

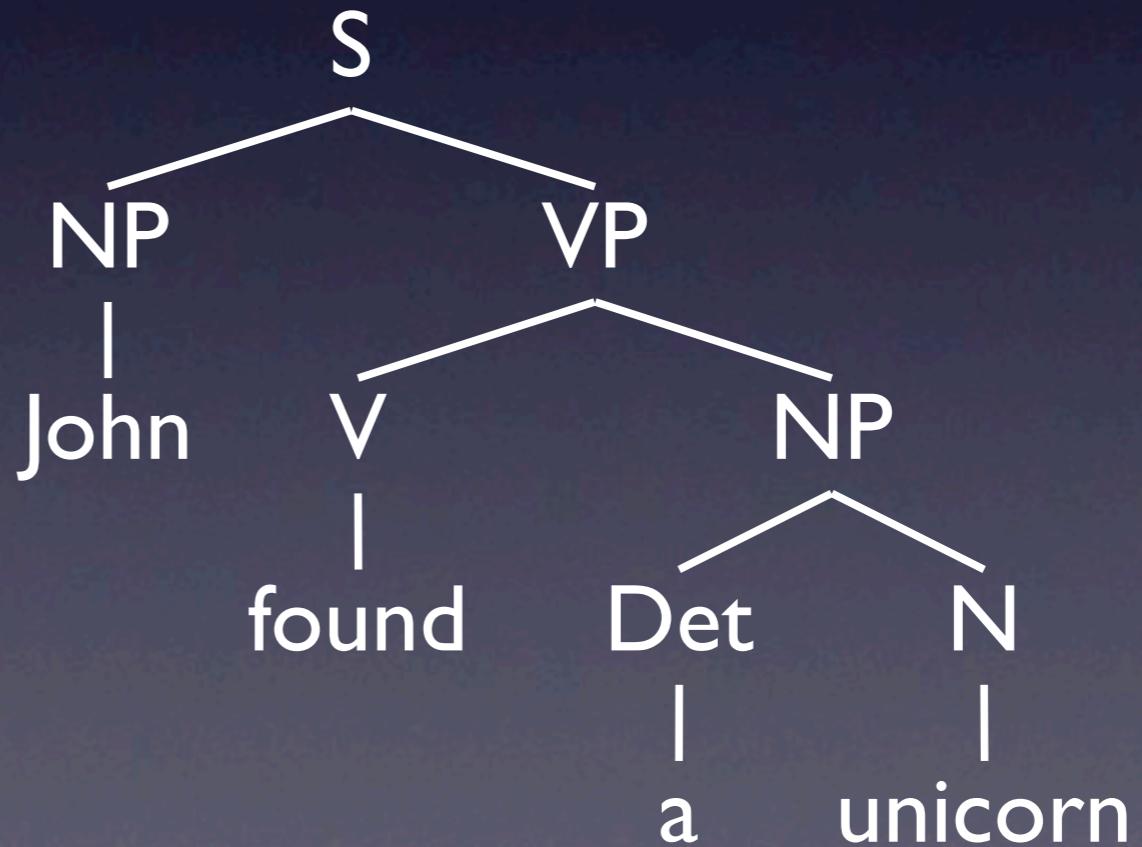
Datalog as a Uniform Framework for Parsing and Generation

Makoto Kanazawa
National Institute of Informatics
Tokyo, Japan

CFG recognition/parsing

```
S → NP VP  
VP → V NP  
NP → Det N  
NP → John  
V → found  
Det → a  
N → unicorn
```

John found a unicorn $\in L(G)$?



To give an idea of what the reduction looks like
Well-known case of CFGs

Datalog query evaluation

```
S(i, k) :- NP(i, j), VP(j, k).  
VP(i, k) :- V(i, j), NP(j, k).  
NP(i, k) :- Det(i, j), N(j, k).  
NP(i, j) :- John(i, j).  
V(i, j) :- found(i, j).  
Det(i, j) :- a(i, j).  
N(i, j) :- unicorn(i, j).
```

program

```
John(0, 1).  
found(1, 2).  
a(2, 3).  
unicorn(3, 4).
```

database

```
?- S(0, 4).
```

query



John(0, 1)

found(1, 2)

a(2, 3)

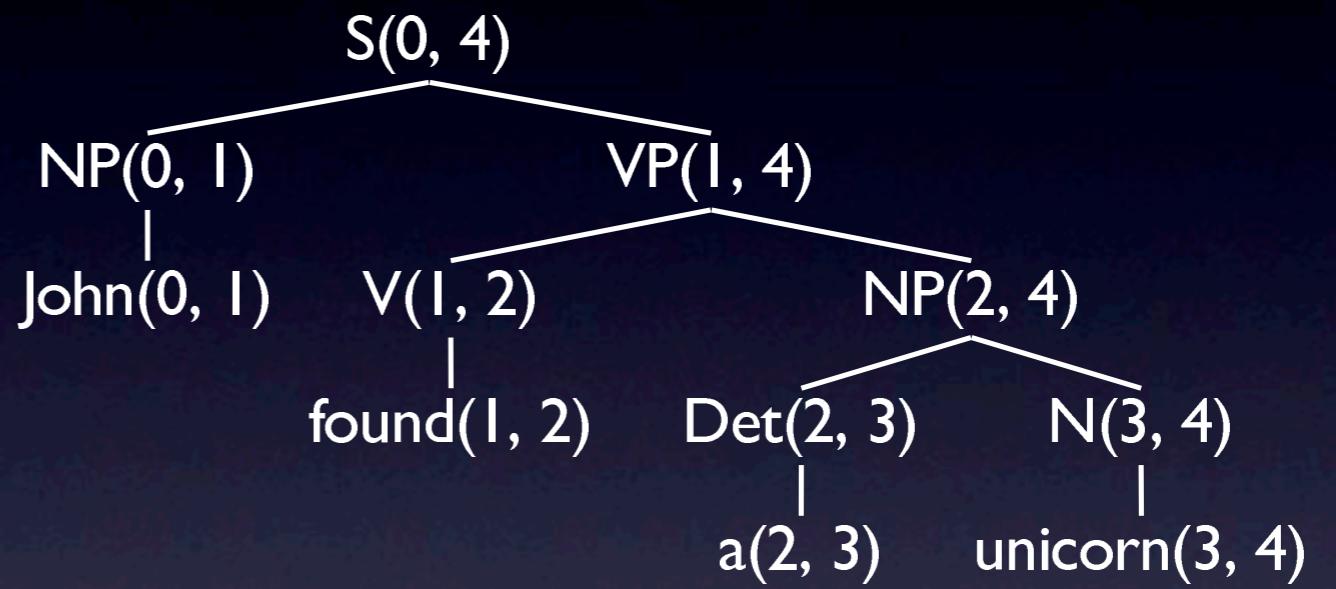
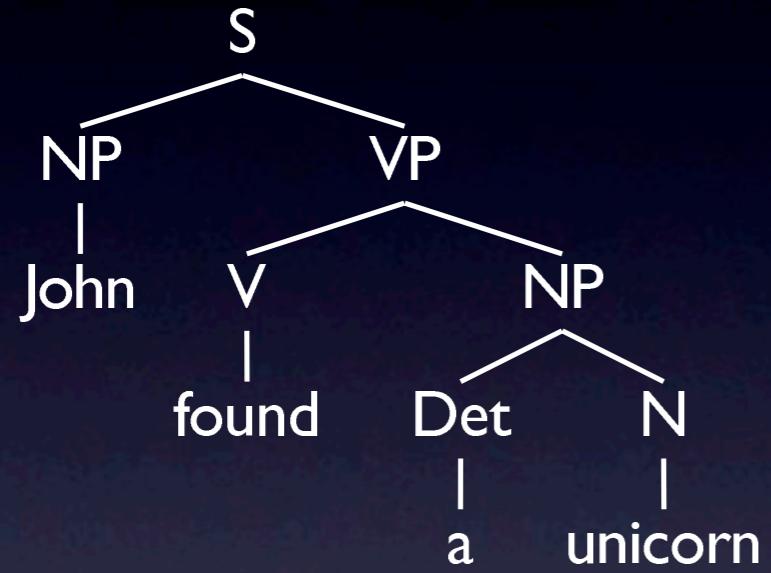
unicorn(3, 4)



