

# WSD Approaches



# WSD Approaches



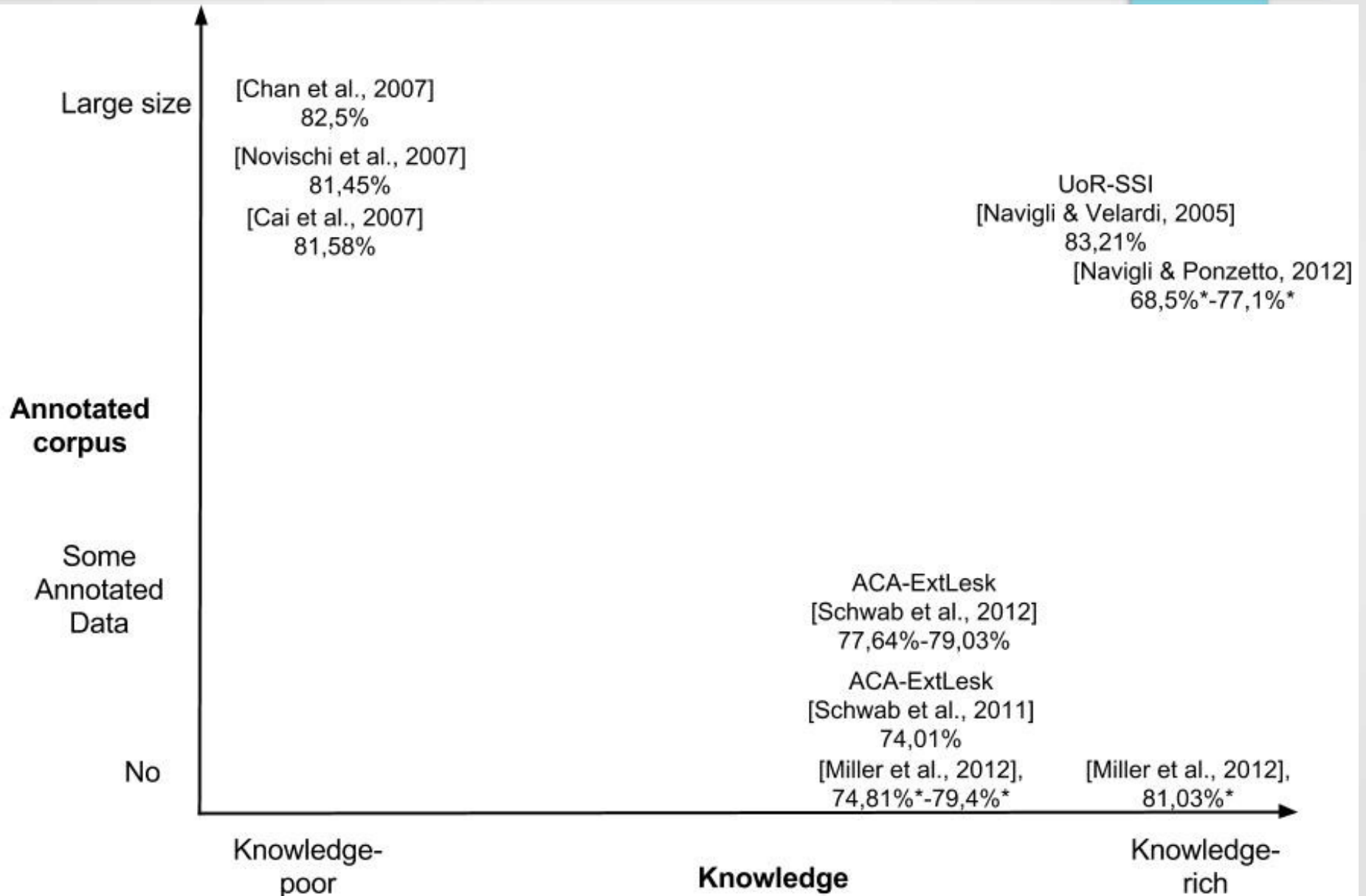
# Word Sense Disambiguation Process

- Composed of 3 steps
  - Build/select raw lexical material(s)
  - Build an elaborate resource
  - Use that resource to lexically disambiguate a text

