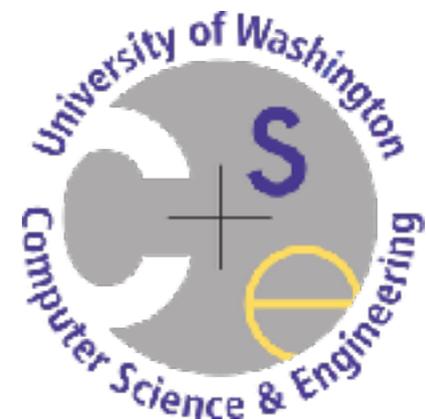


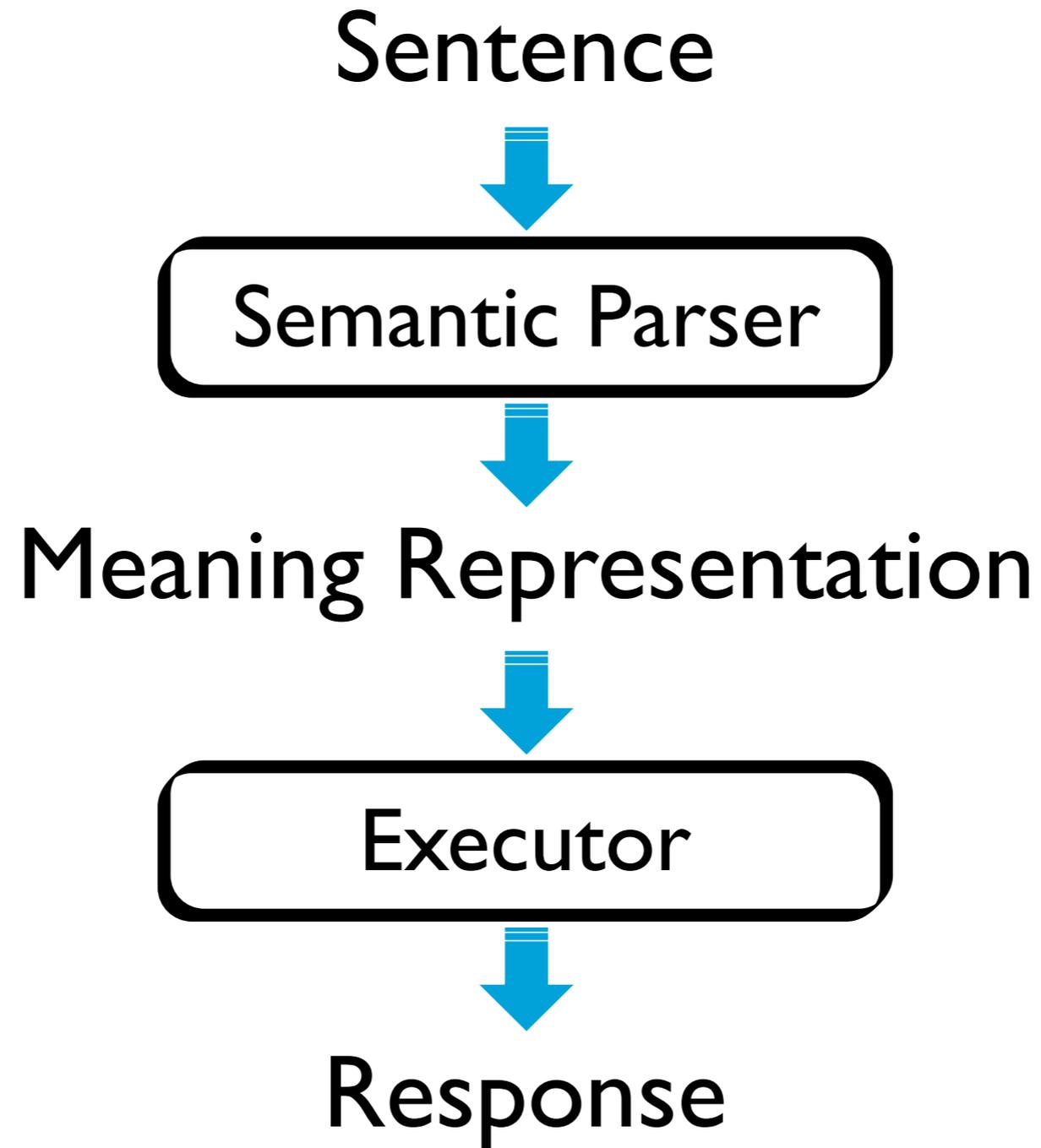
Interactive Learning of Parsers from Weak Supervision

Luke Zettlemoyer

with Luheng He, Kenton Lee, Mike Lewis, Julian Michael

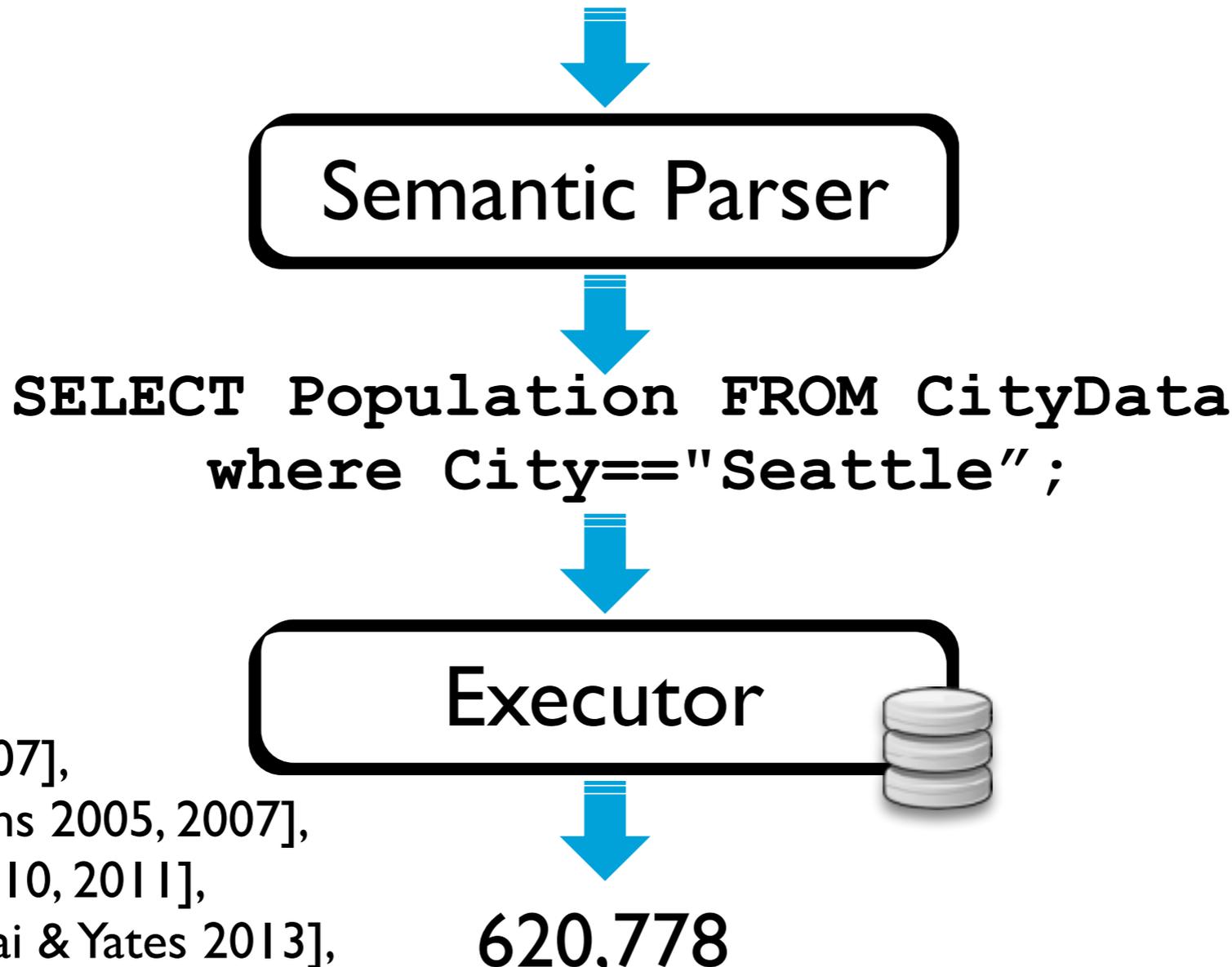


Interpreting Language



Semantic Parsing: QA

How many people live in Seattle?



Semantic Parser

```
SELECT Population FROM CityData  
where City=="Seattle";
```

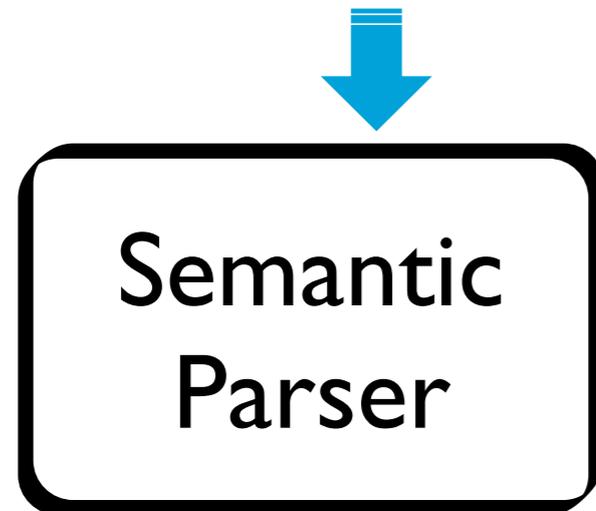
Executor

620,778

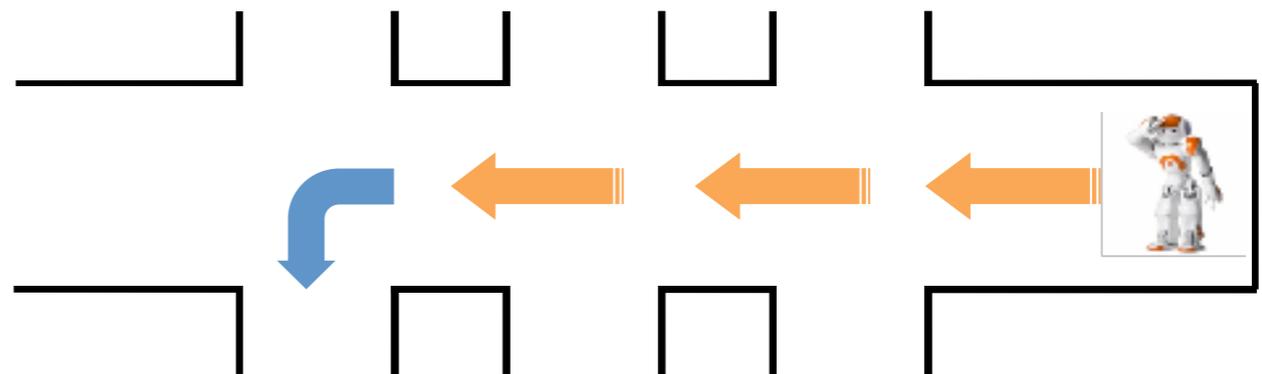
[Wong & Mooney 2007],
[Zettlemoyer & Collins 2005, 2007],
[Kwiatkowski et.al 2010, 2011],
[Liang et.al. 2011], [Cai & Yates 2013],
[Berant et.al. 2013,2014,2015],
[Kwiatkowski et.al. 2013],
[Reddy et.al, 2014,2016]

Semantic Parsing: Instructions

Go to the third junction and take a left



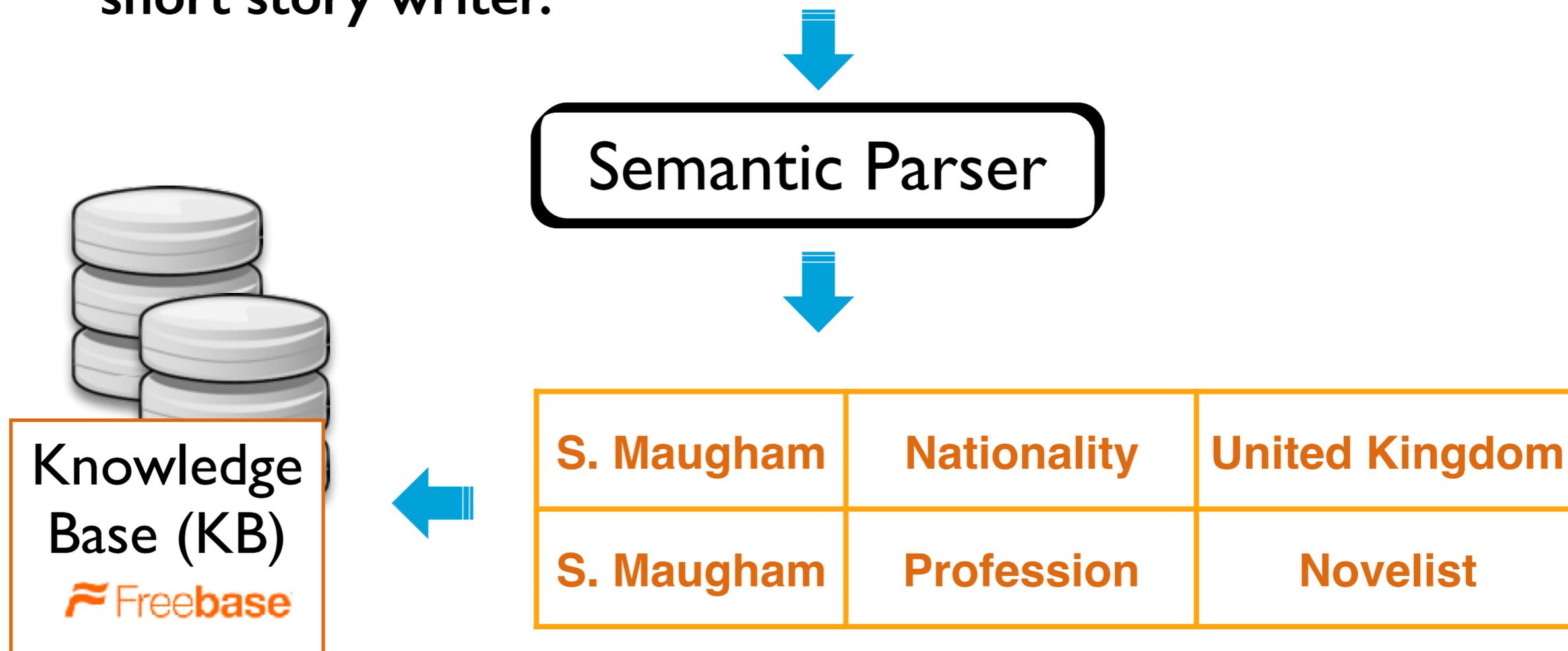
```
(do-seq (do-n-times 3
  (move-to forward-loc
    (do-until
      (junction current-loc
        (move-to forward-loc))))
  (turn-right)))
```



- [Chen & Mooney 2011]
- [Matuszek et.al. 2012]
- [Artzi & Zettlemoyer 2013]
- [Mei et.al. 2015]

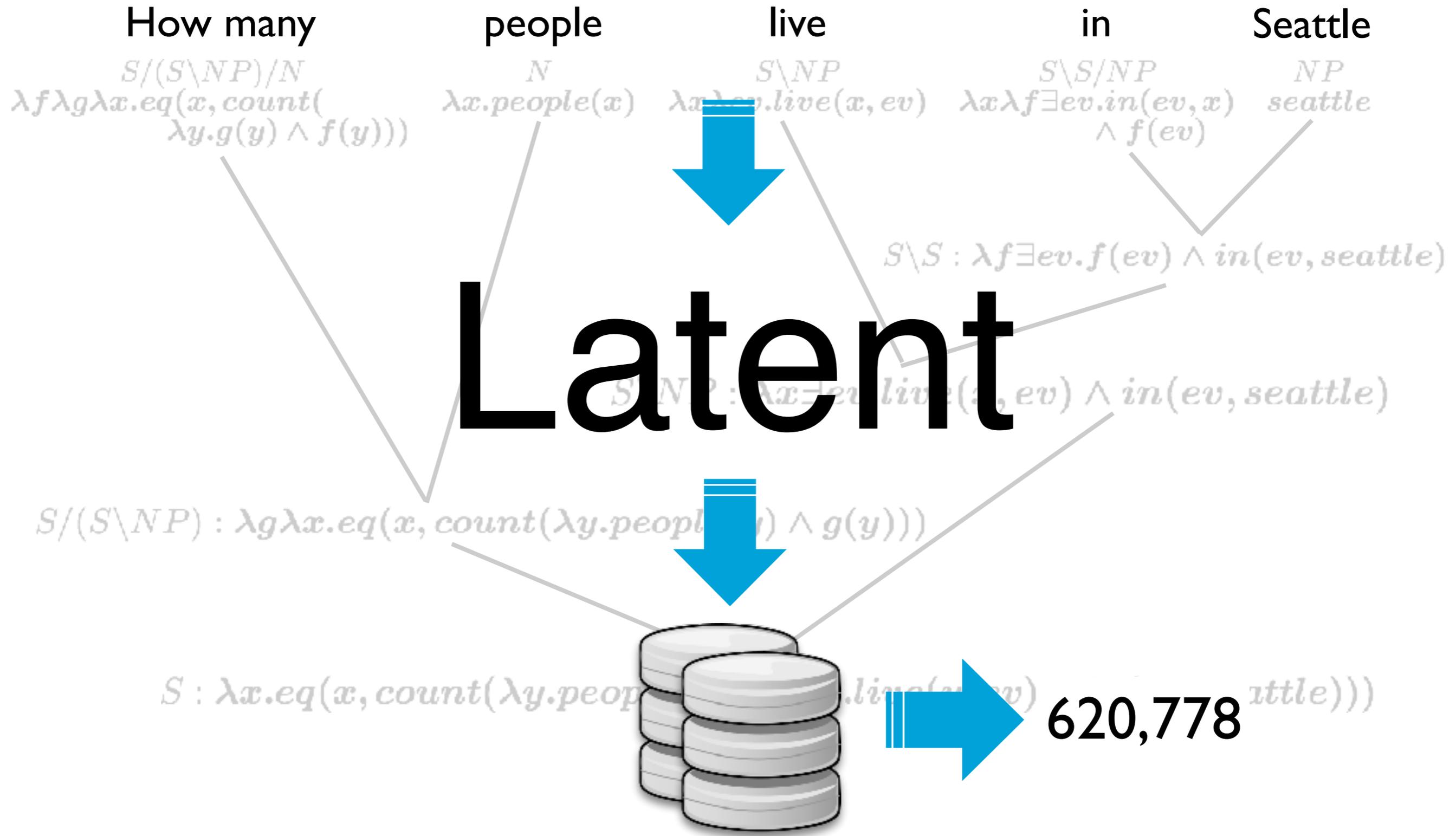
Semantic Parsing: IE

Somerset Maugham was a British playwright, novelist and short story writer.



[Krishnamurthy and Mitchell;
2012, 2014][Choi et al., 2015]

Semantic Parsing: Complex Structure



Lots of Different Applications

We are doing semantic analysis for:

- Visual Semantic Role Labeling [Yatskar et al, 2016]
- Visual Question Answering [FitzGerald et al, in prep]
- Language to Code [Lin et al, in prep]
- Entity-entity sentiment [Choi et al, 2016]
- Understanding Cooking Recipes [Kiddon et al, 2016]
- Zero-shot Relation Extraction [Levy et al, in review]
- Interactive Learning for NLIDBs [Iyer, et al, in review]

Challenge: typically gather data and learn model from scratch in each case...

Understanding Cooking Recipes

Amish Meatloaf (<http://allrecipes.com/recipe/amish-meatloaf/>, recipe condensed)

Ingredients

2 pounds ground beef
2 1/2 cups crushed butter-flavored crackers
1 small onion, chopped
2 eggs
3/4 cup ketchup
1/4 cup brown sugar
2 slices bacon

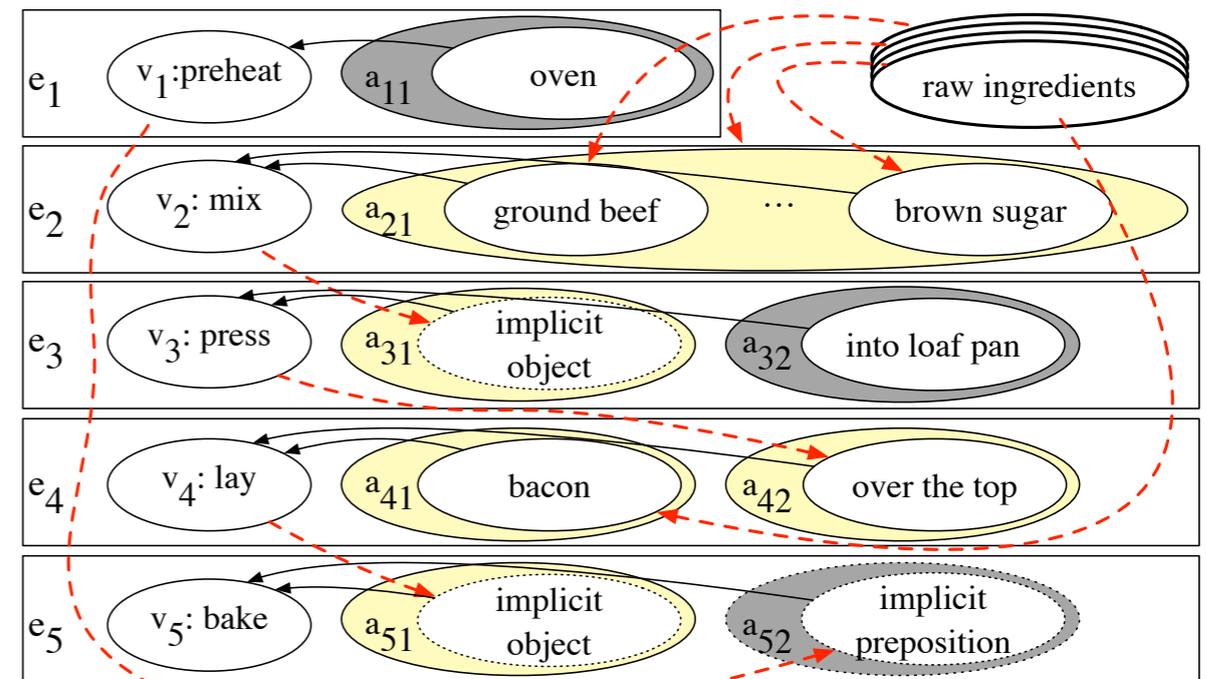
Preheat the oven to 350 degrees F (175 degrees C).

In a medium bowl, mix together ground beef, crushed crackers, onion, eggs, ketchup, and brown sugar until well blended.

Press into a 9x5 inch loaf pan.

Lay the two slices of bacon over the top.

Bake for 1 hour, or until cooked through.



Approach: unsupervised learning for actions and object flow

Open Question:

- Can we build an off-the-shelf parser that would help here?

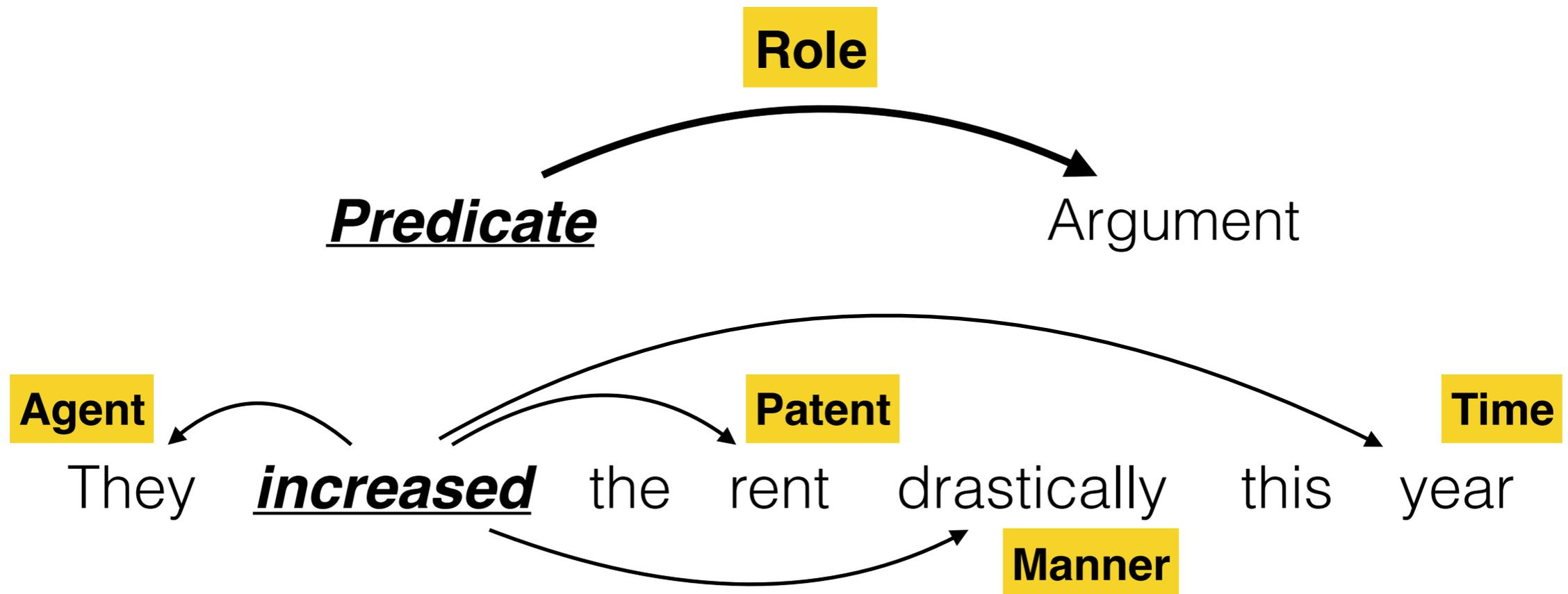
[Kiddon et al 2015, 2016]

Towards Broad Coverage Semantic Parsing

- Can we crowdsource semantics?
- Train with latent syntax?
- Build fast and accurate parsers?
- Actively select which data to label?

Semantic Role Labeling (SRL)

who did **what** to **whom**, **when** and **where**?



- Defining a set of roles can be difficult
- Existing formulations have used different sets

Existing SRL Formulations and Their Frame Inventories

FrameNet

1000+ semantic frames,
roles (frame elements)
shared across frames

Frame: **Change_position_on_a_scale**

This frame consists of words that indicate the change of an **Item**'s position on a scale (the **Attribute**) from a starting point (**Initial_value**) to an end point (**Final_value**). The direction (**Path**) ...

Lexical Units:

..., *reach.v*, *rise.n*, *rise.v*, *rocket.v*, *shift.n*, ...

PropBank

10,000+ frame files
with predicate-specific roles

Roleset Id: **rise.01** , *go up*

Arg1-: *Logical subject, patient, thing rising*

Arg2-EXT: *EXT, amount risen*

Arg3-DIR: *start point*

Arg4-LOC: *end point*

Argm-LOC: *medium*

Unified Verb Index, University of Colorado <http://verbs.colorado.edu/verb-index/>

PropBank Annotation Guidelines, Bonial et al., 2010

FrameNet II: Extended theory and practice, Ruppenhofer et al., 2006

Our Annotation Scheme

Given sentence and a verb:

They ***increased*** the rent this year .

**Step 1: Ask a question
about the verb:**

Who increased something ?

**Step 2: Answer with words
in the sentence:**

They

**Step 3: Repeat, write as many
QA pairs as possible ...**

What is increased ?

