# Word Sense Disambiguation in Queries

Shaung Liu, Clement Yu,

Weiyi Meng

# Objectives

(1) For each content word in a query, find its sense (meaning);

(2) Add terms ( synonyms, hyponyms etc of the determined sense) to the query so as to improve retrieval effectiveness.

# Example

Query: Recycling automobile tire

Recycling:  sense 1: cause to repeat a cycle;

Sense 2 : use again after processing
disambiguated to sense 2:

A synonym: Reuse

Automobile tire has unique sense

A synonym: car tire

Generate phrases: reuse automobile tire,
reuse car tire, recycle car tire

# Our Approach to determine the sense of a content word $t_1$

Find a phrase in the query containing $t_1$. Let the phrase be $(t_1, t_2)$.
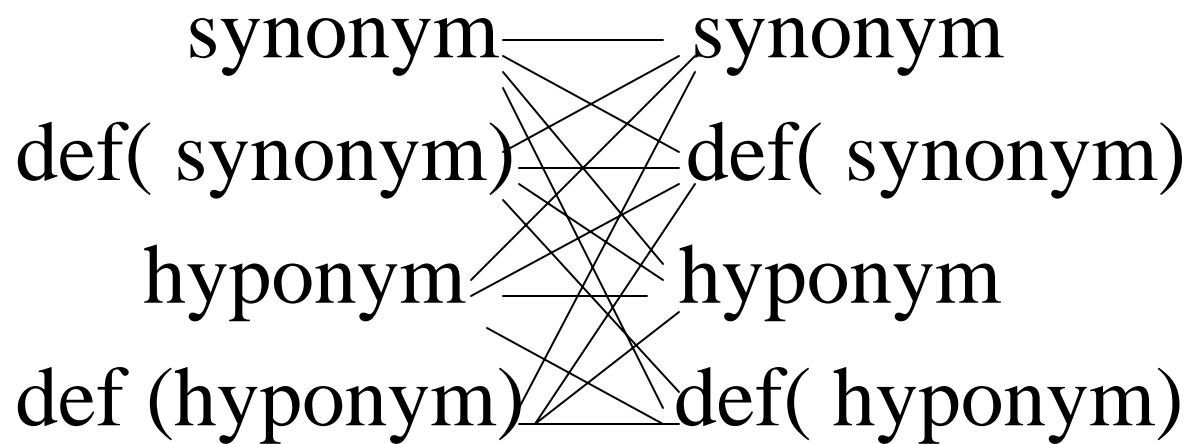
Each $t_i$, $i = 1, 2$, has synonym sets, their definitions, hyponym sets, and their definitions

The sense of $t_1$ is determined by comparing these 4 pieces of information against those of $t_2$

# Comparison of information of t1 against that of t2

t1                              t2

synonym———synonym

def( synonym)———def( synonym)

hyponym———hyponym

def (hyponym)———def( hyponym)

# An Example

Phrase in query: philosophy Stoicism

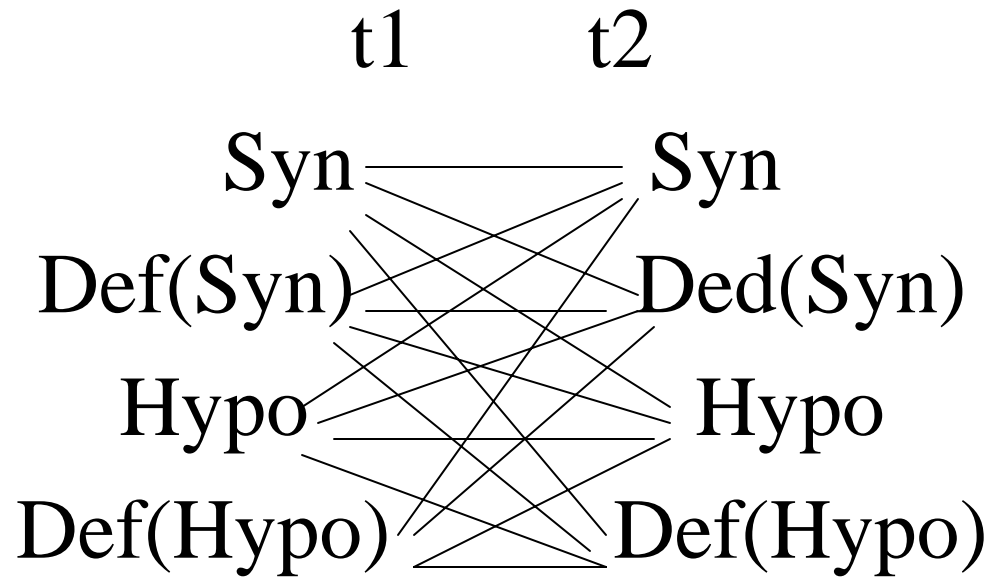A synonym of one sense, S1, of philosophy is "philosophical system"