$S(i, k) :- \text{NP}(i, j), \text{VP}(j, k).$

?- $S(0, 4).$

The conversion is very straightforward
String \rightarrow string graph



CFG recognition/parsing \approx Datalog query evaluation

CFG derivation tree and Datalog derivation tree isomorphic to each other
Finding one amounts to finding the other

Parsing and generation as Datalog query evaluation

Kanazawa 2007

Recognition/Parsing	CFG MCFG RTG MRTG CFTG _{IO} ⋮
Generation	CFG + Montague semantics [†]

[†] with a certain restriction

Parsing and generation as Datalog query evaluation

- Algorithms
 - Seminaive bottom-up \approx CYK
 - Magic-sets rewriting \approx Earley
- Computational complexity
 - Fixed grammar recognition
 - Uniform recognition
 - Parsing

Polynomial-time algorithm

Seminaive(P, D)

```
1 agenda[0]  $\leftarrow D$            holds facts with  
2  $i \leftarrow 0$                       derivation tree of  
3 chart  $\leftarrow \emptyset$           minimal height  $i$   
4 while agenda[ $i$ ]  $\neq \emptyset$     facts that immediately follow  
5   do                            from one fact in agenda[ $i$ ] plus  
6     chart  $\leftarrow \text{chart} \cup \text{agenda}[i]$  some facts in chart  
7     agenda[ $i + 1$ ]  $\leftarrow \text{Conseq}(P, \text{agenda}[i], \text{chart}) - \text{chart}$   
8      $i \leftarrow i + 1$   
9 return chart
```

well-formed
substring table

Computational complexity

- CFG

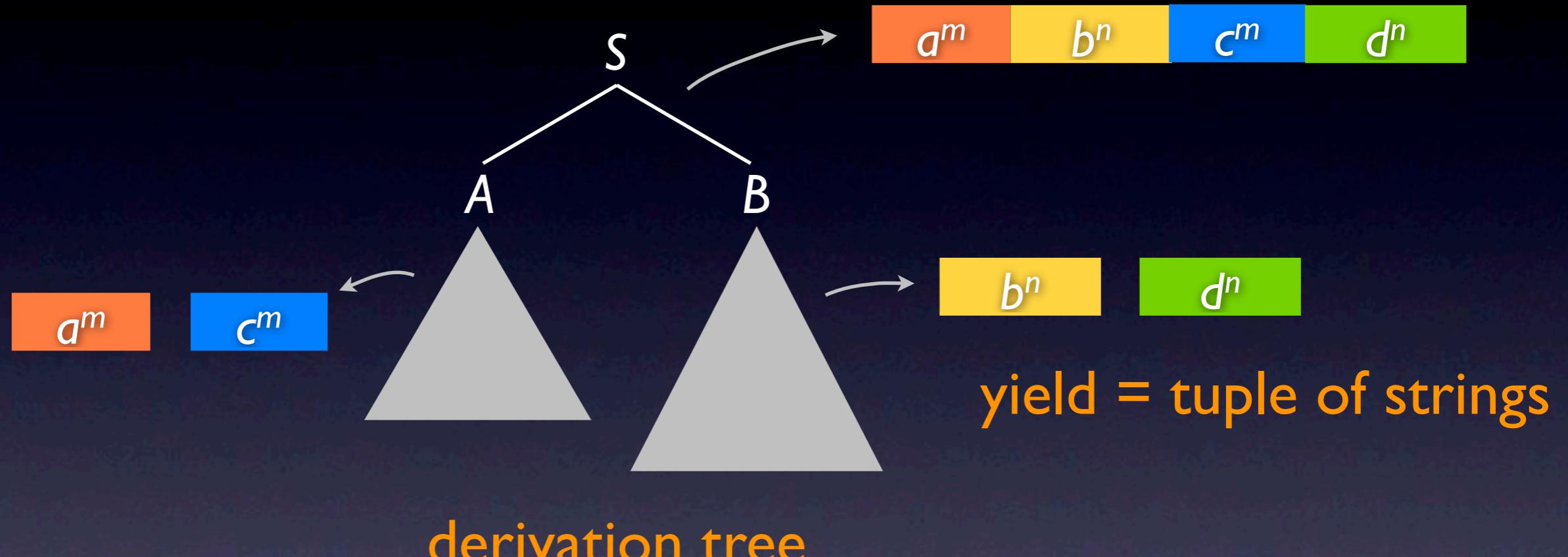
Fixed grammar recognition	LOGCFL-complete
Uniform recognition	P-complete
ϵ -free uniform recognition	LOGCFL-complete

LOGCFL

- The class of problems that reduce to some context-free language in logarithmic space
- The **smallest** computational complexity class that includes the context-free languages

$$\text{AC}^0 \subseteq \text{NC}^1 \subseteq \text{L} \subseteq \text{NL} \subseteq \text{LOGCFL} \subseteq \text{AC}^1 \subseteq \text{NC}^2 \subseteq \text{P} \subseteq \text{NP}$$

Multiple context-free grammars



$S(x_1y_1x_2y_2) \vdash A(x_1, x_2), B(y_1, y_2).$

$A(\varepsilon, \varepsilon).$

$B(\varepsilon, \varepsilon).$

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

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

This is an example of a 2-MCFG.

An m-MCFG allows nonterminals to take up to m arguments.

$S(x_1y_1x_2y_2) :- A(x_1, x_2), B(y_1, y_2).$



$S(i, m) :- A(i, j, k, l), B(j, k, l, m).$

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



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

Range concatenation grammar

- Syntax of EFS (*aka* LMG)
- Variables interpreted as ranges (pairs of string positions)
- Semantics \approx resolution proof of Datalog
- Can be regarded as a subset of Datalog
- Captures P on strings (just like Datalog)
- I focus on a LOGCFL subset of Datalog (on strings, trees, and λ -terms)

CFG + Montague semantics

$S(X_1 X_2) \rightarrow NP(X_1) VP(X_2)$

$VP(\lambda x. X_2 (\lambda y. X_1 yx)) \rightarrow V(X_1) NP(X_2)$

$V(\lambda yx. X_2 (X_1 yx) (X_3 yx)) \rightarrow V(X_1) Conj(X_2) V(X_3)$

$NP(X_1 X_2) \rightarrow Det(X_1) N(X_2)$

$NP(\lambda u. u \mathbf{John}^e) \rightarrow John$

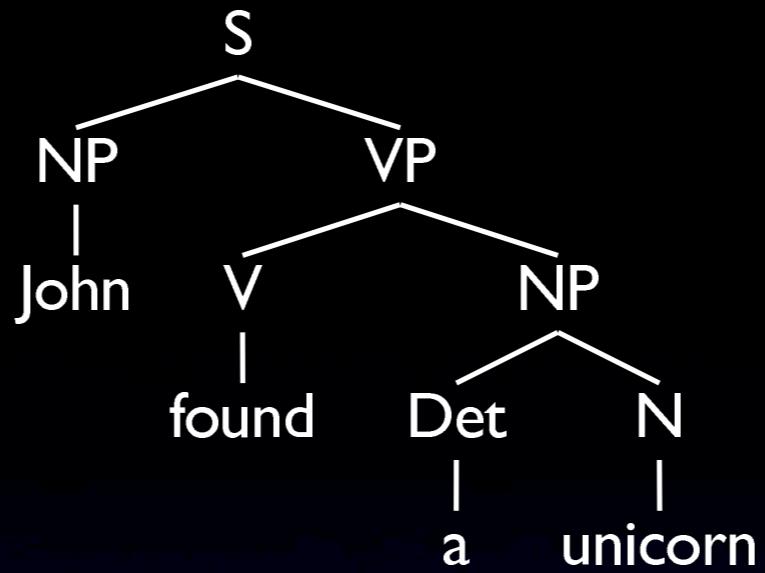
$V(\mathbf{find}^{e \rightarrow e \rightarrow t}) \rightarrow found$