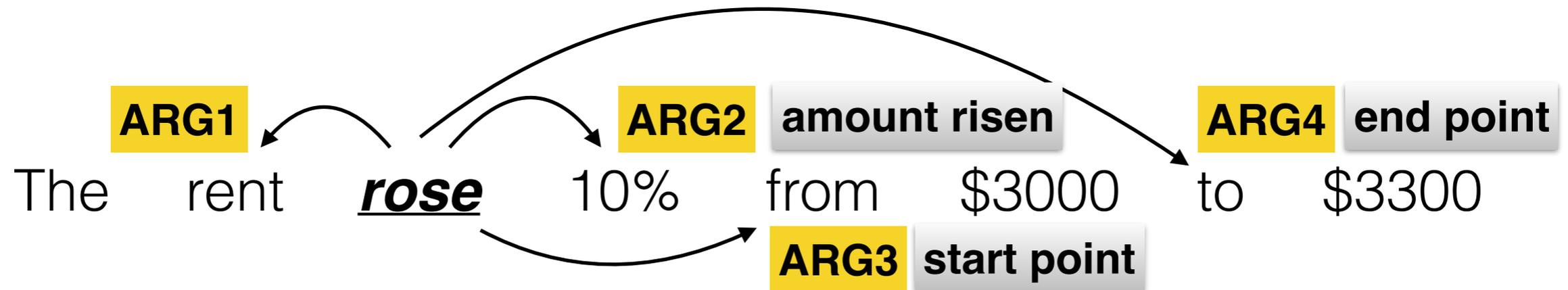
the rent

When is something increased ?

this year

[He et al 2015]

Our Method: Q/A Pairs for Semantic Relations



Wh-Question

Answer

What *rose* ?

the rent

How much did something *rise* ?

10%

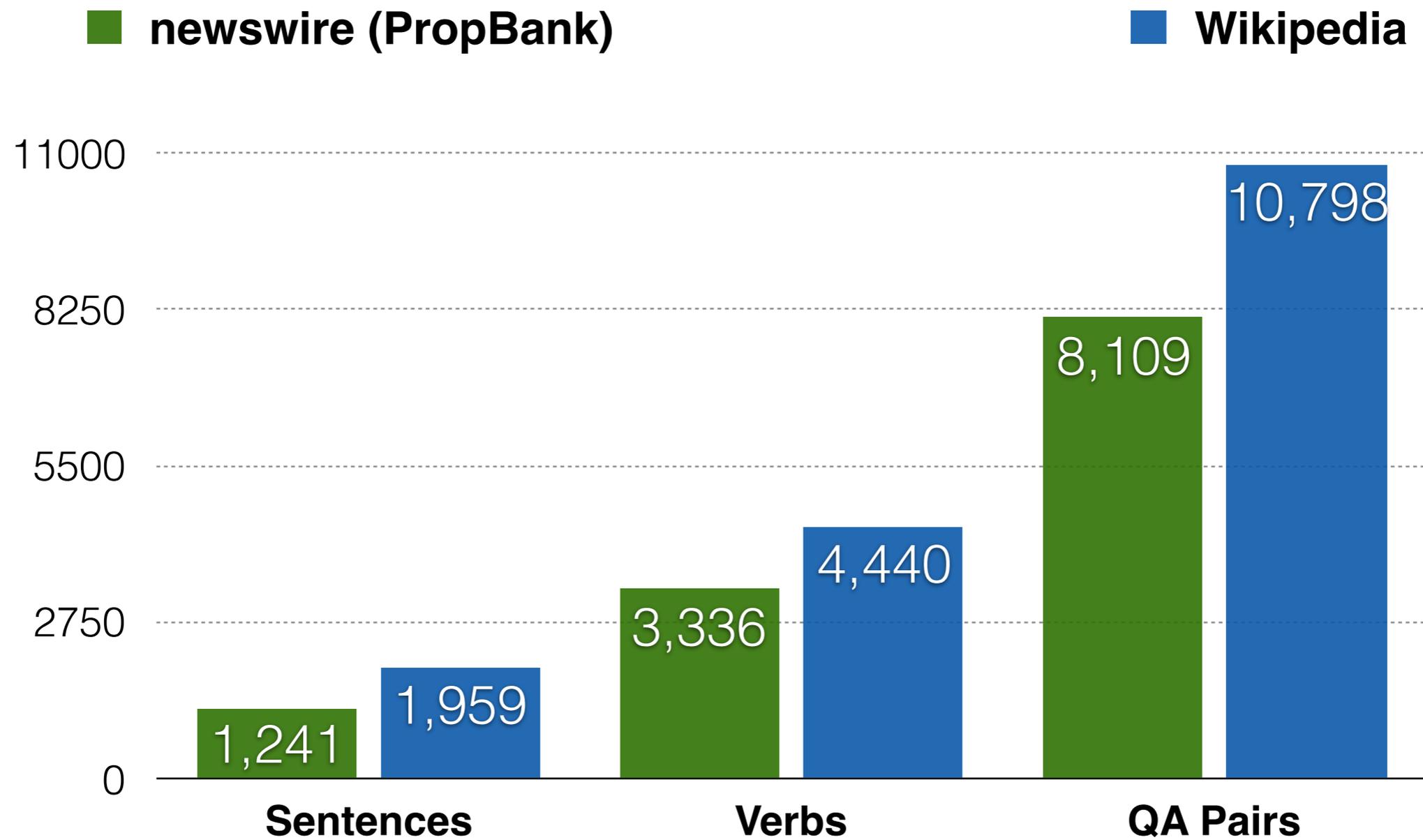
What did something *rise* from ?

\$3000

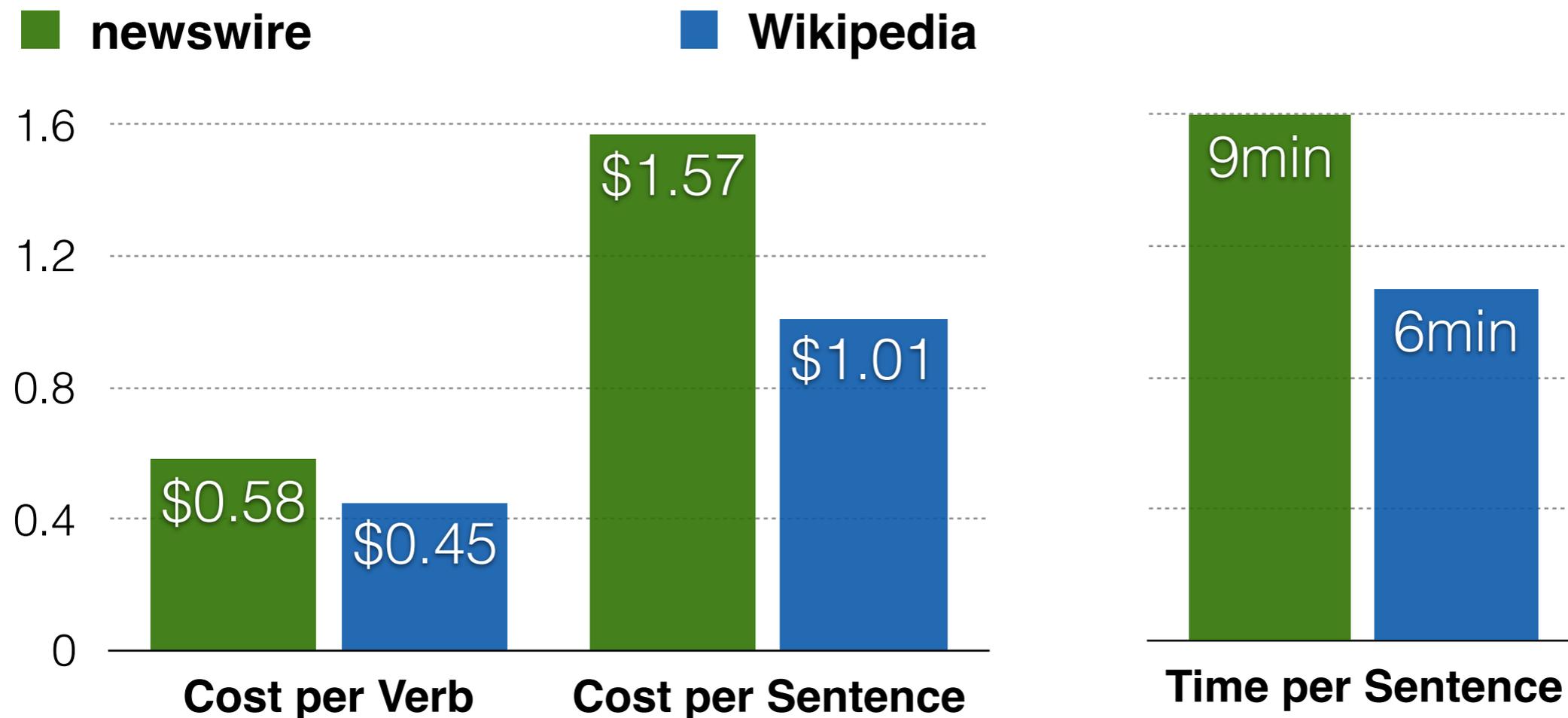
What did something *rise* to ?

\$3300

Dataset Statistics



Cost and Speed



- Part-time freelancers from [upwork.com](https://www.upwork.com) (hourly rate: \$10)
- ~2h screening process for native English proficiency

Wh-words vs. PropBank Roles

	Who	What	When	Where	Why	How	HowMuch
ARG0	1575	414	3	5	17	28	2
ARG1	285	2481	4	25	20	23	95
ARG2	85	364	2	49	17	51	74
ARG3	11	62	7	8	4	16	31
ARG4	2	30	5	11	2	4	30
ARG5	0	0	0	1	0	2	0
AM-ADV	5	44	9	2	25	27	6
AM-CAU	0	3	1	0	23	1	0
AM-DIR	0	6	1	13	0	4	0
AM-EXT	0	4	0	0	0	5	5
AM-LOC	1	35	10	89	0	13	11
AM-MNR	5	47	2	8	4	108	14
AM-PNC	2	21	0	1	39	7	2
AM-PRD	1	1	0	0	0	1	0
AM-TMP	2	51	341	2	11	20	10

Advantages

- Easily explained
- No pre-defined roles, few syntactic assumptions
- Can capture implicit arguments
- Generalizable across domains

Limitations

- Only modeling verbs (for now)
- Not annotating verb senses directly
- Can have multiple equivalent questions

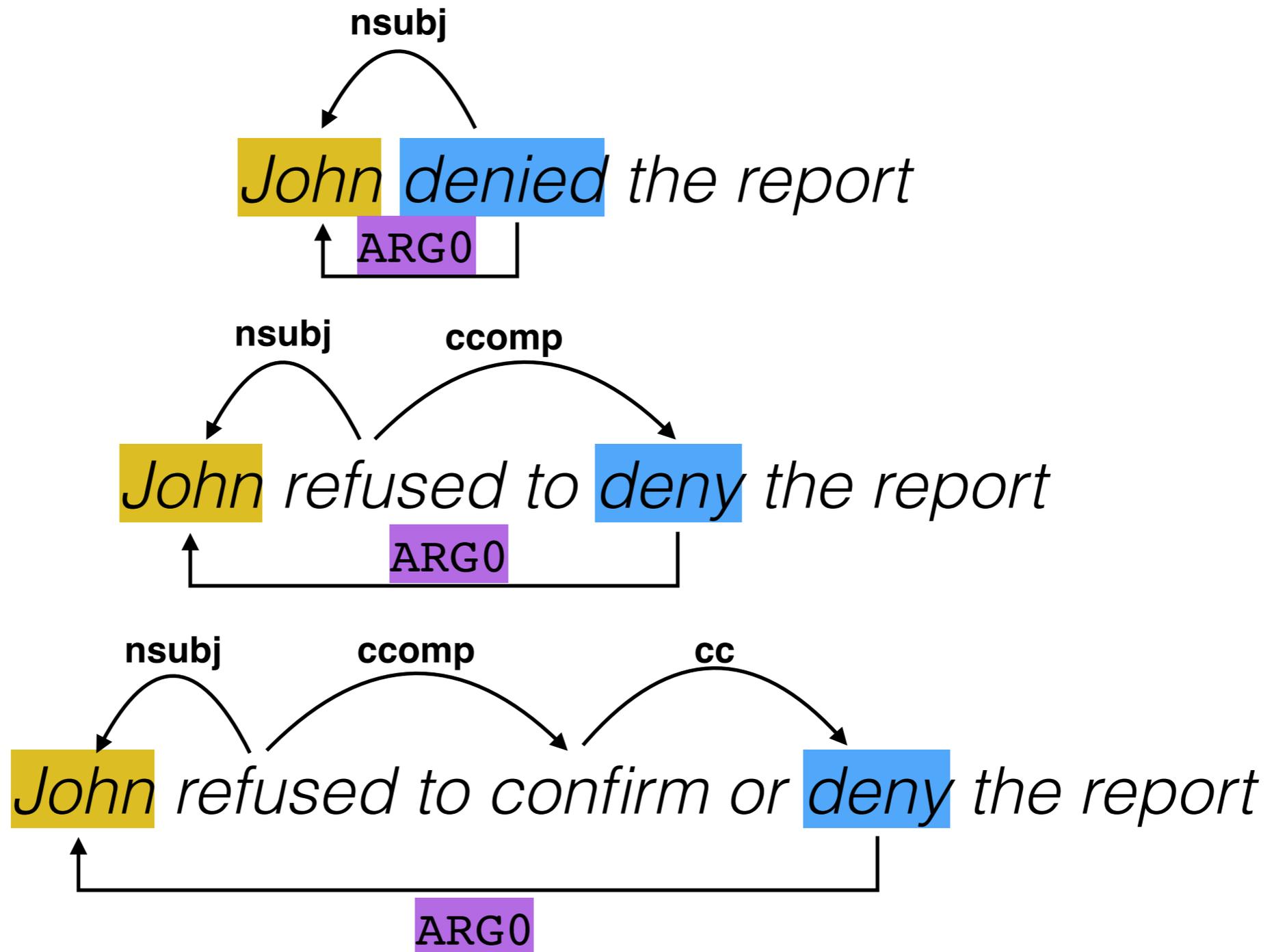
Challenges

- What questions to ask?
- Quality - Can we get good Q/A pairs?
- Coverage - Can we get all the Q/A pairs?

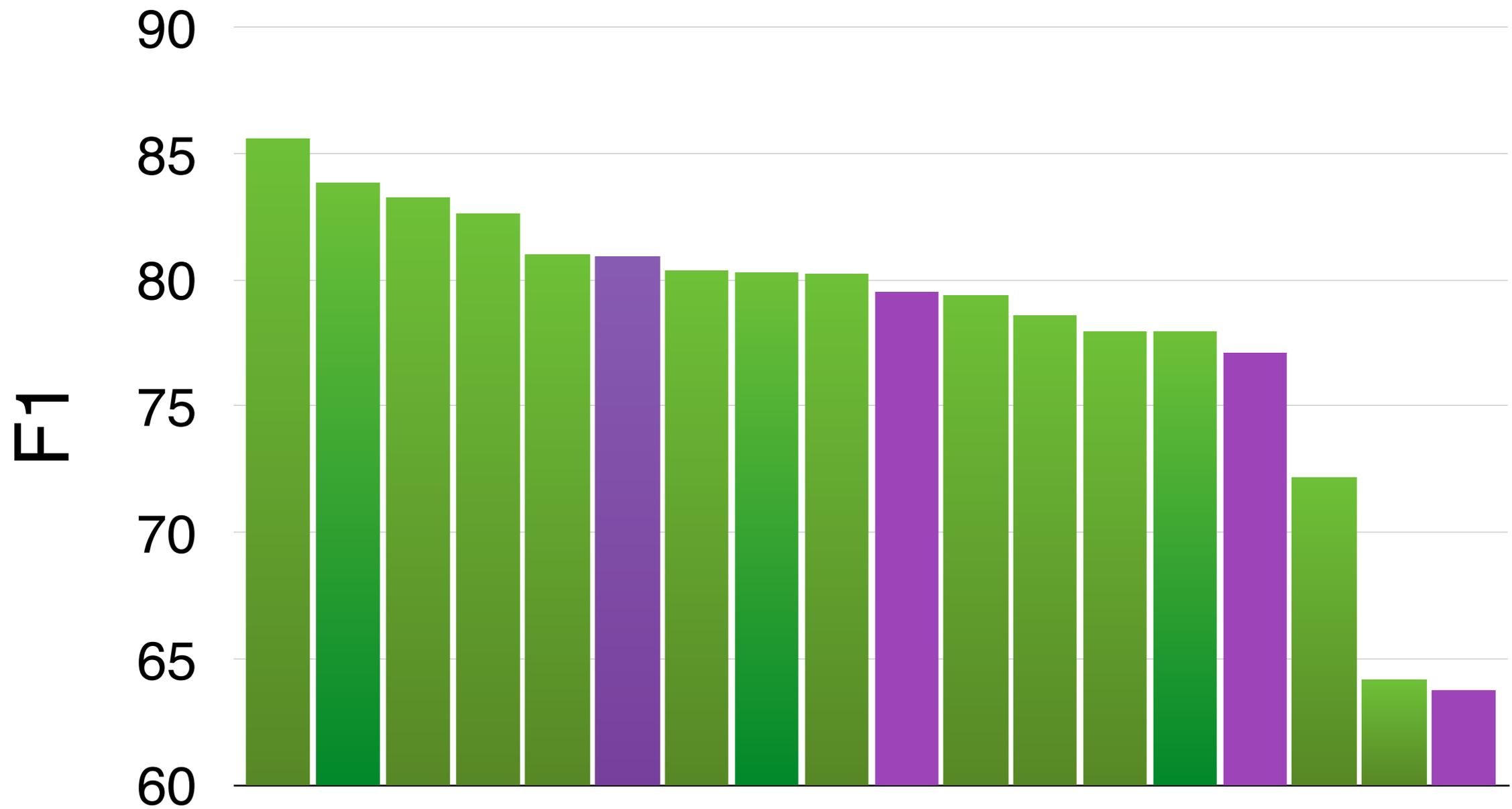
Towards Broad Coverage Semantic Parsing

- Can we crowdsource semantics?
- Train with latent syntax?
- Build fast and accurate parsers?
- Actively select which data to label?

SRL Challenge: Sparsity

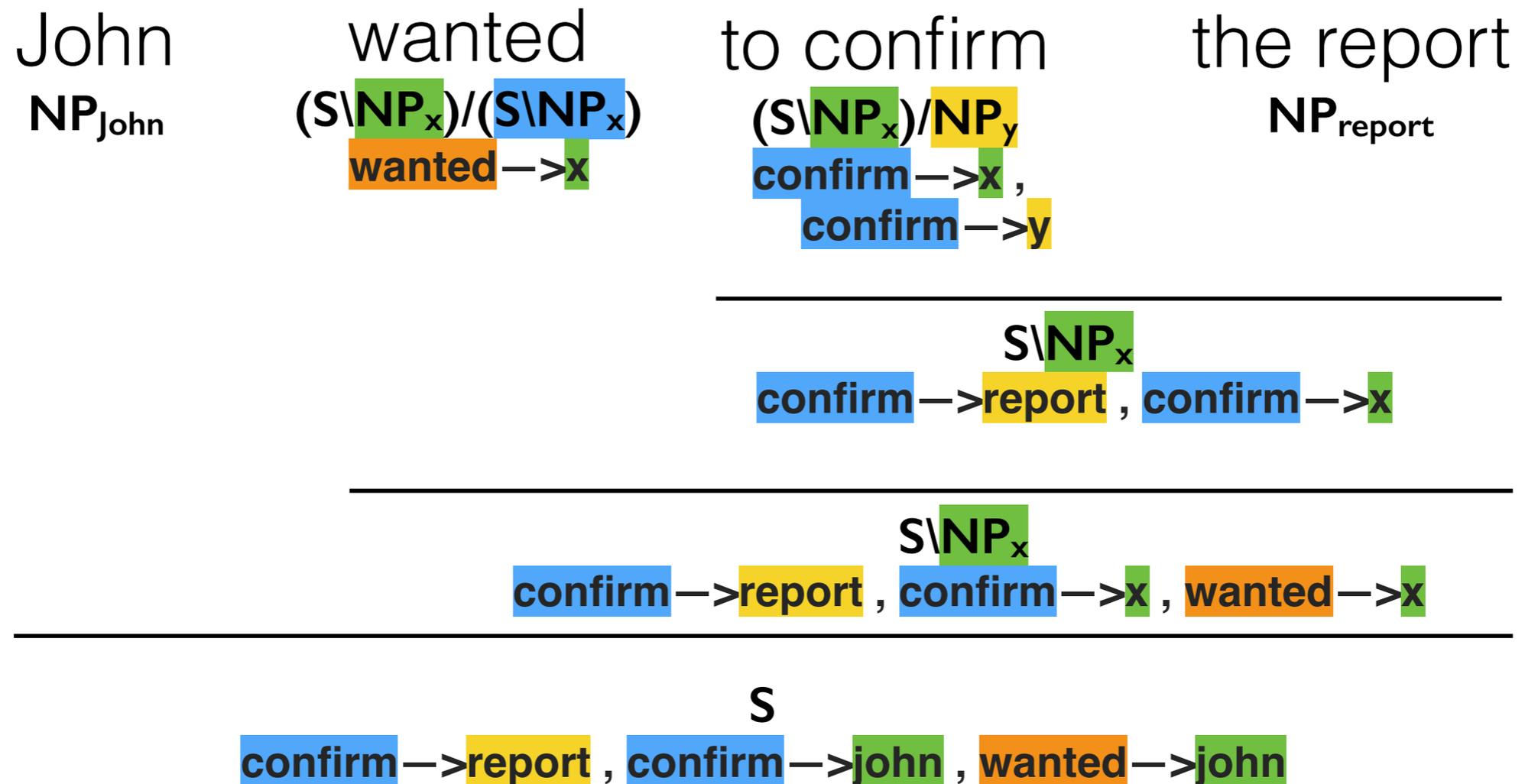


Joint vs. Pipelines



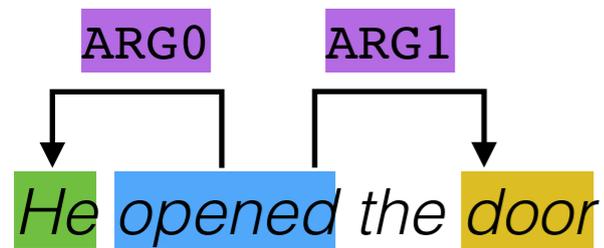
CCG Dependencies

Include nearly all SRL dependencies:



Training

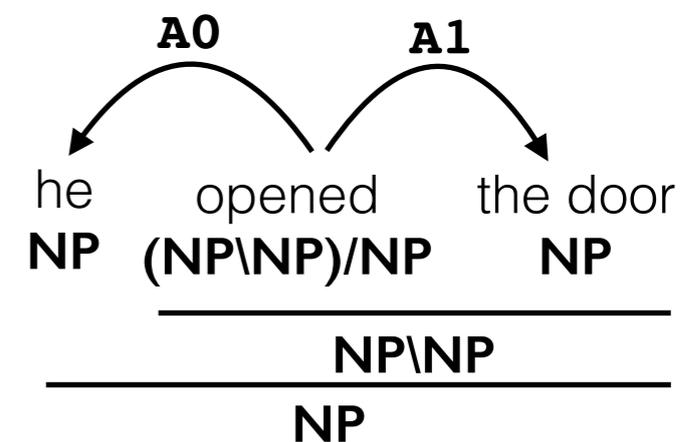
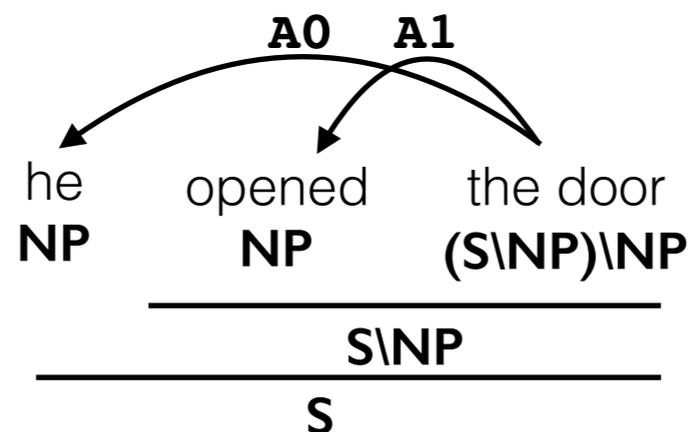
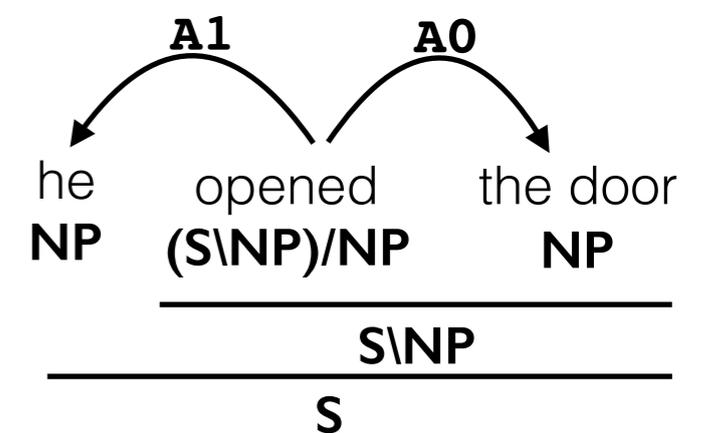
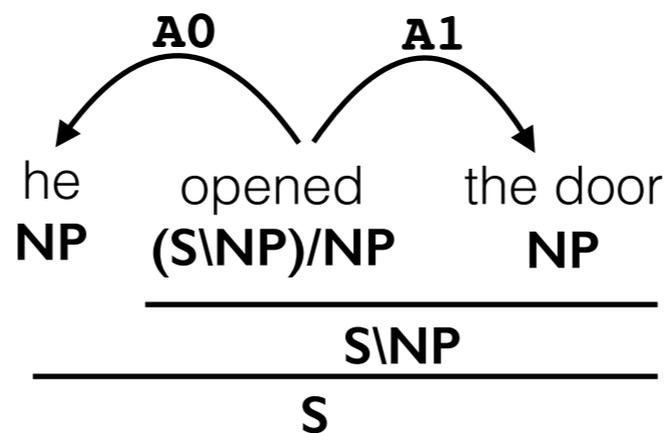
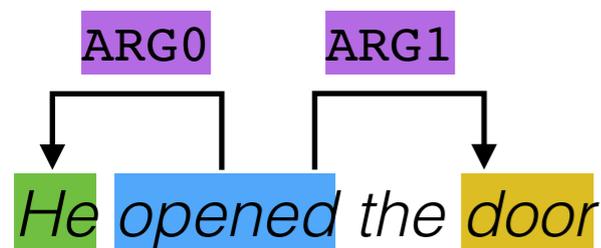
Learn latent CCG that recovers SRL



