# Matt's grand plan for natural language understanding (with CCG)

Matthew Honnibal

September 6, 2010

# NLU is currently divided into several tasks

## Example

*The company was persuaded to buy Power Set.*

- **Syntax:** ((The company) (was (persuaded (to buy (Power Set)))))
- **Semantics:** persuade(x, company, buy(company, Power Set))
- **Sense Disambiguation:** *company* as in *army unit*?
- **Named Entity Recognition and Disambiguation:**
  Power Set → http://en.wikipedia.org/PowerSet
- **Coreference? Sentiment? Discourse? More?**

# Pros and cons of dividing up the task

## Pros

- **Reductionism:** It may be easier to make progress on the tasks in isolation.
- **Modularity:** Don't like one parser? Just plug in another!

## Cons

- **Accuracy:** Information can only flow in one direction.
- **Efficiency:** The same work is repeated many times.
- **Plausibility:** Is a pipeline a realistic model of natural language understanding? Should we be trying to find one?

# Intuition behind joint modelling

- $H(W_s)$: information to disambiguate the words in $s$
- $H(R_s)$: information to assign semantic role labels to $s$
- If word senses are good features for SRL, then $H(R_s|W_s) < H(R_s)$
- But if $H(R_s|W_s) < H(R_s)$, then $H(W_s|R_s) < H(W_s)$
- **If WSD helps SRL, then SRL must be able to help WSD.**
- **So: model $P(R_s, W_s)$ instead of $P(W_s)$ and $P(R_s|W_s)$**
- The grand plan: jointly model all the sentence understanding tasks by bringing all the information into a CCG parse.

# Categorial Grammar: few rules, complex categories

Table: Categorial Grammar has only 2 rule schemas, and 3 atomic types.

| | Rules | | | Types |
|---|---|---|---|---|
| $X$ | $\rightarrow$ | $X/Y$ | $Y$ | $N$ |
| $X$ | $\rightarrow$ | $Y$ | $Y\backslash X$ | $PP$ |
| | | | | $S$ |

Table: Production rules get 'translated' into complex categories.

| | PSG | | | CG |
|---|---|---|---|---|
| $NP$ | $\rightarrow$ | $DT$ | $N'$ | $NP/N$ |
| $PP$ | $\rightarrow$ | $IN$ | $NP$ | $PP/NP$ |
| $S$ | $\rightarrow$ | $NP$ | $VP$ | $S\backslash NP$ |
| $VP$ | $\rightarrow$ | $V$ | $NP$ | $(S\backslash NP)/NP$ |
| $VP$ | $\rightarrow$ | $VP$ | $ADVP$ | $(S\backslash NP)\backslash(S\backslash NP)$ |

Introduction
CCG
CCG SRL
CCG NER
WSD CCG
Conclusion
6

# Example Categorial Grammar Derivation

$$
\begin{array}{ccccc}
\text{The} & \text{company} & \text{bought} & \text{Power} & \text{Set} \\
\hline
NP/N & N & (S\backslash NP)/NP & NP/NP & NP \\
\end{array}
$$

The | company | bought | Power | Set
$NP/N$ | $N$ | $(S\backslash NP)/NP$ | $NP/NP$ | $NP$

$NP$ (>)

$NP$ (>)

$S\backslash NP$ (>)

$S$ (<)

# The pay-off: semantic transparency

$$
\frac{\text{The}}{\substack{NP_w/N_w \\ \text{spec}(w)}} \quad \frac{\text{company}}{\substack{N \\ \text{company}}} \quad \frac{\text{bought}}{\substack{(S\backslash NP_x)/NP_y \\ \text{buy}(x,\ y)}} \quad \frac{\text{Power}}{\substack{NP_z/NP_z \\ \text{Power\_z}}} \quad \frac{\text{Set}}{\substack{NP \\ \text{Set}}}
$$

$$
\frac{NP}{\text{company}} >
$$

$$
\frac{NP}{\text{Power\_Set}} >
$$

$$
\frac{S\backslash NP_x}{\text{buy}(x,\ \text{Power\_Set})} >
$$

$$
\frac{S}{\text{buy}(\text{company},\ \text{Power\_Set})} <
$$

Introduction
CCG
CCG SRL
CCG NER
WSD CCG
Conclusion
8

# CCG adds more rules to reduce category ambiguity

$$
\begin{array}{ccccc}
\text{The} & \text{company} & \text{which} & \text{they} & \text{bought} \\
\hline
NP/N & N & (NP\backslash NP)/(S/NP) & NP & (S/NP)\backslash NP \\
\end{array}
$$

