

Advances in Abstract Categorial Grammars

Language Theory and Linguistic Modeling

Lecture 5

Application of magic-sets rewriting
Prefix-correct Earley recognizer from Datalog program representing MCFG

Early parsing

- From CFGs to MCFGs
- Indirectly applies to string-generating second-order ACGs

Polynomial-time algorithm

Seminaive(P, D)

```
1 agenda[0]  $\leftarrow D$ 
2  $i \leftarrow 0$ 
3 chart  $\leftarrow \emptyset$ 
4 while agenda[ $i$ ]  $\neq \emptyset$ 
5   do
6     chart  $\leftarrow \text{chart} \cup \text{agenda}[i]$ 
7     agenda[ $i + 1$ ]  $\leftarrow \text{Conseq}(P, \text{agenda}[i], \text{chart}) - \text{chart}$ 
8      $i \leftarrow i + 1$ 
9 return chart
```

holds facts with
derivation tree of
minimal height i

facts that immediately follow
from one fact in *agenda*[i] plus
some facts in *chart*

\approx well-formed
substring table

Seminaive is a purely bottom-up tabular algorithm.
This algorithm works for Datalog programs in general.
The size of chart is $O(n^r)$ ($r = \text{maximal arity}$)
Running time $O(n^k)$ depends on the number of variables in rules.

$S \rightarrow NPVP$	$Aux \rightarrow does$
$S \rightarrow Aux \ NP \ VP$	$Det \rightarrow that \mid this \mid a$
$S \rightarrow VP$	$Name \rightarrow Houston \mid TWA$
$NP \rightarrow Det \ NI$	$Pronoun \rightarrow I \mid she \mid me$
$NP \rightarrow Name$	$V \rightarrow book \mid include \mid prefer$
$NP \rightarrow Pronoun$	$N \rightarrow book \mid flight \mid meal \mid money$
$VP \rightarrow V$	$P \rightarrow from \mid to \mid on$
$VP \rightarrow V \ NP$	
$VP \rightarrow V \ NP \ PP$	
$VP \rightarrow VP \ PP$	
$NI \rightarrow N$	
$NI \rightarrow Name \ NI$	
$NI \rightarrow NI \ PP$	
$PP \rightarrow P \ NP$	

An example of a context-free grammar.
Adapted from Jurafsky and Martin.

$S(i, k) :- NP(i, j), VP(j, k).$

$S(i, l) :- Aux(i, j), NP(j, k), VP(k, l).$

$S(i, j) :- VP(i, j).$

$NP(i, k) :- Det(i, j), NI(j, k).$

$NP(i, j) :- Name(i, j).$

$NP(i, j) :- Pronoun(i, j).$

$VP(i, j) :- V(i, j).$

$VP(i, k) :- V(i, j), NP(j, k).$

$VP(i, l) :- V(i, j), NP(j, k), PP(k, l).$

$VP(i, k) :- VP(i, j), PP(j, k).$

$NI(i, j) :- N(i, j).$

$NI(i, k) :- NI(i, j), N(j, k).$

$NI(i, k) :- NI(i, j), PP(j, k).$

$PP(i, k) :- P(i, j), NP(j, k).$

$Aux(i, j) :- does(i, j).$

$Det(i, j) :- that(i, j).$

$Det(i, j) :- this(i, j).$

$Det(i, j) :- a(i, j).$

$Name(i, j) :- Houston(i, j).$

$Name(i, j) :- TWA(i, j).$

$Pronoun(i, j) :- I(i, j).$

$Pronoun(i, j) :- she(i, j).$

$Pronoun(i, j) :- me(i, j).$

$V(i, j) :- book(i, j).$

$V(i, j) :- include(i, j).$

$V(i, j) :- prefer(i, j).$

$N(i, j) :- book(i, j).$

$N(i, j) :- flight(i, j).$

$N(i, j) :- meal(i, j).$

$N(i, j) :- money(i, j).$

$P(i, j) :- from(i, j).$

$P(i, j) :- to(i, j).$

$P(i, j) :- on(i, j).$

Datalog program representing CFG.

Running time of Seminaive is $O(N^4)$

$D = \{ \text{book}(0, 1), \text{the}(1, 2), \text{flight}(2, 3), \text{from}(3, 4), \text{Houston}(4, 5) \}.$

$\text{agenda}[0] = \{ \text{book}(0, 1), \text{the}(1, 2), \text{flight}(2, 3), \text{from}(3, 4), \text{Houston}(4, 5) \}.$

$\text{agenda}[1] = \{ N(0, 1), V(0, 1), \text{Det}(1, 2), N(2, 3), P(3, 4), \text{Name}(4, 5) \}$

$\text{agenda}[2] = \{ NI(0, 1), VP(0, 1), NI(2, 3), NP(4, 5) \}$

$\text{agenda}[3] = \{ NP(1, 3), PP(3, 5) \}$

$\text{agenda}[4] = \{ VP(0, 3), VP(0, 5), NI(2, 5) \}$

$\text{agenda}[5] = \{ S(0, 3), S(0, 5), NP(1, 5) \}$

$\text{agenda}[6] = \emptyset$

$chart = agenda[0] \cup \dots \cup agenda[6]$

The run of Seminaive on “book the flight from Houston”

Outputting shared forest

Seminaive_parse(P, D)

```
1  agenda[0]  $\leftarrow D$            holds instances of  
2   $i \leftarrow 0$                    rules that can be  
3  chart  $\leftarrow \emptyset$           used in derivation  
4  parse  $\leftarrow \emptyset$           trees  
5  while agenda[ $i$ ]  $\neq \emptyset$     facts that immediately follow  
6    do                           from one fact in agenda[ $i$ ] plus  
7      chart  $\leftarrow \text{chart} \cup \text{agenda}[i]$  some facts in chart  
8      agenda[ $i + 1$ ]  $\leftarrow \text{Conseq}(P, \text{agenda}[i], \text{chart}) - \text{chart}$   
9      parse  $\leftarrow \text{parse} \cup \text{Inst}(P, \text{agenda}[i], \text{chart})$   
10      $i \leftarrow i + 1$   
11  return parse
```

shared parse forest

instances of rules with
right-hand side consisting of
one fact in *agenda*[i] and
some facts in *chart*

Parsing algorithms are often not explicitly stated in textbooks.
Lines 8 and 9 should be performed simultaneously.

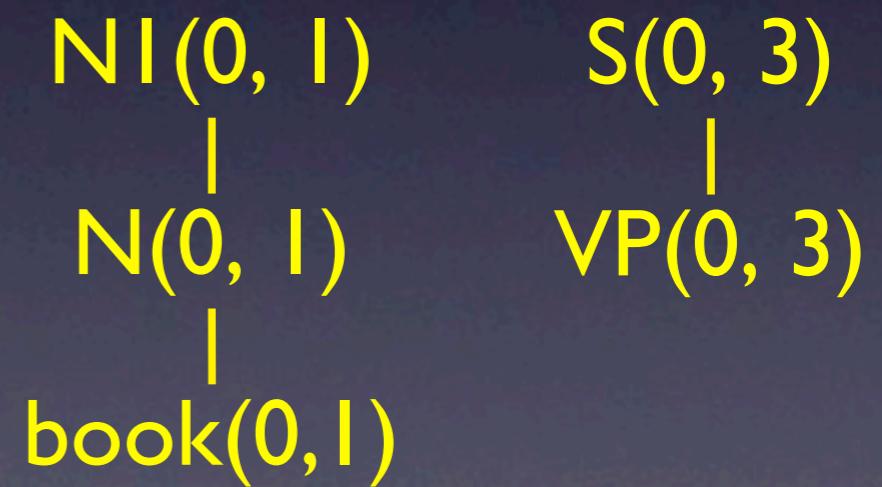
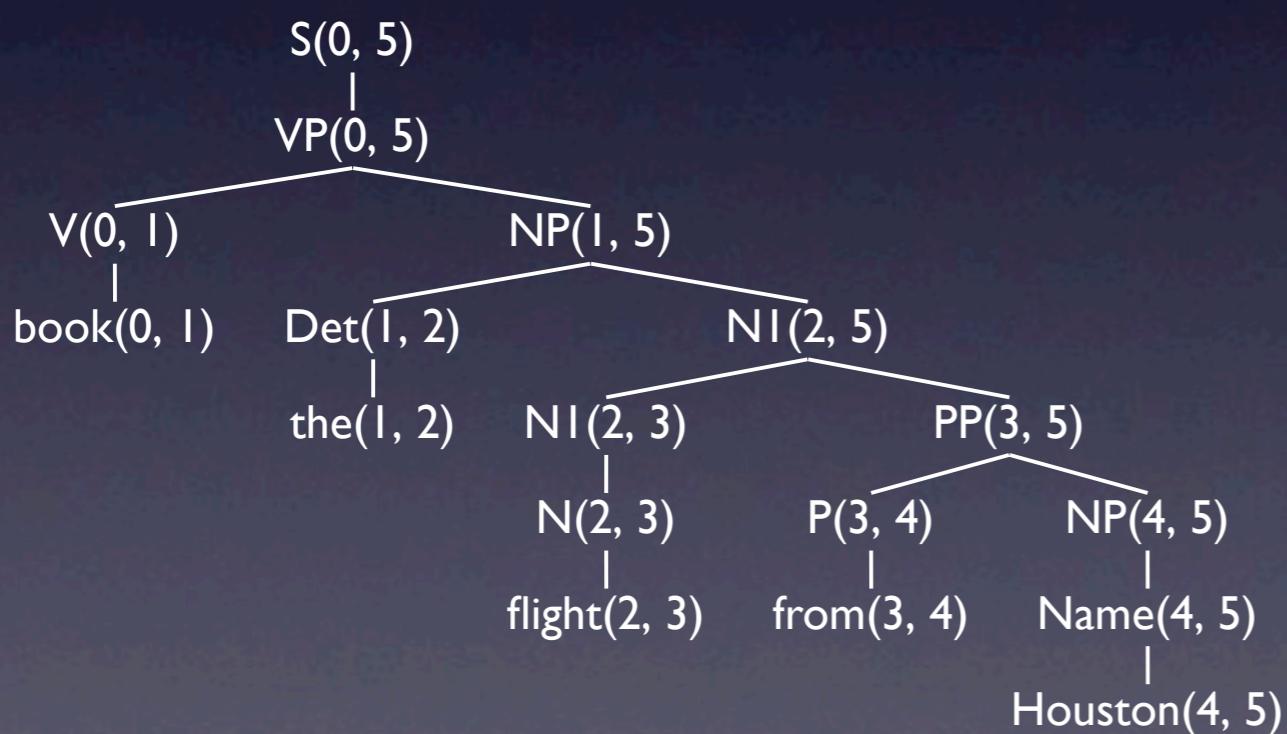
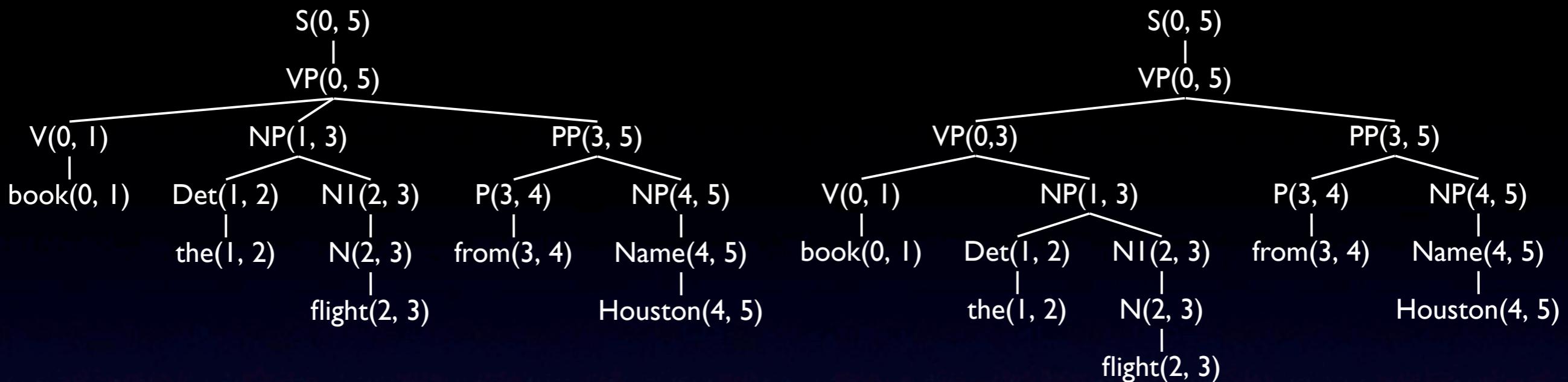
```

parse = { book(0, I). the(I, 2). flight(2, 3). from(3, 4). Houston(4, 5).
          N(0, I) :- book(0, I).
          V(0, I) :- book(0, I).
          Det(I, 2) :- the(I, 2).
          N(2, 3) :- flight(2, 3).
          P(3, 4) :- from(3, 4).
          Name(4, 5) :- Houston(4, 5).
          NI(0, I) :- N(0, I).
          VP(0, I) :- V(0, I).
          NI(2, 3) :- N(2, 3).
          NP(4, 5) :- Name(4, 5).
          NP(I, 3) :- Det(I, 2), NI(2, 3).
          PP(3, 5) :- P(3, 4), NP(4, 5).
          VP(0, 3) :- V(0, I), NP(I, 3).
          VP(0, 5) :- V(0, I), NP(I, 3), PP(3, 5).
          NI(2, 5) :- NI(2, 3), PP(3, 5).
          S(0, 3) :- VP(0, 3).
          S(0, 5) :- VP(0, 5).
          VP(0, 5) :- VP(0, 3), PP(3, 5).
          NP(I, 5) :- Det(I, 2), NI(2, 5).
          VP(0, 5) :- V(0, I), NP(I, 5).
        }