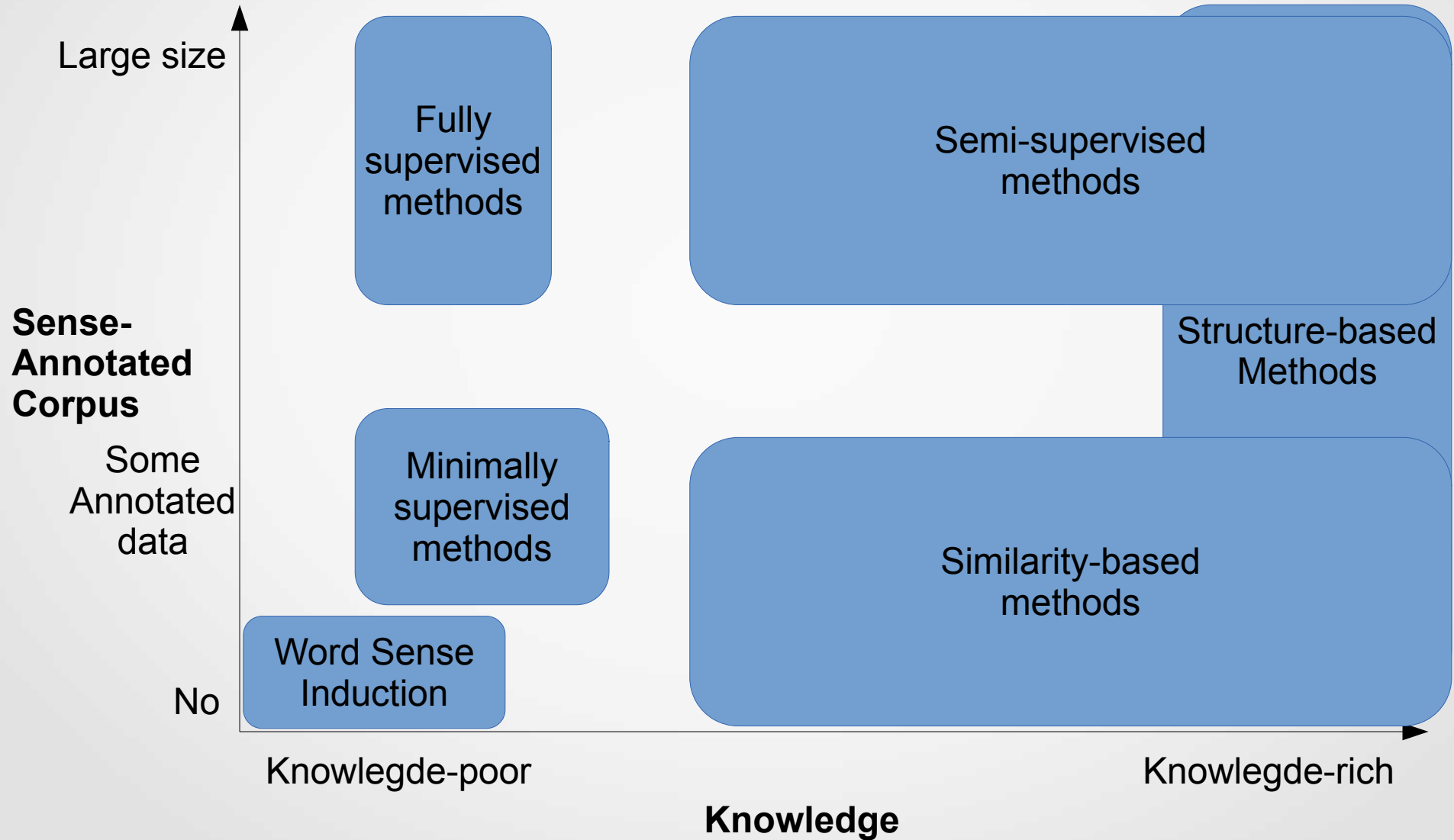
# Word Sense Disambiguation Process

- Composed of 3 steps
  - Build/select of raw lexical material(s)
    - **Mandatory: MRD or Lexical Base**
    - **Optional: corpus (sense-annotated or not)**
  - Build an elaborate resource
    - **Various ways to construct**
  - Use that resource to lexically disambiguate a text
    - **Local algorithm : semantic relatedness between senses**
    - **Global algorithm : Various**

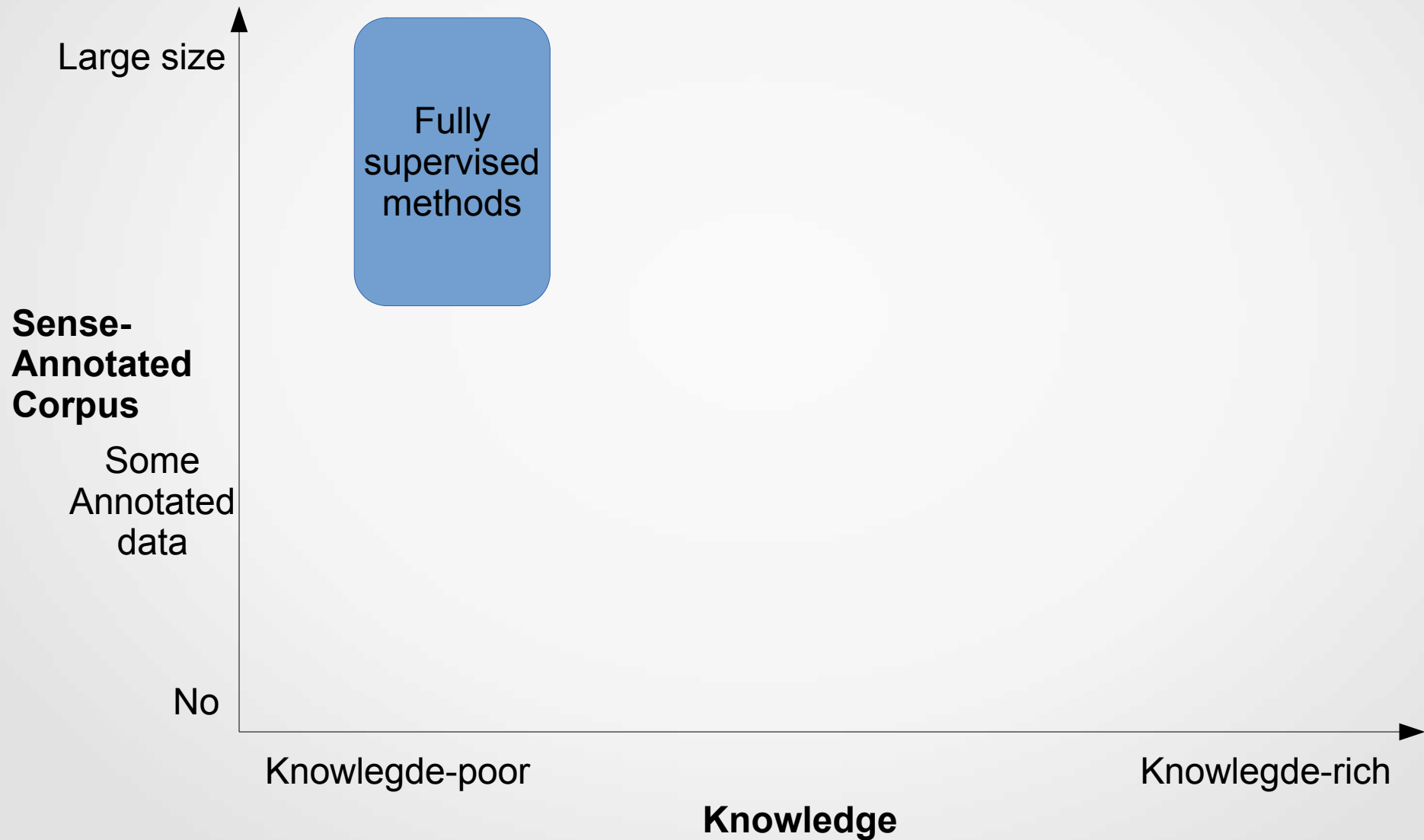
# Semeval 2007 map



# WSD Approaches



# WSD Approaches





# Word Sense Disambiguation Process

- Composed of 3 steps
  - Build/select raw lexical material(s)
  - Build an elaborate resource
  - Use that resource to lexically disambiguate a text

# Supervised WSD

- Build/select raw lexical material(s)
  - Only using sense annotated corpus/corpora
- Build an elaborate resource
  - Learn one classifier per word
- Use that resource to lexically disambiguate a text
  - Use classifiers to find the best sense for each word in texts

# Supervised Word Sense Disambiguation

- Machine Learning techniques
- Learn classical classifiers on sense-tagged corpora
  - Support Vector Machines NUS-PT, (Chan et al., 2007)
  - Naïve Bayes NUS-ML, (Cai et al., 2007)
  - Maximum Entropy / Support Vector Machines LCC-WSD, (Novischi et al., 2007)
- One classifier per word