Training

Learn latent CCG that recovers SRL

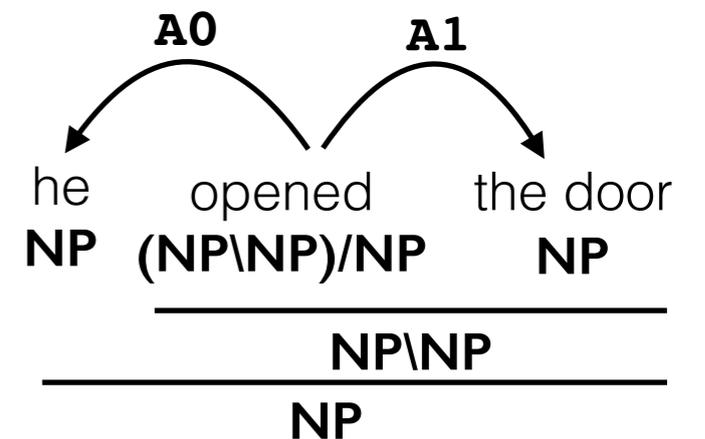
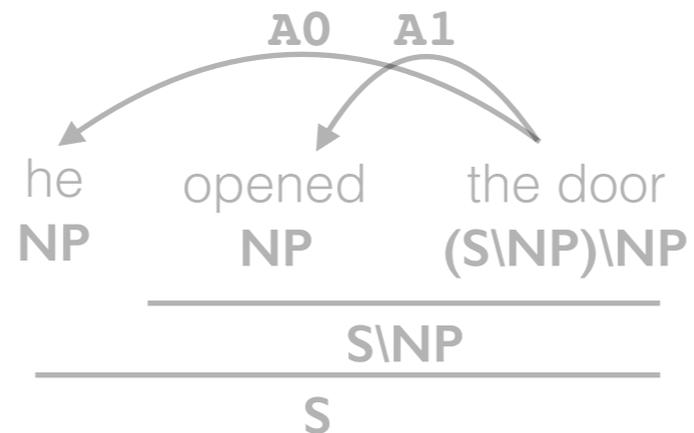
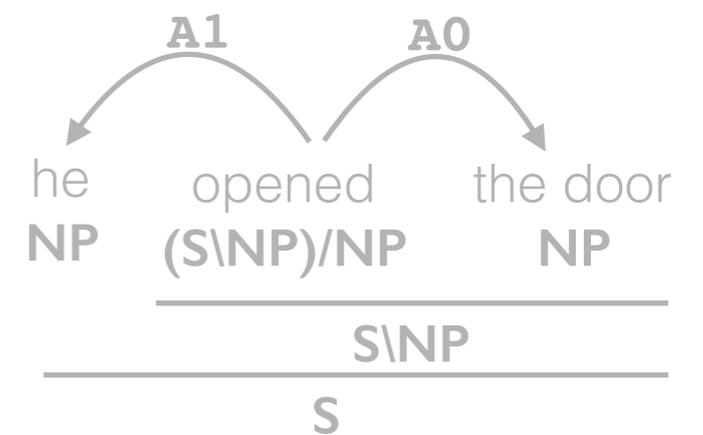
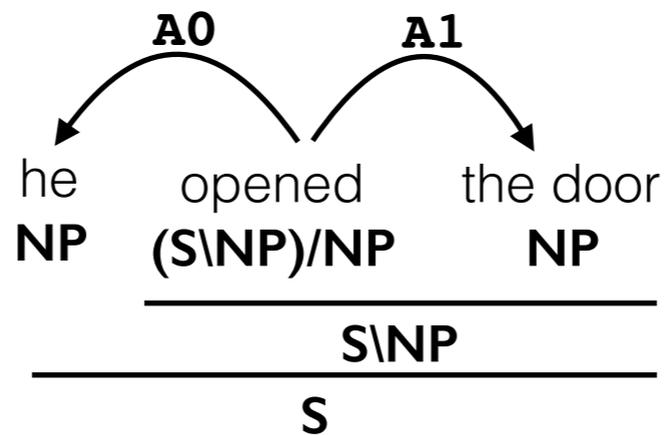
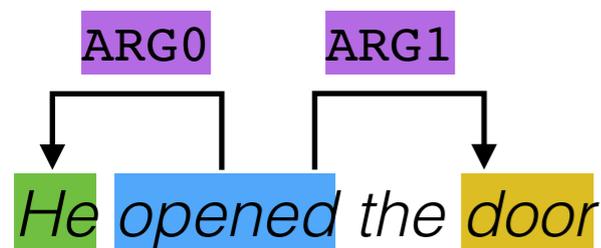
- Generate *consistent* CCG/SRL parses for training sentences



Training

Learn latent CCG that recovers SRL

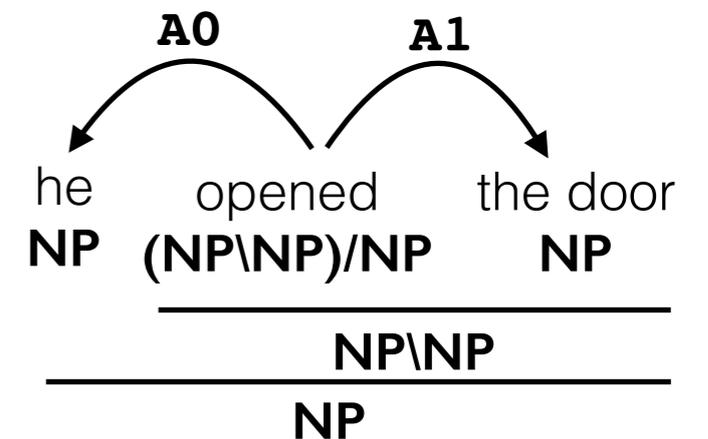
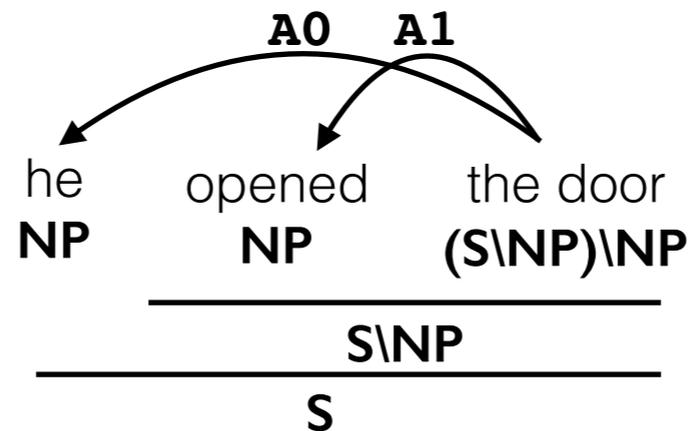
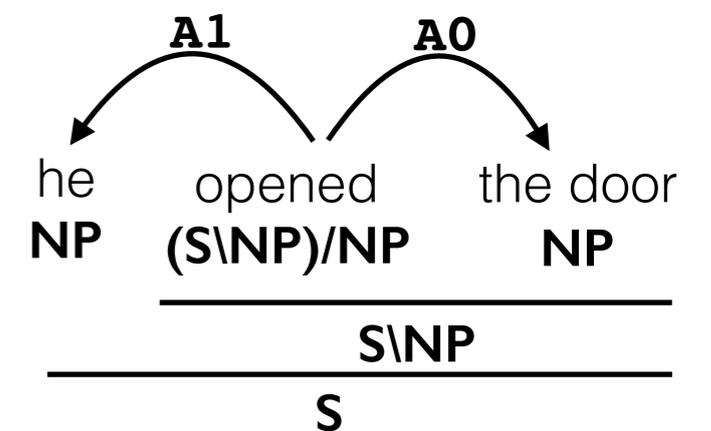
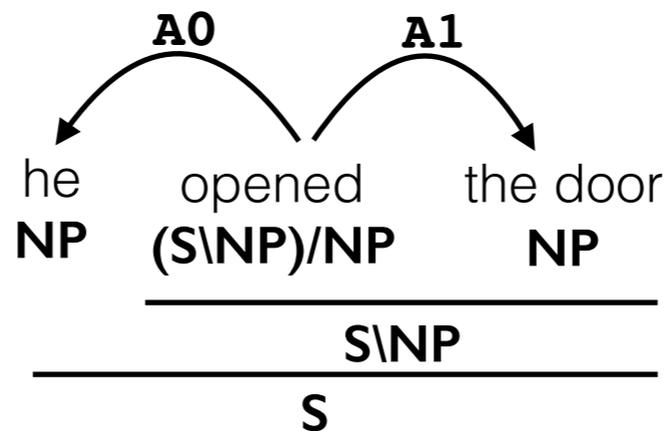
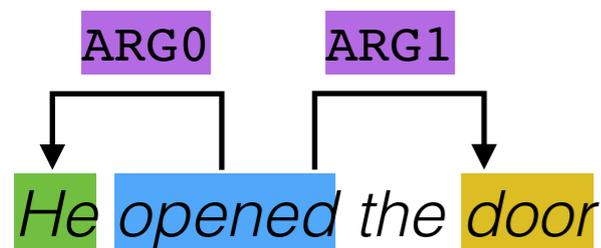
- Mark subset as correct, based on semantic dependencies



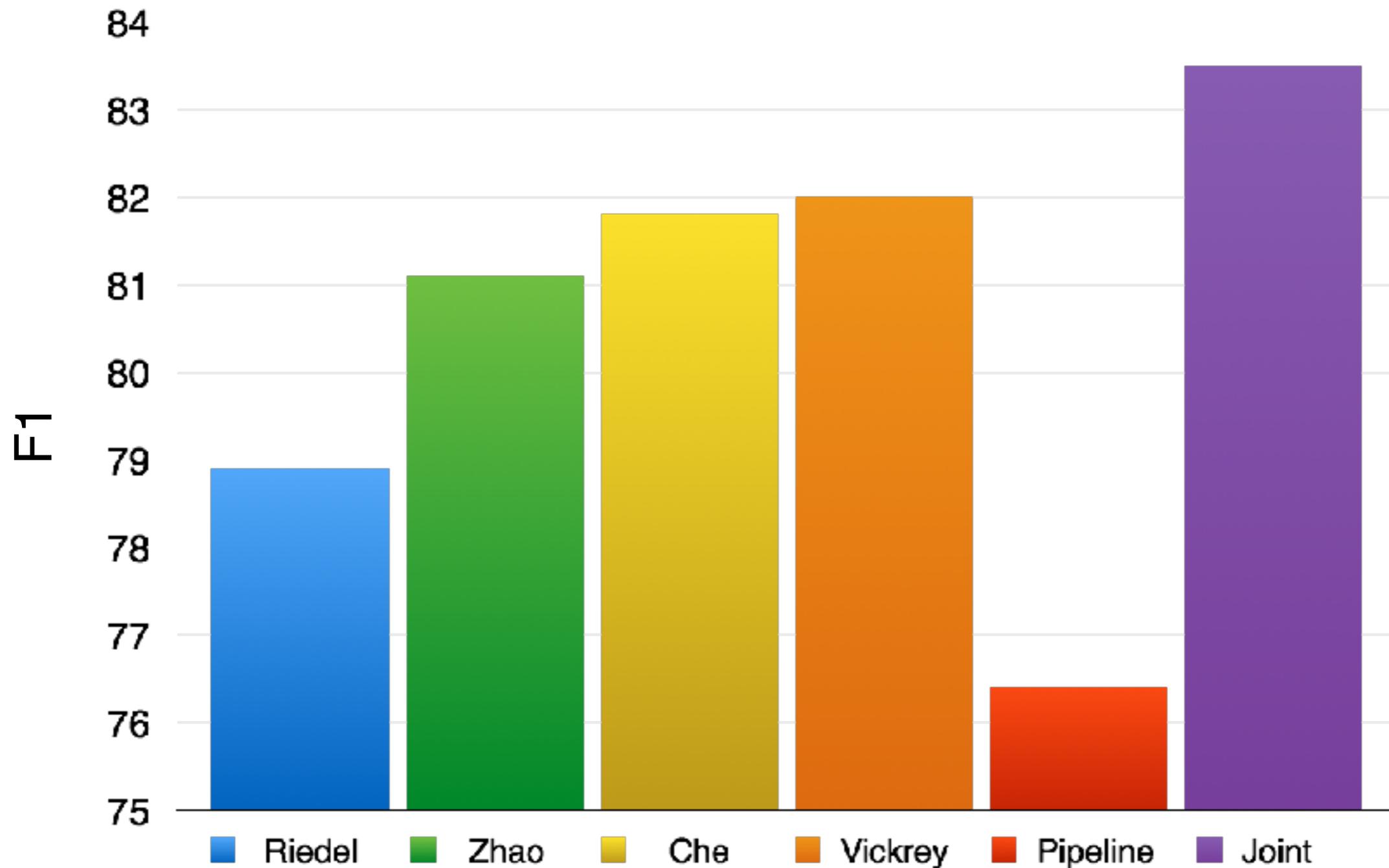
Training

Learn latent CCG that recovers SRL

- Optimize marginal likelihood

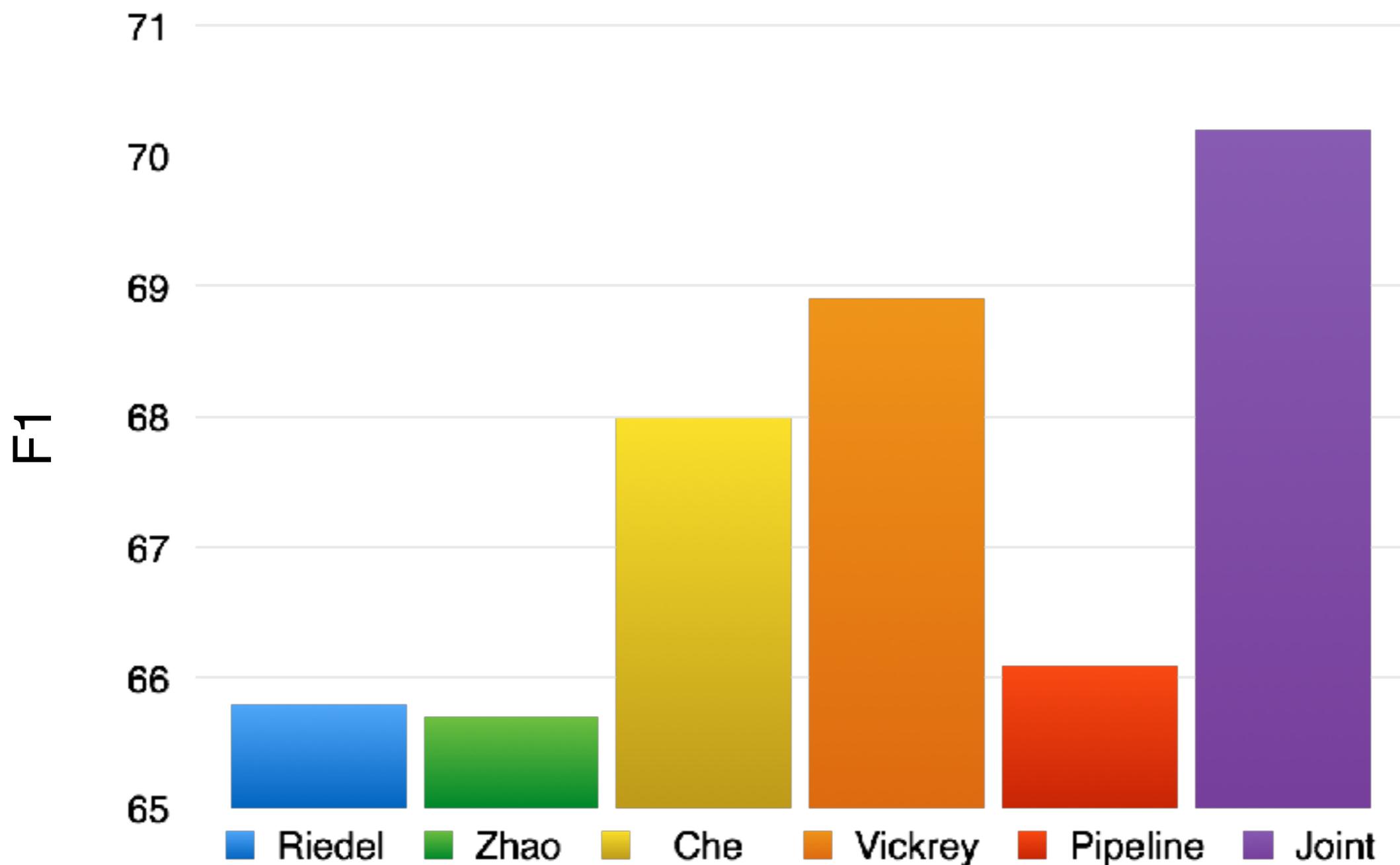


SRL Results



[Lewis et al 2015]

Out-of-domain SRL Results



Towards Broad Coverage Semantic Parsing

- Can we crowdsource semantics?
- Train with latent syntax?
- Build fast and accurate parsers?
- Actively select which data to label?

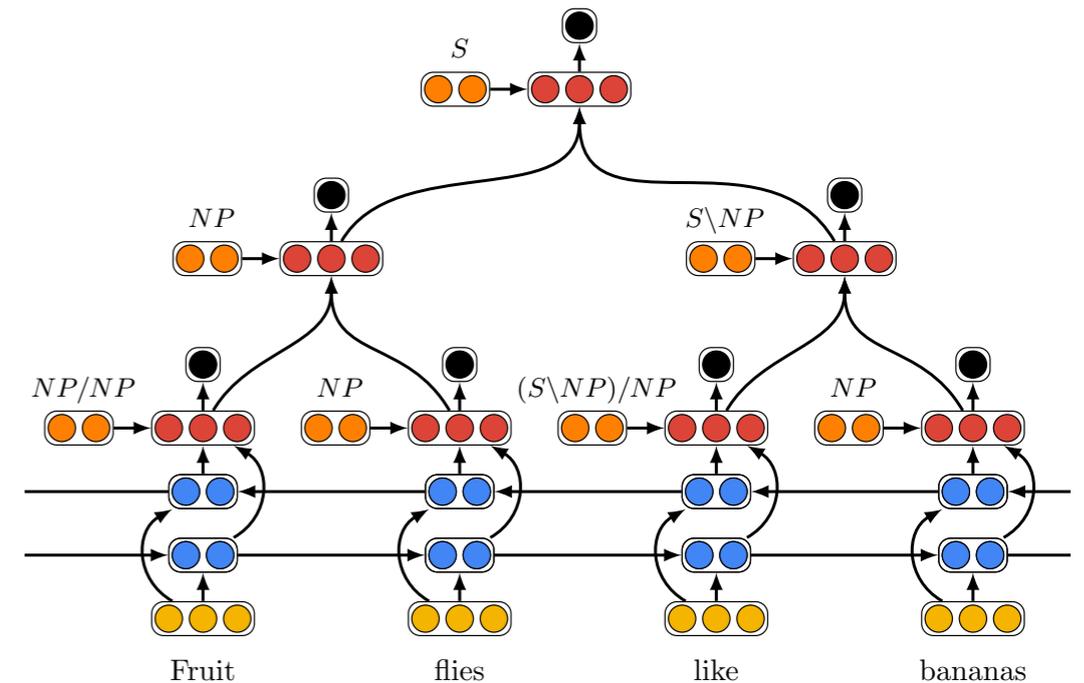
Global A* Parsing

Challenge:

Global models (e.g. Recursive NNs)
break dynamic programs

Our approach:

Combine local and global models in
A* parser



Result:

Accurate models with formal
guarantees

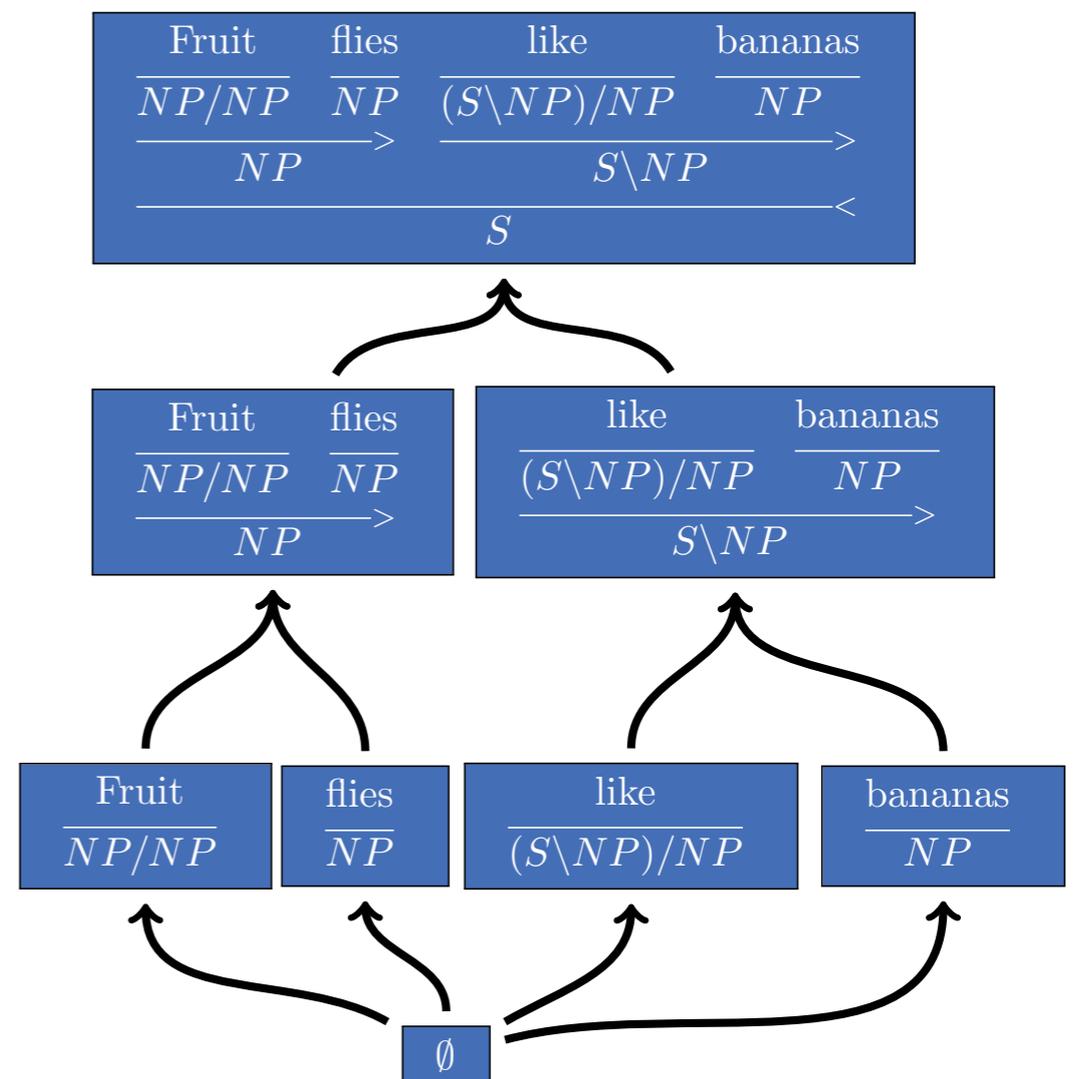
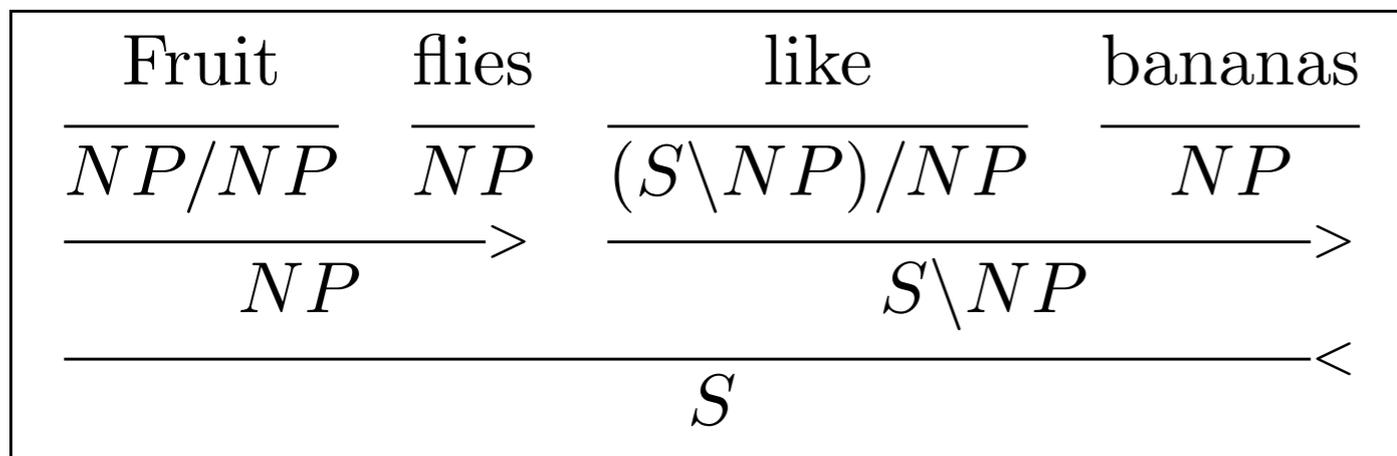
[Lee et al, 2016, EMNLP best paper]

