

# Towards Broad Coverage Surface Realization with CCG

Michael White, Rajakrishnan Rajkumar, Scott Martin

Department of Linguistics  
The Ohio State University  
{mwhite,raja,scott}@ling.ohio-state.edu

UCNLG+MT Workshop  
11 September 2007



# The Talk in a Nutshell

- Progress towards developing the first broad coverage English surface realizer for Combinatory Categorical Grammar (CCG; Steedman, 2000)

# The Talk in a Nutshell

- Progress towards developing the first broad coverage English surface realizer for Combinatory Categorical Grammar (CCG; Steedman, 2000)
- Respectable initial results (after doing many of the tedious and obvious tasks)

# The Talk in a Nutshell

- Progress towards developing the first broad coverage English surface realizer for Combinatory Categorical Grammar (CCG; Steedman, 2000)
- Respectable initial results (after doing many of the tedious and obvious tasks)
- Details on the corpus-based grammar engineering process

# The Talk in a Nutshell

- Progress towards developing the first broad coverage English surface realizer for Combinatory Categorical Grammar (CCG; Steedman, 2000)
- Respectable initial results (after doing many of the tedious and obvious tasks)
- Details on the corpus-based grammar engineering process
- Lessons learned: supertags can help n-grams

## The Talk in a Nutshell

- Progress towards developing the first broad coverage English surface realizer for Combinatory Categorical Grammar (CCG; Steedman, 2000)
- Respectable initial results (after doing many of the tedious and obvious tasks)
- Details on the corpus-based grammar engineering process
- Lessons learned: supertags can help n-grams — **but**, need a generation supertagger!

# The Talk in a Nutshell

- Progress towards developing the first broad coverage English surface realizer for Combinatory Categorical Grammar (CCG; Steedman, 2000)
- Respectable initial results (after doing many of the tedious and obvious tasks)
- Details on the corpus-based grammar engineering process
- Lessons learned: supertags can help n-grams — **but**, need a generation supertagger!
- Much left to explore

## Why CCG?

- State of the art statistical parsing results (Clark and Curran, 2007)



## Why CCG?

- State of the art statistical parsing results (Clark and Curran, 2007)
- CCGbank (Hockenmaier, 2003) — a corpus of CCG derivations created by transforming the Penn Treebank

## Why CCG?

- State of the art statistical parsing results (Clark and Curran, 2007)
- CCGbank (Hockenmaier, 2003) — a corpus of CCG derivations created by transforming the Penn Treebank
- OpenCCG (White, 2006) chart realizer
  - Used with small, precise grammars in various dialogue systems
  - Supports disjunctive logical forms (“packed” inputs)
  - Has API for statistical scoring models
  - Represents words as factor bundles (form, stem, POS, supertag, etc.)

## Why CCG?

- State of the art statistical parsing results (Clark and Curran, 2007)
- CCGbank (Hockenmaier, 2003) — a corpus of CCG derivations created by transforming the Penn Treebank
- OpenCCG (White, 2006) chart realizer
  - Used with small, precise grammars in various dialogue systems
  - Supports disjunctive logical forms (“packed” inputs)
  - Has API for statistical scoring models
  - Represents words as factor bundles (form, stem, POS, supertag, etc.)

⇒ So just turn the crank with the CCGbank, no?

## Why CCG?

- State of the art statistical parsing results (Clark and Curran, 2007)
- CCGbank (Hockenmaier, 2003) — a corpus of CCG derivations created by transforming the Penn Treebank
- OpenCCG (White, 2006) chart realizer
  - Used with small, precise grammars in various dialogue systems
  - Supports disjunctive logical forms (“packed” inputs)
  - Has API for statistical scoring models
  - Represents words as factor bundles (form, stem, POS, supertag, etc.)

⇒ So just turn the crank with the CCGbank, no? (*Actually, need to add semantics — cf. Bos’s Boxer — and improve OpenCCG’s performance with large grammars*)



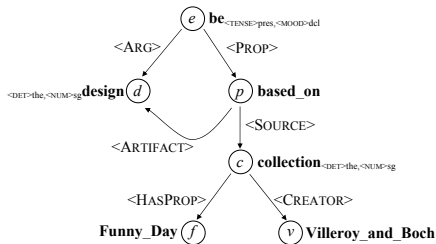
## And this differs how . . .