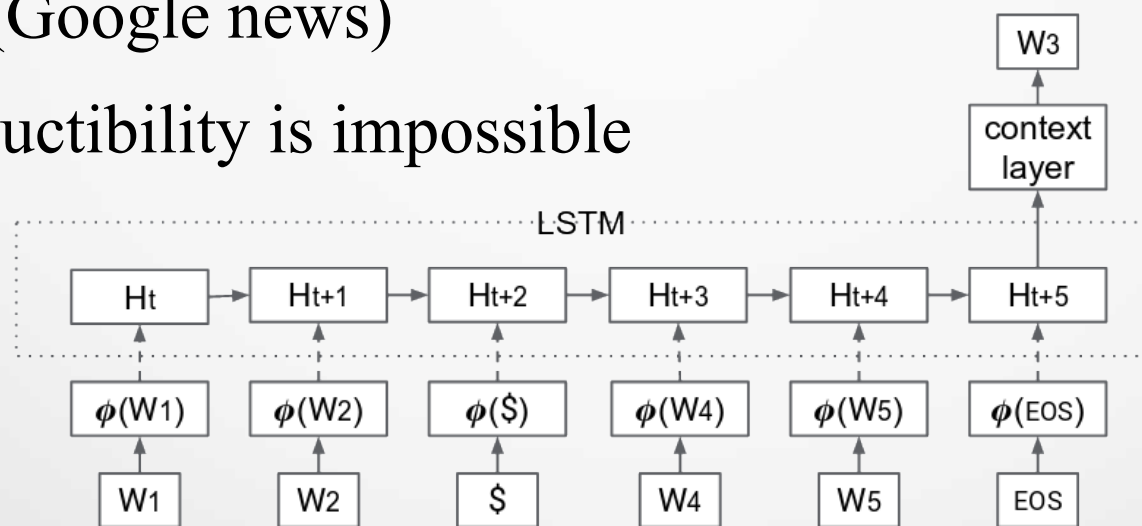
=> **state of the art on WSD 2007 -> 2016**

# Deep Neural Networks

- 2016 → ...
- [Yuan et al., 2016]
- [Raganato et al., 2017]
- [Vial et al., 2018]
- [Vial et al., 2019]

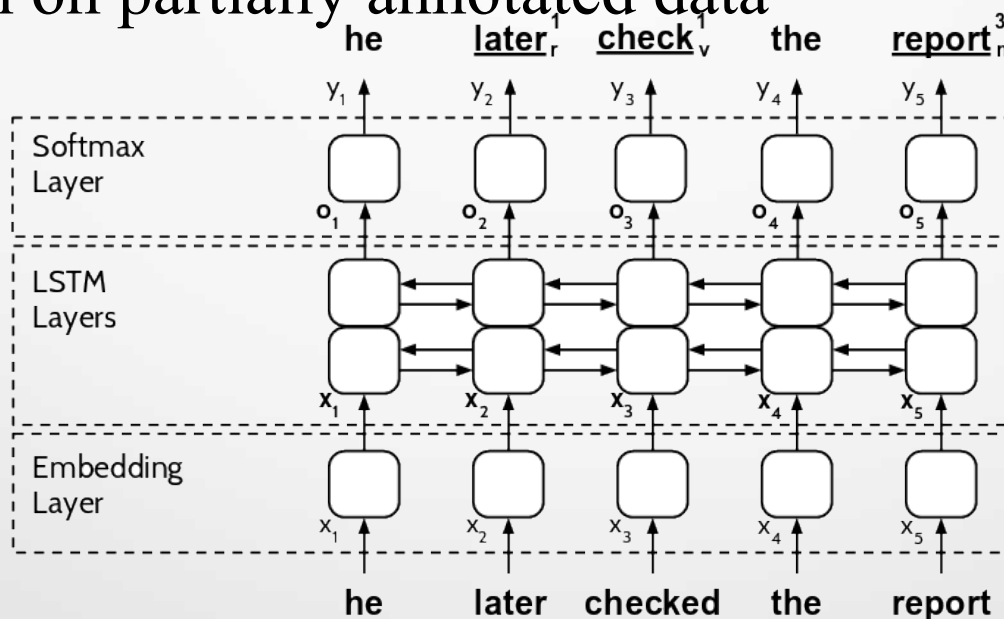
# [Yuan et al., 2016]

- LSTM language Model (Long Short-Term Memory)
- Give a prediction for a target word (classification)
- Closest sense is assigned
- Language model learned on a private corpus of 100 billions words (Google news)
- Reproducibility is impossible



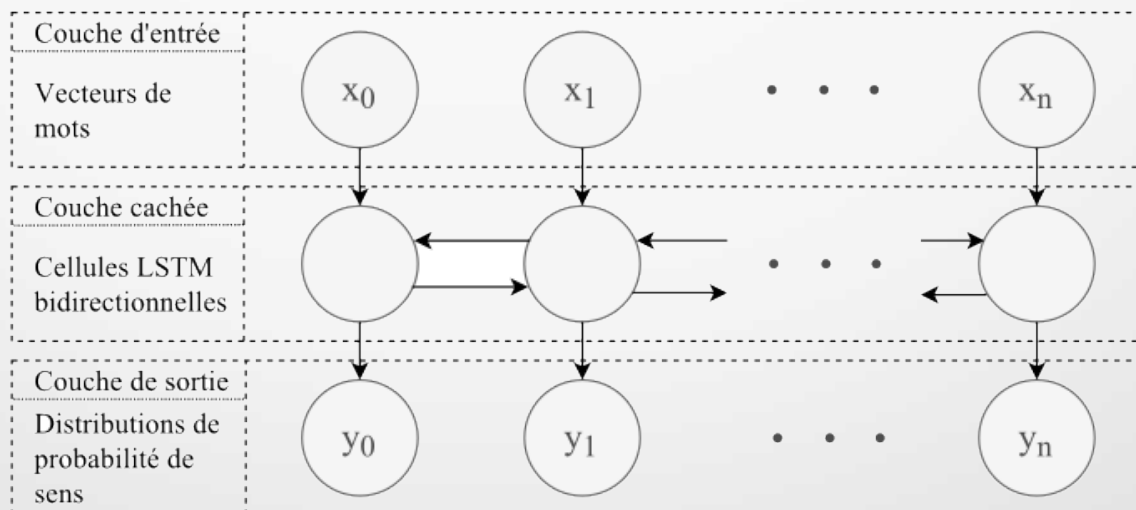
# [Raganato et al., 2017]

- Directly predict sense for each word
- Predict word when no sense can be assigned
- Multi-task learning (POS + WSD)
- Reproducibility is possible
- Can't learn on partially annotated data

