Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
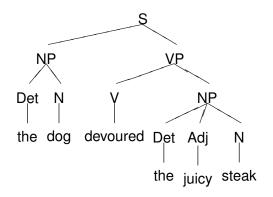Conclusions and Future Work

# Projecting Propbank Roles onto the CCGbank

Stephen A. Boxwell
and
Michael White

The Ohio State University

Three Corpora that Don't Get Along

Distinguishing Arguments from Adjuncts for Fun and Profit

Conclusions and Future Work

The Penn Treebank and the Propbank

CCGbank

Aligning CCGbank and Propbank

# The Penn Treebank WSJ section

- Tens of thousands of sentences from the Wall Street Journal
- Annotated with Part-of-Speech and syntactic structure
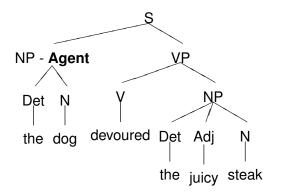- Widely used for a variety of NLP tasks.

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# The Penn Treebank WSJ section

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# The Propbank

- Annotates semantic roles on Penn Treebank trees
- Distinguishes argument roles from modifier roles (manner of action, duration, etc)
- Identifies role-bearing constituents using terminal index and height
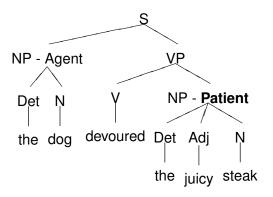- Example: the "Agent" is at terminal index 2, at height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Penn Treebank Tree with Semantic Role annotated



**Agent**: terminal index 2, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Penn Treebank Tree with Semantic Role annotated



Agent: terminal index 2, height 1
**Patient**: terminal index 6, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# The CCGbank

- Combinatory Categorial Grammar is a grammar formalism that treats words as functions and arguments
- A corpus of CCG derivations derived automatically from the Penn Treebank
- CCGbank removes traces and some punctuation
- CCGbank is binary branching, PTB is not.

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# The CCG formalism

- CCG uses syntactically informative lexical categories
- Slash direction ( / or \ ) indicates direction of combinatory potential
  - NP/N = determiner (the, a)
  - PP/NP = preposition (to, with)
  - S\NP = intransitive verb (sleep, die)
  - (S\NP)/NP = transitive verb (devour, love)
  - ((S\NP)/NP)/NP = ditransitive verb (give)
  - ((S\NP)/PP)/NP = ditransitive verb (put)

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# How CCG Categories Make Sentences

*The   dog   devoured   the   juicy   steak*

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# How CCG Categories Make Sentences

| *The* | *dog* | *devoured* | *the* | *juicy* | *steak* |
|-------|-------|------------|-------|---------|---------|
| np/n | n | (s\np)/np | np/n | n/n | n |

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# How CCG Categories Make Sentences

| *The* | *dog* | *devoured* | *the* | *juicy* | *steak* |
|-------|-------|-----------|-------|---------|---------|
| np/n | n | (s\np)/np | np/n | n/n | n |

$$\frac{\text{np/n} \quad \text{n}}{\text{np}} >$$

$$\frac{\text{n/n} \quad \text{n}}{\text{n}} >$$

$$\frac{\text{np/n} \quad \text{n}}{\text{np}} >$$

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# How CCG Categories Make Sentences

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# How CCG Categories Make Sentences

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# The CCGbank and Propbank

- The CCGbank cannot be used directly with the Propbank
- CCGbank terminals $\neq$ PTB terminals
- Binary branching constraint causes tree height mismatch

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Inadvisable Application of Propbank Role to Derivation



Agent: terminal index 2, height 1
**Patient**: terminal index 6, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Trace Annotated with Semantic Role



**Agent** (run): index 3, height 1 AND terminal index 5, height 0

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
**Aligning CCGbank and Propbank**

**Patient** (devour): terminal index 7, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Application of Propbank Role to Derivation Impossible



Patient (devour): terminal index 7, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Application of Propbank Role to Derivation Impossible



Patient (devour): terminal index 7, height 1 **FAIL**

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Aligning the CCGbank and Propbank

- Use a minimum edit distance utility to align the terminals of PTB and CCGB
- Create a mapping of PTB terminals to CCGB terminals
- Find a node in the CCG derivation that covers all and only the correct terminals

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

Theme (want): terminal index 3, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Incorrect Application of Semantic Role to Derivation



| Robin | wants | me | to | eat | steak |
|-------|-------|-----|-----|-----|-------|
| np | ((s\np)/(s\np))/np | np | (s\np)/(s\np) | (s\np)/np | np |

$$\dfrac{\text{wants} \quad me}{(s\backslash np)/(s\backslash np)} >$$

Theme (want): terminal index 3, height 1

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Incorrect Application of Semantic Role to Derivation



Theme (want): terminal index 3, height 1 **FAIL**

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Addressing the Small Clause Mismatch

- Split the role marked on the small clause in two
- Theme → Theme.a, Theme.b
- New notation allows original annotation to be recovered if desired

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Incorrect annotation of theme of "wants"

| Robin | wants | me | to | eat | steak |
|-------|-------|-----|-----|-----|-------|
| np | ((s\np)/(s\np))/np | np | (s\np)/(s\np) | (s\np)/np | np |

$$(s\backslash np)/(s\backslash np)$$ >

$$s\backslash np$$ >

$$s\backslash np$$ >

$$s\backslash np - \textbf{Theme}$$ >

$$s$$ <

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Penn Treebank and the Propbank
CCGbank
Aligning CCGbank and Propbank

