Introduction    Compositional semantics    Lexical semantics    Scientific text processing    Natural and non-natural languages

OO
OOOOOOOOOO
OOOOOOO

OOOOOOO
OOOOO
OOO

OO
OO
OOOO

# What do we mean?

Computational approaches to natural language semantics

## Ann Copestake

Natural Language and Information Processing Group
Computer Laboratory
University of Cambridge

May 2008

# Outline.

Language and language processing

Compositional semantics

Lexical semantics

Scientific text processing

Natural and non-natural languages

Current research in language processing related to semantics, mostly NLIP group, with flashbacks to Karen's work.

# Outline.

## Language and language processing

### Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)

2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)

2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages

○○
○○○○○○○○○○
○○○○○○○

○○○○○○○
○○○○○
○○○

○○
○○
○○○○

## Language and language processing

### Why is automatic language processing difficult?

### Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)

2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)

2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

## Language and language processing

### Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

# Language and language processing

### Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

## Language and language processing

Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

## Language and language processing

Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

# Language and language processing

## Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

## Language and language processing

Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

## Language and language processing

Why is automatic language processing difficult?

Similar strings mean different things:

1. How fast is the TZ? (*fast* CPU speed)
2. How fast will my TZ arrive? (*fast* delivery time)

local ambiguity/vagueness

Different strings mean the same thing:

1. How fast will my TZ arrive? (*my* ordered by me)
2. Please tell me when I can expect the TZ I ordered.

synonymy/near synonymy

## Language and language processing

**So, natural languages are a bad thing, to be replaced wherever possible by precise, well-specified formal languages?**

Natural language properties essential to communication:

- incredibly flexible; learnable while compact
- emergent, evolving systems

Ambiguity/synonymy properties are inherent to flexibility and learnability. (Spärck Jones, 1964, p126–136: 'Model 4 languages')

Language can be indefinitely precise:

- ambiguity is largely local (at least for humans)
- natural languages accommodate (semi-)formal additions

Introduction    Compositional semantics    Lexical semantics    Scientific text processing    Natural and non-natural languages
           oo                        ooooooo            oo
           ooooooooooo               ooooo              oo
           ooooooo                   ooo                oooo

## Language and language processing

So, natural languages are a bad thing, to be replaced wherever possible by precise, well-specified formal languages?

Natural language properties essential to communication:

- incredibly flexible; learnable while compact
- emergent, evolving systems

Ambiguity/synonymy properties are inherent to flexibility and learnability. (Spärck Jones, 1964, p126–136: 'Model 4 languages')

Language can be indefinitely precise:

- ambiguity is largely local (at least for humans)
- natural languages accommodate (semi-)formal additions

## Language and language processing

So, natural languages are a bad thing, to be replaced wherever possible by precise, well-specified formal languages?

Natural language properties essential to communication:

- incredibly flexible; learnable while compact
- emergent, evolving systems

Ambiguity/synonymy properties are inherent to flexibility and learnability. (Spärck Jones, 1964, p126–136: 'Model 4 languages')

Language can be indefinitely precise:

- ambiguity is largely local (at least for humans)
- natural languages accommodate (semi-)formal additions

# Outline.

Introduction   Compositional semantics   Lexical semantics   Scientific text processing   Natural and non-natural languages

●○
○○○○○○○○○○
○○○○○○

○○○○○○○
○○○○○

○○
○○

○○○○

## Natural language interfaces to databases
## (e.g., Copestake and Spärck Jones, 1989)

| OWNER | OOid OSurnam OInits |
| OWNERSHIP | **OWOid OWPid** |
| PARCEL | PPid **PBid** PStrnum PStrnam PLuc |
| | PPark PDwell PFl PCityv PSqft |
| BLOCK | BBid **BWid** |
| WARD | WWid |

- Who owns a house in a street with parcels in Block 3/2?
- Which owners are in Market Place?
  i.e., Which owners own properties which are in Market
  Place?   metonymy

Approach: analyse to produce semantic representation, map to
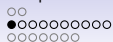domain semantics, map to SQL.

## Natural language interfaces to databases
## (e.g., Copestake and Spärck Jones, 1989)

| | |
|---|---|
| OWNER | OOid OSurnam OInits |
| OWNERSHIP | **OWOid OWPid** |
| PARCEL | PPid **PBid** PStrnum PStrnam PLuc |
| | PPark PDwell PFl PCityv PSqft |
| BLOCK | BBid **BWid** |
| WARD | WWid |

- Who owns a house in a street with parcels in Block 3/2?

- Which owners are in Market Place?
  i.e., Which owners own properties which are in Market Place? metonymy

Approach: analyse to produce semantic representation, map to domain semantics, map to SQL.