- Unlike Halogen (Langkilde-Geary, 2002) and FUF/SURGE (Callaway, 2003), OpenCCG uses a bidirectional grammar

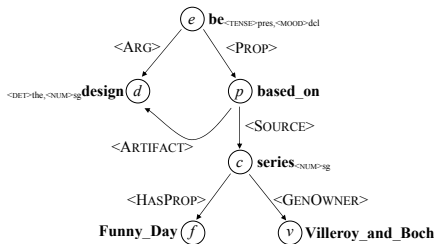
## And this differs how . . .

- Unlike Halogen (Langkilde-Geary, 2002) and FUF/SURGE (Callaway, 2003), OpenCCG uses a bidirectional grammar
- More similar to HPSG/LFG approaches (Carroll and Oepen, 2005; Nakanishi et al., 2005; Cahill and van Genabith, 2006), except for our greater emphasis on more traditional logical forms

## Two Similar Dependency Graphs

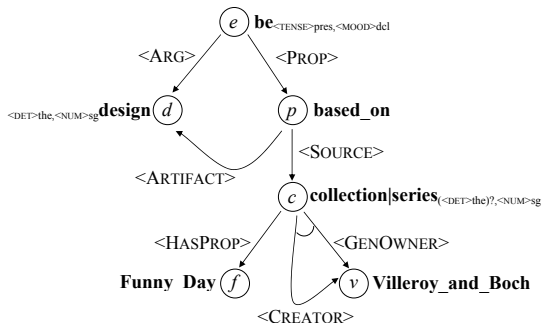


(a) *The design (is|'s) based on the Funny Day collection by Villeroy and Boch.*



(b) *The design (is|'s) based on Villeroy and Boch's Funny Day series.*

# Their Disjunctive Combination



(c) *The design (is|’s) based on (the Funny Day (collection|series) by Villeroy and Boch | Villeroy and Boch’s Funny Day (collection|series)).*



# LF in Hybrid Logic Dependency Semantics (HLDS)

$@_e(\mathbf{be} \wedge \langle \text{TENSE} \rangle \text{pres} \wedge \langle \text{MOOD} \rangle \text{dcl} \wedge$   
 $\langle \text{ARG} \rangle (d \wedge \mathbf{design} \wedge \langle \text{DET} \rangle \text{the} \wedge \langle \text{NUM} \rangle \text{sg}) \wedge$   
 $\langle \text{PROP} \rangle (p \wedge \mathbf{based\_on} \wedge$   
 $\langle \text{ARTIFACT} \rangle d \wedge$   
 $\langle \text{SOURCE} \rangle (c \wedge \mathbf{collection} \wedge \langle \text{DET} \rangle \text{the} \wedge \langle \text{NUM} \rangle \text{sg} \wedge$   
 $\langle \text{HASPROP} \rangle (f \wedge \mathbf{Funny\_Day}) \wedge$   
 $\langle \text{CREATOR} \rangle (v \wedge \mathbf{V\&B}))))$

(a)

## Disjunctive LF in HLDS

$@_e(\mathbf{be} \wedge \langle \text{TENSE} \rangle \text{pres} \wedge \langle \text{MOOD} \rangle \text{dcl} \wedge$   
 $\langle \text{ARG} \rangle (d \wedge \mathbf{design} \wedge \langle \text{DET} \rangle \text{the} \wedge \langle \text{NUM} \rangle \text{sg}) \wedge$   
 $\langle \text{PROP} \rangle (p \wedge \mathbf{based\_on} \wedge$   
 $\langle \text{ARTIFACT} \rangle d \wedge$   
 $\langle \text{SOURCE} \rangle (c \wedge \langle \text{NUM} \rangle \text{sg} \wedge (\langle \text{DET} \rangle \text{the})? \wedge$   
 $\langle \mathbf{collection} \vee \mathbf{series} \rangle) \wedge$   
 $\langle \text{HASPROP} \rangle (f \wedge \mathbf{Funny\_Day}) \wedge$   
 $(\langle \text{CREATOR} \rangle \boxed{\vee} \vee \langle \text{GENOWNER} \rangle \boxed{\vee})))$   
 $\wedge @_v(\mathbf{Villeroy\_and\_Boch})$

