# Finite-state Approximation of Constraint-based Grammars using Left-corner Grammar Transforms
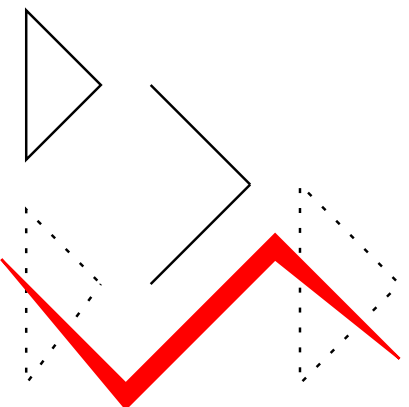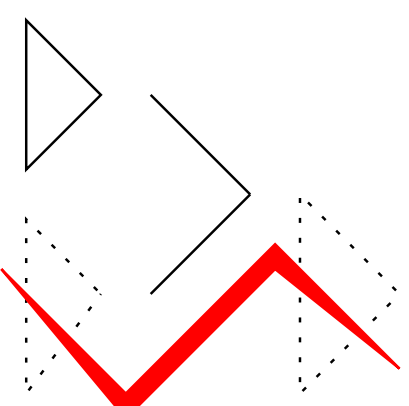
Mark Johnson

Brown University

ACL-COLING 1998

# Summary

- Approximating a Unification Grammar (UG) with a FSM

- FS approximations of top-down parsers

- Grammar transformation

  – Left-Corner (LC)

  – Composition/$\epsilon$-removal

  – Partial evaluation

2

# Why approximate UGs with FSMs?

- FSM processing is faster

  – linear time recognition

  – can be used as oracle to guide UG parser

- LC parsing has some psycholinguistic validity

- UG languages can be manipulated via FS calculus

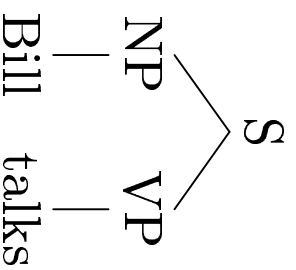# Why use LC approximation?

- LC parsing applies directly to UGs

- LC parsers require only finite stack-depth to parse left linear or right linear grammars

# Non-deterministic top-down parsing

- Parser states are stacks of nonterminals and terminals $(N \cup T)^*$

- State transition function $\delta$:

$$\gamma \in \delta(a\gamma, a) : a \in T, \gamma \in (N \cup T)^*.$$

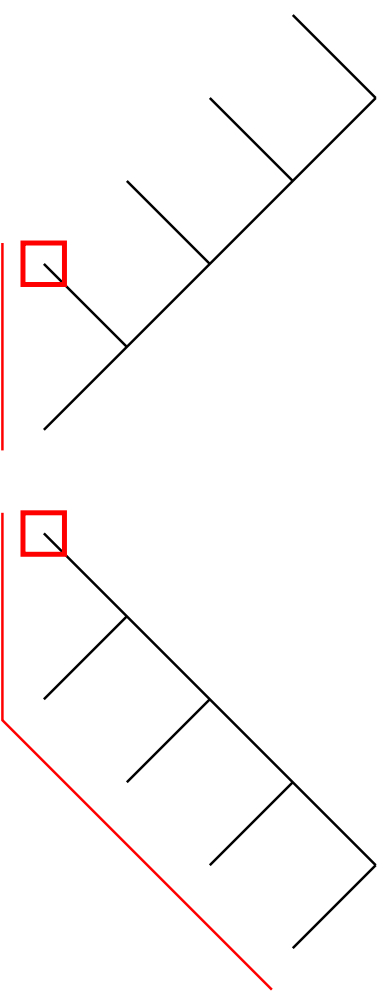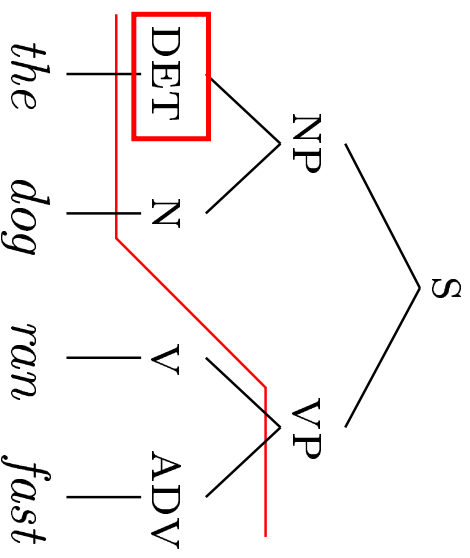$$\beta\gamma \in \delta(A\gamma, \epsilon) : A \in N, \gamma \in (N \cup T)^*, A \to \beta \in P.$$

| State | Remaining input |
|-------|-----------------|
| S     | *NP VP*         |
| NP VP | *NP VP*         |
| VP    | *VP*            |
| $\epsilon$ | $\epsilon$ |

```
        S
       / \
     NP   VP
      |    |
    Bill  talks
```