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages
●○
○○○○○○○○○○
○○○○○○

○○○○○○○
○○○○○

○○
○○

○○
○○○○

## Natural language interfaces to databases
## (e.g., Copestake and Spärck Jones, 1989)

| OWNER | OOid OSurnam OInits |
| --- | --- |
| OWNERSHIP | **OWOid OWPid** |
| PARCEL | PPid **PBid** PStrnum PStrnam PLuc |
| | PPark PDwell PFl PCityv PSqft |
| BLOCK | BBid **BWid** |
| WARD | WWid |

- Who owns a house in a street with parcels in Block 3/2?

- Which owners are in Market Place?
  i.e., Which owners own properties which are in Market Place? metonymy

Approach: analyse to produce semantic representation, map to domain semantics, map to SQL.

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages

●○
○○○○○○○○○○
○○○○○○

○○○○○○○
○○○○○
○○○

○○
○○

○○○○

# Natural language interfaces to databases
# (e.g., Copestake and Spärck Jones, 1989)

| OWNER | OOid OSurnam OInits |
| --- | --- |
| OWNERSHIP | **OWOid OWPid** |
| PARCEL | PPid **PBid** PStrnum PStrnam PLuc |
| | PPark PDwell PFl PCityv PSqft |
| BLOCK | BBid **BWid** |
| WARD | WWid |

- Who owns a house in a street with parcels in Block 3/2?

- Which owners are in Market Place?
  i.e., Which owners own properties which are in Market Place? metonymy

Approach: analyse to produce semantic representation, map to domain semantics, map to SQL.

# Limited domain vs broad coverage language processing

- Until late 1980s: limited domain, often detailed semantics. Systems as agents.
- 1990–2005: broad coverage, information management. Systems as aids to humans.
    - Spoken dialogue systems: limited domain-dependent grammars.
    - Broad coverage text processing: shallow analysis.

    Limited compositional semantics.
- 2005–: question answering (aka 'semantic search'), robust inference.

# Technical progress on broad-coverage compositional semantics

- Better parsing (e.g., PARC/Powerset, DELPH-IN, CCG):
  - Deep parsers incorporating statistical ranking
  - Faster deep parsers
  - More robustness
- Better representations:
  - Language-friendly logical representations (event variables, generalised quantifiers)
  - Underspecification (Alshawi and Crouch (1992): Quasi-logical form (QLF). Copestake, Flickinger, Sag, Pollard (2005): MRS)
  - Semantics from shallower parsers (RMRS)
- Semantics as automatic markup on natural language, not replacement.

# Technical progress on broad-coverage compositional semantics

- Better parsing (e.g., PARC/Powerset, DELPH-IN, CCG):
    - Deep parsers incorporating statistical ranking
    - Faster deep parsers
    - More robustness
- Better representations:
    - Language-friendly logical representations (event variables, generalised quantifiers)
    - Underspecification (Alshawi and Crouch (1992): Quasi-logical form (QLF). Copestake, Flickinger, Sag, Pollard (2005): MRS)
    - Semantics from shallower parsers (RMRS)
- Semantics as automatic markup on natural language, not replacement.

## Logical representations: first order predicate calculus

Every cat chased some dog

$\forall x[\text{cat}'(x) \implies \exists y[\text{dog}'(y) \wedge \text{chase}'(x, y)]]$
$\exists y[\text{dog}'(y) \wedge \forall x[\text{cat}'(x) \implies \text{chase}'(x, y)]]$

Cannot decide between scope on the basis of syntax.

Thus requires full parse and scope disambiguation to produce a valid logical representation.

Underspecification allows useful semantic representation even when this is impossible.

## Underspecification and Sudoku solving

|   |   |   | 7 |   |   |   |   | 8 |
|---|---|---|---|---|---|---|---|---|
|   |   | 9 |   |   |   |   | 2 |   |
|   | 5 |   |   | 3 |   |   | 9 |   |
| 8 |   |   |   |   | 2 |   |   |   |
|   |   | 6 |   |   |   | 7 |   |   |
|   |   |   | 4 |   |   |   |   | 1 |
|   | 3 |   |   | 9 |   |   | 6 |   |
|   | 2 |   |   |   |   | 4 |   |   |
| 7 |   |   |   |   | 1 |   |   |   |

# Solving.

|   |   |   | 7 |   |   |   |   | 8 |
|---|---|---|---|---|---|---|---|---|
|   |   | 9 |   |   |   |   | 2 |   |
|   | 5 |   |   | 3 |   |   | 9 |   |
| 8 |   |   |   |   | 2 |   |   |   |
|   |   | 6 |   |   |   | 7 |   |   |
|   |   |   | 4 |   |   |   |   | 1 |
|   | 3 |   |   | 9 |   |   | 6 |   |
|   | 2 |   |   |   |   | 4 |   |   |
| 7 |   |   |   |   | 1 |   |   |   |