Parsing with Hypergraphs

Input

Fruit flies like bananas

Output

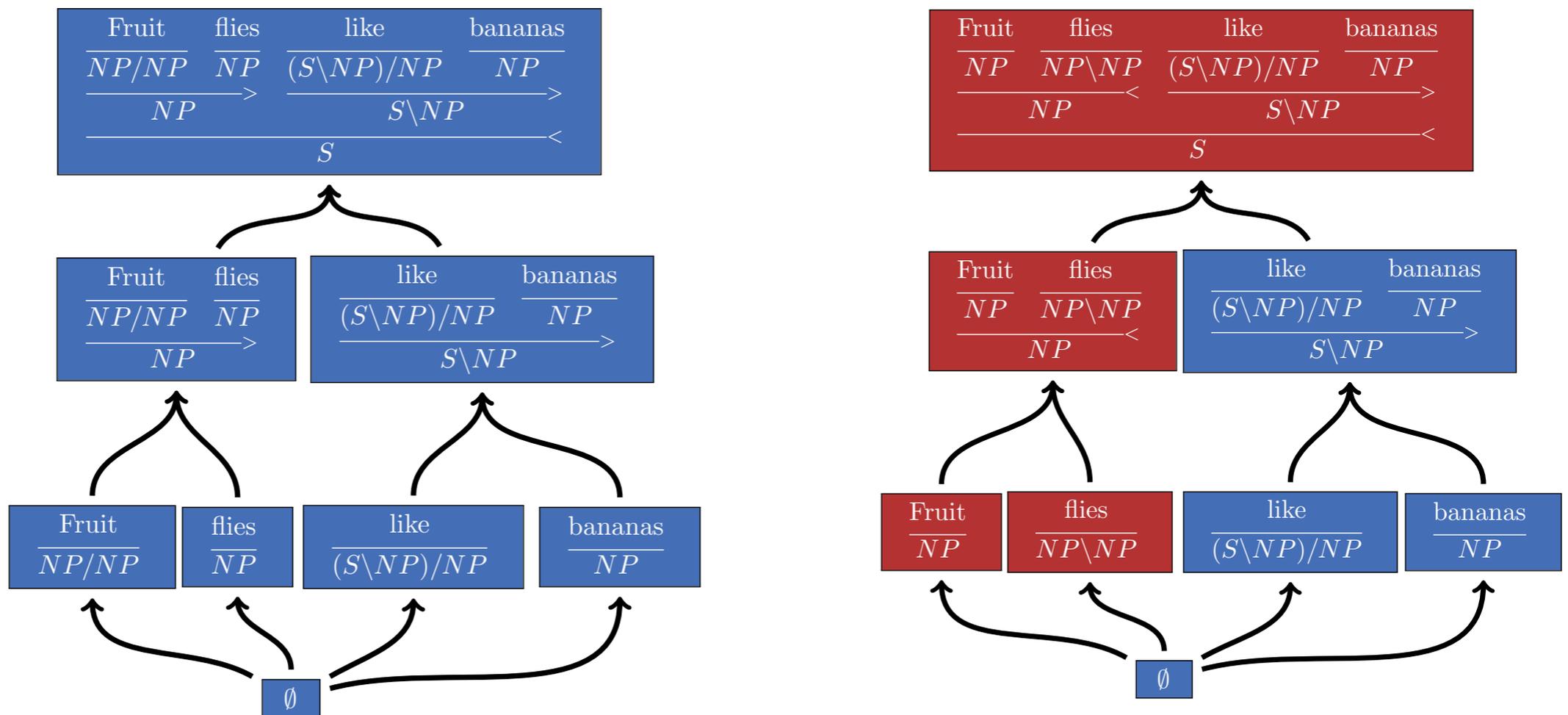


Parsing with Hypergraphs

Input

Fruit flies like bananas

Output

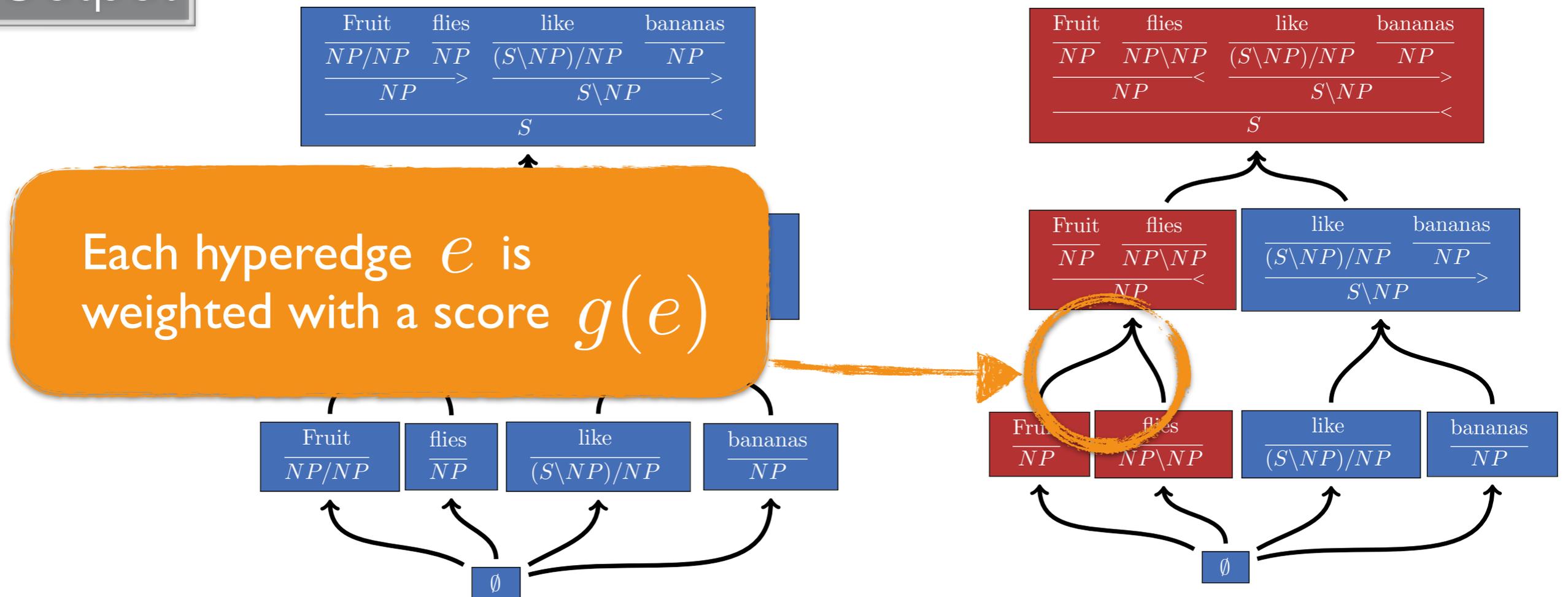


Parsing with Hypergraphs

Input

Fruit flies like bananas

Output

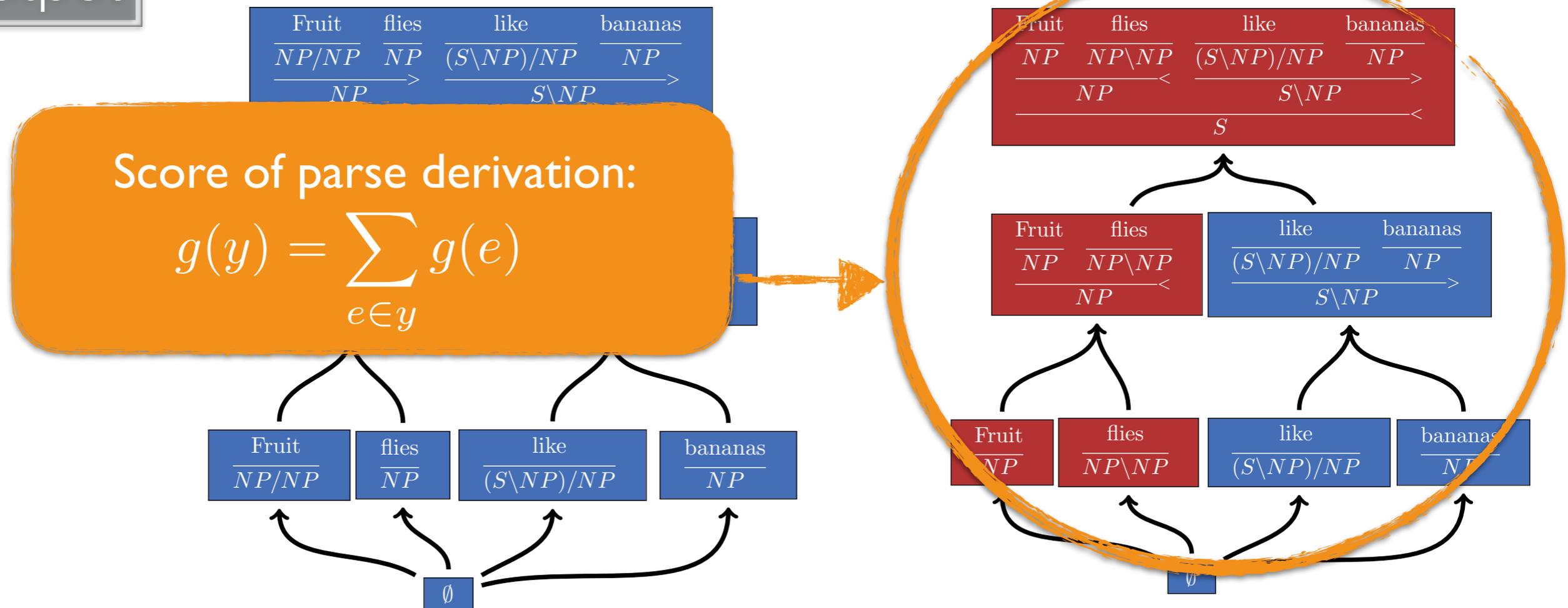


Parsing with Hypergraphs

Input

Fruit flies like bananas

Output



Parsing with Hypergraphs

Fruit	flies	like	bananas
$\overline{NP/NP}$	\overline{NP}	$\overline{(S\backslash NP)/NP}$	\overline{NP}
\overline{NP} >		$\overline{S\backslash NP}$ >	
\overline{S} <			

Fruit	flies	like	bananas
\overline{NP}	$\overline{NP\backslash NP}$	$\overline{(S\backslash NP)/NP}$	\overline{NP}
\overline{NP} <		$\overline{S\backslash NP}$ >	
\overline{S} <			

Fruit	flies	like	bananas
\overline{NP}	$\overline{NP\backslash NP}$	$\overline{(S\backslash NP)/NP}$	\overline{NP}
\overline{NP} <		$\overline{S\backslash NP}$ >	
\overline{S} <			

Fruit	flies	like	bananas
$\overline{NP/NP}$	\overline{NP}	$\overline{(S\backslash NP)/NP}$	\overline{NP}
\overline{NP} >		$\overline{S\backslash NP}$ >	
\overline{S} <			

Fruit	flies
\overline{NP}	$\overline{NP\backslash NP}$
\overline{NP} <	

Fruit	flies
$\overline{NP/NP}$	\overline{NP}
\overline{NP} >	

like	bananas
$\overline{(S\backslash NP)/NP}$	\overline{NP}
$\overline{S\backslash NP}$ >	

flies
$\overline{NP\backslash NP}$

Fruit
$\overline{NP/NP}$

flies
\overline{NP}

like
$\overline{(S\backslash NP)/NP}$

bananas
\overline{NP}

∅

Fruit
\overline{NP}

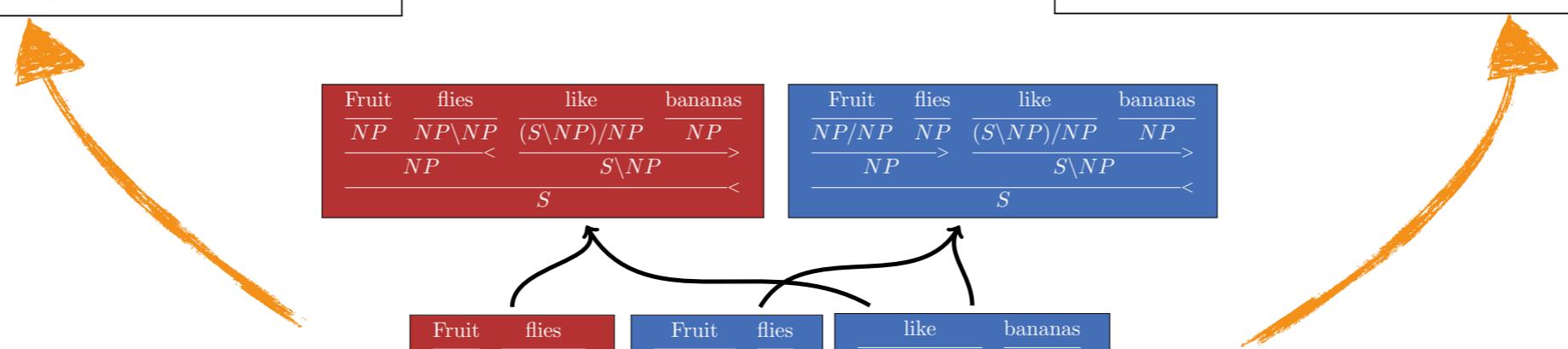
flies
$\overline{S\backslash NP}$

like
$\overline{(S\backslash S)/NP}$

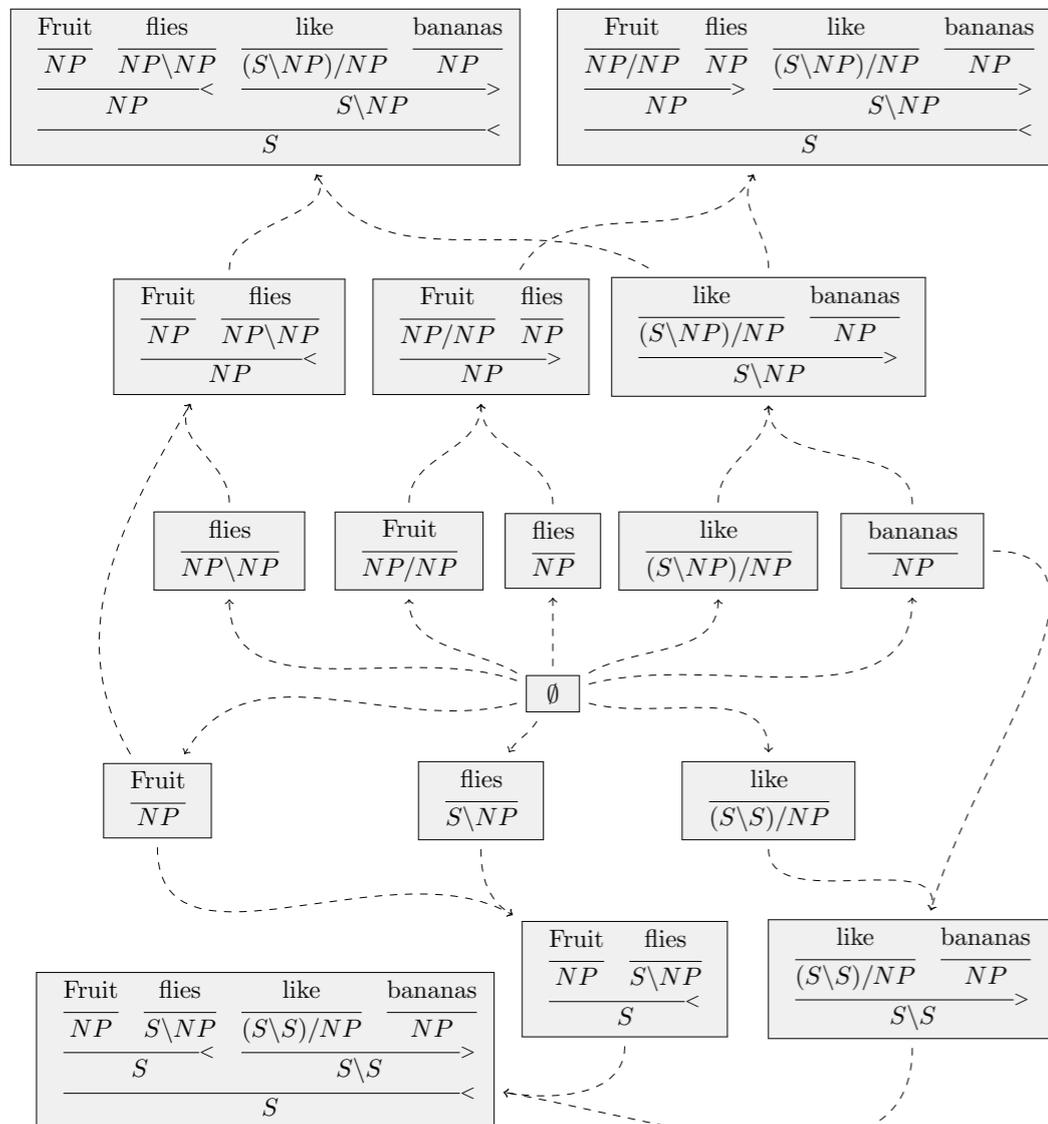
Fruit	flies	like	bananas
\overline{NP}	$\overline{S\backslash NP}$	$\overline{(S\backslash S)/NP}$	\overline{NP}
\overline{S} <		$\overline{S\backslash S}$ >	
\overline{S} <			

Fruit	flies
\overline{NP}	$\overline{S\backslash NP}$
\overline{S} <	

like	bananas
$\overline{(S\backslash S)/NP}$	\overline{NP}
$\overline{S\backslash S}$ >	



Parsing with Hypergraphs

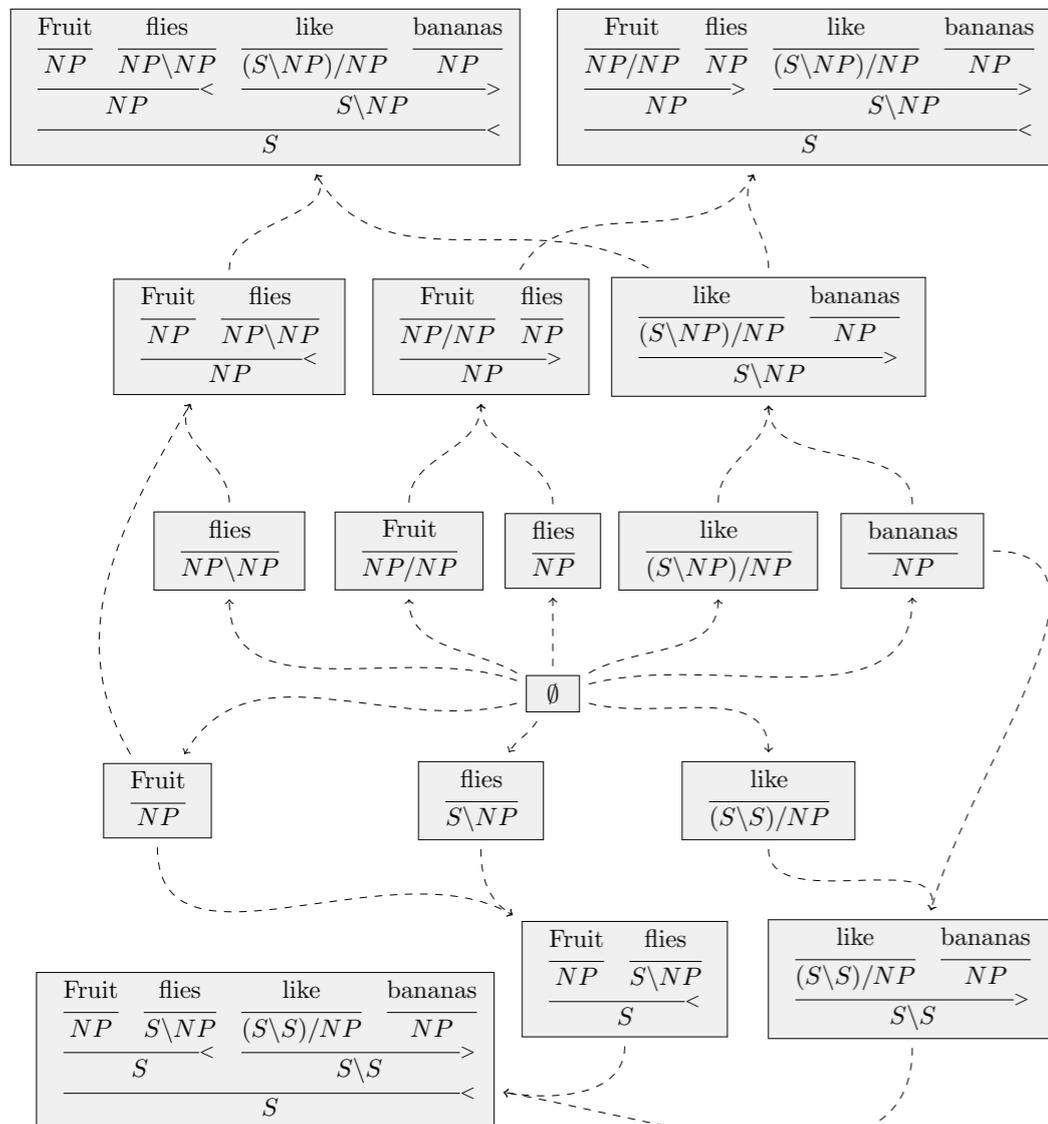


❖ Predicted parse: $y^* = \operatorname{argmax}_{y \in Y} g(y)$

❖ Exponential number of nodes

→ Intractable inference

Managing Intractable Search Spaces



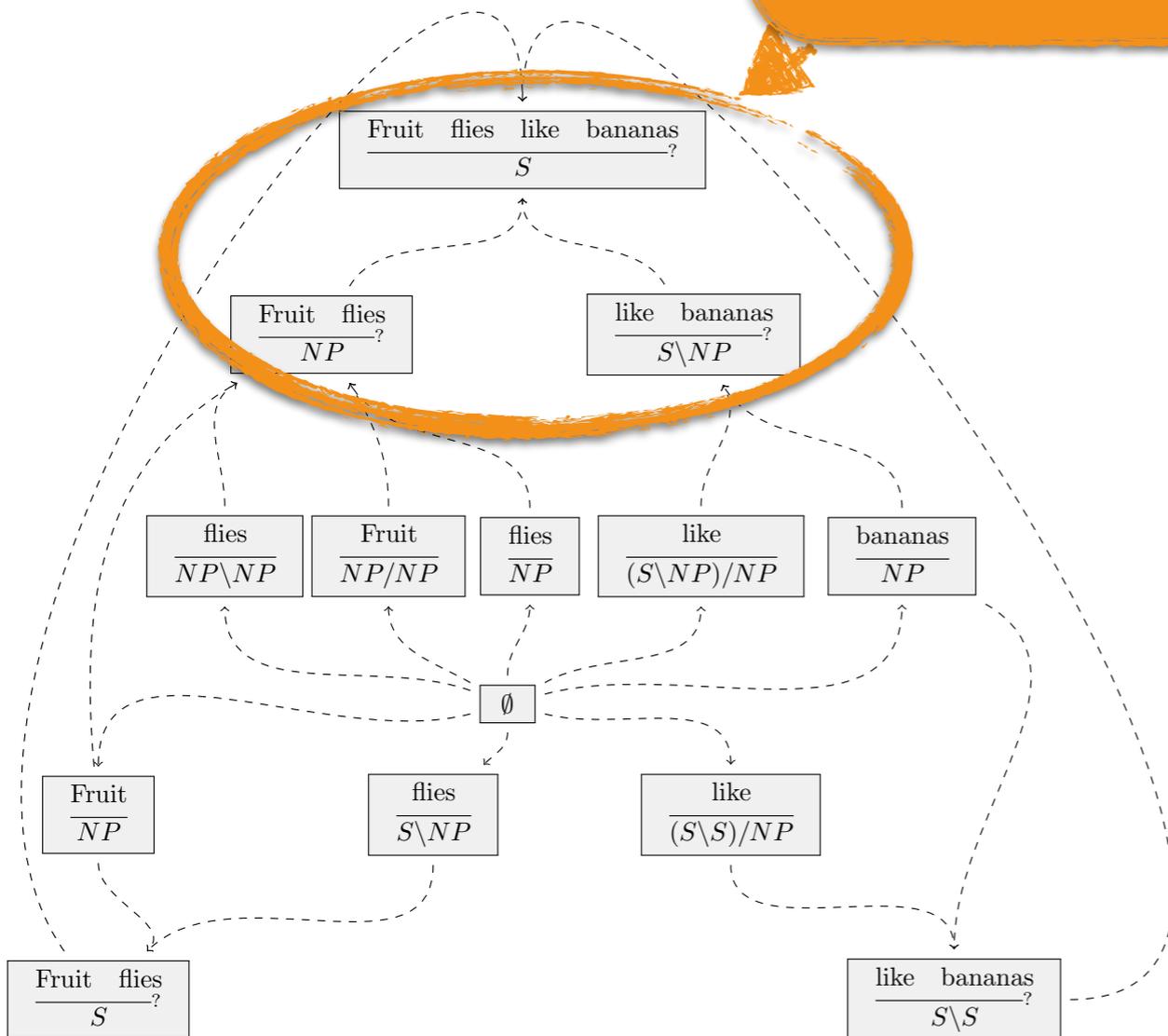
Approximate inference with global expressivity, e.g.

- ❖ Greedy / beam search:
 - ❖ Nivre, 2008
 - ❖ Chen and Manning, 2014
 - ❖ Andor et al., 2016

- ❖ Reranking:
 - ❖ Charniak and Johnson, 2005
 - ❖ Huang, 2008
 - ❖ Socher et al., 2013

Locally Factored Parsing

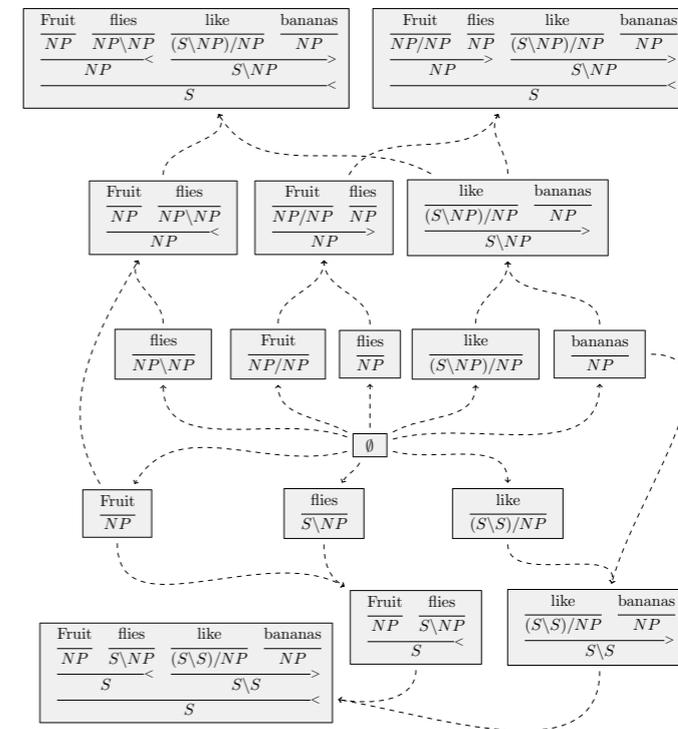
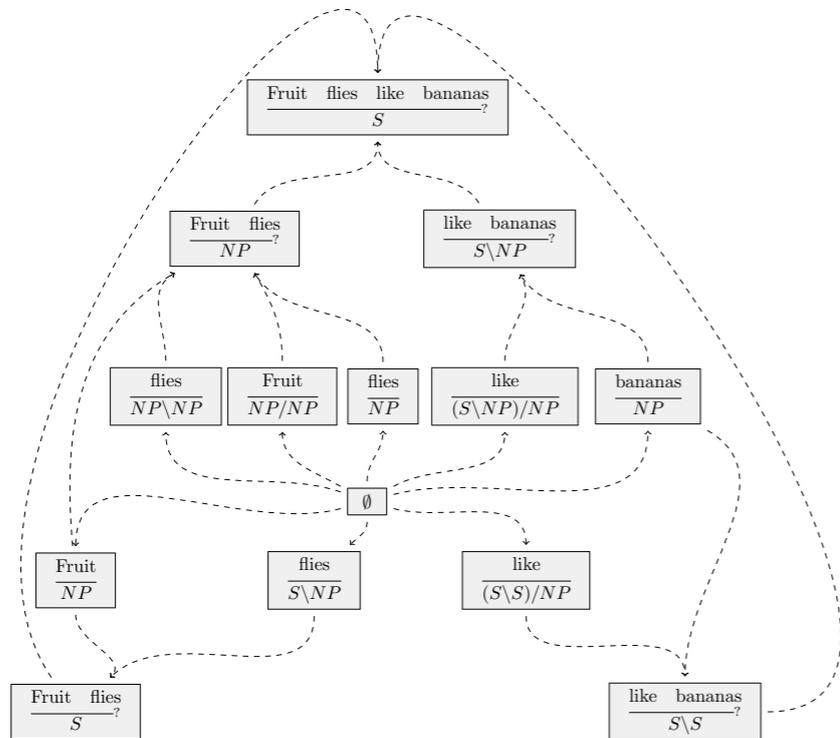
Scores condition on
local structures



Dynamic programs with locally factored models, e.g.