```

Output parse forest.

Can be thought of as a CFG for a singleton language.



Some rule instances found by the algorithm cannot be used in any derivation trees.

```

parse = { book(0, I). the(I, 2). flight(2, 3). from(3, 4). Houston(4, 5).
          N(0, I) :- book(0, I).
          V(0, I) :- book(0, I).                                NI(0, I)
          Det(I, 2) :- the(I, 2).
          N(2, 3) :- flight(2, 3).                            |
          P(3, 4) :- from(3, 4).                            N(0, I)
          Name(4, 5) :- Houston(4, 5).                      |
          NI(0, I) :- N(0, I).                            book(0,I)
          VP(0, I) :- V(0, I).
          NI(2, 3) :- N(2, 3).
          NP(4, 5) :- Name(4, 5).
          NP(I, 3) :- Det(I, 2), NI(2, 3).
          PP(3, 5) :- P(3, 4), NP(4, 5).
          VP(0, 3) :- V(0, I), NP(I, 3).
          VP(0, 5) :- V(0, I), NP(I, 3), PP(3, 5).
          NI(2, 5) :- NI(2, 3), PP(3, 5).                  S(0, 3)
          S(0, 3) :- VP(0, 3).
          S(0, 5) :- VP(0, 5).                            |
          VP(0, 5) :- VP(0, 3), PP(3, 5).                VP(0, 3)
          NP(I, 5) :- Det(I, 2), NI(2, 5).
          VP(0, 5) :- V(0, I), NP(I, 5).
        }

```

This parse forest is not “reduced”.

$D = \{ \text{book}(0, 1), \text{the}(1, 2), \text{flight}(2, 3), \text{from}(3, 4), \text{Houston}(4, 5) \}.$

$agenda[0] = \{ \text{book}(0, 1), \text{the}(1, 2), \text{flight}(2, 3), \text{from}(3, 4), \text{Houston}(4, 5) \}.$

$agenda[1] = \{ N(0, 1), V(0, 1), \text{Det}(1, 2), N(2, 3), P(3, 4), \text{Name}(4, 5) \}$

$agenda[2] = \{ NI(0, 1), VP(0, 1), NI(2, 3), NP(4, 5) \}$

$agenda[3] = \{ NP(1, 3), PP(3, 5) \}$

$agenda[4] = \{ VP(0, 3), VP(0, 5), NI(2, 5) \}$

$agenda[5] = \{ S(0, 3), S(0, 5), VP(0, 5), NP(1, 5) \}$

$agenda[6] = \emptyset$

$chart = agenda[0] \cup \dots \cup agenda[6]$

“Weak parse” as opposed to “parse” (Ruzzo).

Earley parsing

- $O(n^3)$ running time for arbitrary CFG
- Partial reduction of parse forest by top-down filtering
- Correct prefix property

The first two goals independent of each other in principle.
Top-down filtering plus appropriate control achieves partial reduction.

Correct prefix property

- Process input from left to right
- Reject the input as soon as the portion of the input that has been processed so far is not a **correct prefix** (i.e., cannot be extended to an element of the language)

Term comes from stack-based parsing.

Earley parsing

- Binarization of rules
- Additional subgoals

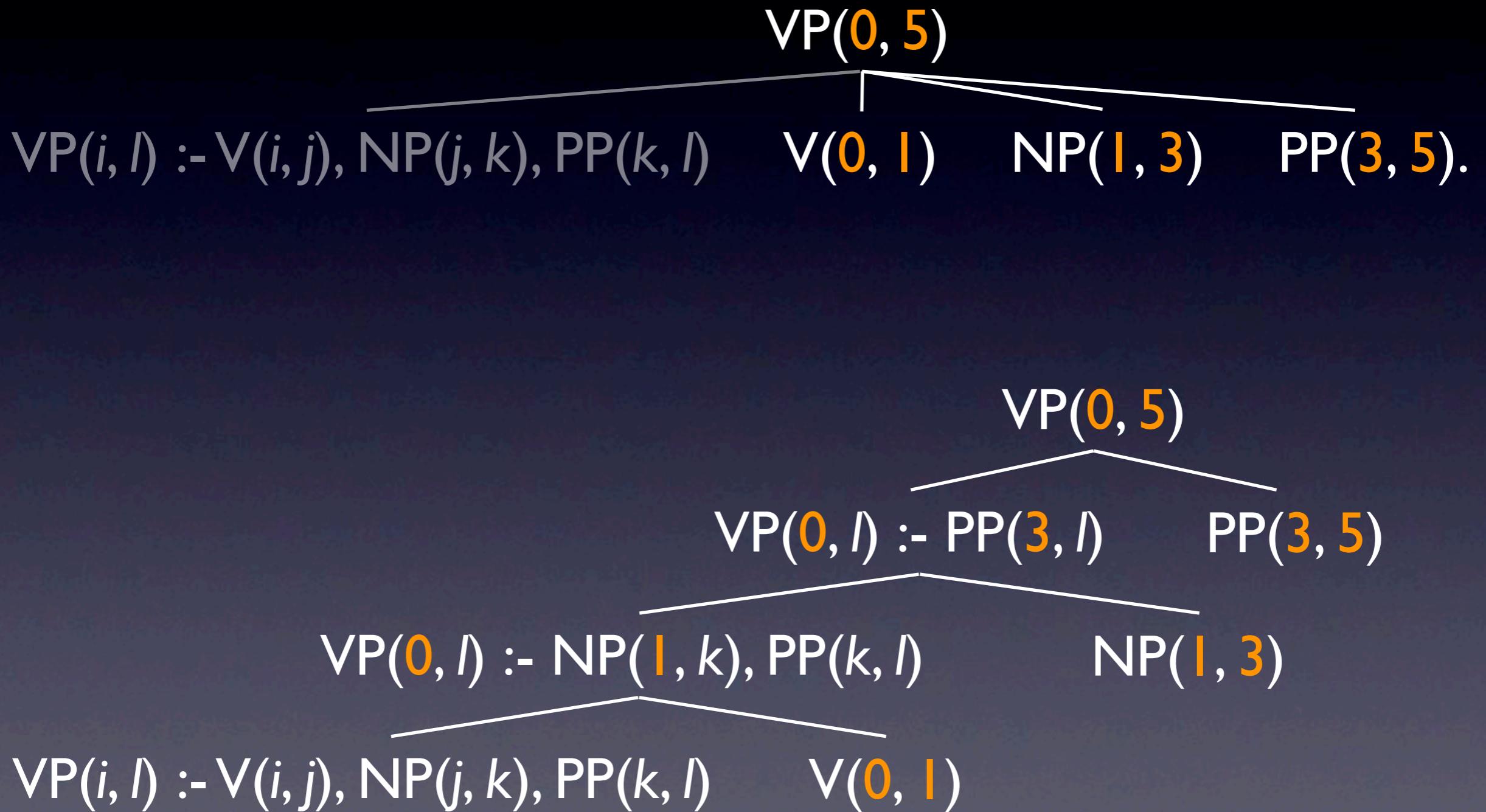
Binarization improves running time.

Additional subgoals achieve top-down filtering.

Binarization helps for top-down filtering as well.

