

Multiword expressions - a pain in the neck for NLP

Sag, Baldwin, Bond, Copestake, Flickinger (2002)

presented by
William Blacoe

Saarland University, Saarbrücken

Outline

- 1 Introduction
- 2 Linguistic Analysis
- 3 Implementation
- 4 Conclusions

Outline

1 Introduction

2 Linguistic Analysis

3 Implementation

4 Conclusions

Current Situation

What is needed for Natural Language Understanding?

- world knowledge?
- disambiguation?
- symbolic or statistical information?
- domain knowledge?

MWE complexity and omnipresence is underappreciated in NLP

Motivation

Problems with understanding MWEs

- **overgeneration**

e.g. "telephone booth" → "telephone closet"

- **idiomaticity**

e.g. how predict that "kick the bucket" is not literal?

- **flexibility**

a words-with-spaces approach is often too rigid

- **lexical proliferation**

listing all possible valid cases for lexical, syntactic or semantic selection

- **tractability**

State of the Art

Current formal approaches

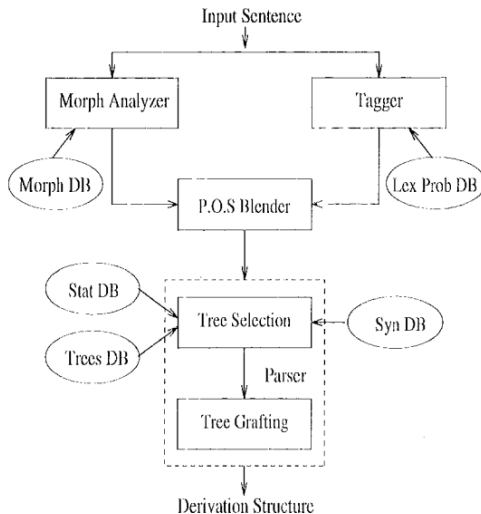
- ParGram
- XTAG
- CCG
- LinGO
- FrameNet

ParallelGrammar (ParGram)

PRED	'see/voir/sehen<(↑ SUBJ),(↑ OBJ)'	
TENSE	FUT	
SUBJ	PRED	'Maria'
	NTYPE	[PROPER NAME]
	PERS	3
	GEND	FEM
	NUM	SG
	CASE	NOM
OBJ	PRED	'Hans'
	NTYPE	[PROPER NAME]
	PERS	3
	GEND	MASC
	NUM	SG
	CASE	ACC
PASSIVE	—	
STMT-TYPE	DECLARATIVE	
VTYPE	MAIN	

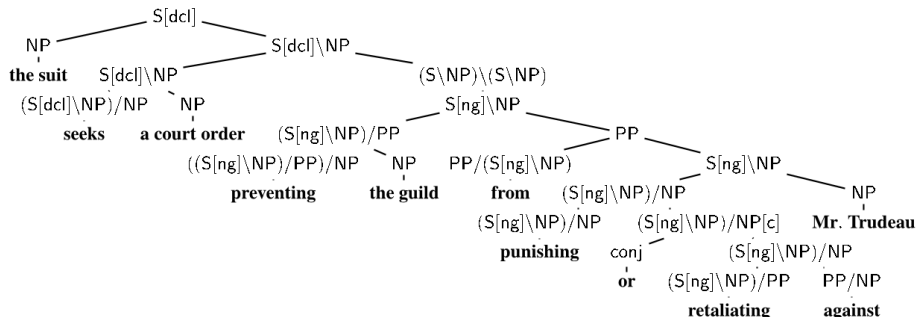
(Butt et al., 1999)

X-Tree Adjoining Grammar (XTAG)



(Doran et al., 1994)

Combinatory Categorical Grammar (CCG)



Categorisation

4 Types of MWEs

- lexicalised phrases
 - fixed expressions
 - semi-fixed expressions
 - syntactically flexible expressions
- institutionalised phrases

Outline

- 1 Introduction
- 2 Linguistic Analysis**
- 3 Implementation
- 4 Conclusions

Type 1

fixed expressions

- examples: *ad hoc*, *Palo Alto*
- no morpho-syntactic variation
- not compositional

