# LCGbank

## A Corpus of Syntactic Analyses
## Based on Proof Nets

Aditya Bhargava, Timothy A. D. Fowler, and Gerald Penn

LREC-COLING 2024

# Parsing with categorial grammars

- Many related members of the categorial grammar family
  - Ajdukiewicz–Bar-Hillel grammar (ABG)
  - Combinatory categorial grammar (CCG)
  - Lambek categorial grammar (LCG)
  - Type-logical grammar (TLG)
- CCG has received (by far) the most attention in CL research
  - Many statistical parsers
  - Corpora in multiple languages

# LCG *versus* CCG

- Generative capacity
  - LCG is weakly context-free
  - CCG is context-sensitive
- Computational complexity
  - LCG parsing is NP-complete
  - CCG has known polynomial-time algorithms

# LCG *versus* CCG

- Generative capacity
  - LCG is weakly context-free, <span style="color:red">but this is rarely a problem</span>
  - CCG is context-sensitive
- Computational complexity
  - LCG parsing is NP-complete
  - CCG has known polynomial-time algorithms

# LCG *versus* CCG

- Generative capacity
  - LCG is weakly context-free, but this is rarely a problem, especially for statistical parsers
  - CCG is context-sensitive
- Computational complexity
  - LCG parsing is NP-complete
  - CCG has known polynomial-time algorithms

# LCG *versus* CCG

- Generative capacity
  - LCG is weakly context-free, but this is rarely a problem, especially for statistical parsers
  - CCG ~~is~~ can be context-sensitive
- Computational complexity
  - LCG parsing is NP-complete
  - CCG has known polynomial-time algorithms

# LCG *versus* CCG

- Generative capacity
  - LCG is weakly context-free, but this is rarely a problem, especially for statistical parsers
  - CCG ~~is~~ can be context-sensitive, but in practice never is
- Computational complexity
  - LCG parsing is NP-complete
  - CCG has known polynomial-time algorithms

# LCG *versus* CCG

- Generative capacity
    - LCG is weakly context-free, but this is rarely a problem, especially for statistical parsers
    - CCG ~~is~~ can be context-sensitive, but in practice never is
- Computational complexity
    - LCG parsing is NP-complete, but polynomial assuming reasonable bounds on categories
    - CCG has known polynomial-time algorithms

# LCG *versus* CCG

- Generative capacity
  - LCG is weakly context-free, but this is rarely a problem, especially for statistical parsers
  - CCG ~~is~~ can be context-sensitive, but in practice never is
- Computational complexity
  - LCG parsing is NP-complete, but polynomial assuming reasonable bounds on categories
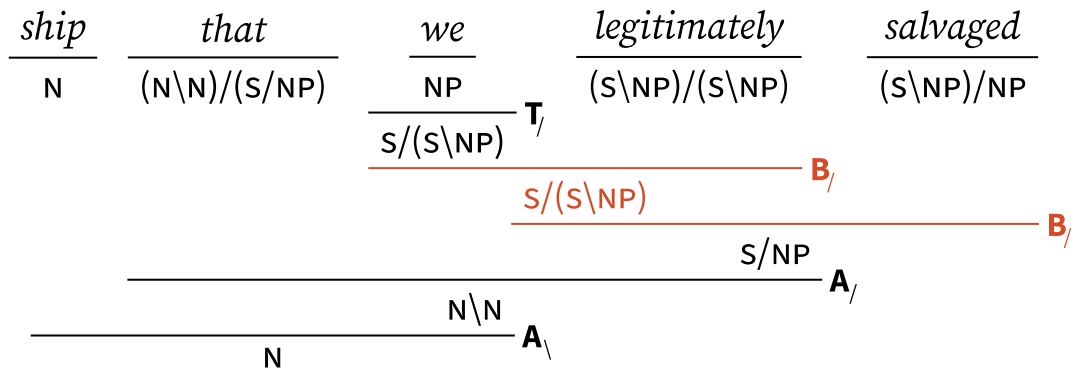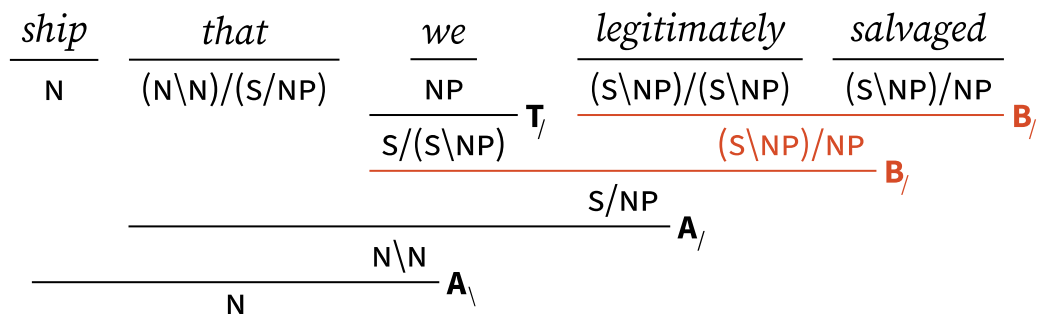  - CCG has known polynomial-time algorithms
  - Not so relevant to statistical parsers