# Possibility 1.

|   |   |   | 7 |   |   |   |   | 8 |
|---|---|---|---|---|---|---|---|---|
|   |   | 9 |   |   |   |   | 2 | 7 |
|   | 5 |   |   | 3 |   |   | 9 |   |
| 8 |   |   |   |   | 2 |   |   |   |
|   |   | 6 |   |   |   | 7 |   |   |
|   |   |   | 4 |   |   |   |   | 1 |
|   | 3 |   |   | 9 |   |   | 6 |   |
|   | 2 |   |   |   |   | 4 |   |   |
| 7 |   |   |   |   | 1 |   |   |   |

## Possibility 2.

|   |   |   | 7 |   |   |   |   | 8 |
|---|---|---|---|---|---|---|---|---|
|   |   | 9 |   |   |   |   | 2 |   |
|   | 5 |   |   | 3 |   |   | 9 | 7 |
| 8 |   |   |   |   | 2 |   |   |   |
|   |   | 6 |   |   |   | 7 |   |   |
|   |   |   | 4 |   |   |   |   | 1 |
|   | 3 |   |   | 9 |   |   | 6 |   |
|   | 2 |   |   |   |   | 4 |   |   |
| 7 |   |   |   |   | 1 |   |   |   |

## Underspecification.

|   |   |   | 7 |   |   |   |   | 8 |
|---|---|---|---|---|---|---|---|---|
|   |   | 9 |   |   |   |   | 2 | 7 |
|   | 5 |   |   | 3 |   |   | 9 | 7 |
| 8 |   |   |   |   | 2 |   |   |   |
|   |   | 6 |   |   |   | 7 |   |   |
|   |   |   | 4 |   |   |   |   | 1 |
|   | 3 |   |   | 9 |   |   | 6 |   |
|   | 2 |   |   |   |   | 4 |   |   |
| 7 |   |   |   |   | 1 |   |   |   |

## Inference on underspecified form.

## Inference on underspecified form.

## Semantics via incremental annotation (RMRS)

Most cats noisily chased a large dog
most_DAT cat_NN2 noisily_RR chase_VVD a_AT1 large_JJ dog_NN1

a1:l1:most_q(x1)
a2:l2:cat_n(x2)
a3:l3:noisy(e3)
a4:l4:chase(e4)
a5:l5:a(x5)
a6:l6:large(e6)
a7:l7:dog(x7)

## Semantics via incremental annotation (RMRS)

Most cats noisily chased a large dog
most_DAT cat_NN2 noisily_RR chase_VVD a_AT1 large_JJ dog_NN1

a1:l1:most_q(x1)      x1=x2
a2:l2:cat_n(x2)
a3:l3:noisy(e3)
a4:l4:chase(e4)
a5:l5:a(x5)           x5=x7
a6:l6:large(e6)       a6:ARG1(x7)  l6=l7
a7:l7:dog(x7)

## Semantics via incremental annotation (RMRS)

Most cats noisily chased a large dog
most_DAT cat_NN2 noisily_RR chase_VVD a_AT1 large_JJ dog_NN1

```
a1:l1:most_q(x1)    x1=x2
a2:l2:cat_n(x2)
a3:l3:noisy(e3)     l3=l4  e3=e4
a4:l4:chase(e4)     a4:ARG1(x1)  a4:ARG2(x5)
a5:l5:a(x5)         x5=x7
a6:l6:large(e6)     a6:ARG1(x7)  l6=l7
a7:l7:dog(x7)
```

## Semantics via incremental annotation (RMRS)

Most cats noisily chased a large dog

most_DAT cat_NN2 noisily_RR chase_VVD a_AT1 large_JJ dog_NN1

a1:l1:most_q(x1)    x1=x2  a1:RSTR(h1)  h1$=_q$l2
a2:l2:cat_n(x2)
a3:l3:noisy(e3)     l3=l4  e3=e4
a4:l4:chase(e4)     a4:ARG1(x1)  a4:ARG2(x5)
a5:l5:a(x5)         x5=x7  a5:RSTR(h5)  h5$=_q$l6
a6:l6:large(e6)     a6:ARG1(x7)  l6=l7
a7:l7:dog(x7)

