



# Algorithms that learn to think on their feet

# What is NLP?



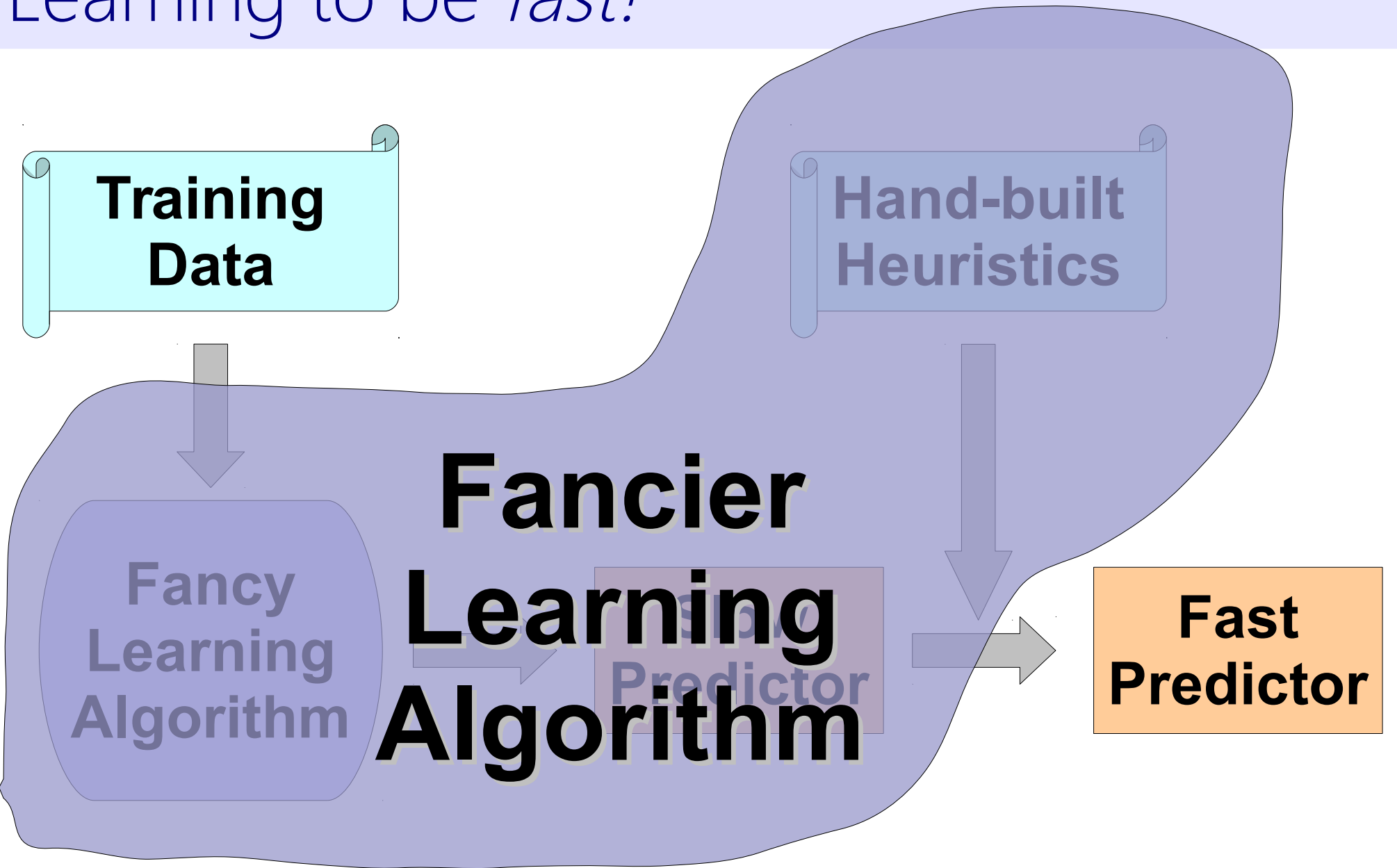
- **Fundamental goal: deep understanding of text**
  - Not just string processing or keyword matching
- **End systems that we want to build**
  - Simple: Spelling correction, text categorization, etc.
  - Complex: Speech recognition, machine translation, information extraction, dialog interfaces, question answering
  - Unknown: human-level comprehension (more than just NLP?)

# Why is language **hard**?

- **Ambiguity abounds (some headlines)**
  - Iraqi Head Seeks Arms
  - Teacher Strikes Idle Kids
  - Kids Make Nutritious Snacks
  - Stolen Painting Found by Tree
  - Local HS Dropouts Cut in Half
  - Enraged Cow Injures Farmer with Ax
  - Hospitals are Sued by 7 Foot Doctors
  - Ban on Nude Dancing on Governor's Desk
  - Scientists study whales from space
- **Why are these funny?**
- **What does ambiguity imply about the role of learning?**

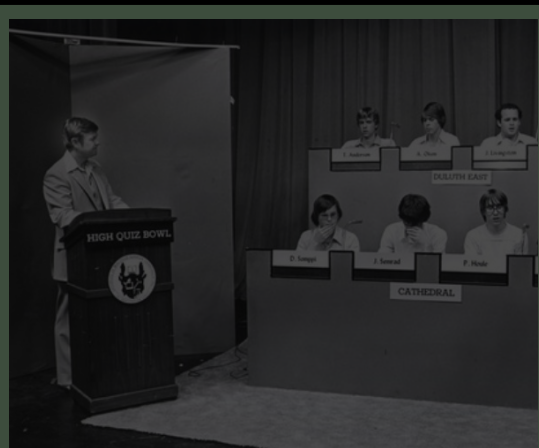
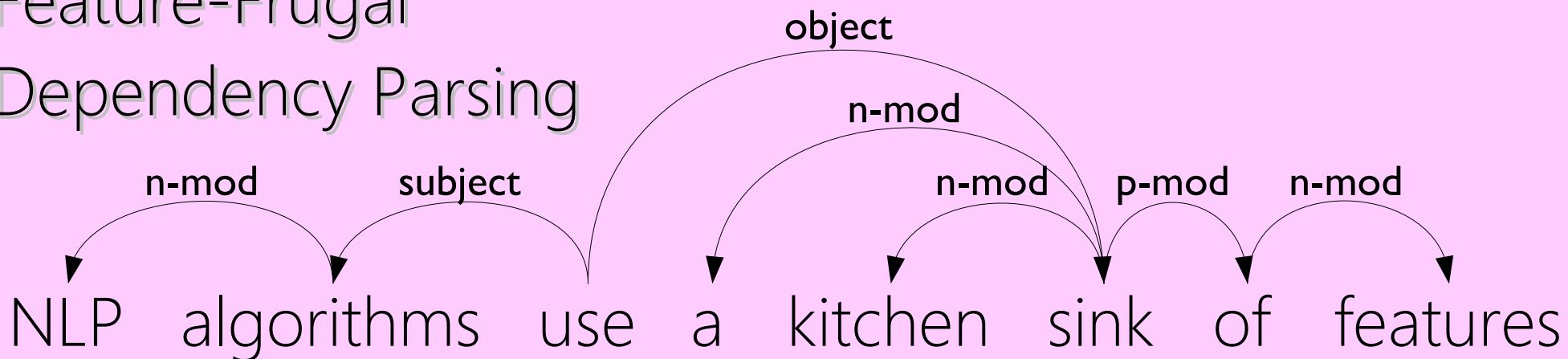


# Learning to be *fast!*



Quality = tradeoff(accuracy, time)

# Feature-Frugal Dependency Parsing



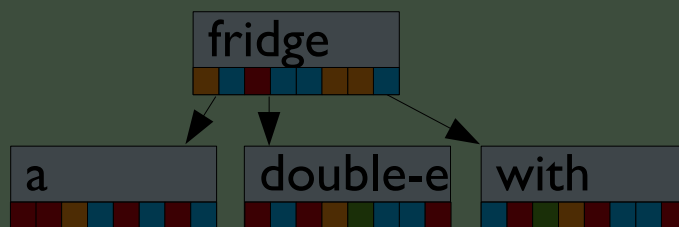
Quizbowl  
(Incremental  
Question  
Answering)