$V(\mathbf{catch}^{e \rightarrow e \rightarrow t}) \rightarrow caught$

$Conj(\wedge^{t \rightarrow t \rightarrow t}) \rightarrow and \quad \lambda uv. \exists y (u(y) \wedge v(y))$

$Det(\lambda uv. \exists^{(e \rightarrow t) \rightarrow t} (\lambda y. \wedge^{t \rightarrow t \rightarrow t} (uy)(vy))) \rightarrow a$

$N(\mathbf{unicorn}^{e \rightarrow t}) \rightarrow unicorn$



$S((\lambda u.u \text{ John})(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uy)(vy))) \text{ unicorn } (\lambda y.\text{find } y x)))$

$\overbrace{\text{NP}(\lambda u.u \text{ John}) \quad \text{VP}(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uy)(vy))) \text{ unicorn } (\lambda y.\text{find } y x))}^{\text{NP}(\lambda u.u \text{ John}) \quad \text{VP}(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uy)(vy))) \text{ unicorn } (\lambda y.\text{find } y x))}$

$\overbrace{\text{NP}(\lambda u.u \text{ John}) \quad \text{VP}(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uy)(vy))) \text{ unicorn } (\lambda y.\text{find } y x))}^{\text{NP}(\lambda u.u \text{ John}) \quad \text{VP}(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uy)(vy))) \text{ unicorn } (\lambda y.\text{find } y x))}$

$(\lambda u.u \text{ John})(\lambda x.(\lambda uv.\exists(\lambda y.\wedge(uy)(vy))) \text{ unicorn } (\lambda y.\text{find } y x))$

$\rightarrow_{\beta} \exists(\lambda y.\wedge(\text{unicorn } y)(\text{find } y \text{ John}))$

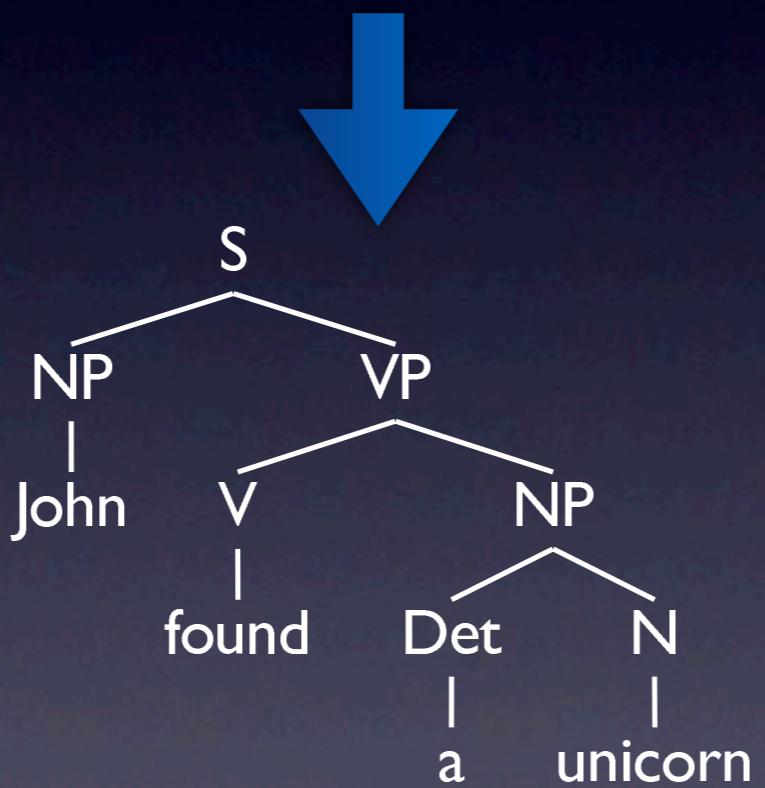
$\approx \exists y(\text{unicorn}(y) \wedge \text{find}(\text{John}, y))$

logical form

Surface realization

$\exists y(\text{unicorn}(y) \wedge \text{find}(\text{John}, y))$

input logical form



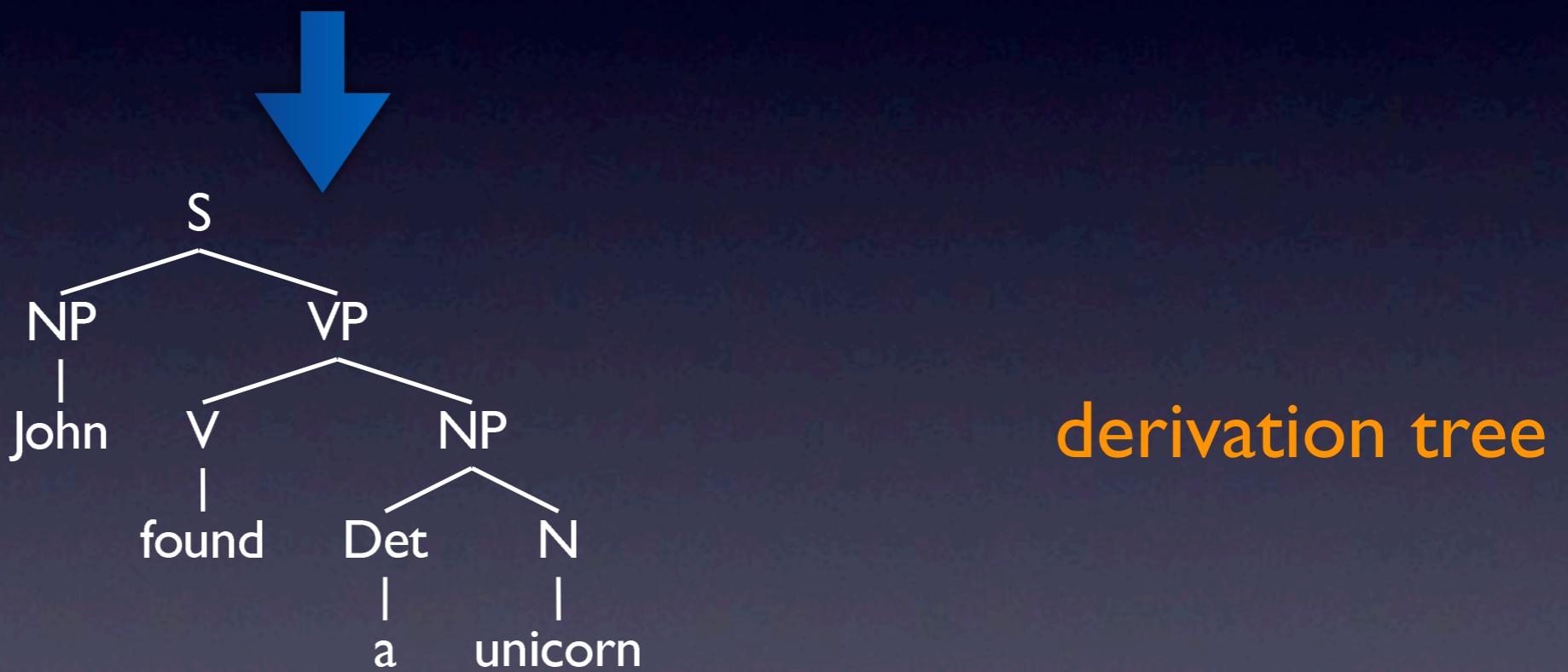
derivation tree

John found a unicorn

output surface form

Parsing of input logical form

$\exists y(\text{unicorn}(y) \wedge \text{find}(\text{John}, y))$ input logical form