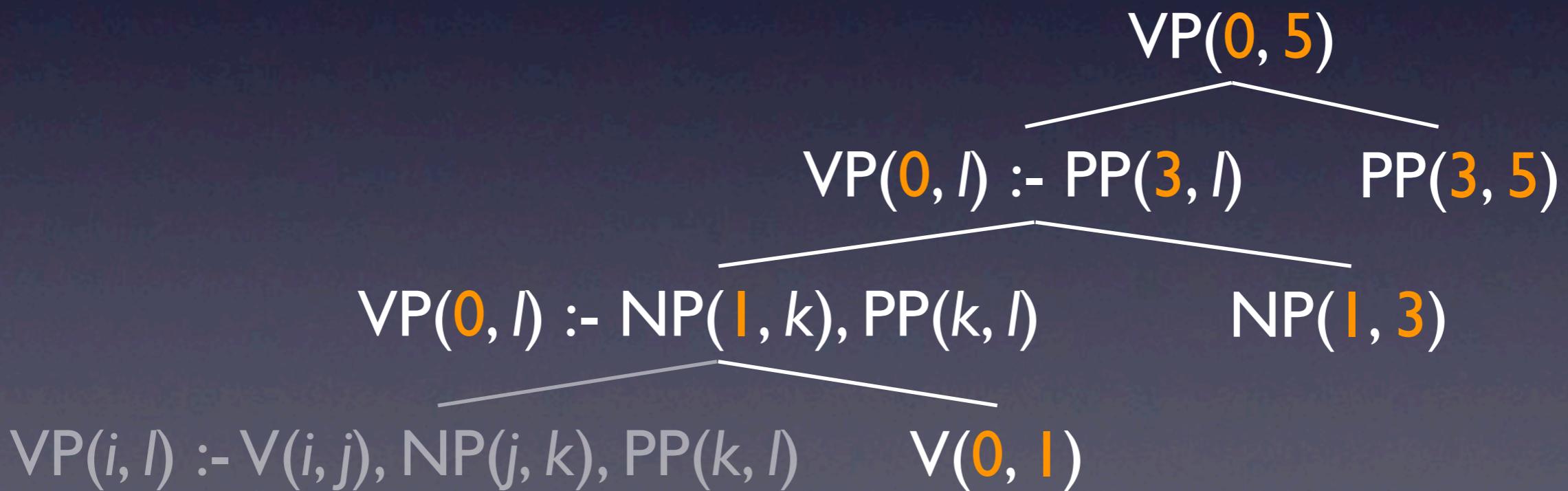
The way the top-down filtering is done achieves prefix-correctness.

Binarization of rules



Derived facts and derived clauses.

Dotted rules

$$VP \rightarrow V \ NP \ PP$$
$$VP \rightarrow V \bullet NP \ PP \quad VP \rightarrow V \ NP \bullet PP$$


Derived clauses \rightarrow derived facts

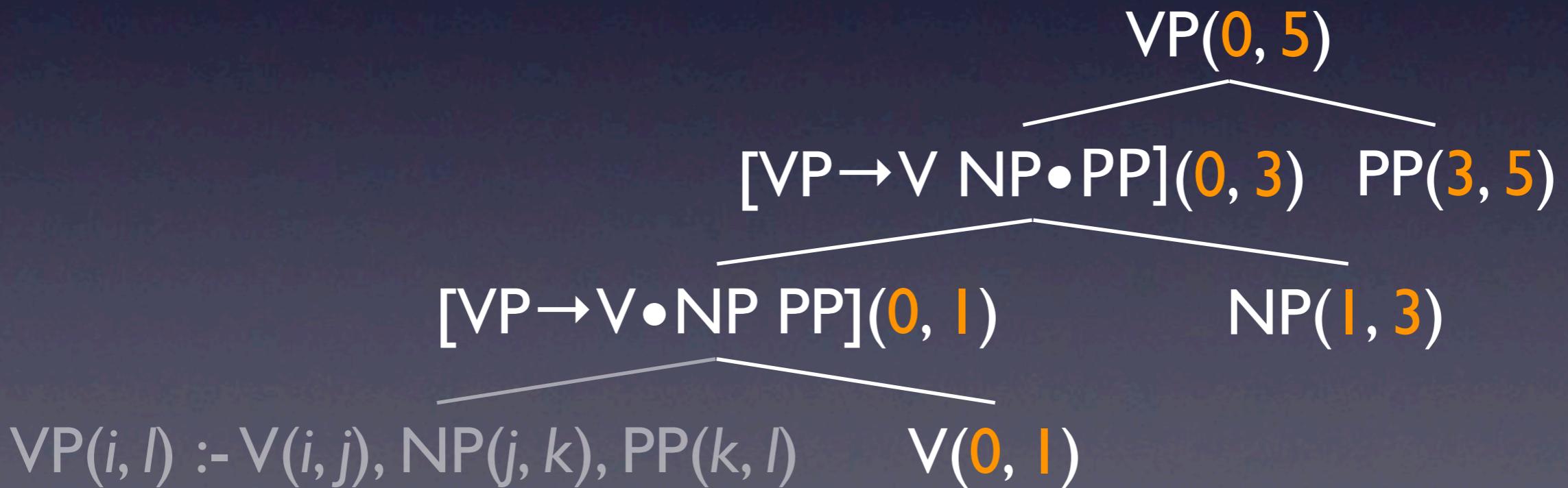
Dotted rules as new predicate names.

Dotted rules

$[VP \rightarrow V \bullet NP \; PP](i, j) :- V(i, j).$

$[VP \rightarrow V \; NP \bullet PP](i, k) :- [VP \rightarrow V \bullet NP \; PP](i, j), NP(j, k).$

$VP(i, l) :- [VP \rightarrow V \; NP \bullet PP](i, k), PP(k, l).$

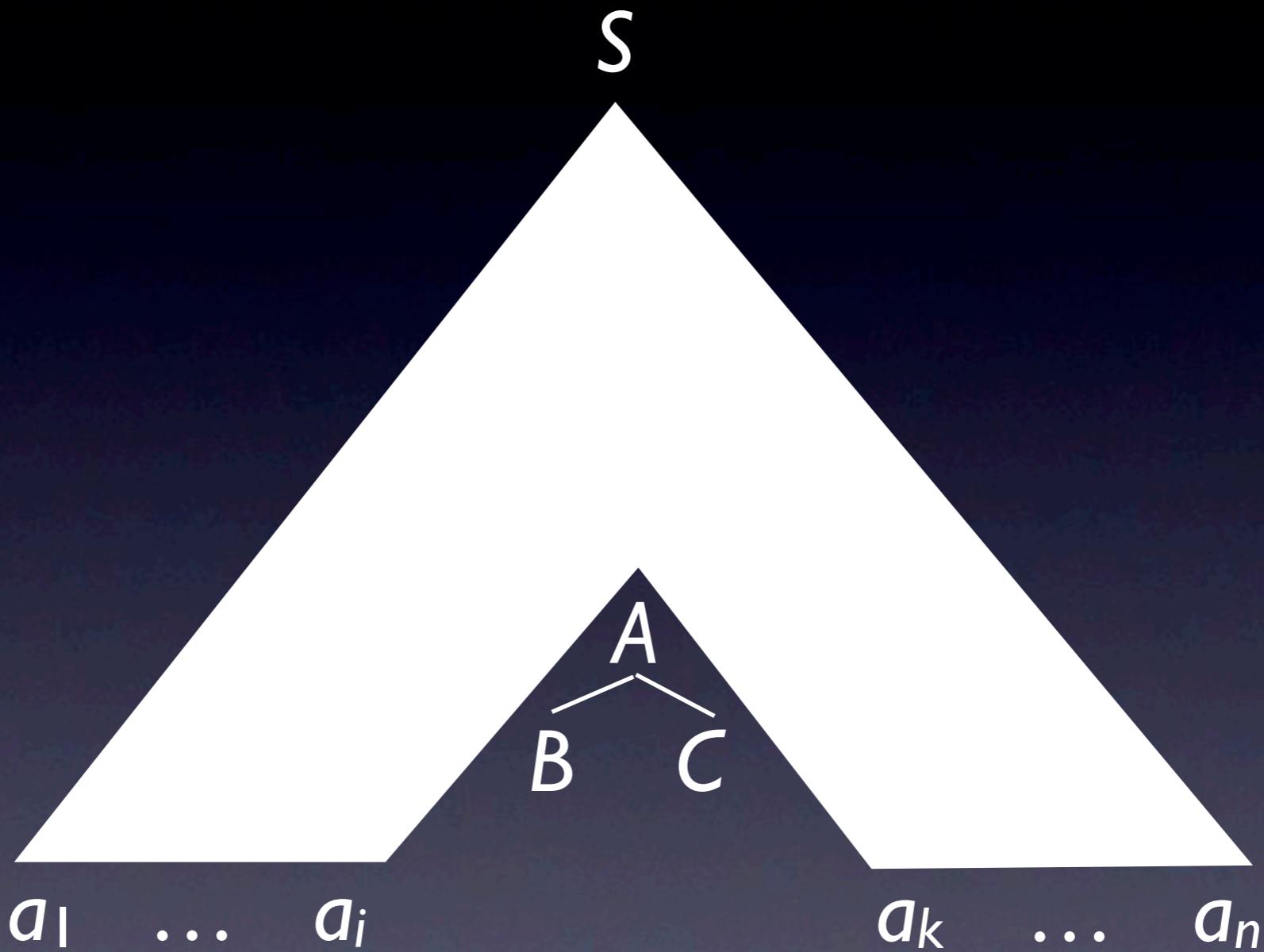


$[S \rightarrow NP \bullet VP](i, j) :- NP(i, j).$
 $S(i, k) :- [S \rightarrow NP \bullet VP](i, j), VP(j, k).$
 $[S \rightarrow Aux \bullet NP VP](i, j) :- Aux(i, j).$
 $[S \rightarrow Aux NP \bullet VP](i, k) :- [S \rightarrow Aux \bullet NP VP](i, j), Aux(j, k).$
 $S(i, l) :- [S \rightarrow Aux NP \bullet VP](i, k), VP(k, l).$
 $S(i, j) :- VP(i, j).$
 $[NP \rightarrow Det \bullet NI](i, j) :- Det(i, j).$
 $NP(i, k) :- [NP \rightarrow Det \bullet NI](i, j), NI(j, k).$
 $NP(i, j) :- Name(i, j).$
 $NP(i, j) :- Pronoun(i, j).$
 $VP(i, j) :- V(i, j).$
 $[VP \rightarrow V \bullet NP](i, j) :- V(i, j).$
 $VP(i, k) :- [VP \rightarrow V \bullet NP](i, j), NP(j, k).$
 $[VP \rightarrow V \bullet NP PP](i, j) :- V(i, j).$
 $[VP \rightarrow V NP \bullet PP](i, k) :- [VP \rightarrow V \bullet NP PP](i, j), NP(j, k).$
 $VP(i, l) :- [VP \rightarrow V NP \bullet PP](i, k), PP(k, l).$