Ich bin mit dem Zug nach Ulm gefahren  
I am with the train to Ulm traveled  
I



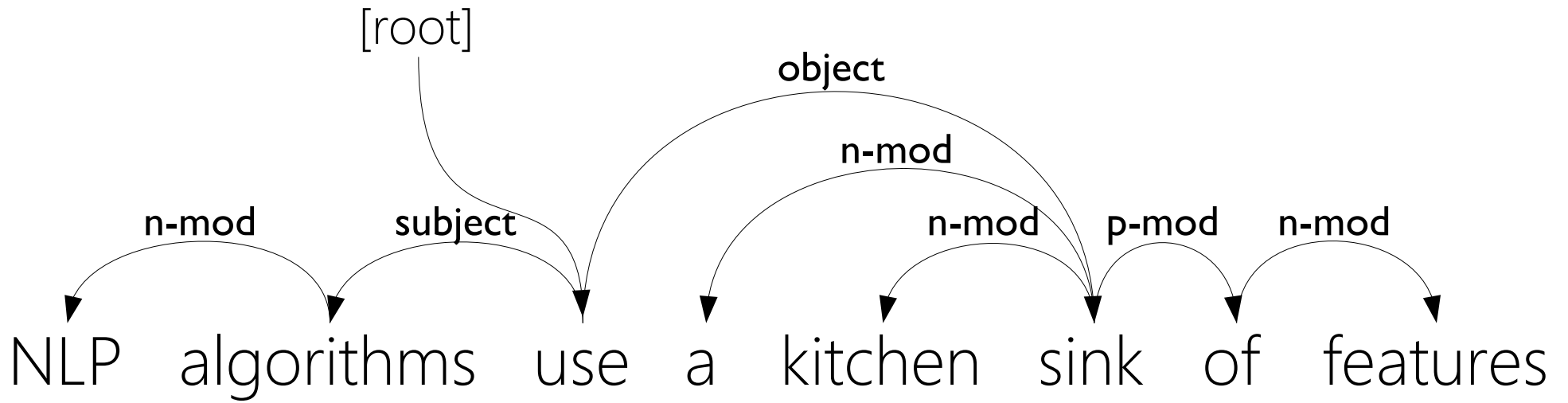
traveled by train to Ulm

Simultaneous  
Machine  
Interpretation

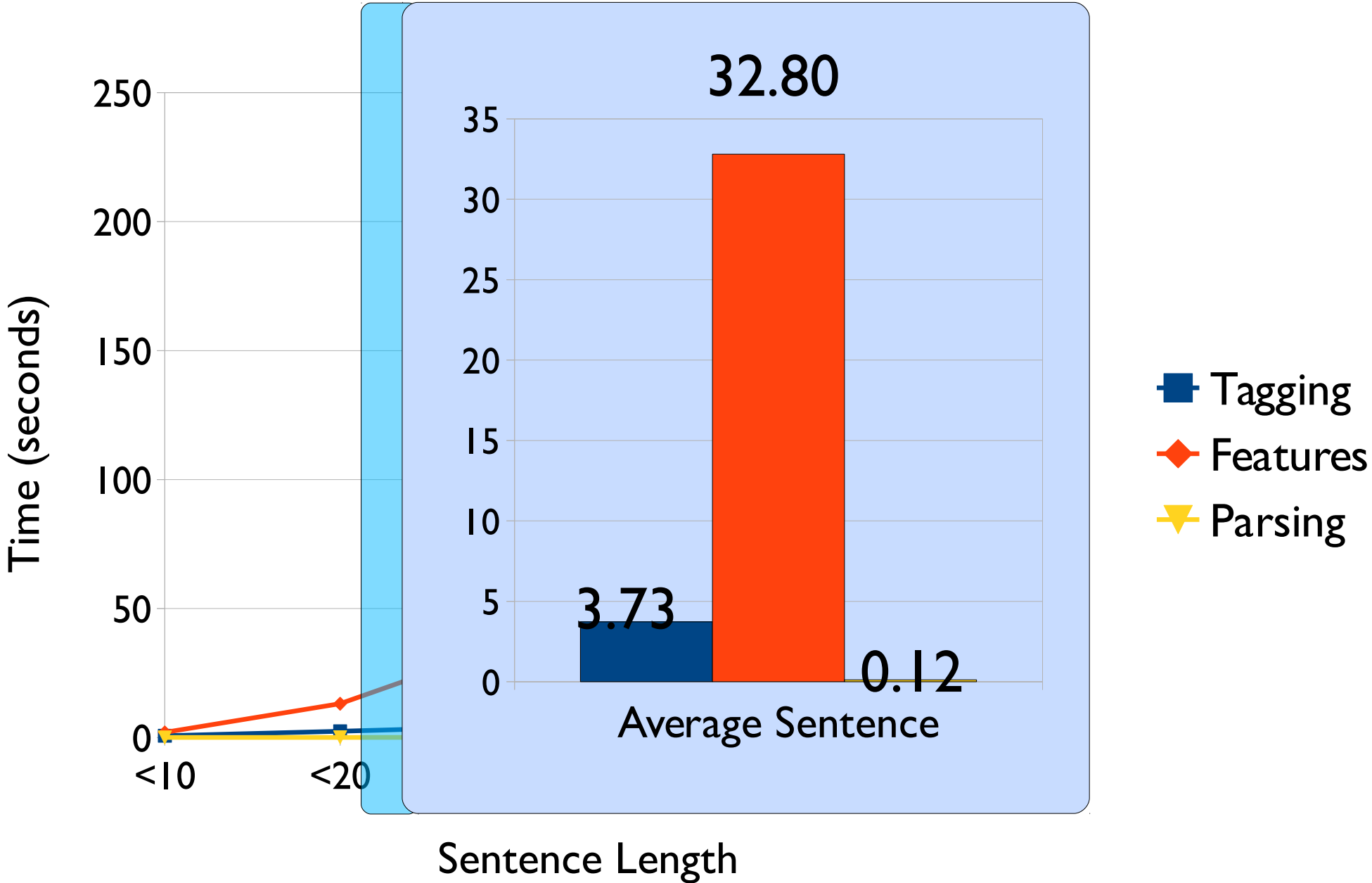


# Outline

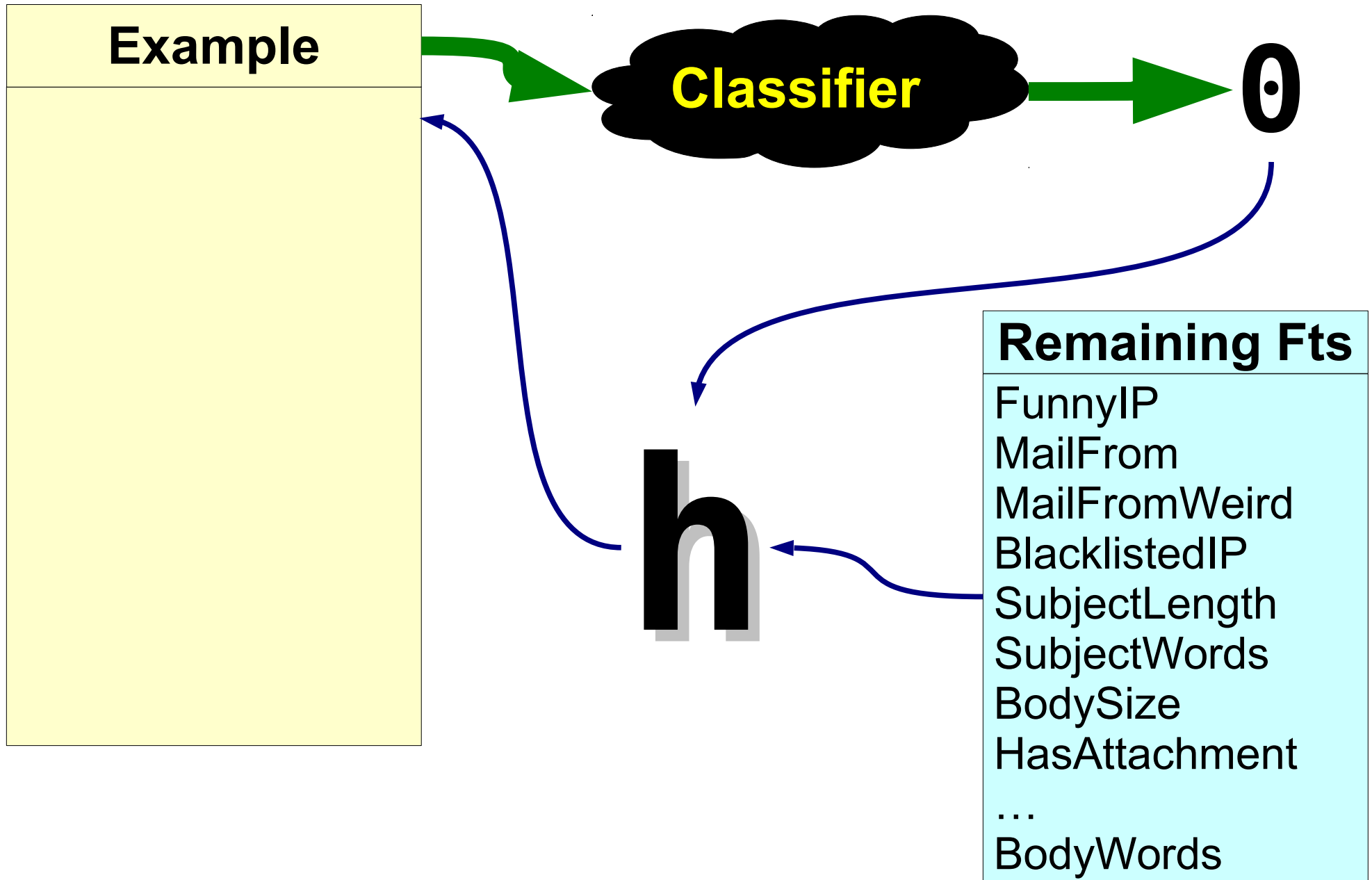
# Dependency parsing



# Case study: dependency parsing

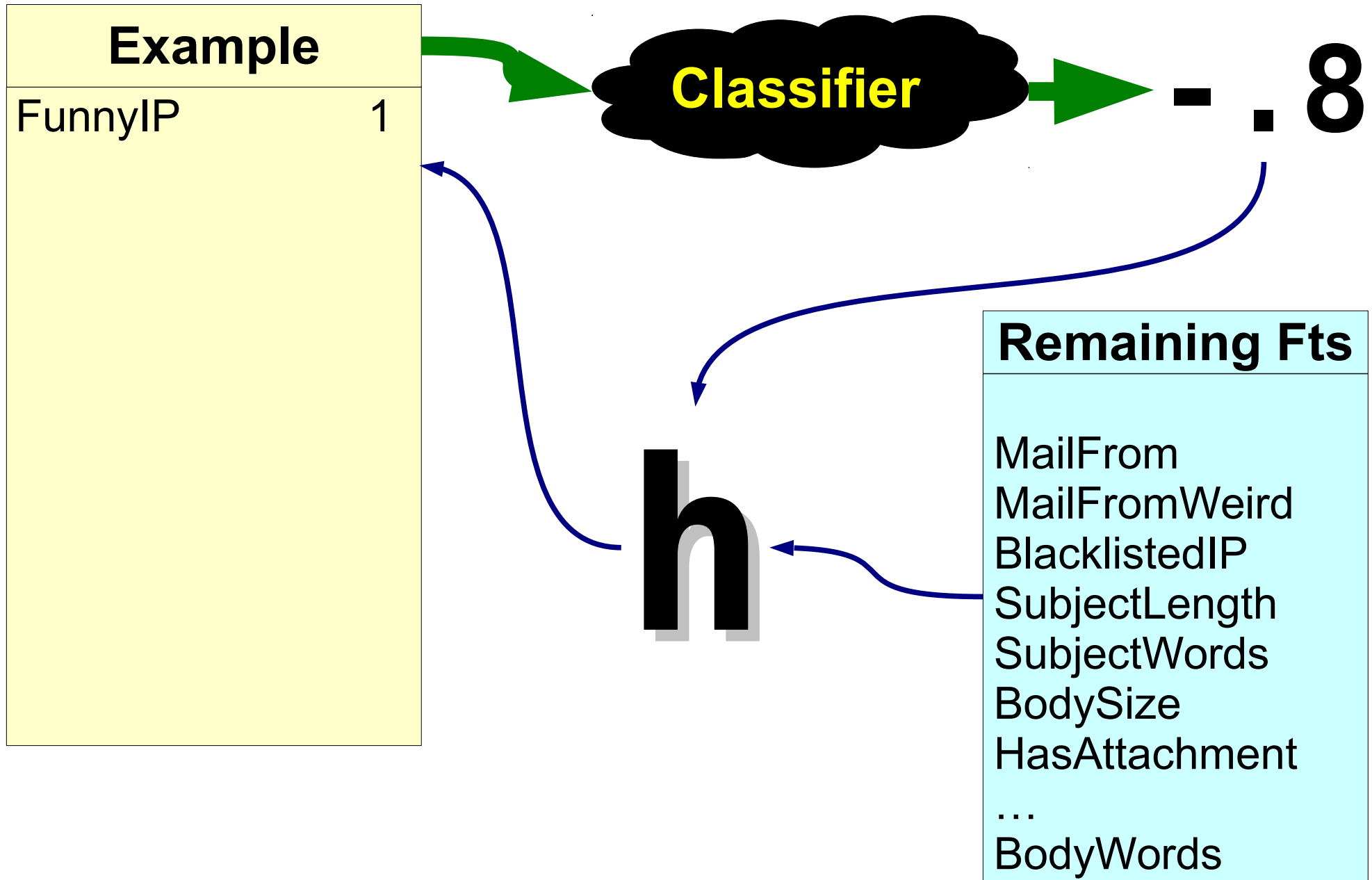


# Dynamic feature selection

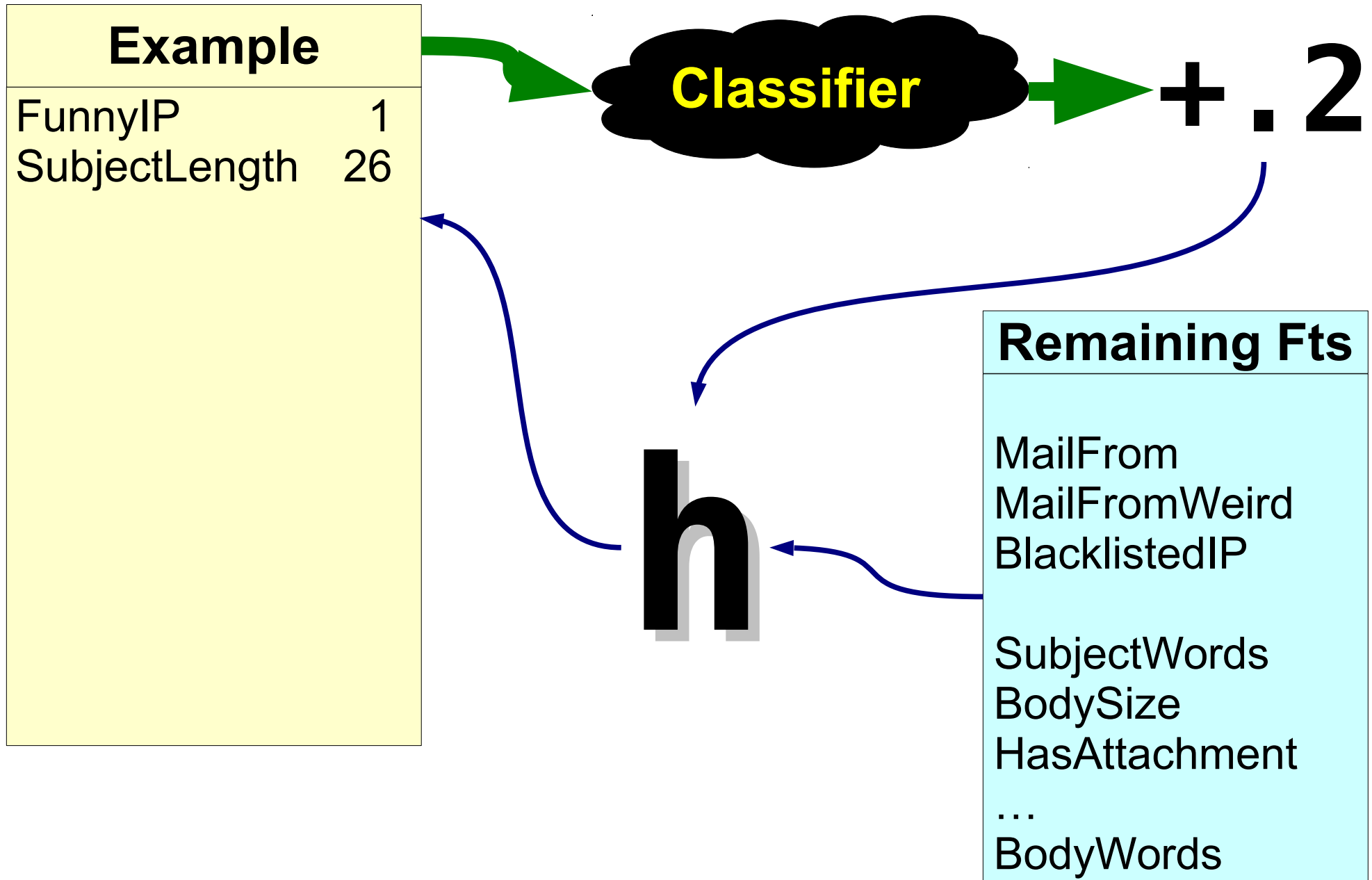




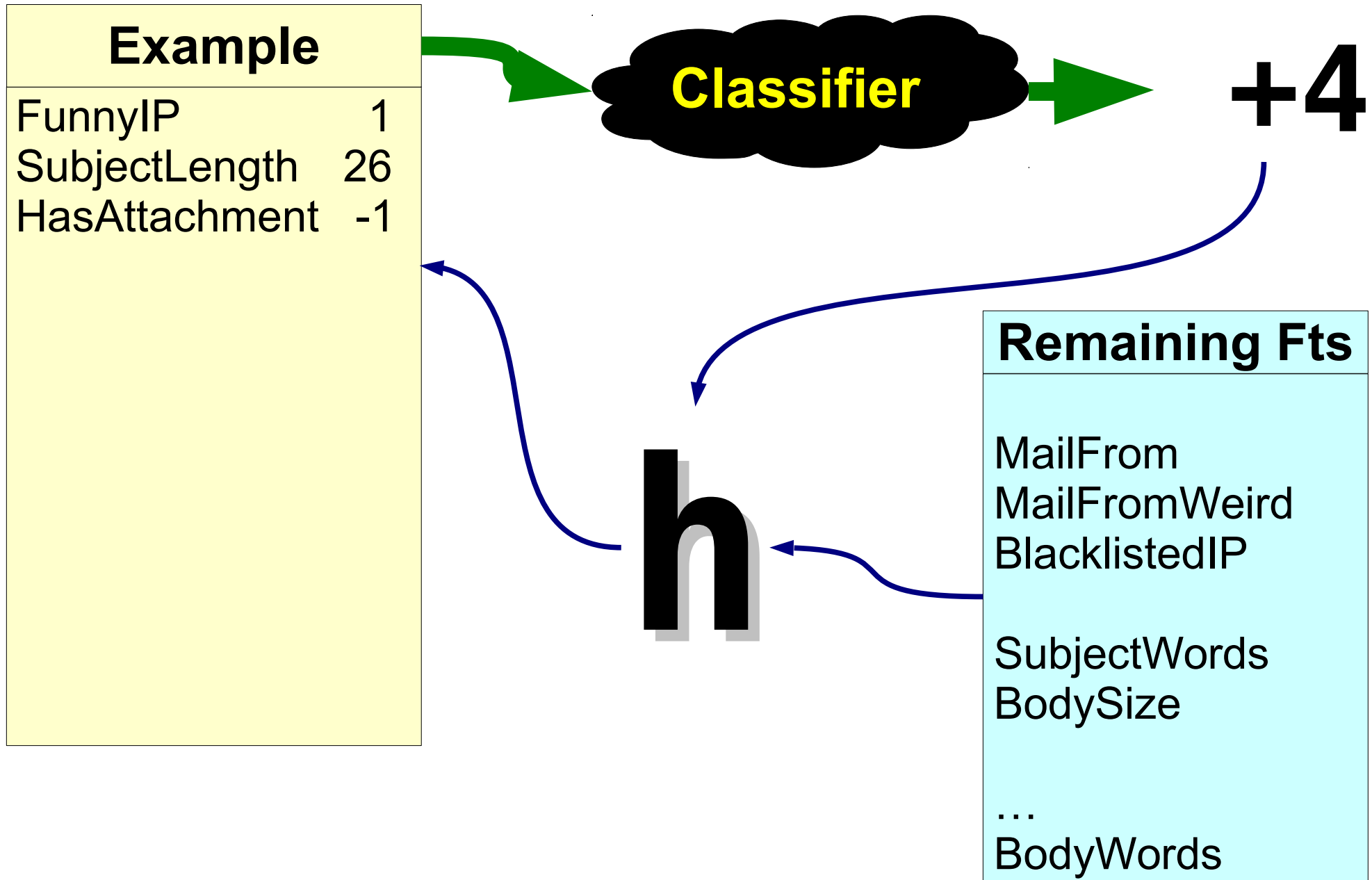
# Dynamic feature selection



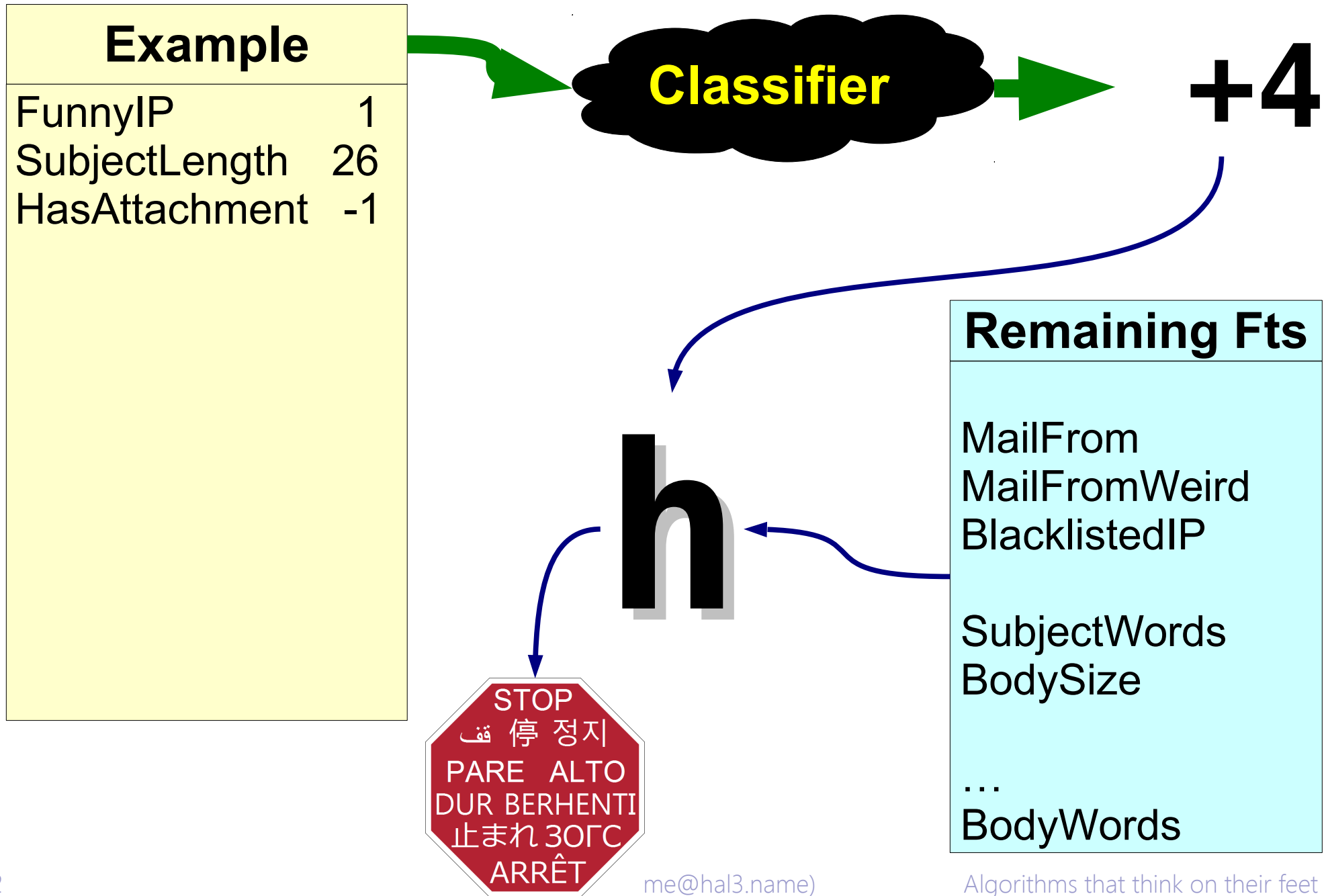
# Dynamic feature selection



# Dynamic feature selection



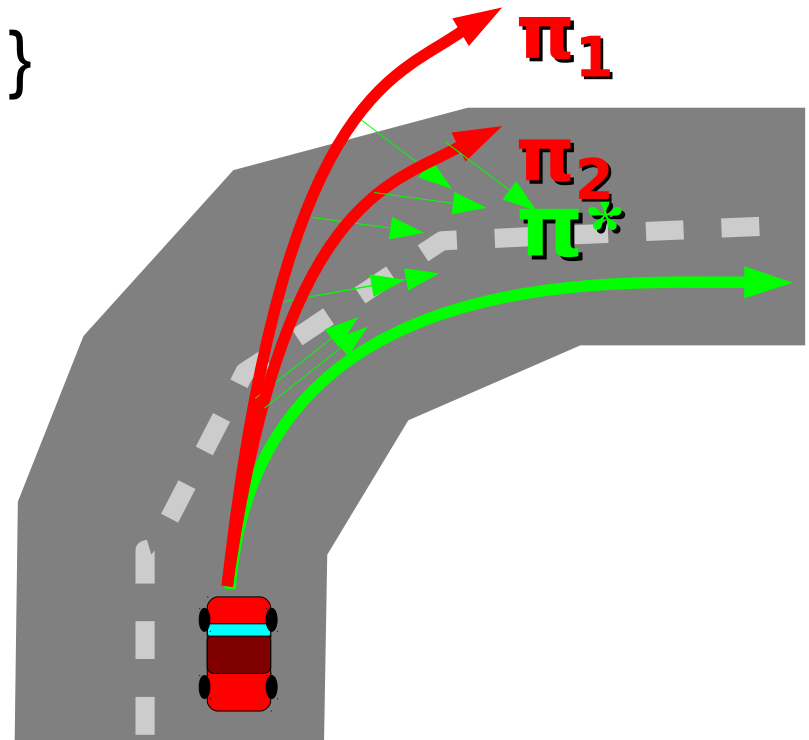
# Dynamic feature selection



# DAgger: Dataset Aggregation

- Collect trajectories from expert  $\pi^*$
- Dataset  $\mathbf{D}_0 = \{ (s, \pi^*(s)) \mid s \sim \pi^* \}$
- Train  $\pi_1$  on  $\mathbf{D}_0$
- Collect new trajectories from  $\pi_1$ 
  - But let the *expert* steer!
- Dataset  $\mathbf{D}_1 = \{ (s, \pi^*(s)) \mid s \sim \pi_1 \}$
- Train  $\pi_2$  on  $\mathbf{D}_0 \cup \mathbf{D}_1$
- In general:
  - $\mathbf{D}_n = \{ (s, \pi^*(s)) \mid s \sim \pi_n \}$
  - Train  $\pi_n$  on  $\mathbf{U}_{i < n} \mathbf{D}_i$

If  $N = T \log T$ ,  
 $L(\pi_n) < T \epsilon_N + O(1)$   
for some  $n$



# The oracle too good!

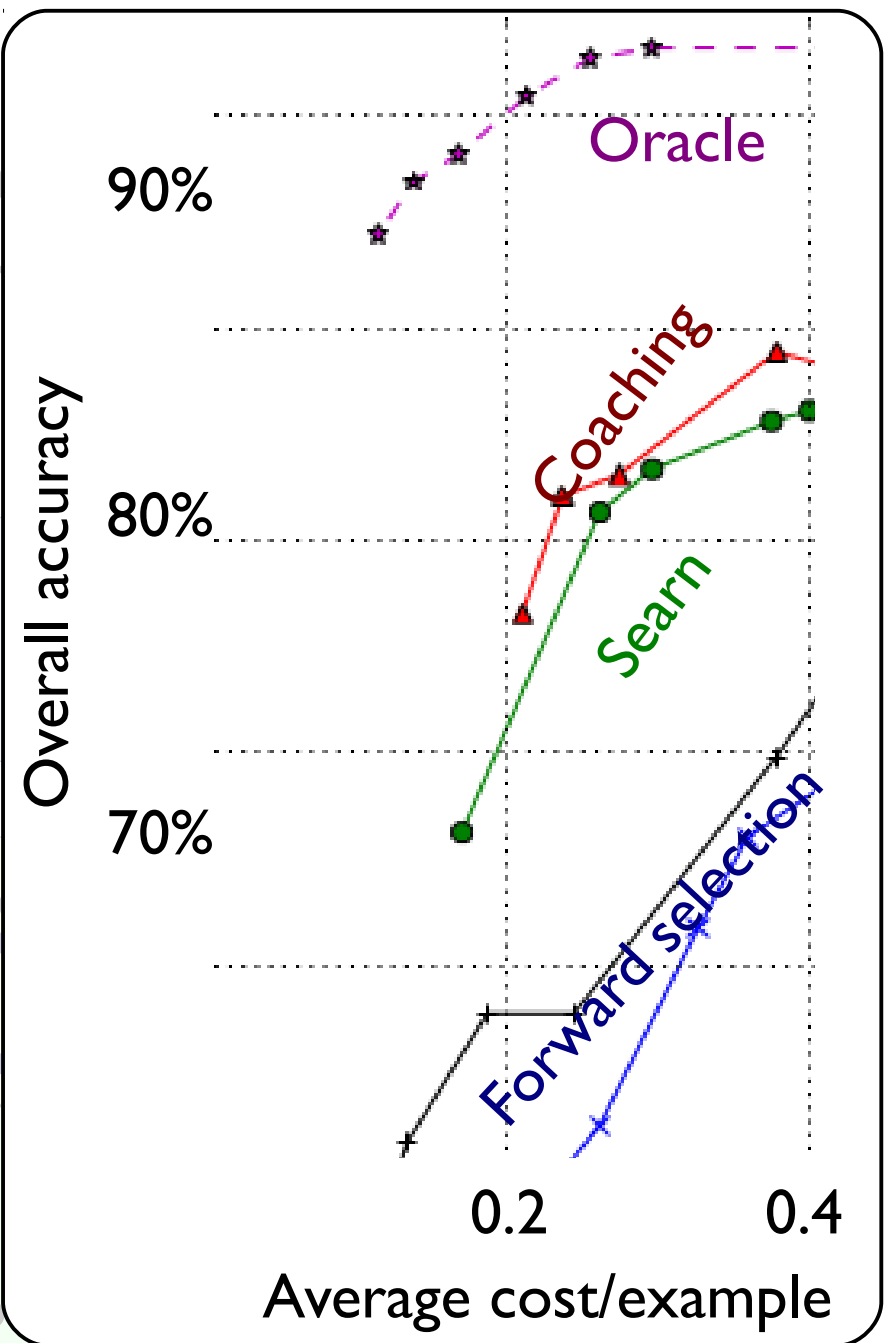
- The oracle *knows the label*
  - Picks feature with highest  $y^* \cdot v_a$
  - Ends after selecting one feature
- Coach says how to improve,

