# Dynamic Feature Selection for Dependency Parsing

## He He, Hal Daumé III and Jason Eisner

EMNLP 2013, Seattle

UNIVERSITY OF MARYLAND

JOHNS HOPKINS UNIVERSITY

# Structured Prediction in NLP
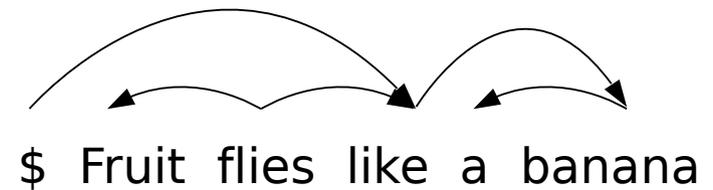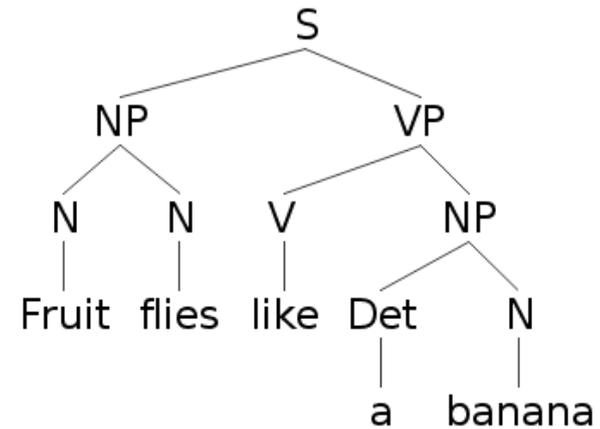
## Part-of-Speech Tagging
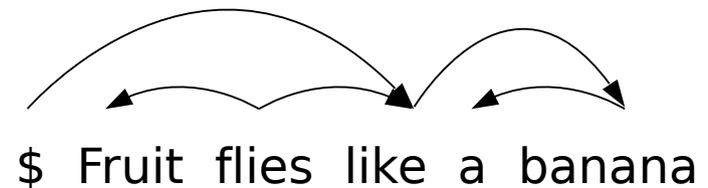
N ⟶ N ⟶ V ⟶ Det ⟶ N

↓      ↓      ↓      ↓      ↓

Fruit flies like    a     banana

## Machine Translation

Fruit flies like a banana .

果     蝇 喜欢    香蕉     。

*summarization, name entity resolution and many more ...*

## Parsing

```
              S
          ┌───┴───┐
         NP       VP
        ┌─┴─┐   ┌──┴──┐
        N   N   V     NP
        │   │   │    ┌─┴─┐
      Fruit flies like Det  N
                      │    │
                      a  banana
```

$ Fruit flies like a banana

# Structured Prediction in NLP



**Part-of-Speech Tagging**
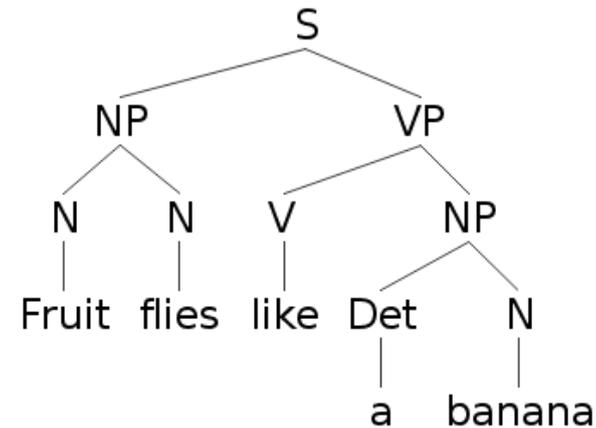
N ⟶ N ⟶ V ⟶ Det ⟶ N

↓    ↓    ↓    ↓    ↓

Fruit flies  like    a    banana

**Machine Translation**

Fruit flies like a banana  .

果　　蝇 喜欢　 香蕉　 。

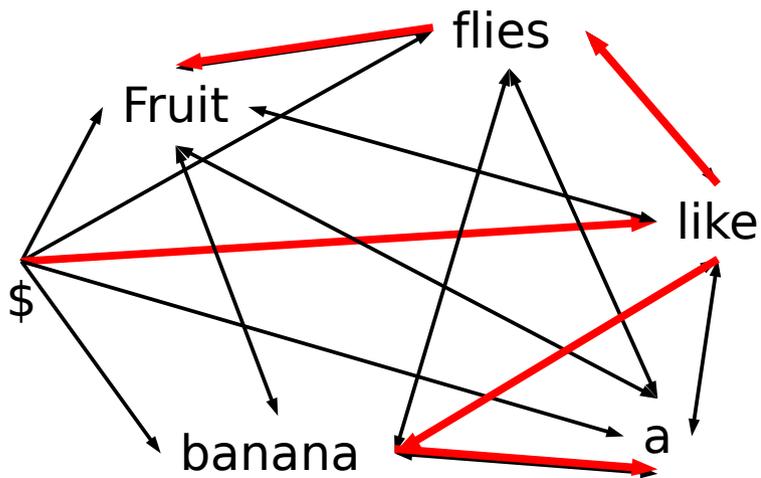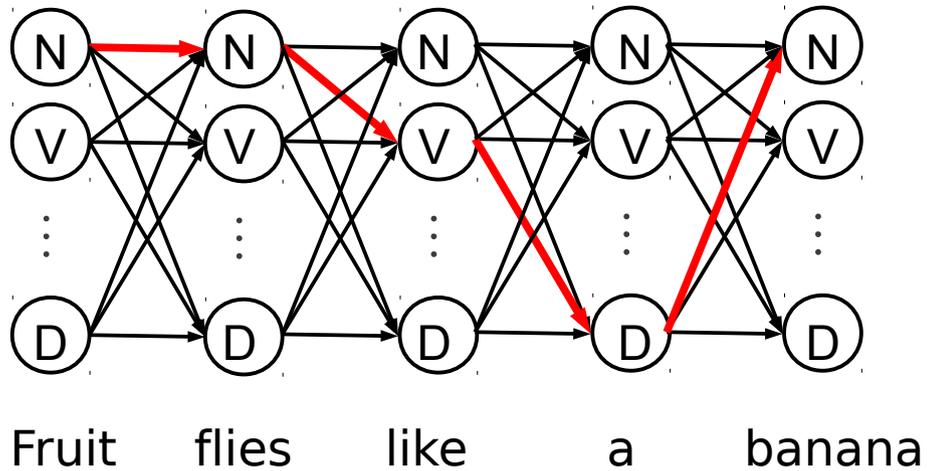*summarization, name entity resolution and many more ...*

**Parsing**



$ Fruit  flies  like  a  banana

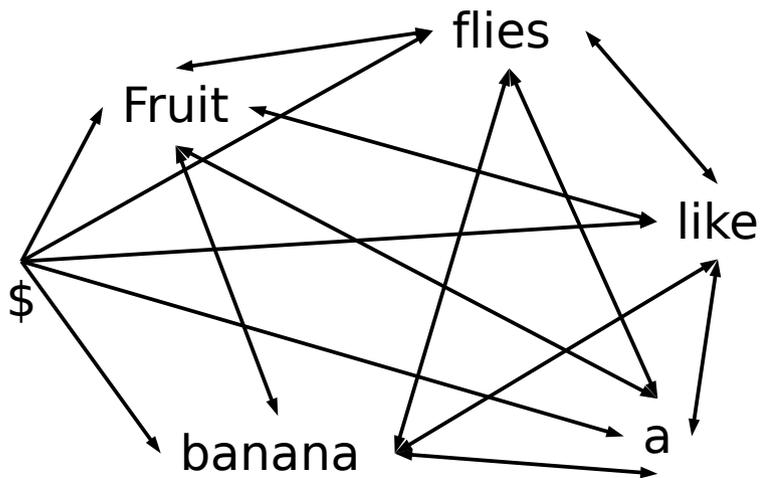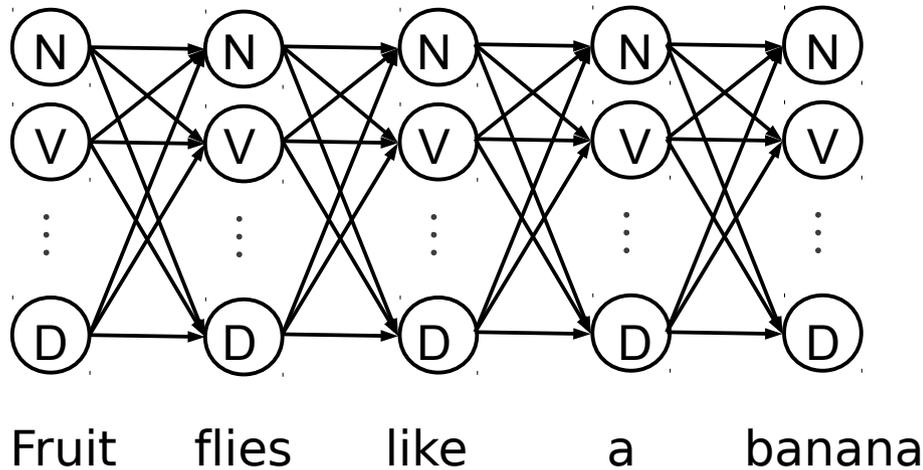**<span style="color:red">Exponentially increasing search space</span>**

**<span style="color:red">Millions of features for scoring</span>**
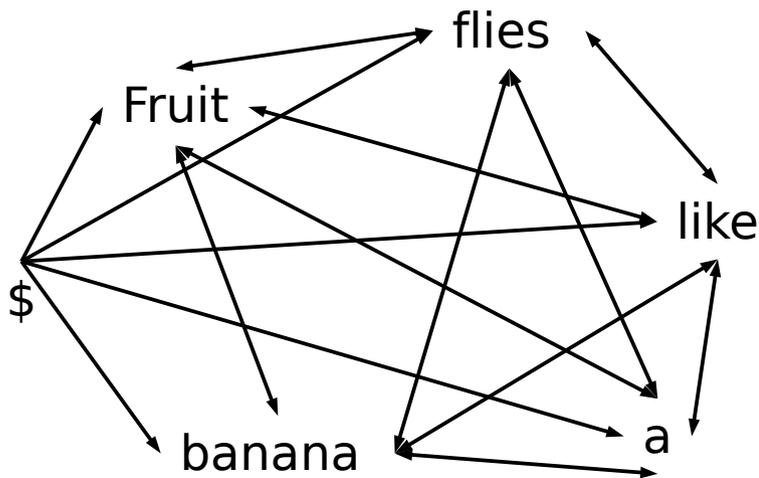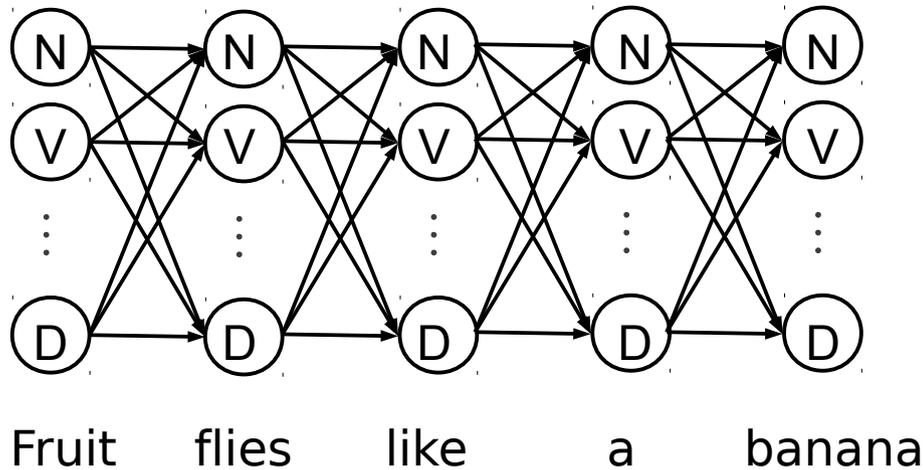
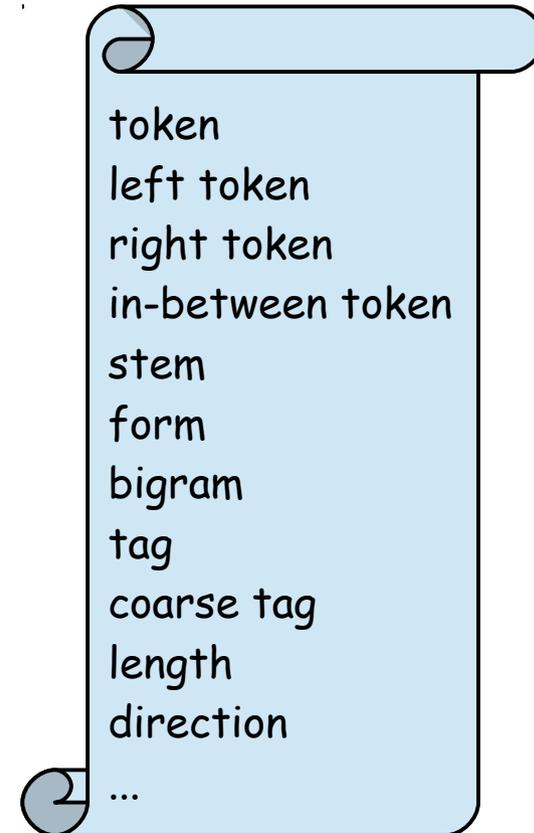# Structured Prediction in NLP