 \vdots
 \vdots

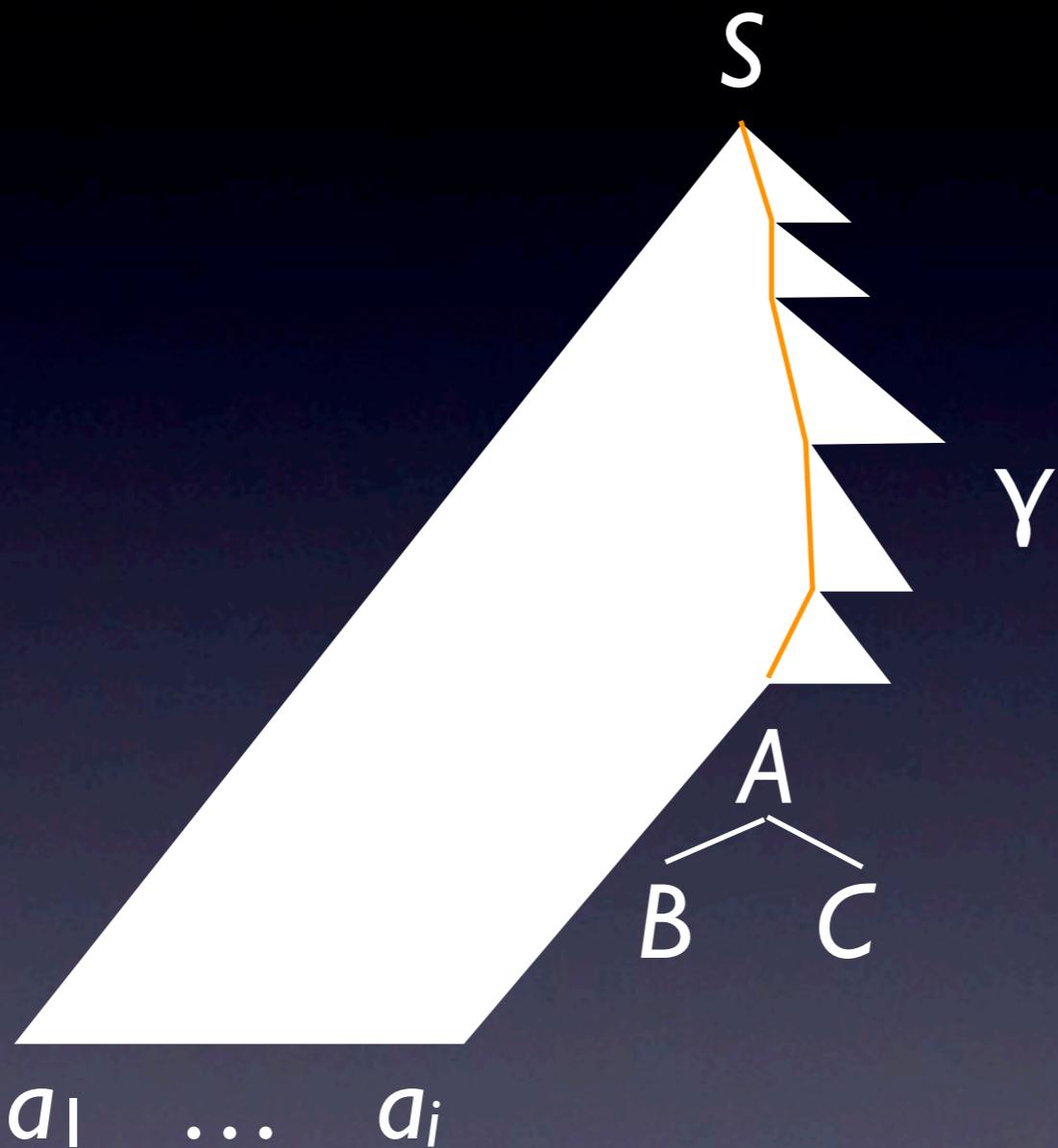
Binarized program: each rule contains at most 3 variables.
 Running time is $O(N^3)$
 “Bottom-up chart” recognizer

Top-down filtering



The existence of the white region needs to be checked in order to produce a complete parse/reduced parse forest.
But this is too much and inconsistent with correct prefix property.

Top-down filtering



$S \Rightarrow^* a_1 \dots a_i A \gamma$ leftmost derivation

Require there to be a path from the root and the part to the left of the path to be realized. This has the dual effect of partial reduction of parse forest and correct prefix property.

Top-down prediction

$m_VP(i) \Leftrightarrow S \Rightarrow^* a_1 \dots a_i \ VP \ \gamma$

“VP is predicted to start at i ”

$[VP \rightarrow V \bullet NP \ PP](i, j) :- m_VP(i), V(i, j).$

$[VP \rightarrow V \ NP \bullet PP](i, k) :- [VP \rightarrow V \bullet NP \ PP](i, j), NP(j, k).$

$VP(i, l) :- [VP \rightarrow V \ NP \bullet PP](i, k), PP(k, l).$

$m_S(0).$ seed

$S(i, j) :- VP(i, j).$ $m_VP(i) :- m_S(i).$

$m_V(i) :- m_VP(i).$

$m_NP(j) :- [VP \rightarrow V \bullet NP \ PP](i, j).$

$m_PP(k) :- [VP \rightarrow V \ NP \bullet PP](i, k).$

Earley deduction system

$m_S(0)$.

INITIALIZE

$m_B(j) :- [A \rightarrow \alpha \bullet B\beta](i, j)$.

$m_B(i) :- m_A(i). \quad (A \rightarrow B\beta)$

PREDICT

$[A \rightarrow B \bullet \beta](i, j) :- m_A(i), B(i, j)$.

$[A \rightarrow \alpha B \bullet \beta](i, k) :- [A \rightarrow \alpha \bullet B\beta](i, j), B(j, k)$.

$A(i, k) :- [A \rightarrow \alpha \bullet B](i, j), B(j, k)$.

COMPLETE

$[A \rightarrow b \bullet \beta](i, j) :- m_A(i), b(i, j)$.

$[A \rightarrow \alpha b \bullet \beta](i, k) :- [A \rightarrow \alpha \bullet b\beta](i, j), b(j, k)$.

$A(i, k) :- [A \rightarrow \alpha \bullet b](i, j), b(j, k)$.

SCAN

A minor variation, expressed in the style of Datalog.

“ $m_$ ” comes from magic-sets rewriting.

$m_NP(i) :- m_S(i).$
 $[S \rightarrow NP \bullet VP](i, j) :- m_S(i), NP(i, j).$
 $m_VP(j) :- [S \rightarrow NP \bullet VP](i, j).$
 $S(i, k) :- [S \rightarrow NP \bullet VP](i, j), VP(j, k).$
 $m_Aux(i) :- m_S(i).$
 $[S \rightarrow Aux \bullet NP \bullet VP](i, j) :- m_S(i), Aux(i, j).$
 $m_NP(j) :- [S \rightarrow Aux \bullet NP \bullet VP](i, j).$
 $[S \rightarrow Aux \bullet NP \bullet VP](i, k) :- [S \rightarrow Aux \bullet NP \bullet VP](i, j), Aux(j, k).$
 $m_VP(k) :- [S \rightarrow Aux \bullet NP \bullet VP](i, k).$
 $S(i, l) :- [S \rightarrow Aux \bullet NP \bullet VP](i, k), VP(k, l).$
 $m_VP(i) :- m_S(i).$
 $S(i, j) :- m_S(i), VP(i, j).$
 $m_Det(i) :- m_NP(i).$
 $[NP \rightarrow Det \bullet NI](i, j) :- m_NP(i), Det(i, j).$
 $m_NI(j) :- [NP \rightarrow Det \bullet NI](i, j).$
 $NP(i, k) :- [NP \rightarrow Det \bullet NI](i, j), NI(j, k).$

 \vdots
 \vdots

Earley deduction system as Datalog program.

Prefix-correct control algorithm

Chart_recognize(P, D)

agenda $\leftarrow D$

i $\leftarrow 0$

chart $\leftarrow \emptyset$

while *agenda* $\neq \emptyset$

do *trigger* $\leftarrow \text{Pop}(\text{agenda})$

chart $\leftarrow \text{chart} \cup \{\text{trigger}\}$

new $\leftarrow \text{Conseq}(P, \{\text{trigger}\}, \text{chart}) - \text{chart}$

if *trigger* $= a_j(j-1, j)$ and *new* $= \emptyset$

then reject

foreach *item* $\in \text{new}$

do *Push(agenda, item)*

return *chart*

Agenda is now first-in first-out.

Push and pop one item at a time.

Earley = magic-sets rewriting

$m_S(0)$.

$m_B(j) :- sup_{r,n}(i,j). \quad (r = A \rightarrow \alpha B \beta, n = |\alpha|)$

$m_B(i) :- m_A(i). \quad (A \rightarrow B \beta)$

$sup_{r,1}(i,j) :- m_A(i), B(i,j). \quad (r = A \rightarrow B \beta)$

$sup_{r,n+1}(i,k) :- sup_{r,n}(i,j), B(j,k). \quad (r = A \rightarrow \alpha B \beta, n = |\alpha|)$

$A(i,k) :- sup_{r,n}(i,j), B(j,k). \quad (r = A \rightarrow \alpha \bullet B, n = |\alpha|)$

$sup_{r,1}(i,j) :- m_A(i), b(i,j). \quad (r = A \rightarrow b \beta)$

$sup_{r,n+1}(i,k) :- sup_{r,n}(i,j), b(j,k). \quad (r = A \rightarrow \alpha b \beta, n = |\alpha|)$

$A(i,k) :- sup_{r,n}(i,j), b(j,k). \quad (r = A \rightarrow \alpha \bullet b, n = |\alpha|)$