- ❖ CKY:
 - ❖ Collins, 1997
 - ❖ Durrett and Klein, 2015
- ❖ Minimum spanning tree:
 - ❖ McDonald et al., 2005
 - ❖ Kiperwasser and Goldberg, 2016

Local vs. Global Models



Local model:

$$y^* = \underset{y \in Y}{\operatorname{argmax}} (g_{\text{local}}(y))$$

Efficient

Inexpressive

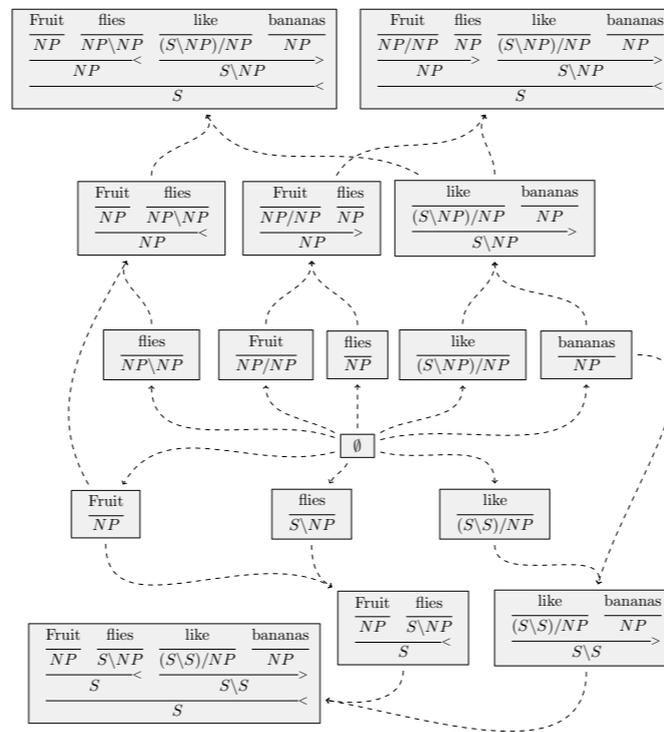
Global model:

$$y^* = \underset{y \in Y}{\operatorname{argmax}} (g_{\text{global}}(y))$$

Intractable

Expressive

This Work



Combined model:

$$y^* = \underset{y \in Y}{\operatorname{argmax}} (g_{\text{local}}(y) + g_{\text{global}}(y))$$

Efficient

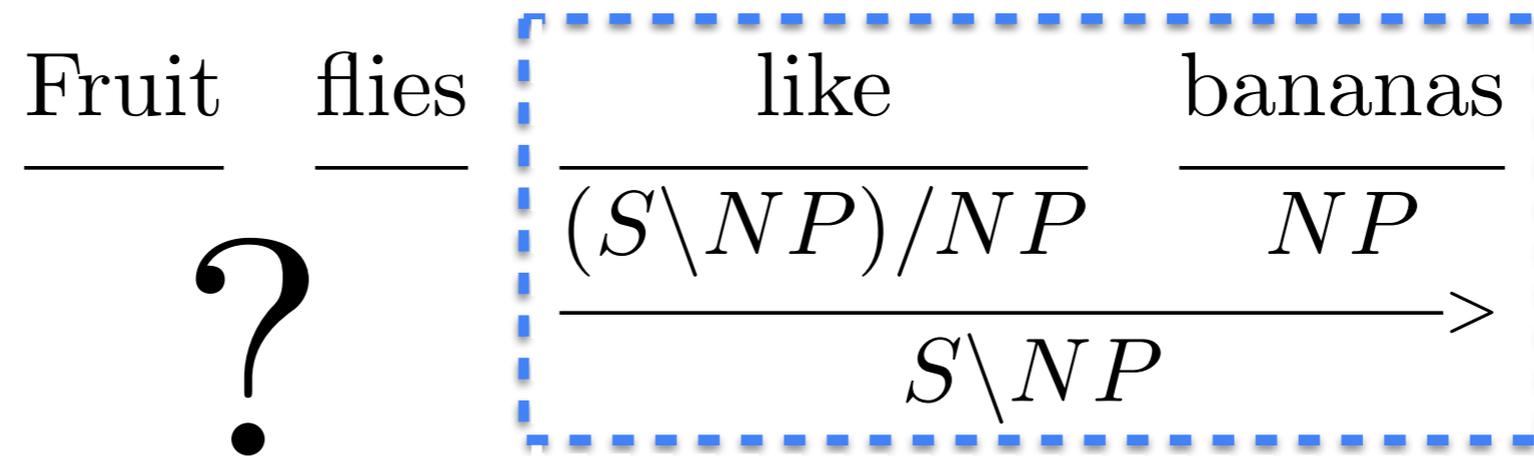
Expressive

A* Parsing

$$y^* = \operatorname{argmax}_{y \in Y} g(y)$$

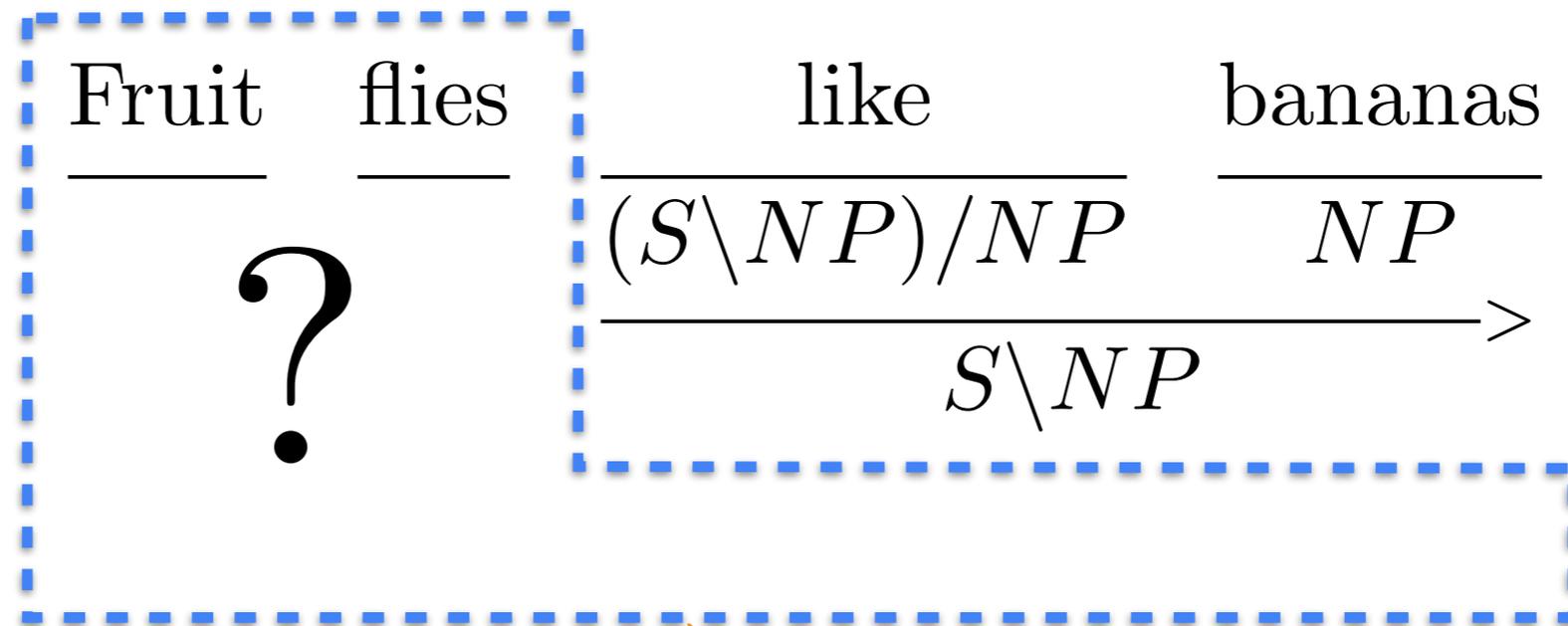
- ❖ Search in the space of partial parses
- ❖ First explored full parse **guaranteed to be optimal**

A* Parsing



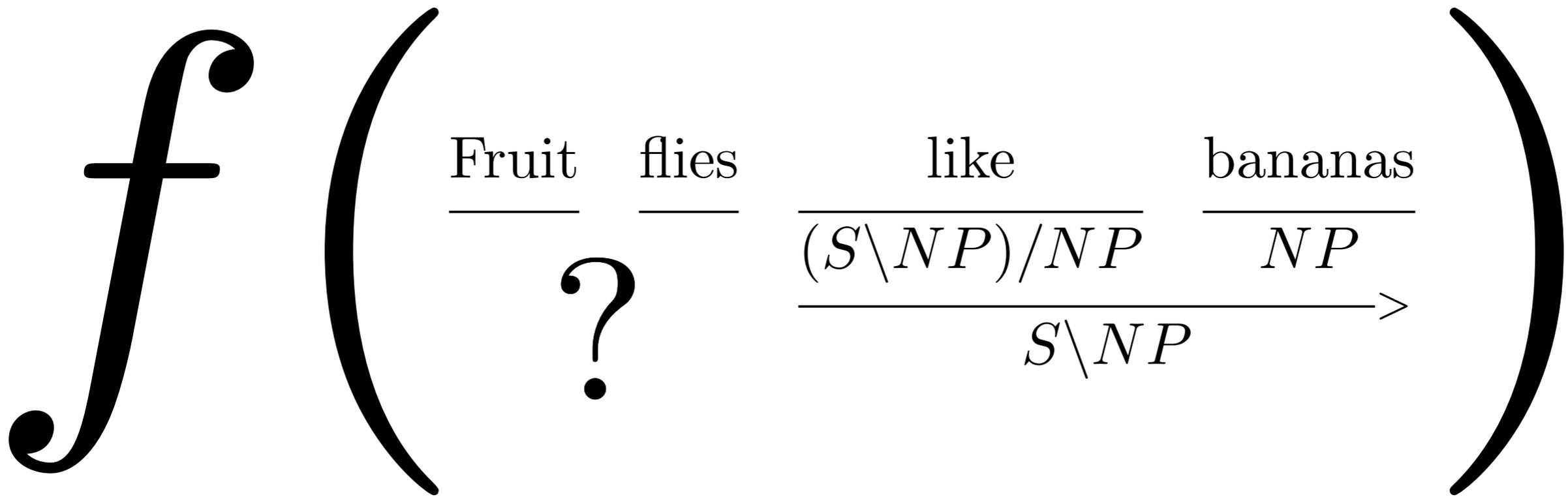
Partial parse

A* Parsing



Partial parse

A* Parsing



Exploration priority

Partial parse

A* Parsing

Exploration priority

$$f\left(\begin{array}{c} \text{Fruit} \quad \text{flies} \quad \text{like} \quad \text{bananas} \\ \frac{\quad}{?} \quad \frac{\quad}{(S \setminus NP) / NP} \quad \frac{\quad}{NP} \\ \hline S \setminus NP \end{array}\right) = g\left(\begin{array}{c} \text{Fruit} \quad \text{flies} \quad \text{like} \quad \text{bananas} \\ \frac{\quad}{?} \quad \frac{\quad}{(S \setminus NP) / NP} \quad \frac{\quad}{NP} \\ \hline S \setminus NP \end{array}\right) + h\left(\begin{array}{c} \text{Fruit} \quad \text{flies} \quad \text{like} \quad \text{bananas} \\ \frac{\quad}{?} \quad \frac{\quad}{(S \setminus NP) / NP} \quad \frac{\quad}{NP} \\ \hline S \setminus NP \end{array}\right)$$

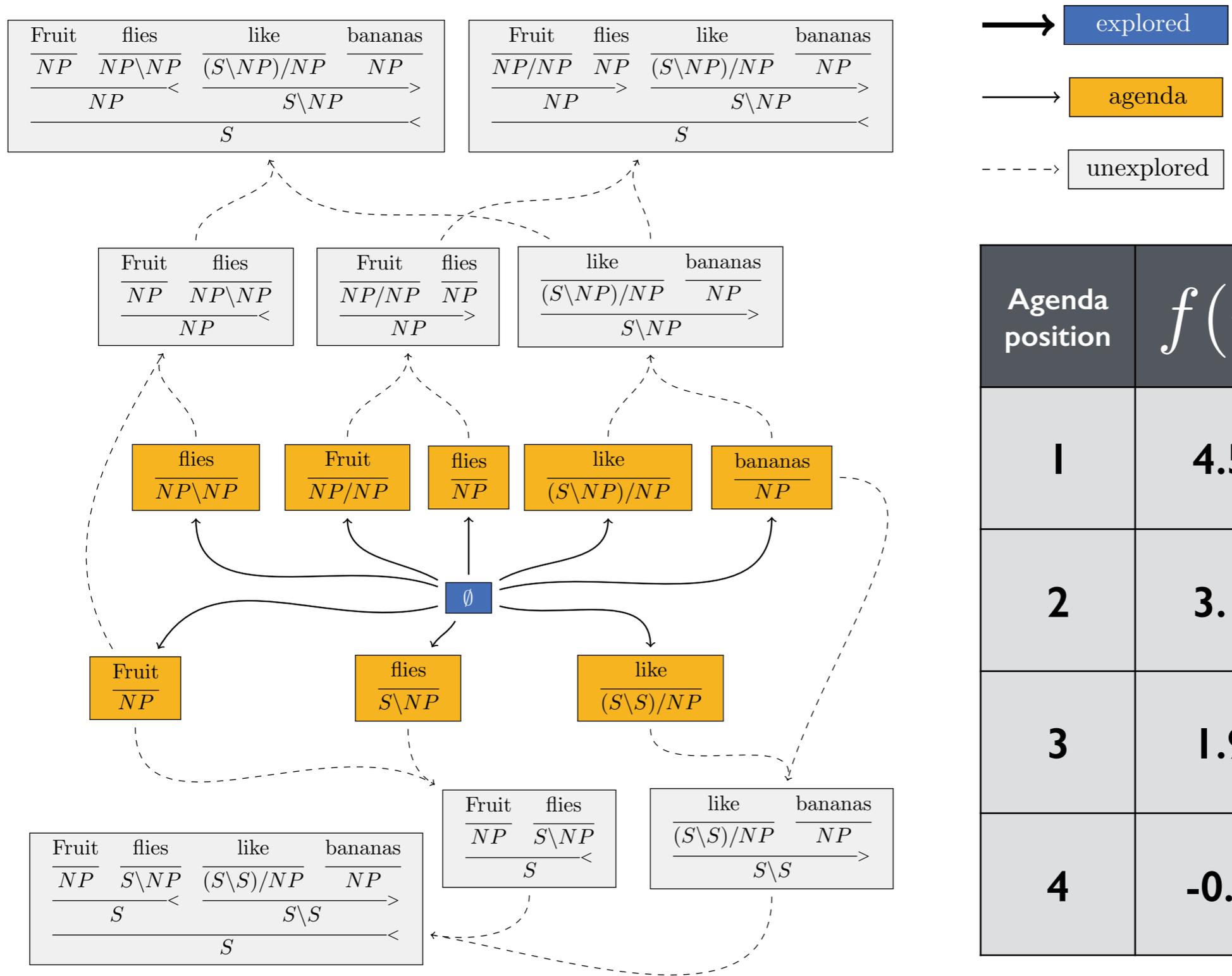
Inside score

$$\begin{array}{c} \text{Fruit} \quad \text{flies} \quad \text{like} \quad \text{bananas} \\ \frac{\quad}{?} \quad \frac{\quad}{(S \setminus NP) / NP} \quad \frac{\quad}{NP} \\ \hline S \setminus NP \end{array}$$

Admissible A* heuristic

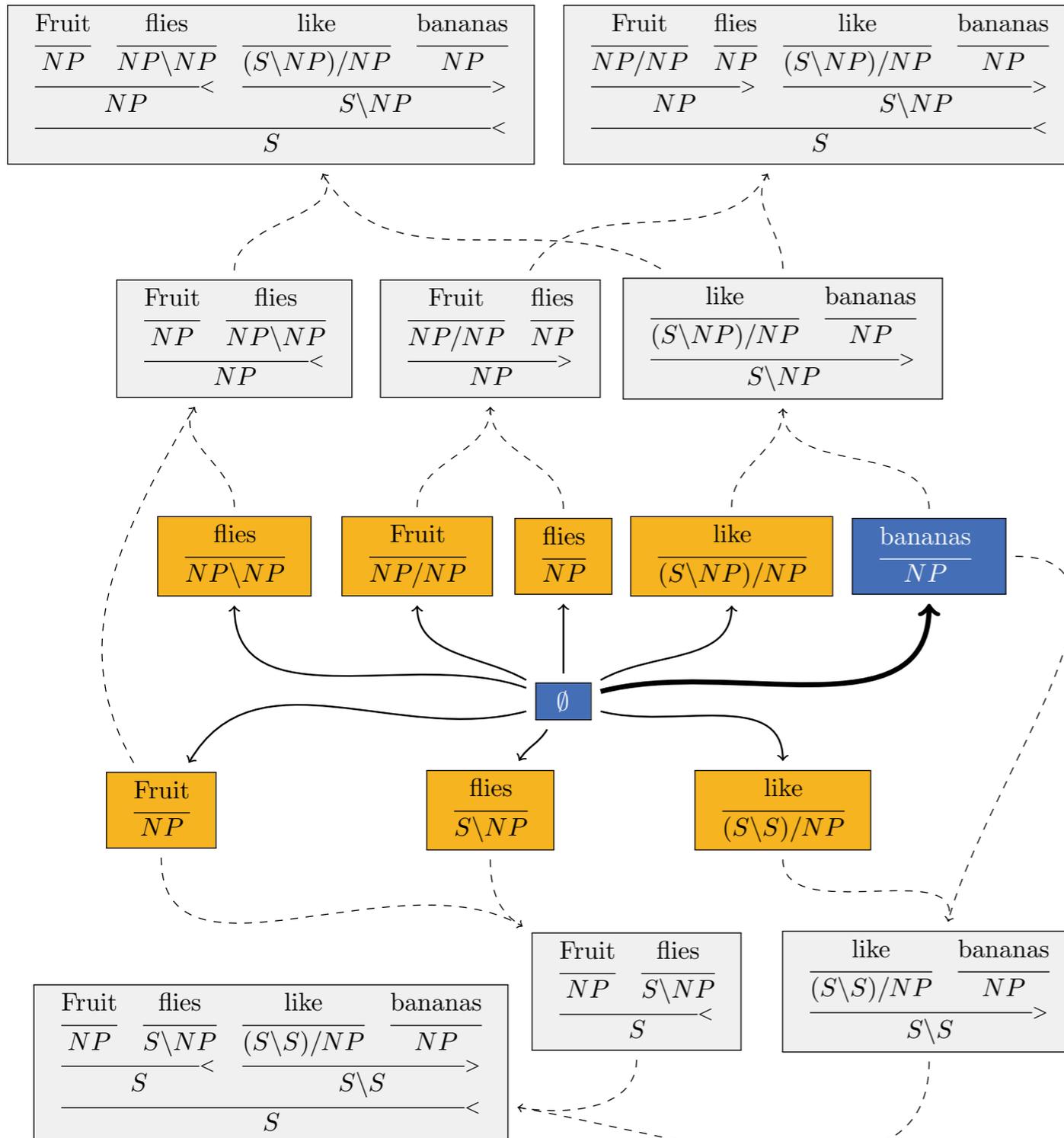
$$\begin{array}{c} \text{Fruit} \quad \text{flies} \quad \text{like} \quad \text{bananas} \\ \frac{\quad}{?} \quad \frac{\quad}{(S \setminus NP) / NP} \quad \frac{\quad}{NP} \\ \hline S \setminus NP \end{array}$$

A* Parsing



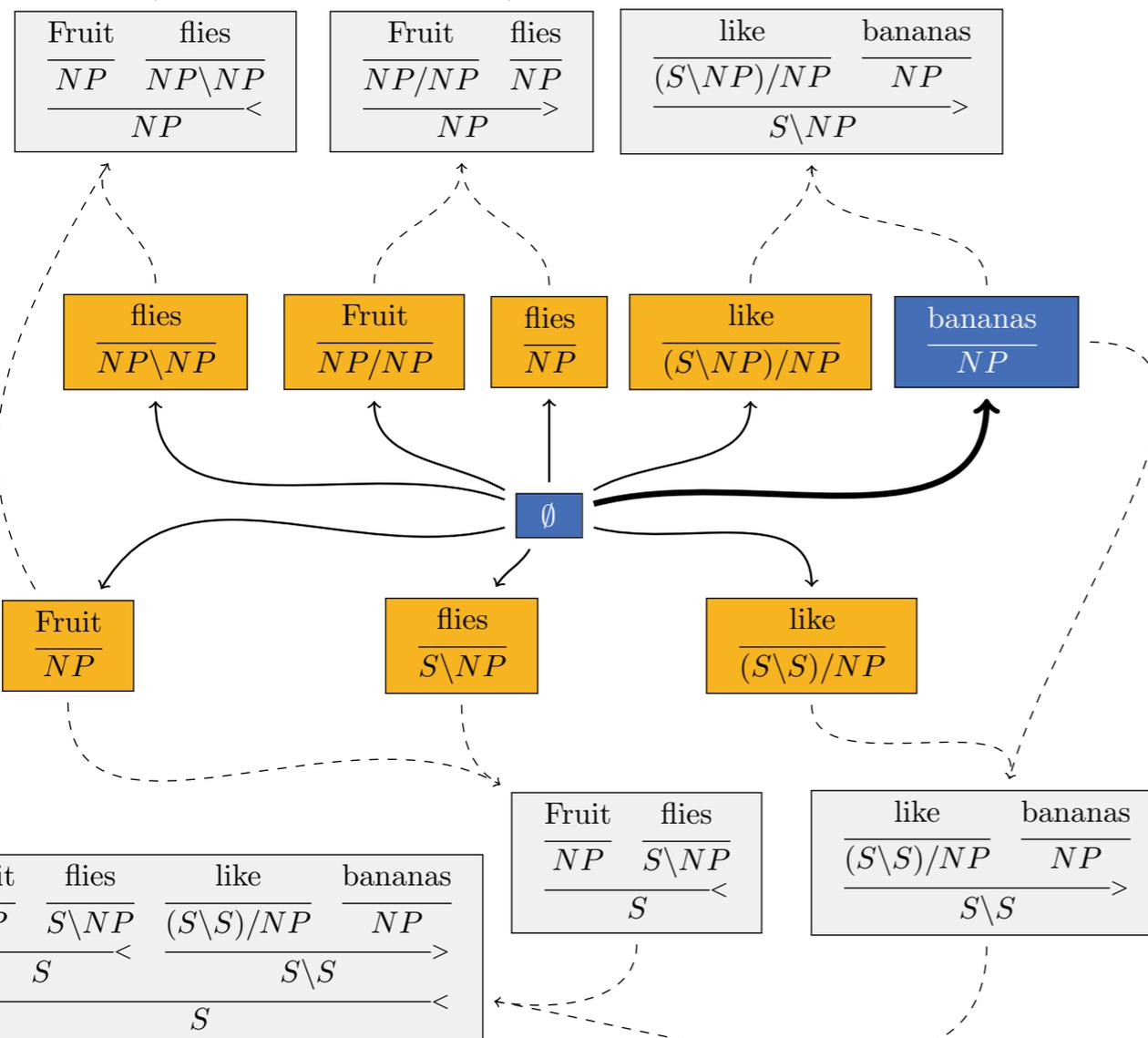
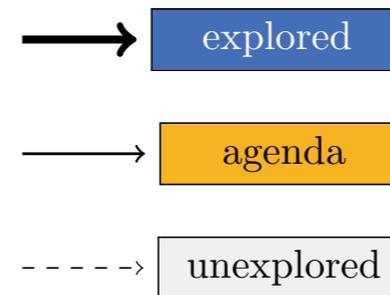
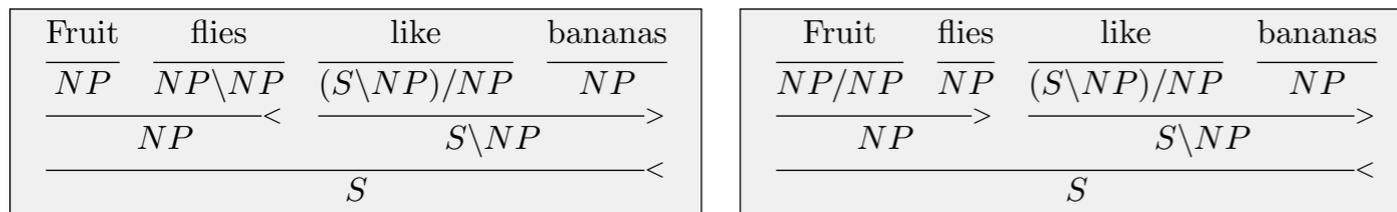
Agenda position	$f(y)$	y
1	4.5	$\frac{\text{bananas}}{NP}$
2	3.1	$\frac{\text{like}}{(S \backslash NP) / NP}$
3	1.9	$\frac{\text{Fruit}}{NP}$
4	-0.5	$\frac{\text{Fruit}}{NP / NP}$

A* Parsing



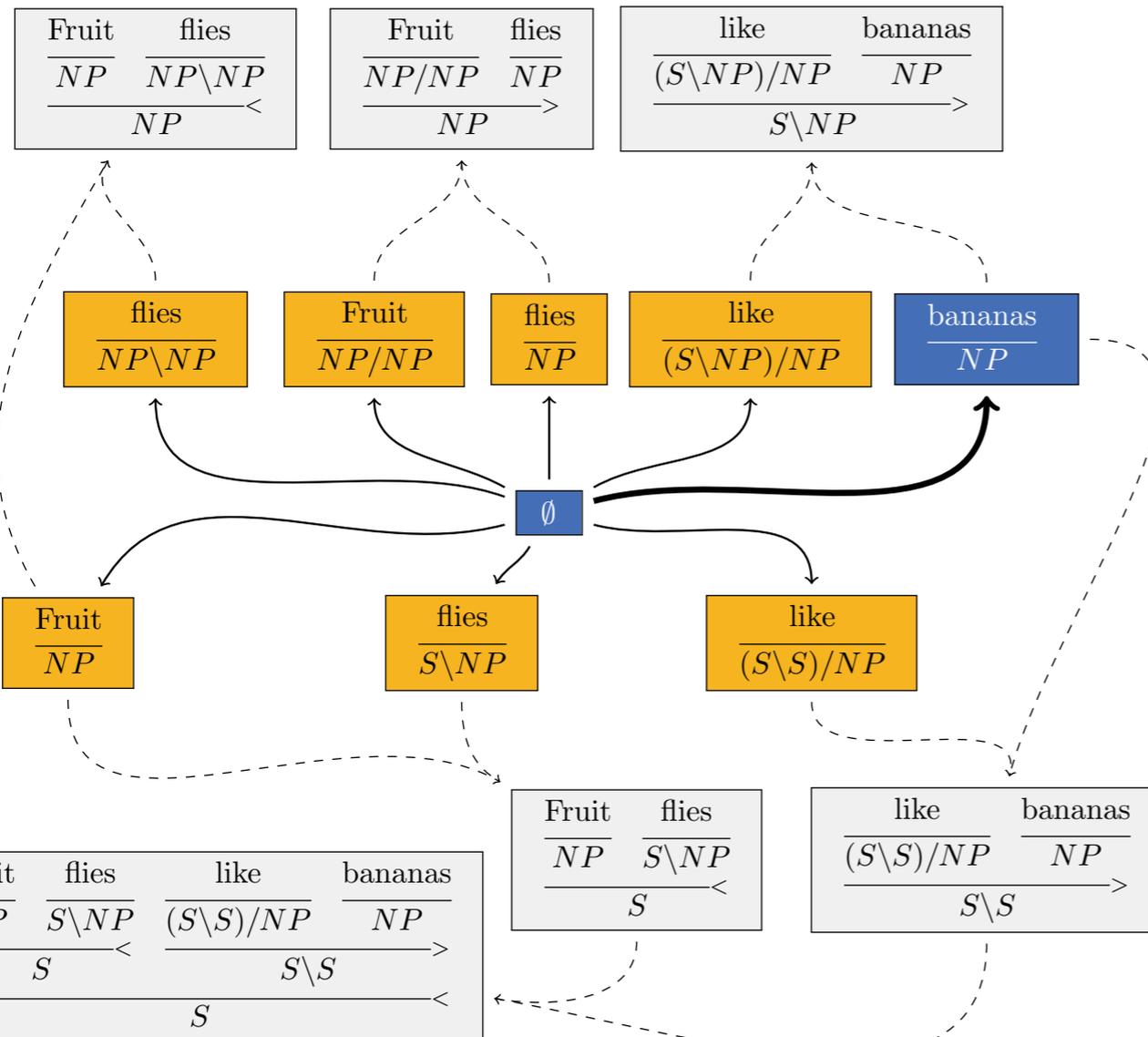
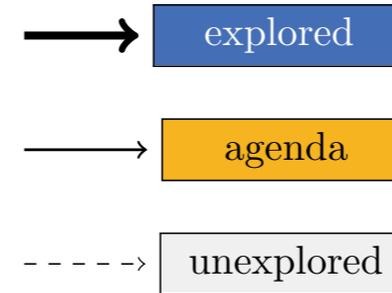
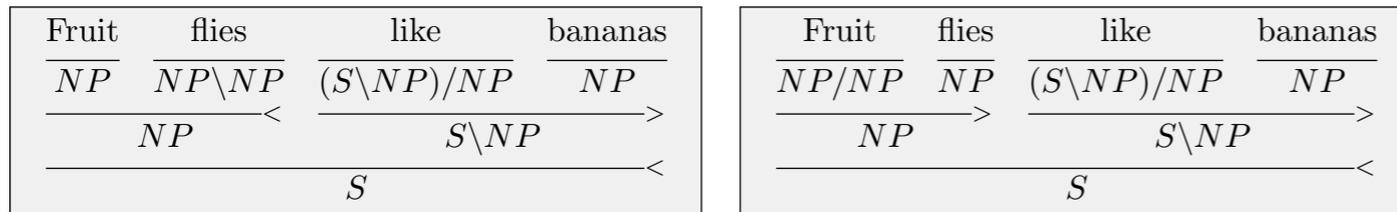
Agenda position	$f(y)$	y
1	4.5	\overline{NP} / bananas
2	3.1	$\overline{(S \setminus NP) / NP}$ / like
3	1.9	\overline{NP} / Fruit
4	-0.5	$\overline{NP / NP}$ / Fruit