Pssst! You should choose  
 $\operatorname{argmin}_a \mathbf{E}[l(a)]$



If  $N = T \log$   
 $L(\pi_n) < T \epsilon_N$   
 for some

Provably smaller  
 oracle's epsilon

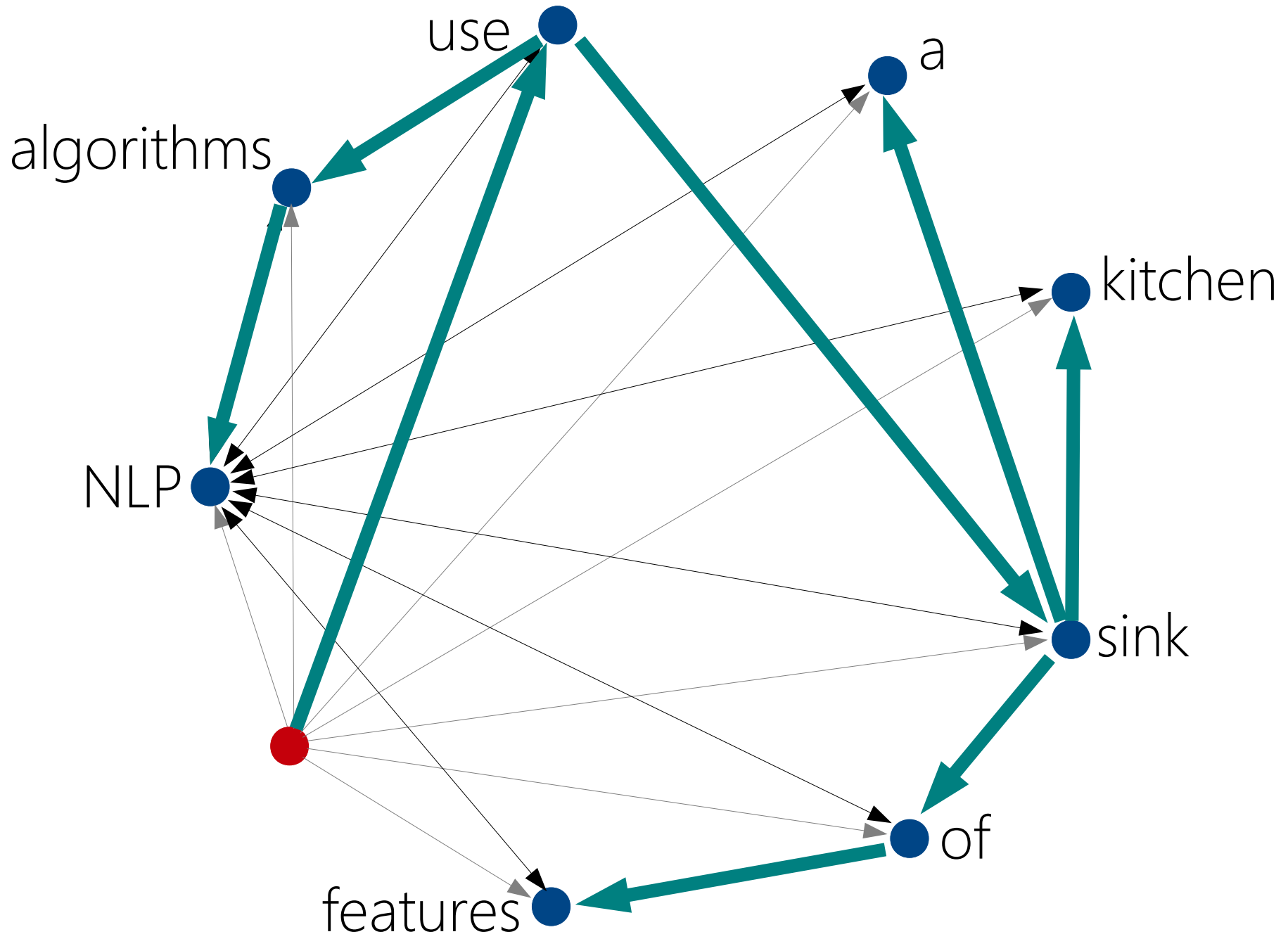


# Dependency parsing

[root]

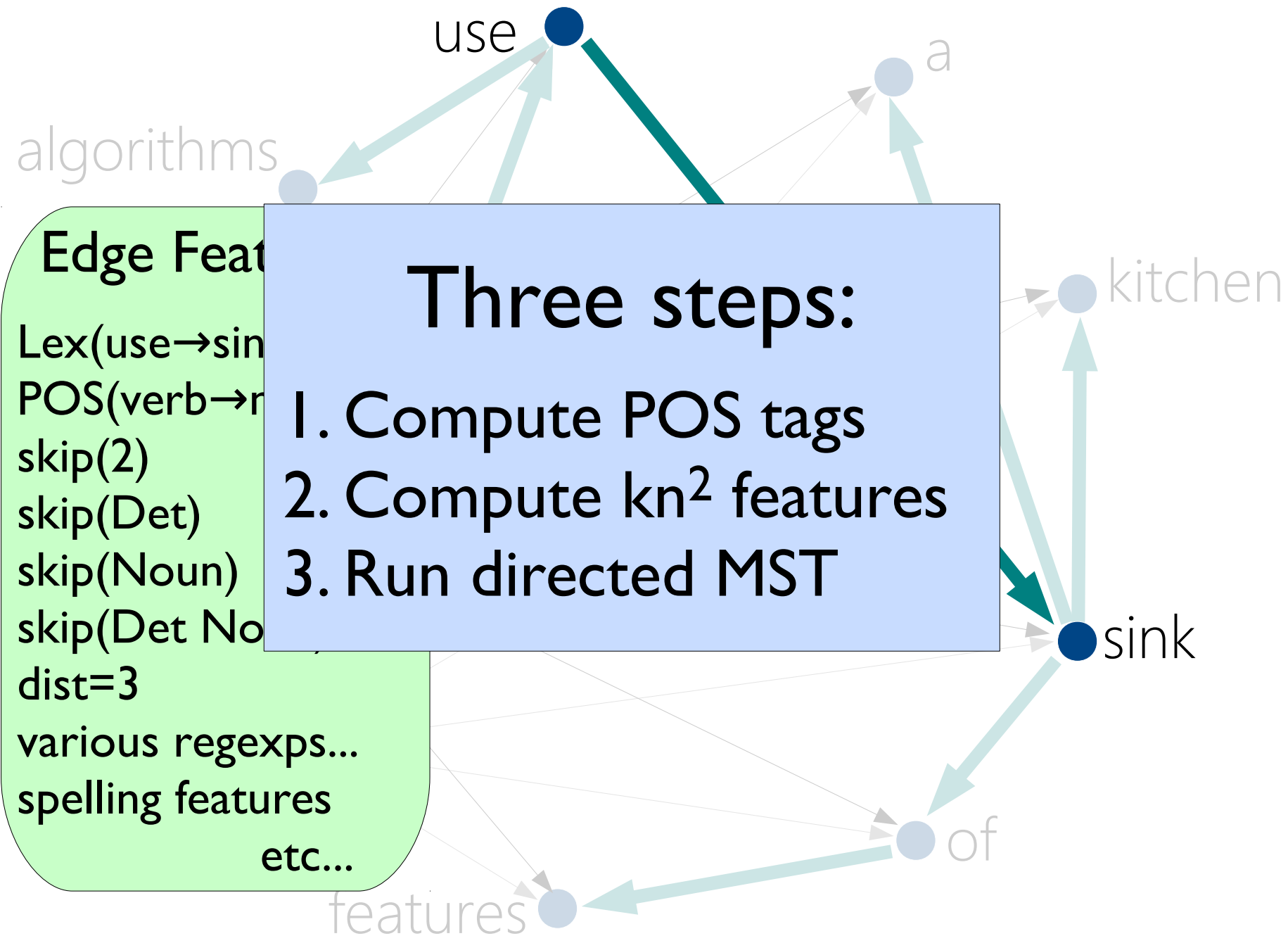
NLP algorithms use a kitchen sink of features

# Dependency parsing





# Dependency parsing

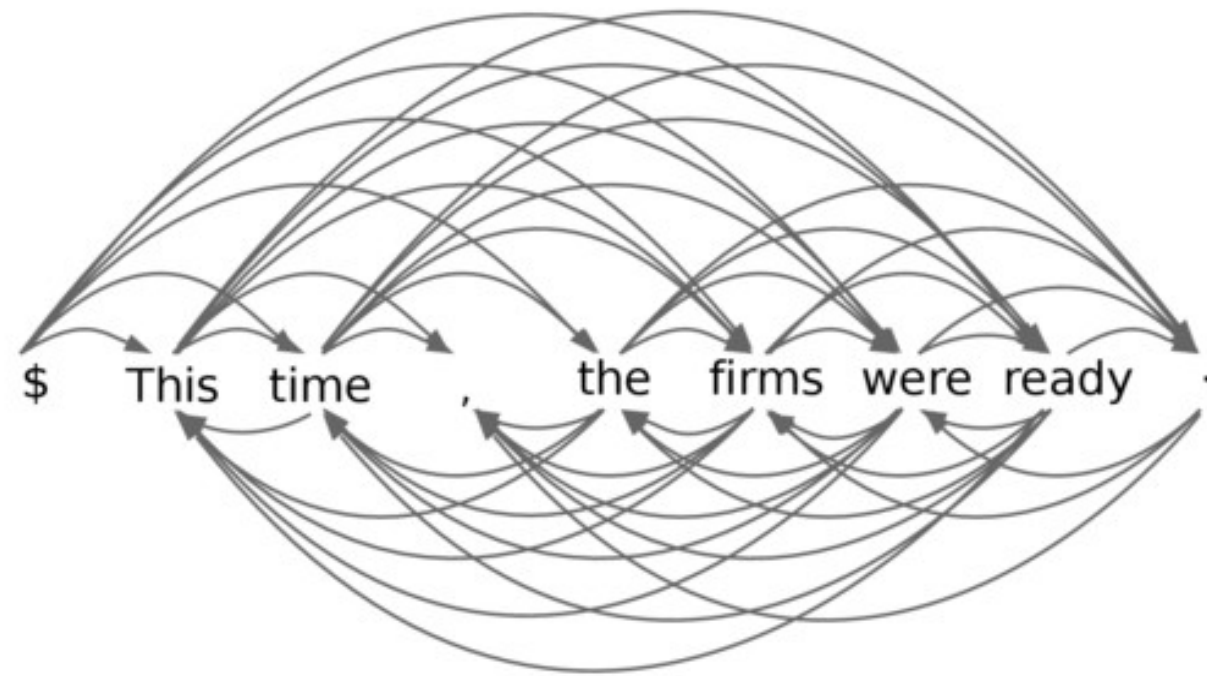


# The system we learn to control

+ first feature group

5  
51

features per gray edge  
gray edge with unknown fate...

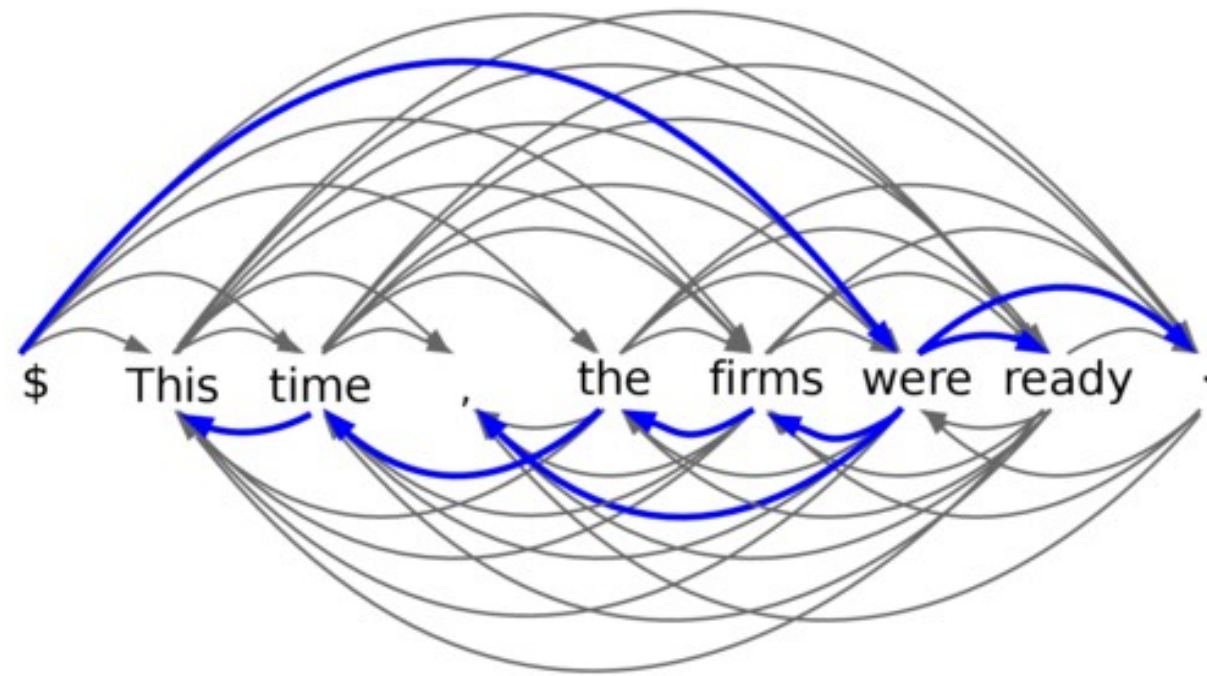


- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

# The system we learn to control



**5** features per gray edge  
**51** gray edge with unknown fate...



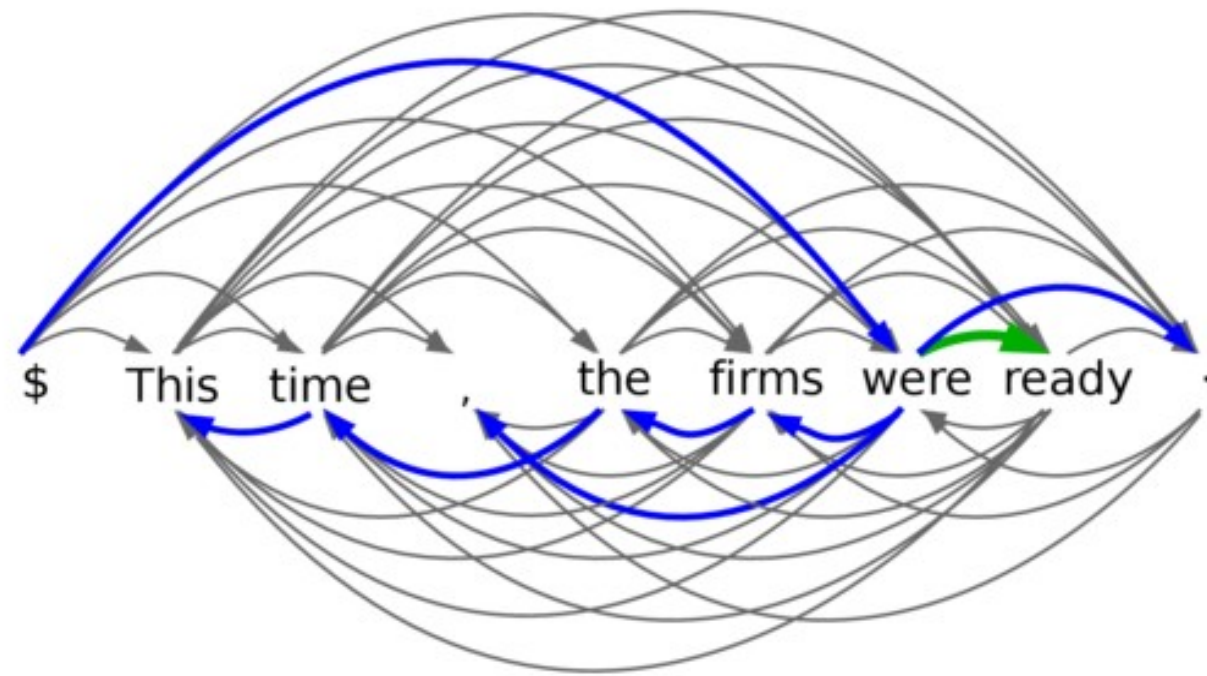
- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

*Non-projective decoding*

# The system we learn to control



5 features per gray edge  
50 gray edge with unknown fate...

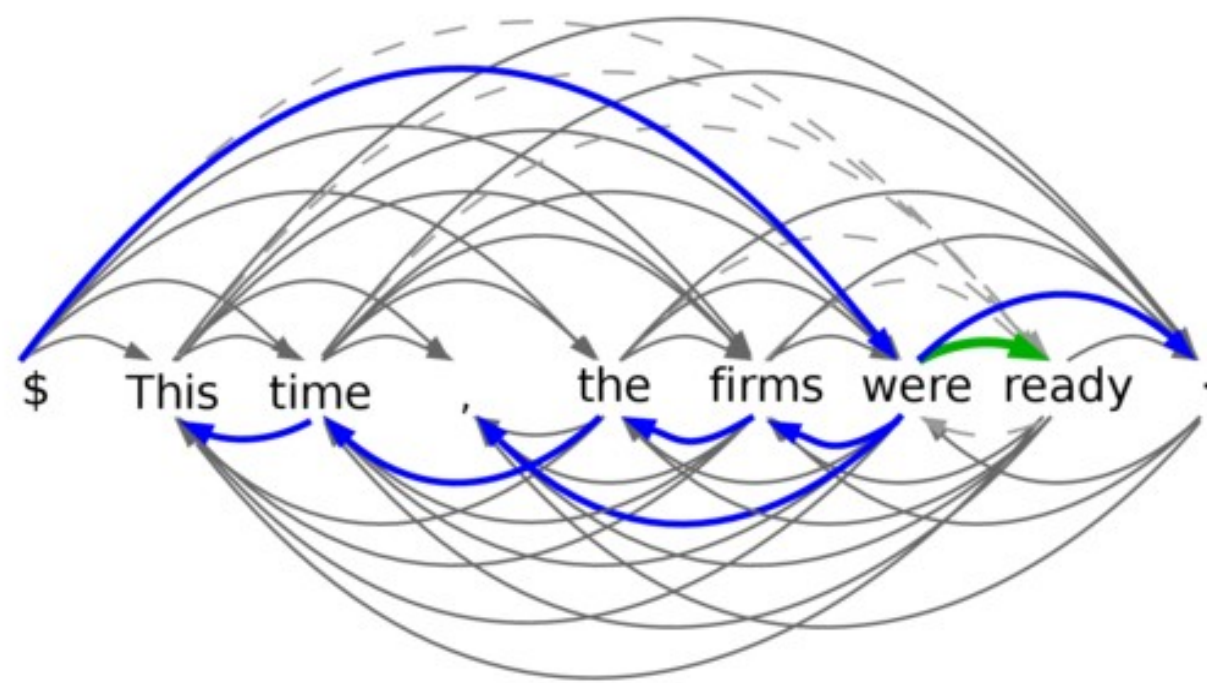


- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

Decide *winners* among the *blue edges*

# The system we learn to control

5 features per gray edge  
44 gray edge with unknown fate...



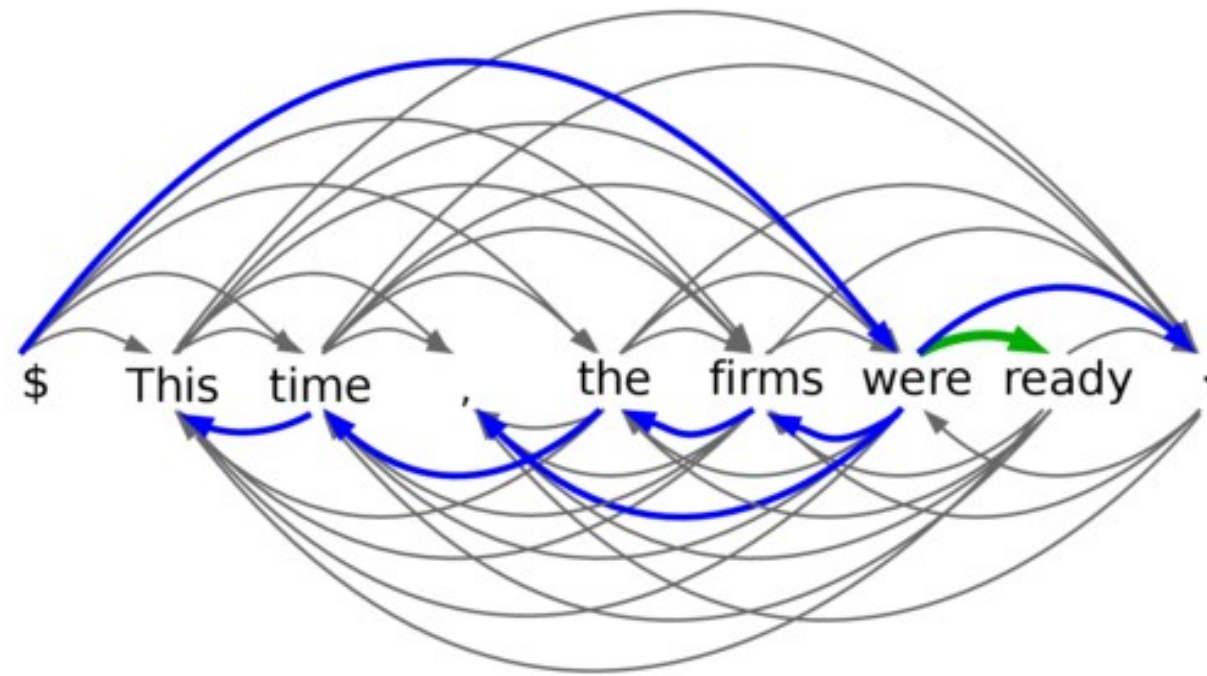
- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

*Remove losers in conflict with the winners*



# The system we learn to control

5 features per gray edge  
44 gray edge with unknown fate...