The   company   which   they   bought
$NP/N$   $N$   $(NP\backslash NP)/(S/NP)$   $NP$   $(S/NP)\backslash NP$
$NP$ >
$S/NP$ <
$NP\backslash NP$ >
$NP$ <

The   company   which   they   bought
$NP/N$   $N$   $(NP\backslash NP)/(S/NP)$   $NP$   $(S\backslash NP)/NP$
$NP$ >
$S/(S\backslash NP)$ >**T**
$S/NP$ >**B**
$NP\backslash NP$ >
$NP$ <

# PropBank and NomBank: Penn Treebank SRL layers

- A **predicate** heads a proposition (but might not assert it)
- **Arguments** can be **core** or **peripheral**.

## Example predicate-argument structures

(1) *Google*   *bought*   *YouTube*   *October 2006*   *for 1.6bn*
    Arg-0    Predicate   Arg-1      Arg-TMP        Arg-3

(2) *Google*   *paid*    *1.6bn*   *for YouTube*   *October 2006*
    Arg-0    Predicate   Arg-1    Arg-3          Arg-TMP
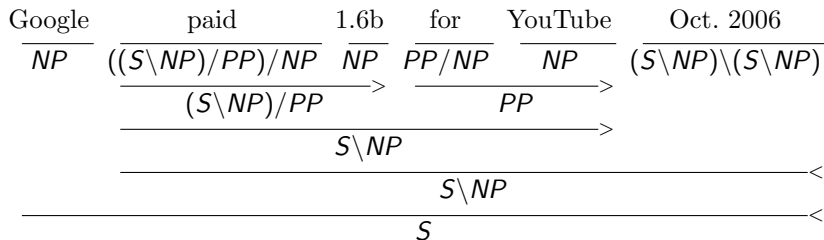
(3) *Google's   1.6bn   acquisition   of YouTube   October 2006*
    Arg-0     Arg-1    Predicate     Arg-3        Arg-TMP

- **PropBank:** Propositions headed by **verbs** in the PTB.
- **NomBank:** Propositions headed by **nouns** in the PTB.

# Integrating PropBank annotation into CCG
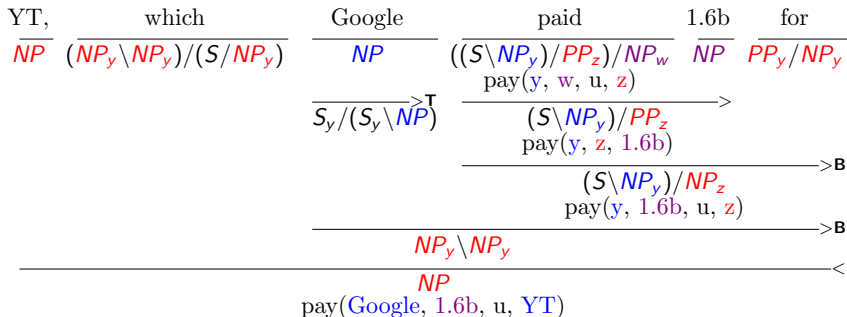
- Target: CCG derivations that map unambiguously to PropBank analyses.

- **Predicates** will be identified by the **semantic category** assigned to them.

- **Core arguments** will be **syntactic complements**. Argument labels will be assigned by the syntax-semantics mapping.

- **Peripheral arguments** will be **syntactic adjuncts**. Their type will be specified in their semantics.