A* Parsing



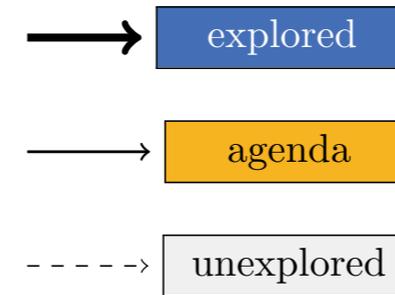
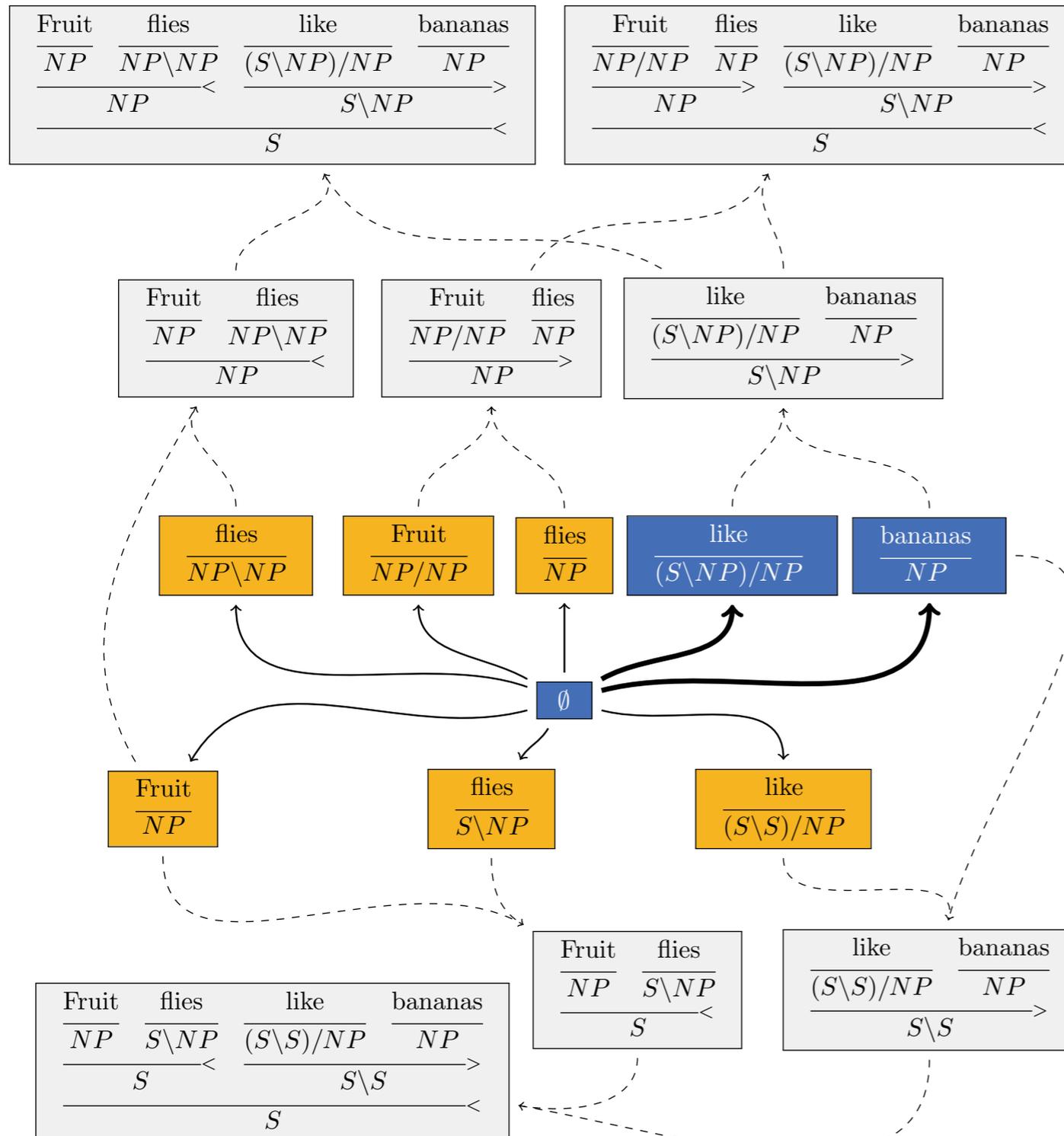
Agenda position	$f(y)$	y
2	3.1	$\frac{\text{like}}{(\text{S}\backslash\text{NP})/\text{NP}}$
3	1.9	$\frac{\text{Fruit}}{\text{NP}}$
4	-0.5	$\frac{\text{Fruit}}{\text{NP}/\text{NP}}$

A* Parsing



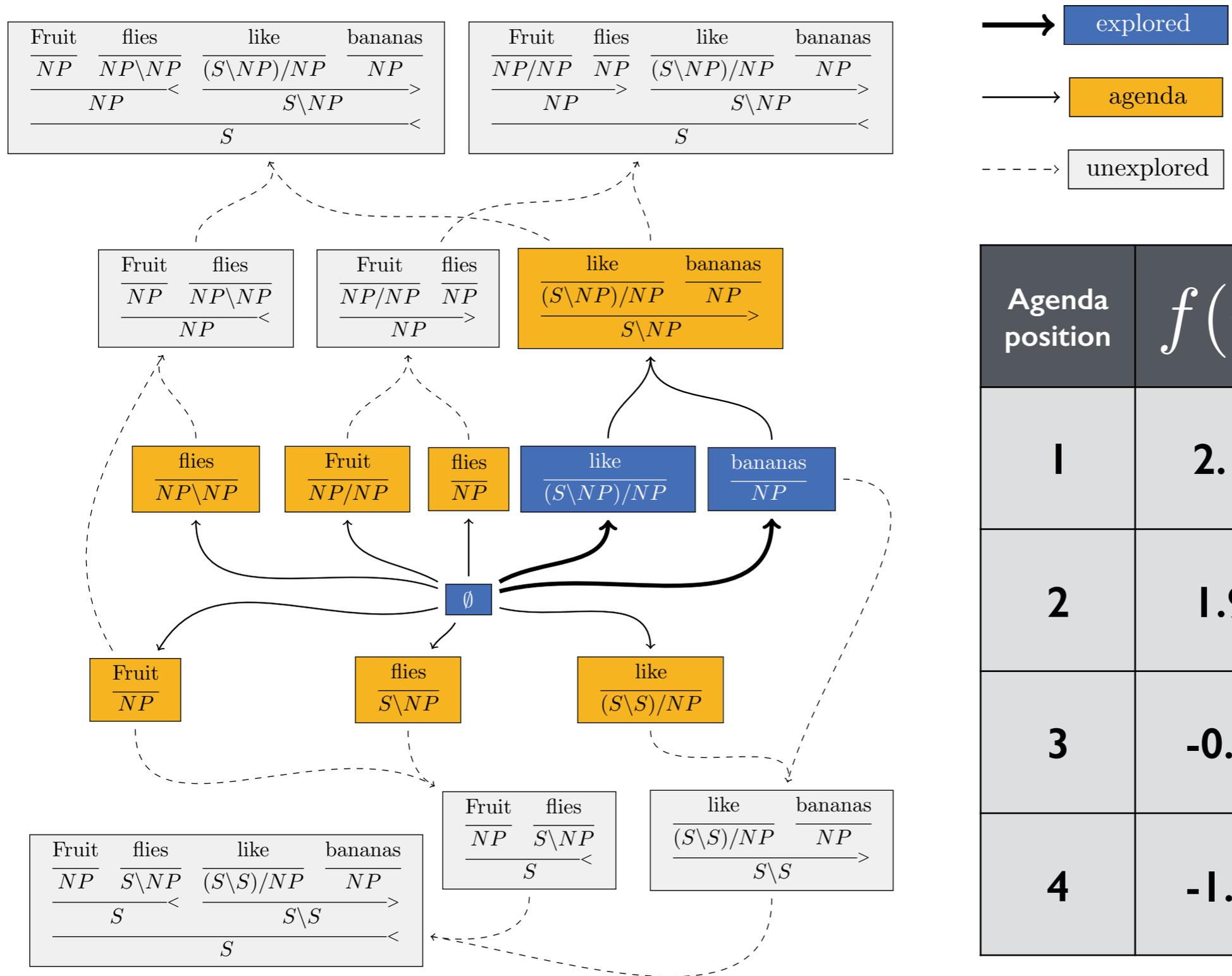
Agenda position	$f(y)$	y
1	3.1	$\frac{\text{like}}{(\text{S}\backslash\text{NP})/\text{NP}}$
2	1.9	$\frac{\text{Fruit}}{\text{NP}}$
3	-0.5	$\frac{\text{Fruit}}{\text{NP}/\text{NP}}$
4	-1.3	$\frac{\text{flies}}{\text{NP}}$

A* Parsing



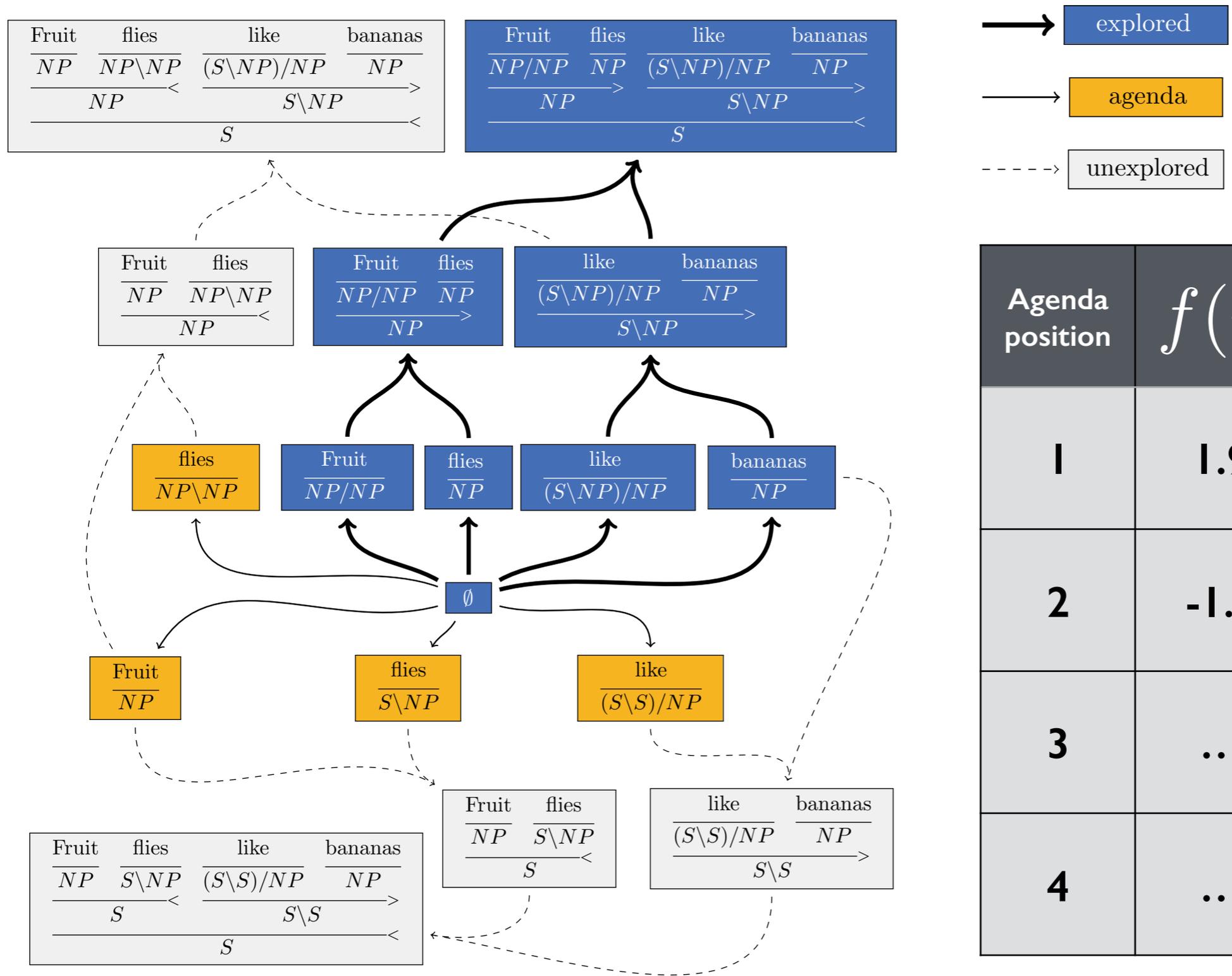
Agenda position	$f(y)$	y
1	3.1	$\frac{\text{like}}{(S \setminus NP) / NP}$
2	1.9	$\frac{\text{Fruit}}{NP}$
3	-0.5	$\frac{\text{Fruit}}{NP / NP}$
4	-1.3	$\frac{\text{flies}}{NP}$

A* Parsing



Agenda position	$f(y)$	y
1	2.1	$\frac{\text{like}}{(S \setminus NP) / NP} \frac{\text{bananas}}{NP} \rightarrow$ $S \setminus NP$
2	1.9	$\frac{\text{Fruit}}{NP}$
3	-0.5	$\frac{\text{Fruit}}{NP / NP}$
4	-1.3	$\frac{\text{flies}}{NP}$

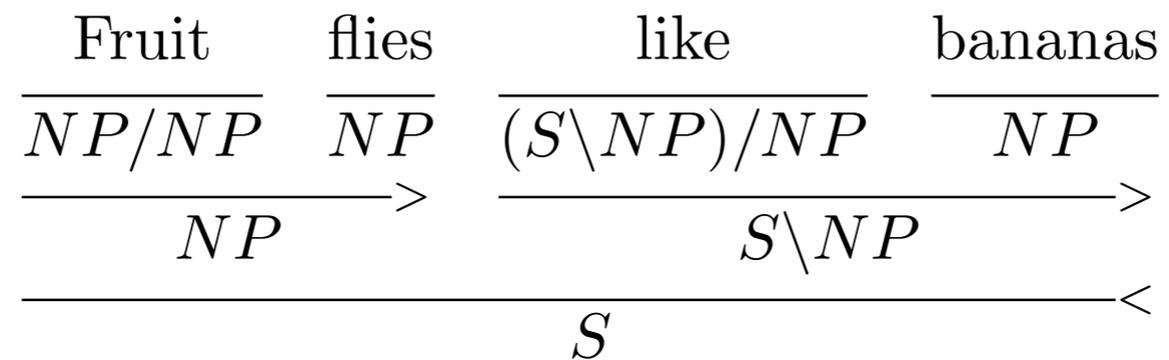
A* Parsing



Agenda position	$f(y)$	y
1	1.9	$\overline{Fruit / NP}$
2	-1.5	$\overline{like / (S \setminus S) / NP}$
3
4

Locally Factored Model

Supertag-factored A* CCG Parser (Lewis et al, 2016):



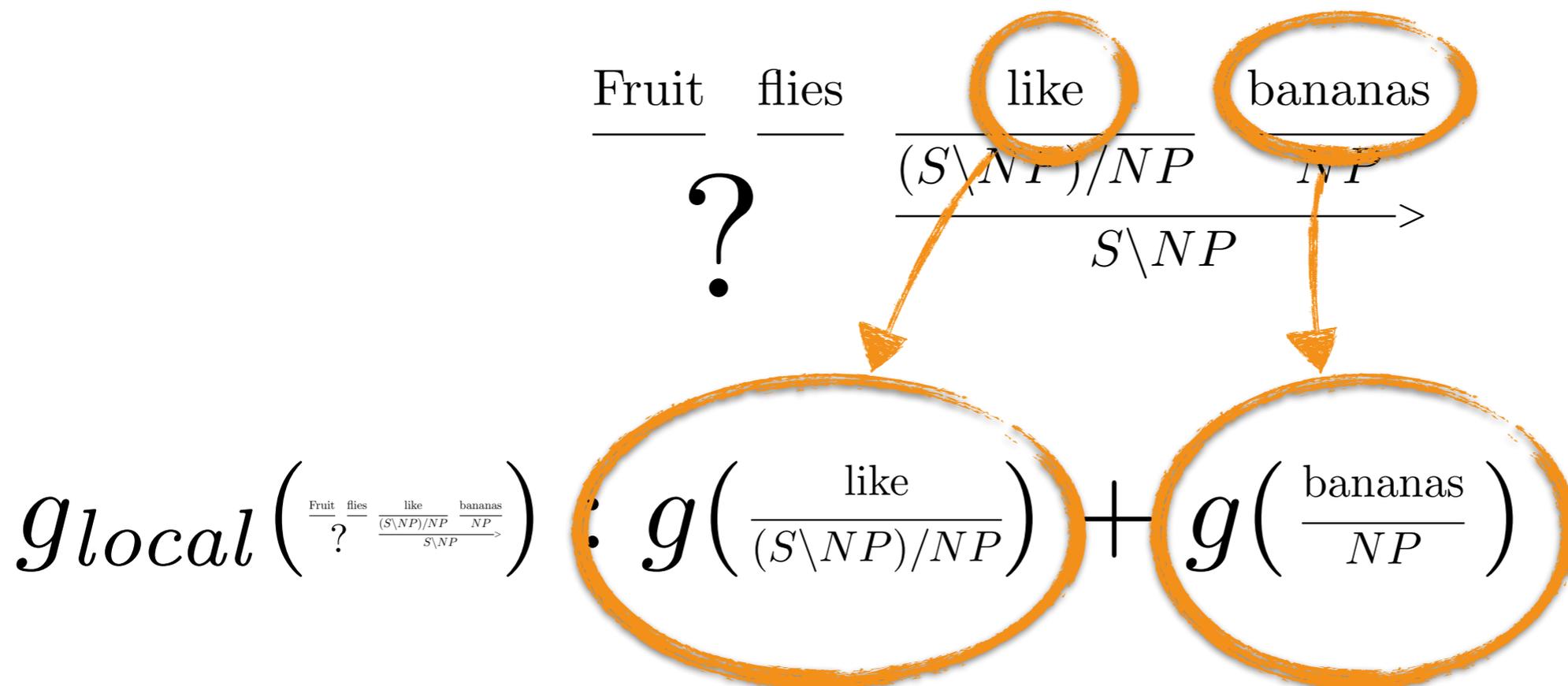
Locally Factored Model

Supertag-factored A* CCG Parser (Lewis et al, 2016):

$$\begin{array}{ccc} \text{Fruit} & \text{flies} & \text{like} & \text{bananas} \\ \hline & & (S \setminus NP) / NP & NP \\ \hline ? & & S \setminus NP & \rightarrow \end{array}$$

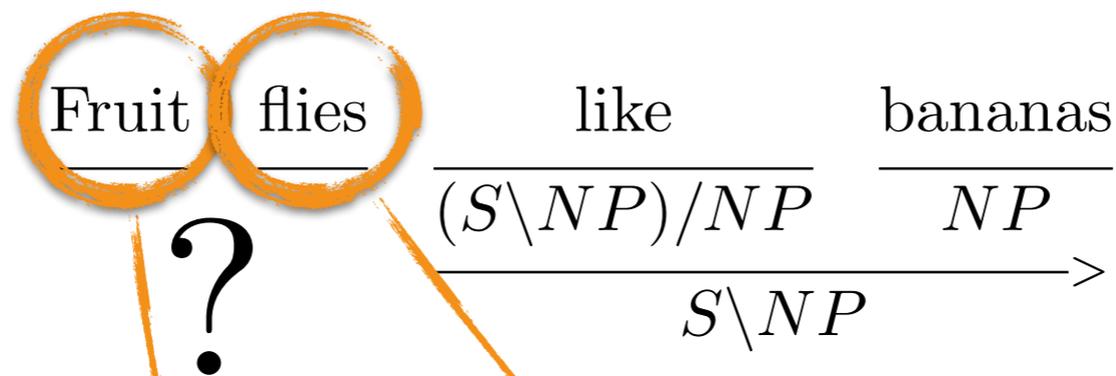
Locally Factored Model

Supertag-factored A* CCG Parser (Lewis et al, 2016):



Locally Factored Model

Supertag-factored A* CCG Parser (Lewis et al, 2016):



$$g_{local} \left(\frac{\text{Fruit flies}}{?} \frac{\text{like bananas}}{\frac{(S \setminus NP) / NP}{S \setminus NP} \quad NP} \right) : g \left(\frac{\text{like}}{(S \setminus NP) / NP} \right) + g \left(\frac{\text{bananas}}{NP} \right)$$

$$h_{local} \left(\frac{\text{Fruit flies}}{?} \frac{\text{like bananas}}{\frac{(S \setminus NP) / NP}{S \setminus NP} \quad NP} \right) : \max_{\text{tag}} g \left(\frac{\text{Fruit}}{\text{tag}} \right) + \max_{\text{tag}} g \left(\frac{\text{flies}}{\text{tag}} \right)$$

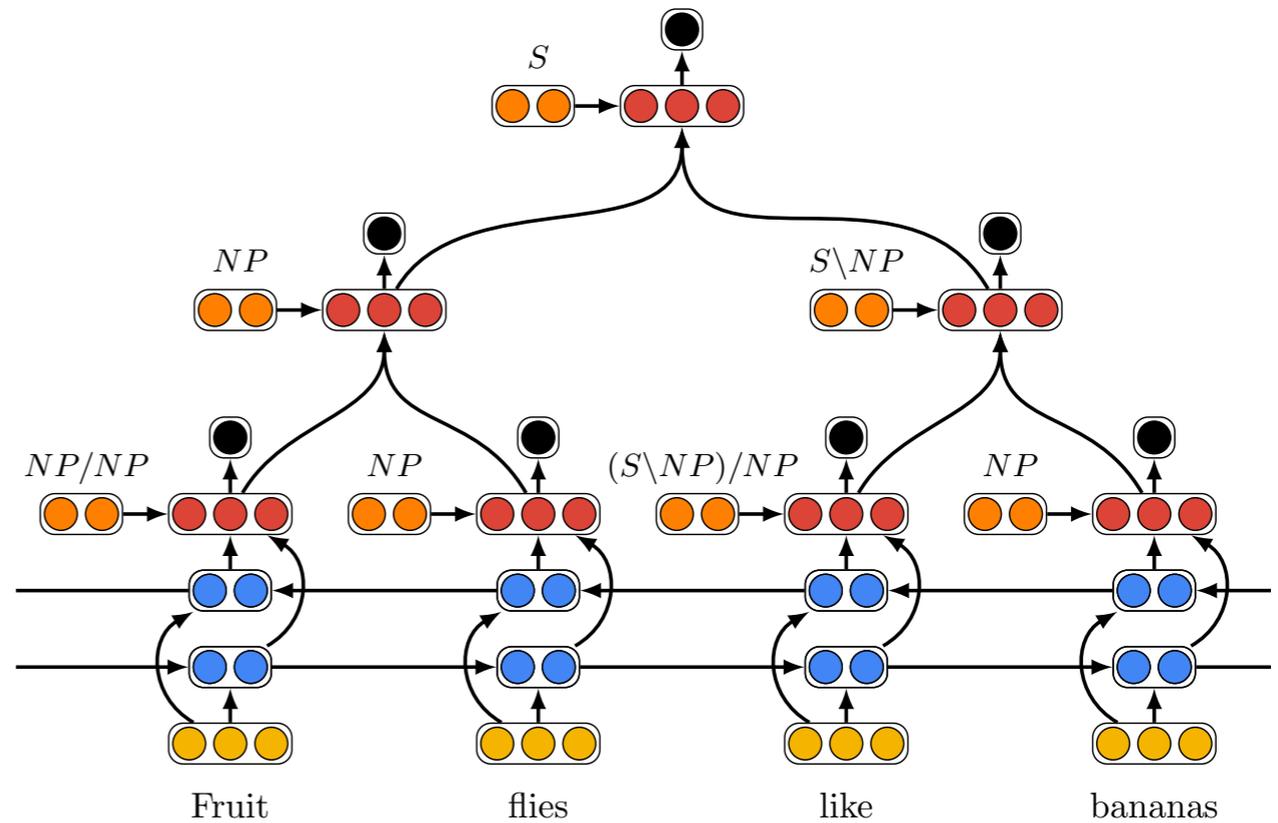
Global A* Parsing

$$y^* = \operatorname{argmax}_{y \in Y} g(y)$$

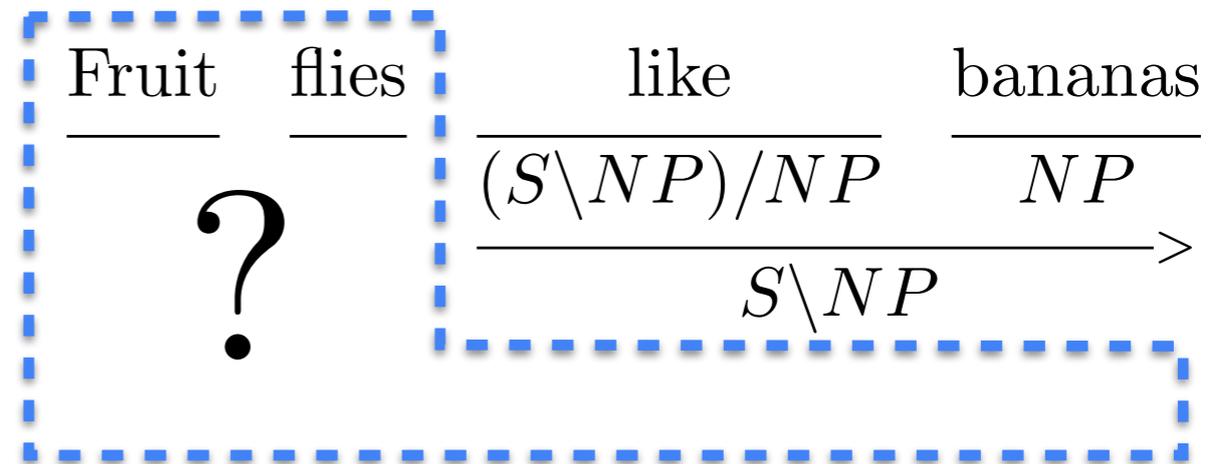
- ❖ First explored full parse **guaranteed to be optimal**
- ❖ Global search graph is **exponential** in sentence length
- ❖ Open question: Can we still **learn to search** efficiently?