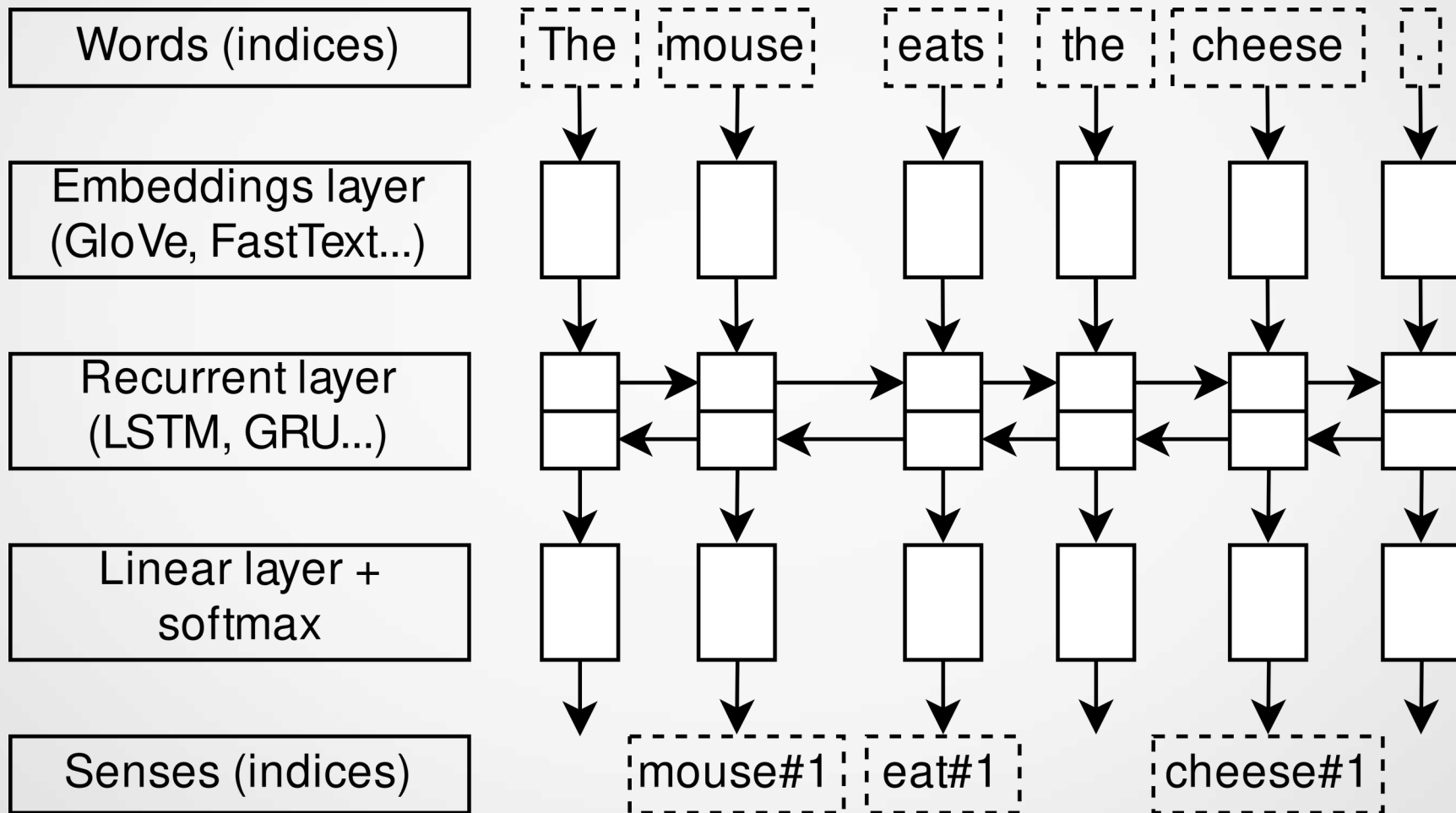


# [Vial et al., 2018]

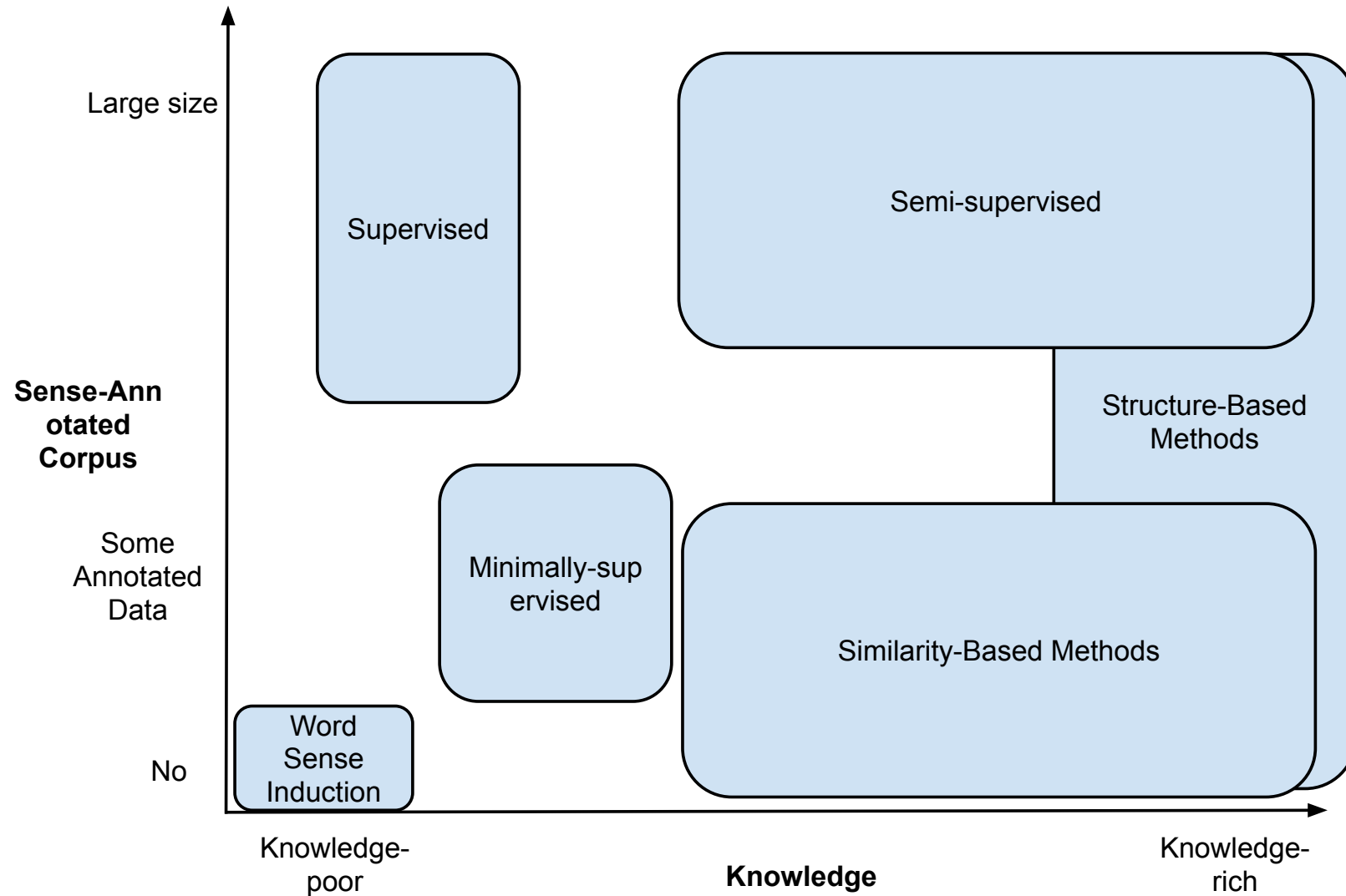
- Input layer : pre-trained vectors (Glove (Pennington et al., 2014))
- Hidden layer : Bidirectional LSTM (size : 1000)
- Output layer : size number of senses ( $\sim 100\ 000$ )
- Dropout : 50%