The definition of one sense, S2, of Stoicism contains "philosophical system". Thus, the sense of philosophy is S1 and that of Stoicism is S2.

# Another example

Query: induction, deduction

The definition of one sense, S1, of induction and that of one sense, S2, of deduction have the common words "reasoning, general". Thus, the sense of induction is determined to be S1 and that of deduction is determined to be S2.

# What happens if multiple senses of a content word are obtained?

t1　　　t2

Syn —————— Syn

Def(Syn)　　　Ded(Syn)

Hypo　　　Hypo

Def(Hypo)　　　Def(Hypo)

16 cases

Two or more cases yield different senses

# Resolve Mutiple senses

2 key parameters:

(1) Historical accuracies of the Cases:

Determined by experiments

(2) Likelihood that a word has a given sense: given by Wordnet (frequency)

# What happens if the technique yields no sense

(1) Choose the most likely sense, if it is at least 50% chance of being correct.

(2) Use Web search to determine the sense.

# Web search to determine sense of a term t

Suppose t has two senses.

From the definition of each sense of t,

form a vector of content words, say V1, V2.

Submit the query containing t to Google.
From the top 20 documents, extract the
content words around t to form a vector V.

Choose sense i, if sim( V, Vi) is maximum.

# Experimental Results

- TREC 2004 queries, robust track

-  250 queries

- 258 unique sense terms, 333 ambiguous terms

|                 | Case  | Frequency | Web |
|-----------------|-------|-----------|-----|
| Applicability   | 65%   | 30%       | 5%  |
| Accuracy        | 89.4% | 93%       | 81% |
| Overall accuracy |      | 90%       |     |

# Similarity function of our system

- Similarity( Q, D) =

  ( phrase similarity, term similarity);

  phrase similarity = sum of idfs of phrases;

  term similarity = Okapi similarity

  D1 is ranked ahead of D2 if phrase-sim 1

  > phrase-sim 2 or if phrase-sim1 =phrase-
  sim 2 and term-sim 1 > term-sim 2

# Recognition of phrases in queries

A phrase, say p,  is  recognized in a query as
   (a) named entity: eg name of person or
   (b) dictionary phrase: in Wordnet or
   (c)  simple phrase: containing two words or
   (d) complex phrase: more than 2 words

# Recognition of phrases in documents

A phrase p, say (term 1, term 2) appears in a document if the terms are within a certain distance.

named entity: terms need to be adjacent

dictionary phrase: terms within distance d1

simple phrase: terms within d2;

complex phrase: d3;  d1 < d2 < d3;

d1, d2, d3 determined by decision tree

# Impact of WSD on effectiveness

|         | No-WSD | WSD | improvement |
|---------|--------|-----|-------------|
| TREC6   | .28    | .32 | 17%         |
| TREC7   | .25    | .31 | 22.6%       |
| TREC8   | .29    | .32 | 11.4%       |
| TREC12  | .37    | .41 | 10.5%       |
| TREC13  | .38    | .42 | 10%         |
| Hard 50 | .18    | .20 | 14.7%       |
| Old 200 | .30    | .34 | 14.9%       |
| Overall | .31    | .35 | 13.7%       |

(previous best known result: .33)

# Summary

- Utilizes 3 methods for word sense disambiguation.
- Case analysis, guessing based on frequency, Web search
- Yields 100% coverage and 90% accuracy
- Improves retrieval effectiveness

# Comparison with other word sense disambiguation algorithm

- Earlier works mostly disambiguates words in documents rather than in queries


- Previous "best" result is around 71% accuracy.

# Conclusion

- Accuracy of our current system is around 90%.

- Yields improvement in retrieval effectiveness

- Will attempt to improve both accuracy in word sense disambiguation and retrieval effectiveness