Modeling Global Structure

$g_{global}(y) :$



$h_{global}(y) :$



Modeling Global Structure

$$g(y) =$$

$$g_{global}(y)$$

Non-positive
global model



$$h(y) =$$

0

Modeling Global Structure

$$g(y) = g_{local}(y) + g_{global}(y)$$

Any locally factored model with
an admissible A^* heuristic

Non-positive
global model

$$h(y) = h_{local}(y) + 0$$

Division of Labor

$$g(y) = g_{local}(y) + g_{global}(y)$$

- ❖ Limited expressivity
- ❖ Provides guidance with an A^* heuristic

- ❖ Global expressivity
- ❖ Discriminative only when necessary

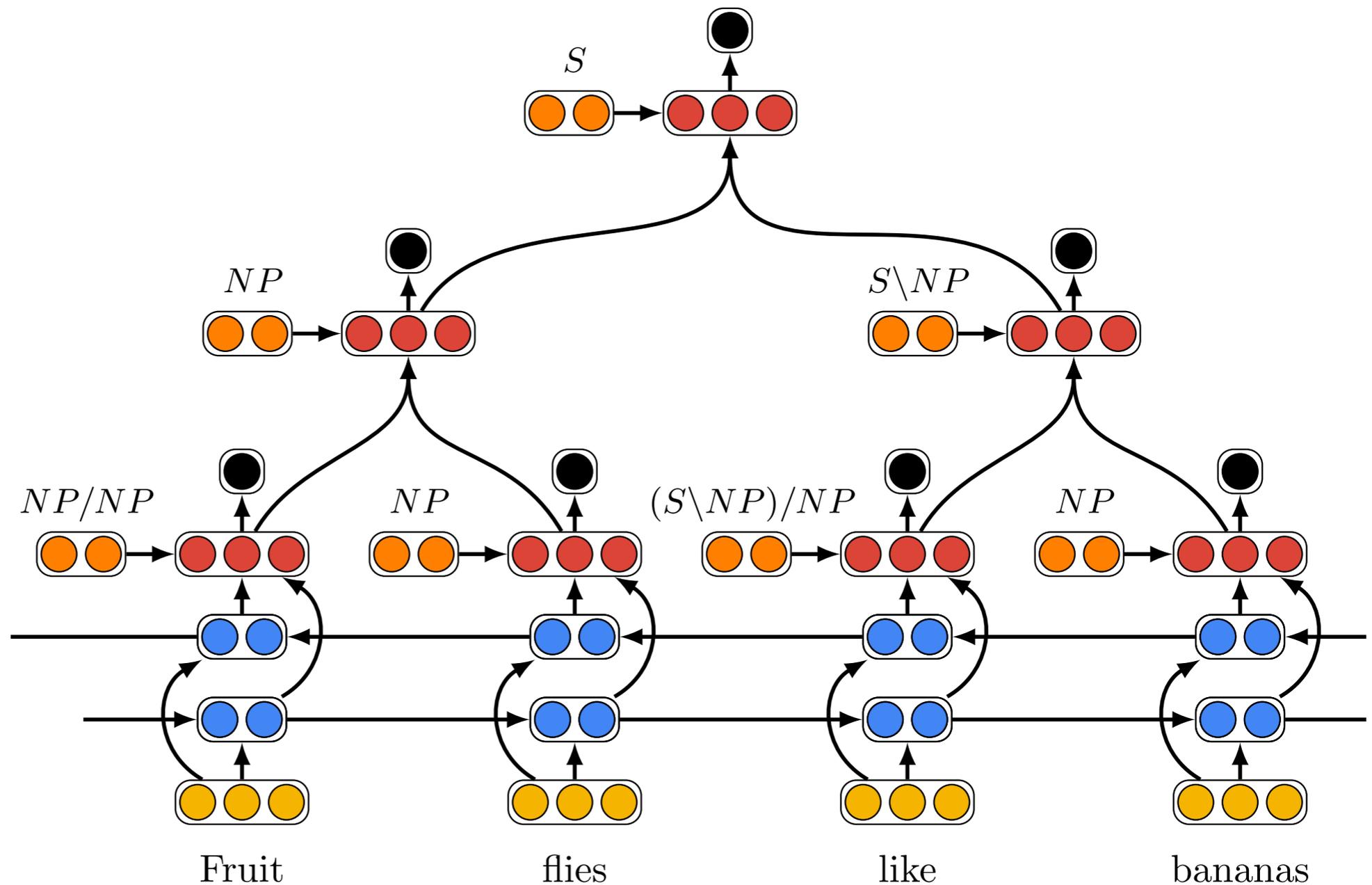
Global Model: $g_{global}(y)$

Parse Scores

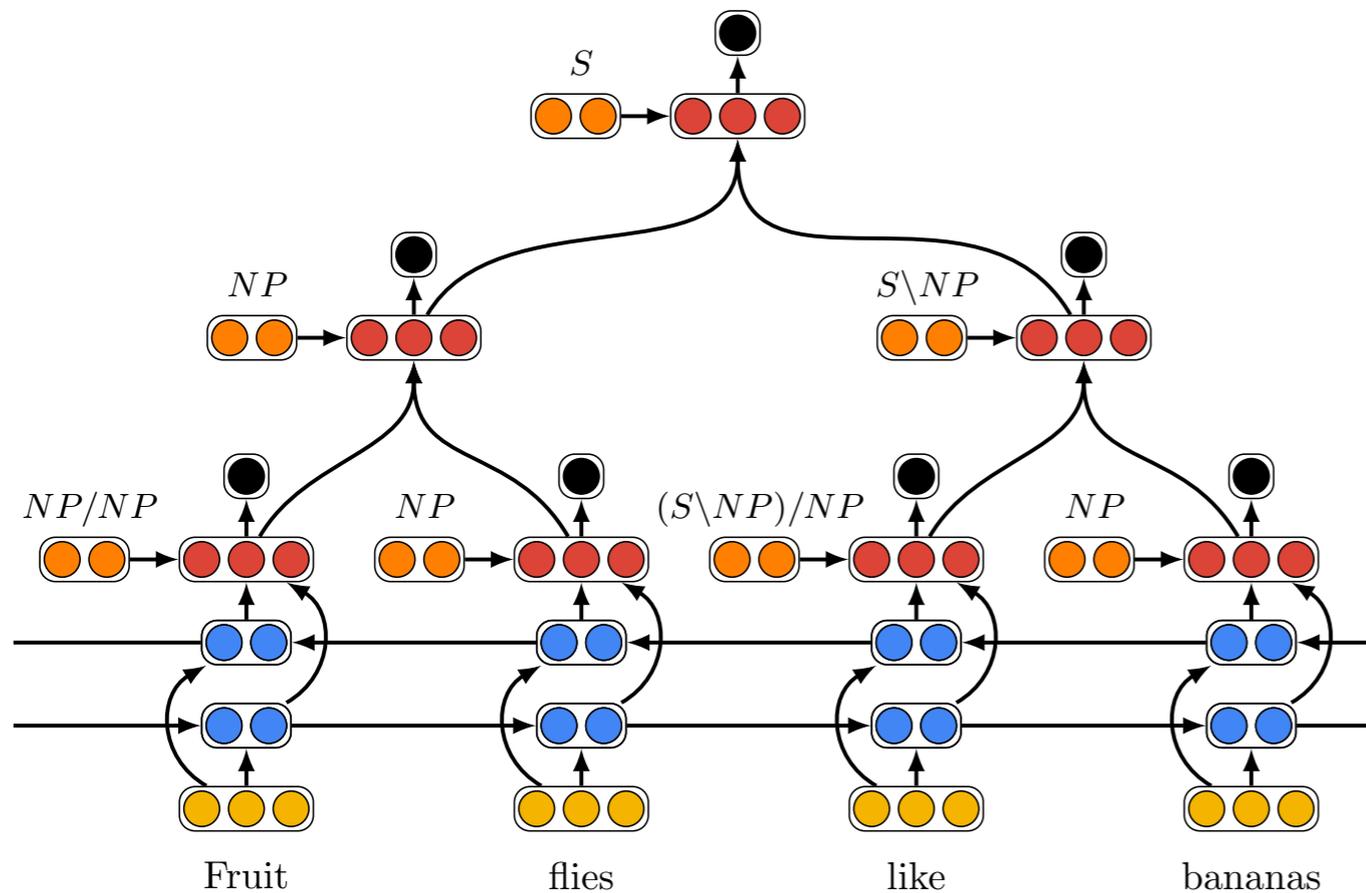
Tree-LSTM

Bidirectional LSTM

Word embeddings



Non-positive Global Model



Log-probability of a logistic regression layer

$$g_{global}(\text{red nodes}) = \log(\sigma(w \cdot \text{red nodes}))$$

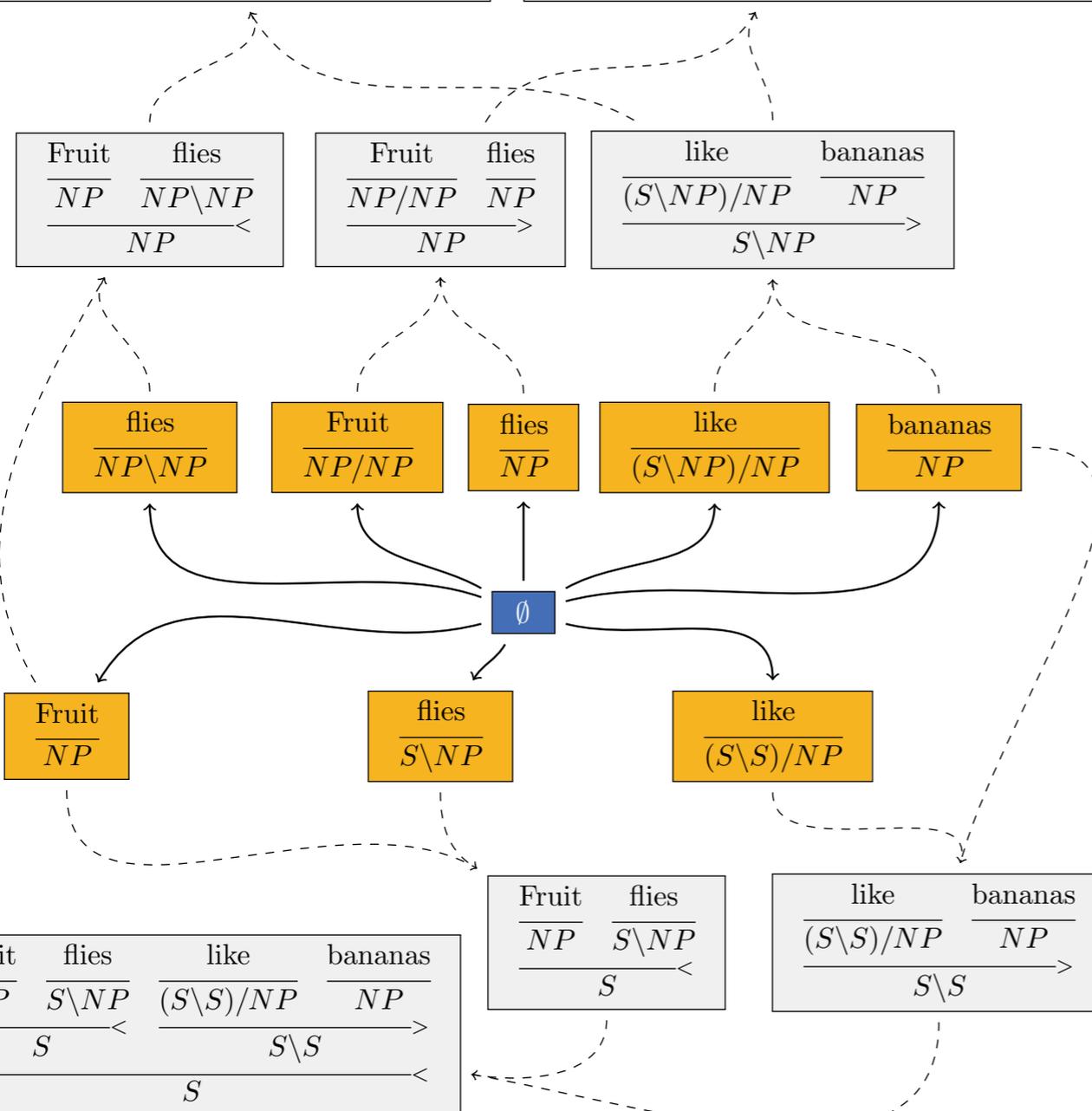
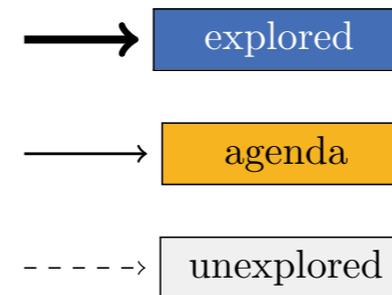
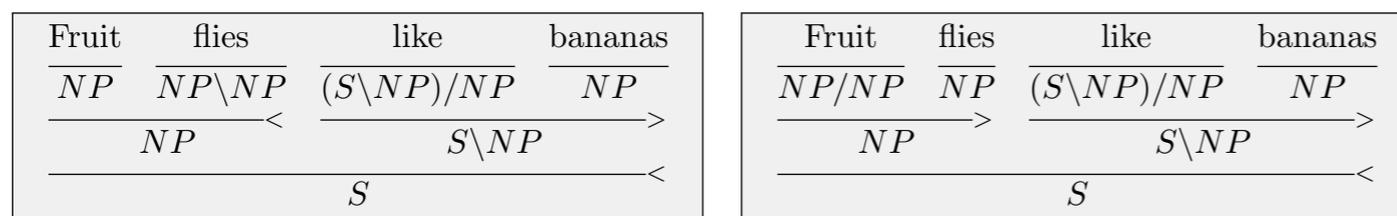
Division of Labor

$$g(y) = g_{local}(y) + g_{global}(y)$$

- ❖ Limited expressivity
- ❖ Provides guidance with an A^* heuristic

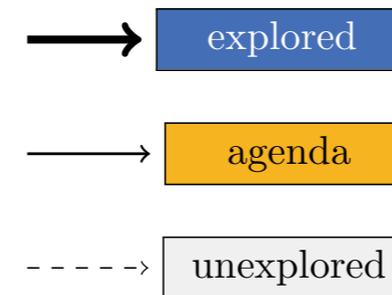
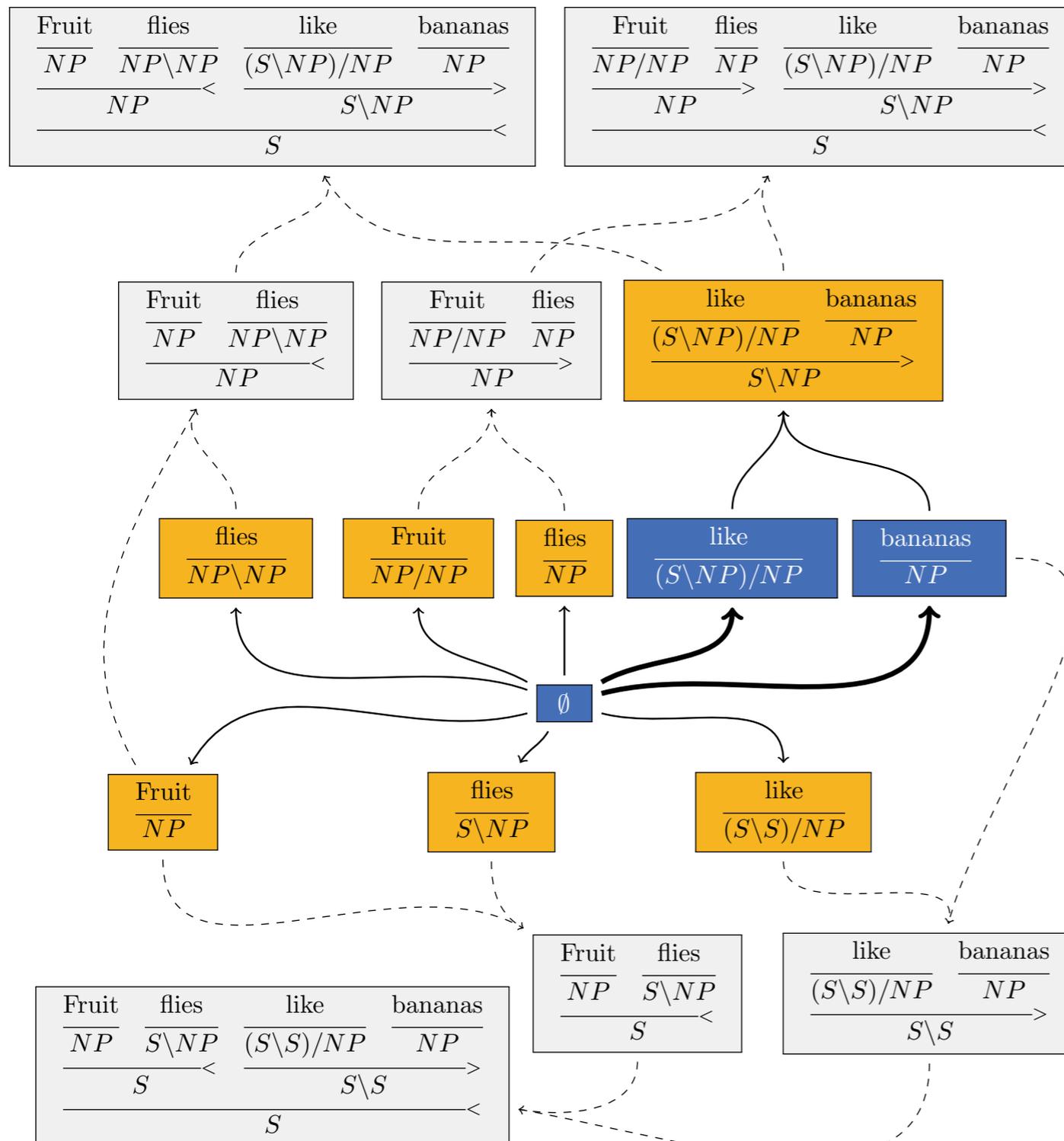
- ❖ Global expressivity
- ❖ **Discriminative only when necessary**

Learning with A*



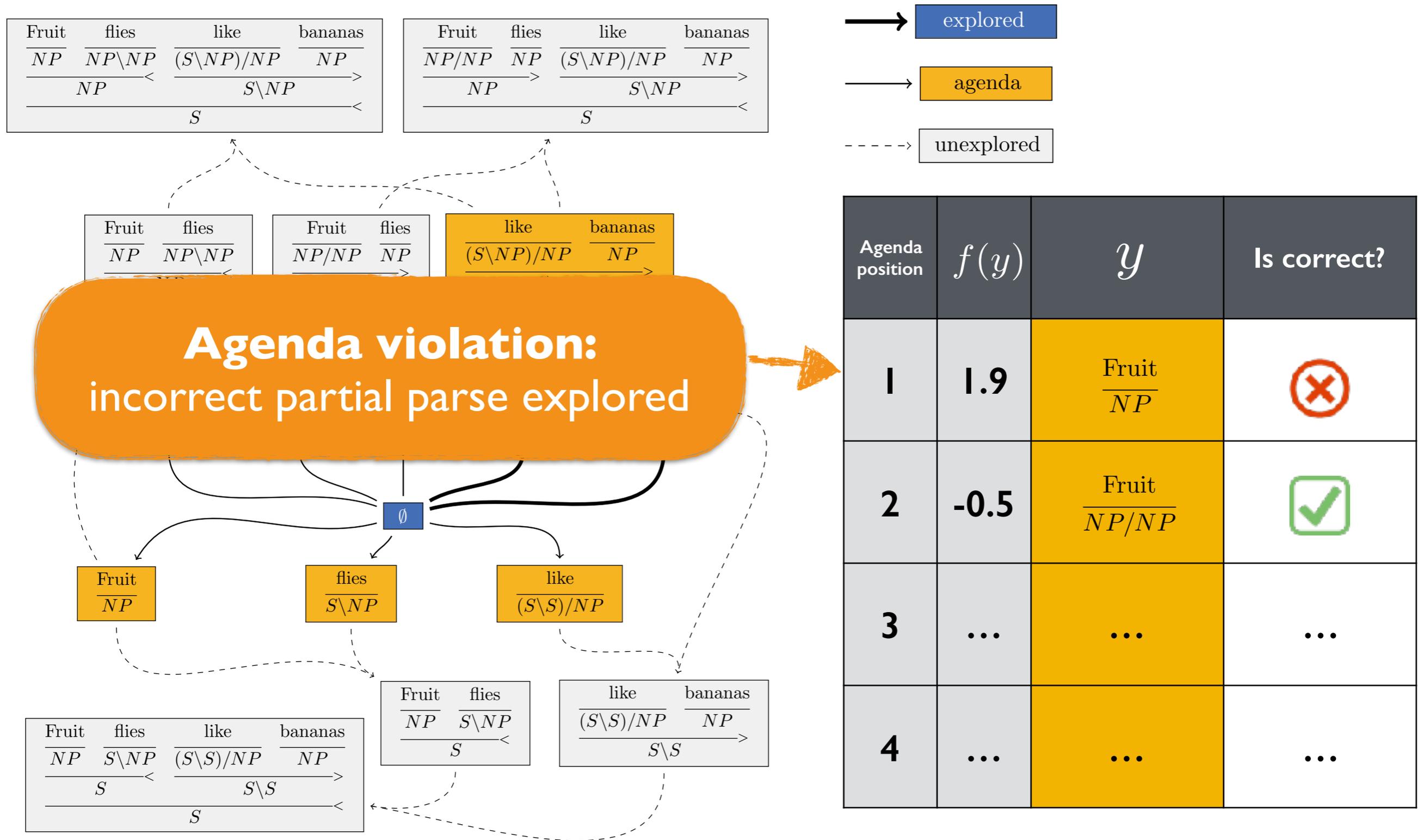
Agenda position	$f(y)$	y	Is correct?
1	4.5	$\frac{\text{bananas}}{NP}$	✓
2	3.1	$\frac{\text{like}}{(S \setminus NP) / NP}$	✓
3	1.9	$\frac{\text{Fruit}}{NP}$	✗
4	-0.5	$\frac{\text{Fruit}}{NP / NP}$	✓

Learning with A^*



Agenda position	$f(y)$	y	Is correct?
1	1.9	Fruit $\frac{NP}{NP}$	✗
2	-0.5	Fruit $\frac{NP}{NP/NP}$	✓
3
4

Learning with A*



Violation-based Loss

A

⋮ [

Agenda position	$f(y)$	y	Is correct?
1	4.5	bananas <u>NP</u>	✓
2	3.1	Is <u>(S, NP) NP</u>	✓
3	1.9	Fruit <u>NP</u>	✗
4	-0.5	Fruit <u>NP NP</u>	✓

■ ■ ■

Agenda position	$f(y)$	y	Is correct?
1	1.9	Fruit <u>NP</u>	✗
2	-0.5	Fruit <u>NP NP</u>	✓
3
4

■ ■ ■

]

Violation-based Loss

$\mathcal{A} : [$

Agenda position	$f(y)$	y	Is correct?
1	4.5	bananas NP	✓
2	3.1	Is (S,NP)NP	✓
3	1.9	Fruit NP	✗
4	-0.5	Fruit NP/NP	✓

... [

Agenda position	$f(y)$	y	Is correct?
1	1.9	Fruit NP	✗
2	-0.5	Fruit NP/NP	✓
3
4

]

$$L(\mathcal{A}) = \sum_{t=1}^T \underbrace{\max_{y \in \mathcal{A}_t} f(y)} - \underbrace{\max_{y \in \text{GOLD}(\mathcal{A}_t)} f(y)}$$

Top of agenda

Best gold partial parse

Jointly Optimizing Accuracy and Efficiency

Correct partial parse can still be predicted via backtracking

Agenda position	$f(y)$	y	Is correct?
1	1.9	$\frac{\text{Fruit}}{NP}$	
2	-0.5	$\frac{\text{Fruit}}{NP/NP}$	
3
4

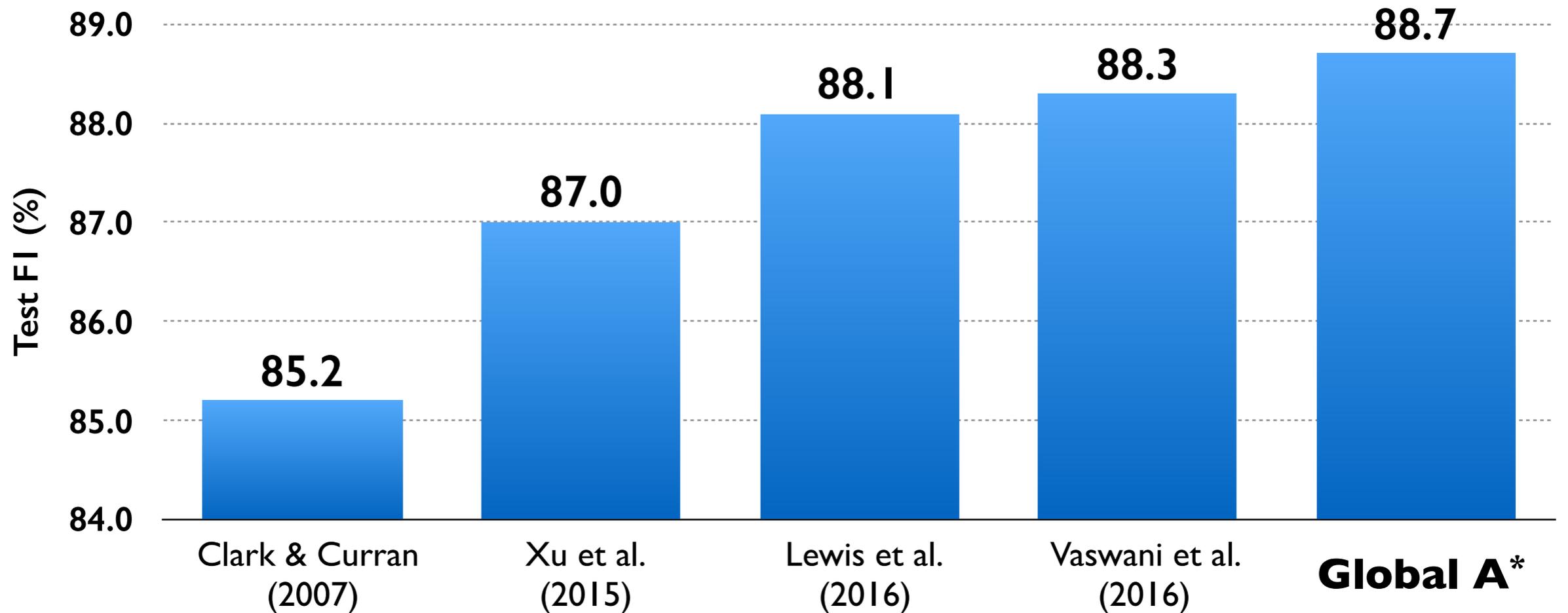
Jointly Optimizing Accuracy and Efficiency

Agenda position	$f(y)$	y	Is correct?

Explicitly optimize for search efficiency!

3
4

CCG Parsing Results

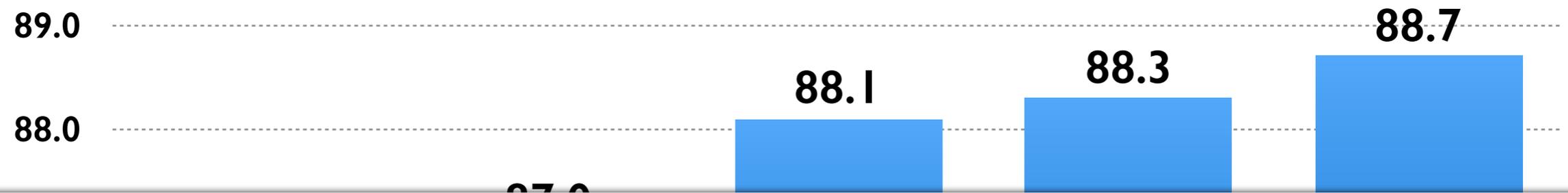


Is global?

Is exact?

	✓		✓	✓
		✓		✓

CCG Parsing Results



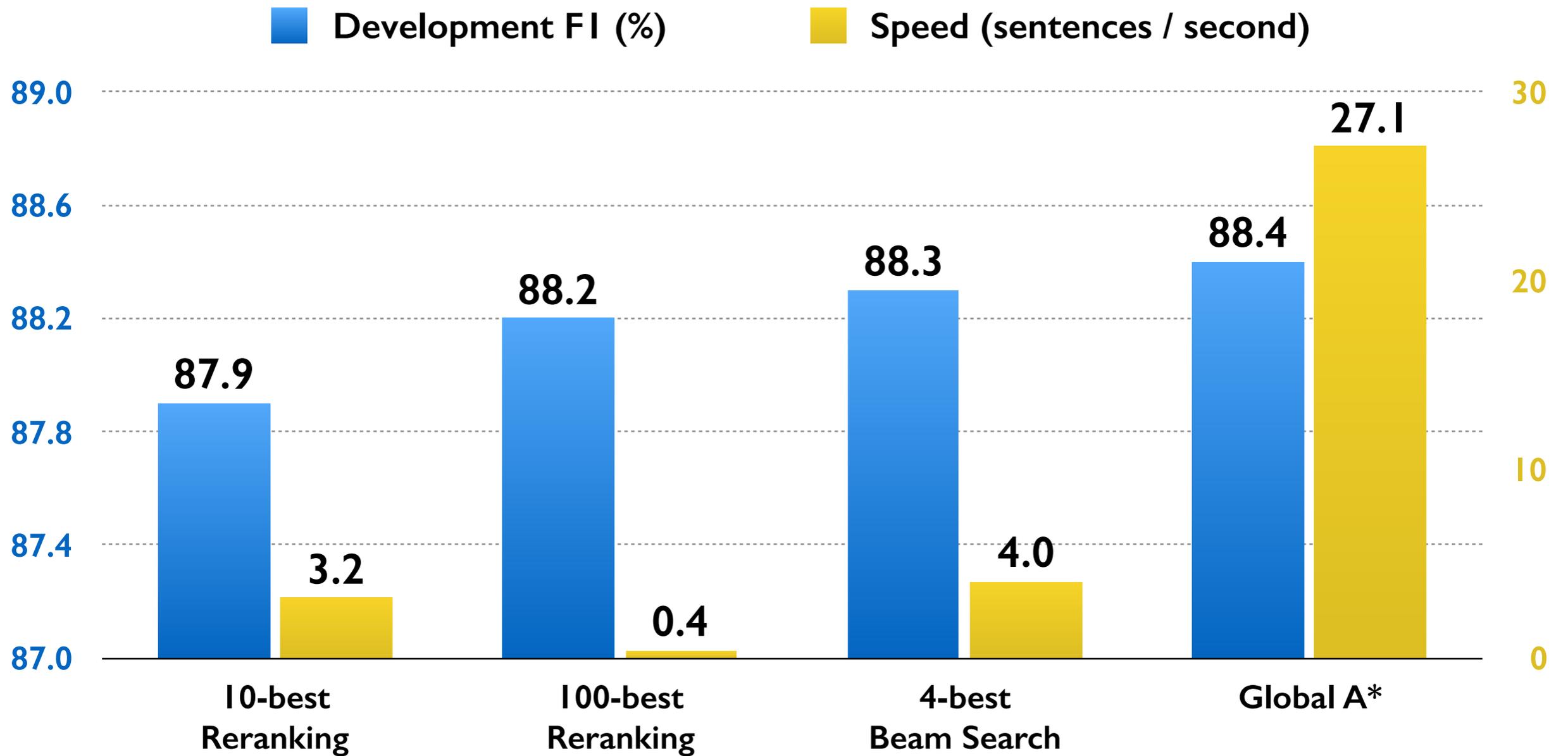
- ❖ Optimal parse found for 99.9% of sentences
- ❖ Explores only 190 partial parses on average

Is global:

Is exact?