Introduction
000

CCG
0000

CCG SRL
000●0000

CCG NER
000

WSD CCG
000

Conclusion
000

11

# Distinguishing core and peripheral arguments in CCG

$$
\begin{array}{c}
\dfrac{\text{Google}}{NP} \quad
\dfrac{\text{paid}}{((S\backslash NP)/PP)/NP} \quad
\dfrac{\text{1.6b}}{NP} \quad
\dfrac{\text{for}}{PP/NP} \quad
\dfrac{\text{YouTube}}{NP} \quad
\dfrac{\text{Oct. 2006}}{(S\backslash NP)\backslash(S\backslash NP)}
\end{array}
$$

$$
\cfrac{(S\backslash NP)/PP}{} {\scriptstyle >}
$$

$$
\cfrac{PP}{} {\scriptstyle >}
$$

$$
\cfrac{S\backslash NP}{} {\scriptstyle >}
$$

$$
\cfrac{S\backslash NP}{} {\scriptstyle <}
$$

$$
\cfrac{S}{} {\scriptstyle <}
$$

THE UNIVERSITY OF SYDNEY

Introduction
○○○

CCG
○○○○

CCG SRL
○○○●○○○

CCG NER
○○○

WSD CCG
○○○

Conclusion
○○○

12

# Compositional semantics for SRL with CCG

# Compositional semantics for nominal predicates in CCG

$$\frac{\frac{\text{acquisition}}{N_{acquisition}/PP_b}}{\text{acquire}(a,\ b)} \quad \frac{\frac{\text{of}}{PP_c/NP_c} \quad \frac{\text{YouTube}}{NP_{\text{YouTube}}}}{\frac{PP_{\text{YouTube}}}{}>}$$

$$\frac{N_{\text{acquisition}}}{\text{acquire}(a,\ \text{YouTube})}>$$

Introduction
CCG
**CCG SRL**
CCG NER
WSD CCG
Conclusion
14

# Nominal predicates require some creative analyses

$$
\begin{array}{ccc}
\text{YouTube} & \text{'s} & \text{acquisition} \\
\hline
NP_{\text{YouTube}} & (NP_a/(N_a/PP_b))\backslash NP_b & N_{\text{acquisition}}/PP_d \\
& & \text{acquire}(c,\ d)
\end{array}
$$

$$
\underline{\qquad NP_a/(N_a/PP_{\text{YouTube}}) \qquad} <
$$

$$
\underline{\qquad\qquad\qquad NP_{\text{acquisition}} \qquad\qquad\qquad} >
$$
$$
\text{acquire}(c,\ \text{YouTube})
$$

## Nominal predicates with support verbs

$$
\begin{array}{cccc}
\text{Google} & \text{made} & \text{a} & \text{decision} \\
\hline
NP_{\text{Google}} & (S \backslash NP_a)/(NP_b/PP_a) & NP_c/N_c & N_{\text{decision}}/PP_d \\
& & & \text{decide}(d,\ e)
\end{array}
$$

$$
\cfrac{NP_{\text{decision}}/PP_d \quad\quad >\mathbf{B}}{\text{decide}(d,\ e)}
$$

$$
\cfrac{S \backslash NP_d \quad\quad >}{\text{decide}(d,\ e)}
$$

$$
\cfrac{S \quad\quad <}{\text{decide}(\text{Google},\ e)}
$$

# Joint Named Entity Recognition and PTB parsing

- Named entities recognition is usually modelled as a **sequence tagging** task, e.g.
  - **Power|ORG Set|ORG**
- This makes it difficult to account for **nested named entities**, e.g.
  - **New York Stock Exchange**
  - **Sydney, Australia**
  - **David and Melissa Smith**
- Finkel and Manning (2009) joint NER and parsing:
  - Up to 1.36% F-measure parsing improvement;
  - Up to 9% F-measure NER improvement.

# Named entities screw up CCG parses if handled naively
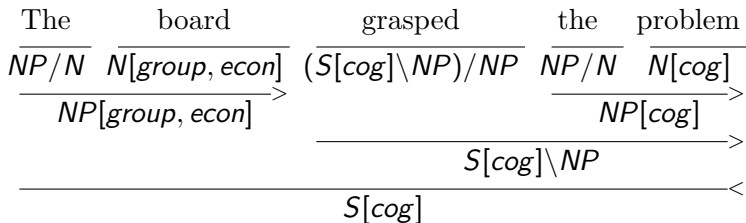
$$
\frac{
\frac{October}{(((VP\backslash VP)/(VP\backslash VP)))/(((VP\backslash VP)/(VP\backslash VP)))} \quad \frac{26}{(VP\backslash VP)/(VP\backslash VP)}
}{
\frac{(VP\backslash VP)/(VP\backslash VP)}{}
} \quad \frac{2006}{VP\backslash VP}
$$