# FS approximations to TD states

- Unbounded state stack size

  − ignore state stacks larger than some fixed bound

  ⇒ approximation accepts a *subset* of UG language

  − collapse all states sharing a common prefix

  ⇒ approximation accepts a *superset* of UG language

- Unbounded UG categories

  − Restriction (a.k.a. abstraction) (Shieber 1985)

  ⇒ approximation accepts a *superset* of UG language

  − Ignore categories whose complexity exceeds some bound

  ⇒ approximation accepts a *subset* of UG language

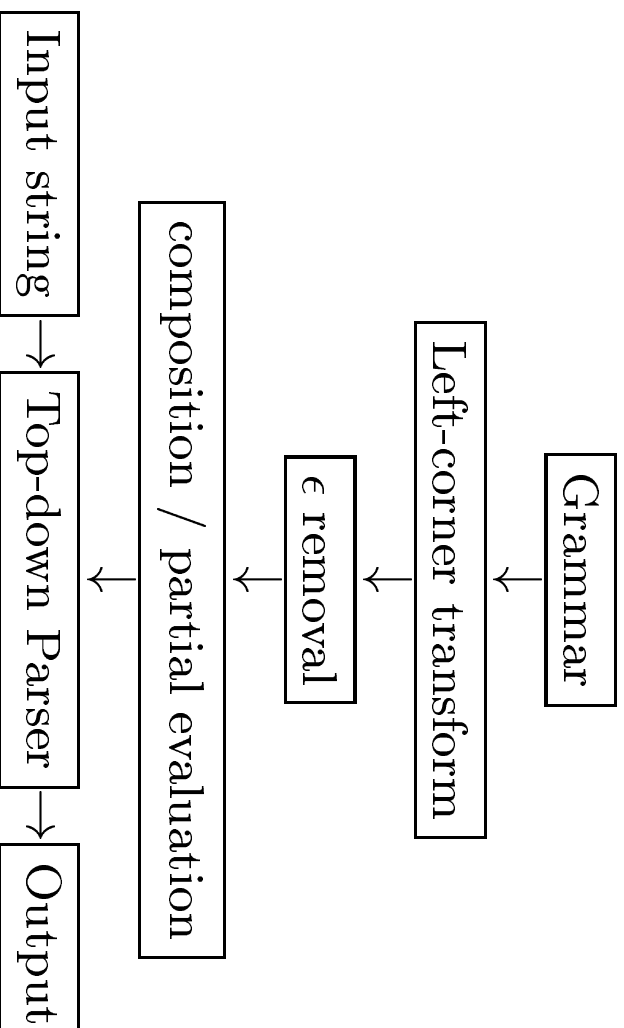  ∗ In many UGs, the *syntactically potent features* range over finite values

# States of a TD parser

- Just before $X$ is expanded, the TD parser's state consists of $X$ followed by *the right siblings of it and all its ancestors.*

$\Rightarrow$ Right-linear grammars ($A \to w\,B$) require *finite* state size

$\Rightarrow$ Left-linear grammars ($A \to B\,w$) require *unbounded* state size

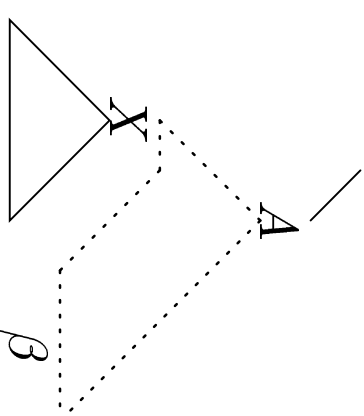# Left-corner grammar transforms

- A Left-Corner (LC) parser exhibits finite state size on both left-linear and right-linear CFGs (*)

- A LC parser for grammar $G$ acts isomorphically to a top-down parser using $\mathcal{LC}(G)$.

```
Grammar → Left-corner transform → ε removal → composition / partial evaluation →

Input string → Top-down Parser → Output
```

# Left-corner grammar transform

- Left-corner of each production is recognized bottom-up, everything else is predicted top-down