# Why LCG?

# Why LCG?
# Proof nets!

# Spurious ambiguities

| ship | that | we | legitimately | salvaged |
|------|------|-----|--------------|----------|
| N | (N\N)/(S/NP) | NP | (S\NP)/(S\NP) | (S\NP)/NP |

$\dfrac{\text{NP}}{\text{S/(S\NP)}}$ **T**$_/$

$\dfrac{\text{(S\NP)/(S\NP)} \quad \text{(S\NP)/NP}}{\text{(S\NP)/NP}}$ **B**$_/$

$\dfrac{}{\text{(S\NP)/NP}}$ **B**$_/$

$\dfrac{}{\text{S/NP}}$ **A**$_/$

$\dfrac{}{\text{N\N}}$ **A**$_\backslash$

$\dfrac{}{\text{N}}$

---

| ship | that | we | legitimately | salvaged |
|------|------|-----|--------------|----------|
| N | (N\N)/(S/NP) | NP | (S\NP)/(S\NP) | (S\NP)/NP |

$\dfrac{\text{NP}}{\text{S/(S\NP)}}$ **T**$_/$

$\dfrac{}{\text{S/(S\NP)}}$ **B**$_/$

$\dfrac{}{\text{S/(S\NP)}}$ **B**$_/$

$\dfrac{}{\text{S/NP}}$ **A**$_/$

$\dfrac{}{\text{N\N}}$ **A**$_\backslash$

$\dfrac{}{\text{N}}$

# Spurious ambiguities

# Spurious ambiguities

- CCG parsers use normal-form constraints during parsing
- Proof nets represent LCG derivations such that semantically-equivalent derivations correspond to the same proof net
- CCG is incompatible

# Statistical parsing with proof nets

- Recent interest in statistical parsing with proof nets
  - Involves a more structured treatment of lexical categories
- Can even train *without* ground-truth derivations

# Statistical parsing with proof nets

- Recent interest in statistical parsing with proof nets
  - Involves a more structured treatment of lexical categories
- Can even train *without* ground-truth derivations

…but very few corpora, and none in English!

# LCGbank

# LCGbank development process

- Start from CCGbank
- Discard category features
- Lexicalize incompatible CCG rules for LCG compatibility
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- **Discard category features**
- Lexicalize incompatible CCG rules for LCG compatibility
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- Discard category features
- **Lexicalize incompatible CCG rules for LCG compatibility**
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- Discard category features
- Lexicalize incompatible CCG rules for LCG compatibility
- **Lexicalize incompatible CCG*bank* rules (e.g., type-changing)**
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- Discard category features
- Lexicalize incompatible CCG rules for LCG compatibility
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- Discard category features
- Lexicalize incompatible CCG rules for LCG compatibility
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- Discard category features
- Lexicalize incompatible CCG rules for LCG compatibility
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# LCGbank development process

- Start from CCGbank
- Discard category features
- Lexicalize incompatible CCG rules for LCG compatibility
- Lexicalize incompatible CCG*bank* rules (e.g., type-changing)
- Converted (LCG) derivations then specify the proof net
- New derivations for left-out sentences
- Result: fully compositional semantics