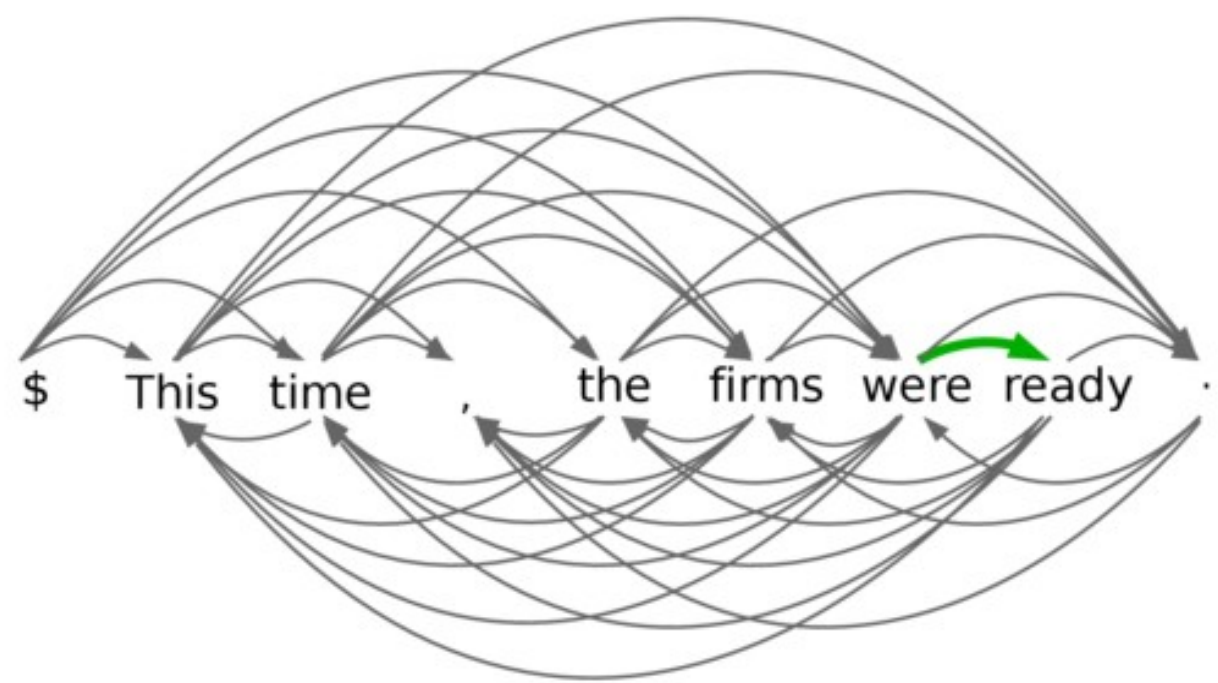
- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

*Remove losers in conflict with the winners*

# The system we learn to control

+ next feature group **27**

features per gray edge  
gray edge with unknown fate...

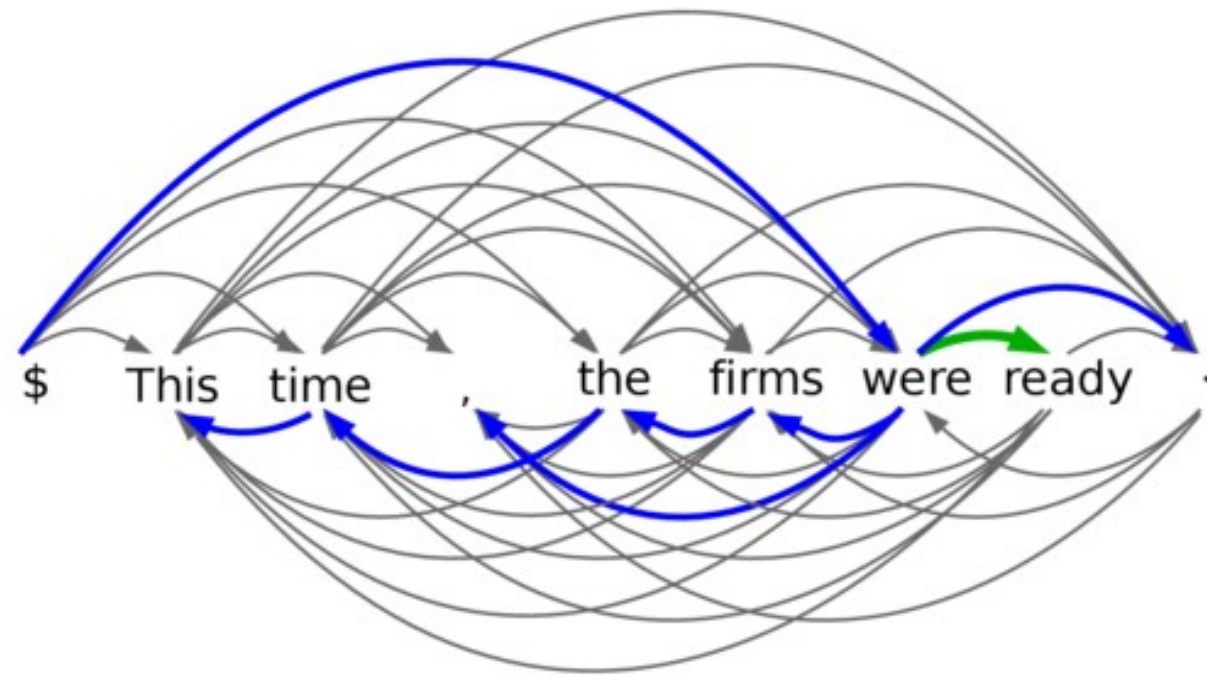


- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

# The system we learn to control

+ next feature group **27**  
**44**

features per gray edge  
gray edge with unknown fate...

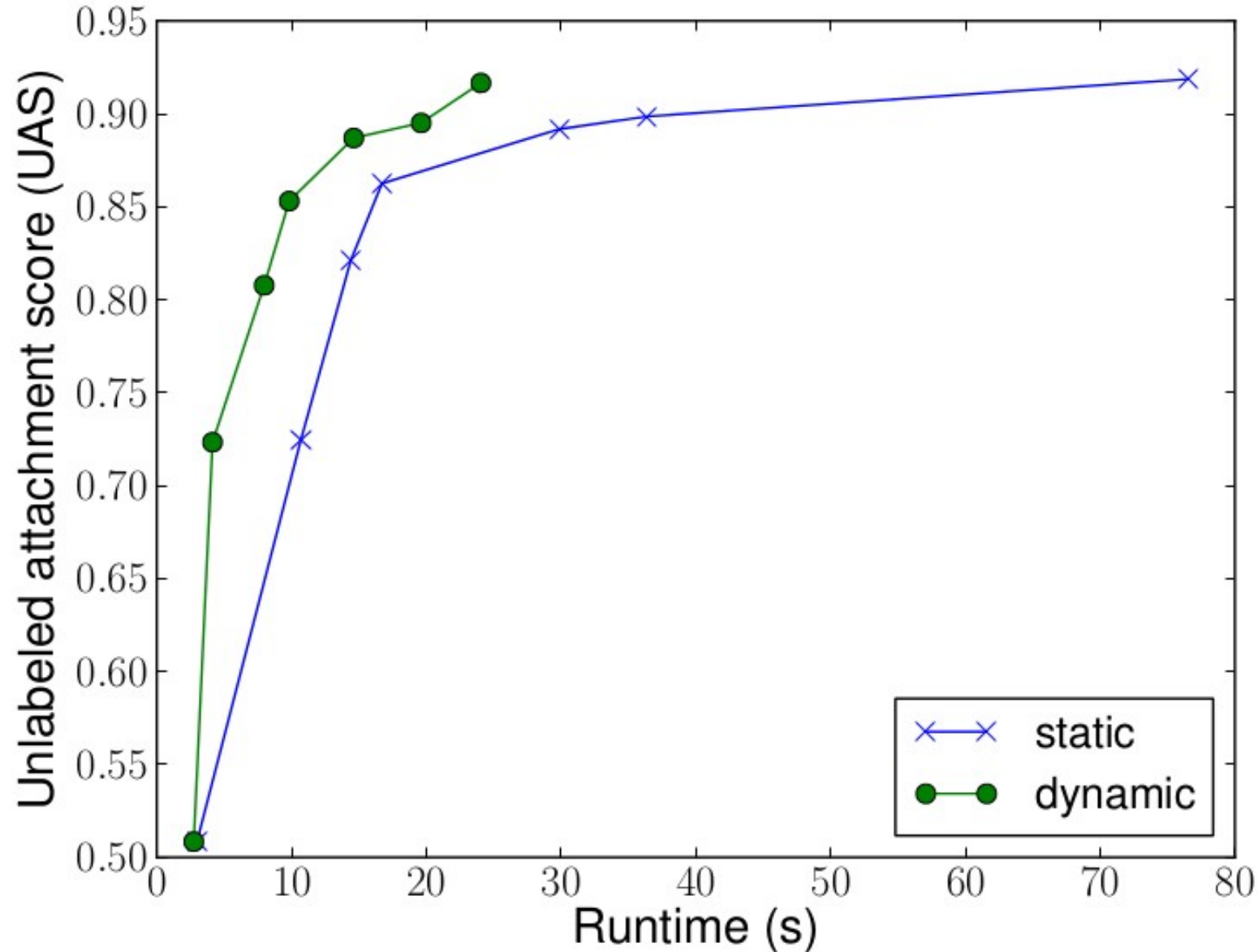


- Undetermined edge
- Current 1-best tree
- Winner edge
- - Loser edge

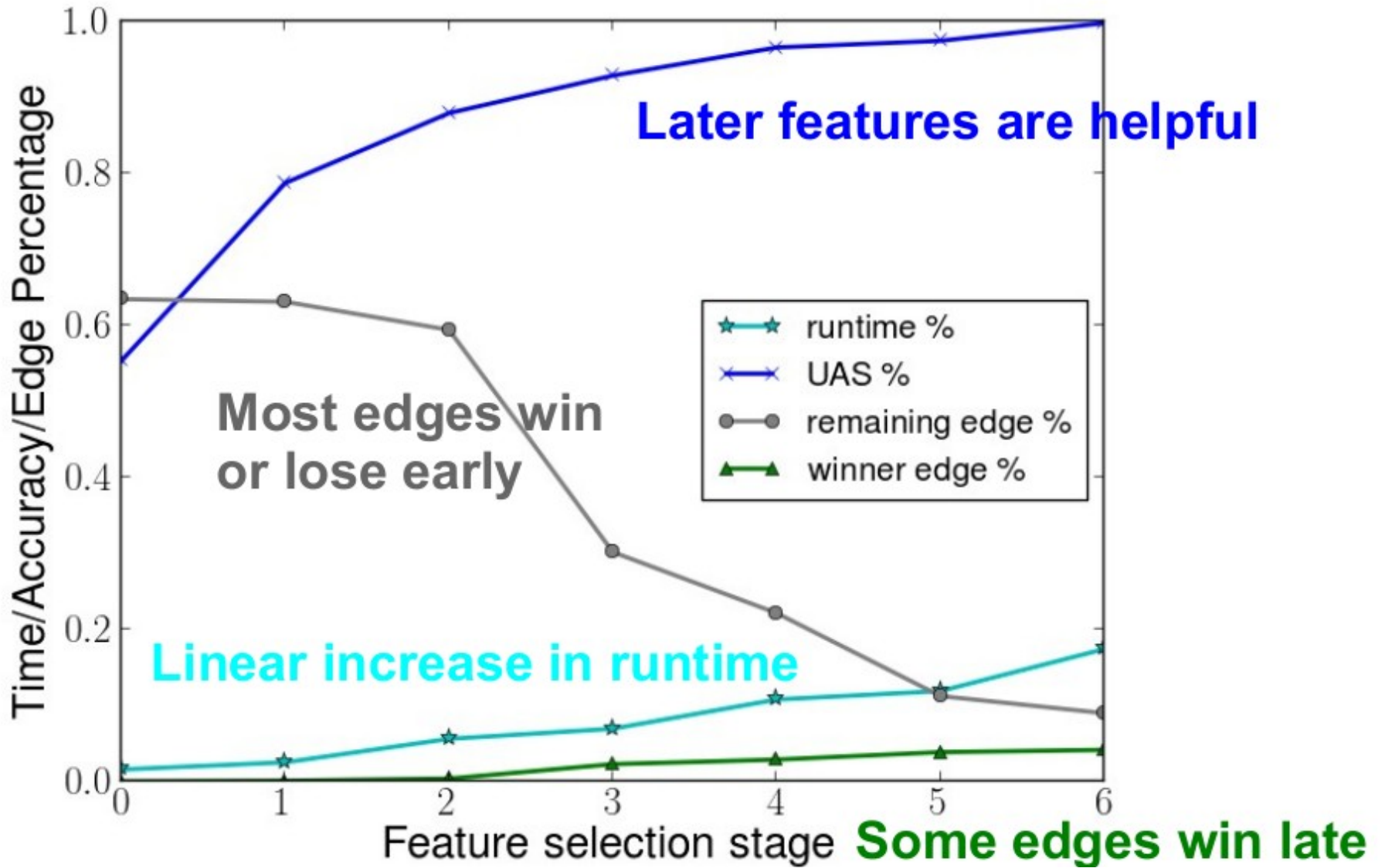
*Non-projective decoding*



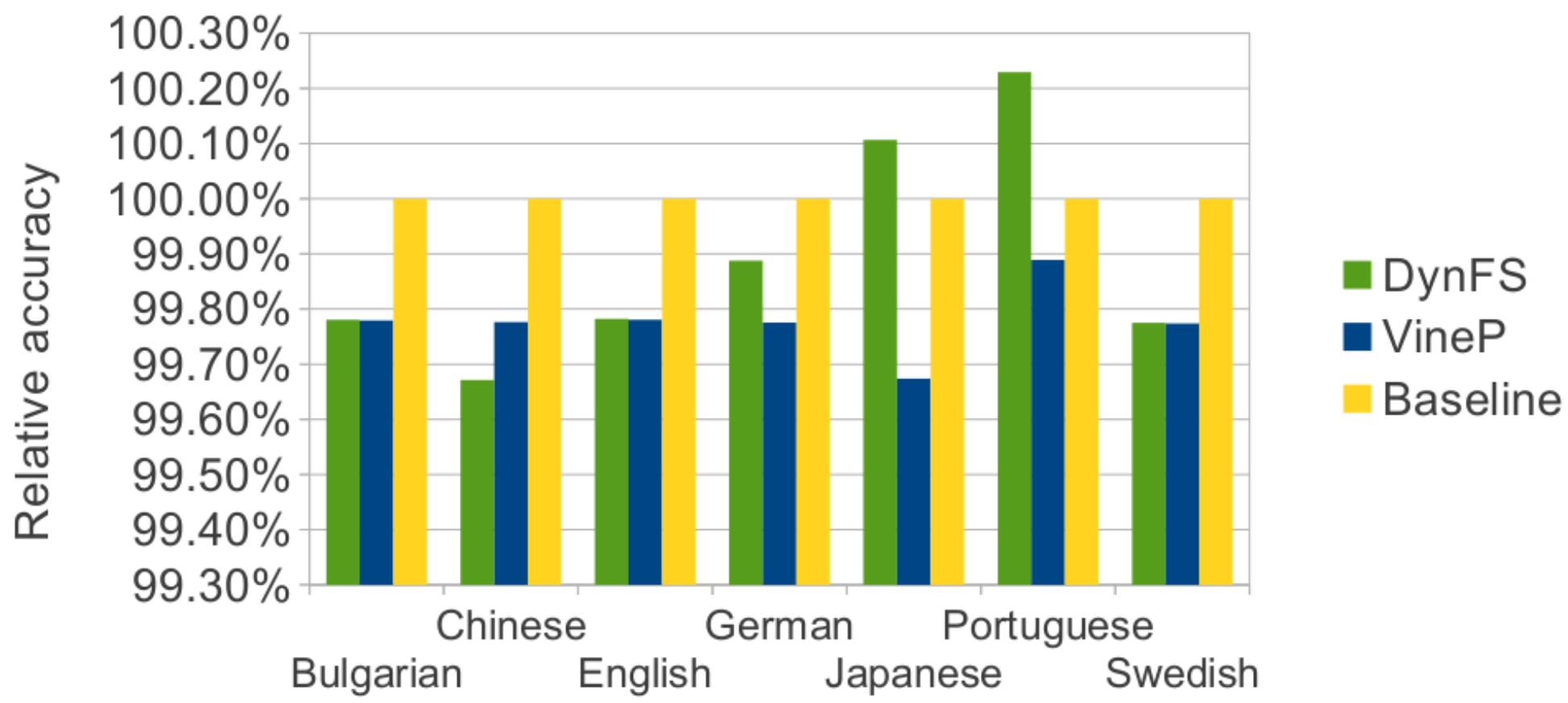
# Static versus dynamic feature selection



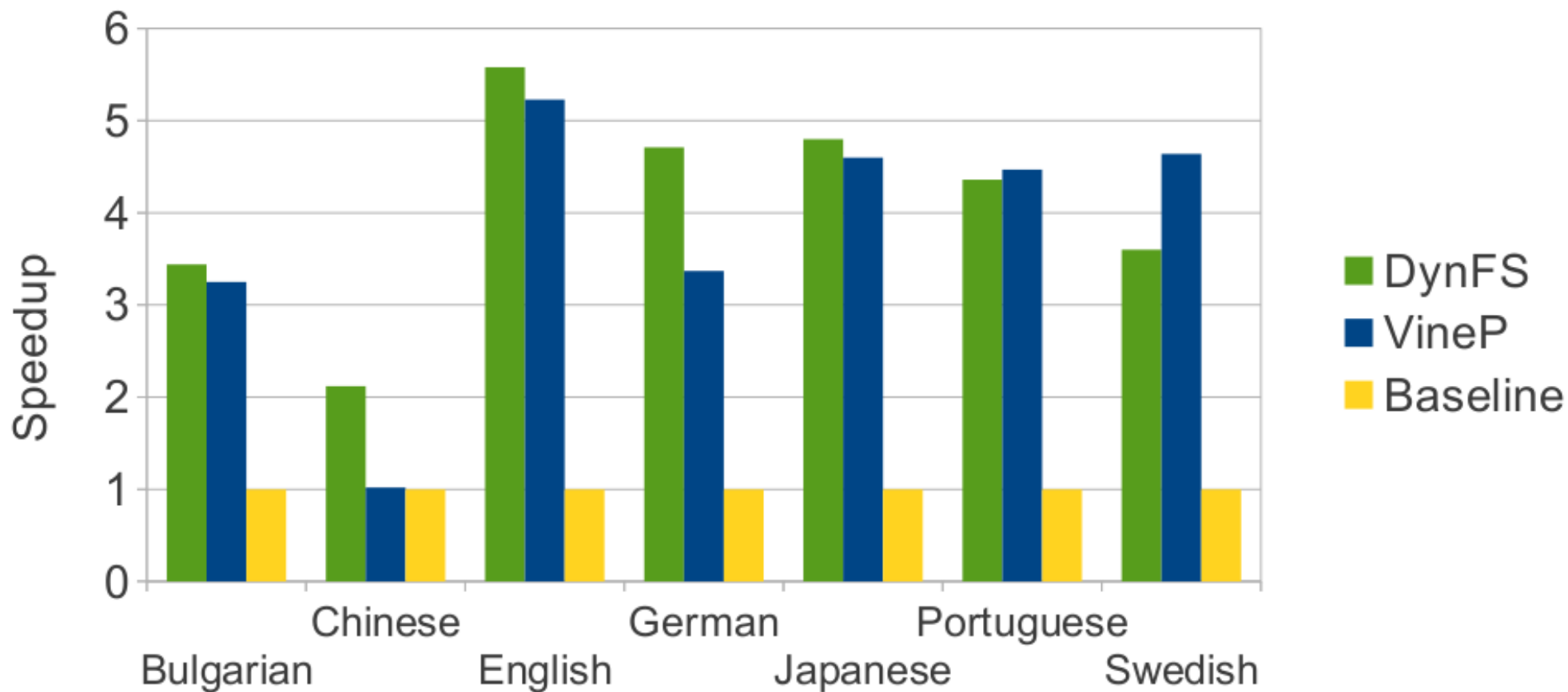
# Looking inside the box...