## Semantics via incremental annotation (RMRS)

Most cats noisily chased a large dog
most_DAT cat_NN2 noisily_RR chase_VVD a_AT1 large_JJ dog_NN1

a1:l1:most_q(x1)　　x1=x2　a1:RSTR(h1)　h1$=_q$l2　　a1:BODY(l5)
a2:l2:cat_n(x2)
a3:l3:noisy(e3)　　　l3=l4　e3=e4
a4:l4:chase(e4)　　　a4:ARG1(x1)　a4:ARG2(x5)
a5:l5:a(x5)　　　　　x5=x7　a5:RSTR(h5)　h5$=_q$l6　　a1:BODY(l3)
a6:l6:large(e6)　　　a6:ARG1(x7)　l6=l7
a7:l7:dog(x7)

## Semantics via incremental annotation (RMRS)

Most cats noisily chased a large dog
most_DAT cat_NN2 noisily_RR chase_VVD a_AT1 large_JJ dog_NN1

a1:l1:most_q(x1)     x1=x2   a1:RSTR(h1)   h1=$_q$l2     a1:BODY(l3)
a2:l2:cat_n(x2)
a3:l3:noisy(e3)      l3=l4   e3=e4
a4:l4:chase(e4)      a4:ARG1(x1)   a4:ARG2(x5)
a5:l5:a(x5)          x5=x7   a5:RSTR(h5)   h5=$_q$l6     a1:BODY(l1)
a6:l6:large(e6)      a6:ARG1(x7)   l6=l7
a7:l7:dog(x7)

## A real example

Very few of the Chinese construction companies consulted were even remotely interested in entering into such an arrangement with a local partner.

# A real example

Very few of the Chinese construction companies consulted were even remotely interested in entering into such an arrangement with a local partner.

## modified quantifier

# A real example

Very few of the Chinese construction companies consulted were even remotely interested in entering into such an arrangement with a local partner.

## partitive

# A real example

Very few of the Chinese construction companies consulted were even remotely interested in entering into such an arrangement with a local partner.

## compound nominal

Introduction | **Compositional semantics** | Lexical semantics | Scientific text processing | Natural and non-natural languages

oo
oooooooooo●
ooooooo

ooooooo
ooooo
ooo

oo
oo
oooo

# A real example

Very few of the Chinese construction companies consulted
were even remotely interested in entering into such an
arrangement with a local partner.

## reduced relative

# A real example

Very few of the Chinese construction companies consulted were <span style="color:red">even remotely</span> interested in entering into such an arrangement with a local partner.

## <span style="color:red">modified modifier</span>

# A real example

Very few of the Chinese construction companies consulted were even remotely interested in entering into such an arrangement with a local partner.

## predeterminer

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages

○○ | ○○○○○○○ | ○○ | ○○
○○○○○○○○○○ | ○○○○○ | ○○
●○○○○○○ | ○○○ | ○○○○

# Question Answering by semantic pattern matching

What eats jellyfish?

Match robust semantics of question with semantics of possible answer:

[ ?x, a:eat(e), a:ARG1(x), a:ARG2(y), jellyfish(y) ] (simplified)

Matches on *turtles eat jellyfish*, *jellyfish are eaten by turtles*
[ turtle(x), a:eat(e), a:ARG1(x), a:ARG2(y), jellyfish(y) ]

But won't match on *jellyfish eat fish*
[ jellyfish(x), a:eat(e), a:ARG1(x), a:ARG2(y), fish(y) ]

Introduction    **Compositional semantics**    Lexical semantics    Scientific text processing    Natural and non-natural languages
oo                oo                            oooooo               oo
                  oooooooooo                    ooooo                oo
                  o●ooooo                       ooo                  oooo

## Jellyfish eaters: pattern matching and inference

Turtles eat jellyfish and they have special hooks in their throats to help them swallow these slimy animals.

Semantic pattern matches
Inference: $P \wedge Q$ entails $P$

Introduction   **Compositional semantics**   Lexical semantics   Scientific text processing   Natural and non-natural languages

○○          ○○○○○○○     ○○
○○○○○○○○○○   ○○○○○      ○○
○○●○○○○       ○○○       ○○○○

## Jellyfish eaters: pattern matching and inference

Sea turtles, ocean sunfish (Mola mola) and blue rockfish all are able to eat large jellyfish, seemingly without being affected by the nematocysts.

Semantic pattern matching: contexts have to be specified to block.
Inference: axioms have to be specified to license.

Negative context may exist in another document, especially in scientific text.

# Jellyfish eaters: pattern matching and inference

Sea turtles, ocean sunfish (Mola mola) and blue rockfish all are able to eat large jellyfish, seemingly without being affected by the nematocysts.

Semantic pattern matching: contexts have to be specified to block.
Inference: axioms have to be specified to license.

Negative context may exist in another document, especially in scientific text.

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages

○○ | ○○○○○○○ | ○○ |
○○○○○○○○○○ | ○○○○○ | ○○ |
○○○●○○○ | ○○○ | ○○○○

## Compositional semantics: summary

- Broad coverage grammars for English and other languages exist which can provide quite detailed compositional semantic representations.