# CCGbank's type-changing

$$
\begin{array}{c}
\dfrac{issues}{\text{NP}} \qquad
\dfrac{\dfrac{they}{\text{NP}}}{\text{S/(S\textbackslash NP)}}\,\mathbf{T}_/ \qquad
\dfrac{floated}{\text{(S\textbackslash NP)/NP}} \qquad
\dfrac{in}{\text{((S\textbackslash NP)\textbackslash (S\textbackslash NP))/NP}} \qquad
\dfrac{\dfrac{1980}{\dfrac{\text{N}}{\text{NP}}\,\text{TC}}}{\ }
\end{array}
$$

# CCGbank's type-changing

$$
\begin{array}{c}
\dfrac{\textit{issues}}{\text{NP}} \qquad
\dfrac{\dfrac{\textit{they}}{\text{NP}}}{\text{S/(S\backslash NP)}}\,\mathbf{T}_{/} \qquad
\dfrac{\textit{floated}}{\text{(S\backslash NP)/NP}} \qquad
\dfrac{\textit{in}}{\text{((S\backslash NP)\backslash(S\backslash NP))/NP}} \qquad
\dfrac{\dfrac{\textit{1980}}{\dfrac{\text{N}}{\text{NP}}\,\text{TC}}}{}
\end{array}
$$

$$
\dfrac{\phantom{(S\backslash NP)}}{\text{(S\backslash NP)\backslash(S\backslash NP)}}\,\mathbf{A}_{/}
$$

$$
\dfrac{\phantom{xxxxxxxxxxxxxx}}{\text{(S\backslash NP)\backslash(S\backslash NP)}}\times\!\mathbf{B}_{\backslash}
$$

$$
\dfrac{\phantom{xxxxxxxxxxxxxxxxxx}}{\text{(S\backslash NP)/NP}}\,\mathbf{B}_{/}
$$

$$
\dfrac{\text{S/NP}}{\dfrac{\text{NP\backslash NP}}{\text{NP}}\,\mathbf{A}_{\backslash}}\,\text{TC}
$$

# CCGbank's type-changing

$$
\frac{issues}{\text{NP}} \quad
\frac{\dfrac{they}{\text{NP}}}{\text{S/(S\textbackslash NP)}}\ \mathbf{T}_{/} \quad
\frac{floated}{\text{(S\textbackslash NP)/NP}} \quad
\frac{in}{\text{((S\textbackslash NP)\textbackslash(S\textbackslash NP))/NP}} \quad
\frac{1980}{\text{NP}}
$$

(S\NP)\(S\NP)  ×**B**\

(S\NP)/NP  **B**/

S/NP  TC

NP\NP  **A**\

NP

$$\frac{1980}{\text{NP}}\ \mathbf{A}_{/}$$

# CCGbank's type-changing

$$\frac{\textit{issues}}{\text{NP}} \quad \frac{\dfrac{\textit{they}}{\text{NP}}}{\text{S/(S\textbackslash NP)}}\mathbf{T}_/ \quad \frac{\textit{floated}}{\text{(S\textbackslash NP)/NP}} \quad \frac{\dfrac{\textit{in}}{\text{((S\textbackslash NP)\textbackslash (S\textbackslash NP))/NP}} \quad \dfrac{\textit{1980}}{\text{NP}}}{\text{(S\textbackslash NP)\textbackslash (S\textbackslash NP)}}\mathbf{A}_/$$

issues    NP

they    NP    S/(S\NP)    **T**/

floated    (S\NP)/NP

in    ((S\NP)\(S\NP))/NP

1980    NP    **A**/

(S\NP)\(S\NP)    ×**B**\

(S\NP)/NP    **B**/

S/NP    TC
NP\NP    **A**\

NP

# CCGbank's type-changing

$$
\begin{array}{c}
\dfrac{ship}{\text{N}} \quad \dfrac{that}{\text{(N\textbackslash N)/(S/NP)}} \quad \dfrac{\dfrac{\dfrac{we}{\text{NP}}}{\text{S/(S\textbackslash NP)}}\,\mathbf{T}_{/}}{} \quad \dfrac{\dfrac{legitimately}{\text{(S\textbackslash NP)/(S\textbackslash NP)}} \quad \dfrac{salvaged}{\text{(S\textbackslash NP)/NP}}}{\text{(S\textbackslash NP)/NP}}\,\mathbf{B}_{/}
\end{array}
$$

ship — N

that — (N\N)/(S/NP)

we — NP → S/(S\NP) **T**/

legitimately — (S\NP)/(S\NP)