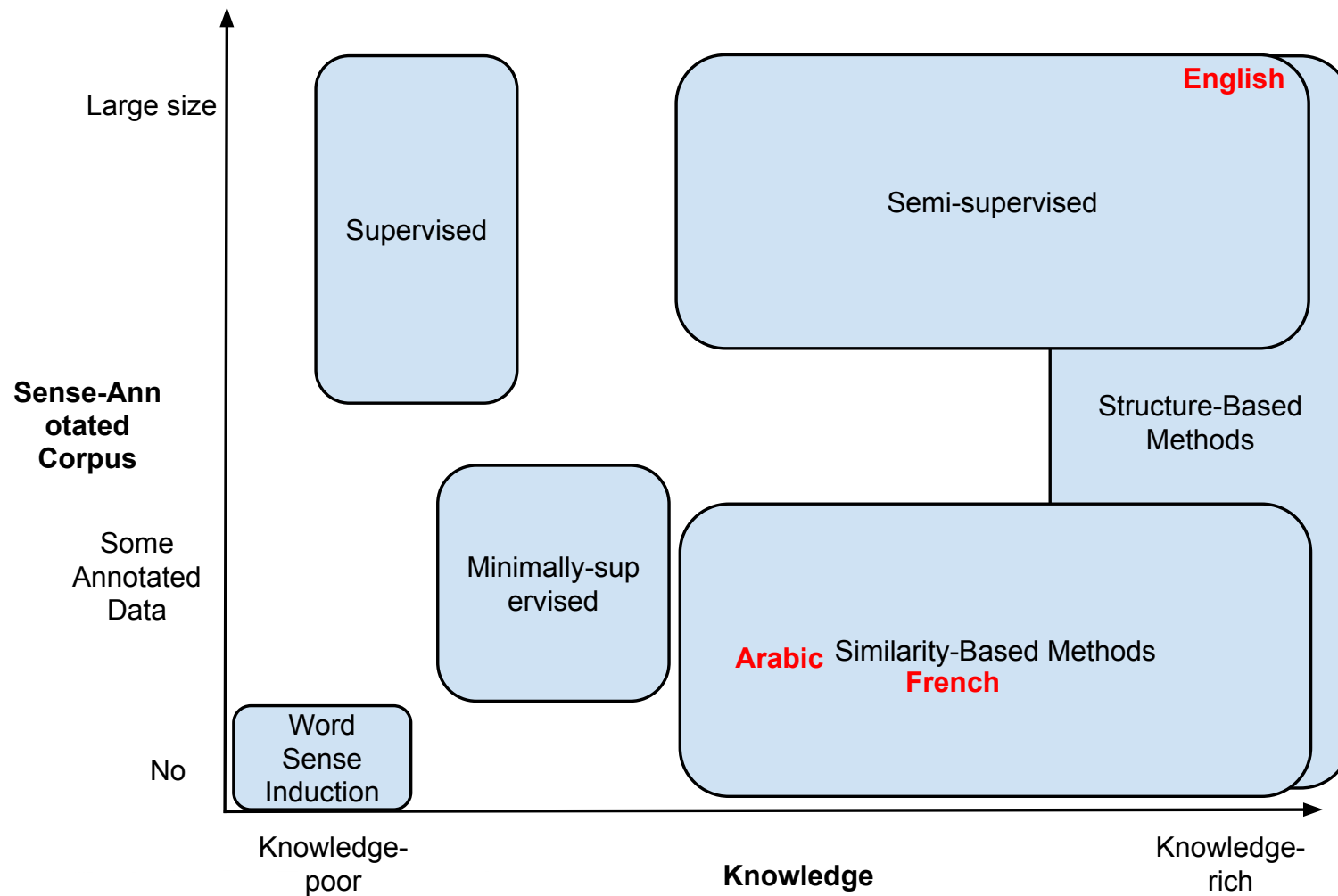


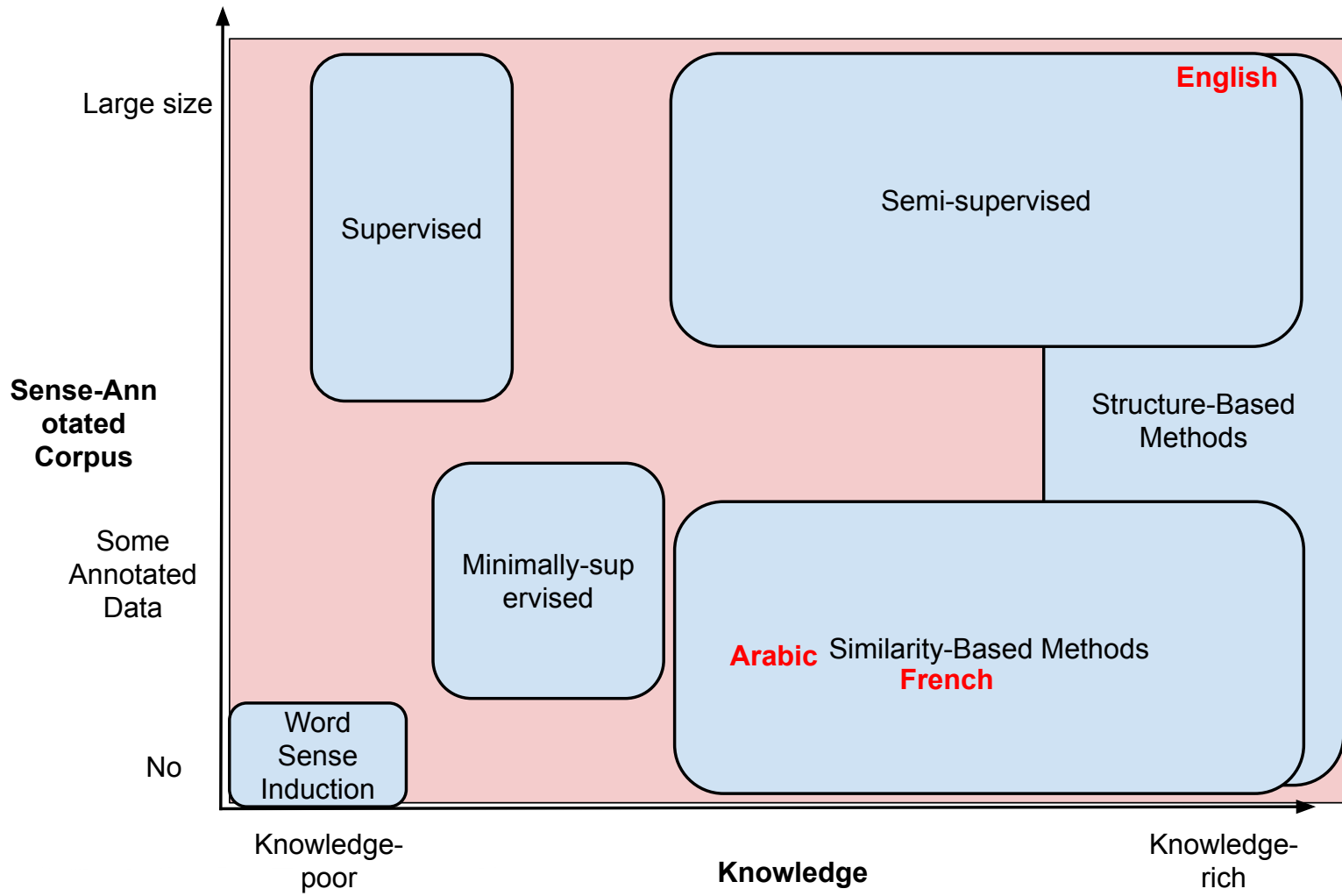
# State of the art neural approach for supervised Word Sense Disambiguation