(c)

# Flattening

- (2) 0: @<sub>e</sub>(**be**), 1: @<sub>e</sub>(⟨TENSE⟩pres),  
2: @<sub>e</sub>(⟨MOOD⟩dcl), 3: @<sub>e</sub>(⟨ARG⟩d),  
4: @<sub>d</sub>(**design**), 5: @<sub>d</sub>(⟨DET⟩the),  
6: @<sub>d</sub>(⟨NUM⟩sg),  
7: @<sub>e</sub>(⟨PROP⟩p), 8: @<sub>p</sub>(**based\_on**),  
9: @<sub>p</sub>(⟨ARTIFACT⟩d), 10: @<sub>p</sub>(⟨SOURCE⟩c),  
11: @<sub>c</sub>(⟨NUM⟩sg), 12: @<sub>c</sub>(⟨DET⟩the),  
13: @<sub>c</sub>(**collection**), 14: @<sub>c</sub>(**series**),  
15: @<sub>c</sub>(⟨HASPROP⟩f), 16: @<sub>f</sub>(**Funny\_Day**),  
17: @<sub>c</sub>(⟨CREATOR⟩v), 18: @<sub>c</sub>(⟨GENOWNER⟩v),  
19: @<sub>v</sub>(**Villeroy\_and\_Boch**)

- (3) alt<sub>0,0</sub> = {13}; alt<sub>0,1</sub> = {14}  
alt<sub>1,0</sub> = {17, 19}; alt<sub>1,1</sub> = {18, 19}  
opt<sub>0</sub> = {12}

# Edges

## Packed Edges

- In packing mode, a *representative* edge maintains a list of alternative edges whose signs have the same category (but different word sequences)
- Representative edges stand in for their alternative edges during chart construction

# Edges

## Packed Edges

- In packing mode, a *representative* edge maintains a list of alternative edges whose signs have the same category (but different word sequences)
- Representative edges stand in for their alternative edges during chart construction

## The Disjunctive Case

Inspired by Shemtov (1997); see INLG-06 paper for details . . .

## Lexical Instantiation

- (4) a.  $\{11,13,14\}$  *collection*  $\vdash n_c$  :  
 $@_c(\mathbf{collection}) \wedge @_c(\langle \text{NUM} \rangle \text{sg})$
- b.  $\{11,13,14\}$  *series*  $\vdash n_c$  :  
 $@_c(\mathbf{series}) \wedge @_c(\langle \text{NUM} \rangle \text{sg})$
- c.  $\{17\}$   $\text{alt}_{1,0}$  *by*  $\vdash n_c \setminus n_c / np_v$  :  
 $@_c(\langle \text{CREATOR} \rangle v)$
- d.  $\{18\}$   $\text{alt}_{1,1}$  *'s*  $\vdash np_c / n_c \setminus np_v$  :  
 $@_c(\langle \text{GENOWNER} \rangle v)$
- e.  $\{19\}$   $\text{alt}_{1,0}; \text{alt}_{1,1}$  *Villeroy\_and\_Boch*  $\vdash np_v$  :  
 $@_v(\mathbf{V\&B})$

# Derivation

1. {8-10} *based\_on*  $\vdash s_p \setminus np_d / np_c$
2. {12} *the*  $\vdash np_c / n_c$
3. {15, 16} *Funny\_Day*  $\vdash n_c / n_c$
4. {11, 13, 14} *collection*  $\vdash n_c$   
   {11, 13, 14} *series*  $\vdash n_c$
5. {17} *alt*<sub>1,0</sub> *by*  $\vdash n_c \setminus n_c / np_v$
6. {18} *alt*<sub>1,1</sub> *'s*  $\vdash np_c / n_c \setminus np_v$
7. {19} *alt*<sub>1,0</sub>; *alt*<sub>1,1</sub> *Villeroy\_and\_Boch*  $\vdash np_v$