- Logics are relatively 'language friendly' and support underspecification.

- Compositional semantics seen as annotation of text rather than replacement.

- Robust inference and semantic pattern matching (NB ongoing work by Bergmair)

## Karen on compositional semantics

### Spärck Jones, 1985

More recent developments in the theory of grammar, for example Generalized Phrase Structure Grammar (Gazdar et al, 1985) are much more hospitable to exploitation for automatic language processing, though as far as the semantic content necessary for effective language processing goes, one view is that they are essentially still empty vessels, awaiting the water of life in an account of word meanings.

## 'They all had a use once'

# 'They all had a use once'

# Outline.

## Lexical semantics in language applications

- The Information Retrieval approach: no explicit semantic representation.
- Domain-specific semantics: e.g., interfaces to databases.
- Hand code: e.g., WordNet, specialist terminology resources/ontologies.
- Supervised and unsupervised machine learning.

## You shall know a word by the company it keeps! (Firth, 1957)

Words represented as vectors of features:

|          | $feature_1$ | $feature_2$ | ... | $feature_n$ |
|----------|-------------|-------------|-----|-------------|
| $word_1$ | $f_{1,1}$   | $f_{2,1}$   |     | $f_{n,1}$   |
| $word_2$ | $f_{1,2}$   | $f_{2,2}$   |     | $f_{n,2}$   |
| ...      |             |             |     |             |
| $word_m$ | $f_{1,m}$   | $f_{2,m}$   |     | $f_{n,m}$   |

**Features:** co-occur with $word_n$ in some window, co-occur with $word_n$ as a syntactic dependent, occur in $paragraph_n$, occur in $document_n$ ...

First computational application: Spärck Jones (1964)

## Words co-occurring with words

|  | arts | boil | data | function | large | sugar | summarized | water |
|---|---|---|---|---|---|---|---|---|
| apricot | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| pineapple | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| digital | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
| information | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |

(from Jurafsky and Martin, 2008)

apricot: { boil, large, sugar, water }
pineapple: { boil, large, sugar, water }
digital: { arts, data, function, summarized }
information: { arts, data, function, summarized }

Clustering: group together words with 'similar' vectors.

## Words co-occurring with words

|             | arts | boil | data | function | large | sugar | summarized | water |
|-------------|------|------|------|----------|-------|-------|------------|-------|
| apricot     | 0    | 1    | 0    | 0        | 1     | 1     | 0          | 1     |
| pineapple   | 0    | 1    | 0    | 0        | 1     | 1     | 0          | 1     |
| digital     | 1    | 0    | 1    | 1        | 0     | 0     | 1          | 0     |
| information | 1    | 0    | 1    | 1        | 0     | 0     | 1          | 0     |

(from Jurafsky and Martin, 2008)

apricot: { boil, large, sugar, water }
pineapple: { boil, large, sugar, water }
digital: { arts, data, function, summarized }
information: { arts, data, function, summarized }

Clustering: group together words with 'similar' vectors.

## Words co-occurring with words

|  | arts | boil | data | function | large | sugar | summarized | water |
|---|---|---|---|---|---|---|---|---|
| apricot | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| pineapple | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| digital | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
| information | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |

(from Jurafsky and Martin, 2008)

apricot: { boil, large, sugar, water }
pineapple: { boil, large, sugar, water }
digital: { arts, data, function, summarized }
information: { arts, data, function, summarized }

Clustering: group together words with 'similar' vectors.

| Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages |
|---|---|---|---|---|
| oo | oooooooooo | ooooooo | oo | oo |
| | oooooooo | ooooo | oo | oo |
| | | ooo | | oooo |

## Early clustering: Spärck Jones (1967)

Harper (1965): cooccurrence data for 40 nouns from 120,000 words of Russian scientific text: adjective dependents, noun dependents, noun governors.

Harper clustered by:

$$\frac{|V_1 \cap V_2|}{F_1 F_2}$$

where $V_1$, $V_2$ are cooccurring sets, $F_1$, $F_2$ are the frequencies of the nouns in the corpus.

Spärck Jones (1967): Harper's similarity coefficient is 'of doubtful propriety'. Instead clustered ('clumped') by Jaccard:

$$\frac{|V_1 \cap V_2|}{|V_1 \cup V_2|}$$

## Early clustering: Spärck Jones (1967)

Harper (1965): cooccurrence data for 40 nouns from 120,000 words of Russian scientific text: adjective dependents, noun dependents, noun governors.

Harper clustered by:

$$\frac{|V_1 \cap V_2|}{F_1 F_2}$$

where $V_1$, $V_2$ are cooccurring sets, $F_1$, $F_2$ are the frequencies of the nouns in the corpus.