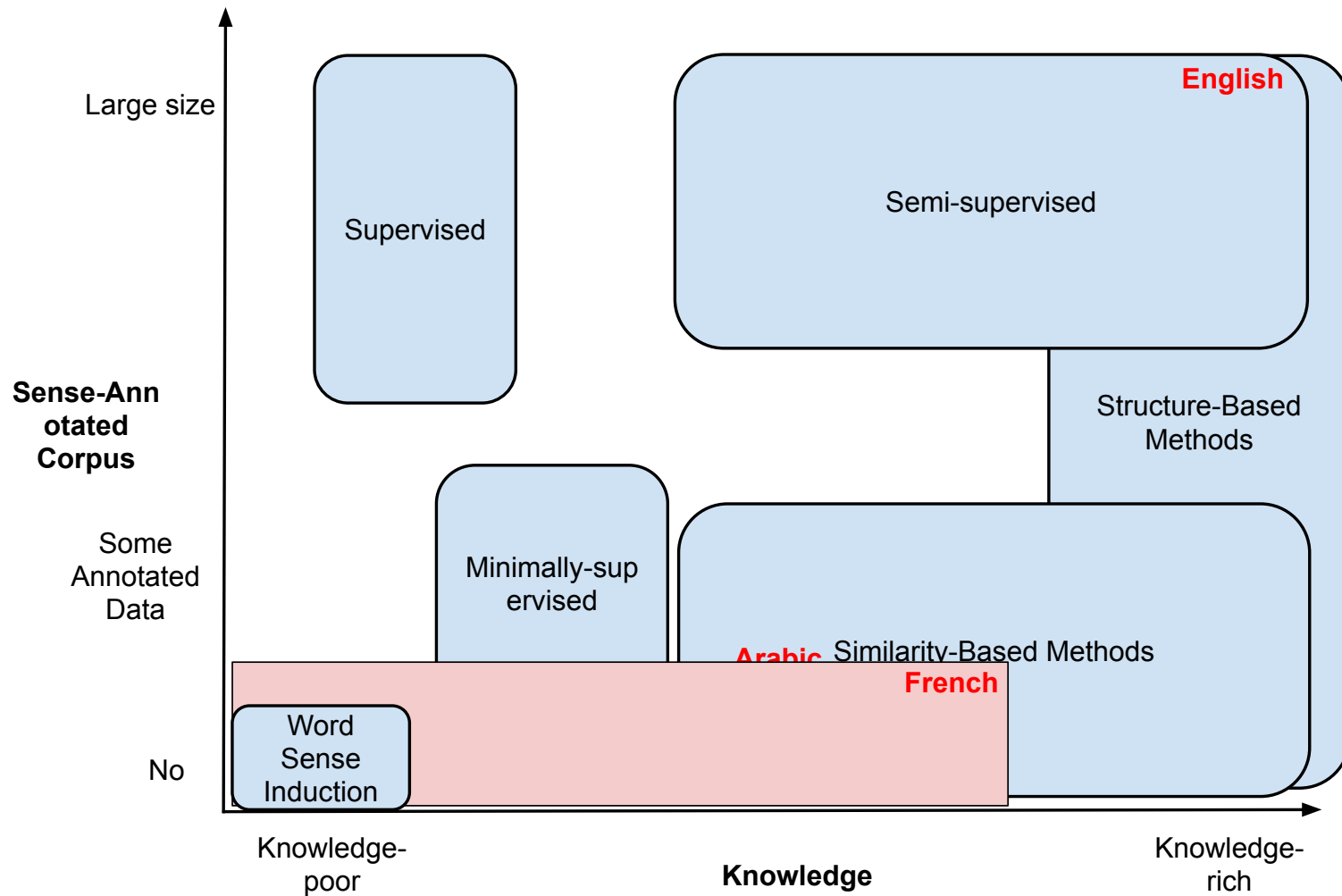




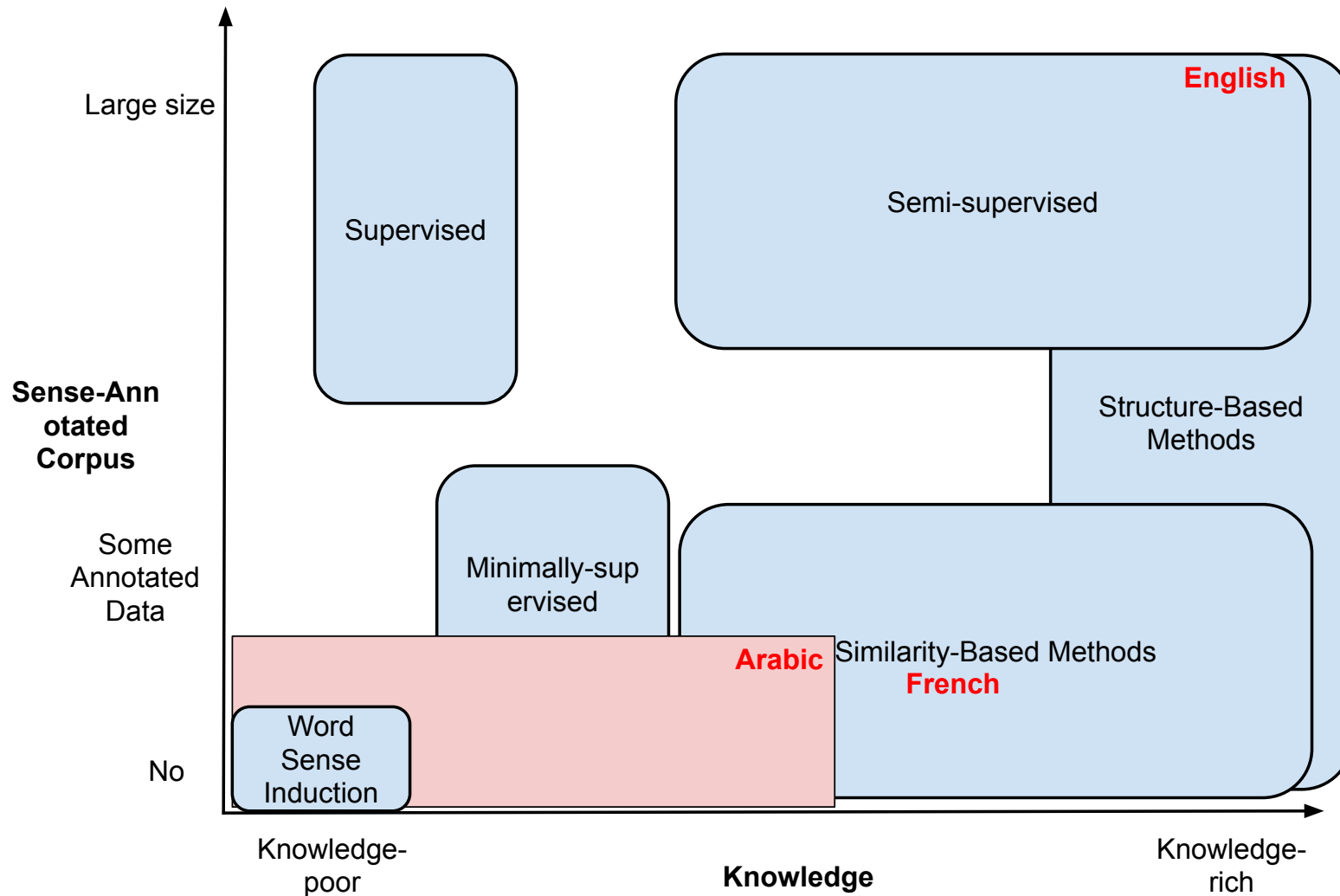
# Situation in 2015 (Methods and Languages)