Type 2

semi-fixed expressions

- non-decomposable idioms:
kick the bucket, shoot the breeze, but not *spill the beans*
- compound nominals: *attorney general, part of speech*
pluralisation is more complex than just adding an s at the end
- proper names: *(the/those) (San Francisco) 49ers*, but not *the Oakland 49ers*
- no syntactic variation, e.g. passivisation
- limited morphological variation, e.g. number
- words-with-spaces is no real solution

Type 3

syntactically flexible expressions

- verb-particle constructions: *write up*, *look up*, *brush up on*
→ semi-compositional, but the semantics of *up* depends highly on the verb
- semi-compositional: *eat up* → *gobble up*
is particle-initial possible with transitive VPCs?
fall off a truck but not *fall a truck off*,
however *call Kim up* and *call up Kim*

Type 3

syntactically flexible expressions

- decomposeable idioms: *spill the beans*, *let the cat out of the bag*
somewhat compositional
- light verbs: *make a mistake*, *give a demo*
highly idiosyncratic
hard to predict which noun can be selected
syntactically very variable, but not fully compositional
- not representable as words-with-spaces
- only partially compositional due to problems with overgeneration and idiomaticity

Type 4

institutionalised phrases

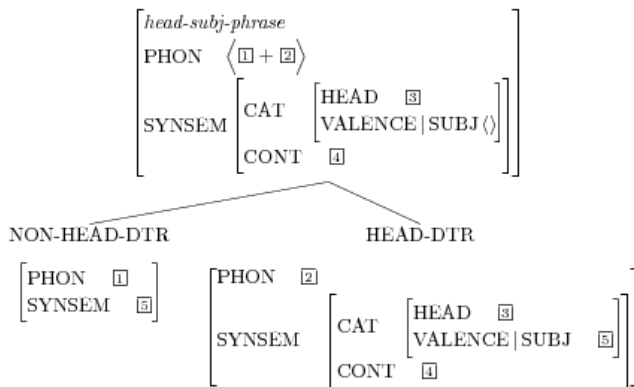
- examples: *traffic light*, *telephone booth*
- semantically and syntactically fully compositional but statistically idiosyncratic
- very high frequency, lexical variants have particularly low frequency

Outline

- 1 Introduction
- 2 Linguistic Analysis
- 3 Implementation**
- 4 Conclusions

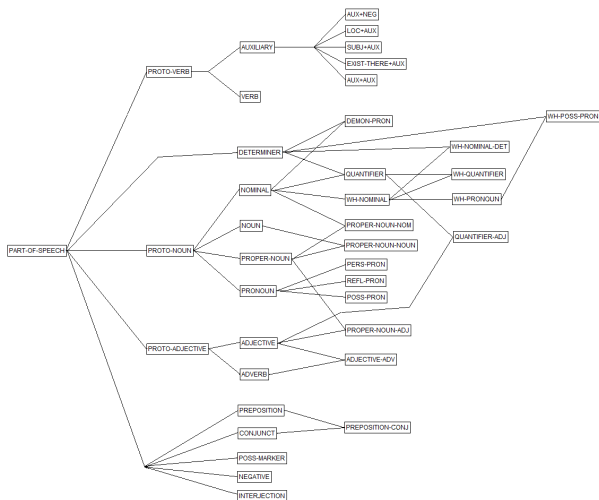
Formal Frameworks

constraint-based HPSG (head-driven phase structure grammar)



Formal Frameworks

constraint-based HPSG (head-driven phase structure grammar)



Formal Frameworks

constraint-based HPSG (head-driven phase structure grammar)
implemented as LKB (lexical knowledge base)

```
scissor := pair-noun-lxm &  
[ ORTH.LIST.FIRST "scissor",  
  SEM.RELS.LIST.FIRST.PRED "scissor_rel" ].
```