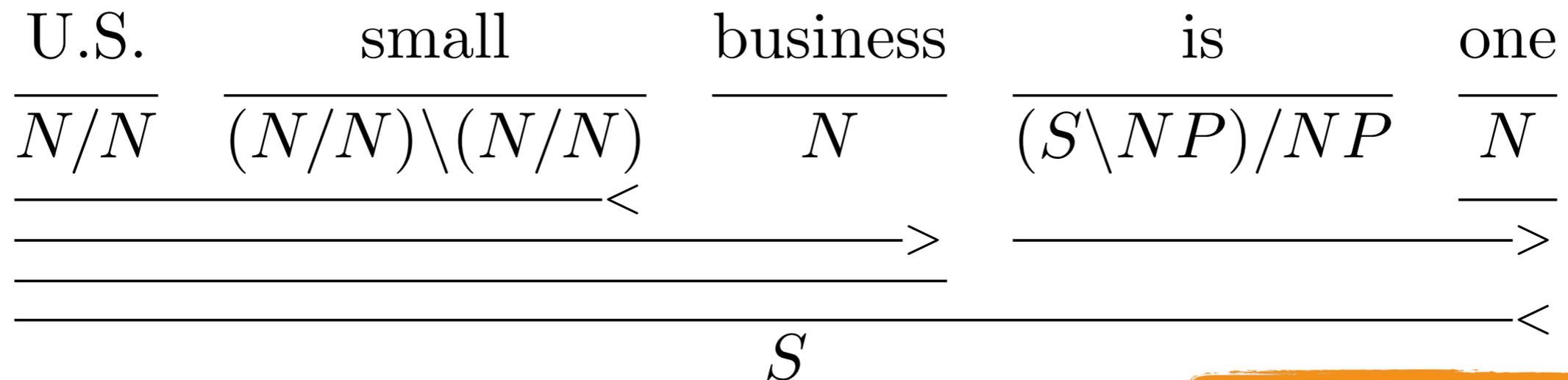
	✓		✓	✓
		✓		✓

Decoder Comparisons



Garden Paths

Incorrect partial parse (syntactically plausible in isolation):



Heavily penalized by
the global model

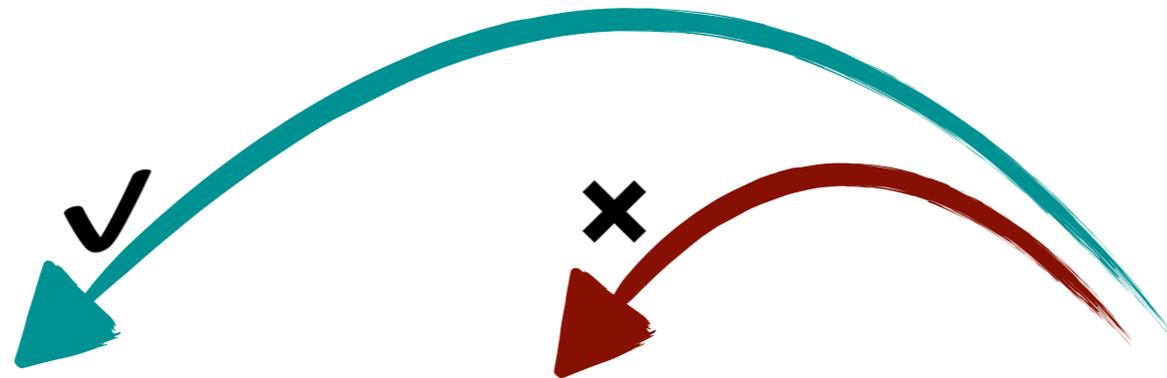
Input sentence:

The favorite **U.S. small business is one** whose research and development can be milked for future Japanese use.

Towards Broad Coverage Semantic Parsing

- Can we crowdsource semantics?
- Train with latent syntax?
- Build fast and accurate parsers?
- Actively select which data to label?

Our key hypothesis:
Anyone who **understands the meaning of a sentence**
should be able to correct **parser mistakes**.



Pat ate the cake on the table that I ***baked*** last night.

Parser: I baked **table**

Human understanding: I baked **cake**

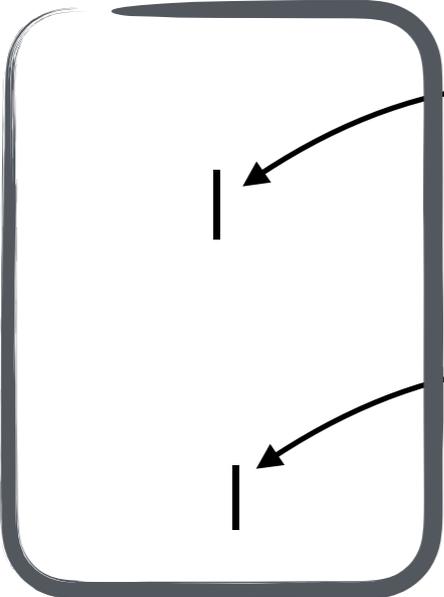
Can we use human judgements to improve parse?

[He et al, 2016]

Pat ate the cake on the table that I **baked** last night.

Q: What did someone **bake**?
1. table 2. cake

**parses from
the n-best list**

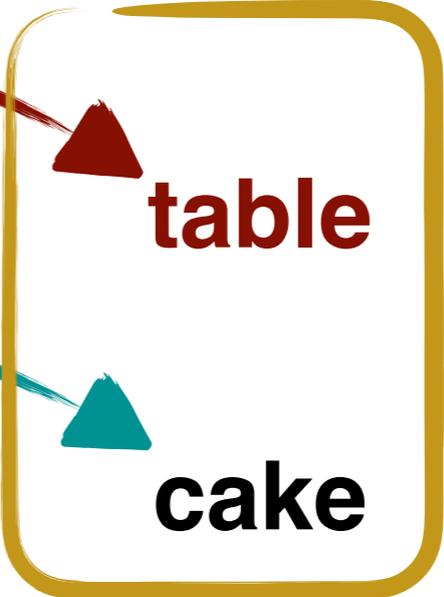


**Confident
attachment
decision**

baked

baked

.... ..



**Uncertain
attachment
decision**

**human
judgment**

Workflow

**Candidate dependencies
from the n-best list:**

baked → table
baked → cake

**CCG
Parser**

**Question
Generator**

Q: “What did
someone ***bake***?”
1) table **2)** cake

**Crowdsourcing
Platform**



**Re-parsed CCG
Dependency
Tree**

**Re-parse
w/ Constraints**

C_pos (bake → cake)
C_neg (bake → table)

cake (4 votes)
table (1 vote)

**Not re-training
the model**

Generate Q/A Pairs from CCG Dependencies

Predicted CCG category of **baked**: $(S \setminus NP_1) / NP_2$

Convert to template: NP₁ bake NP₂

Filling-in the Slots:

what bake sth.

I baked table
What baked something?
— I

I baked cake
What baked something?
— I

sth. bake what

What did someone bake?
— **the table**

What did someone bake?
— **the cake**

Infer **someone/something** and the **answer spans** based on the n-best parses

Used “**what**” for all questions

Group Q/A Pairs into Queries

Questions	Answers	Scores	Question Confidence	Answer Uncertainty (Entropy)
What baked something?	I	1.0	1.0	0.0
What did someone bake?	the table	0.7	1.0	0.88
	the cake	0.3		
What was baked something something?	the table	0.1	0.1	0.0

Non-sensical question

No uncertainty

Our Annotation Task

Sentence:

Pat ate the cake on the table that I **baked** last night.

Question:

What did someone bake?

Check one or more

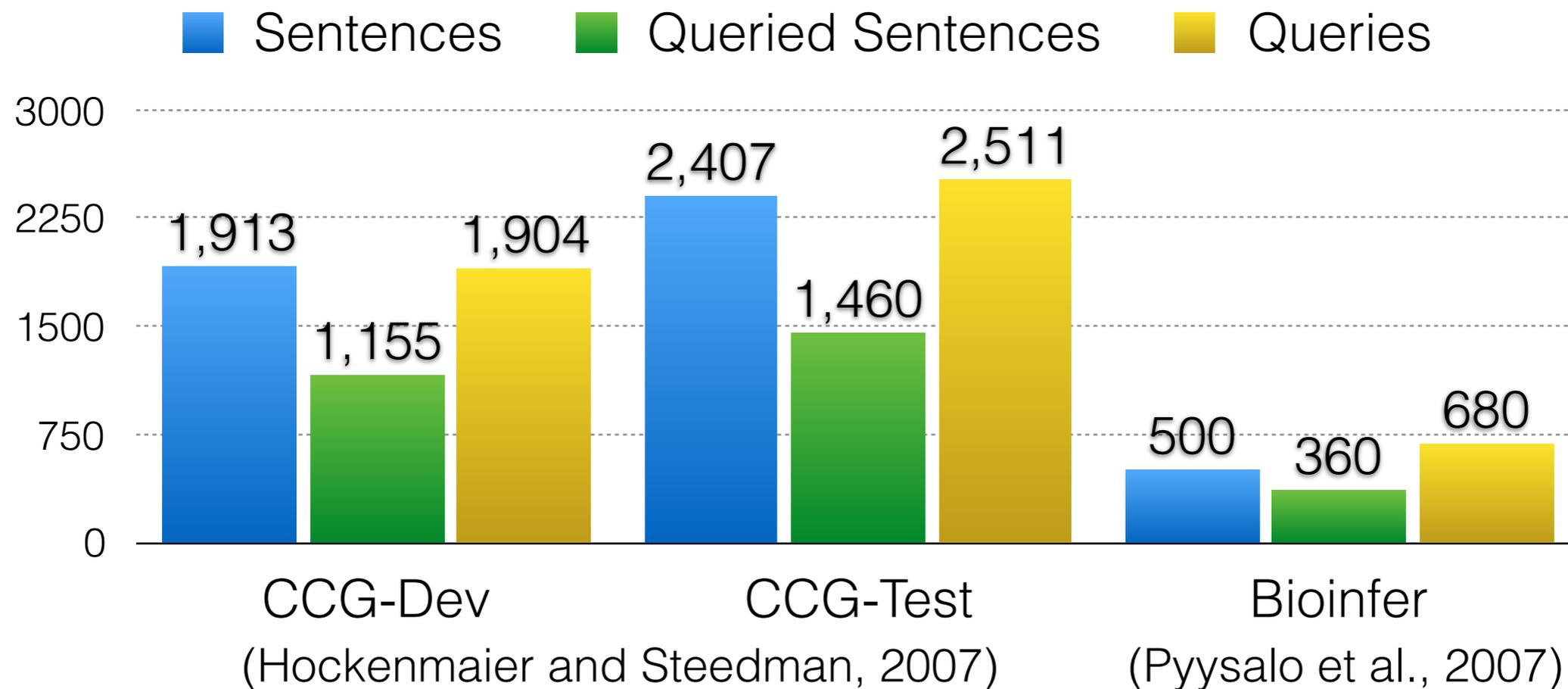
- the cake
- the table
- None of the above.

- Annotators are instructed to choose options that “***explicitly and directly***” answer the question.
- Multiple answers are allowed.
- 5 judgements per query.

Comment:

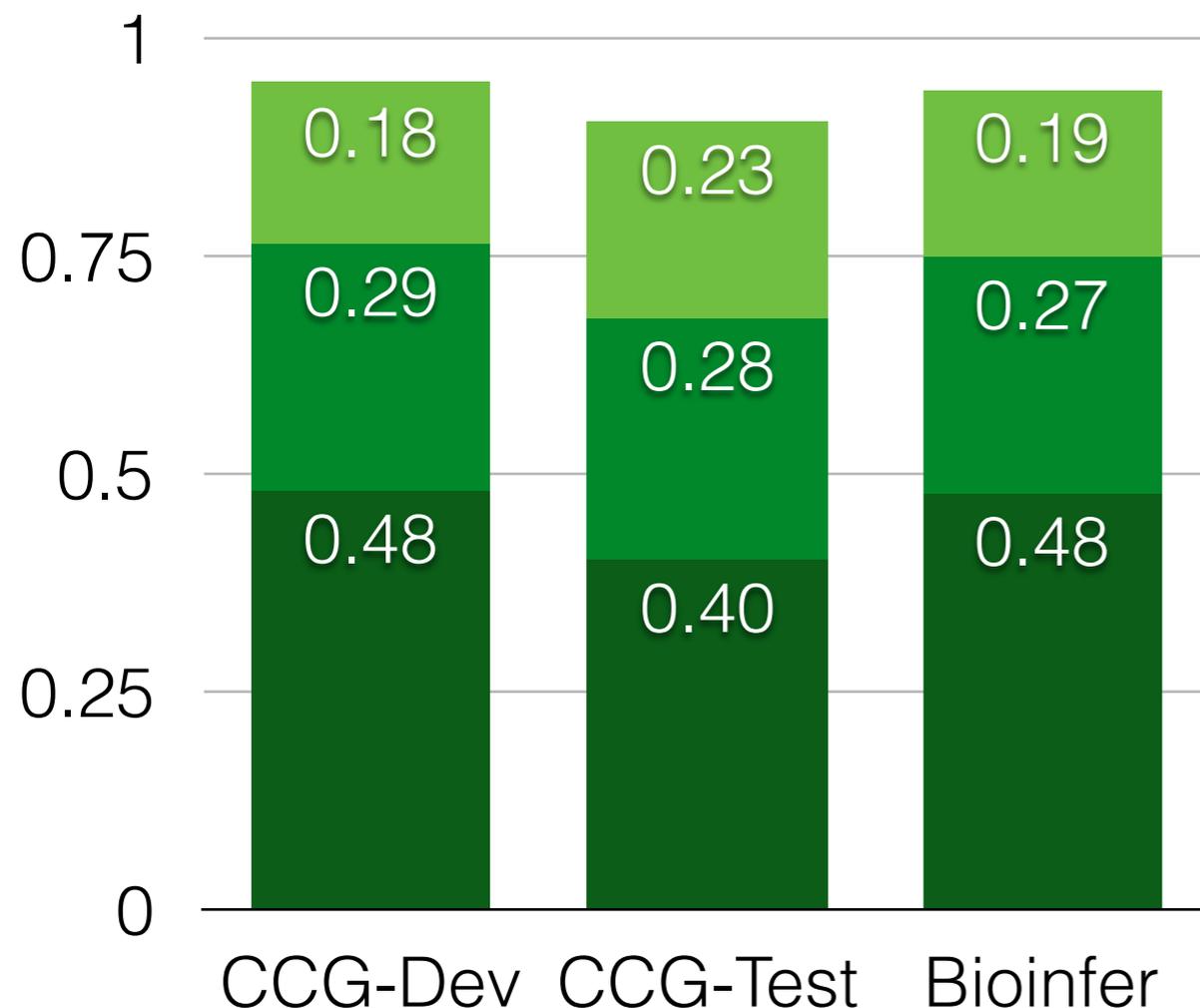
* Crowdsourcing platform: <https://www.crowdflower.com/>.

Data Collection with Crowdsourcing



- All developments are done on CCG-Dev only.
- Less than 2 queries per sentence, for about 60% of the sentences.
- **Cost:** 46 cents per query.
- **Speed:** 200 queries per hour.

Inter-Annotator Agreement



- Agreement is computed only for matching the exact set of answers. i.e. (A, B) and (B) are considered disagreement.
- Unanimous agreement for over 40% of the queries.
- Over 90% absolute majority.

Putting our hypothesis to the test:
How well does annotators' **human understanding**
align with the **gold syntax**?

- Successes: Long-range attachment decisions
- Challenges: **Syntax-semantics mismatch**
- Use heuristics to fix the mismatch problems at re-parsing time.

Success - Long-range Dependency

Temple also said Sea Containers' plan raises numerous legal, regulatory, financial and fairness issues, but didn't ***elaborate***.

What *didn't* ***elaborate*** something?

4

Temple

1

Sea Containers' plan

0

None of the above.

Success - Coordination

To **avoid** these costs, and a possible default, immediate action is imperative.

What would something **avoid**?

- 4** these costs
- 3** a possible default
- 0** None of the above.

Challenge - Coreference

Kalipharma is a New Jersey-based pharmaceuticals concern that ***sells*** products under the Purepac label.

What ***sells*** something?

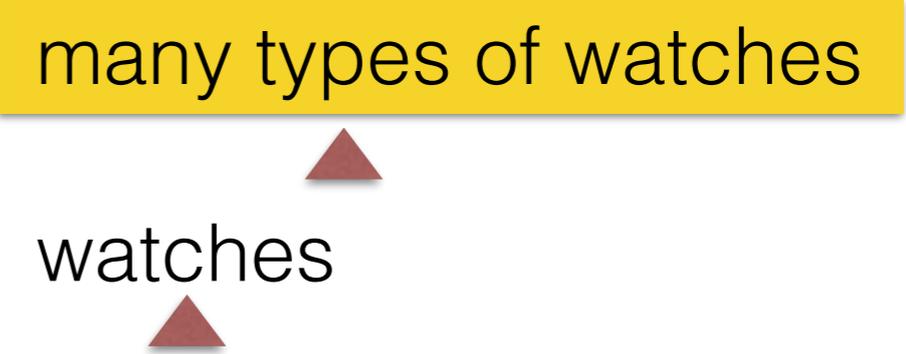
- 5 Kalipharma
- 0 a New Jersey-based pharmaceuticals concern
- 0 None of the above.

- Syntax-semantics mismatch
- Also happens with pronouns and appositives.
- Some cases are heuristically fixed during reparsing.

Challenge - Headedness

Timex had requested duty-free treatment for many types of watches,
covered by 58 different U.S. tariff classifications.

What would be **covered**?

- | | | | |
|----------|---------------------|----------|-----------------------|
| 0 | Timex | 2 | many types of watches |
| 0 | duty-free treatment | 3 | watches |
| 0 | None of the above. | | |
- 

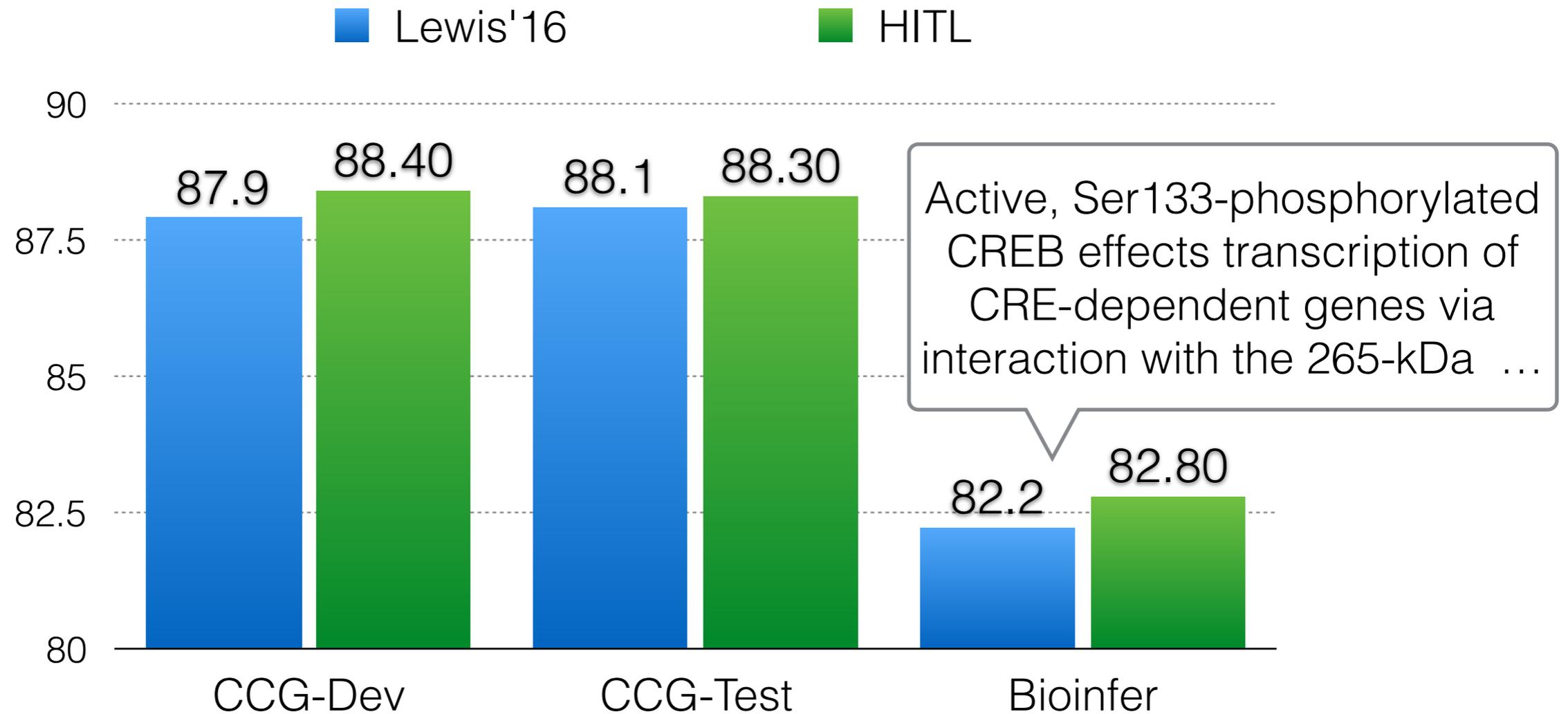
- Annotators tend to struggle with headedness.
- We add “disjunctive constraint”, forcing the re-parser to produce either of the two dependencies.

Re-Parsing with Crowdsourced Constraints

Q1: What did someone **bake**? $y^{\text{new}} = \arg \max_y \text{base_parser_score}(y)$
votes(cake) = 4 $-T^+ \times \mathbb{1}(\text{baked} \rightarrow \text{cake} \in y)$
votes(table) = 1 $-T^- \times \mathbb{1}(\text{baked} \rightarrow \text{table} \in y)$
votes(*None of the above*) = 0

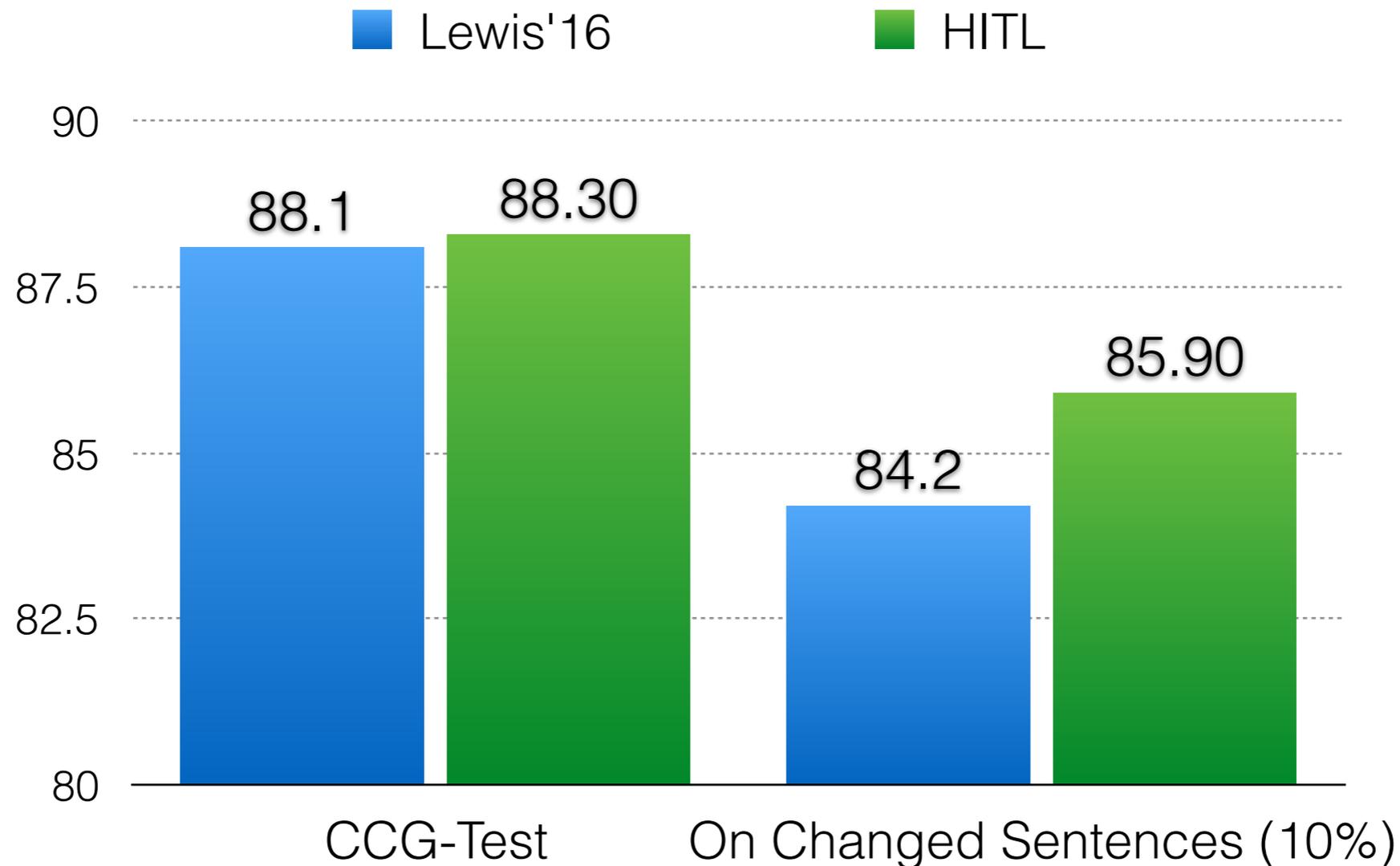
- Penalizes parses that disagree with crowdsourced judgments.
- Constraints are decomposed by dependencies.
- Thresholds and penalties are tuned on CCG-Dev.

Re-parsing Results (Labeled F1)



- Modest improvement due to syntax-semantics mismatch.
- Larger improvement on out-of-domain data.

Re-parsing Results



- Modified parse trees for about 10% of the sentences after incorporating human judgments.
- Larger gain on changed sentences.
- Changed sentences are “more difficult” on average.

Towards Broad Coverage Semantic Parsing

- Can we crowdsource semantics?
 - Yes, but need more than verbs....
- Train with latent syntax?
 - Yes, but must extend to QA supervision...
- Build fast and accurate parsers?
 - Yes, but need to extend to latent-variable case...
- Actively select which data to label?
 - Yes, but need to scale up...

Questions

?