## Derivation (2)

8. {11, 13-16} *FD [collection]*  $\vdash n_c$  (3 4 >)
9. {17-19} *by V&B*  $\vdash n_c \setminus n_c$  (5 7 >)
10. {17-19} *V&B 's*  $\vdash np_c / n_c$  (7 6 <)
11. {11, 13-19} *FD [coll.] by V&B*  $\vdash n_c$  (8 9 <)
12. {11, 13-19} *V&B 's FD [coll.]*  $\vdash np_c$  (10 8 >)
13. {11-19} *the FD [coll.] by V&B*  $\vdash np_c$  (2 11 >)  
    {11-19} *V&B 's FD [coll.]*  $\vdash np_c$  (12 optC)
14. {8-19} *b.\_on [the FD [coll.] ...]*  $\vdash s_p \setminus np_d$  (1 13 >)



# Unpacking

- Complete edges are unpacked bottom-up, a la Langkilde (2000)
- Pruning and scoring configured via API
- At present, edges are pruned only within equivalence classes, during the unpacking stage

# Unpacking

- Complete edges are unpacked bottom-up, a la Langkilde (2000)
- Pruning and scoring configured via API
- At present, edges are pruned only within equivalence classes, during the unpacking stage — this ensures that pruning does not cause the realizer to fail (i.e., fail to find a complete derivation)

# Unpacking

- Complete edges are unpacked bottom-up, a la Langkilde (2000)
- Pruning and scoring configured via API
- At present, edges are pruned only within equivalence classes, during the unpacking stage — this ensures that pruning does not cause the realizer to fail (i.e., fail to find a complete derivation)
- **But**, with large grammars, considering all lexical category assignments often leads to inordinately large charts



# Anytime Best-First Search

- In the anytime best-first mode, the packing and unpacking stages are essentially interleaved
- The search can be cut off after configurable time limits, without first completing the packed chart

# Anytime Best-First Search

- In the anytime best-first mode, the packing and unpacking stages are essentially interleaved
- The search can be cut off after configurable time limits, without first completing the packed chart
- If no complete realization is found within the time limit, fragments are greedily assembled

# Greedy Fragment Assembly

- 1 Start with the best partial realization (by semantic coverage)

# Greedy Fragment Assembly

- 1 Start with the best partial realization (by semantic coverage)
- 2 Successively select the best partial realization whose semantic coverage is disjoint from those selected so far

# Greedy Fragment Assembly

- 1 Start with the best partial realization (by semantic coverage)
- 2 Successively select the best partial realization whose semantic coverage is disjoint from those selected so far
- 3 Again starting with the original best partial realization, greedily concatenate the remaining selected edges (by score), trying both orders



# Corpus Conversion

- Transform CCGbank to reflect lexicalized treatment of coordination assumed in newer, multimodal version of CCG (Baldrige and Kruijff, 2003)



# Corpus Conversion

- Transform CCGbank to reflect lexicalized treatment of coordination assumed in newer, multimodal version of CCG (Baldrige and Kruijff, 2003)
- Reanalyze punctuation along the lines of Doran's (1998) lexicalized TAG analysis



# Corpus Conversion

- Transform CCGbank to reflect lexicalized treatment of coordination assumed in newer, multimodal version of CCG (Baldrige and Kruijff, 2003)
- Reanalyze punctuation along the lines of Doran's (1998) lexicalized TAG analysis
- Change unification constraints to support semantic rather than surface syntactic dependencies (complementizers, infinitival-*to*, expletive subjects, case-marking prepositions)

## Corpus Conversion

- Transform CCGbank to reflect lexicalized treatment of coordination assumed in newer, multimodal version of CCG (Baldrige and Kruijff, 2003)
- Reanalyze punctuation along the lines of Doran's (1998) lexicalized TAG analysis
- Change unification constraints to support semantic rather than surface syntactic dependencies (complementizers, infinitival-*to*, expletive subjects, case-marking prepositions)

⇒ Viewed as a grammar engineering process! (And accordingly, implemented in a general fashion as an XSLT pipeline)



# Grammar Extraction

From the converted CCGbank, a lexico-grammar is extracted and augmented with logical forms

- Extracted categories, unary rules and lexical assignments must meet specified frequency thresholds