- Nonterminals of $\mathcal{LC}(G) = N \cup N \times (N \cup T)$

$$A \overset{*}{\Rightarrow}_G X\beta \text{ iff } A\text{-}X \overset{*}{\Rightarrow}_{\mathcal{LC}(G)} \beta$$



- Productions of $\mathcal{LC}(G) =$

$$
\begin{aligned}
A &\to a\,A\text{-}a & &: & &A \in N, a \in T. \\
A\text{-}X &\to \beta\,A\text{-}B & &: & &A \in N, B \to X\beta \in P. \\
A\text{-}A &\to \epsilon & &: & &A \in N.
\end{aligned}
$$

# Parsing with $\mathcal{LC}(G)$: start

NP
DET — *the*
N — *dog*

S

VP
V — *ran*
ADV — *fast*

S

S—DET
DET — *the*
N — *dog*
S—NP
VP
V — *ran*
VP—V
ADV — *fast*
VP—VP
S—S

$A \to aA{-}a \ : \ A \in N, a \in T.$

# Parsing with $\mathcal{LC}(G)$: shift DET



$$A \to aA\text{-}a \quad : \quad A \in N, a \in T.$$

$$A\text{-}X \to \beta A\text{-}B \quad : \quad B \to X\beta \in P.$$

# Parsing with $\mathcal{LC}(G)$: NP



$A\text{-}X \to \beta A\text{-}B \ : \ B \to X\beta \in P.$

# Parsing with $\mathcal{LC}(G)$: S



$A\text{-}X \rightarrow \beta A\text{-}B \quad : \quad B \rightarrow X\beta \in P.$

$A\text{-}A \rightarrow \epsilon \qquad : \quad A \in N.$

# States of an LC parser

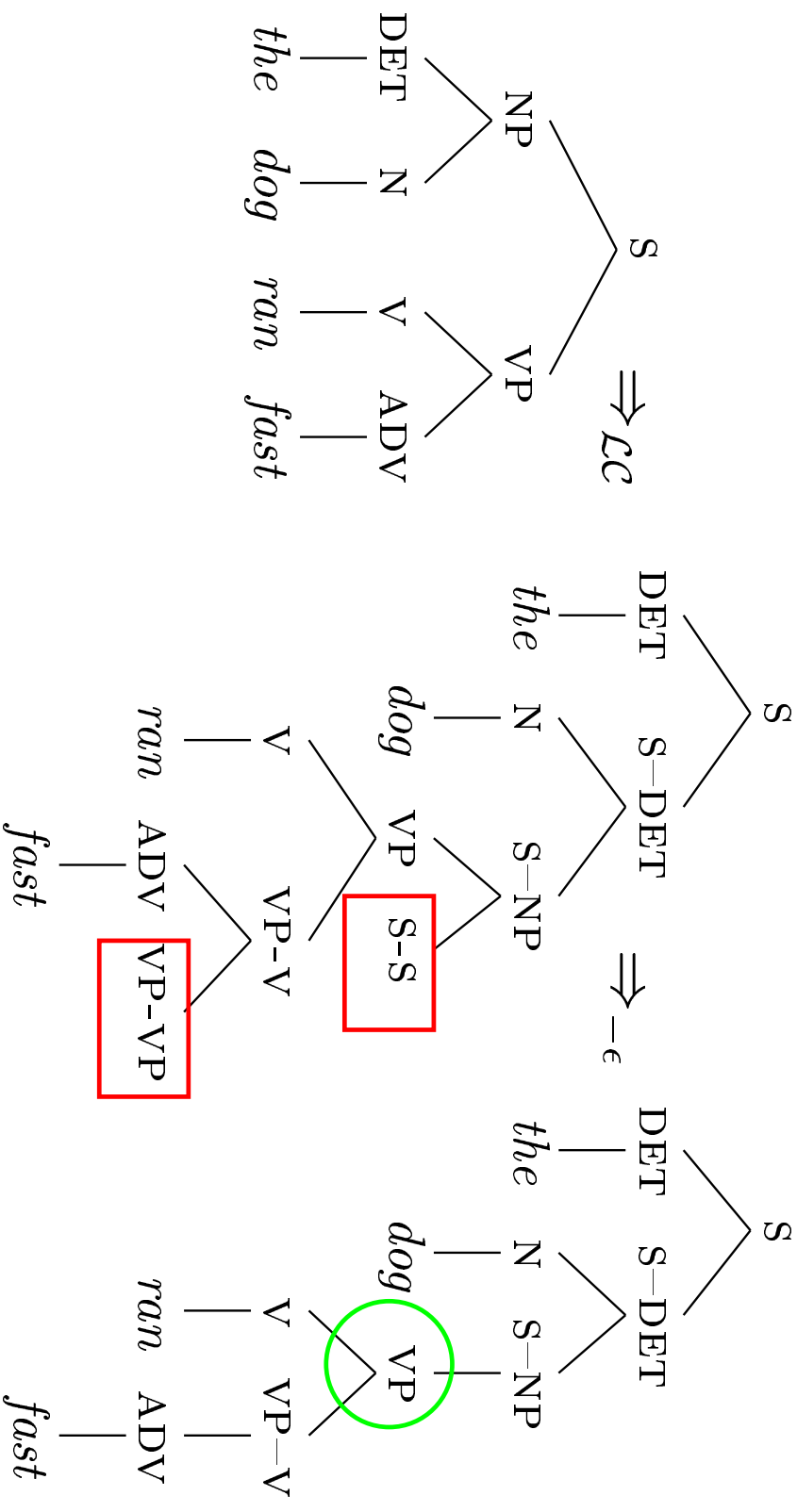- Left-linear $G$ $\Rightarrow$ right-linear $\mathcal{LC}(G)$ $\Rightarrow$ finite states