# Structured Prediction in NLP



Fruit    flies    like    a    banana

### Feature templates per edge

token
left token
right token
in-between token
stem
form
bigram
tag
coarse tag
length
direction
...

# Structured Prediction in NLP



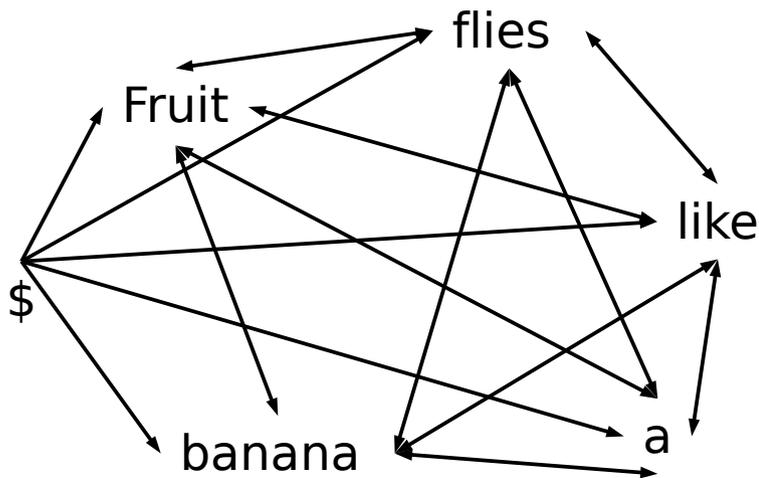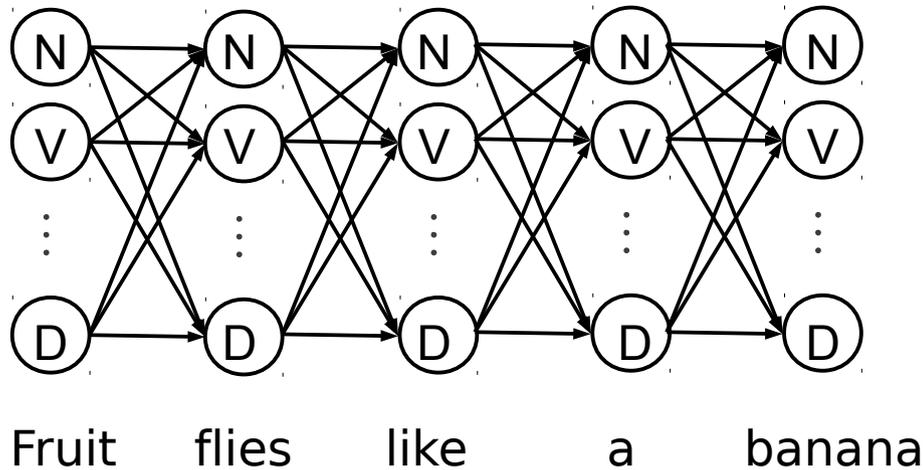Fruit    flies    like    a    banana

Feature templates per edge

token
left token
right token
in-between token
stem
form
bigram
tag
coarse tag
length
direction
...

(head_token + mod_token)
X
(head_tag + mod_tag)

6

# Structured Prediction in NLP



Fruit    flies    like    a    banana

Feature templates per edge

token
left token
right token
in-between token
stem
form
bigram
tag
coarse tag
length
direction
...

**HUGE**

(head_token + mod_token)
X
(head_tag + mod_tag)

# Structured Prediction in NLP



N N N N N
V V V V V
D D D D D

Fruit    flies    like    a    banana

Feature templates per edge

token

tag
coarse tag
length
direction
...

Do you need all features everywhere ?

# Structured Prediction in NLP



Feature templates per edge

token

tag
coarse tag
length
direction
...

Do you need all features everywhere ?

# Structured Prediction in NLP



Feature templates per edge
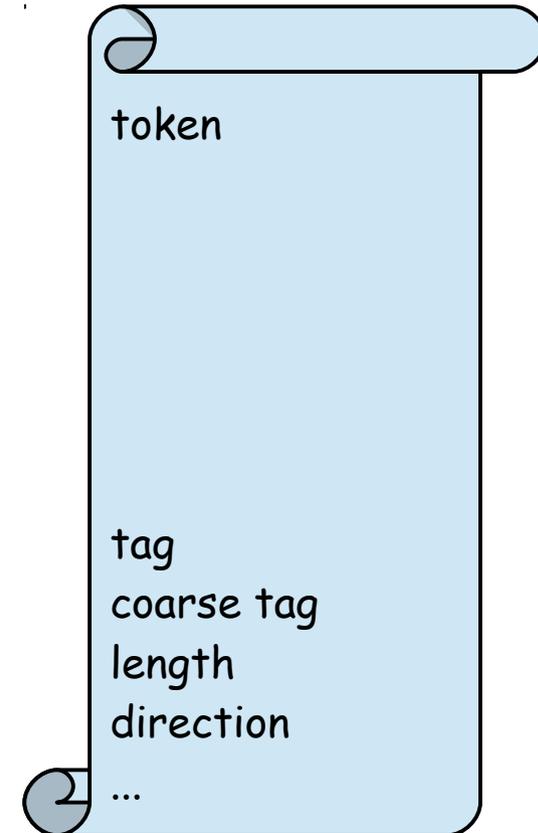
token

tag
coarse tag
length
direction
...

Do you need all features everywhere ?

# Structured Prediction in NLP



Feature templates per edge
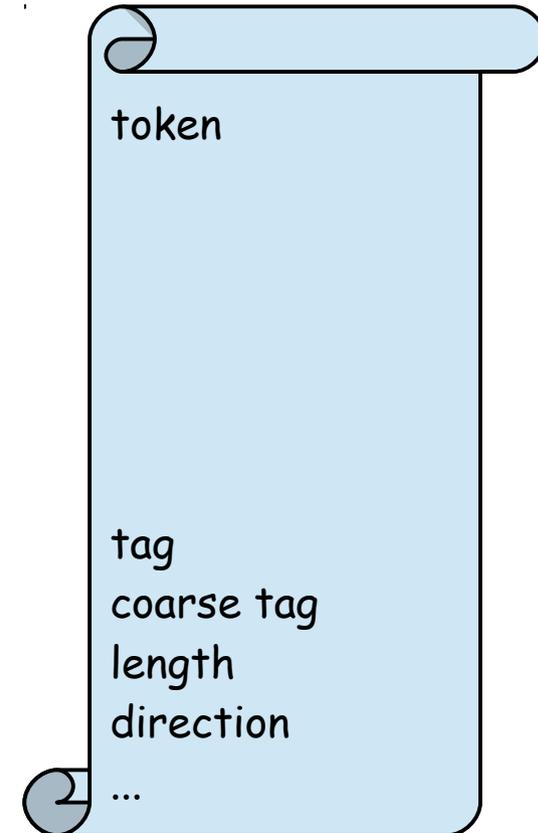
token

**Dynamic Decisions**

tag
coarse tag
length
direction
...

# Case Study: Dependency Parsing



2x to 6x speedup with little loss in accuracy

# Graph-based Dependency Parsing



$ This time , the firms were ready .

Scoring: $\phi(E) \cdot w$

# Graph-based Dependency Parsing



$ This time , the firms were ready .

firms ⟶ were

length:              1
direction:           right
modifier_token:      were
head_token:          firms
head_tag:            noun
⋮

And hundreds more!

Scoring: $\phi(E) \cdot w$

# Graph-based Dependency Parsing



Decoding: find the highest-scoring tree

# MST Dependency Parsing
## (1st-order projective)

# MST Dependency Parsing
## (1st-order projective)



**Find highest-scoring tree O($n^3$)**

# MST Dependency Parsing
## (1st-order projective)

~268 feature templates
~76M features

**Find edge scores**

**Find highest-scoring tree O(n³)**

Average Sentence

# Add features only when necessary!



This       the    firms      ready

score(This → ready) =

score(the → firms) =

# Add features only when necessary!



score(This → ready) = -0.23

score(the → firms) = 0.63

# Add features only when necessary!



+0.1

-0.23

+0.63  +0.7

This            the  firms          ready

score(This → ready) = -0.13

score(the → firms) = 1.33

# Add features only when necessary!



score(This → ready) = -0.13

score(the → firms) = 1.33

# Add features only when necessary!



score(This → ready) = -1.33

score(the → firms) = 1.33

# Add features only when necessary!



**LOSER**
+0.1
- 0.55
-1.2
-0.23

**WINNER**
+0.63
+0.7

This        the    firms        ready

score(This → ready) = -1.88

score(the → firms) = 1.33

# Add features only when necessary!

**LOSER**

+0.1  - 0.55

-0.23  -1.2

**WINNER**  +0.63  +0.7

This          the    firms            ready

score(This → ready) = -1.88

score(the → firms) = 1.33

This is a **structured** problem!
Should not look at scores independently.