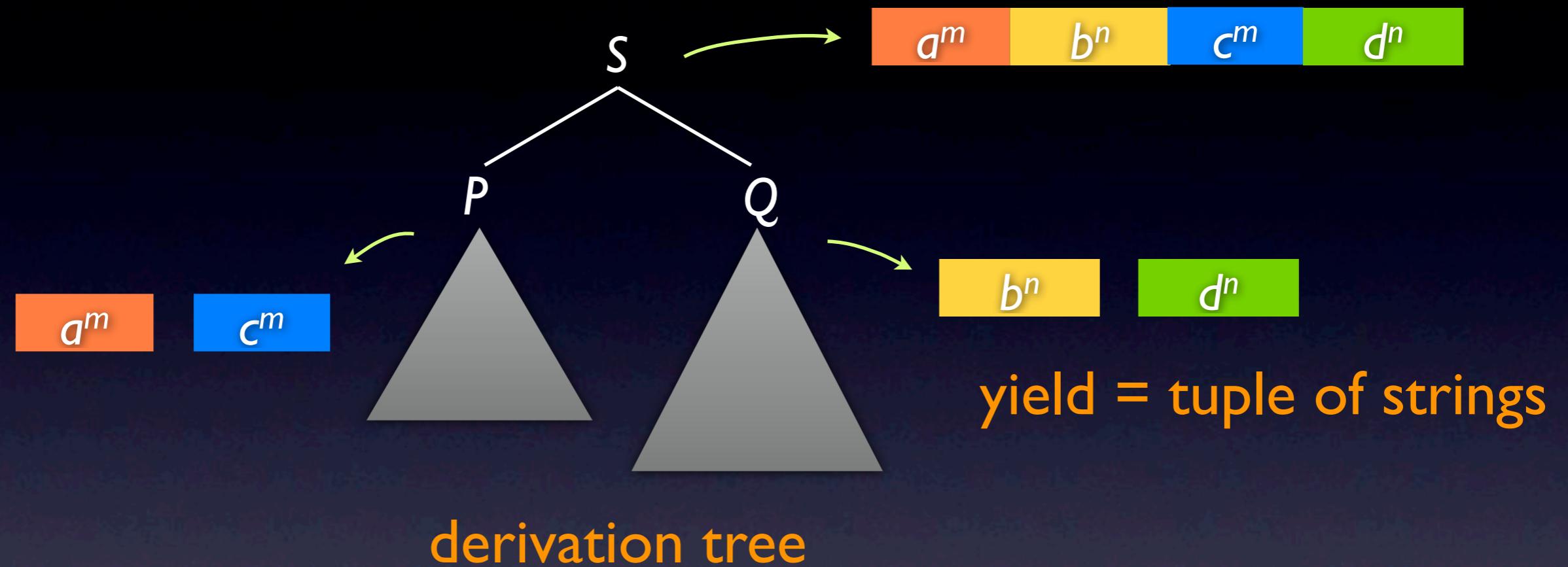
Notation for magic-sets rewriting.

Supplementary predicates correspond to dotted rules.

Magic-sets rewriting

- Logic program transformation particularly effective for Datalog
- Roughly equivalent to “memoized” top-down evaluation methods (including Earley deduction)
- Binarization of rules, top-down filtering

Multiple context-free grammars



$S(x_1y_1x_2y_2) \vdash P(x_1, x_2), Q(y_1, y_2).$

bottom-up

$P(a, c).$

$Q(b, d).$

$P(ax_1, cx_2) \vdash P(x_1, x_2).$

$Q(by_1, dy_2) \vdash Q(y_1, y_2).$

A slight variation of a by-now-familiar 2-MCFG.

Generates $\{ a^m b^n c^m d^n \mid m, n \geq 1 \}$

Directly convert to Datalog.

$S(x_1y_1x_2y_2) :- P(x_1, x_2), Q(y_1, y_2).$



$S(i, m) :- P(i, j, k, l), Q(j, k, l, m).$

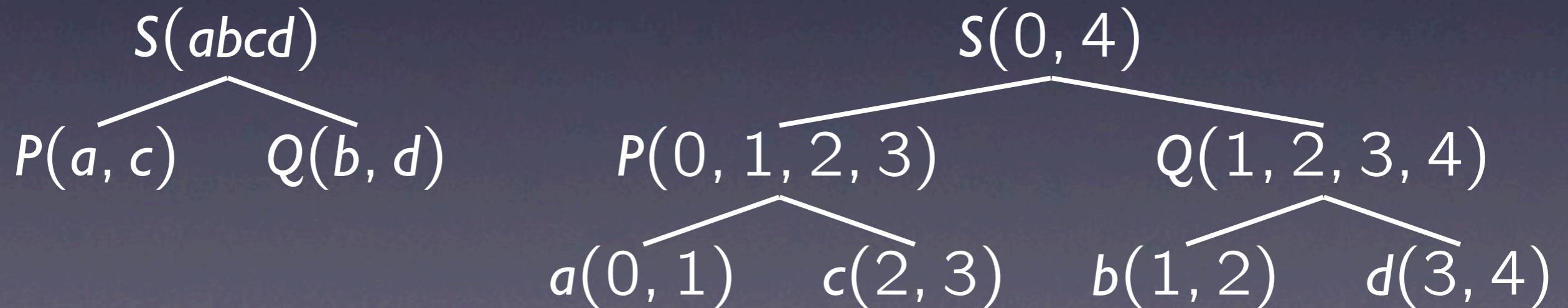
$P(ax_1, cx_2) :- P(x_1, x_2).$



$P(i, k, l, n) :- P(j, k, m, n), a(i, j), a(l, m).$

From MCFG to Datalog

$S(x_1 \textcolor{red}{y_1} x_2 y_2) :- P(x_1, x_2), Q(y_1, y_2).$	$S(i, m) :- P(i, j, k, l), Q(j, k, l, m).$
$P(a \textcolor{blue}{c}).$	$P(i, j, k, l) :- a(i, j), c(k, l).$
$P(a \textcolor{blue}{x}_1, c \textcolor{blue}{x}_2) :- P(x_1, x_2).$	$P(i, k, l, n) :- a(i, j), P(j, k, m, n), c(l, m).$
$Q(b, d).$	$Q(i, j, k, l) :- b(i, j), d(k, l).$
$Q(b y_1, d y_2) :- Q(y_1, y_2).$	$Q(i, k, l, n) :- b(i, j), d(l, m), Q(j, k, m, n).$



Adornment

- 1: $S^{bf}(i, m) :- P^{bfff}(i, j, k, l), Q^{bbbf}(j, k, l, m).$
- 2: $P^{bfff}(i, j, k, l) :- a^{bf}(i, j), c^{ff}(k, l).$
- 3: $P^{bfff}(i, k, l, n) :- a^{bf}(i, j), P^{bfff}(j, k, m, n), c^{fb}(l, m).$
- 4: $Q^{bbbf}(i, j, k, l) :- b^{bb}(i, j), d^{bf}(k, l).$
- 5: $Q^{bbbf}(i, k, l, n) :- b^{bf}(i, j), d^{bf}(l, m), Q^{bbbf}(j, k, m, n).$

?– $S(0, x).$

SLD derivation

$\xrightarrow{1} ?- P(0, j_1, k_1, l_1), Q(j_1, k_1, l_1, x).$

$\xrightarrow{2} ?- a(0, j_1), c(k_1, l_1), Q(j_1, k_1, l_1, x).$

$\xrightarrow{a(0,1)} ?- c(k_1, l_1), Q(1, k_1, l_1, x).$

$\xrightarrow{c(2,3)} ?- Q(1, 2, 3, x).$

$\xrightarrow{4} ?- b(1, 2), d(3, x).$

$\xrightarrow{b(1,2)} ?- d(3, x).$

$\xrightarrow{d(3,4)} ?- [].$

Applying magic-sets rewriting.

Each argument position is marked as “free” (uninstantiated) or “bound” (instantiated), according to the status of arguments in SLD derivations.

Initial query is assumed to be $S(0, x).$

SLD derivations generalize leftmost derivations of CFGs to logic programs.

Magic predicates

1: $S^{bf}(i, m) :- P^{bfff}(i, j, k, l), Q^{bbbf}(j, k, l, m).$

2: $P^{bfff}(i, j, k, l) :- a^{bf}(i, j), c^{ff}(k, l).$

3: $P^{bfff}(i, k, l, n) :- a^{bf}(i, j), P^{bfff}(j, k, m, n), c^{fb}(l, m).$

4: $Q^{bbbf}(i, j, k, l) :- b^{bb}(i, j), d^{bf}(k, l).$

5: $Q^{bbbf}(i, k, l, n) :- b^{bf}(i, j), d^{bf}(l, m), Q^{bbbf}(j, k, m, n).$

?– $S(0, x).$

$m_S(0)$

$\xrightarrow{1} ?- P(0, j_1, k_1, l_1), Q(j_1, k_1, l_1, x).$ $m_P(0)$

$\xrightarrow{2} ?- a(0, j_1), c(k_1, l_1), Q(j_1, k_1, l_1, x).$

$\xrightarrow{a(0,1)} ?- c(k_1, l_1), Q(1, k_1, l_1, x).$

$\xrightarrow{c(2,3)} ?- Q(1, 2, 3, x).$ $m_Q(1, 2, 3)$

$\xrightarrow{4} ?- b(1, 2), d(3, x).$

$\xrightarrow{b(1,2)} ?- d(3, x).$

$\xrightarrow{d(3,4)} ?- [].$

The fact that $S(0, x)$ is the first goal is represented by the “magic” predicate $m_S(0)$, etc.
 Magic predicates take only bound arguments.

Supplementary predicates

5: $Q^{bbbf}(i, k, l, n) :- b^{bf}(i, j), d^{bf}(l, m), Q^{bbbf}(j, k, m, n).$

?– $Q(2, 4, 6, x).$

SLD derivation

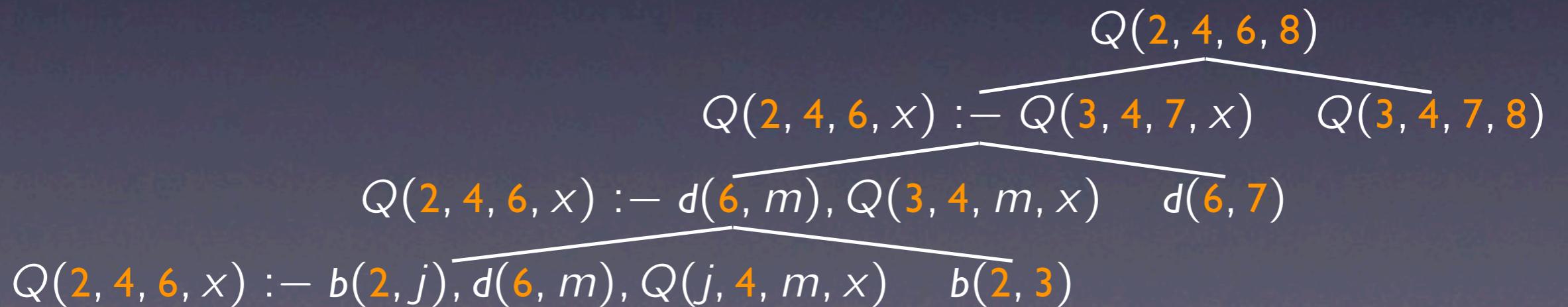
$\xrightarrow{5} ?- b(2, j), d(6, m), Q(j, 4, m, x).$

$\xrightarrow{b(2,3)} ?- d(6, m), Q(3, 4, m, x).$

$\xrightarrow{d(6,7)} ?- Q(3, 4, 7, x).$

:
:

$\Rightarrow ?- [].$



Here's how rule 4 is converted to Datalog rules.
A Datalog derivation corresponding to an SLD derivation.