Formal Frameworks

ERG (English resource grammar) is compatible with LKB
and uses MRS (minimal recursion semantics) for semantic entries

every dog probably chased some white cat

$$\begin{array}{l}
 \text{TOP } \mathbf{h1} \\
 \text{LZT} < \left[\begin{array}{l} \text{prpstn_rel} \\ \text{HNL } \mathbf{h1} \\ \text{SOA } \mathbf{h21} \end{array} \right] \cdot \left[\begin{array}{l} \text{every_rel} \\ \text{HNL } \mathbf{h3} \\ \text{BV } \mathbf{x4} \\ \text{RESTR } \mathbf{h5} \\ \text{BODY } \mathbf{h6} \end{array} \right] \cdot \left[\begin{array}{l} \text{dog_rel} \\ \text{HNL } \mathbf{h8} \\ \text{INST } \mathbf{x4} \end{array} \right] \cdot \left[\begin{array}{l} \text{probably_rel} \\ \text{HNL } \mathbf{h9} \\ \text{ARG } \mathbf{h10} \end{array} \right] \cdot \left[\begin{array}{l} \text{chase_v_rel} \\ \text{HNL } \mathbf{h12} \\ \text{EVENT } \mathbf{e2} \\ \text{ARG1 } \mathbf{x4} \\ \text{ARG2 } \mathbf{x13} \end{array} \right] \left[\begin{array}{l} \text{TENSE } \text{past} \\ \text{MOOD } \text{indic} \end{array} \right] \cdot \left[\begin{array}{l} \text{some_rel} \\ \text{HNL } \mathbf{h14} \\ \text{BV } \mathbf{x13} \\ \text{RESTR } \mathbf{h15} \\ \text{BODY } \mathbf{h16} \end{array} \right] \cdot \left[\begin{array}{l} \text{white_rel} \\ \text{HNL } \mathbf{h18} \\ \text{ARG } \mathbf{x13} \end{array} \right] \cdot \left[\begin{array}{l} \text{cat_rel} \\ \text{HNL } \mathbf{h18} \\ \text{INST } \mathbf{x13} \end{array} \right] > \\
 \text{H-CONS} < \mathbf{h5} \text{ qeq } \mathbf{h8}, \mathbf{h10} \text{ qeq } \mathbf{h12}, \mathbf{h15} \text{ qeq } \mathbf{h18}, \mathbf{h21} \text{ qeq } \mathbf{h9} >
 \end{array}$$

```

prpstn(probably(every(x, dog(x), some(y, white(y) ∧ cat(y), chase(x, y)))))
prpstn(every(x, dog(x), probably(some(y, white(y) ∧ cat(y), chase(x, y)))))
prpstn(every(x, dog(x), some(y, white(y) ∧ cat(y), probably(chase(x, y)))))
prpstn(probably(some(y, white(y) ∧ cat(y), every(x, dog(x), chase(x, y)))))
prpstn(some(y, white(y) ∧ cat(y), probably(every(x, dog(x), chase(x, y)))))
prpstn(some(y, white(y) ∧ cat(y), every(x, dog(x), probably(chase(x, y)))))
  
```

(Copestake and Flickinger, 2000)

Formal Frameworks

ERG + MRS

running

<div> <div>VBG</div> <div>running</div> </div>	
# 0	<div> <div>XP</div> <div>NP</div> <div>N</div> <div>N</div> <div>V</div> <div>V</div> <div>running</div> </div> <div> <div>e3:</div> <div>e3:unknown(0:7)[ARG x4]</div> <div>_1:udef_q(0:7)[BV x4]</div> <div>e9:_run_v_1(0:7)[]</div> <div>x4.nominalization(0:7)[ARG1 e9]</div> </div>

running

# 1	<div> <div>XP</div> <div>VP</div> <div>V</div> <div>V</div> <div>running</div> </div> <div> <div>e3:</div> <div>e3:_run_v_1(0:7)[]</div> </div>
-----	---

running

# 2	<div> <div>XP</div> <div>NP</div> <div>N</div> <div>N</div> <div>V</div> <div>V</div> </div> <div> <div>e3:</div> <div>e3:unknown(0:7)[ARG x4]</div> <div>_1:udef_q(0:7)[BV x4]</div> <div>e9:_run_v_1(0:7)[]</div> <div>x4.nominalization(0:7)[ARG1 e9]</div> </div>
-----	---

try

# 0	<div> <div>VB</div> <div>try</div> </div> <div> <div>S</div> <div>VP</div> <div>V</div> <div>V</div> <div>try</div> </div> <div> <div>e3:</div> <div>_1:pronoun_q(0:3)[BV x6]</div> <div>x6.pron(0:3)[]</div> <div>e3:_try_v_1(0:3)[ARG1 x6]</div> </div>
-----	---