23



# Situation in 2015 (Methods and Languages)



# Work on WSD for anything else but English: Sense-annotated corpus wanted!



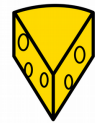
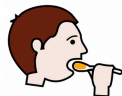
25

## From English Sense-Annotated Corpus [Hadj Salah et al., 2018]

- Only need English-to-target-language machine translation system
- Method:
  - Translation of corpus to target language (home-made machine translation system or external tool)
  - Word alignment (FastAlign)
  - Post-processing (word reordering, duplication correction,...)

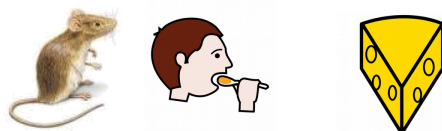


The mouse ate the cheese





The mouse ate the cheese



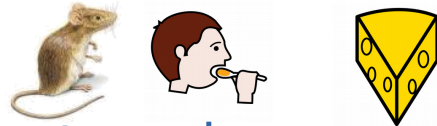
La souris mangea le fromage

Translation





The mouse ate the cheese



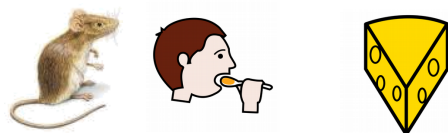
La souris mangea le fromage



Word Alignment



The mouse ate the cheese



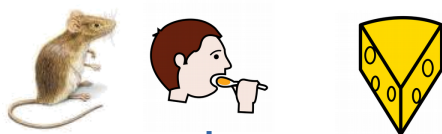
La souris mangea le fromage



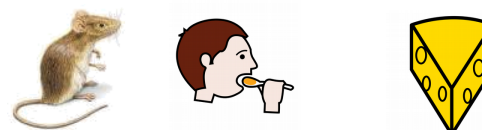
Annotation Transfert



The mouse ate the cheese



La souris mangea le fromage



Annotation transfert

# Work on WSD for anything else but English: sense-annotated corpus wanted!



31

## From English Sense-Annotated Corpus [Hadj Salah et al., 2018]

- Just need an English to target language Machine Translation system
- Method:
  - Translation of corpus to target language (home-made machine translation system or external tool)
  - Word alignment (FastAlign)
  - Post-processing (word reordering, duplication correction,...)

## UFSAC [Vial et al., 2018]

- Unification of Sense Annotated Corpora and Tools
- 12 English Corpora

# Situation in 2019 (resources)



32

resource	Sentence	Words		Part of Speech			
		Overall	Annotated	Nouns	Verbs	Adjectives	Adverbs
SemCor	37 176	778 587	229 533	87 581	89 051	33 752	19 149
DSO	101 004	2 705 190	176 197	105 245	70 952	0	0
WNGT	117 659	1 634 691	496 776	287 798	77 234	107 135	24 609
MASC	31 760	585 354	113 546	49 474	39 356	12 894	11 822
OMSTI	820 084	35 800 061	920 357	476 692	253 555	190 110	0
OntoNotes	124 851	2 475 926	233 616	79 765	153 851	0	0
SemEval 2007 task 07	245	5 637	2 261	1 108	591	356	206
SemEval 2007 task 17	126	3 438	455	159	296	0	0
SemEval 2 013 task 12	306	8 142	1 644	1 644	0	0	0
SemEval 2015 task 13	138	2 637	1 053	554	251	166	82
Senseval 2	238	5 589	2 301	1 061	541	422	277
Senseval 3 task 1	300	5 507	1 957	886	723	336	12
Total (UFSAC)	1 233 649	44 010 759	2 179 696	1 091 967	686 401	345 171	56 157
Total (UFSAC-Ara)	1 233 649	36 213 777	2 001 918	1 011 258	624 771	314 449	51 440
Total (UFSAC-Fra)	1 233 649	41 447 346	1 661 726	949 304	526 715	149 306	36 401

UFSAC [Vial et al., 2018], UFSAC-ARA [Hadj Salah et al., 2018]

- <https://github.com/getalp/UFSAC>
- Marwa Hadj Salah, Loïc Vial, Mounir Zrigui, Hervé Blanchon, Benjamin Lecouteux, Didier Schwab



## Drawbacks of current supervised systems

- Output vocabulary (number of sense tags) is large
  - WordNet 3.0 = 206 941 senses
  - These are too many outputs for the softmax layer of a typical NN
- Sense annotated corpora = costly resource ; SemCor: largest manually annotated corpus but only 16% of all WordNet senses are represented

## Sense Vocabulary Compression

- Form groups of similar senses, for instance:
  - group n1 : {mouse1, rat1, rodent1...}
  - group n2 : {mouse4, keyboard1, click4...}
- Learn to predict group tags instead of sense tags during training
- Find back the “true” sense at disambiguation time



At training time:

"I plug my **keyboard** on my PC"

"The **mouse** ate my salad"



At disambiguation time:

"The **rat** eats the cheese"

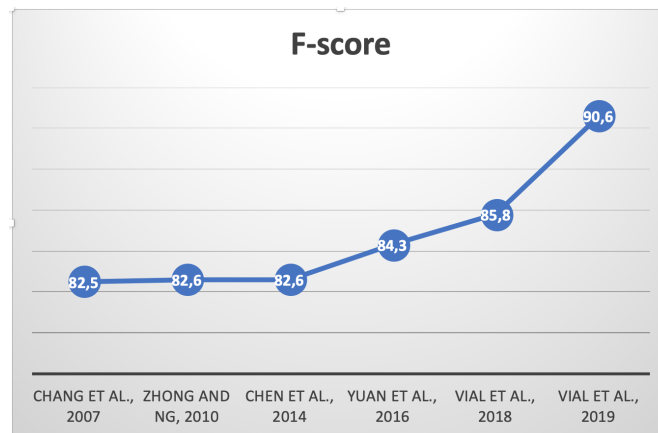


Method	Vocabulary size	Compression rate	SemCor coverage
Senses	206 941	0 %	16 %
Synsets	117 659	43 %	22 %
Hypernyms	39 147	81 %	32 %
All relations	11 885	94 %	39 %



## DNNs and lexical database : best of both worlds

- Smaller number of senses
  - Smaller size of neural models
  - Shorter training time
- Increase coverage
- Better generalization
- Improve results (see below for Eng.)