Spärck Jones (1967): Harper's similarity coefficient is 'of doubtful propriety'. Instead clustered ('clumped') by Jaccard:

$$\frac{|V_1 \cap V_2|}{|V_1 \cup V_2|}$$

# Spärck Jones (1967)

## IR (Robertson and Spärck Jones, 1976, 1994)

Term Frequency:

```
TF(i,j) = number of terms t(i) in document d(j)
```

Collection Frequency Weight (inverse document frequency):

```
CFW(i) = log N - log n
where n is the number of documents t(i) occurs in,
N is the total number of documents
```

Document length:

```
NDL = number of terms in d(j) / average number terms
```

Combined weight:

```
CW(i,j) = [CFW(i)*TF(i,j)*(K+1)] / [K*NDL(j)+TF(i,j)]
```

## Verbs in biomedical text (Korhonen et al, 2006)

Gold standard clusters:

| 1 Have an effect on activity (BIO/29) |
|---|
| **1.1 Activate / Inactivate** |
| 1.1.1 Change activity: *activate, inhibit* |
| 1.1.2 Suppress: *suppress, repress* |
| 1.1.3 Stimulate: *stimulate* |
| 1.1.4 Inactivate: *delay, diminish* |
| **1.2 Affect** |
| 1.2.1 Modulate: *stabilize, modulate* |
| 1.2.2 Regulate: *control, support* |
| **1.3 Increase / decrease:** *increase,* |
| decrease |
| **1.4 Modify:** *modify, catalyze* |

| 4 Experimental Procedures (BIO/30) |
|---|
| **4.1 Prepare** |
| 4.1.1 Wash: *wash, rinse* |
| 4.1.2 Mix: *mix* |
| 4.1.3 Label: *stain, immunoblot* |
| 4.1.4 Incubate: *preincubate, incubate* |
| 4.1.5 Elute: *elute* |
| **4.2 Precipitate:** *coprecipitate* |
| *coimmunoprecipitate* |
| **4.3 Solubilize:** *solubilize,lyse* |
| **4.4 Dissolve:** *homogenize, dissolve* |
| **4.5 Place:** *load, mount* |

Verb clustering using a range of features derived via robust parsing (Briscoe and Carroll, 2002).

## Distributional differences (Copestake, 2005)

Magnitude adjectives and non-physical-solid nouns.
Distributional data from the British National Corpus (100 million words)

|       | importance | success | majority | number | proportion | quality | role | problem | part | winds | support | rain |
|-------|-----------|---------|----------|--------|------------|---------|------|---------|------|-------|---------|------|
| great | 310 | 360 | 382 | 172 | 9 | 11 | 3 | 44 | 71 | 0 | 22 | 0 |
| large | 1 | 1 | 112 | 1790 | 404 | 0 | 13 | 10 | 533 | 0 | 1 | 0 |
| high | 8 | 0 | 0 | 92 | 501 | 799 | 1 | 0 | 3 | 90 | 2 | 0 |
| major | 62 | 60 | 0 | 0 | 7 | 0 | 272 | 356 | 408 | 1 | 8 | 0 |
| big | 0 | 40 | 5 | 11 | 1 | 0 | 3 | 79 | 79 | 3 | 1 | 1 |
| strong | 0 | 0 | 2 | 0 | 0 | 1 | 8 | 0 | 3 | 132 | 147 | 0 |
| heavy | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 2 | 4 | 198 |

Andersen: evidence from error corpus that language learners overuse *big*.

Introduction | Compositional semantics | **Lexical semantics** | Scientific text processing | Natural and non-natural languages

○○
○○○○○○○○○○
○○○○○○○

○○○○○○○
●○○○○
○○○

○○
○○
○○○○

## Compound noun relations

- *cheese knife*: knife for cutting cheese
- *steel knife*: knife made of steel
- *kitchen knife*: knife characteristically used in the kitchen

(Spärck Jones (1983) on compound nouns: implications for overall processing architecture.)

- Syntactic parsers can't distinguish: N1(x), N2(y), compound(x,y)
- One approach: human annotation of compounds, machine learning of unseen examples.

# Compound noun relation learning
## (Ó Séaghdha, 2007)

# Compound noun relation learning
## (Ó Séaghdha)



honey bee

company president

**HAVE**

pork pie

**BE**

pine tree

car door

tuna fish

**ACTOR**

pine cone

steel knife

**ABOUT**  fairy tale

**IN**

crime investigation

forest hut

midnight mass

**INST**

cheese knife  rice cooker

home secretary

machine learning

**LEX**

squirrel pasty?