Supplementary predicates

5: $Q^{bbbf}(i, k, l, n) :- b^{bf}(i, j), d^{bf}(l, m), Q^{bbbf}(j, k, m, n).$

?– $Q(2, 4, 6, x).$

SLD derivation

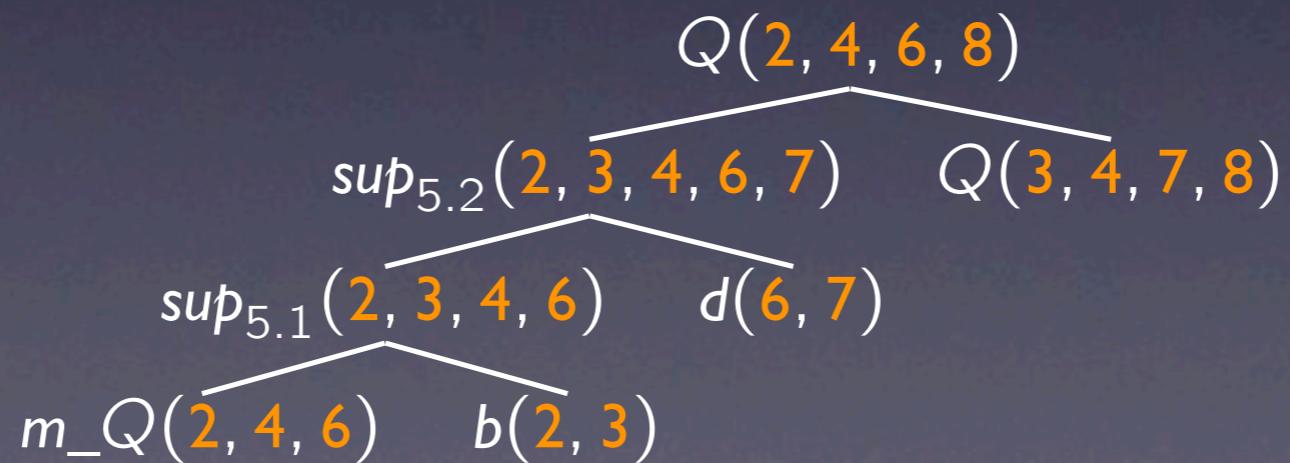
$\xrightarrow{5} ?- b(2, j), d(6, m), Q(j, 4, m, x).$

$\xrightarrow{b(2,3)} ?- d(6, m), Q(3, 4, m, x).$

$\xrightarrow{d(6,7)} ?- Q(3, 4, 7, x).$

:
:

$\Rightarrow ?- [].$



A partially instantiated clause is represented with a magic predicate.
A derived clause is represented with a supplementary predicate.

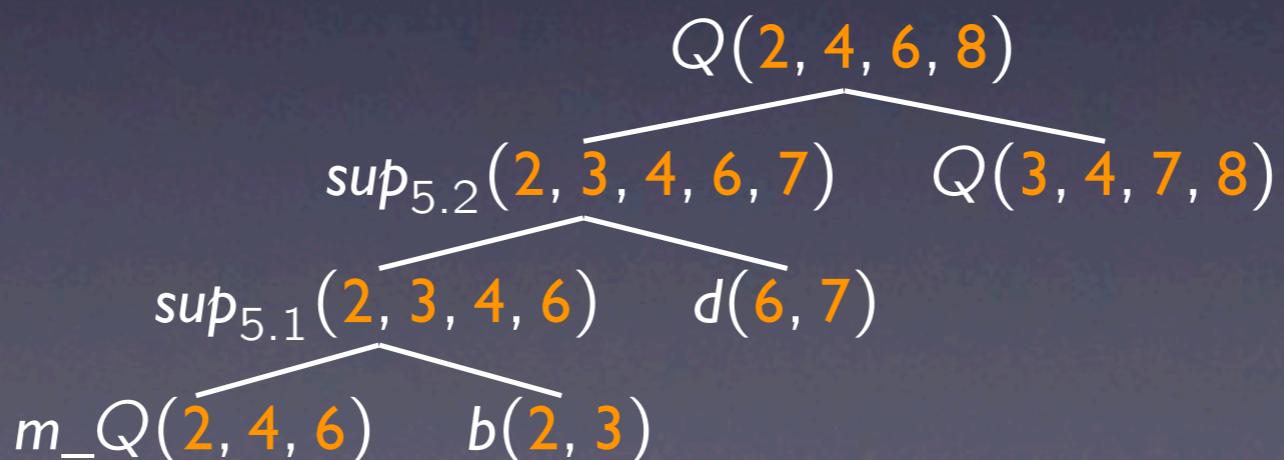
Rules for supplementary predicates

5: $Q^{bbbf}(i, k, l, n) :- b^{bf}(i, j), d^{bf}(l, m), Q^{bbbf}(j, k, m, n).$

$sup_{5.1}(i, j, k, l) :- m_Q(i, k, l), b(i, j).$

$sup_{5.2}(i, j, k, l, m) :- sup_{5.1}(i, j, k, l), d(l, m).$

$Q(i, k, l, n) :- sup_{5.2}(i, j, k, l, m), Q(j, k, m, n).$



Rules for magic predicates

- 1: $S^{bf}(i, m) :- P^{bfff}(i, j, k, l), Q^{bbbf}(j, k, l, m).$
- 2: $P^{bfff}(i, j, k, l) :- a^{bf}(i, j), c^{ff}(k, l).$
- 3: $P^{bfff}(i, k, l, n) :- a^{bf}(i, j), P^{bfff}(j, k, m, n), c^{fb}(l, m).$
- 4: $Q^{bbbf}(i, j, k, l) :- b^{bb}(i, j), d^{bf}(k, l).$
- 5: $Q^{bbbf}(i, k, l, n) :- b^{bf}(i, j), d^{bf}(l, m), Q^{bbbf}(j, k, m, n).$

$m_P(i) :- m_S(i).$

$m_Q(j, k, l) :- sup_{1.1}(i, j, k, l).$

$m_P(j) :- sup_{3.1}(i, j).$

$m_Q(j, k, m) :- sup_{5.2}(i, j, k, l, m).$

	a	c	
$r_1: m_P(i) :- m_S(i).$	0		2
$r_2: m_Q(j, k, l) :- sup_{1.1}(i, j, k, l).$			
$r_3: m_P(j) :- sup_{3.1}(i, j).$	$m_S(0)$	$sup_{2.1}(0, 1)$	$P(0, 1, 1, 2)$
$r_4: m_Q(j, k, m) :- sup_{5.2}(i, j, k, l, m).$	$m_P(0)$	$sup_{3.1}(0, 1)$	$sup_{1.1}(0, 1, 1, 2)$
$r_5: sup_{1.1}(i, j, k, l) :- m_S(i), P(i, j, k, l).$		$m_P(1)$	$m_Q(1, 1, 2)$
$r_6: S(i, m) :- sup_{1.1}(i, j, k, l), Q(j, k, l, m).$			
$r_7: sup_{2.1}(i, j) :- m_P(i), a(i, j).$			
$r_8: P(i, j, k, l) :- sup_{2.1}(i, j), c(k, l).$			not prefix-correct!
$r_9: sup_{3.1}(i, j) :- m_P(i), a(i, j).$			
$r_{10}: sup_{3.2}(i, k, m, n) :- sup_{3.1}(i, j), P(j, k, m, n).$			
$r_{11}: P(i, k, l, n) :- sup_{3.2}(i, k, m, n), c(l, m).$			
$r_{12}: sup_{4.1}(i, j, k) :- m_Q(i, j, k), b(i, j).$			
$r_{13}: Q(i, j, k, l) :- sup_{4.1}(i, j, k), d(k, l).$			
$r_{14}: sup_{5.1}(i, j, k, l) :- m_Q(i, k, l), b(i, j).$			
$r_{15}: sup_{5.2}(i, j, k, l, m) :- sup_{5.1}(i, j, k, l), d(l, m).$			
$r_{16}: Q(i, k, l, n) :- sup_{5.2}(i, j, k, l, m), Q(j, k, m, n).$			

Magic-sets corresponds to Earley for CFGs, but in general, does not imply prefix-correctness.

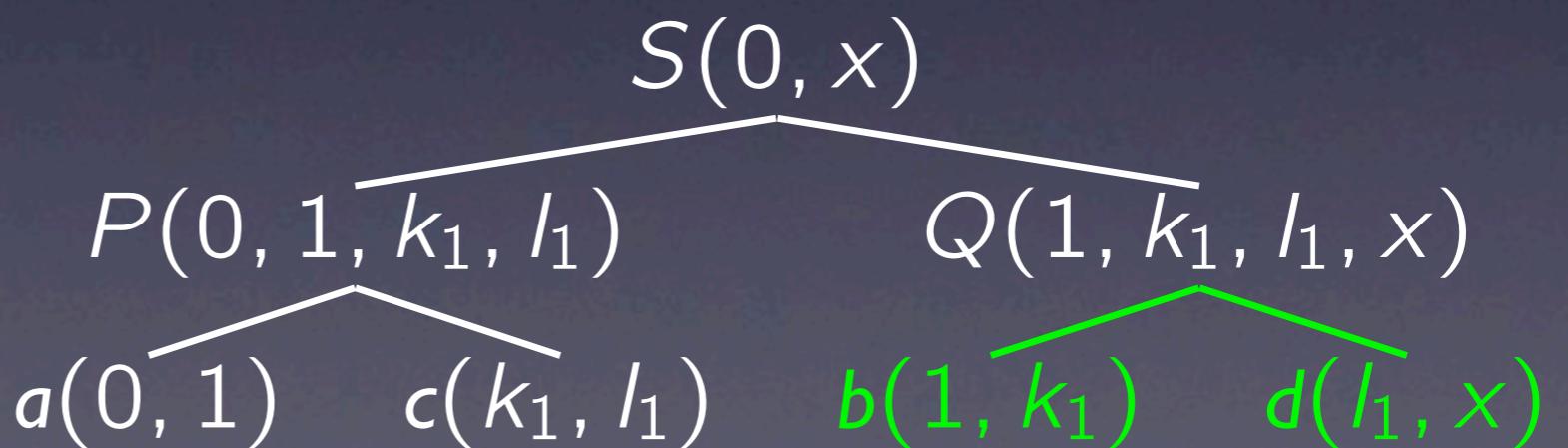
Cause of non-prefix-correctness

SLD derivation

$$\begin{aligned} & ?- S(0, x). \\ \xrightarrow{1} & ?- P(0, j_1, k_1, l_1), Q(j_1, k_1, l_1, x). \\ \xrightarrow{2} & ?- a(0, j_1), c(k_1, l_1), Q(j_1, k_1, l_1, x). \\ \xrightarrow{a(0,1)} & ?- c(k_1, l_1), Q(1, k_1, l_1, x). \\ \xrightarrow{c(1,2)} & ?- Q(1, 1, 2, x). \\ \xrightarrow{4} & ?- b(1, 1), d(2, x). \end{aligned}$$

adornment

$$2: P^{bfff}(i, j, k, l) :- a^{bf}(i, j), c^{ff}(k, l).$$



incomplete derivation tree

SLD derivation with the original program is not prefix-correct.