# Modified annotation of theme of "wants"

| Robin | wants | me | to | eat | steak |
|-------|-------|-----|-----|-----|-------|
| np | $((s\backslash np)/(s\backslash np))/np$ | $np - $ **Theme.a** | $(s\backslash np)/(s\backslash np)$ | $(s\backslash np)/np$ | np |

$$(s\backslash np)/(s\backslash np) \qquad \rangle$$

$$s\backslash np \qquad \rangle$$

$$s\backslash np - \textbf{Theme.b} \qquad \rangle$$

$$s\backslash np \qquad \rangle$$

$$s \qquad \langle$$

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# The Argument-Adjunct Distinction

- Penn Treebank does not make a strong distinction between arguments and adjuncts
- Argument - adjunct distinction can make a big difference in word-word dependencies, which has implications for generation and semantic role prediction
- CCG theory requires distinction

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# A Verb Consuming an Argument

| *Shakespeare* | *wrote* | *Macbeth* |
|:---:|:---:|:---:|
| np | $(s{\backslash}np)/np$ | np |

$$\frac{(s{\backslash}np)/np \quad np}{s{\backslash}np} >$$

$$\frac{np \quad s{\backslash}np}{s} <$$

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# A Verb Modified by an Adjunct

$$
\frac{
\underset{\text{np}}{\underline{\textit{Shakespeare}}}
\quad
\underset{\text{s\textbackslash np}}{\underline{\textit{wrote}}}
\quad
\frac{
\underset{((\text{s\textbackslash np})\textbackslash(\text{s\textbackslash np}))/\text{np}}{\underline{\textit{in}}}
\quad
\underset{\text{np}}{\underline{\textit{1605}}}
}{
\frac{
(\text{s\textbackslash np})\textbackslash(\text{s\textbackslash np})
}{
\frac{\text{s\textbackslash np}}{\text{s}}
}
}
}{}
$$

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# The Argument-Adjunct Distinction

- Because PTB does not make a good distinction between arguments and adjuncts, CCGbank must make its best guess
- Sometimes CCGbank gets it wrong
- These errors can be identified by discrepencies between Propbank roles and CCGbank categories

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# An Argument that should be an Adjunct

| *join* | *the* | *board* | *as* | *a* | *director* |
|---|---|---|---|---|---|
| $((s\backslash np)/pp)/np$ | $np/n$ | $n$ | $pp/np$ | $np/n$ | $n$ |

$$np \ >$$

$$(s\backslash np)/pp \ >$$

$$pp \ >$$

$$s\backslash np \ >$$

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# An Adjunct that should be an Argument

| bring | new | attention | to | the | problem |
|-------|-----|-----------|-----|-----|---------|
| (s\np)/np | n/n | n | ((s\np)\(s\np))/np | np/n | n |

$$\frac{bring}{(s\backslash np)/np} \quad \frac{\frac{new}{n/n} \quad \frac{attention}{n}}{\frac{n}{np}>} \quad \frac{to}{((s\backslash np)\backslash(s\backslash np))/np} \quad \frac{\frac{the}{np/n} \quad \frac{problem}{n}}{np}>$$

$$\frac{s\backslash np}{} \quad > \quad \frac{(s\backslash np)\backslash(s\backslash np)}{} \quad >$$

$$\frac{s\backslash np}{} \quad >$$

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# Repairing the CCGbank

- 11569 adjuncts converted to arguments
- 1543 arguments converted to adjuncts
- Modifications reflect the judgement of propbank annotators rather than educated guesses from automatic CCGbank generation algorithm

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

# Why Is This Useful?

- We can use syntactic dependencies to annotate verbal categories with semantic roles
- Creates a mapping from CCG syntactic categories to semantic role frames
- Strong implications for semantic role labeling

Three Corpora that Don't Get Along
**Distinguishing Arguments from Adjuncts for Fun and Profit**
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
Why Is This Useful?

| *The* | *dog* | *devoured* | *the* | *juicy* | *steak* |
|---|---|---|---|---|---|
| np/n | n | (s\np)/np | np/n | n/n | n |

np − **Agent** (>)

n (>)

np − **Patient** (>)

s\np (>)

s (<)

Three Corpora that Don't Get Along
**Distinguishing Arguments from Adjuncts for Fun and Profit**
Conclusions and Future Work

The Argument Adjunct Distinction
CCGbank errors
**Why Is This Useful?**

$$
\begin{array}{ccccccc}
\textit{The} & \textit{dog} & \textit{devoured} & \textit{the} & \textit{juicy} & \textit{steak} \\
\hline
\mathrm{np/n} & \mathrm{n} & \mathrm{(s\backslash np_{agent})/np_{patient}} & \mathrm{np/n} & \mathrm{n/n} & \mathrm{n} \\
\end{array}
$$

np − **Agent** >

n >

np − **Patient** >

s\np >

s <

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

How the Argument / Adjunct Repair Improves Performance
Current and Future Work
Acknowledgements

# How Argument / Adjunct Repair Improves Performance

- 96.85% of syntactic arguments found a numbered role (up from 96.13%)
- 89.24% of semantic roles found a syntactic argument (up from 85.71%)
- differences in improvement reflect the relative number of arguments that are converted to adjuncts, and vice versa.

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
Conclusions and Future Work

How the Argument / Adjunct Repair Improves Performance
Current and Future Work
Acknowledgements

# Current and Future Work

- Current work using the modified CCGbank:
  - Hypertagging - generating surface realizations from a logical form (Espinosa, White, and Mehay, ACL 2008)
  - More precise punctuation analysis for CCG realization (White and Rajkumar)

Three Corpora that Don't Get Along
Distinguishing Arguments from Adjuncts for Fun and Profit
**Conclusions and Future Work**

How the Argument / Adjunct Repair Improves Performance
Current and Future Work
Acknowledgements

## Acknowledgements