# Compound noun relation learning
## (Ó Séaghdha)

- Treat compounds as single words: doesn't work!

- Constituent similarity: compounds x1 x2 and y1 y2, compare x1 vs y1 and x2 vs y2.
  *squirrel* vs *pork*, *pasty* vs *pie*

- Relational similarity: **sentences** with x1 and x2 vs sentences with y1 and y2.
  *squirrel is very tasty, especially in a pasty* vs
  *pies are filled with tasty pork*

## Human annotation

- Preliminary to supervised machine learning, evaluation of unsupervised techniques.

- Methodology: define categories, develop guidelines, multiple annotators, measure annotator agreement, refine categories and guidelines . . .

- Agreement of 70% quite usual in semantic annotation.

- What's going on?
  Sometimes, local effects: *sponsorship cash*. Cash gained through sponsorship (INST) or sponsorship in form of cash (BE)?

## Human annotation

- Preliminary to supervised machine learning, evaluation of unsupervised techniques.

- Methodology: define categories, develop guidelines, multiple annotators, measure annotator agreement, refine categories and guidelines . . .

- Agreement of 70% quite usual in semantic annotation.

- What's going on?
  Sometimes, local effects: *sponsorship cash*. Cash gained through sponsorship (INST) or sponsorship in form of cash (BE)?

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | Natural and non-natural languages

○○
○○○○○○○○○○
○○○○○○○

○○○○○
○○○○○

○●○

○○
○○

○○○○

# Ontology extraction



**Extraction of hyponymies**

**A beetle is an insect**

**A tibia is a bone**

## Ontology extraction (Herbelot, 2007, 2008)

- Improving recall by extracting complex examples with
  robust semantic patterns:
  *Opah (also known colloquially as moonfish, sunfish,
  kingfish, and Jerusalem haddock) are large, colourful,
  deep-bodied pelagic Lampriform fish comprising the small
  family Lampridae (also spelt Lamprididae).*
- Learning difference between generic and individual uses:
  - A whale is a mammal.
  - A whale escaped from a zoo yesterday.

## Computational lexical semantics

- Karen was a pioneer of many of the basic methods.
- Research really took off in the 1990s with the availability of corpora (and disk space).
- Many linguistic phenomena involved: generics, compounds, polysemy, metonymy.
- Semantic annotation requires considerable thought about phenomenon and experimentation to be successful: even then, quite low agreement.
- Unsupervised methods, such as clustering, are very attractive, but evaluation can be a problem (especially soft clustering).

# Outline.

Introduction   Compositional semantics   Lexical semantics   Scientific text processing   Natural and non-natural languages
oo              oooooooooo               ooooooo             ●o                          
                ooooooo                  ooooo               oo
                                         ooo                 oooo

# FlySlip: aiding manual curation

- FlyBase: database for Drosophila genetics, manually constructed from literature.
- FlySlip: using NLP to improve the process: NLIP group and Dept of Genetics (Karamanis, Seal, Lewin, McQuilton, Vlachos, Gasperin, Drysdale, Briscoe)

## FlySlip: PaperBrowser



- Entity view: anaphorically-linked gene references highlighted (focus determined by curator).
- Base NPs identified: more useful than just gene names.

# Hedge terms: Medlock and Briscoe (2007)

Hedge: *a word or phrase used to allow for additional possibilities or to avoid over-precise commitment.* (OED)

Hedge classification is the task of identifying and labeling the use of speculative language in written text.

Speculative: This unusual substrate specificity may explain why Dronc is resistant to inhibition by the pan-caspase inhibitor.

Non-speculative: These results demonstrate that ADGF-A overexpression can partially rescue the effects of constitutively active Toll signaling in larvae

Weakly-supervised machine learning technique.

# Citations in IR: Ritchie (2008)

# SciBorg: extracting the science from scientific publications

- Use RMRS language as semantic annotation on chemistry papers (standoff annotation on SciXML).
- Support ontology extraction, discourse markup and information extraction.
- NLIP group, Chemistry dept, CeSC (Copestake, Teufel, Murray-Rust, Parker, Corbett, Rupp, Siddharthan, Waldron) with IUCr, Nature, Royal Society of Chemistry (Batchelor).

# SciBorg: information extraction

Paper 1: The synthesis of 2,8-dimethyl-6H,12H-5,11
methanodibenzo[b,f][1,5]diazocine (Troger's base) from
p-toluidine and of two Troger's base analogs from other anilines

Paper 2: . . . Tröger's base (TB) . . . The TBs are usually
prepared from para-substituted anilines

Eventually, robust inference: e.g., search for papers describing
Tröger's base syntheses which don't involve anilines?