$\Rightarrow_{\mathcal{LC}}$

# States of an LC parser (cont.)

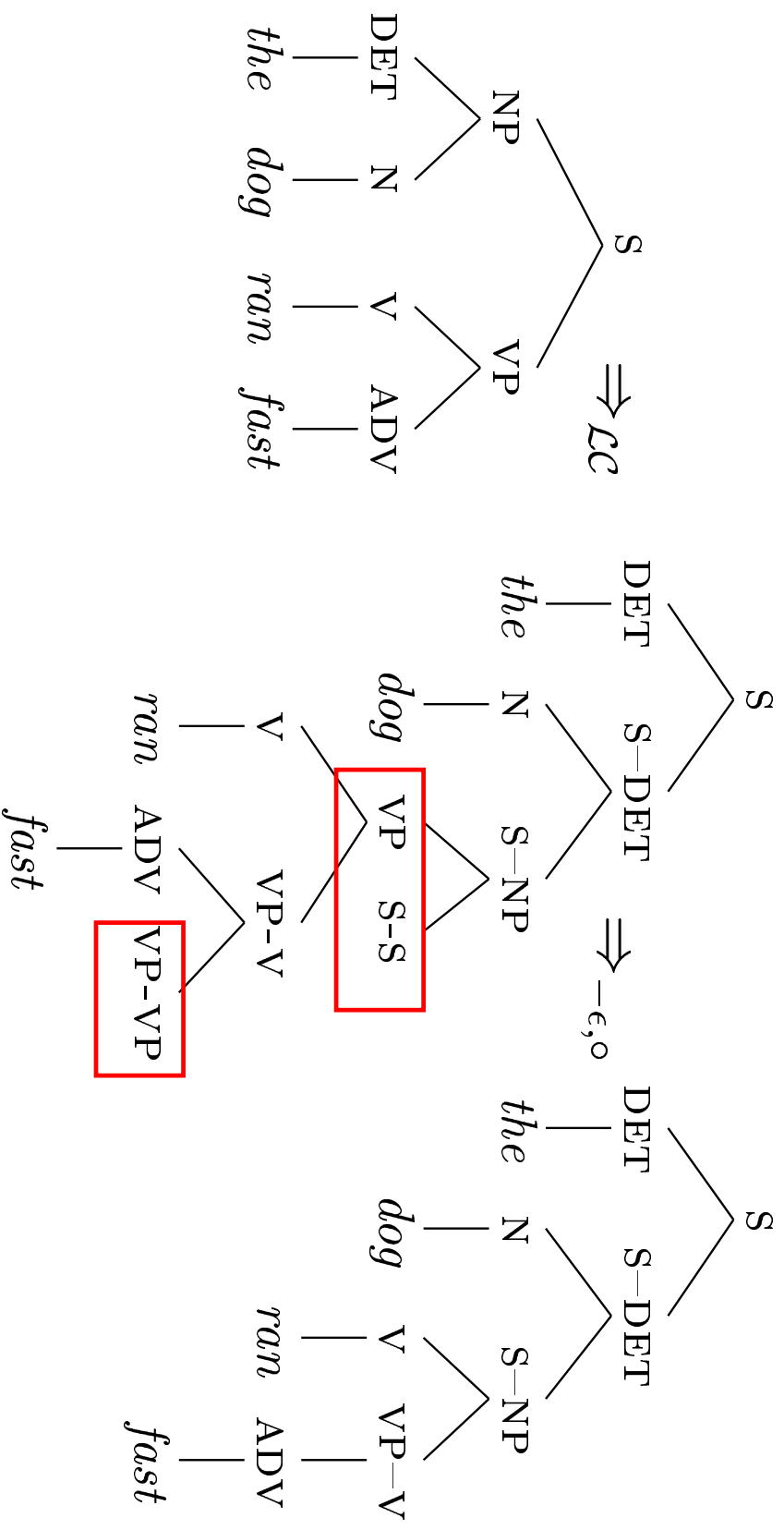- Right-linear $G$ ⇒ unbounded TD states in $\mathcal{LC}(G)$