Recognition of surface realizability

$\exists y(\mathbf{unicorn}(y) \wedge \mathbf{find}(\mathbf{John}, y))$

input logical form

surface realizable?



yes/no

Context-free grammars on λ -terms

$S(X_1 X_2) \vdash NP(X_1) VP(X_2)$

$VP(\lambda x. X_2 (\lambda y. X_1 yx)) \vdash V(X_1), NP(X_2).$

$V(\lambda yx. X_2 (X_1 yx) (X_3 yx)) \vdash V(X_1), Conj(X_2), V(X_3).$

$NP(X_1 X_2) \vdash Det(X_1), N(X_2).$

$NP(\lambda u. u \mathbf{John}^e).$

$V(\mathbf{find}^{e \rightarrow e \rightarrow t}).$

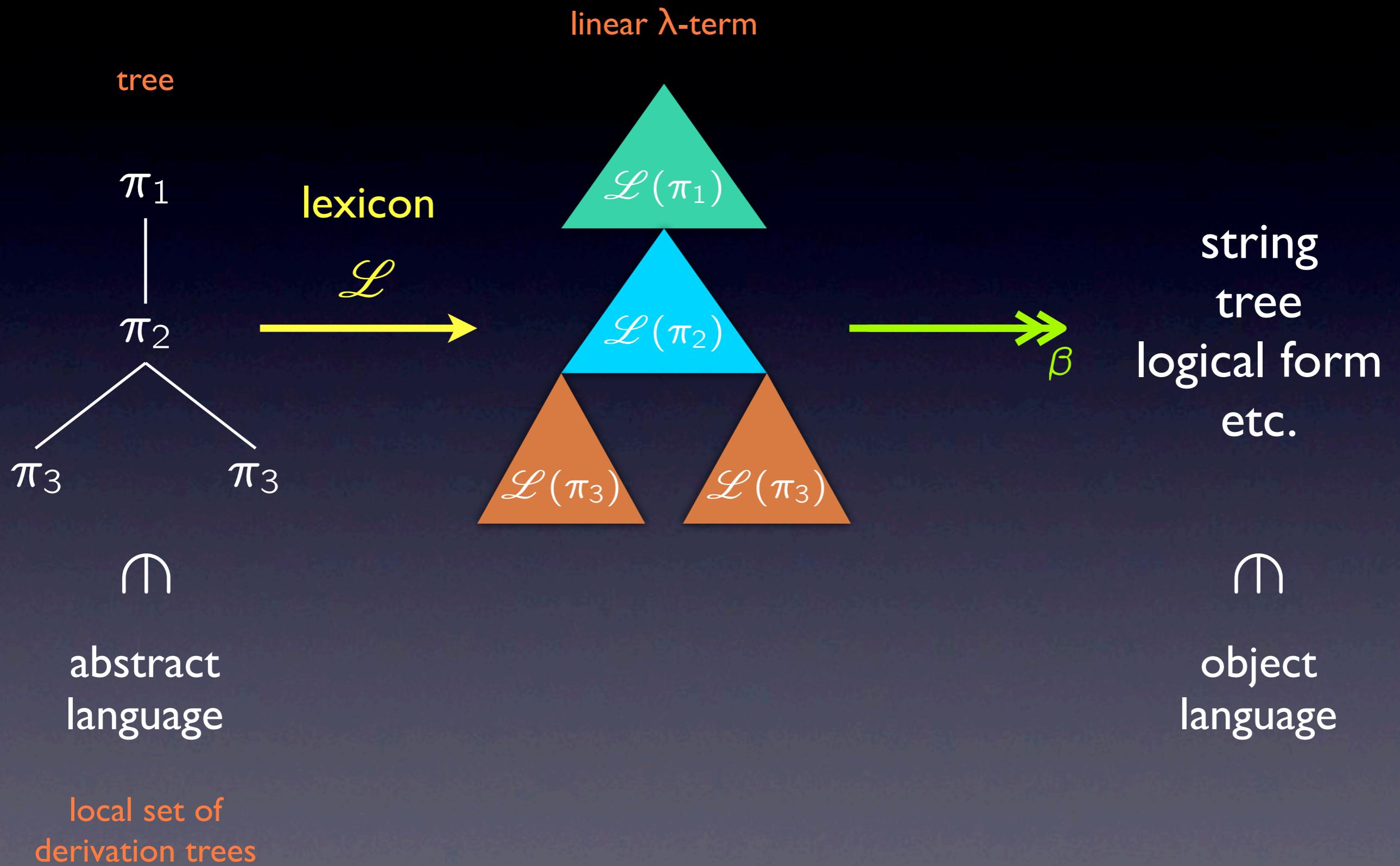
$V(\mathbf{catch}^{e \rightarrow e \rightarrow t}).$

$Conj(\wedge^{t \rightarrow t \rightarrow t}).$

$Det(\lambda uv. \exists^{(e \rightarrow t) \rightarrow t} (\lambda y. \wedge^{t \rightarrow t \rightarrow t} (uy)(vy))).$

$N(\mathbf{unicorn}^{e \rightarrow t}).$

Second-order ACGs



Context-free grammars on λ -terms = second-order **non-linear** ACGs

$$\pi : B(M) \coloneq B_1(X_1), \dots, B_n(X_n).$$
$$\pi : B_1 \rightarrow \dots \rightarrow B_n \rightarrow B : \lambda X_1 \dots X_n . M$$

abstract
constant

type of π

object
realization of
 π

From CFLG recognition to Datalog query evaluation

$S(X_1 X_2) \vdash NP(X_1) \vee NP(X_2)$

$VP(\lambda x. X_2 (\lambda y. X_1 yx)) \vdash V(X_1), NP(X_2).$

$V(\lambda yx. X_2 (X_1 yx) (X_3 yx)) \vdash V(X_1), Conj(X_2), V(X_3).$

$NP(X_1 X_2) \vdash Det(X_1), N(X_2).$

$NP(\lambda u. u \text{ John}^e).$

$V(\text{find}^{e \rightarrow e \rightarrow t}).$

$V(\text{catch}^{e \rightarrow e \rightarrow t}).$

$Conj(\wedge^{t \rightarrow t \rightarrow t}).$

$Det(\lambda uv. \exists^{(e \rightarrow t) \rightarrow t} (\lambda y. \wedge^{t \rightarrow t \rightarrow t} (uy)(vy))).$

$N(\text{unicorn}^{e \rightarrow t}).$

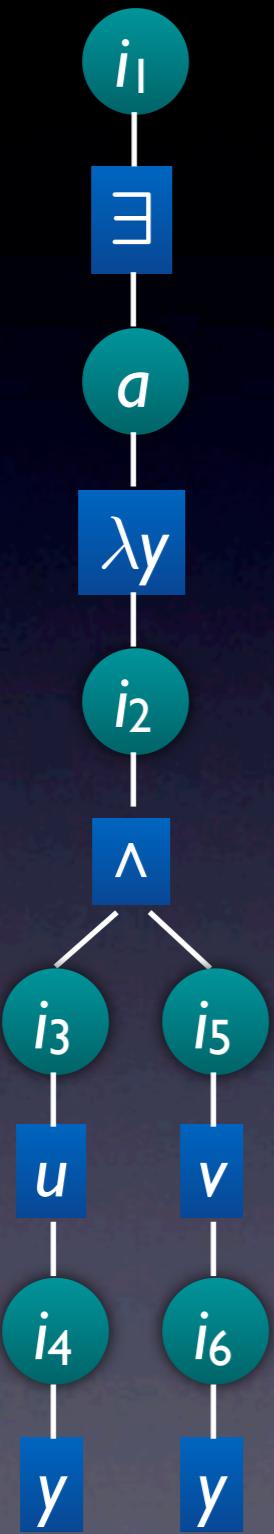
$\exists(\lambda y. \wedge(\text{unicorn } y)(\text{find } y \text{ John})) \in L(G)?$