# Grammar Extraction

From the converted CCGbank, a lexico-grammar is extracted and augmented with logical forms

- Extracted categories, unary rules and lexical assignments must meet specified frequency thresholds
- For unseen open class words, lexical smoothing assigns most frequent categories for POS

# Grammar Extraction

From the converted CCGbank, a lexico-grammar is extracted and augmented with logical forms

- Extracted categories, unary rules and lexical assignments must meet specified frequency thresholds
- For unseen open class words, lexical smoothing assigns most frequent categories for POS
- Logical forms are inserted using around two dozen XSLT templates, with numbered semantic roles (a la PropBank)

## Example Logical Form Insertion Templates

- (1)  $s_{1:dcl} \setminus np_2 / np_3 \implies$   
 $s_{1:dcl,x1} \setminus np_{2:x2} / np_{3:x3} : @_{x1}(*\mathbf{pred}^* \wedge \langle \text{ARG0} \rangle_{x2} \wedge \langle \text{ARG1} \rangle_{x3})$
- (2)  $s_{1:pss} \setminus np_2 \implies s_{1:pss,x1} \setminus np_{2:x2} : @_{x1}(*\mathbf{pred}^* \wedge \langle \text{ARG1} \rangle_{x2})$
- (3)  $np_1 / n_1 \implies np_{1:x1} / n_{1:x1} : @_{x1}(\langle \text{DET} \rangle (d \wedge *\mathbf{pred}^*))$
- (4)  $np_1 / n_1 \setminus np_2 \implies np_{1:x1} / n_{1:x1} \setminus np_{2:x2} : @_{x1}(\langle \text{GENOWN} \rangle_{x2})$



## Creating Dev/Train/Test Files

- To obtain logical form inputs for the realizer, the extracted grammar is used to constrain parse the corpus files
- When the gold standard derivation succeeds, the resulting logical form is saved
- Sentence-internal punctuation is skipped when necessary

# Coverage

Paper/Current:

test set	LF created	single root
dev (00)	95.1% / 98.0%	67.4.% / 76.4%
test (23)	94.3% / 96.0%	69.7% / 77.2%

## Coverage

Paper/Current:

test set	LF created	single root
dev (00)	95.1% / 98.0%	67.4.% / 76.4%
test (23)	94.3% / 96.0%	69.7% / 77.2%

- LFs with multiple roots have missing dependencies

## Coverage

Paper/Current:

test set	LF created	single root
dev (00)	95.1% / 98.0%	67.4.% / 76.4%
test (23)	94.3% / 96.0%	69.7% / 77.2%

- LFs with multiple roots have missing dependencies
- Problems usually due to LF templates, but have found bugs in CCGbank

## Exact regeneration

⇒ Also helpful to look at whether sentence can be exactly regenerated with oracle n-grams, from target string

test set	grammar	complete	exact
wsj_0003	wsj_0003	86.7%	86.7%
	dev	76.7%	70.0%
	train	63.3%	56.7%
dev (00)	dev	59.1%	53.4%
	train	46.6%	37.7%

# N-gram Models

- Factored trigram models over words, part-of-speech tags and supertags
- Data from standard training sections (02–21)
- SRILM toolkit
- Null (arbitrary choice) baseline

# Word, POS and Supertag Models

$$p(\vec{F}_1^n) \approx \prod_{i=1}^n p(\vec{F}_i \mid \vec{F}_{i-2}, \vec{F}_{i-1}) \quad (1)$$

$$\begin{aligned} p^W(\vec{F}_i \mid \vec{F}_{i-2}^{i-1}) &= p(W_i \mid W_{i-2}, W_{i-1}) \\ p^P(\vec{F}_i \mid \vec{F}_{i-2}^{i-1}) &= p(P_i \mid P_{i-2}, P_{i-1}) \\ p^S(\vec{F}_i \mid \vec{F}_{i-2}^{i-1}) &= p(S_i \mid P_{i-2}, P_{i-1}) \end{aligned} \quad (2)$$