# Accuracy is (essentially) unaffected

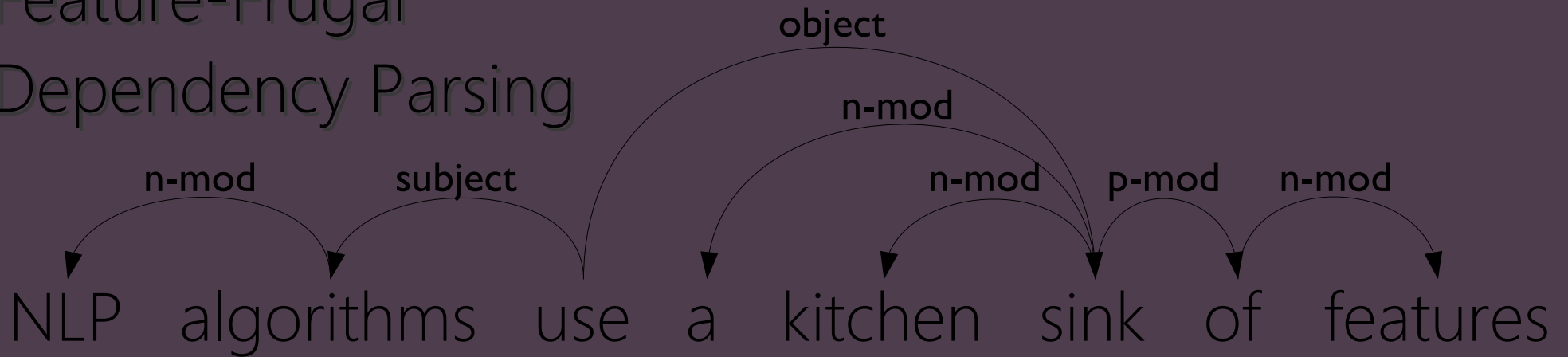


# ...but we get a lot faster



(He+Eisner+D, EMNLP 2013)

# Feature-Frugal Dependency Parsing



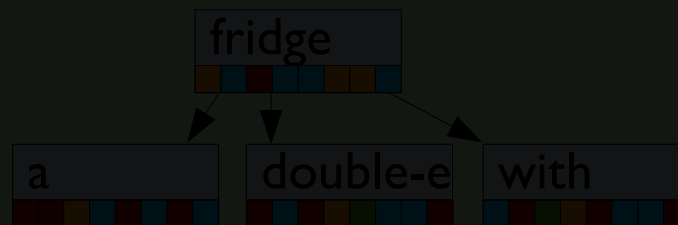
Quizbowl  
(Incremental  
Question  
Answering)

Ich bin mit dem Zug nach Ulm gefahren  
I am with the train to Ulm traveled  
I



traveled by train to Ulm

Simultaneous  
Machine  
Interpretation



# Outline

# Simultaneous (machine) interpretation



## Nuremberg Trials

- Dozens of defendants
- Judges from four nations (three languages)
- Status quo: speak, then translate
- After Nuremberg, simultaneous translations became the norm
- Long wait → bad conversation



# Why simultaneous interpretation is hard

- Human languages have vastly different word orders
  - About half are OV, the other half are VO
  - This comes with a lot more baggage than just verb-final

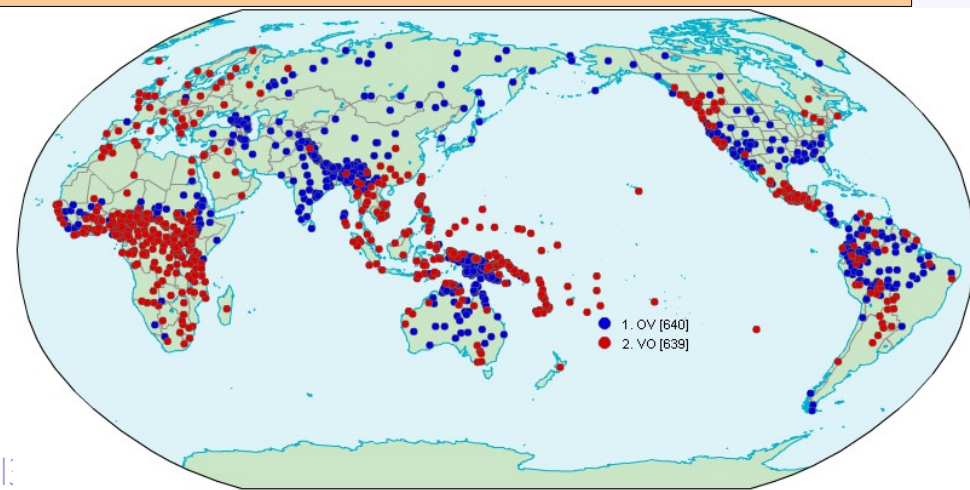
## Running (German/English) Example:

Ich bin mit dem Zug nach Ulm gefahren

I am with the train to Ulm traveled

I (..... *waiting*.....)

traveled by train to Ulm



# Model for interpretation decisions

- **We have a set of actions (predict / translate)**
  - Wait
  - Predict clause-verb
  - Predict next word
  - Commit (“speak”)
- **In a changing environment (state)**
  - The words we've seen so far
  - Our models' internal predictions
- **With a well defined oracle**