$$\Rightarrow_{\mathcal{LC}}$$

# Epsilon-removal after $\mathcal{LC}$ transform

- Linear $G \Rightarrow$ right-linear $\mathcal{LC}'(G) \Rightarrow$ finite TD states

$$\Rightarrow_{\mathcal{LC}}$$

$$\Rightarrow_{-\epsilon}$$

# Partial evaluation/composition

- Converts binary branches into (almost) binary branches

# Special case of binary productions

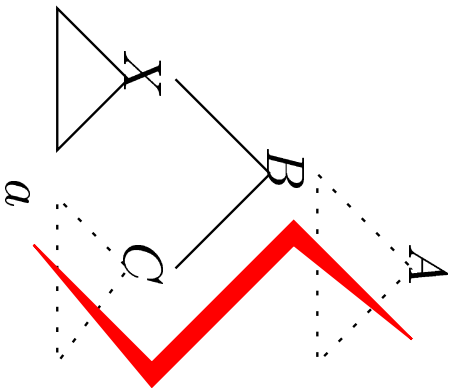$$S \rightarrow aSa \qquad : a \in T.$$
$$A\text{-}X \rightarrow aA\text{-}B \quad : A \in N, B \rightarrow Xa \in P.$$
$$A\text{-}X \rightarrow a \qquad : A \rightarrow Xa \in P.$$
$$A\text{-}X \rightarrow aC\text{-}a \qquad : A \rightarrow XC \in P.$$
$$A\text{-}X \rightarrow aC\text{-}aA\text{-}B \quad : A \in N, B \rightarrow XC \in P.$$

- All but one schema are right-linear

- Exactly one transformed rule per input item

- Such productions can be implemented as FSM arcs, e.g.:

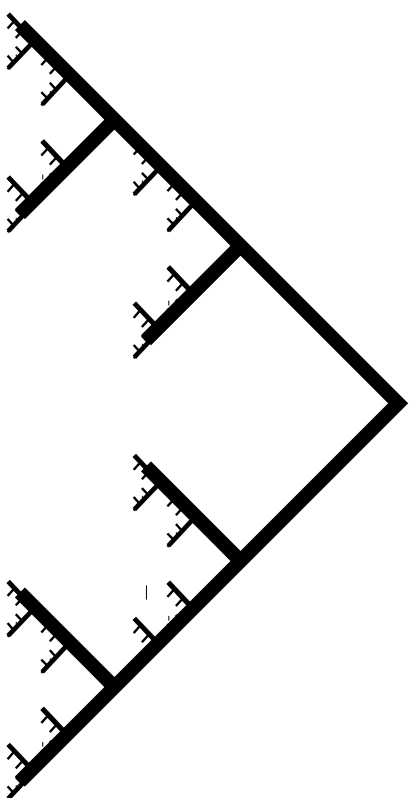$$A\text{-}B\beta \in \delta(A\text{-}X\beta, a) : B \rightarrow Xa \in P.$$

# Geometry of LC state complexity

- Because only one production schema increases the stack state size, LC state complexity is associated with a specific tree geometry

- Helps characterize the errors in a FS approximation

$A\text{-}X \rightarrow a\,C\text{-}a\,A\text{-}B$

# Odds and ends

- Classifying unification grammar categories

- Identifying useless productions in $\mathcal{LC}(G)$ (link table)

- Obtaining parse trees from FSM transitions

  – FS transducer emits rule schema used at each transition,

  – which guides LC parser for $G'$

# Conclusion

- Left-corner grammar transforms convert left recursion into right recursion

- A finite-state approximation can be directly constructed from transformed unification grammars

- The approximation is exact for left linear and right linear CFGs

- A characterization of LC state complexity identifies constructions for which the approximation is inexact