salvaged — (S\NP)/NP

(S\NP)/NP **B**/

(S\NP)/NP → S/NP **B**/

S/NP **A**/

N\N **A**\

N

# CCGbank's type-changing

$$
\frac{ship}{N} \quad \frac{that}{\text{(N\textbackslash N)/(S/NP)}} \quad \frac{\dfrac{\dfrac{we}{NP}}{S/(S\textbackslash NP)} \, \textbf{T}_{/} \quad \dfrac{\dfrac{legitimately}{\text{(S\textbackslash NP)/(S\textbackslash NP)}} \quad \dfrac{salvaged}{\text{(S\textbackslash NP)/NP}}}{\text{(S\textbackslash NP)/NP}} \, \textbf{B}_{/}}{\dfrac{\dfrac{S/NP}{} \, \textbf{A}_{/}}{\dfrac{N\textbackslash N}{N} \, \textbf{A}_{\backslash}} \, \textbf{B}_{/}}
$$

# CCGbank's type-changing

$$\frac{issues}{NP} \quad \frac{\dfrac{they}{NP}}{S/(S\backslash NP)} \textbf{T}_{/} \quad \frac{floated}{(S\backslash NP)/NP} \quad \frac{\dfrac{in}{((S\backslash NP)\backslash(S\backslash NP))/NP} \quad \dfrac{1980}{NP}}{(S\backslash NP)\backslash(S\backslash NP)} \textbf{A}_{/}$$

issues — NP

they — NP → S/(S\NP)   **T**/

floated — (S\NP)/NP

in — ((S\NP)\(S\NP))/NP    1980 — NP    **A**/

(S\NP)\(S\NP)    ×**B**\

(S\NP)/NP    **B**/

<span style="color:#c0392b">S/NP</span>

<span style="color:#c0392b">NP\NP</span>  TC

**A**\

NP

# CCGbank's type-changing

$$\frac{issues}{NP} \qquad \frac{\dfrac{they}{NP}}{S/(S\backslash NP)}\textbf{T}_{/} \qquad \frac{floated}{(S\backslash NP)/NP} \qquad \frac{\dfrac{in}{((S\backslash NP)\backslash(S\backslash NP))/NP} \quad \dfrac{1980}{NP}}{(S\backslash NP)\backslash(S\backslash NP)}\textbf{A}_{/}$$

issues — NP

they — NP — S/(S\NP) **T**/

floated — (S\NP)/NP

in — ((S\NP)\(S\NP))/NP

1980 — NP — **A**/

(S\NP)\(S\NP) ×**B**\

(S\NP)/NP **B**/

S/NP — TC

NP\NP — **A**\

NP

# CCGbank's type-changing



$$issues \quad they \quad floated \quad in \quad 1980$$

$$NP \quad NP \quad (S\backslash NP)/NP \quad ((S\backslash NP)\backslash(S\backslash NP))/NP \quad NP$$

$$NP^- \quad NP^- \quad NP^+ \quad NP^+ \quad NP^- \quad NP^+ \quad NP^- \quad NP^+ \quad NP^- \quad NP^+ \quad NP^- \quad NP^+$$

$$NP \quad NP \quad (NP\backslash NP)\backslash NP \quad ((NP\backslash NP)\backslash(NP\backslash NP))/NP \quad NP \quad NP$$

$$issues \quad they \quad floated \quad in \quad 1980$$

# Crossed composition

$$\frac{\textit{issues}}{\text{NP}} \qquad \frac{\frac{\textit{they}}{\text{NP}}}{\text{S/(S\textbackslash NP)}} \textbf{T}_/ \qquad \frac{\textit{floated}}{\text{(S\textbackslash NP)/NP}} \qquad \frac{\textit{in}}{\text{((S\textbackslash NP)\textbackslash(S\textbackslash NP))/NP}} \qquad \frac{\frac{\textit{1980}}{\text{N}}}{\text{NP}} \text{TC}$$

issues — NP

they — NP — S/(S\NP) **T**/

floated — (S\NP)/NP

in — ((S\NP)\(S\NP))/NP

1980 — N — NP — TC

(S\NP)\(S\NP) — **A**/

(S\NP)/NP — ×**B**\

S/NP — **B**/

NP\NP — TC

S/NP

NP — **A**\

# Crossed composition

$$\frac{\cfrac{\textit{would}}{(S\backslash NP)/(S\backslash NP)} \quad \cfrac{\textit{temporarily}}{(S\backslash NP)\backslash(S\backslash NP)}}{(S\backslash NP)/(S\backslash NP)} \times\mathbf{B}_{\backslash} \qquad \cfrac{\cfrac{\textit{dilute}}{(S\backslash NP)/NP} \quad \cfrac{\cfrac{\textit{earnings}}{\cfrac{N}{NP}\; TC}}{S\backslash NP}\mathbf{A}_{/}}{} }{S\backslash NP}\mathbf{A}_{/}$$