From CFLG recognition to Datalog query evaluation

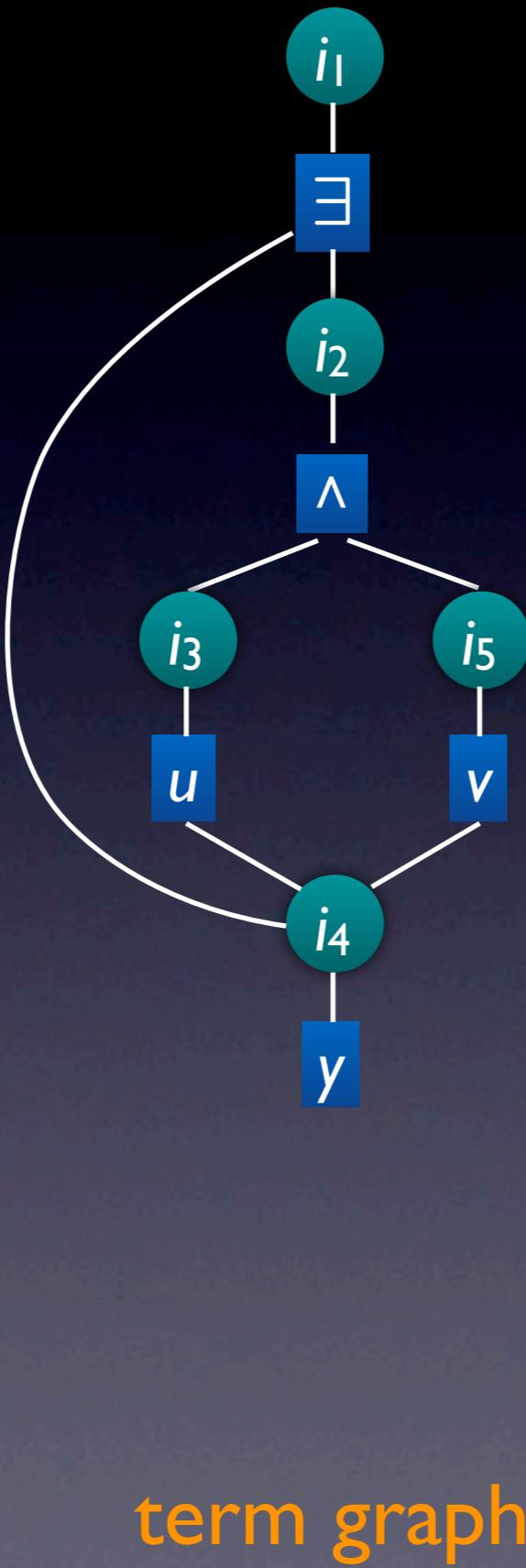
```
S( $i_1$ ) :- NP( $i_1, i_2, i_3$ ), VP( $i_2, i_3$ ).  
VP( $i_1, i_4$ ) :- V( $i_2, i_4, i_3$ ), NP( $i_1, i_2, i_3$ ).  
V( $i_1, i_4, i_3$ ) :- V( $i_2, i_4, i_3$ ), Conj( $i_1, i_5, i_2$ ), V( $i_5, i_4, i_3$ ).  
NP( $i_1, i_4, i_5$ ) :- Det( $i_1, i_4, i_5, i_2, i_3$ ), N( $i_2, i_3$ ).  
NP( $i_1, i_1, i_2$ ) :- John( $i_2$ ).  
V( $i_1, i_3, i_2$ ) :- find( $i_1, i_3, i_2$ ).  
V( $i_1, i_3, i_2$ ) :- catch( $i_1, i_3, i_2$ ).  
Conj( $i_1, i_3, i_2$ ) :-  $\wedge(i_1, i_3, i_2)$ .  
Det( $i_1, i_5, i_4, i_3, i_4$ ) :-  $\exists(i_1, i_2, i_4)$ ,  $\wedge(i_2, i_5, i_3)$ .  
N( $i_1, i_2$ ) :- unicorn( $i_1, i_2$ ).
```

```
 $\exists(i_1, i_2, i_4)$ .  
 $\wedge(i_2, i_5, i_3)$ .  
unicorn( $i_3, i_4$ ).  
find( $i_5, i_6, i_4$ ).  
John( $i_6$ ).
```