# Dynamic Dependency Parsing

1. Find the highest-scoring tree after adding some features   *fast non-projective decoding*

# Dynamic Dependency Parsing

1. Find the highest-scoring tree after adding some features   *fast non-projective decoding*

2. Only edges in the current best tree can win

# Dynamic Dependency Parsing

1. Find the highest-scoring tree after adding some features   *fast non-projective decoding*

2. Only edges in the current best tree can win

   😃   are chosen by a classifier   ≤ n *decisions*

   ☹   are killed because they fight with the winners

# Dynamic Dependency Parsing

1. Find the highest-scoring tree after adding some features   *fast non-projective decoding*

2. Only edges in the current best tree can win

   🙂   are chosen by a classifier   ≤ n *decisions*

   🙁   are killed because they fight with the winners

3. Add features to undetermined edges   *by group*

# Dynamic Dependency Parsing

1. Find the highest-scoring tree after adding some features   *fast non-projective decoding*

2. Only edges in the current best tree can win

    🙂   are chosen by a classifier   *≤ n decisions*

    🙁   are killed because they fight with the winners

3. Add features to undetermined edges   *by group*

   Max # of iterations = # of feature groups

# + first feature group



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



$ This time , the firms were ready .

— Undetermined edge
— Current 1-best tree
— Winner edge
(permanently in 1-best tree)
- - Loser edge

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

$ This time , the firms were ready .

Undetermined edge
Current 1-best tree
Winner edge
(permanently in 1-best tree)
Loser edge

*Non-projective decoding to find new 1-best tree*

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

$ This time , the firms were ready .

Undetermined edge
Current 1-best tree
Winner edge
(permanently in 1-best tree)
Loser edge

*Classifier picks winners among the blue edges*

33

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



$ This time , the firms were ready .

Undetermined edge
Current 1-best tree
Winner edge
(permanently in 1-best tree)
Loser edge

*Remove losers in conflict with the winners*

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

$ This time , the firms were ready .

Remove losers in conflict
with *the winners*

—— Undetermined edge
—— Current 1-best tree
—— Winner edge
(permanently in 1-best tree)
- - Loser edge

+ next feature group       44    gray edges with unknown fate...
                           27    features per gray edge



$ This time , the firms were ready .

— Undetermined edge
— Current 1-best tree
— Winner edge
(permanently in 1-best tree)
- - Loser edge

+ next feature group   44   gray edges with unknown fate...
                       27   features per gray edge
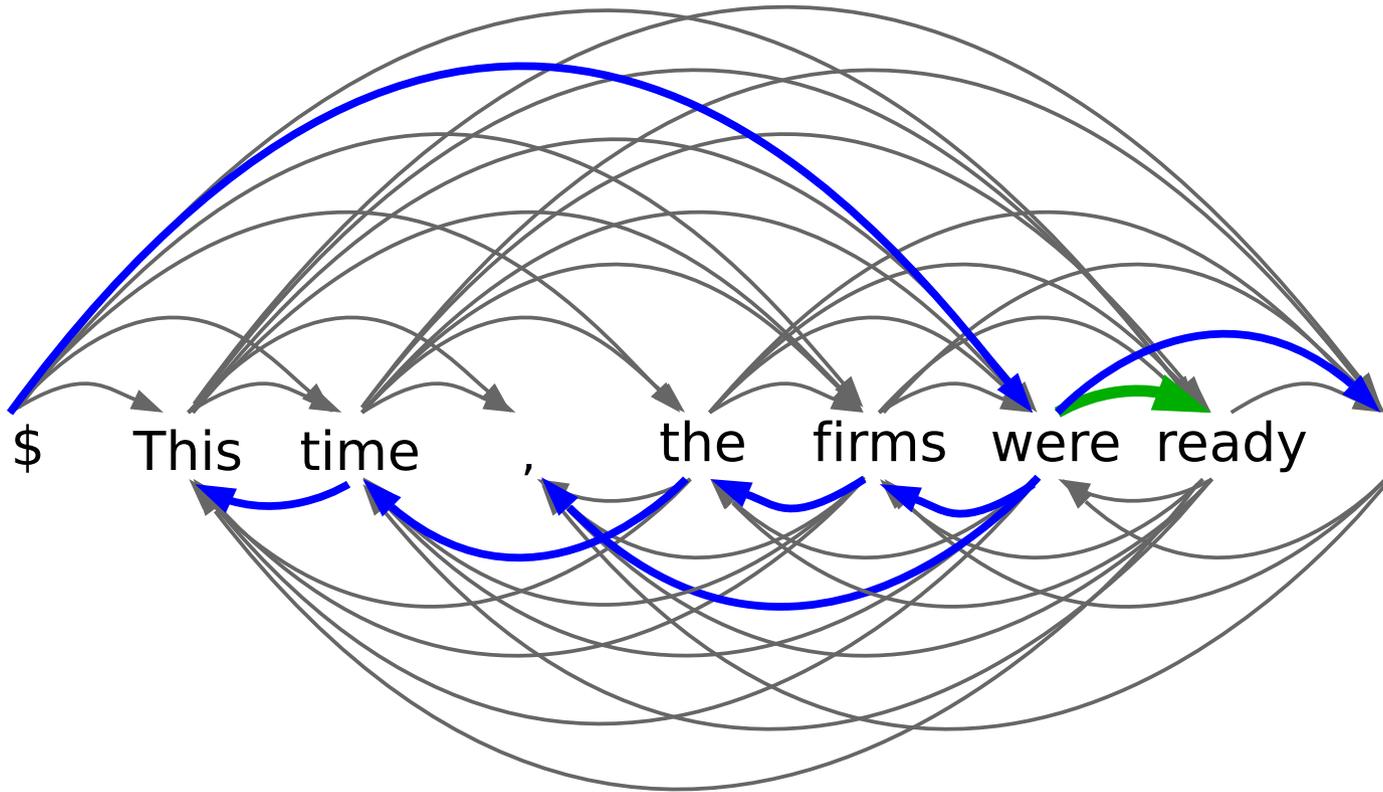


$ This time , the firms were ready .

— Undetermined edge
— Current 1-best tree
— Winner edge
  (permanently in 1-best tree)
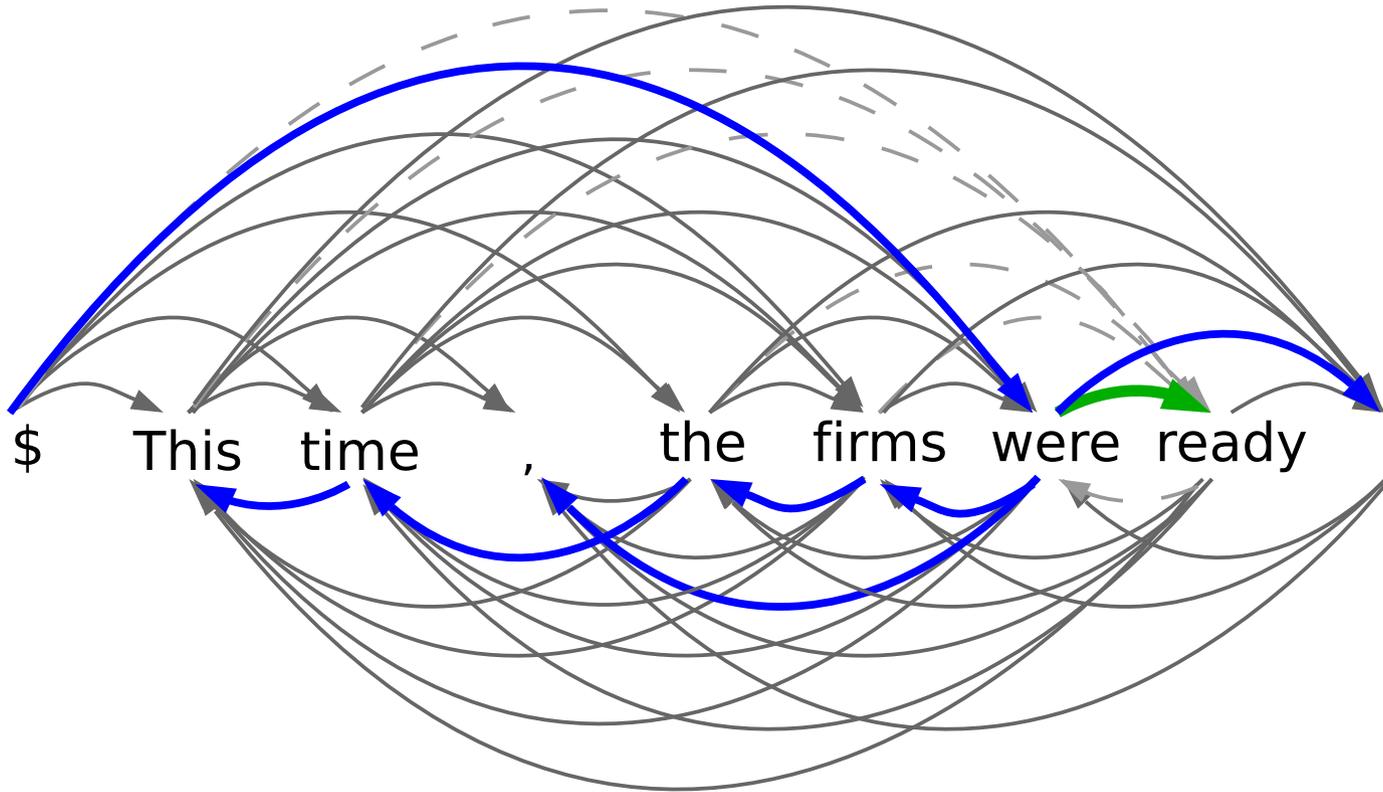- - Loser edge

*Non-projective decoding to find new 1-best tree*

**42** gray edges with unknown fate...
**27** features per gray edge

$ This time , the firms were ready .

Undetermined edge
Current 1-best tree
Winner edge
(permanently in 1-best tree)
Loser edge

*Classifier picks winners among the blue edges*

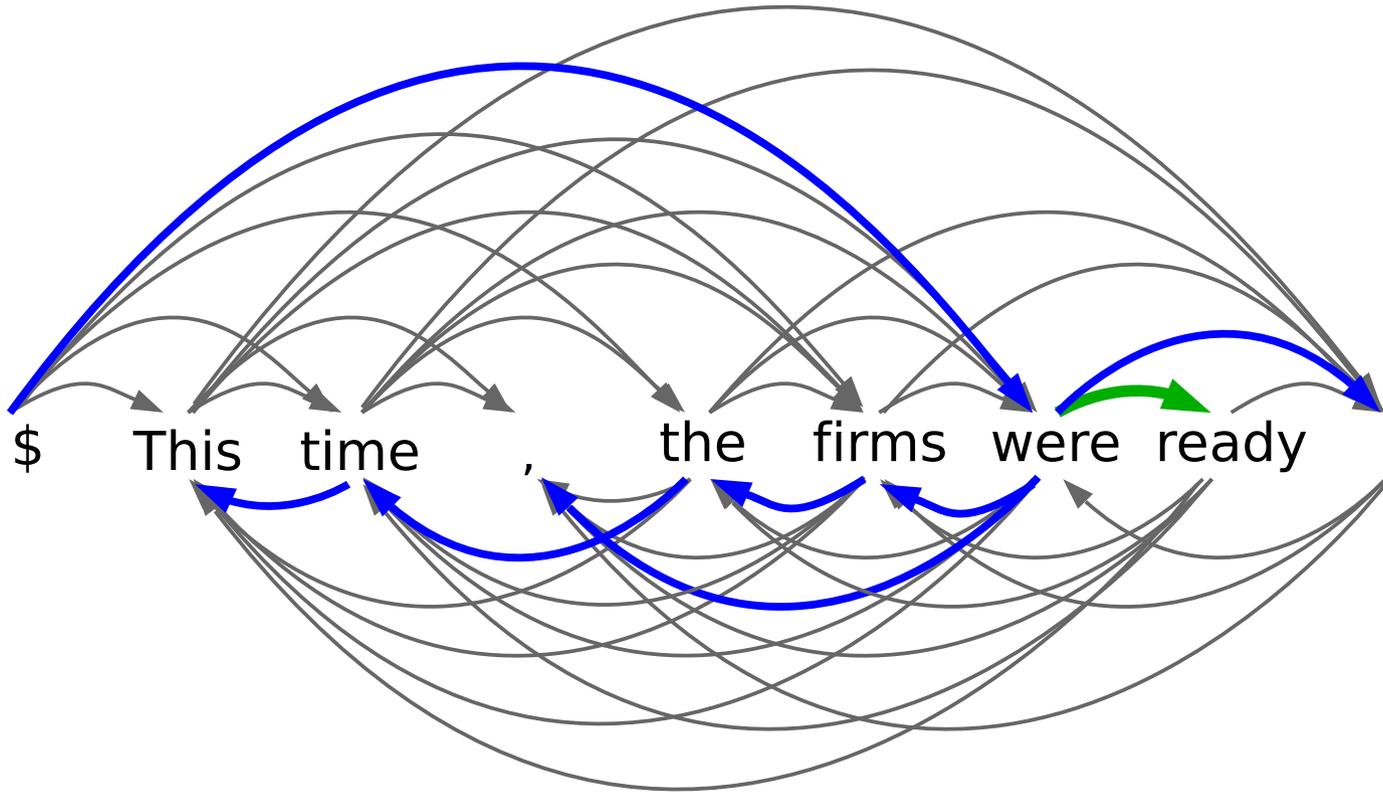31 gray edges with unknown fate...
27 features per gray edge



$ This time , the firms were ready .