The order of extensional facts at the leaves of the derivation tree does not match the order of terminals in the input.

Securing prefix-correctness

$S(x_1 \ y_1 \ x_2 \ y_2) :- P(x_1, x_2), Q(y_1, y_2).$
 $i \ j \ k \ l \ m$

~~$S(i, m) :- P(i, j, k, l), Q(j, k, l, m).$~~

$S(i, m) :- P_1(i, j), Q_1(j, k), P(i, j, k, l), Q(j, k, l, m).$

~~$P(i, j, k, l) :- a(i, j), c(k, l).$~~ $P(i, j, k, l) :- aux(i, j), c(k, l).$

~~$P_1(i, j) :- a(i, j).$~~ $P_1(i, j) :- aux(i, j).$

$aux(i, j) :- a(i, j).$

$S^{bf}(i, m) :- P_1^{bf}(i, j), Q_1^{bf}(j, k), P^{bbbbf}(i, j, k, l), Q^{bbbbf}(j, k, l, m).$

$P^{bbbbf}(i, j, k, l) :- aux^{bb}(i, j), c^{bf}(k, l).$

$P_1^{bf}(i, j) :- aux^{bf}(i, j).$

$aux^{bf}(i, j) :- a^{bf}(i, j).$

Add “redundant” subgoals.

The aux predicate “folds” common part of two rules.

New starting point

- 1: $S^{bf}(i, m) :- P_1^{bf}(i, j), Q_1^{bf}(j, k), P^{bbbbf}(i, j, k, l), Q^{bbbbf}(j, k, l, m).$
- 2: $P_1^{bf}(i, j) :- aux_2^{bf}(i, j).$
- 3: $P^{bbbbf}(i, j, k, l) :- aux_2^{bf}(i, j), c^{bf}(k, l).$
- 4: $aux_2^{bf}(i, j) :- a^{bf}(i, j).$
- 5: $P_1^{bf}(i, k) :- aux_3^{bff}(i, j, k).$
- 6: $P^{bbbbf}(i, k, l, n) :- aux_3^{bff}(i, j, k), c^{bf}(l, m), P^{bbbbf}(j, k, m, n).$
- 7: $aux_3^{bff}(i, j, k) :- a^{bf}(i, j), P_1^{bf}(j, k).$
- 8: $Q_1^{bf}(i, j) :- aux_4^{bf}(i, j).$
- 9: $Q^{bbbbf}(i, j, k, l) :- aux_4^{bf}(i, j), d^{bf}(k, l).$
- 10: $aux_4^{bf}(i, j) :- b^{bf}(i, j).$
- 11: $Q_1^{bf}(i, k) :- aux_5^{bff}(i, j, k).$
- 12: $Q^{bbbbf}(i, k, l, n) :- aux_5^{bff}(i, j, k), d^{bf}(l, m), Q^{bbbbf}(j, k, m, n).$
- 13: $aux_5^{bff}(i, j, k) :- b^{bf}(i, j), Q_1^{bf}(j, k).$

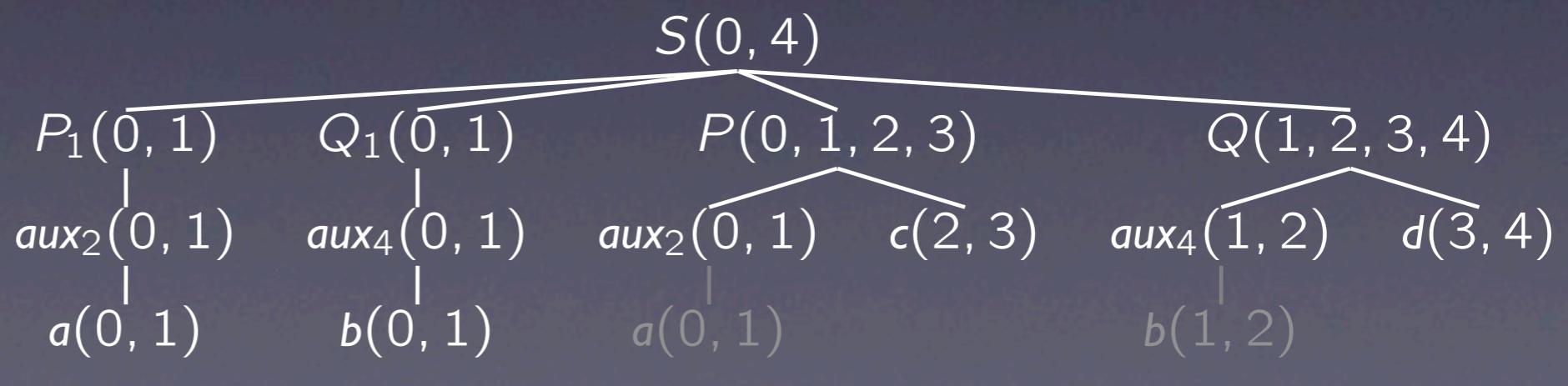
The result of introducing redundant subgoals to the original program.

Prefix-correctness

SLD derivation

$? - S(0, x).$
 $\xrightarrow{1} ? - P_1(0, j_1), Q_1(j_1, k_1), P(0, j_1, k_1, l_1), Q(j_1, k_1, l_1, x).$
 $\xrightarrow{2} ? - aux_2(0, j_i), Q_1(j_1, k_1), P(0, j_1, k_1, l_1), Q(j_1, k_1, l_1, x).$
 $\xrightarrow{4} ? - a(0, j_1), Q_1(j_1, k_1), P(0, j_1, k_1, l_1), Q(j_1, k_1, l_1, x).$
 $\xrightarrow{a(0,1)} ? - Q_1(1, k_1), P(0, 1, k_1, l_1), Q(1, k_1, l_1, x).$
 $\xrightarrow{8} ? - aux_4(1, k_1), P(0, 1, k_1, l_1), Q(1, k_1, l_1, x).$
 $\xrightarrow{10} ? - b(1, k_1), P(0, j_1, k_1, l_1), Q(1, k_1, l_1, x).$
 $\xrightarrow{b(1,2)} ? - P(0, 1, 2, l_1), Q(1, 2, l_1, x).$
 $\xrightarrow{3} ? - aux_3(0, 1), c(2, l_1), Q(1, 2, l_1, x).$
 \vdots
 $\xrightarrow{} ? - c(2, l_1), Q(1, 2, l_1, x).$
 $\xrightarrow{c(2,3)} ? - Q(1, 2, 3, x).$
 $\xrightarrow{9} ? - aux_4(1, 2), d(3, x).$
 \vdots
 $\xrightarrow{} ? - d(3, x).$
 $\xrightarrow{d(3,4)} ? - []$

derivation tree

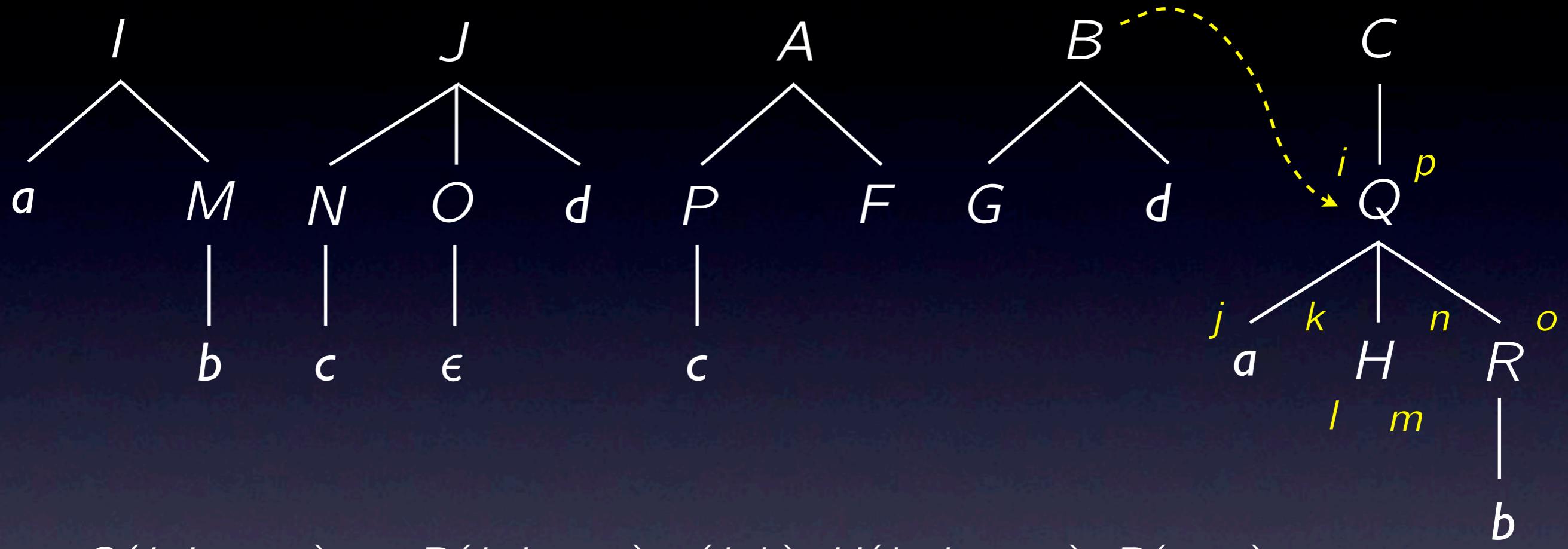


SLD derivation is now prefix-correct.