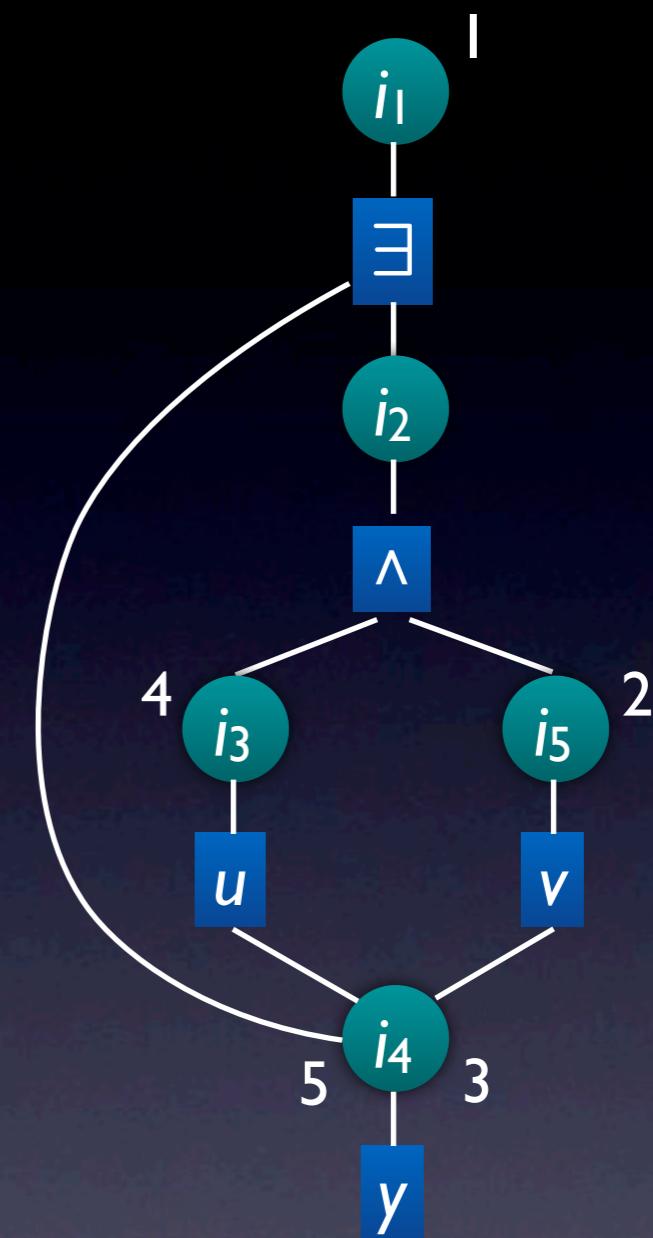
```
?- S( $i_1$ ).
```

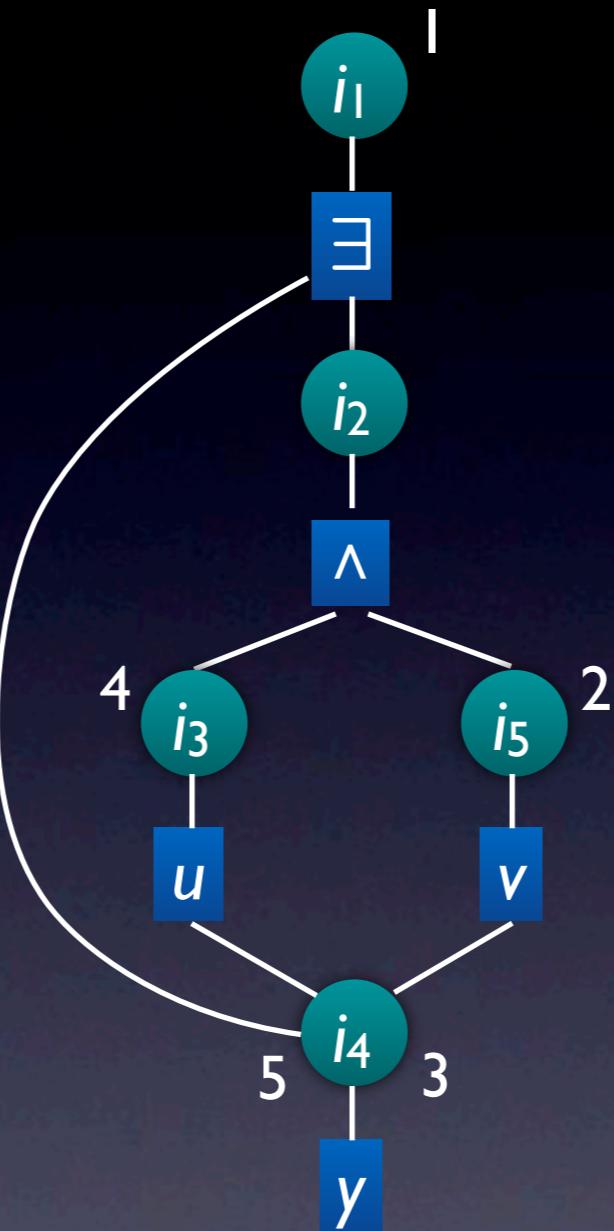
$\exists(\lambda y.\wedge(uy)(vy))$ 

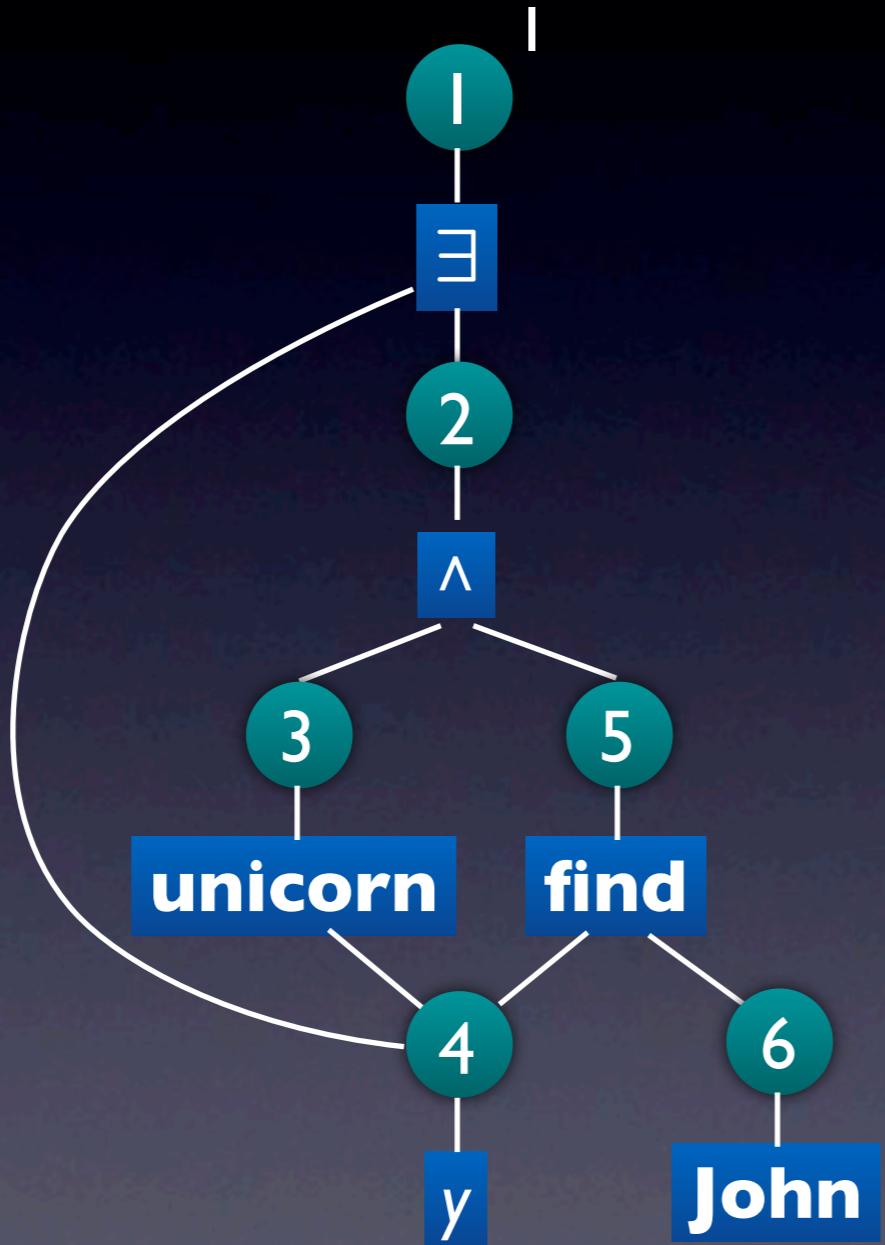
tree graph

 $\lambda uv.\exists(\lambda y.\wedge(uy)(vy))$ 

term graph

term graph with
external nodes

$$\text{Det}(\lambda uv. \exists^{(e \rightarrow t) \rightarrow t} (\lambda y. \wedge^{t \rightarrow t \rightarrow t} (uy)(vy))).$$

$$\text{Det}(i_1, i_5, i_4, i_3, i_4) \leftarrow \exists(i_1, i_2, i_4), \wedge(i_2, i_5, i_3).$$

$$\exists(\lambda y. \wedge(\mathbf{unicorn}y)(\mathbf{find}y \mathbf{John}))$$


$\exists(1, 2, 4).$
 $\wedge(2, 5, 3).$
unicorn(3, 4).
find(5, 6, 4).
John(6).

?– S(I).

From CFLG recognition to Datalog query evaluation

- The reduction is correct when all λ -terms in the grammar are **almost linear**.