## Chained, Interpolated Models

$$p^{PS}(\vec{F}_i | \vec{F}_{i-2}^{i-1}) = p(P_i | P_{i-2}^{i-1})p(S_i | P_{i-2}^i) \quad (3)$$

$$\begin{aligned} p^{W+P}(\vec{F}_i | \vec{F}_{i-2}^{i-1}) &= \lambda_1 p^W(\vec{F}_i | \vec{F}_{i-2}^{i-1}) + \\ &\quad \lambda_2 p^P(\vec{F}_i | \vec{F}_{i-2}^{i-1}) \\ p^{W+PS}(\vec{F}_i | \vec{F}_{i-2}^{i-1}) &= \lambda_1 p^W(\vec{F}_i | \vec{F}_{i-2}^{i-1}) + \\ &\quad \lambda_2 p^{PS}(\vec{F}_i | \vec{F}_{i-2}^{i-1}) \end{aligned} \quad (4)$$



## Initial Non-Blind Development Results

scoring model	exact	BLEU
word 3g + pos 3g * stag 3g	14.8%	0.6615
word 3g + pos 3g	13.7%	0.6407
word 3g, interp. Kneser-Ney	12.2%	0.6247
word 3g, Good-Turing	11.7%	0.6219
pos 3g * supertag 3g	10.6%	0.6042
supertag 3g	10.0%	0.5886
pos 3g	8.0%	0.5413
null	5.1%	0.5251

## Initial Results With Usual Splits

test set	scoring model	exact	BLEU
dev	w3g + pos3g * stag3g	8.1%	0.5578
	word 3g + pos 3g	7.1%	0.5210
	word 3g, Kneser-Ney	6.5%	0.4872
	null	2.2%	0.3697
test	w3g + pos3g * stag3g	9.8%	0.5768
	word 3g, Kneser-Ney	6.9%	0.5178

# Updated Results With Best Model

Paper/Current:

test set	condition	exact	BLEU
dev	non-blind	14.8% / 24.3%	0.6615 / 0.7317
	usual	8.1% / 12.4%	0.5578 / 0.6101
test	usual	9.8% / 13.0%	0.5768 / 0.6223

## BLEU Comparison (PTB 23)

(N.B.: direct comparison difficult!)

	coverage	BLEU
OpenCCG (07)	96.0%	0.6223
Cahill & van Genabith (06)	98.5%	0.6651
Langkilde-Geary (02), 'Permute, no dir'	83%	0.757
Nakanishi et al. (05), $\leq 20w$	90.8%	0.7733

# Need a Supertagger for Realization!

⇒ Search errors revealed by generating with oracle n-grams (from target string), vs. best model

Oracle/Best:

test set	grammar	complete
wsj_0003	wsj_0003	86.7% / 86.7%
	dev	76.7% / 76.7%
	train	63.3% / 10.0%
dev (00)	dev	59.1% / 57.0%
	train	46.6% / 21.3%

## Continued Relevance of BLEU Scores?

⇒ Once realizations are generally satisfactory, BLEU scores may no longer be useful in measuring progress

## Continued Relevance of BLEU Scores?

⇒ Once realizations are generally satisfactory, BLEU scores may no longer be useful in measuring progress

- With MT outputs, Callison-Burch et al. (2006) contend that improved BLEU scores are neither necessary nor sufficient to achieve better human evaluation scores
- Stent et al. (2005) suggest that BLEU is a poor judge of fluency with generators that aim to produce desirable variation (e.g., in discourse)

## Example: Good

- ref.1 four of the five surviving workers have asbestos-related diseases , including three with recently diagnosed cancer .
- 0.52 four of the surviving five workers have asbestos-related diseases including three with recently diagnosed cancer .

(Score is BLEU approximation using rank order centroid weights)



## Example: Bad

- ref.2 although preliminary findings were reported more than a year ago , the latest results appear in today 's New England Journal of Medicine , a forum likely to bring new attention to the problem .
- 0.65 likely to bring new attention to the problem , today's New England Journal of Medicine in a forum the latest results appear in although preliminary findings were reported more than a year ago .

## Future Work: The Laundry List

- 1 Supertagger!
- 2 Google 1T 5-gram counts
- 3 Grammar improvements (stemming, agreement, punctuation)
- 4 PropBank integration
- 5 Perceptron tree models
- 6 Targeted human evaluations