Undetermined edge

Current 1-best tree

Winner edge
(permanently in 1-best tree)

Loser edge

*Remove losers in conflict with the winners*

**31** gray edges with unknown fate...
**27** features per gray edge



$ This time , the firms were ready .

— Undetermined edge
— Current 1-best tree
— Winner edge
(permanently in 1-best tree)
- - Loser edge

*Remove losers in conflict with the winners*

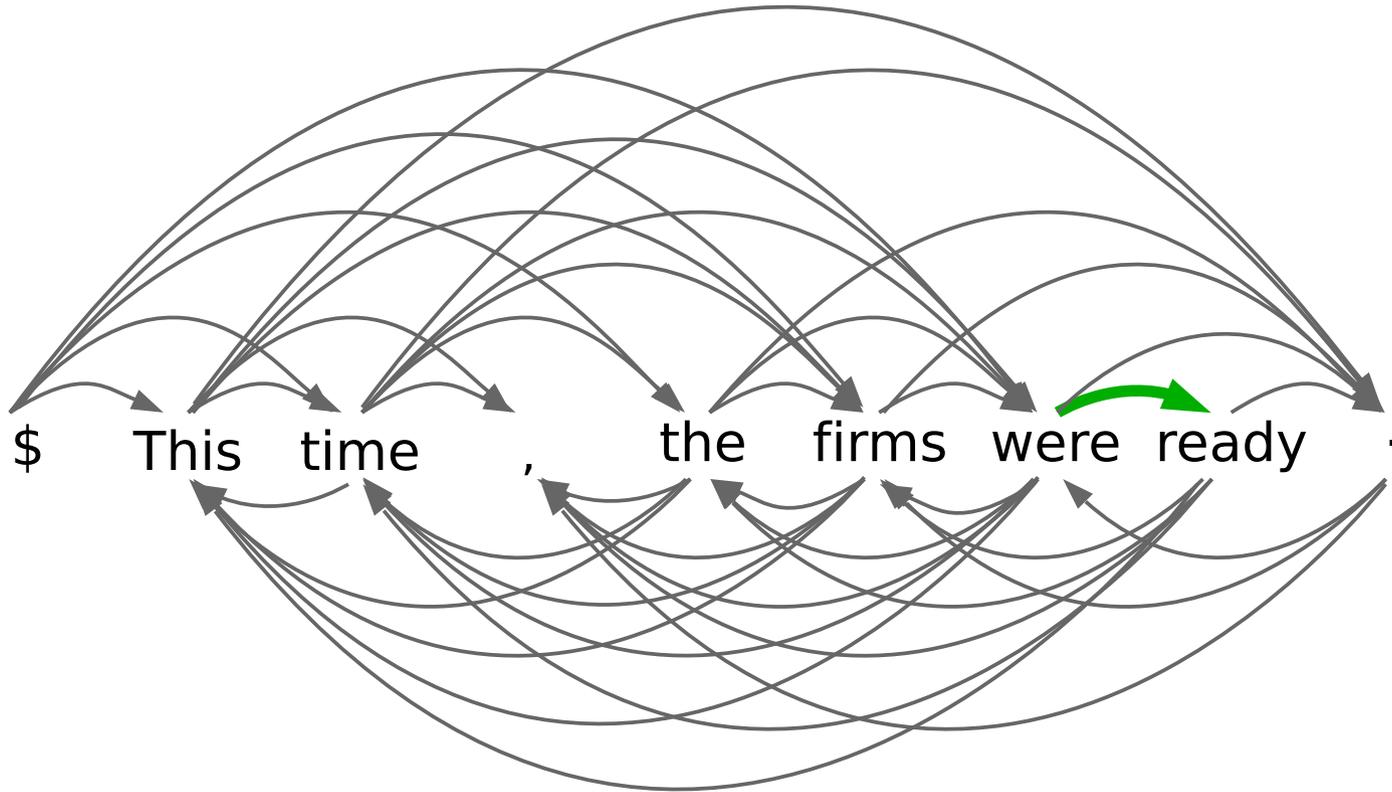40

+ next feature group     31     gray edges with unknown fate...
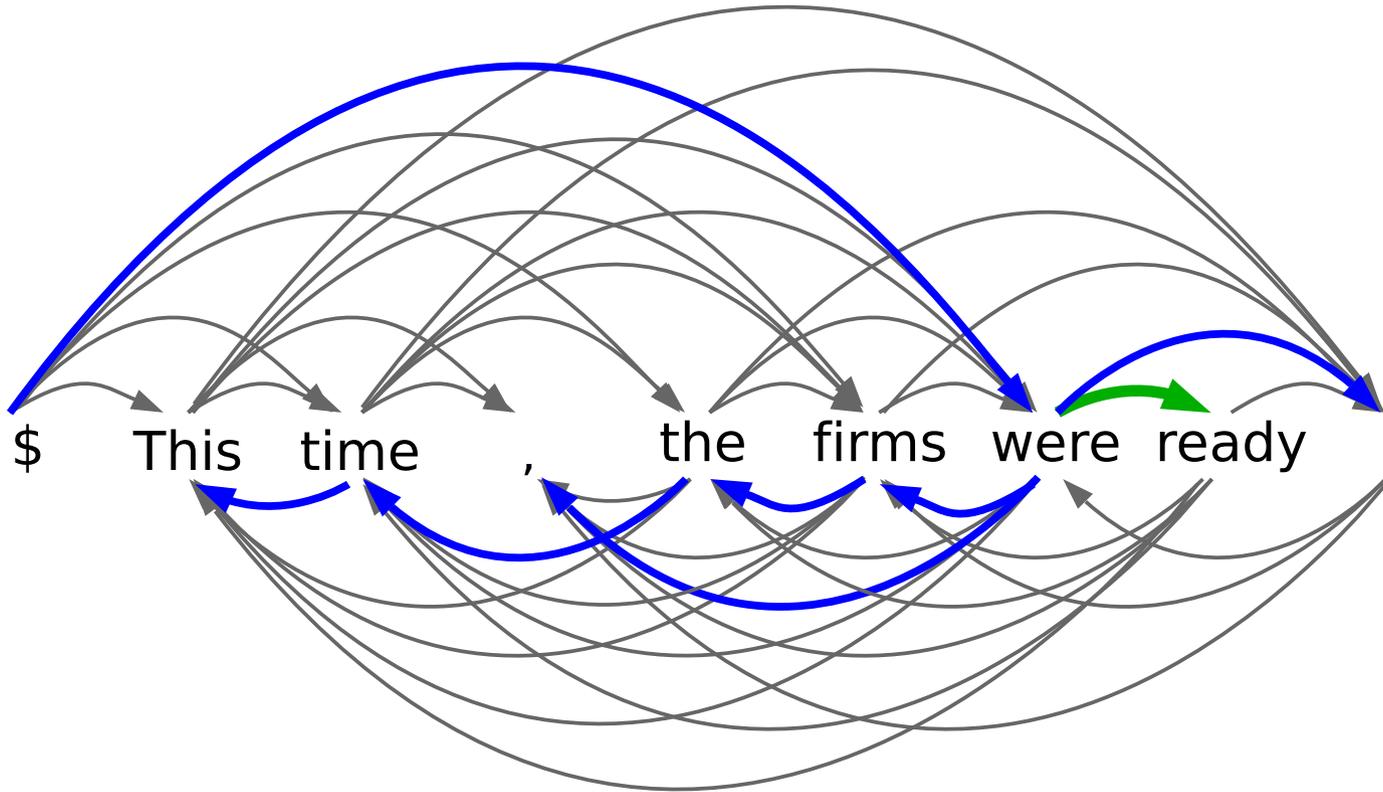                         74     features per gray edge

$   This   time   ,   the   firms   were   ready   .

—— Undetermined edge
—— Current 1-best tree
—— Winner edge
(permanently in 1-best tree)
- - Loser edge

31 gray edges with unknown fate...
74 features per gray edge

$ This time , the firms were ready .

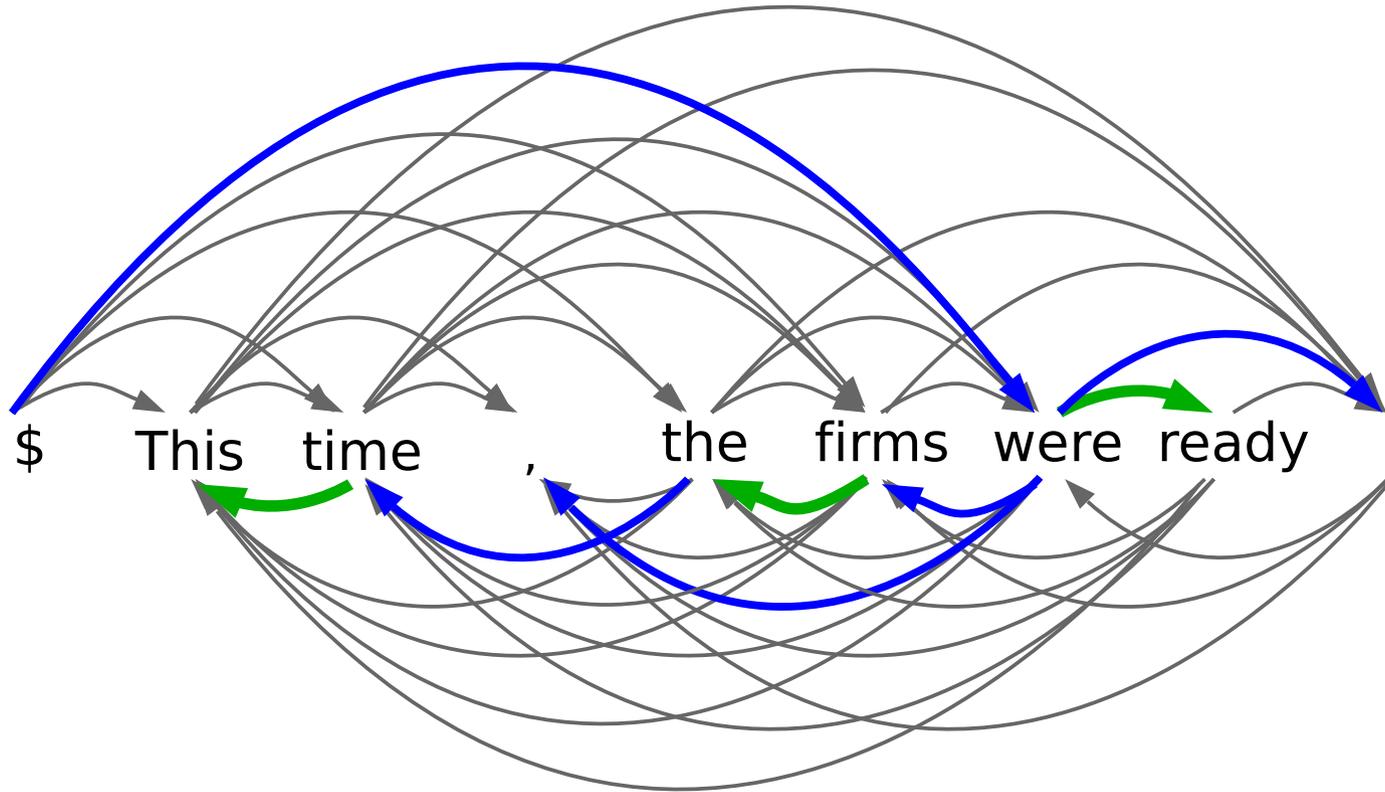— Undetermined edge
— Current 1-best tree
— Winner edge
(permanently in 1-best tree)
- - Loser edge