Magic-sets rewriting

$r_1 : m_P_1(i) :- m_S(i).$	$r_{21} : P_1(i, j) :- m_P_1(i), aux_2(i, j).$
$r_2 : m_Q_1(j) :- sup_{1.1}(i, j).$	$r_{22} : sup_{3.1}(i, j, k) :- m_P(i, j, k), aux_2(i, j).$
$r_3 : m_P(i, j, k) :- sup_{1.2}(i, j, k).$	$r_{23} : P(i, j, k, l) :- sup_{3.1}(i, j, k), c(k, l).$
$r_4 : m_Q(j, k, l) :- sup_{1.3}(i, j, k, l).$	$r_{24} : aux_2(i, j) :- m_aux_2(i), a(i, j).$
$r_5 : m_aux_2(i) :- m_P_1(i).$	$r_{25} : P_1(i, k) :- m_P_1(i), aux_3(i, j, k).$
$r_6 : m_aux_2(i) :- m_P(i, j, k).$	$r_{26} : sup_{6.1}(i, j, k, l) :- m_P(i, k, l), aux_3(i, j, k).$
$r_7 : m_aux_3(j) :- m_P_1(i).$	$r_{27} : sup_{6.2}(i, j, k, l, m) :- sup_{6.1}(i, j, k, l), c(l, m).$
$r_8 : m_aux_3(i) :- m_P(i, k, l).$	$r_{28} : P(i, k, l, n) :- sup_{6.2}(i, j, k, l, m), P(j, k, m, n).$
$r_9 : m_P(j, k, m) :- sup_{6.2}(i, j, k, l, m).$	$r_{29} : sup_{7.1}(i, j) :- m_aux_3(i), a(i, j).$
$r_{10} : m_P_1(j) :- sup_{7.1}(i, j).$	$r_{30} : aux_3(i, j, k) :- sup_{7.1}(i, j), P_1(j, k).$
$r_{11} : m_aux_4(i) :- m_Q_1(i).$	$r_{31} : Q_1(i, j) :- m_Q_1(i), aux_4(i, j).$
$r_{12} : m_aux_4(i) :- m_Q(i, j, k).$	$r_{32} : sup_{9.1}(i, j, k) :- m_Q(i, j, k), aux_4(i, j).$
$r_{13} : m_aux_5(i) :- m_Q_1(i).$	$r_{33} : Q(i, j, k, l) :- sup_{9.1}(i, j, k), d(k, l).$
$r_{14} : m_aux_5(i) :- m_Q(i, k, l).$	$r_{34} : aux_4(i, j) :- m_aux_4(i), b(i, j).$
$r_{15} : m_Q(j, k, m) :- sup_{12.2}(i, j, k, l, m).$	$r_{35} : Q_1(i, k) :- m_Q_1(i), aux_5(i, j, k).$
$r_{16} : m_Q_1(j) :- sup_{13.1}(i, j).$	$r_{36} : sup_{12.1}(i, j, k, l) :- m_Q(i, k, l), aux_5(i, j, k).$
$r_{17} : sup_{1.1}(i, j) :- m_S(i), P_1(i, j).$	$r_{37} : sup_{12.2}(i, j, k, l, m) :- sup_{12.1}(i, j, k, l), d(l, m).$
$r_{18} : sup_{1.2}(i, j, k) :- sup_{1.1}(i, j), Q_1(j, k).$	$r_{38} : Q(i, k, l, n) :- sup_{12.2}(i, j, k, l, m), Q(j, k, m, n).$
$r_{19} : sup_{1.3}(i, j, k, l) :- sup_{1.2}(i, j, k), P(i, j, k, l).$	$r_{39} : sup_{13.1}(i, j) :- m_aux_5(i), b(i, j).$
$r_{20} : S(i, m) :- sup_{1.3}(i, j, k, l), Q(j, k, l, m).$	$r_{40} : aux_5(i, j, k) :- sup_{13.1}(i, j), Q_1(j, k).$

Application to TAG



$Q(i, l, m, p) :- B(i, j, o, p), a(j, k), H(k, l, m, n), R(n, o).$



$\text{aux}(i, j, k, l) :- B_1(i, j), a(j, k), H_1(k, l).$

$Q_1(i, l) :- \text{aux}(i, j, k, l).$

$Q(i, l, m, p) :- \text{aux}(i, j, k, l), H(k, l, m, n), R(n, o), B(i, j, o, p).$

Direct conversion to Datalog.
Essentially the method of Lang.

Magic-sets rewriting

$m_B_1(i) :- m_aux(i).$

$m_H_1(k) :- sup_{1.2}(i, j, k).$

$m_aux(i) :- m_Q_1(i).$

$m_aux(i) :- m_Q(i, l, m).$

$m_H(k, l, m) :- sup_{3.1}(i, j, k, l, m).$

$m_R(n) :- sup_{3.2}(i, j, l, m, n).$

$m_B(i, j, o) :- sup_{3.3}(i, j, l, m, o).$

$sup_{1.1}(i, j) :- m_aux(i), B_1(i, j).$

$sup_{1.2}(i, j, k) :- sup_{1.1}(i, j), a(j, k).$

$aux(i, j, k, l) :- sup_{1.2}(i, j, k), H_1(k, l).$

$Q_1(i, l) :- m_Q_1(i, l), aux(i, j, k, l).$

$sup_{3.1}(i, j, k, l, m) :- m_Q(i, l, m), aux(i, j, k, l).$

$sup_{3.2}(i, j, l, m, n) :- sup_{3.1}(i, j, k, l, m), H(k, l, m, n).$

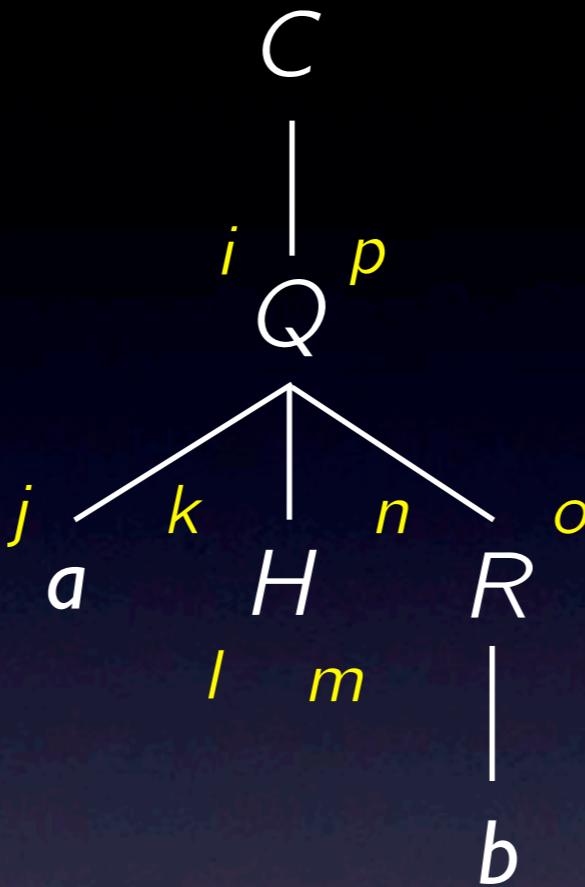
$sup_{3.3}(i, j, l, m, o) :- sup_{3.2}(i, j, l, m, n), R(n, o).$

$Q(i, l, m, p) :- sup_{3.3}(i, j, l, m, o), B(i, j, o, p).$

$O(n^6)$

Apply the technique for MCFGs.

Time and space complexity



$Q(i, l, m, p) :- \text{aux}(i, j, k, l), H(k, l, m, n), R(n, o), B(i, j, o, p).$ $Q(i, l, m, p)$



Optimal complexity bounds are obtained without any fine-tuning

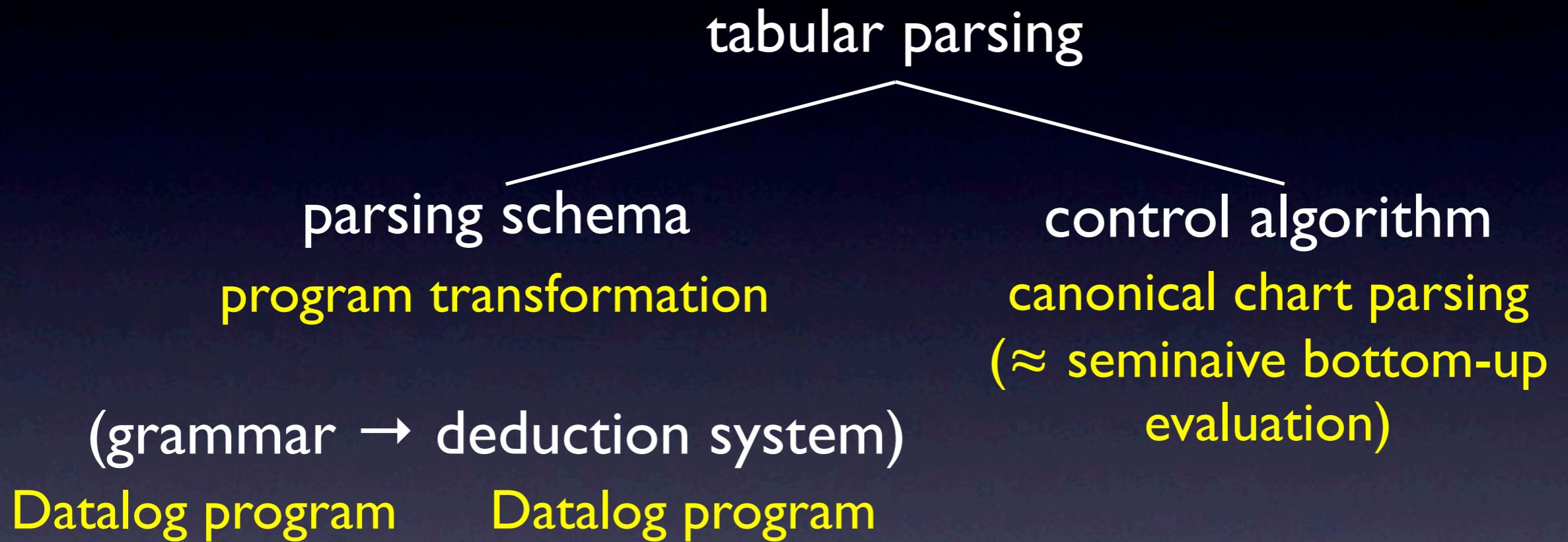
The arity of supplementary predicate indicated by the number of yellow lines crossing orange.

The same time and space complexity as Nederhof's prefix-correct Earley recognizer for TAG, previously known optimal bound.

Nothing ad hoc except the introduction of redundant subgoals.

A special case of general technique for MCFGs.

How parsing should be approached



- Use well-established, general formalism/method
- Avoids ad hoc techniques as much as possible
- Easier to understand and easier to prove correct

Some people even divide parsing schema into conversion to PDA and “PDA tabulation”. This is not necessarily a natural approach.

How parsing should be approached

Earley parsing

parsing schema

introduction of redundant predicates
+ magic-sets rewriting

(grammar → deduction system)

Datalog program

Datalog program

control algorithm

canonical chart parsing
(≈ seminaive bottom-up
evaluation)

- Use well-established, general formalism/method
- Avoids ad hoc techniques as much as possible
- Easier to understand and easier to prove correct

In the case of Earley, introduction of redundant predicates plus magic-sets rewriting.