# OSCAR: chemistry terms (Corbett, Murray-Rust)

# OSCAR: chemistry terms (Corbett, Murray-Rust)

# OSCAR: chemistry terms (Corbett, Murray-Rust)

# Citation classification (Teufel, Siddharthan, Batchelor)

# Outline.

## Semantic web, scientific text and language processing

- Description logics, OWL etc.
- Ontologies/terminology resources.
- Chemistry Markup Language (CML: Murray-Rust).
- Availability of texts in XML for language processing.
- Publishing as mixture of texts and structured output (e.g., spectra).

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | **Natural and non-natural languages**

○○ | ○○○○○○○○○○ | ○○○○○○○ | ○○
○○○○○○○ | ○○○○○ | ○○ | ○○○
| | ○○○ | ○○○○

## Semantic web publishing

- Claim: Language processing will soon just be needed for legacy texts. All new scientific publication will use semantic markup.

- Scientific publishing is not simply about facts slotting into an agreed framework.

- Counter-claim 1: where we understand what's going on in scientific text, we can learn to annotate it automatically. But most aspects cannot currently be formalised.

- Counter-claim 2: we need language processing experiments and methodology to work out how to do semantic markup.

## Semantic web publishing

- Claim: Language processing will soon just be needed for legacy texts. All new scientific publication will use semantic markup.

- Scientific publishing is not simply about facts slotting into an agreed framework.

- Counter-claim 1: where we understand what's going on in scientific text, we can learn to annotate it automatically. But most aspects cannot currently be formalised.

- Counter-claim 2: we need language processing experiments and methodology to work out how to do semantic markup.

Introduction | Compositional semantics | Lexical semantics | Scientific text processing | **Natural and non-natural languages**

○○
○○○○○○○○○○
○○○○○○○

○○○○○○○
○○○○○
○○○

○○
○○
○○○○

## Semantic web publishing

- Claim: Language processing will soon just be needed for legacy texts. All new scientific publication will use semantic markup.

- Scientific publishing is not simply about facts slotting into an agreed framework.

- Counter-claim 1: where we understand what's going on in scientific text, we can learn to annotate it automatically. But most aspects cannot currently be formalised.

- Counter-claim 2: we need language processing experiments and methodology to work out how to do semantic markup.

## Semantic web publishing

- Claim: Language processing will soon just be needed for legacy texts. All new scientific publication will use semantic markup.

- Scientific publishing is not simply about facts slotting into an agreed framework.

- Counter-claim 1: where we understand what's going on in scientific text, we can learn to annotate it automatically. But most aspects cannot currently be formalised.

- Counter-claim 2: we need language processing experiments and methodology to work out how to do semantic markup.

## Information Layer and scientific publishing

- 'Information Layer' (Spärck Jones 2007): connection via words may be good enough for many computing system tasks.
- Semantic publishing best seen as an addition to natural language, not a replacement. One objective should be to make scientific publications more accessible to humans.
- Natural language is flexible and adaptable: can this be emulated in formal languages?

## Maths texts and natural languages (Ganesalingam)

Then $V = U \cap H$ for some $U$ in $\mathcal{T}$, by definition of $\mathcal{T}_H$, and
$U \cap H = i^{-1}(U)$, so $g^{-1}(V) = g^{-1}(i^{-1}(U)) = (i \circ g)^{-1}(U)$.

Sutherland, W. A., Introduction to Metric and Topological Spaces, OUP 1975, p. 52.

Analogous to 'donkey sentence' in linguistics.

Every farmer who owns a donkey beats it.

$\forall x[farmer(x) \wedge \exists y[donkey(y) \wedge own(x,y)]] \implies beat(x,y)]$

## Maths texts and natural languages (Ganesalingam)

Then $\underline{V = U \cap H}$ for some $\underline{U}$ in $\mathcal{T}$, by definition of $\mathcal{T}_H$, and $U \cap H = i^{-1}(U)$, so $g^{-1}(V) = g^{-1}(i^{-1}(U)) = (i \circ g)^{-1}(U)$.

Sutherland, W. A., Introduction to Metric and Topological Spaces, OUP 1975, p. 52.

Analogous to 'donkey sentence' in linguistics.

Every farmer who owns a donkey beats it.

$$\forall x[farmer(x) \wedge \exists y[donkey(y) \wedge own(x, y)]] \implies beat(x, y)]$$

## Concluding comments

- Computational semantics: enrich texts to make aspects of meaning more accessible to subsequent processing.

- Underspecifiable, 'surfacy' representations of compositional semantics: logically defined, but robustness, reasonable processing speed.

- Lexical semantics by distributional methods can (partially) model ambiguity/synonymy behaviour (though evaluation still a problem).

- Practical applications to scientific text processing.

- Karen's 'Information Layer' challenges us to take natural language's properties seriously.