# Example of interpretation trajectory

Observation

1. Mit dem Zug

state

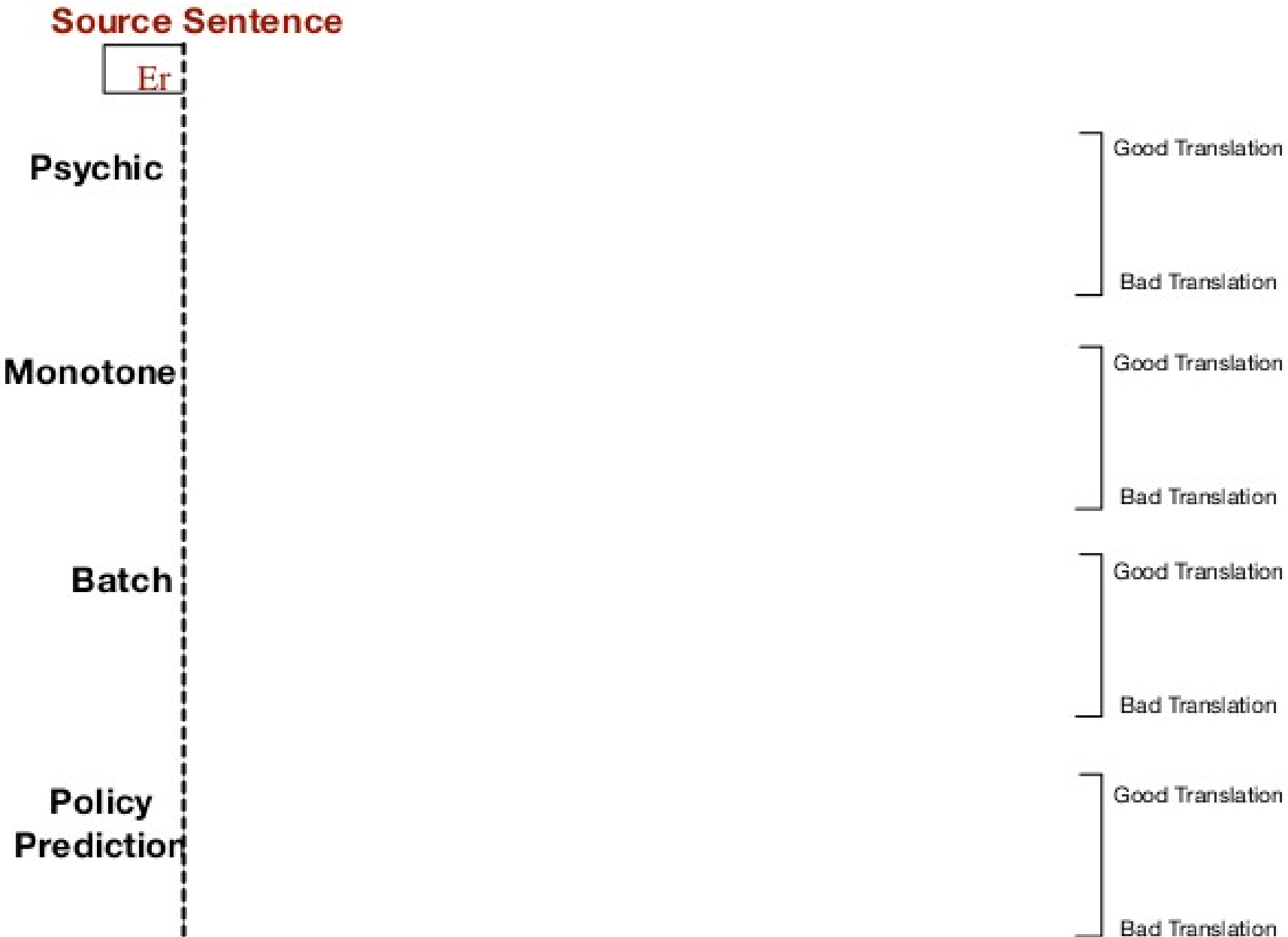
Verb: **gewesen**  
Next: **und**

Ich bin mit dem Zug nach Ulm gefahren

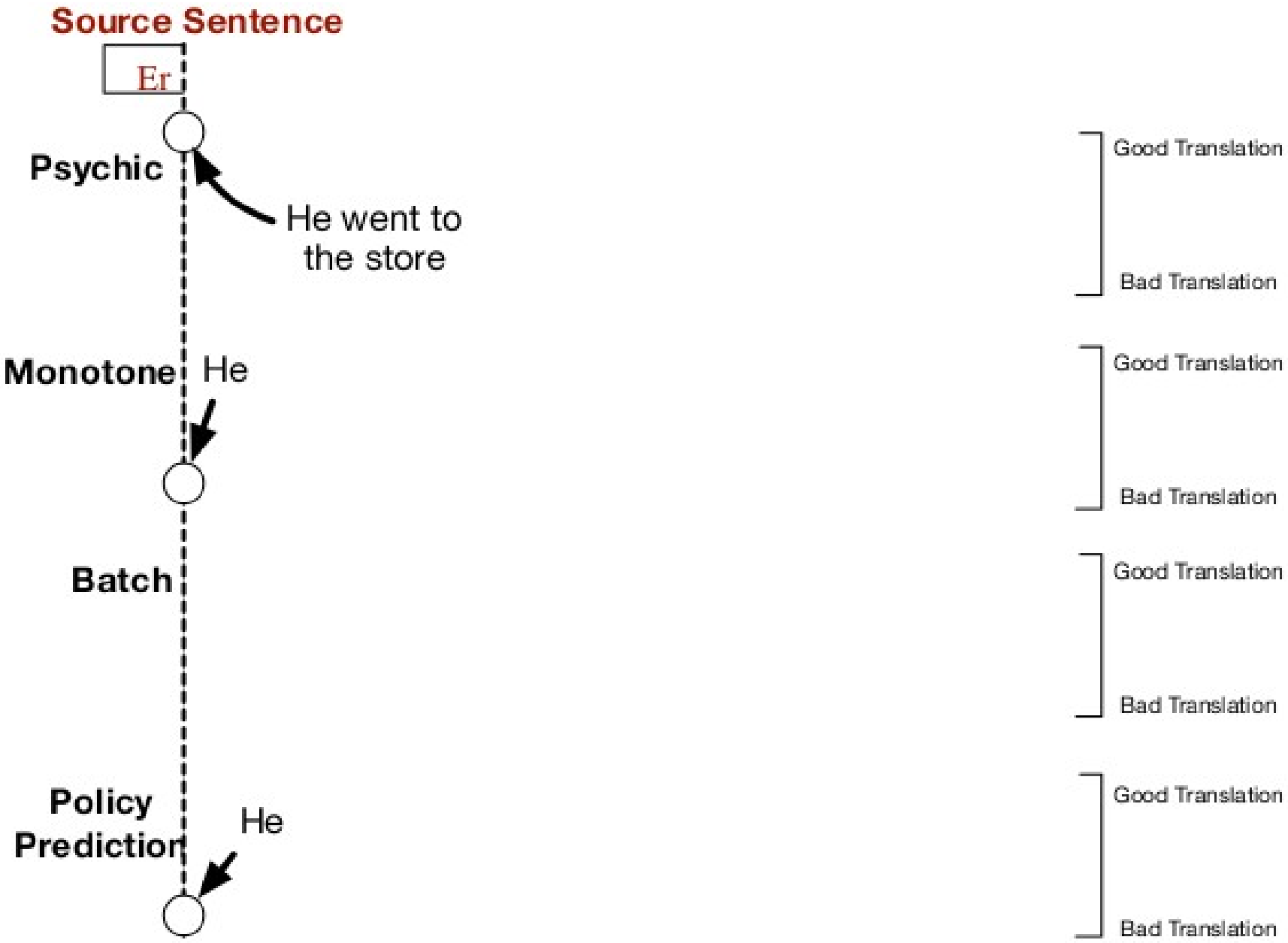
I am with the train to Ulm traveled

I (..... *waiting*.....) traveled by train to Ulm

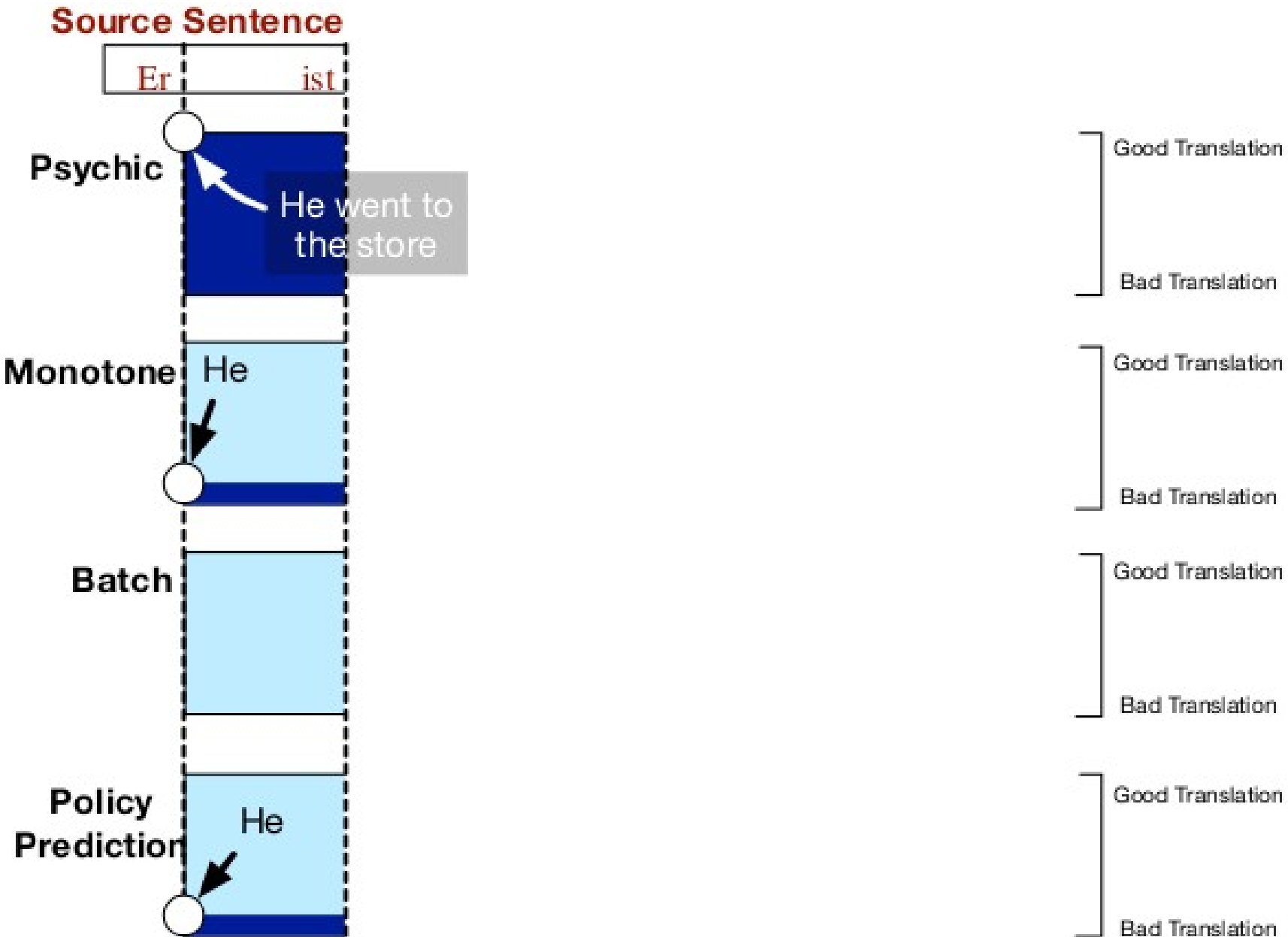
# Evaluating performance and baselines



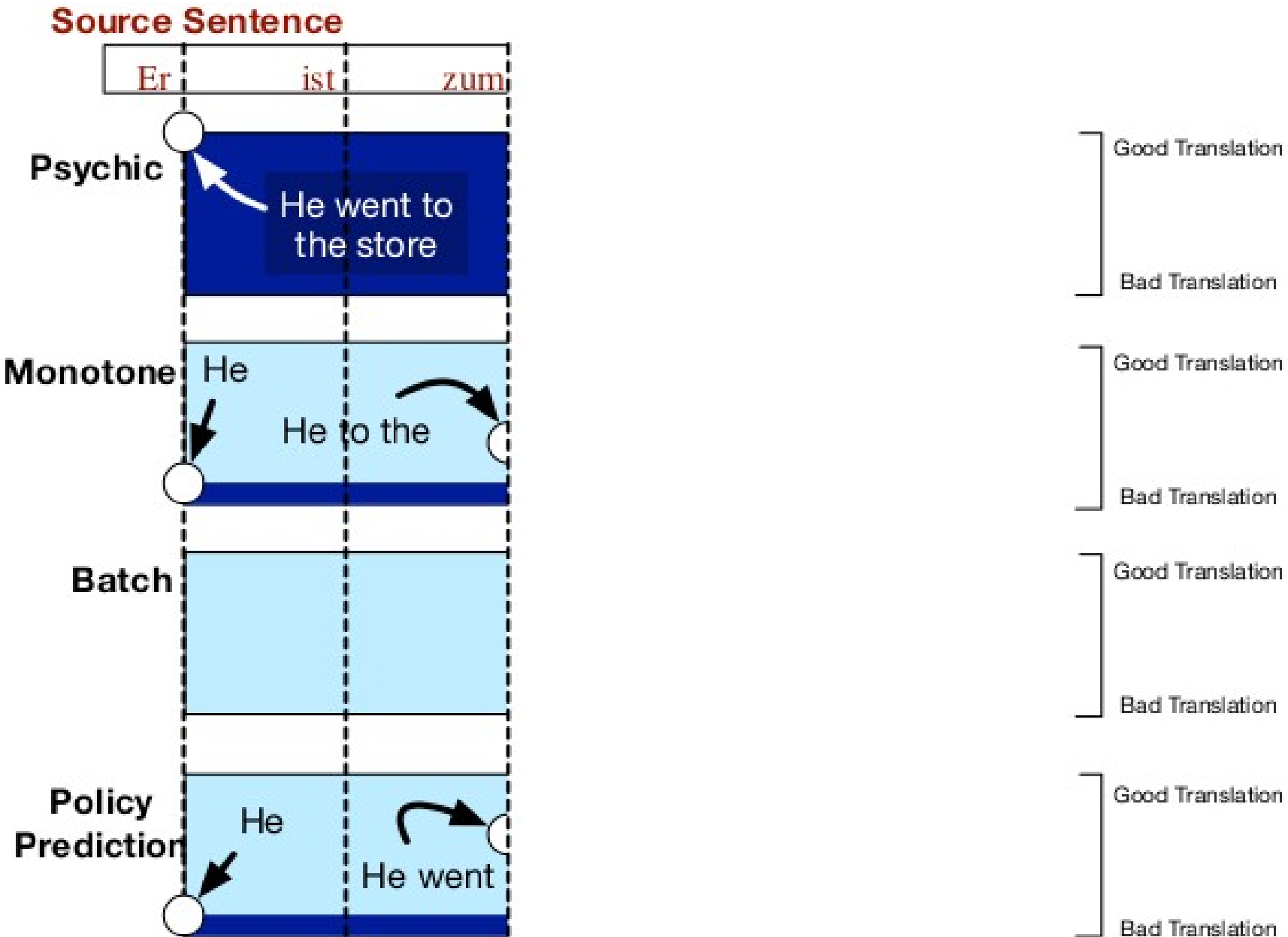
# Evaluating performance and baselines



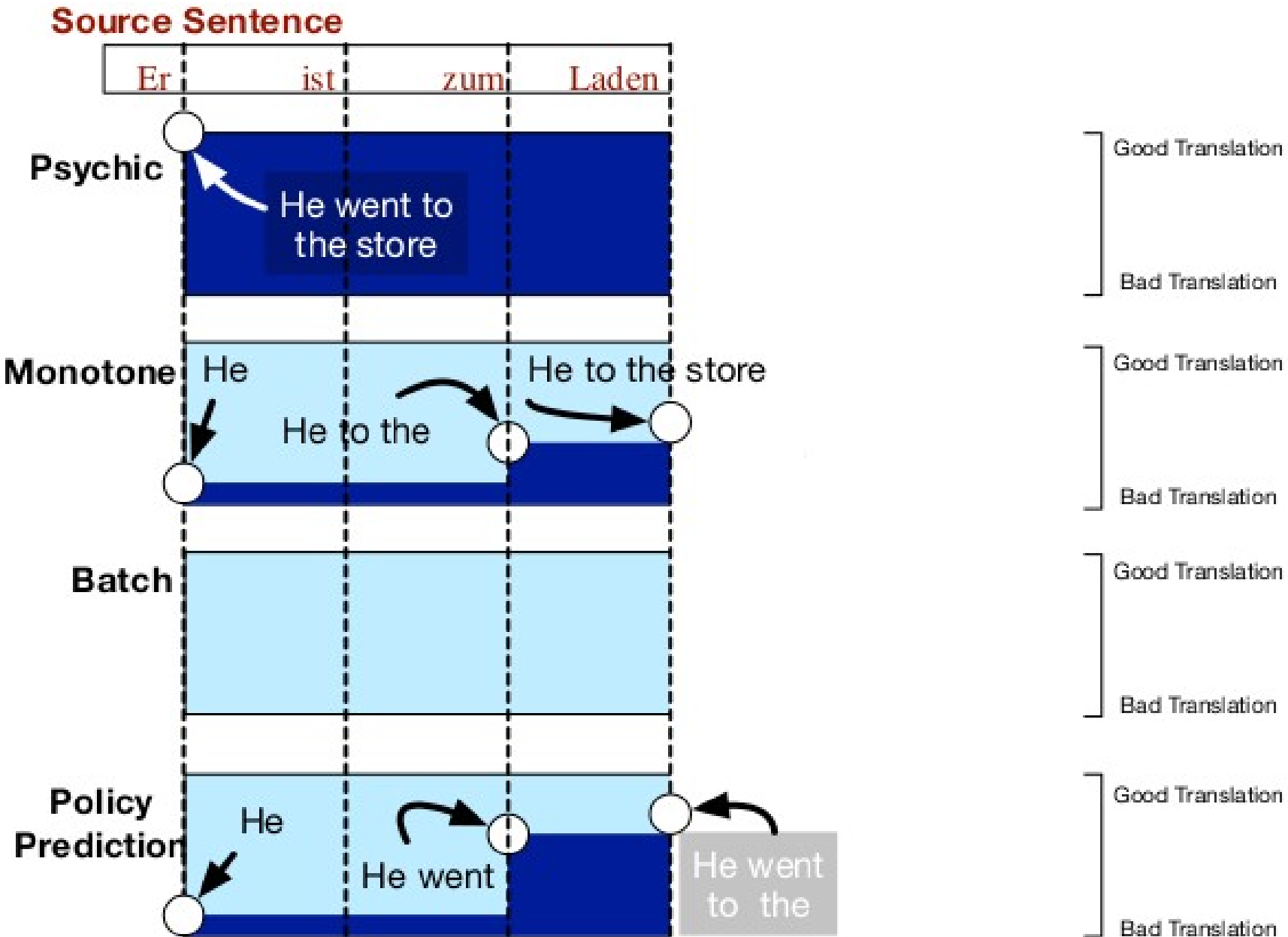
# Evaluating performance and baselines



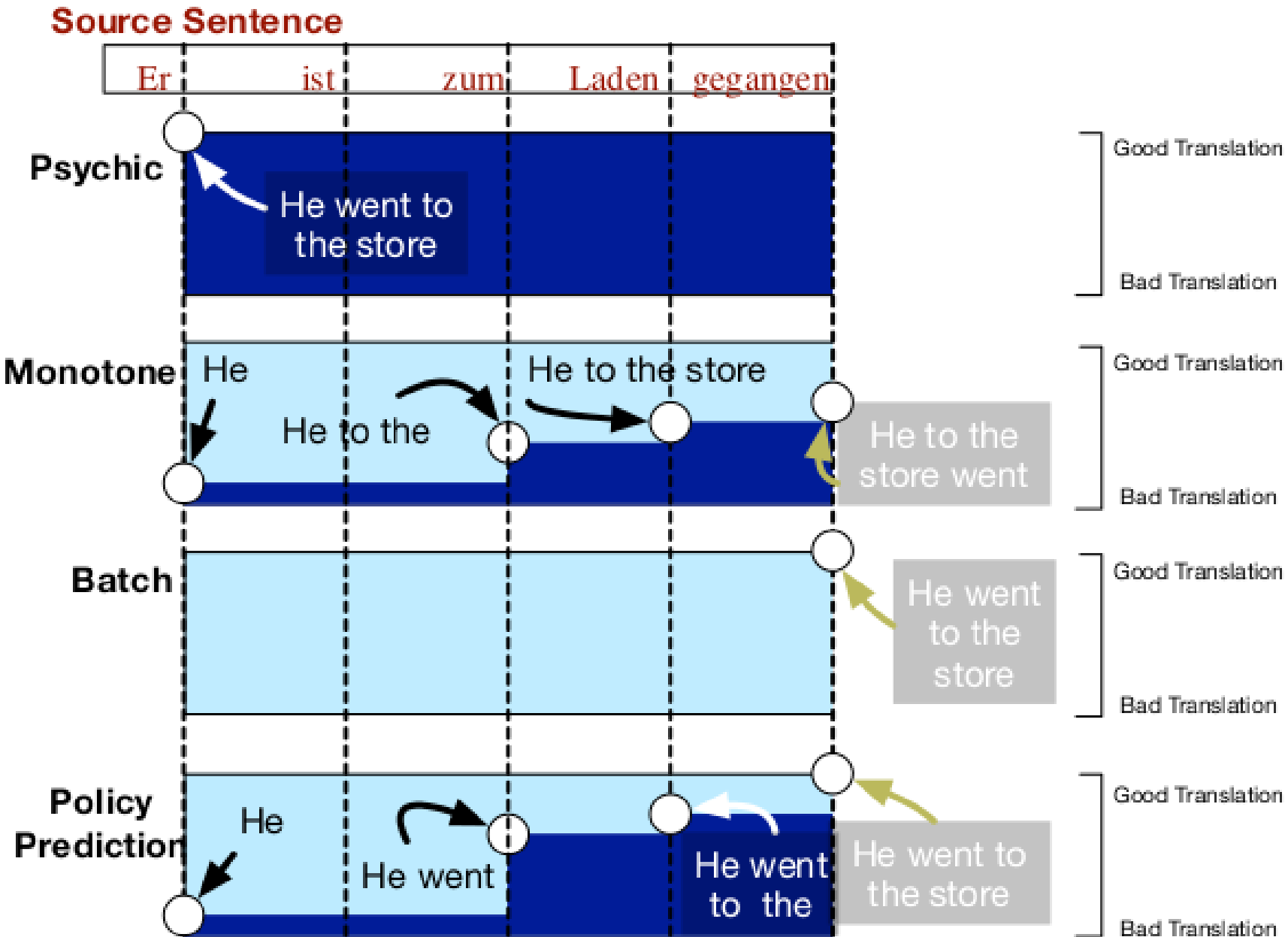
# Evaluating performance and baselines



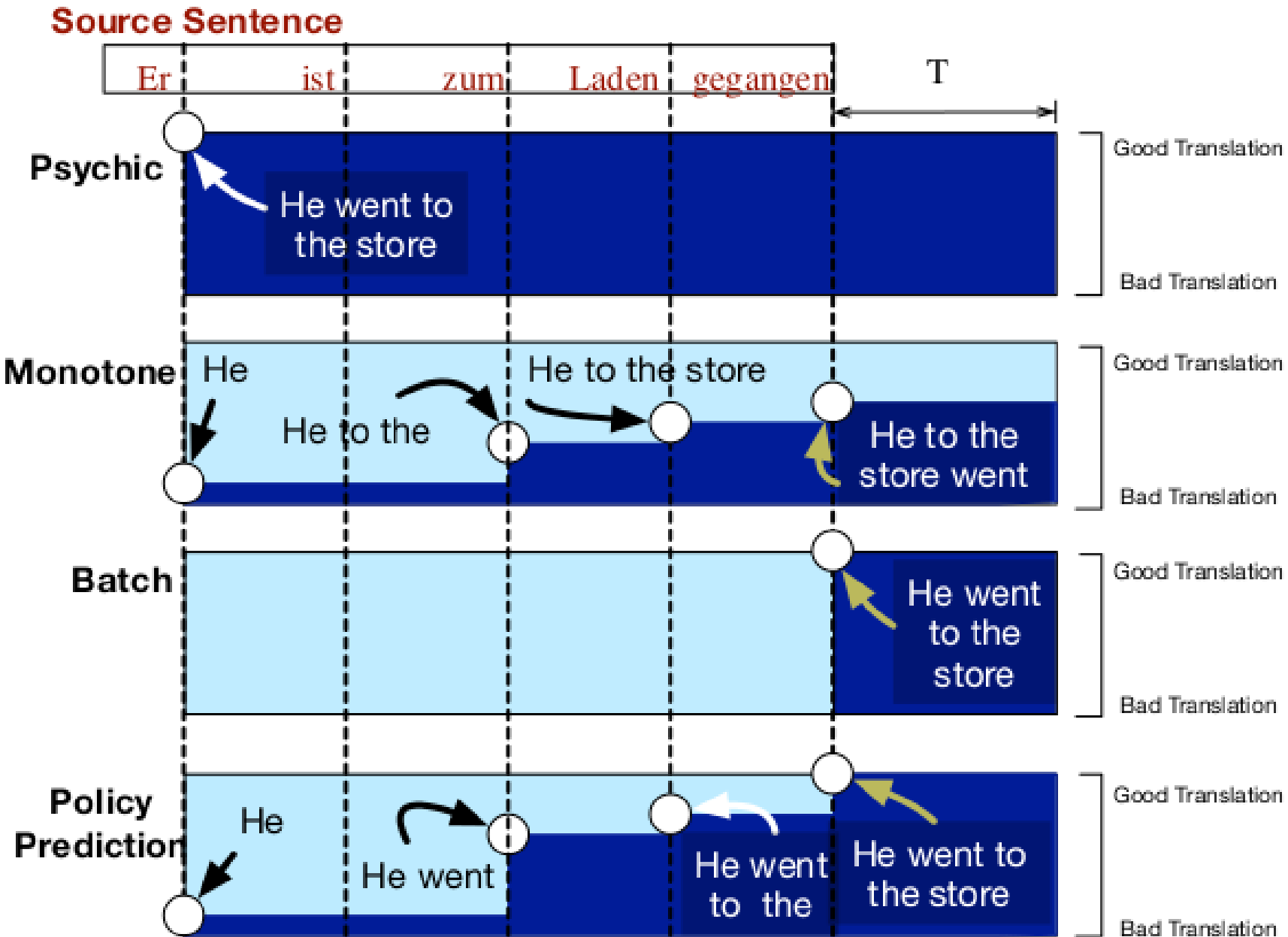
# Evaluating performance and baselines



# Evaluating performance and baselines

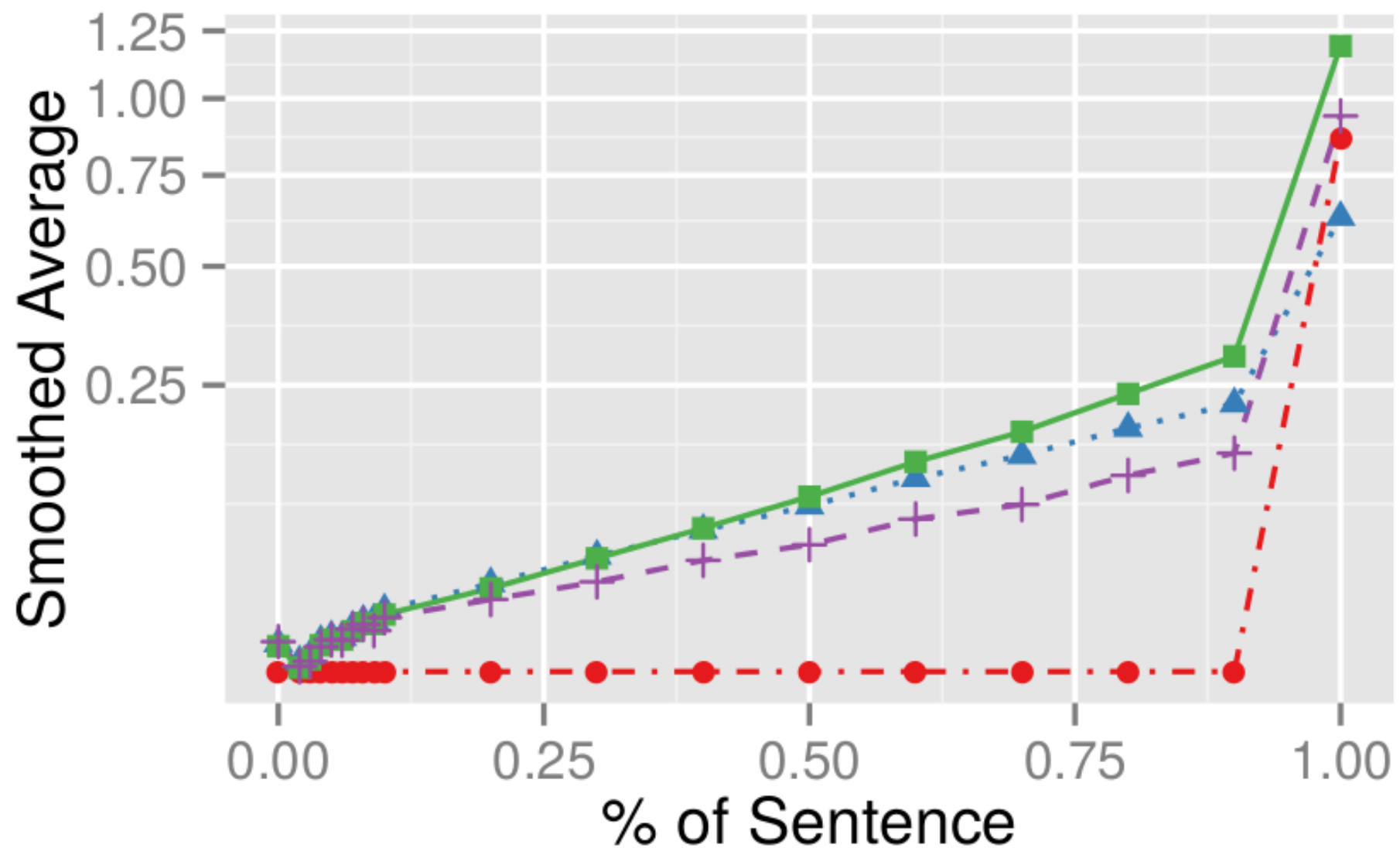


# Evaluating performance and baselines



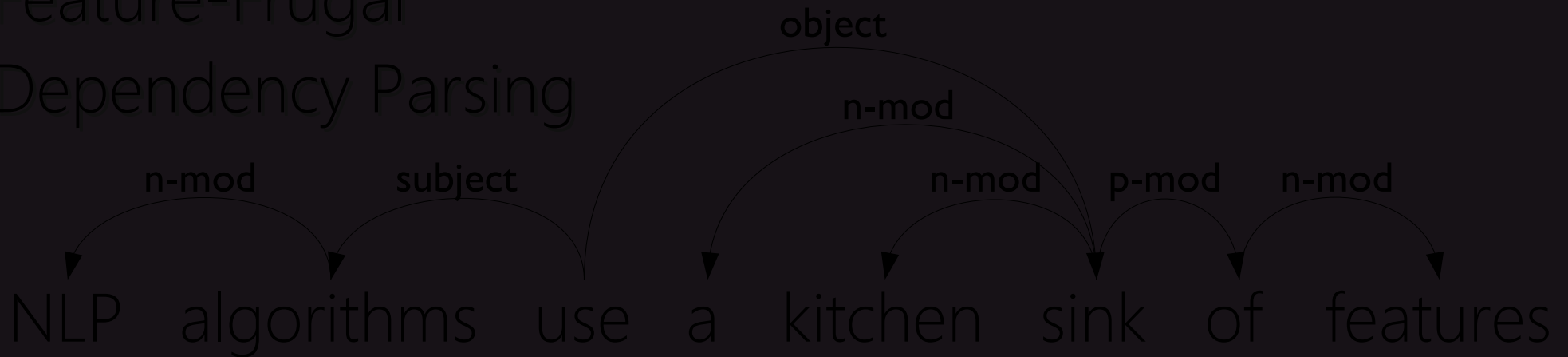


# Evaluating performance and baselines



● Batch ▲ Monotone ■ Optimal + Learned

# Feature-Frugal Dependency Parsing

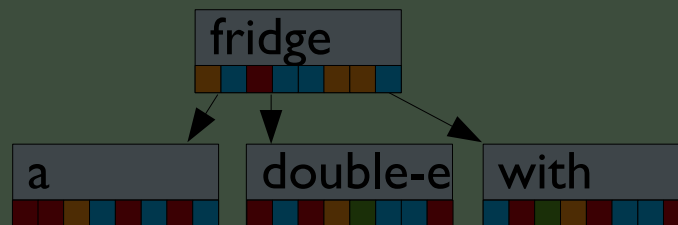


Quizbowl  
(Incremental  
Question  
Answering)

Ich bin mit dem Zug nach Ulm gefahren  
I am with the train to Ulm traveled  
I traveled by train to Ulm



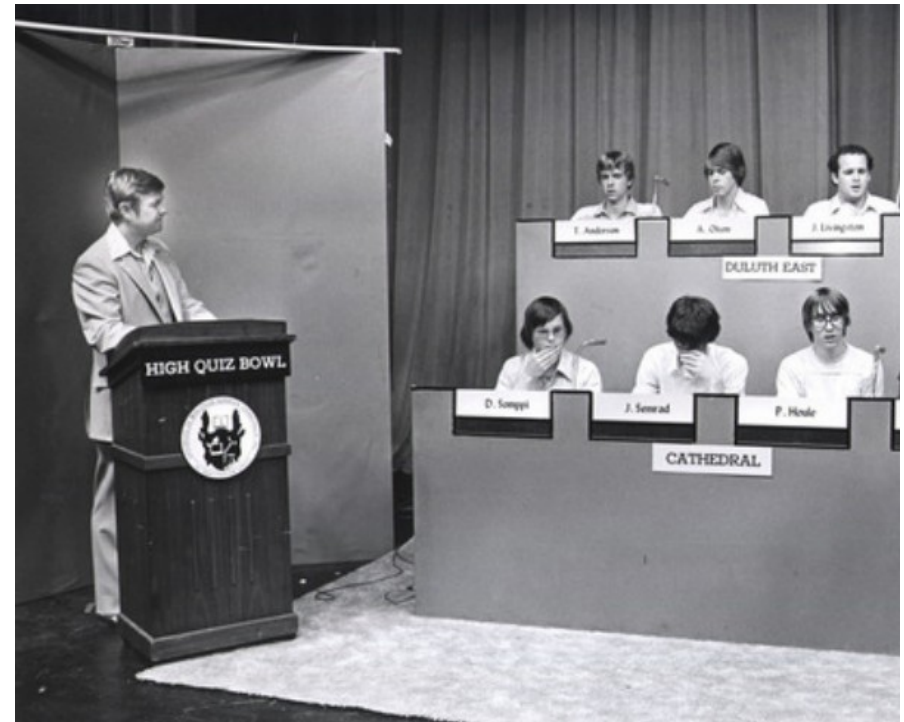
Simultaneous  
Machine  
Interpretation



# Outline

# Humans doing incremental prediction

- Game called “quiz bowl”
- Two teams play each other
  - Moderator reads a question
  - When a team knows the answer, they buzz in
  - If **right**, they get points; **otherwise**, rest of the question is read to the other team
- Hundreds of teams in the US alone
- Example ...