$$\frac{}{x : \alpha \vdash x : \alpha} \quad \frac{}{\vdash c : \tau(c)}$$

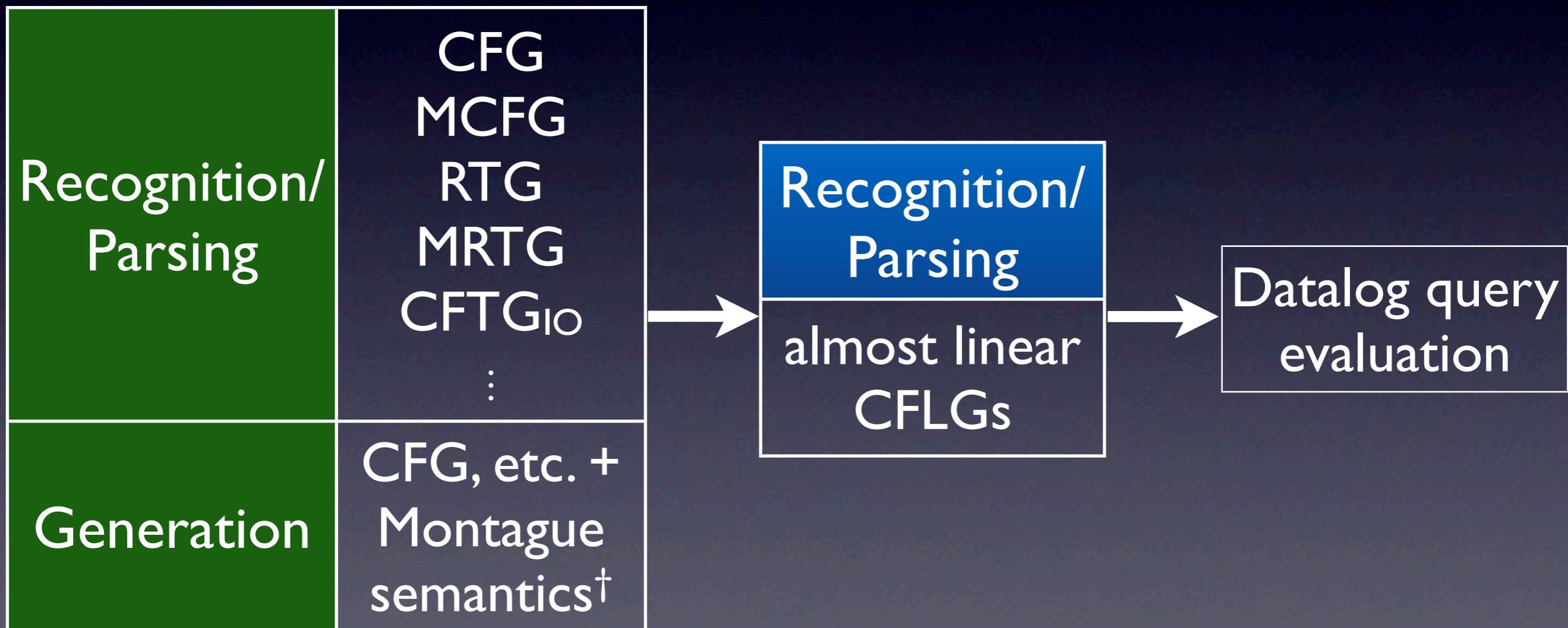
$$\frac{\Gamma \vdash M : \alpha \quad \Delta \vdash N : \beta}{\Gamma \cup \Delta \vdash MN : \beta} \text{ if } \Gamma \cap \Delta \subseteq \text{At} \quad \frac{\Gamma, x : \alpha \vdash M : \beta}{\Gamma \vdash \lambda x. M : \alpha \rightarrow \beta}$$

Almost linear λ -terms

- **Coherence Theorem:** An almost linear λ -term is determined by its principal typing up to $\beta\eta$ -equality.
- **Subject Expansion Theorem:** An almost linear λ -term has the same principal typing as its β -normal form.

Parsing and generation as Datalog query evaluation

Kanazawa 2007



[†]when almost linear

Computational complexity

- Almost linear CFLGs with **bounded width** and **rank**

Fixed grammar recognition	LOGCFL-complete
Uniform recognition	P-complete
ϵ -free uniform recognition	LOGCFL-complete

width = arity of Datalog predicates

rank = number of subgoals in rules

ϵ -rule = Datalog rule with empty right-hand side

LOGCFL

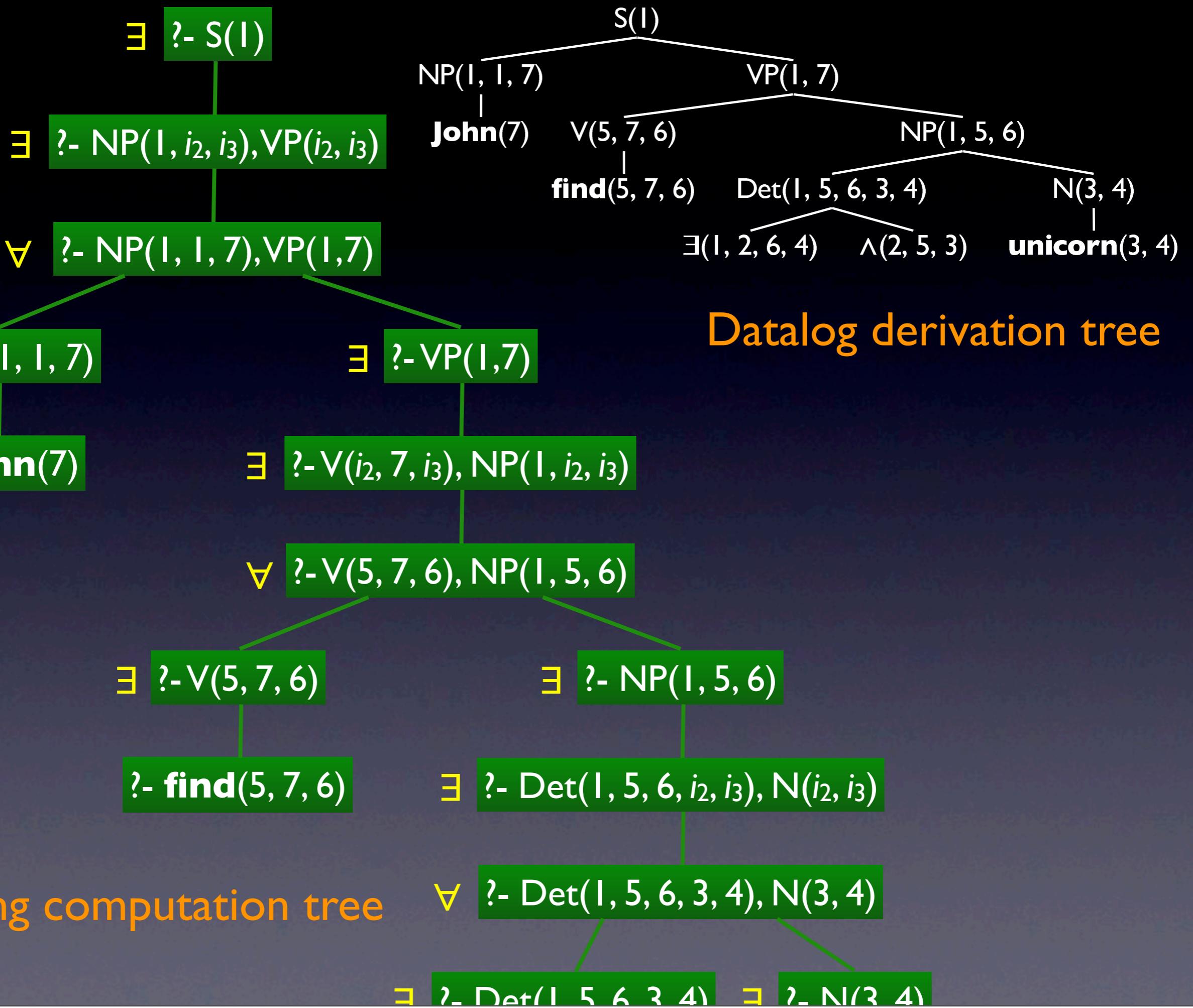
- Characterized in terms of **alternating** Turing machines operating in **log space** with **polynomial-size accepting computation trees**

Ruzzo 1980



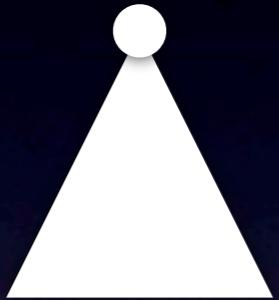
Alternating algorithm for Datalog query evaluation

- Each configuration is an attempt to prove some ground fact
- Ǝ
 - Pick a matching rule
 - Pick constants that instantiate the rule
- ∀
 - Prove all subgoals of the instantiated rule
 - Can operate in log space if width and rank are bounded