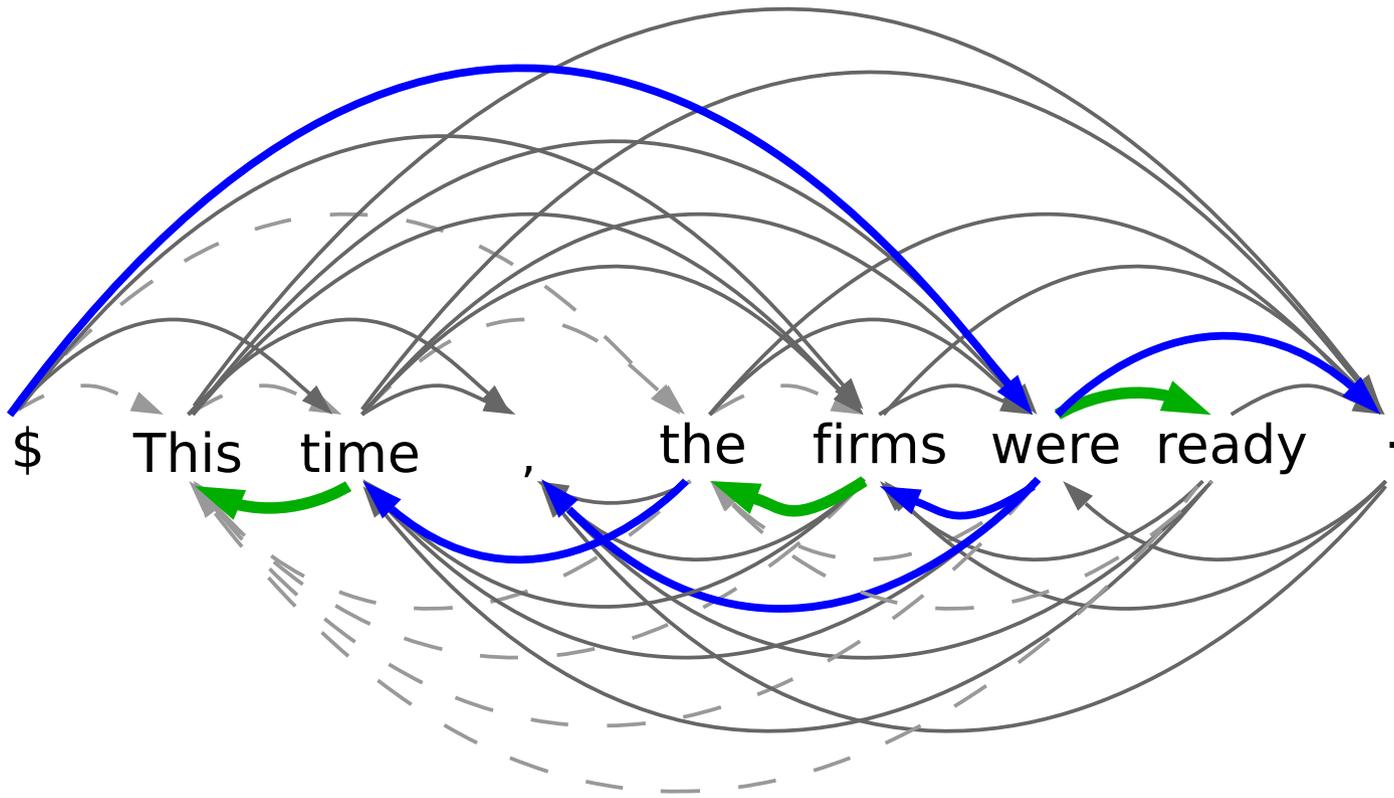*Non-projective decoding to find new 1-best tree*

42

**28** gray edges with unknown fate...
**74** features per gray edge

$ This time , the firms were ready .

—— Undetermined edge
—— Current 1-best tree
—— Winner edge
(permanently in 1-best tree)
- - Loser edge

*Classifier picks winners among the blue edges*

43

**8** gray edges with unknown fate...
**74** features per gray edge
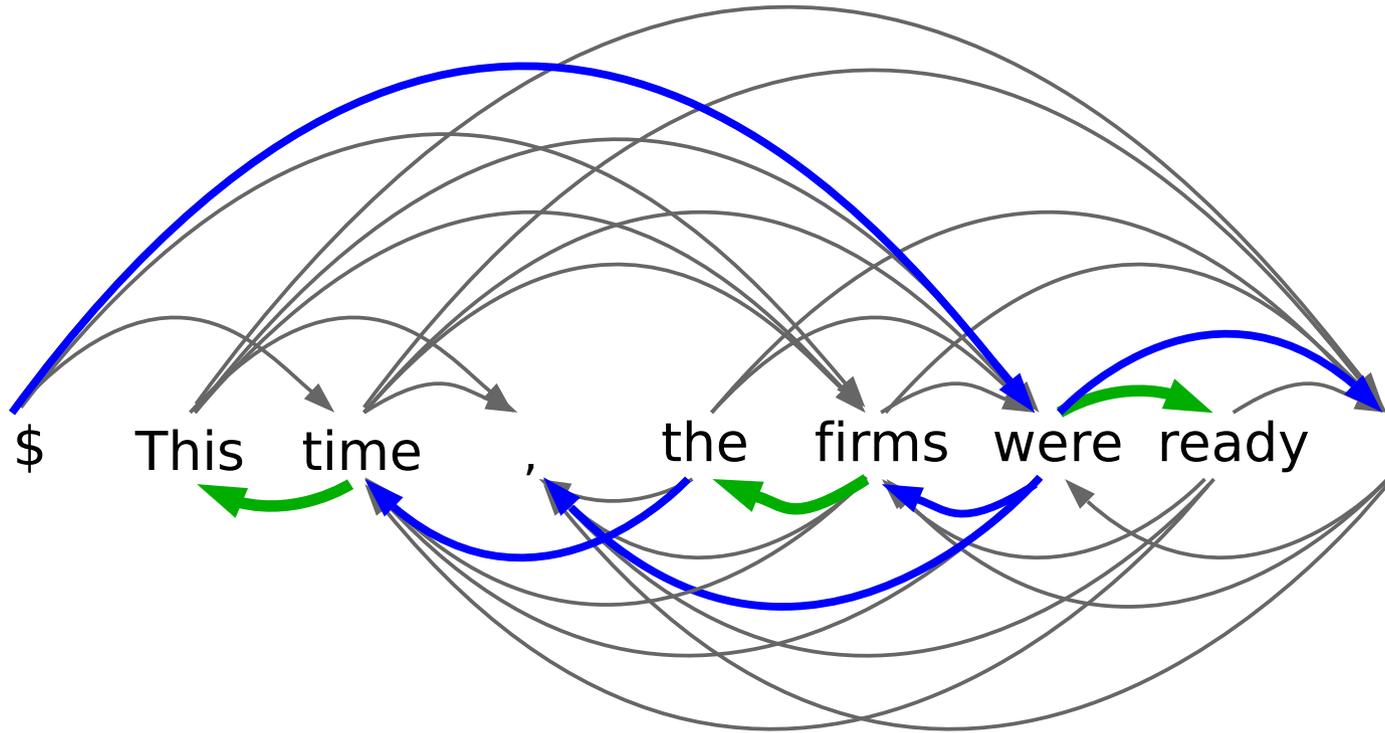


$ This time , the firms were ready .

— Undetermined edge
— Current 1-best tree
— Winner edge
(permanently in 1-best tree)
- - Loser edge

*Remove losers in conflict with the winners*

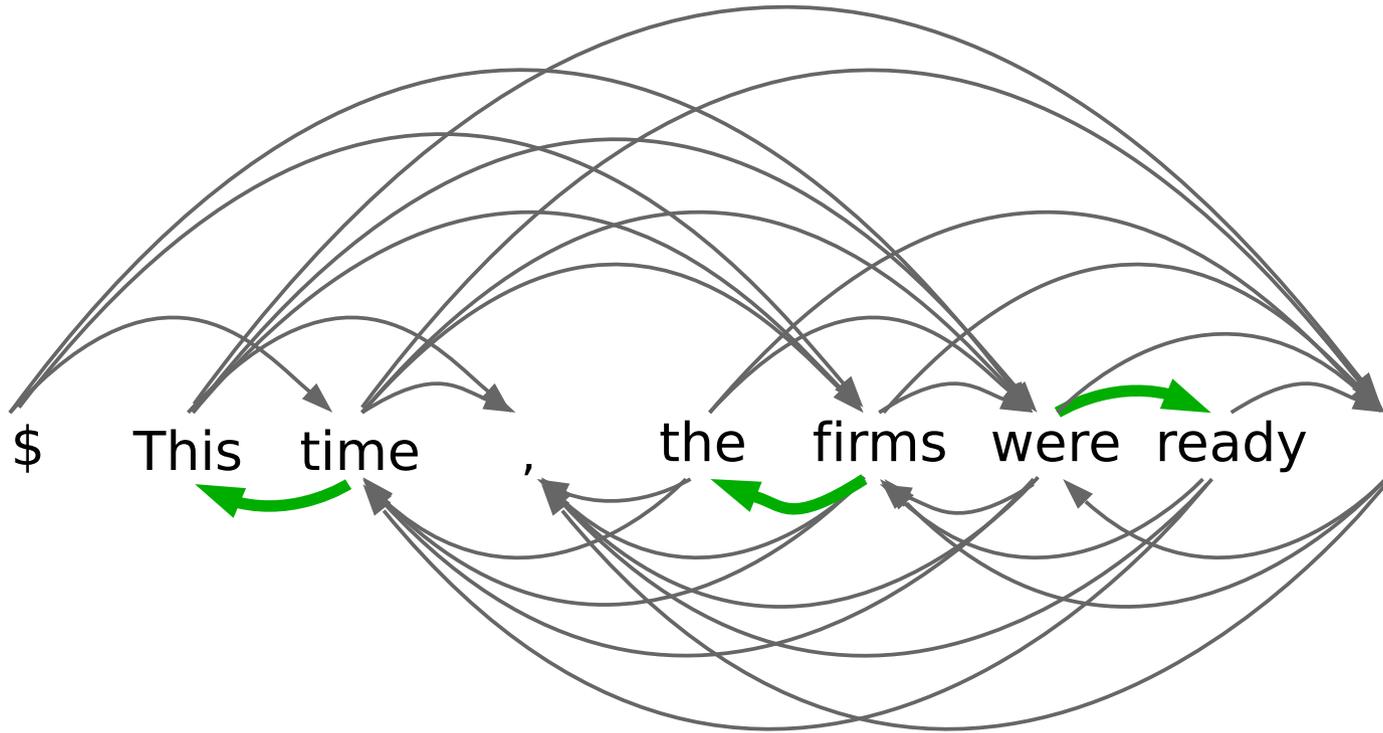8    gray edges with unknown fate...
74   features per gray edge

$    This   time    ,    the   firms   were   ready   .