# Quizbowl example

**With Leo Szilard, he invented a doubly-eponymous**

# Quizbowl example

With Leo Szilard, he invented a doubly-eponymous refrigerator with no moving parts. He did not take interaction with neighbors into account when formulating his theory

# Quizbowl example

With Leo Szilard, he invented a doubly-eponymous refrigerator with no moving parts. He did not take interaction with neighbors into account when formulating his theory of heat capacity, so

# Quizbowl example

With Leo Szilard, he invented a doubly-eponymous refrigerator with no moving parts. He did not take interaction with neighbors into account when formulating his theory of heat capacity, so Debye adjusted the theory for low temperatures. His summation convention automatically sums repeated indices in tensor products. His name is attached to the A and B coefficients



# Quizbowl example

With Leo Szilard, he invented a doubly-eponymous refrigerator with no moving parts. He did not take interaction with neighbors into account when formulating his theory of heat capacity, so Debye adjusted the theory for low temperatures. His summation convention automatically sums repeated indices in tensor products. His name is attached to the A and B coefficients for spontaneous and stimulated emission, the subject of one of his multiple groundbreaking 1905 papers. He further developed the model of statistics sent to him by

# Quizbowl example

With Leo Szilard, he invented a doubly-eponymous refrigerator with no moving parts. He did not take interaction with neighbors into account when formulating his theory of heat capacity, so Debye adjusted the theory for low temperatures. His summation convention automatically sums repeated indices in tensor products. His name is attached to the A and B coefficients for spontaneous and stimulated emission, the subject of one of his multiple groundbreaking 1905 papers. He further developed the model of statistics sent to him by Bose to describe particles with integer spin. For 10 points, who is this German physicist best known for formulating

# Quizbowl example

With Leo Szilard, he invented a doubly-eponymous refrigerator with no moving parts. He did not take interaction with neighbors into account when formulating his theory of heat capacity, so Debye adjusted the theory for low temperatures. His summation convention automatically sums repeated indices in tensor products. His name is attached to the A and B coefficients for spontaneous and stimulated emission, the subject of one of his multiple groundbreaking 1905 papers. He further developed the model of statistics sent to him by Bose to describe particles with integer spin. For 10 points, who is this German physicist best known for formulating the special and general theories of relativity?

# Solving incrementally

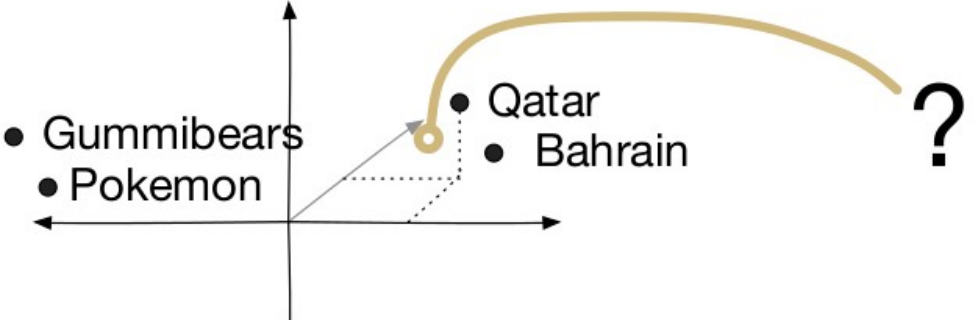
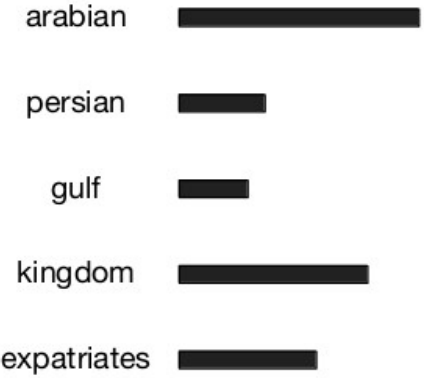
- Action: buzz now or wait
  - Content Model is constantly generating guesses
  - Oracle provides examples where it is correct
  - The Policy generalizes to test data
  - Features represent our state

## Qatar

From Wikipedia, the free encyclopedia

*For other places with the same name, see [Qatar \(disambiguation\)](#).*

**Qatar** (ⓘ/ˈkɑːtɑːr/, ⓘ/kɑːtər/ or ⓘ/kəˈtɑːr/<sup>[6]</sup> Arabic: قطر *Qatar* [ˈqɑtˤɑr]; local the **State of Qatar** (Arabic: دولة قطر *Dawlat Qatar*), is a sovereign Arab the small Qatar Peninsula on the northeastern coast of the **Arabian Penir** to the south, with the rest of its territory surrounded by the **Persian Gulf**. from the nearby island kingdom of **Bahrain**. In 2013, Qatar's total populat and 1.5 million **expatriates**.<sup>[8]</sup>



# Evaluation methodology

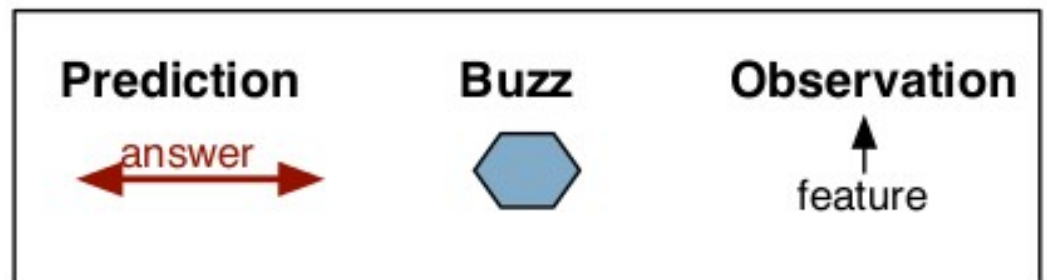
- Mechanical Turk to collect human data
- 7000 questions were



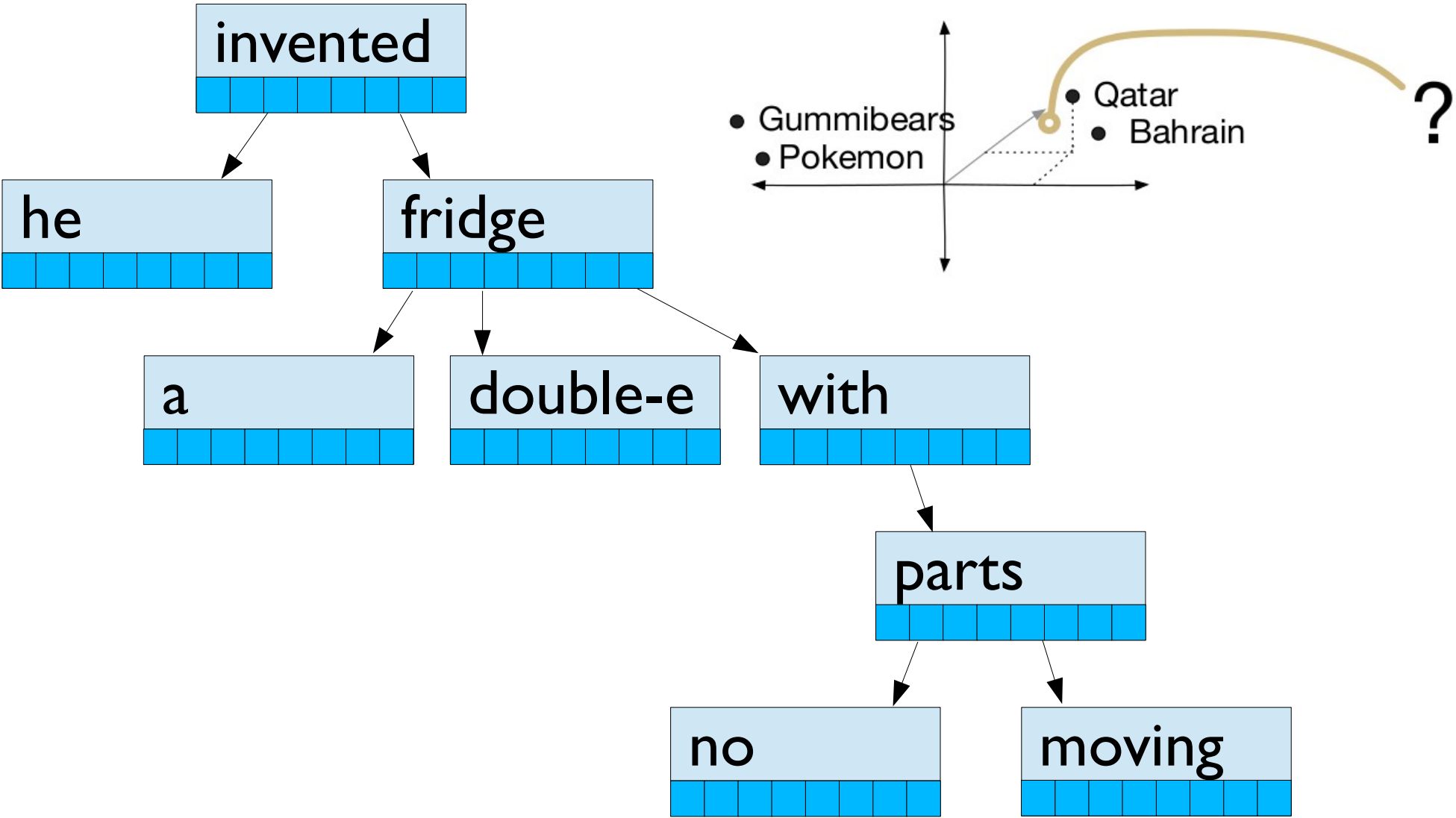
## Big problem:

“this man shot at Aaron Burr”  
*is very different from*  
“Aaron Burr shot at this man”

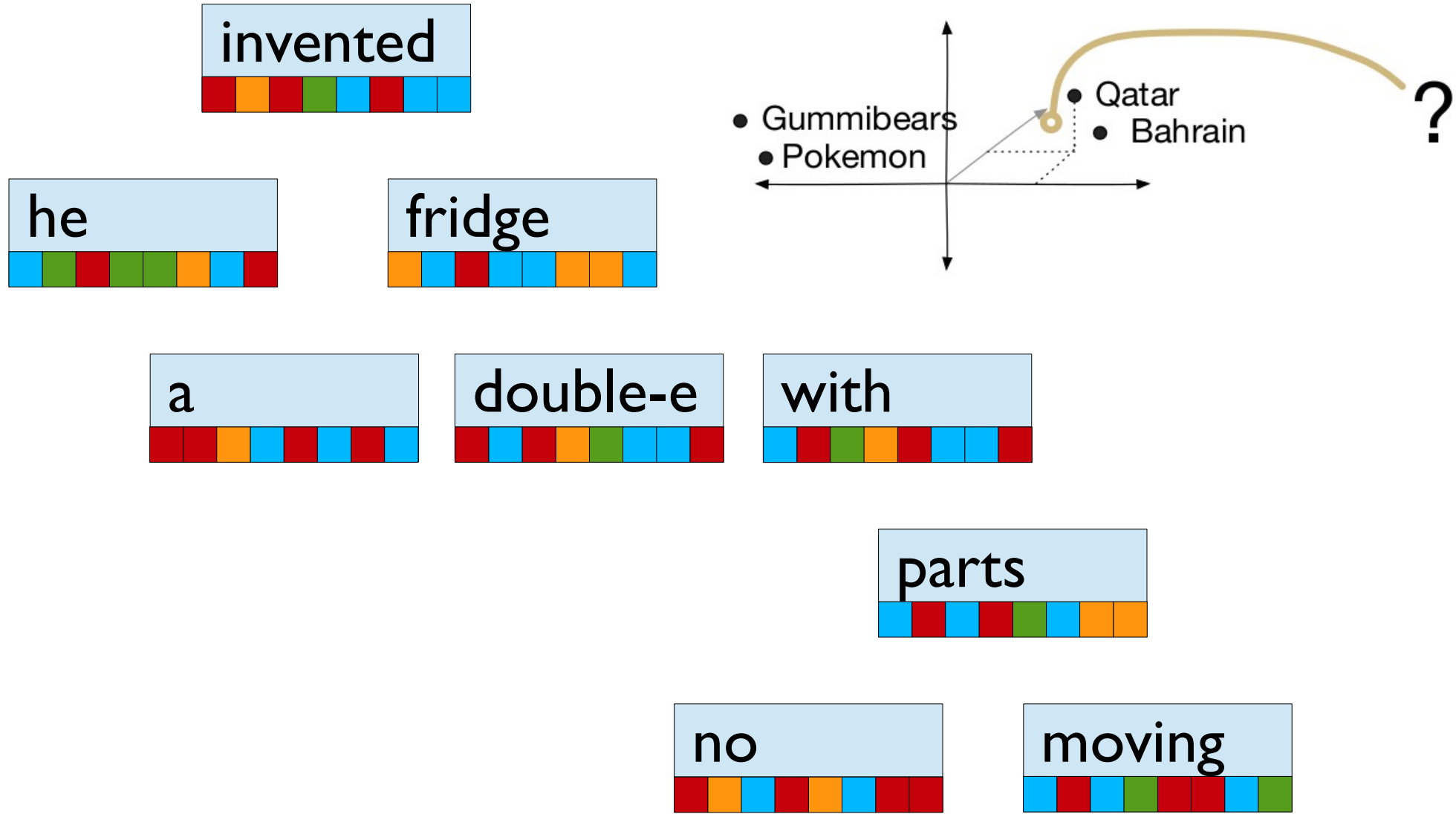
- Total of 461 unique users
- Leaderboard to encourage users