# Crossed composition

$$\frac{\displaystyle \frac{\displaystyle \frac{would}{(S\backslash NP)/(S\backslash NP)} \quad \frac{temporarily}{{\color{orange}(S\backslash NP)\backslash(S\backslash NP)}}}{(S\backslash NP)/(S\backslash NP)} \times \mathbf{B}_\backslash \quad \frac{\displaystyle \frac{dilute}{(S\backslash NP)/NP} \quad \frac{\displaystyle \frac{earnings}{N}}{NP}\, TC}{S\backslash NP}\, \mathbf{A}_/}{S\backslash NP}\, \mathbf{A}_/$$

# Crossed composition

# Coordination

- CCGbank assigns a special conj category to coordinators ($and$, $or$, etc.)
- Usual interpretation is as polymorphic $(x\backslash x)/x$
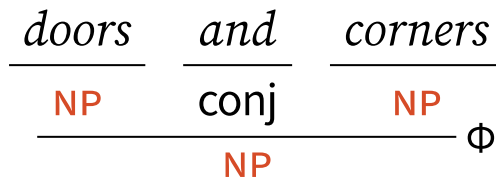    - But only for *like* coordination

# Coordination

- CCGbank assigns a special conj category to coordinators (*and*, *or*, etc.)
- Usual interpretation is as polymorphic $(x \backslash x)/x$
  - But only for *like* coordination

$$\frac{\displaystyle \frac{doors}{NP} \quad \frac{and}{conj} \quad \frac{corners}{NP}}{NP} \Phi$$
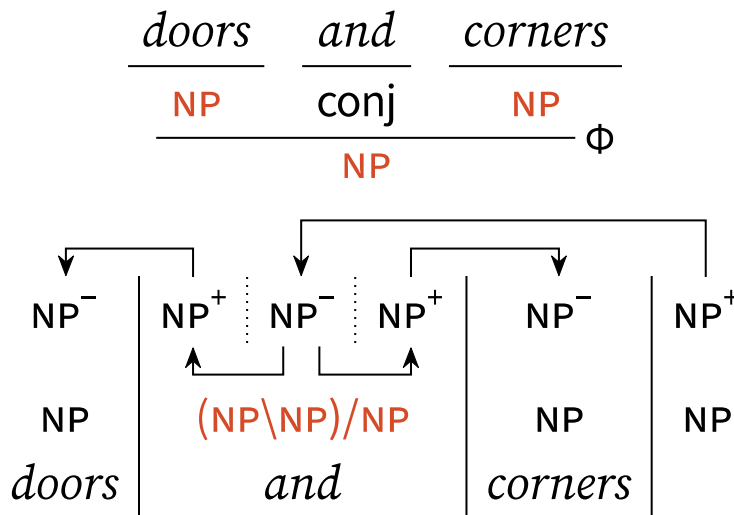
# Coordination

- CCGbank assigns a special conj category to coordinators (*and*, *or*, etc.)
- Usual interpretation is as polymorphic $(x \backslash x)/x$
  - But only for *like* coordination

$$\frac{\displaystyle \frac{doors}{NP} \quad \frac{and}{conj} \quad \frac{corners}{NP}}{NP} \Phi$$