— Undetermined edge
— Current 1-best tree
— Winner edge
  (permanently in 1-best tree)
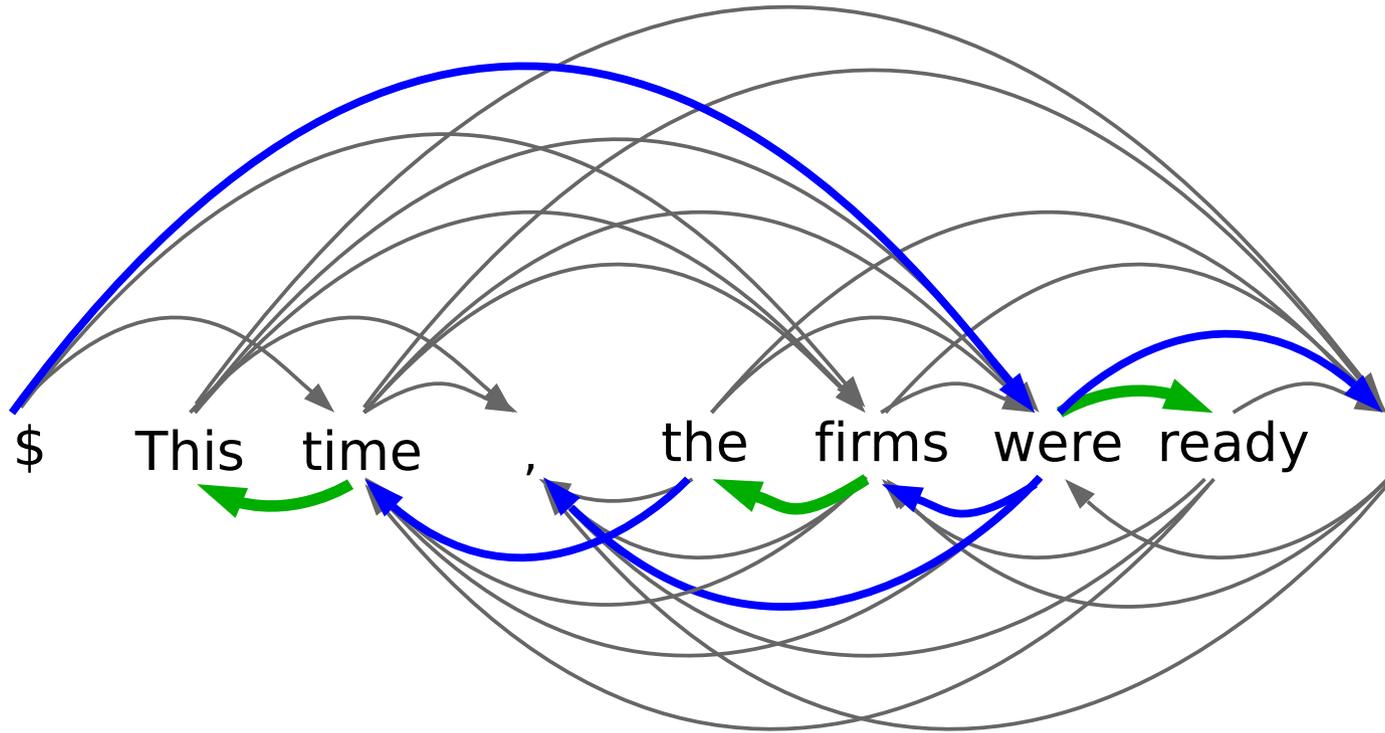- - Loser edge

*Remove losers in conflict with the winners*

45

+ next feature group   8   gray edges with unknown fate...
107   features per gray edge

$ This time , the firms were ready .

—— Undetermined edge
—— Current 1-best tree
—— Winner edge
(permanently in 1-best tree)
- - Loser edge

46

**8** gray edges with unknown fate...

**107** features per gray edge



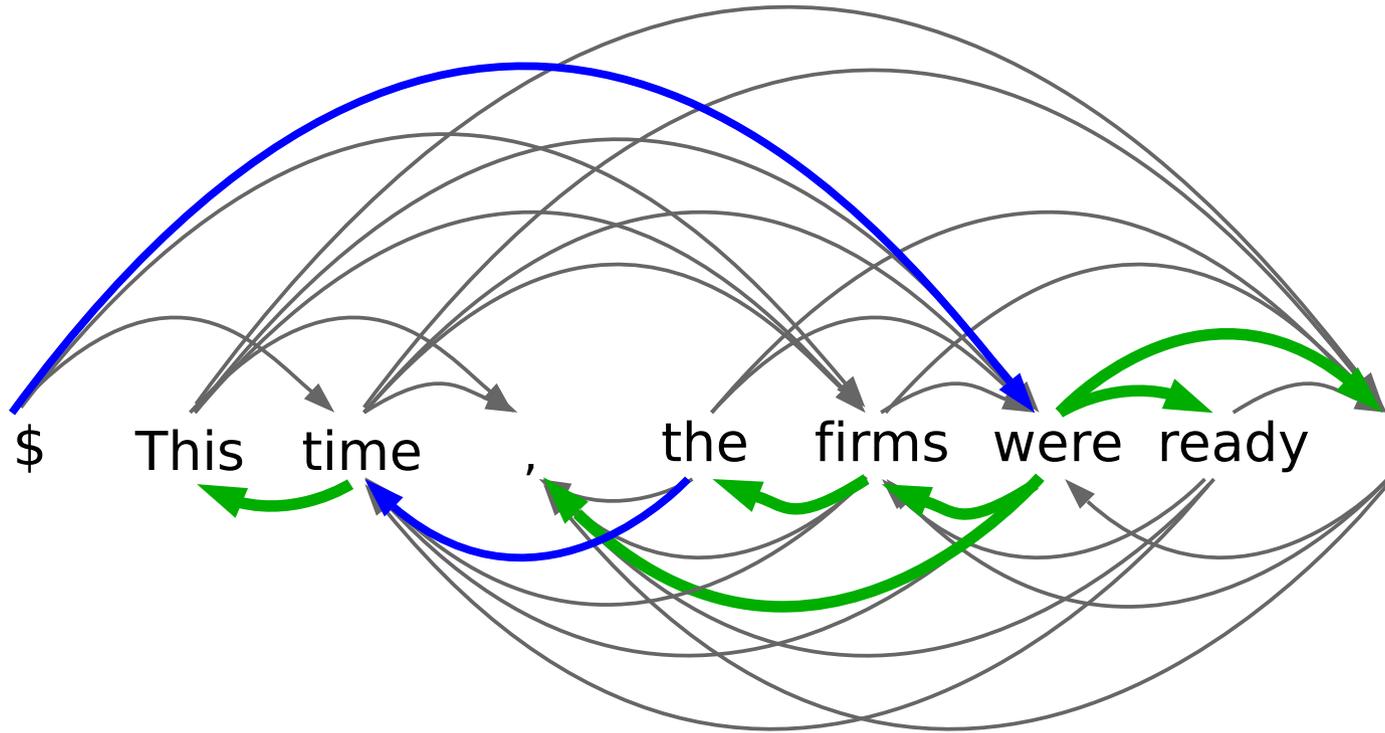$   This   time   ,   the   firms   were   ready   .

— Undetermined edge

— Current 1-best tree

— Winner edge
(permanently in 1-best tree)