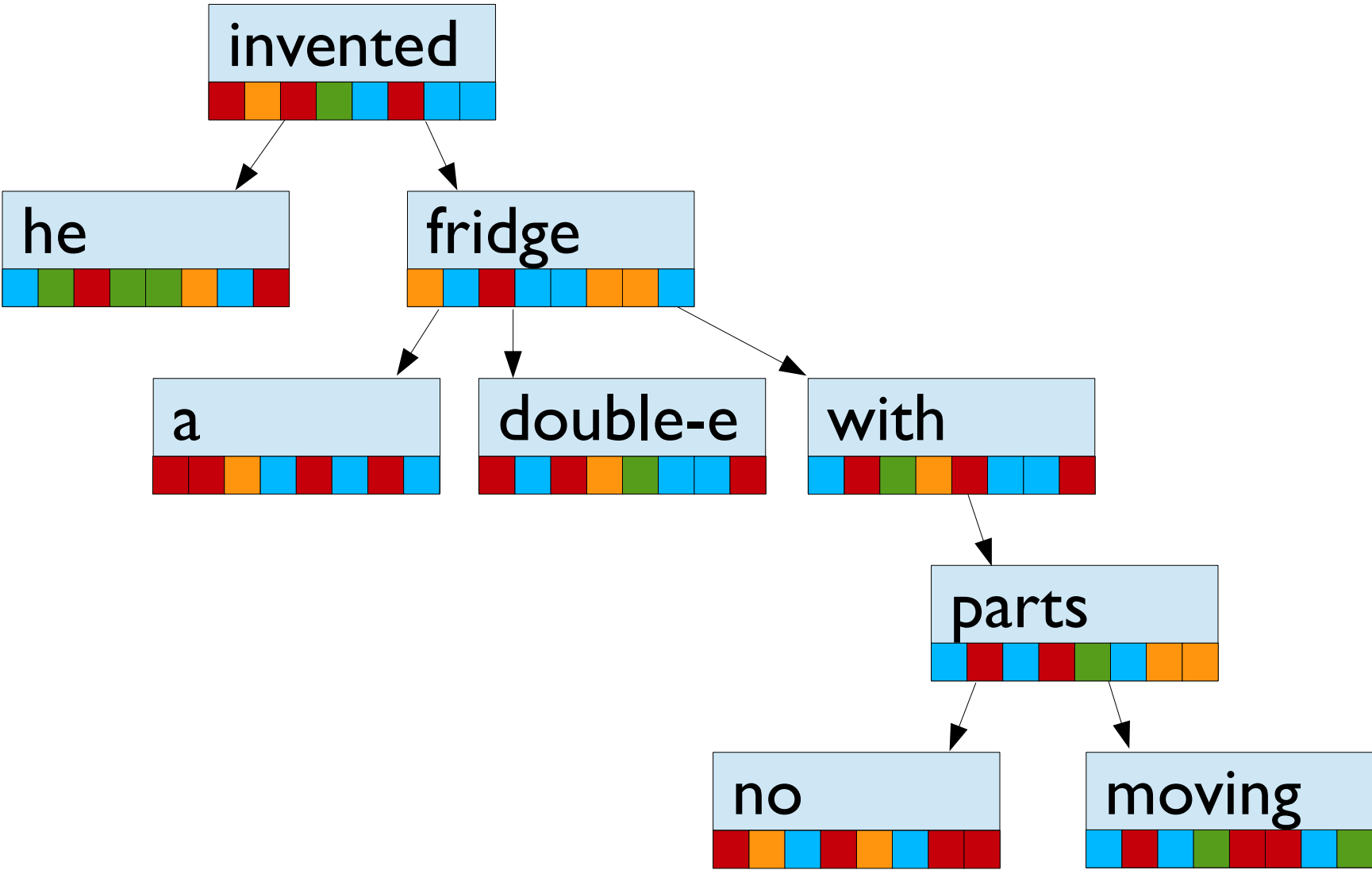
# Challenge: modeling compositionality



# Challenge: modeling compositionality

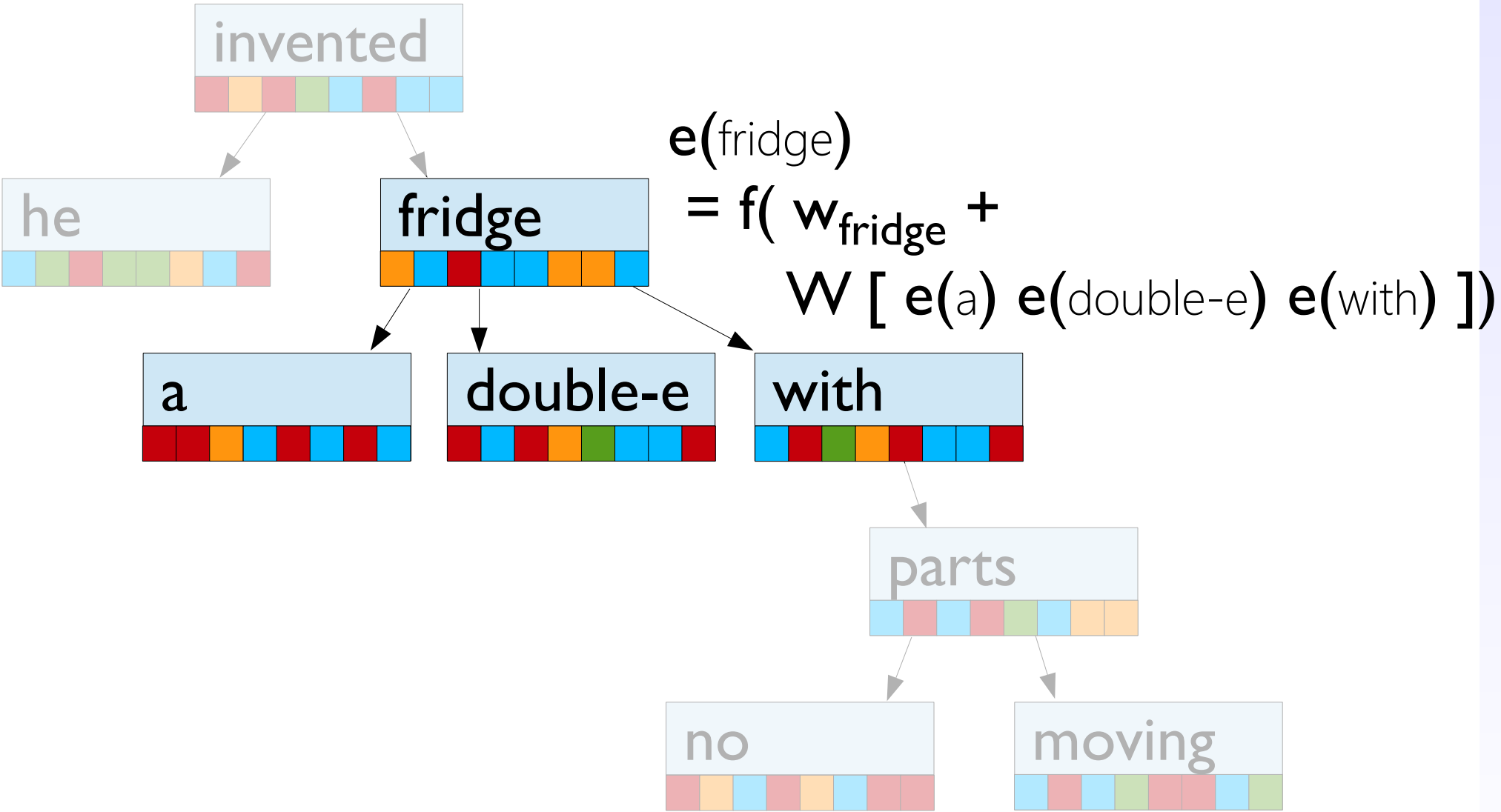


# Challenge: modeling compositionality

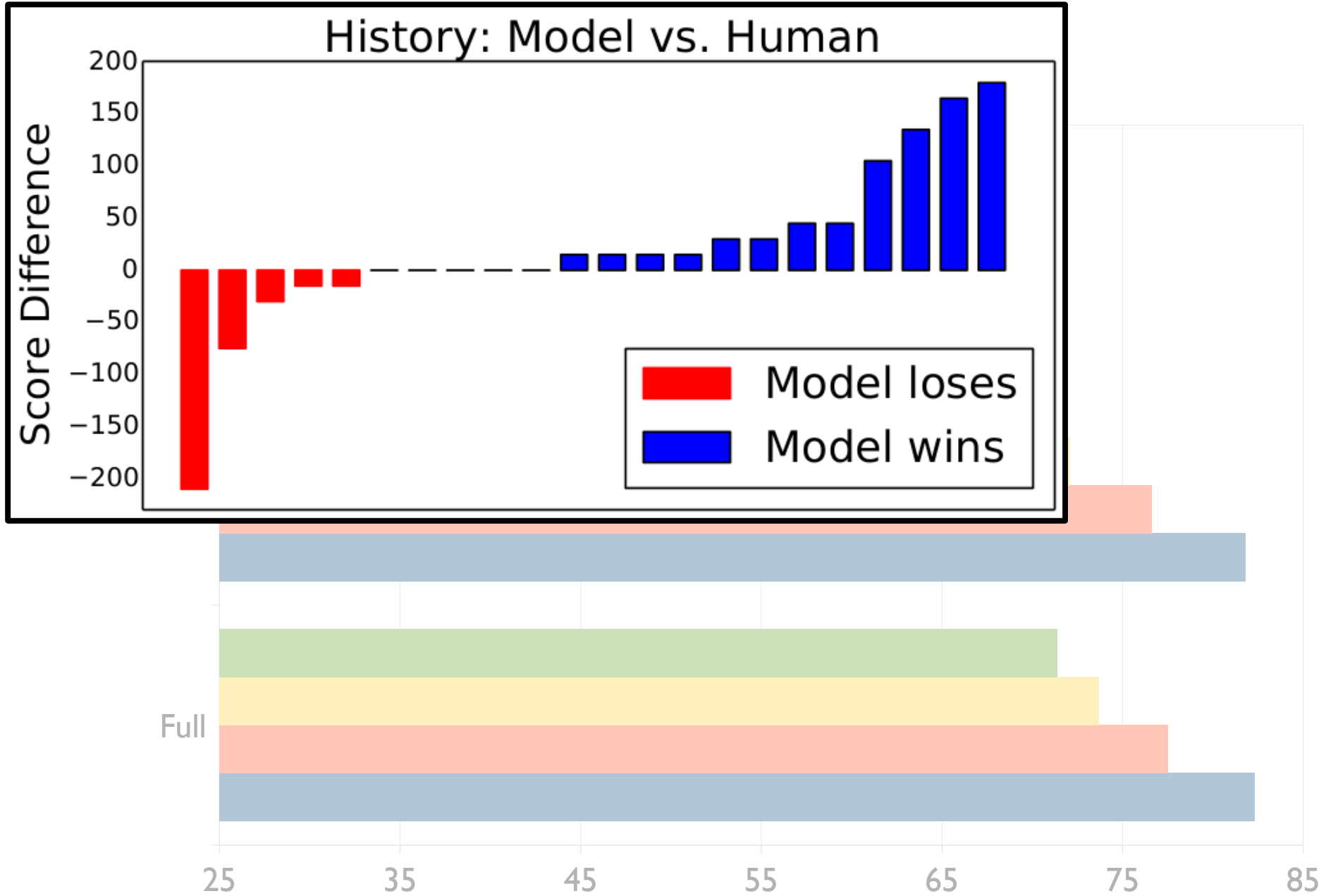




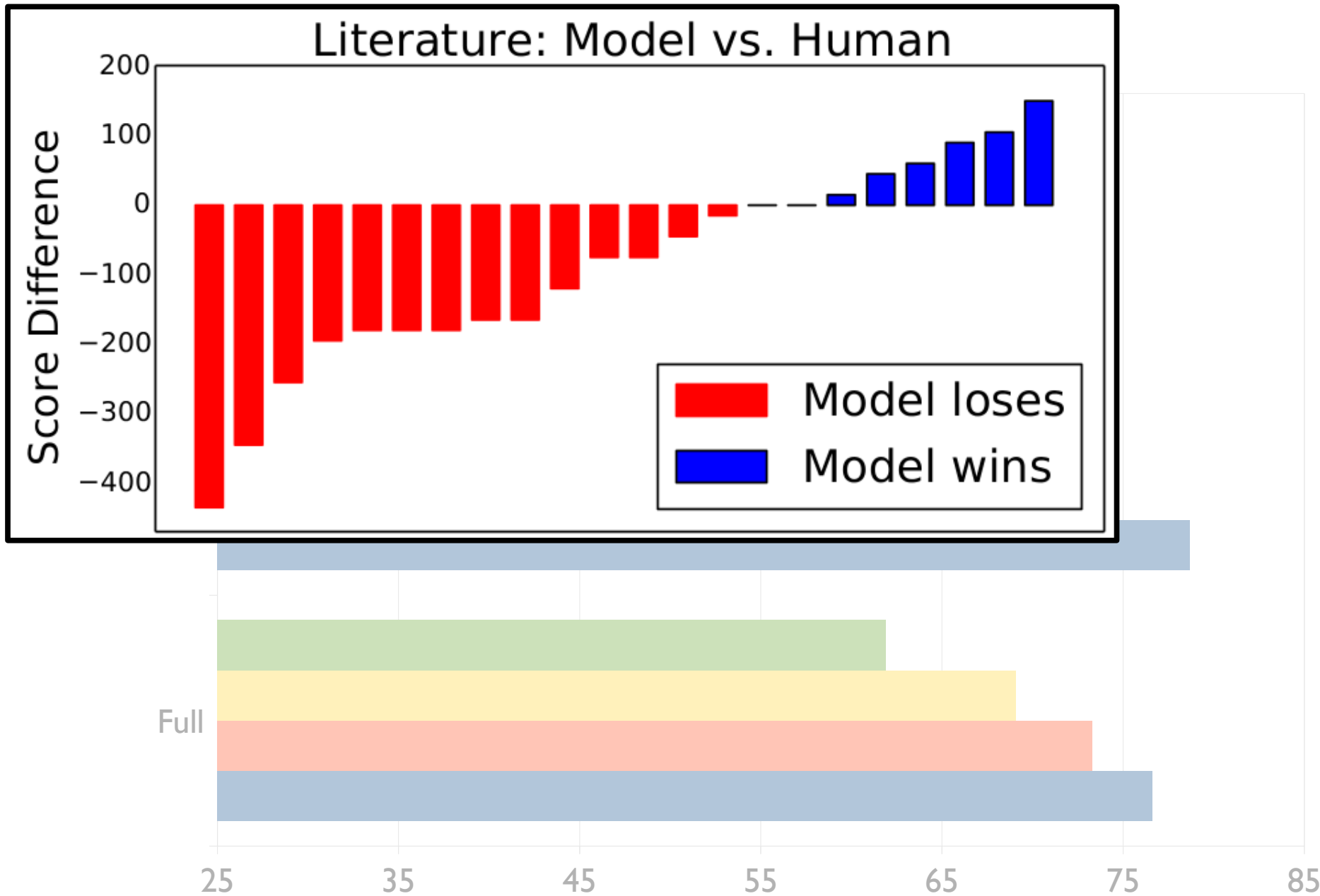
# Challenge: modeling compositionality



# Results on question-answering task



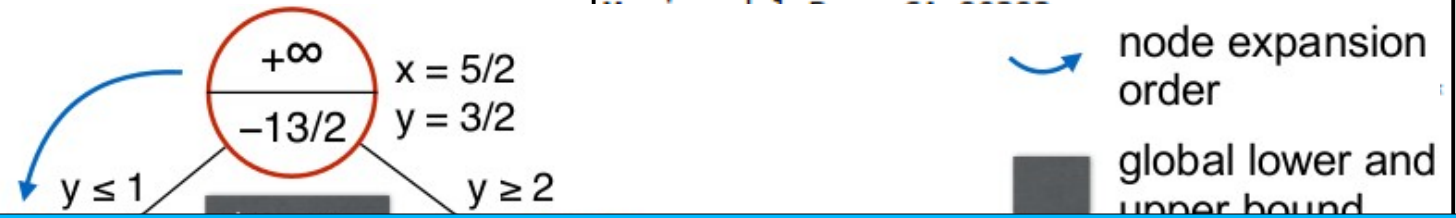
# Results on question-answering task



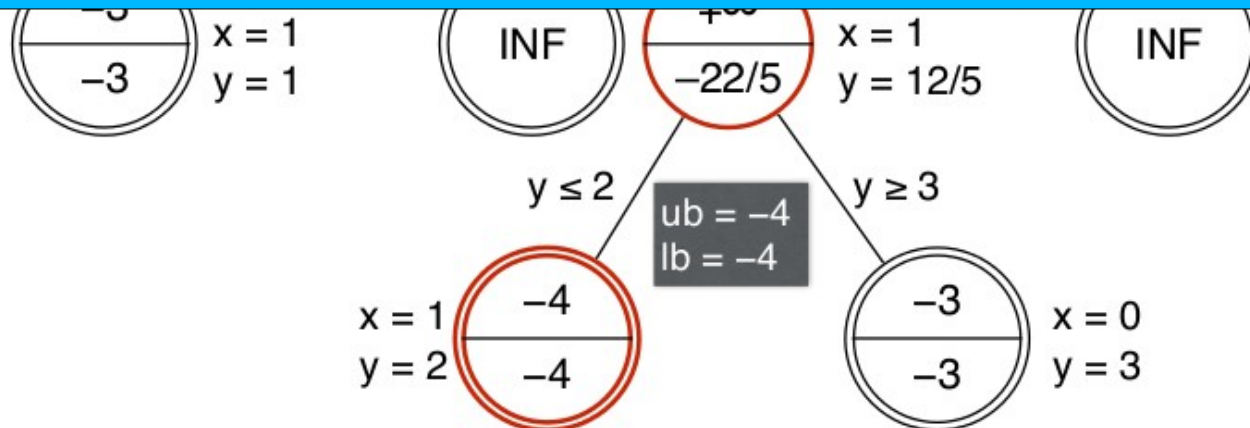
# Moving to more general frameworks

- Lots of NLP (+al) problems can be cast *at test time* as integer linear programs
- ILPs are usually solved using

Thursday, March 6, 2003  
 10:30am - 12:00pm  
 11th Floor Large Conference Room  
 USC/Information Sciences Institute  
 4676 Admiralty Way, Suite 1001



Branch and bound involves a complex heuristic search  
 Can we learn to perform this search efficiently?



min  $-2x - y$   
 s.t.  $3x - 5y \leq 0$   
 $3x + 5y \leq 15$   
 $x \geq 0, y \geq 0$   
 $x, y \in \mathbb{Z}$

(He+Eisner+D, NIPS 2014)

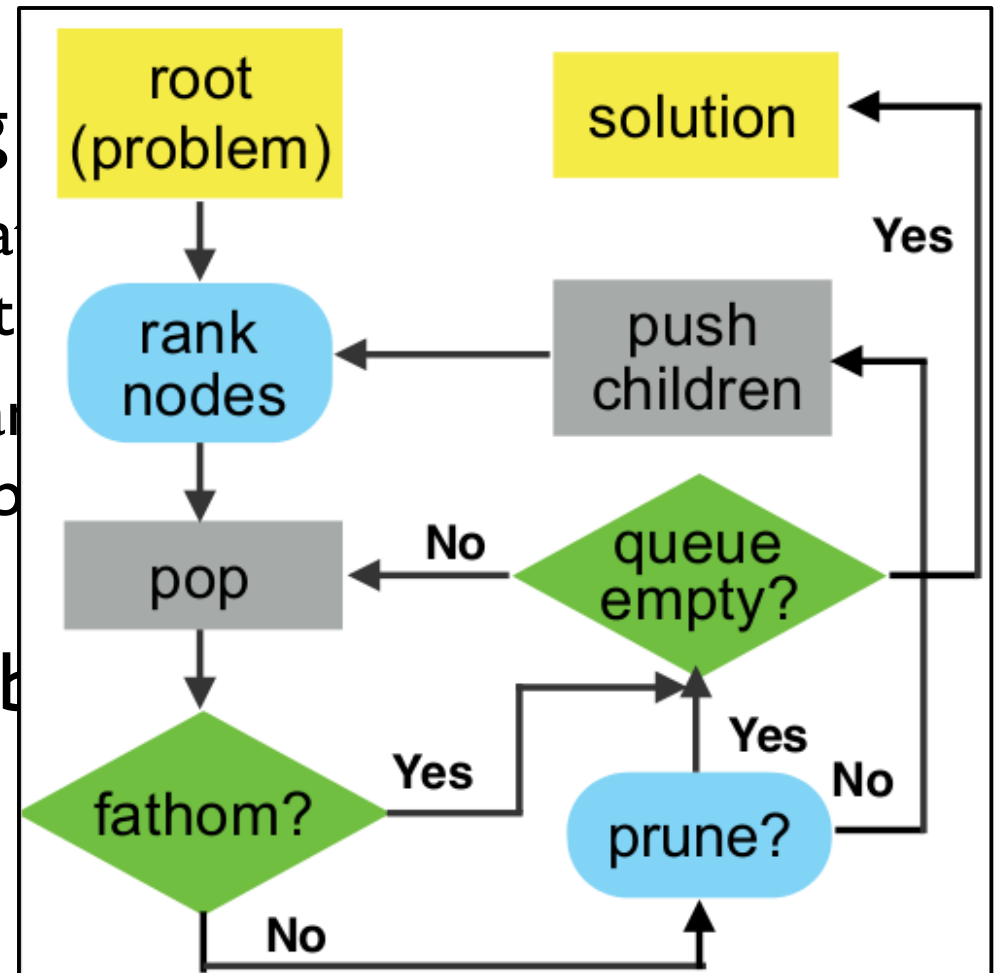
# Some intuition

- A good search strategy should:
  - find a good incumbent solution early
  - identify non-promising nodes before expansion

- “Good” varies depending

- DFS should only be used as a good feasible solution to
- Best-bound-first search can visit many nodes, but should not be

- We will learn a heuristic to capture this intuition



# Training and experiments

- Same algorithm
- Four (standard deviation) while exploring 0.05%, 1.5%, 5.1% and 47% of the nodes explored by Gurobi!
- Comparison
  - DFS (baseline)
  - Gurobi (thousands of person-hours of effort)
- Measures:
  - Optimality Gap, Integrality Gap, and improvement from initial heuristic solution

Dataset	Ours(DAgger training)			DFS			Gurobi		
	OGap	IGap	Impr	OGap	IGap	Impr	OGap	IGap	Impr
MIK	<b>0.23</b>	16.63	<b>4.39</b>	6.74	35.48	0.00	<b>0.17</b>	<b>15.24</b>	0.36
Regions1	<b>0.54</b>	<b>4.53</b>	<b>10.57</b>	3.07	8.48	8.61	2.24	7.20	0.60
Regions2	<b>1.22</b>	<b>6.76</b>	<b>19.36</b>	4.75	11.38	15.12	1.65	7.48	2.15
Hybrid	<b>0.87</b>	<b>20.28</b>	<b>24.46</b>	1.69	23.08	23.53	<b>1.37</b>	23.49	1.58

# Some other fun stuff...

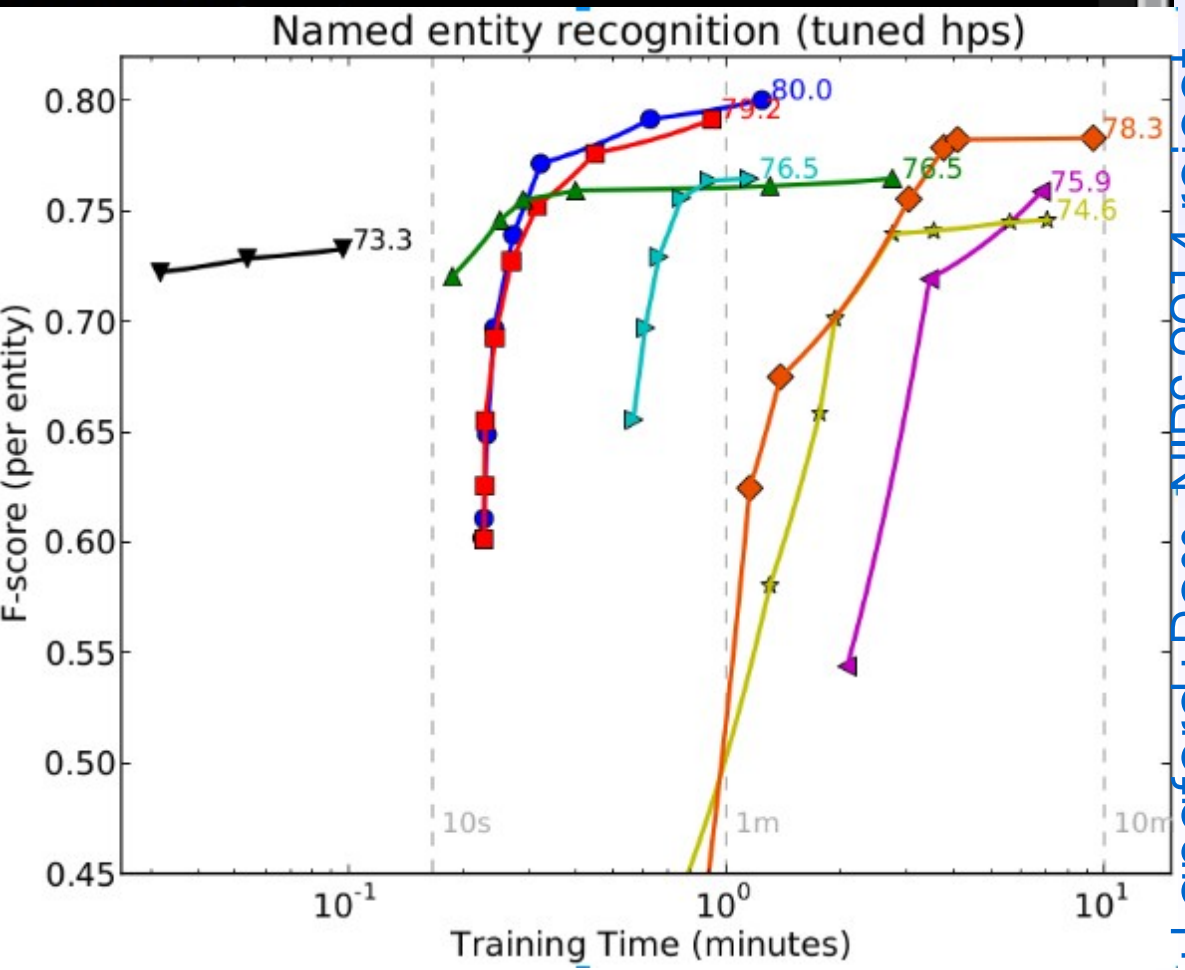
```
test_search.py
File Edit Options Buffers Tools Python YASnippet Development Help

import sys
import pyvw

class SequenceLabeler(pyvw.SearchTask):
    def __init__(self, vw, srn):
        # you must must must initialize srn
        # this will automatically initialize srn
        pyvw.SearchTask.__init__(self, vw)
        srn.set_options(srn.AllOptions)

    def _run(self, sentence):
        output = []
        for tag, word in sentence:
            with self.vw.example(word):
                pred = self.srn.predict()
                output.append(pred)
        return output

for curPass in range(10):
    sequenceLabeler.learn(my_data)
```



(D+Langford+Ross, NIPS 2014 reject - on arxiv)





Jordan B-G



Jason Eisner



Alvin Grissom



He He



Mohit Iyer



John Langford

- Reasoning with incomplete information is useful for *speed* and *modeling*
- *Imitation learning* can help us build such systems
- Wide range of new, interesting problems to work on!

Thanks! Questions?