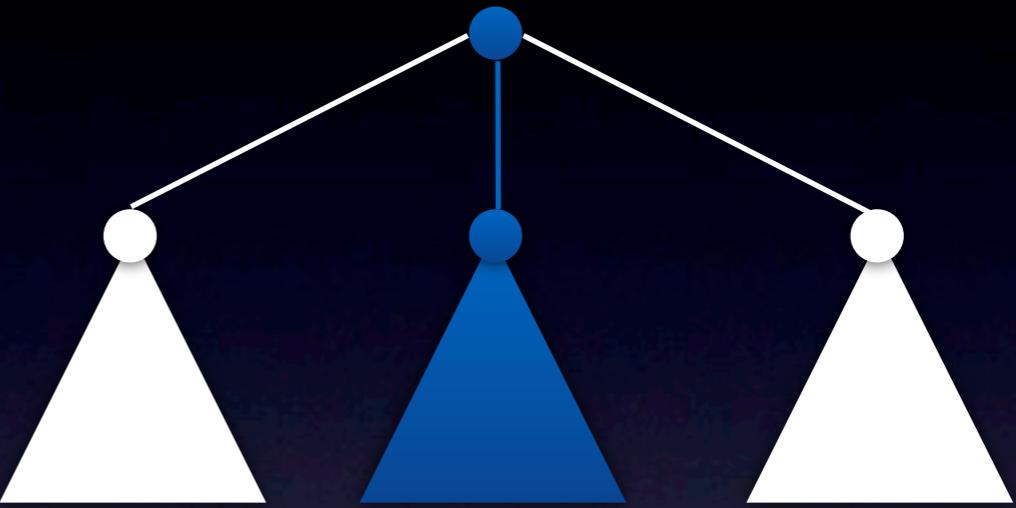
Size of Datalog derivation tree

Type 0



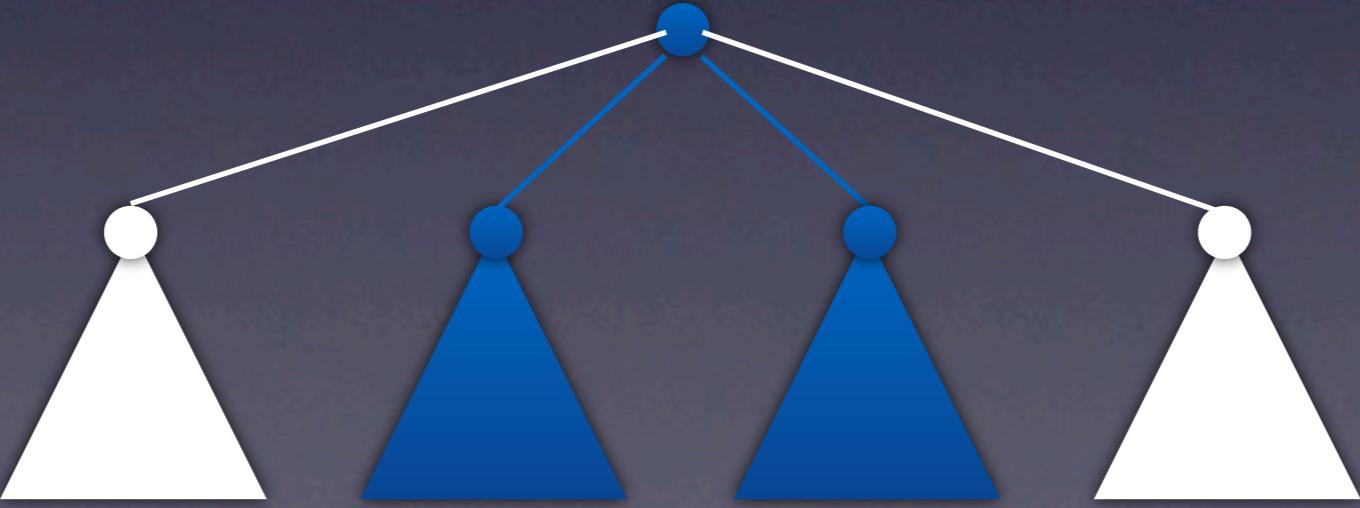
no database facts below

Type I



all children **except one** are Type 0

Type 2



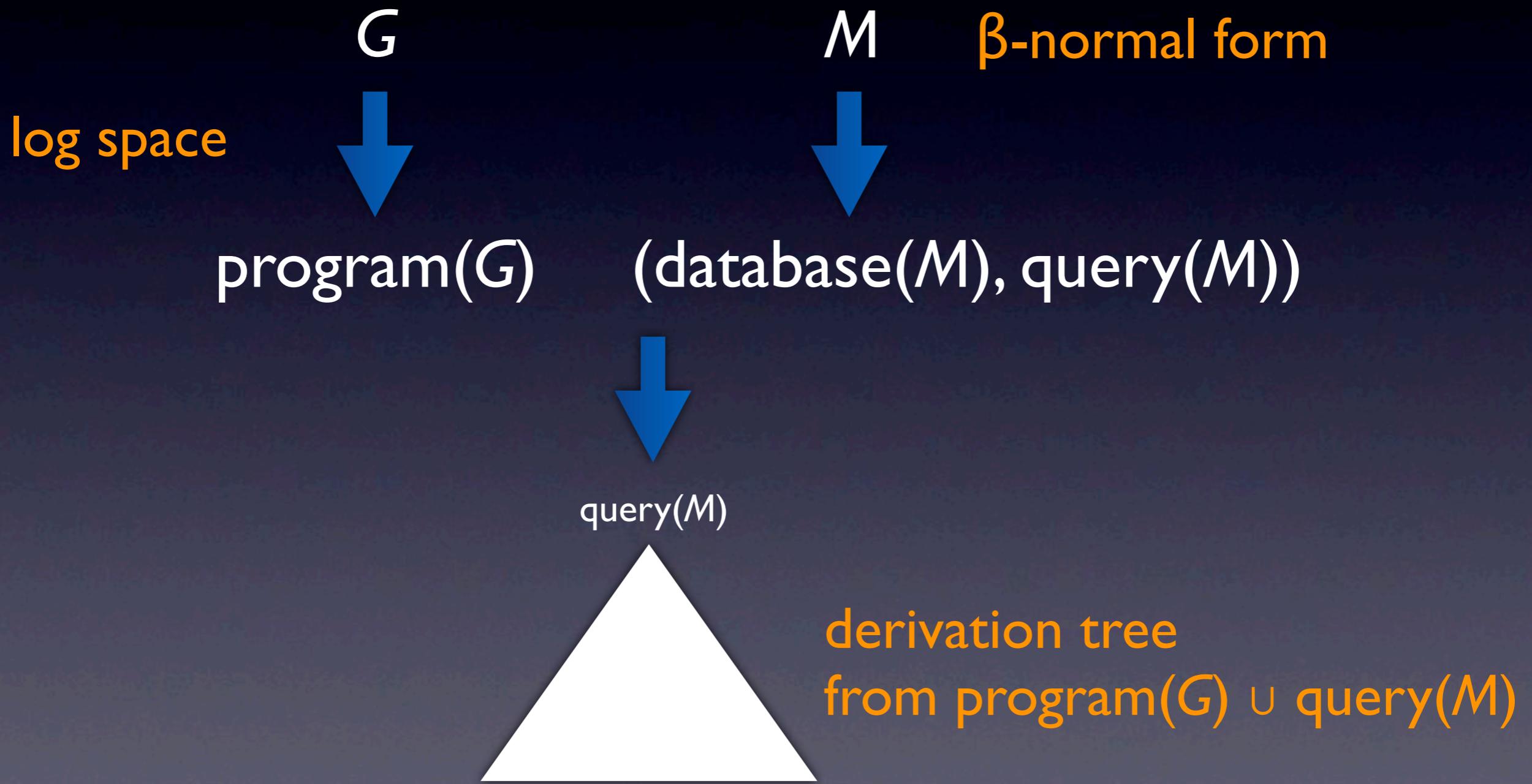
at least two children that are not Type 0

Number of nodes

n = number of extensional nodes
(nodes labeled by database facts)

- Type 0: **exponential** in the number of predicates, polynomial in n
- Type 1: linear in the number of predicates, polynomial in n
- Type 2: at most $n-1$

Alternating algorithm for almost linear CFLG recognition



The number of extensional nodes bounded by the size of M

Complexity of parsing

- Almost linear CFLGs with **bounded width** and **rank**
- Output: parse forest = set of all ground rule instances useful for deriving the input query

Fixed grammar	$\text{FL}^{\text{LOGCFL}}$
Uniform	FP
ϵ -free uniform	$\text{FL}^{\text{LOGCFL}}$