- - Loser edge

*Non-projective decoding to find new 1-best tree*

**7** gray edges with unknown fate...
**107** features per gray edge
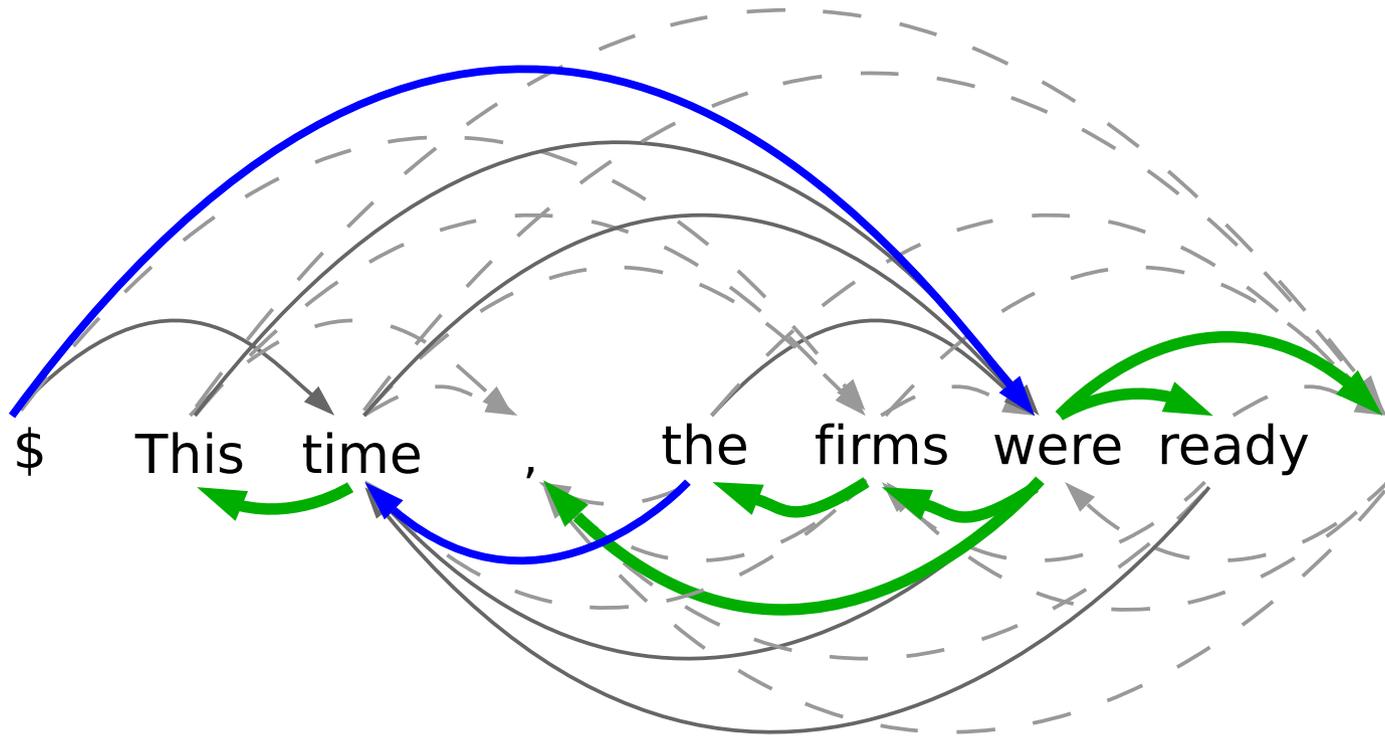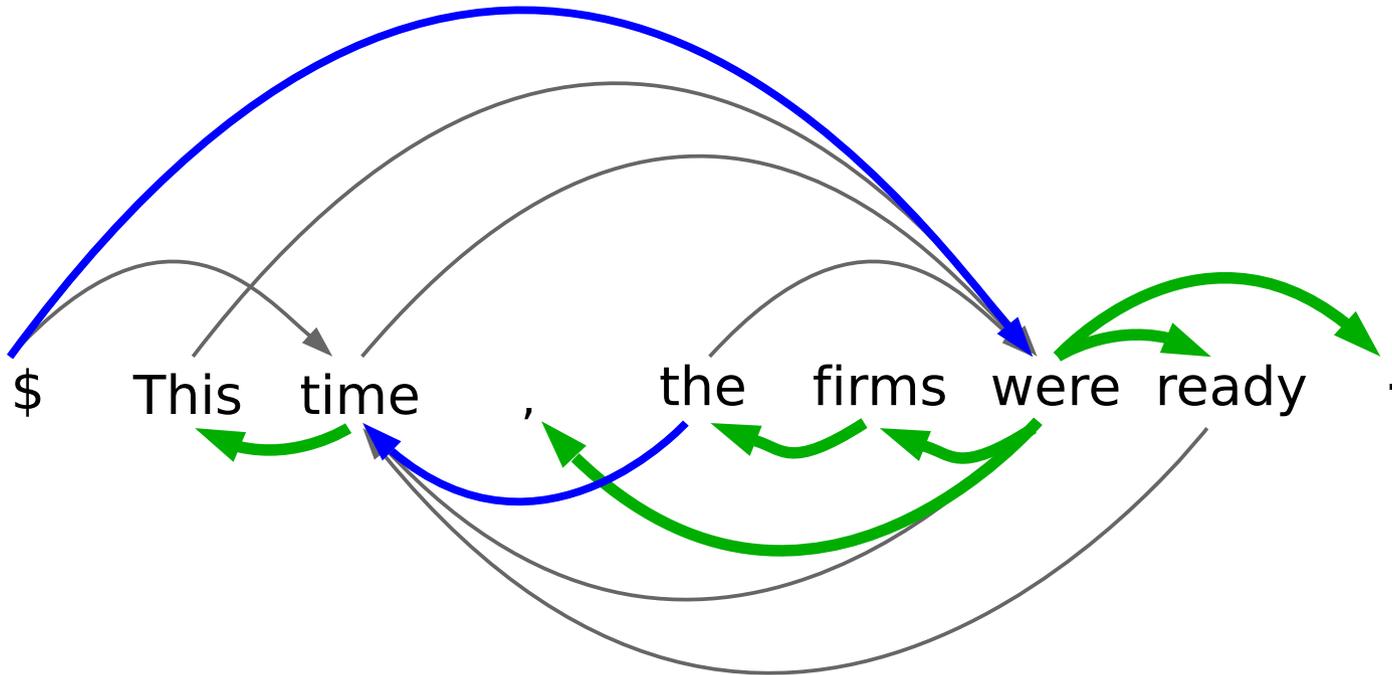
$ This time , the firms were ready .

— Undetermined edge
— Current 1-best tree
— Winner edge
(permanently in 1-best tree)
- - Loser edge

*Classifier picks winners among the blue edges*

48

3 gray edges with unknown fate...
107 features per gray edge

$ This time , the firms were ready .

Undetermined edge
Current 1-best tree
Winner edge
(permanently in 1-best tree)
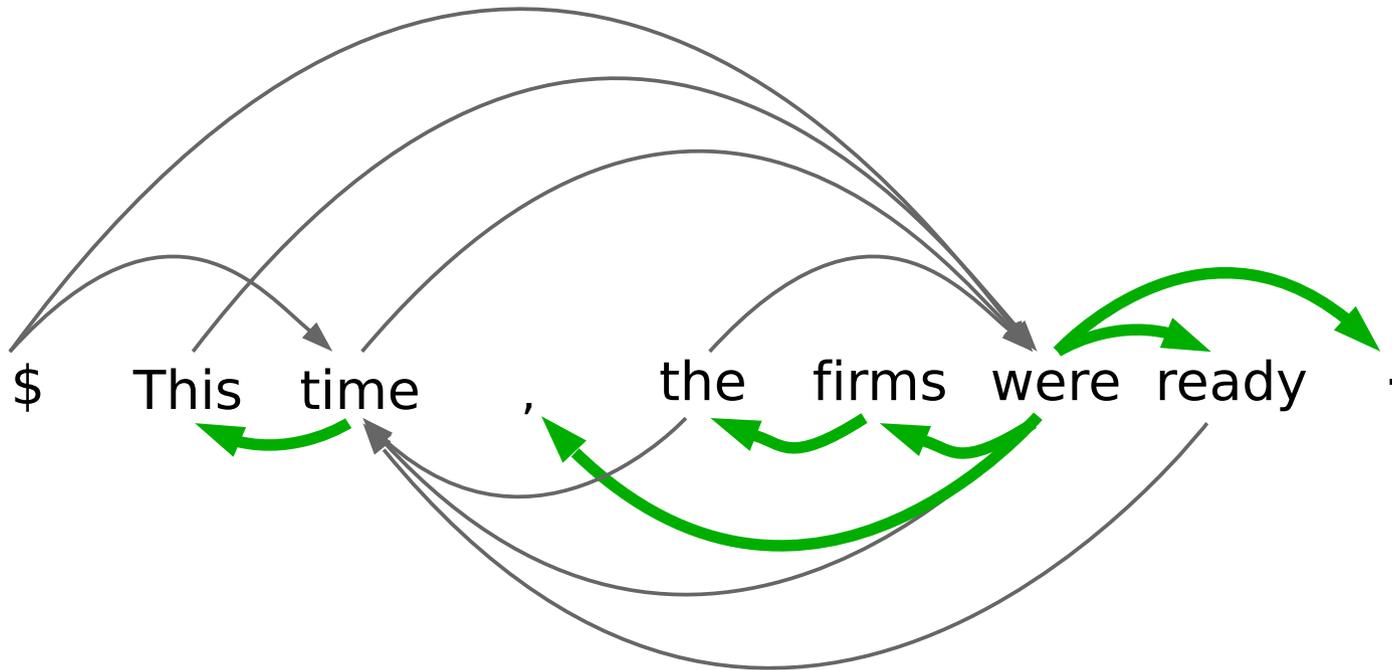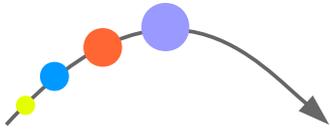Loser edge

*Remove losers in conflict with the winners*

49

3    gray edges with unknown fate...
107   features per gray edge

$ This time , the firms were ready .

── Undetermined edge
── Current 1-best tree
── Winner edge
(permanently in 1-best tree)
- - Loser edge

*Remove losers in conflict with the winners*

+ last feature group    3    gray edges with unknown fate...
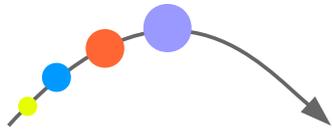                      268    features per gray edge



$   This   time   ,   the   firms   were   ready   .

—— Undetermined edge
—— Current 1-best tree
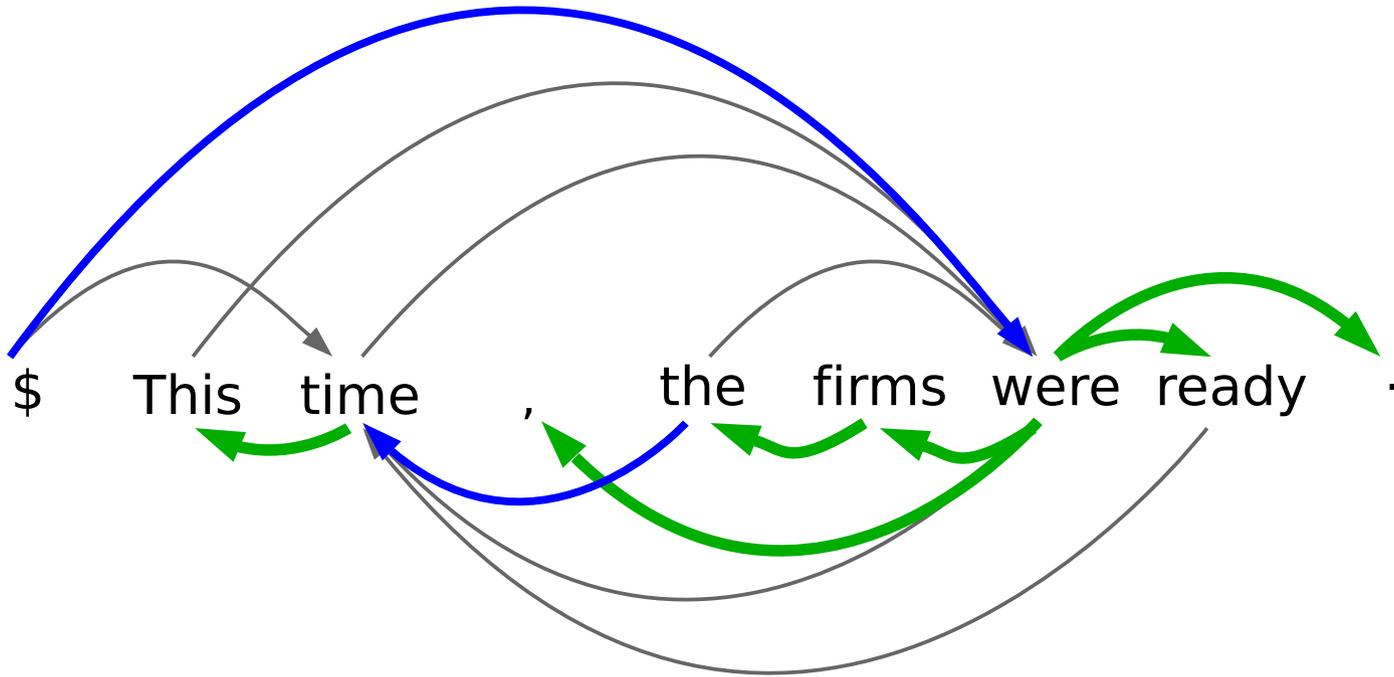—— Winner edge
(permanently in 1-best tree)
- - Loser edge