$$
\frac{VP\backslash VP}{}
$$

$$
\frac{The}{NP/N} \quad \frac{Grand}{((N/N)/(N/N))/((N/N)/(N/N))} \quad \frac{Rapids,}{(N/N)/(N/N)} \quad \frac{MI}{N/N} \quad \frac{man}{N}
$$

$$
\frac{(N/N)/(N/N)}{}
$$

$$
\frac{N/N}{}
$$

$$
\frac{N}{}
$$

$$
\frac{NP}{}
$$

# Integrating NER into CCG with Hat Categories

$$\frac{\cfrac{\text{October}}{MON/DAY} \quad \cfrac{\text{26}}{DAY} \quad \cfrac{\text{2006}}{DATE^{VP\backslash VP}\backslash MON}}{}$$

October — $MON/DAY$
26 — $DAY$
2006 — $DATE^{VP\backslash VP}\backslash MON$

$MON$ (>)

$DATE^{VP\backslash VP}$ (<)

$VP\backslash VP$ **H**

The — $NP/N$
Grand — $CITY/CITY$
Rapids — $CITY^{N/N}/STATE$
MI — $STATE$
man — $N$

$CITY^{N/N}$ (>)

$CITY^{N/N}$ (>)

$N/N$ **H**

$N$ (>)

$NP$ (>)

# Tentative thoughts on Word Sense Disambiguation

- Full word sense disambiguation involves many fine-grained labels
- Integrating these labels into CCG category sets may cause sparse data problems
- What if I just use super senses and WordNet Domains?
- 
  - 41 supersenses e.g. noun.food, noun.group, verb.cognition.
  - 46 domains, e.g. economy, sport, fashion, sexuality

# Adding SuperSenses and domains as category features

$$
\frac{\dfrac{\text{The}}{NP/N} \quad \dfrac{\text{board}}{N[group, econ]}}{\dfrac{NP[group, econ]}{}>} \quad \dfrac{\text{grasped}}{(S[cog]\backslash NP)/NP} \quad \dfrac{\dfrac{\text{the}}{NP/N} \quad \dfrac{\text{problem}}{N[cog]}}{NP[cog]}>
$$

$$
\frac{S[cog]\backslash NP}{S[cog]}>
$$

$$
S[cog]<
$$

# WordNet senses for 'board' and 'problem'

| Sense | Super sense | Definition |
|---|---|---|
| 1 | noun.group | A committee having supervisory powers *the board has seven members* |
| 2 | noun.substance | A stout length of sawn timber; made in a wide variety of sizes and used for many purposes |
| 4 | noun.food | Food or meals in general *room and board* |
| 9 | noun.artifact | A flat portable surface (usually rectangular) designed for board games. *he got out the board and set up the pieces* |
| 1 | noun.state | A state of difficulty that needs to be resolved: *she and her husband are having problems* |
| 2 | noun.communication | A question raised for consideration or solution: *our homework consisted of ten problems to solve* |
| 3 | noun.cognition | A source of difficulty *one trouble after another delayed the job* |

# Progress so far

- PropBank/CCGbank integration complete
- Most difficult NomBank/CCGbank integration complete
- Preliminary parsing experiments on modified corpora
- Nicky has BBN and CCG aligned and is working on the integration
- Mike White's group have done something with CCG and discourse parsing

# Current priorities

- Get oracle figures for CCGbank-to-SRL
- Error analysis over oracle errors. Further improvements? Problems with CCG?
- Parse with SRL-CCGbank to get joint model performance.
- Tinker with WSD/CCG ideas at some point.

# Conclusion

- I am focussing on a representation problem, rather than the learning problem. But do these tasks all fit in one hypothesis space? Will the task be tractable?

- It's currently very difficult to deploy a system that makes use of all the NLU modules.

- If my approach works, it will produce a very efficient all-singing-all-dancing NLU solution.

- The project also raises a lot of questions about our current theories of compositional semantics.