try

# 1	<div> <div>XP</div> <div>VP</div> <div>V</div> <div>V</div> </div> <div> <div>e3:</div> <div>e3:_try_v_1(0:3)[]</div> </div>
-----	--

try running

# 0	<div> <div>VB</div> <div>try</div> </div> <div> <div>S</div> <div>VP</div> <div>V</div> <div>V</div> <div>V</div> <div>try</div> <div>running</div> </div> <div> <div>e3:</div> <div>_1:pronoun_q(0:11)[BV x6]</div> <div>x6.pron(0:11)[]</div> <div>e3:_try_v_1(0:3)[ARG1 x6, ARG2 e11]</div> <div>e11:_run_v_1(4:11)[ARG1 x6]</div> </div>
-----	--

try running

# 1	<div> <div>XP</div> <div>VP</div> <div>V</div> <div>V</div> <div>V</div> <div>try</div> <div>running</div> </div> <div> <div>e3:</div> <div>e3:_try_v_1(0:3)[ARG2 e7]</div> <div>e7:_run_v_1(4:11)[]</div> </div>
-----	---

try running

# 2	<div> <div>XP</div> <div>AP</div> <div>N</div> <div>V</div> <div>V</div> <div>try</div> </div> <div> <div>e3:</div> <div>e3:unknown(0:11)[]</div> <div>e7:compound(0:11)[ARG1 e5, ARG2 x6]</div> <div>_1:udef_q(0:3)[BV x6]</div> </div>
-----	--

Encoding MWEs

for fixed expressions use words-with-spaces

```
ad_hoc_1 := intr_adj_1 &
  [ STEM < "ad", "hoc" >,
    SEMANTICS [KEY ad-hoc_rel ]].
```

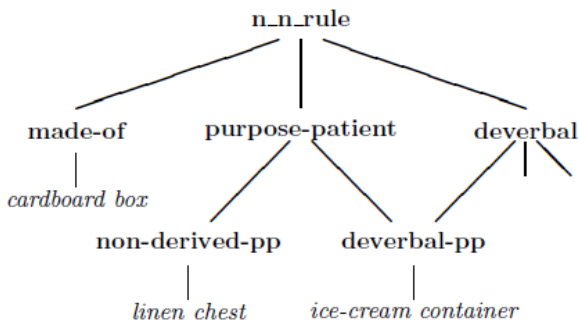
for semi-fixed expressions

- encode internal inflection
- hierarchical lexicon

```
part_of_speech_1 := intr_noun_1 &
  [ STEM < "part", "of", "speech" >,
    INFL-POS "1",
    SEMANTICS [KEY part_of_speech_rel ]].
```

Encoding MWEs

for syntactically flexible expressions encode semi-productivity in a type graph (here: noun compounds)



Encoding MWEs

for syntactically flexible expressions

- encode lexical selection in co-occurring word's representation
- use semantic relations for light verbs
- for semantics of decomposeable idioms combine predicates with idiomatic interpretation of involved words

for institutionalised phrases use frequency information

Outline

- 1 Introduction
- 2 Linguistic Analysis
- 3 Implementation
- 4 Conclusions**

Conclusions

- useful classification of MWEs
- disambiguation is key (not treated in this paper)
- existing approaches help, but there is much work to be done
- symbolic treatment of grammar is not enough
statistical information is necessary

References

- Miriam Butt, Stefanie Dipper, Anette Frank, and Tracy Holloway King. Writing large-scale parallel grammars for english, french, and german. In Miriam Butt and Tracy Holloway King, editors, *Proceedings of the Lexical Functional Grammar Conference*, Manchester, UK, 1999.
- Ann Copestake and Dan Flickinger. An open-source grammar development environment and broad-coverage english grammar using hpsg. In *Proceedings of the Second conference on Language Resources and Evaluation (LREC)*, Athens, Greece, 2000.
- Christy Doran, Dania Egedi, Beth Ann Hockey, B Srinivas, and Martin Zaidel. Xtag system - a wide coverage grammar for english. In *Proceedings of the 15th International Conference on Computational Linguistics (COLING)*, pages 922–928, 1994.