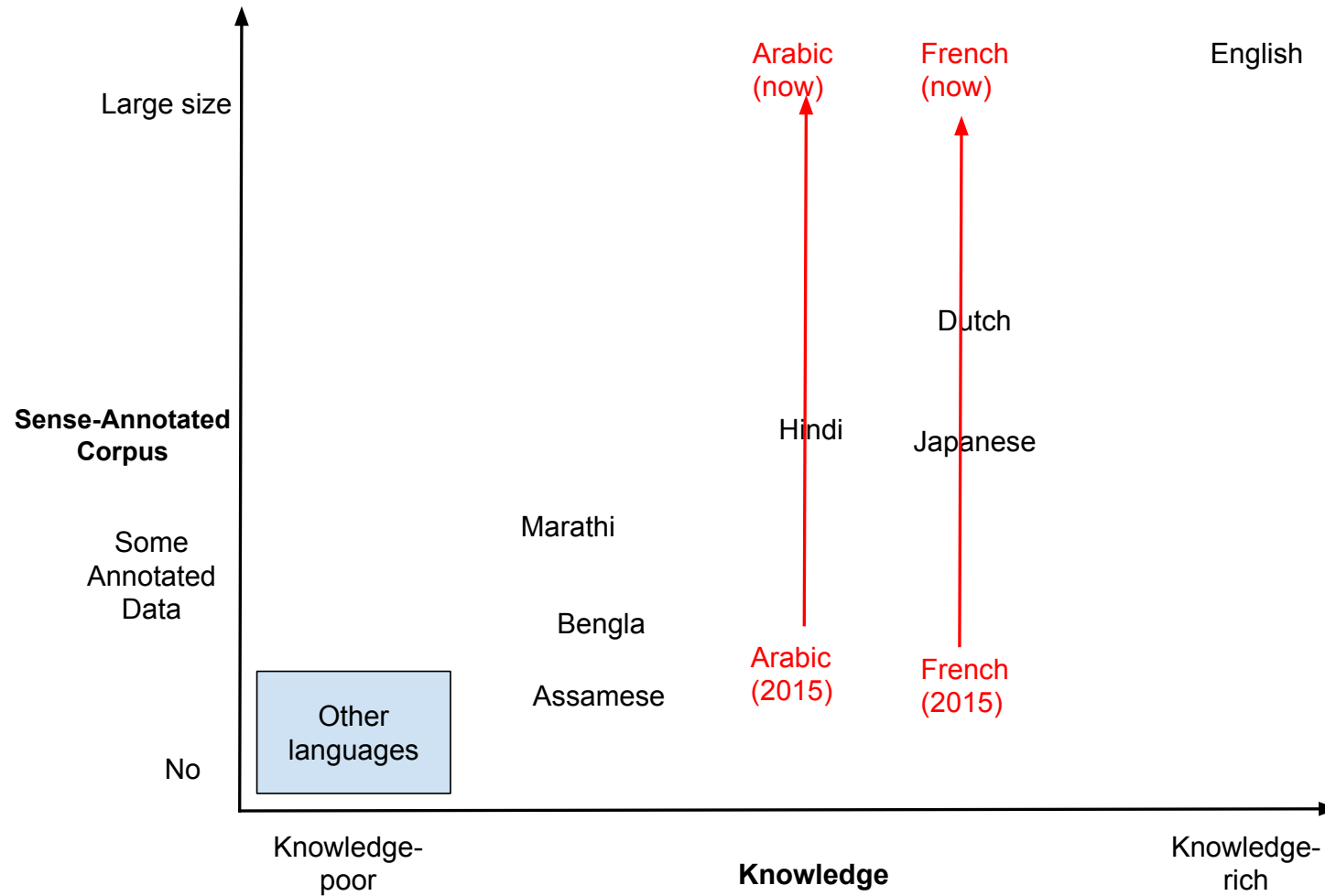




# Situation in 2019 (resources)



36





## Conclusions

- On English, state-of-the-art performance (Best paper TALN 2019)
  - Unification of sense-annotated corpora
  - Use of knowledge to obtain better generalisation in Neural WSD
- On other languages (where there is MT from English)
  - Supervised WSD is now possible
  - State of the art with the same method
- Joint models with neural machine translation: 2 PhDs (Marwa Hadj Salah (Ara), Loïc Vial)

## Perspectives

Automatic Generation of pictograms from speech or text: for cognitive disable people and allophone population - Geneva Hospital, Univ. Geneva (Suisse) and Univ. Louvain-la-Neuve (Belgium)



# Conclusion

- Sense Vocabulary Compression :
  - Easy to implement method
  - Improves the coverage and generalization ability of neural WSD systems
  - Reduces the number of parameters of neural models
- New “contextualized” word embeddings (ELMo, BERT) :
  - Greatly improve the performance of neural WSD systems
  - Improve the state of the art by almost 10 points
- Our code and our pre-trained models are available:  
<https://github.com/getalp/disambiguate>

# Bibliography

- [Gale et al., 1991] One sense per discourse, Gale, William A. and Church, Kenneth W. and Yarowsky, David, HLT '91 Proceedings of the workshop on Speech and Natural Language, Pages 233-237, Harriman, New York, 1991
- [Lesk, 1986] Lesk, M. (1986). Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In SIGDOC '86: Proceedings of the 5th annual international conference on Systems documentation, pages 24-26, New York, NY, USA. ACM
- [Navigli & Lapata, 2010] R. Navigli, M. Lapata, An experimental study on graph connectivity for unsupervised Word Sense Disambiguation, IEEE Transactions on Pattern, 2010  
Analysis and Machine Intelligence 32 (2010) 678–692.
- [Navigli & Ponzetto, 2012] R. Navigli and S. Ponzetto. BabelNet: The Automatic Construction, Evaluation and Application of a Wide-Coverage Multilingual Semantic Network. Artificial Intelligence, 193, Elsevier, 2012, pp. 217-250
- [Rivest, 1987] Learning Decision Lists, Ronald Rivest, Machine Learning, 2, pp 229-246, 1987
- [Sánchez-de-Madariaga & Fernández-del-Castillo, 2009] The bootstrapping of the Yarowsky algorithm in real corpora Ricardo Sánchez-de-Madariaga & José R. Fernández-del-Castillo, Journal Information Processing and Management: an International Journal, Volume 45 Issue 1, Pages 55-69, January, 2009
- [Sarkar, 2009] Bootstrapping a Classifier Using the Yarowsky Algorithm, Anoop Sarkar
- [Schütze, 1992] Dimensions of Meaning, H. Schütze Proceedings of the 1992 ACM/IEEE Conference on Supercomputing, Supercomputing '92, Minneapolis, Minnesota, USA, 1992
- [Yarowsky, 1993] One Sense Per Collocation, Yarowsky, David, Proceedings of the Workshop on Human Language Technology, HLT '93, pages 266-271, 1993