# Coordination

- CCGbank assigns a special conj category to coordinators (*and*, *or*, etc.)
- Usual interpretation is as polymorphic $(x\backslash x)/x$
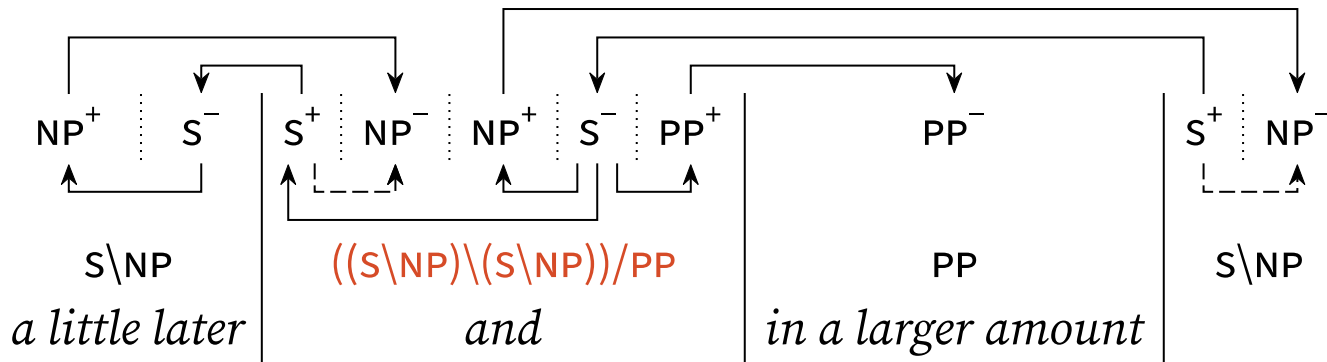  - But only for *like* coordination

$$
\begin{array}{ccc}
\underline{doors} & \underline{and} & \underline{corners} \\
\text{NP} & \text{conj} & \text{NP} \\
\end{array}
$$

$$\text{NP} \quad \Phi$$

$$
\begin{array}{cccccc}
\text{NP}^- & \text{NP}^+ & \text{NP}^- & \text{NP}^+ & \text{NP}^- & \text{NP}^+ \\
\text{NP} & & \text{(NP\backslash NP)/NP} & & \text{NP} & \text{NP} \\
doors & & and & & corners & \\
\end{array}
$$

# Coordination: *unlike*

$$\frac{\frac{a \quad little \quad later}{S\backslash NP} \quad \frac{and}{conj} \quad \frac{in \quad a \quad larger \quad amount}{PP}}{S\backslash NP} \text{magic}$$

# Coordination: *unlike*

$$a \quad little \quad later \quad\quad and \quad\quad in \quad a \quad larger \quad amount$$

| S\NP | conj | PP |

S\NP      magic



NP⁺ ┆ S⁻ ┆ S⁺ NP⁻ ┆ NP⁺ S⁻ ┆ PP⁺      PP⁻      ┆ S⁺ ┆ NP⁻

S\NP          ((S\NP)\(S\NP))/PP              PP              S\NP

*a little later*              *and*              *in a larger amount*

# Manual annotations

- 274 sentences from PTB omitted from CCGbank (e.g., sentential gapping)
  - We parsed these with a CCG parser and adjusted manually

# Manual annotations

- 274 sentences from PTB omitted from CCGbank (e.g., sentential gapping)
  - We parsed these with a CCG parser and adjusted manually
- ~500 rules not accounted for by our conversion rules
  - We annotated these manually; most were annotation errors
- 40 sentences with links *within* lexical categories
  - We provide *additional* analyses without using these links

# LCGbank: the corpus

- ~49k sentences with analyses
- Use all sections for data splits
- Released as a set of conversion scripts & data
  - Apache 2.0 license
  - CCGbank is required

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---:|---:|---:|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---|---|---|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---|---|---|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---|---|---|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---|---|---|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---|---|---|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# Training set statistics

| | CCGbank | CCGbank w/o feats | LCGbank |
|---|---|---|---|
| Sentences | 44,614 | 44,614 | 44,870 |
| Atomic categories | 34 | 11 | 5 |
| Lexical categories | 1,327 | 487 | 1,071 |
| Avg. cat. order | 1.748 | 1.916 | 2.317 |
| Avg. cats/word | 1.701 | 1.577 | 1.947 |
| Exp. cats/word | 20.083 | 14.958 | 29.731 |

# What we didn't get to do

- Porting CCGbank dependencies
    - Motivation partially taken care of by proof nets
    - But still useful for evaluation
    - No good evaluation in place for statistical LCG parsing

# What we didn't get to do

- CCGrebank
  - Numerous improvements over CCGbank
  - Not readily available 😢

# LCGbank

## A Corpus of Syntactic Analyses Based on Proof Nets

Aditya Bhargava, Timothy A. D. Fowler, and Gerald Penn
LREC-COLING 2024