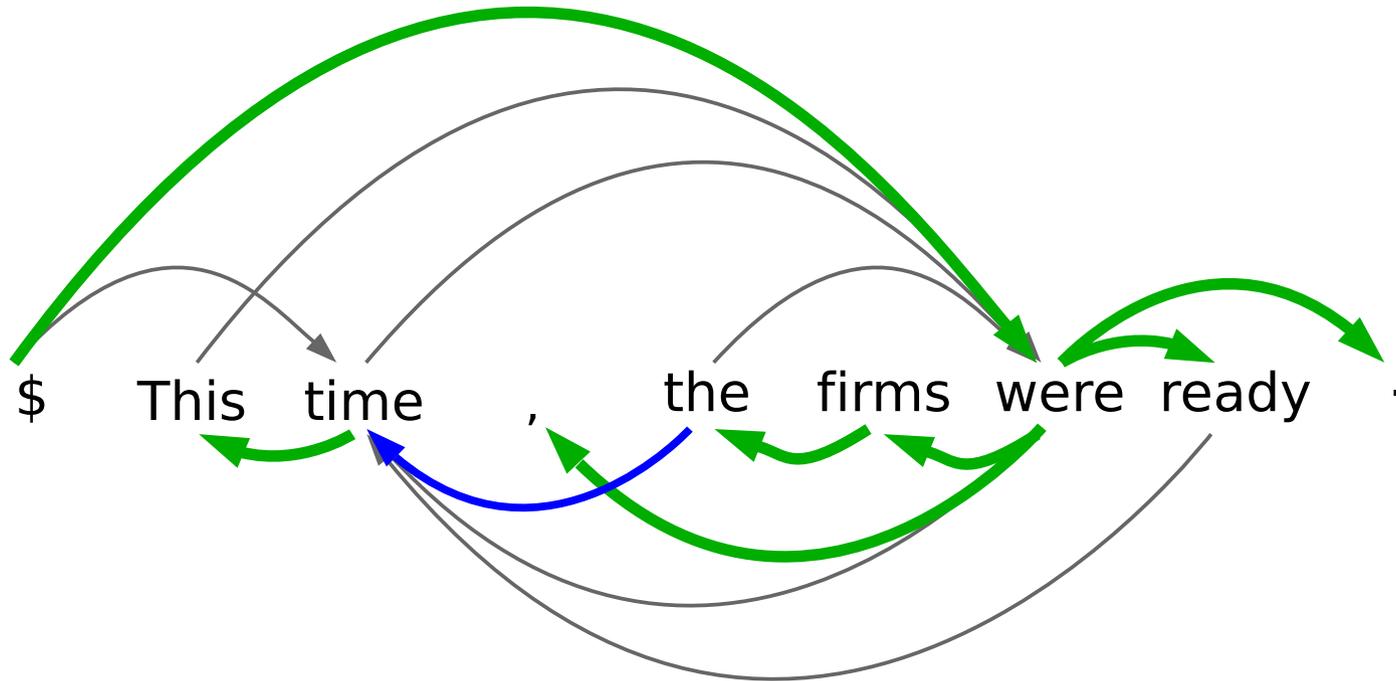0     gray edge with unknown fate...
268   features per gray edge

$    This   time   ,    the   firms   were   ready   .

Undetermined edge
Current 1-best tree
Winner edge
(permanently in 1-best tree)
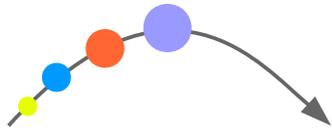Loser edge

*Projective decoding to find final 1-best tree*

52

# What Happens During the Average Parse?

# What Happens During the Average Parse?

# What Happens During the Average Parse?

# What Happens During the Average Parse?



Later features are helpful

Most edges win
or lose early

| | | runtime % |
| | | UAS % |
| | | remaining edge % |
| | | winner edge % |

Some edges win late

# What Happens During the Average Parse?



Later features are helpful

Most edges win or lose early

Linear increase in runtime

Some edges win late

Legend:
- runtime %
- UAS %
- remaining edge %
- winner edge %

# Summary: How Early Decisions Are Made

- Winners 🙂
  - Will definitely appear in the 1-best tree
- Losers ☹
  - Have the same child as a winning edge
  - Form cycle with winning edges
  - Cross a winning edge *(optional)*
  - Share root ($) with a winning edge *(optional)*
- Undetermined
  - Add the next feature group to the remaining gray edges

# Feature Template Ranking

- Forward selection

# Feature Template Ranking

- Forward selection

<span style="color:red">A   0.60</span>
B   0.49
C   0.55

# Feature Template Ranking

- Forward selection

<div>

A   0.60

B   0.49      A →

C   0.55

</div>

1   A

# Feature Template Ranking

- Forward selection

A   0.60
B   0.49
C   0.55

A →

A&B   0.80
A&C   0.85

1   A

# Feature Template Ranking

- Forward selection

A   0.60
B   0.49  →(A)→  A&B  0.80  →(C)→
C   0.55                 A&C  0.85

```
1    A
2    C
```

# Feature Template Ranking

- Forward selection

A   0.60
B   0.49    A   →   A&B   0.80    C   →   A&C&B   0.9
C   0.55                A&C   0.85

| 1 | A |
| 2 | C |
| 3 | B |

# Feature Template Ranking

- ## Forward selection

A   0.60
B   0.49          A          A&B   0.80          C          A&C&B   0.9
C   0.55        →           A&C   0.85        →

```
1    A
2    C
3    B
```

- ## Grouping

head cPOS+ mod cPOS + in-between punct #          0.49
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
in-between cPOS                                                        0.59
⋮
head POS + mod POS + in-between conj #          0.71
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
head POS + mod POS + in-between POS + dist   0.72
⋮
head token + mod cPOS + dist                              0.80
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
⋮

# Feature Template Ranking

- Forward selection

A   0.60
B   0.49    A →
C   0.55

A&B  0.80    C →
A&C  0.85

A&C&B  0.9

1   A
2   C
3   B

- Grouping

| | |
|---|---|
| head cPOS+ mod cPOS + in-between punct # | 0.49 |
| in-between cPOS | 0.59 |
| head POS + mod POS + in-between conj # | 0.71 |
| head POS + mod POS + in-between POS + dist | 0.72 |
| head token + mod cPOS + dist | 0.80 |

$+ \sim 0.1$

$+ \sim 0.1$

# Partition Feature List Into Groups

# How to pick the winners?

# How to pick the winners?

- Learn a classifier

# How to pick the winners?

- Learn a classifier

- Features

  - Currently added parsing features
  - Meta-features -- confidence of a prediction

# How to pick the <span style="color:green">winners</span>?

- Learn a classifier
- Features
  - Currently added parsing features
  - Meta-features -- confidence of a prediction
- Training examples
  - Input: each <span style="color:blue">blue edge</span> in current 1-best tree
  - Output: is the edge in the gold tree? If so, we want it to <span style="color:green">win</span>!

# Classifier Features

- Currently added parsing features

- Meta-features

  - the    firms    : …, 0.5, 0.8, 0.85

    (scores are normalized by the sigmoid function)

  - Margins to the highest-scoring competing edge

    0.72
    0.65    0.30
    0.23
    $    This    time        the    firms    were        .
    0.12

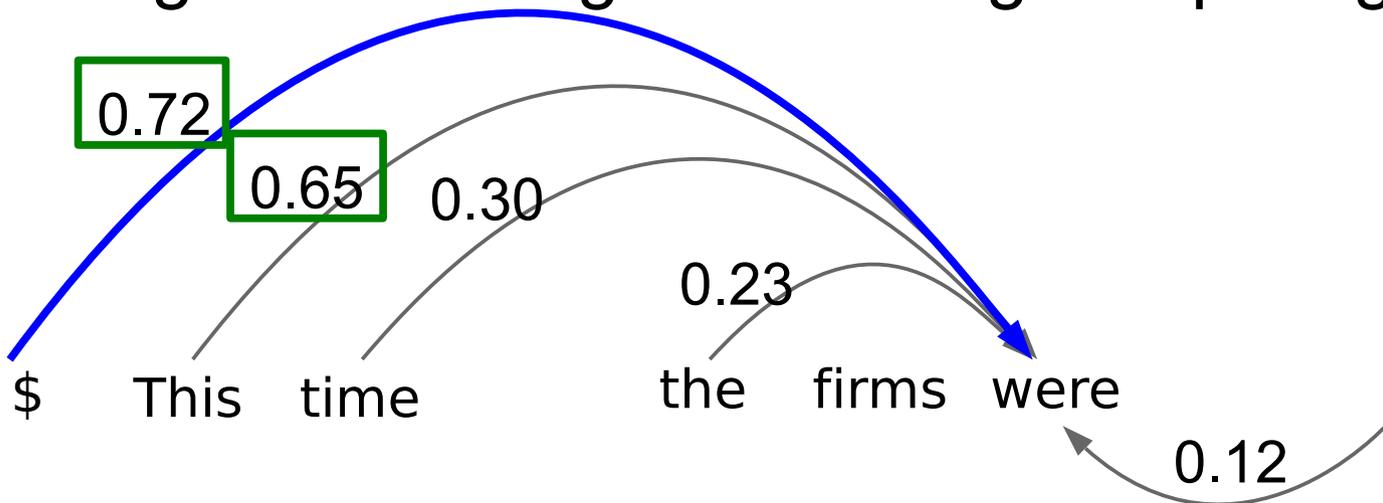  - Index of the next feature group

# Classifier Features

- Currently added parsing features

- Meta-features

  - the   firms   : …, 0.5, 0.8, 0.85

    (scores are normalized by the sigmoid function)
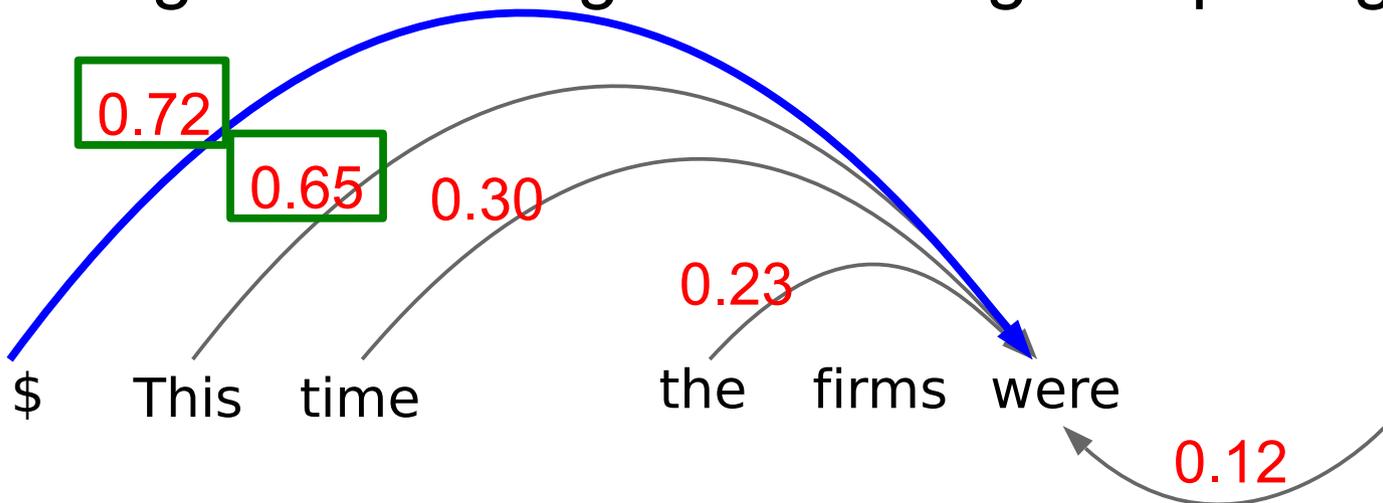
  - Margins to the highest-scoring competing edge

    0.72   0.65   0.30   0.23

    $    This   time         the   firms  were        .

    0.12

  - Index of the next feature group

# Classifier Features

- Currently added parsing features

- Meta-features                              **Dynamic Features**
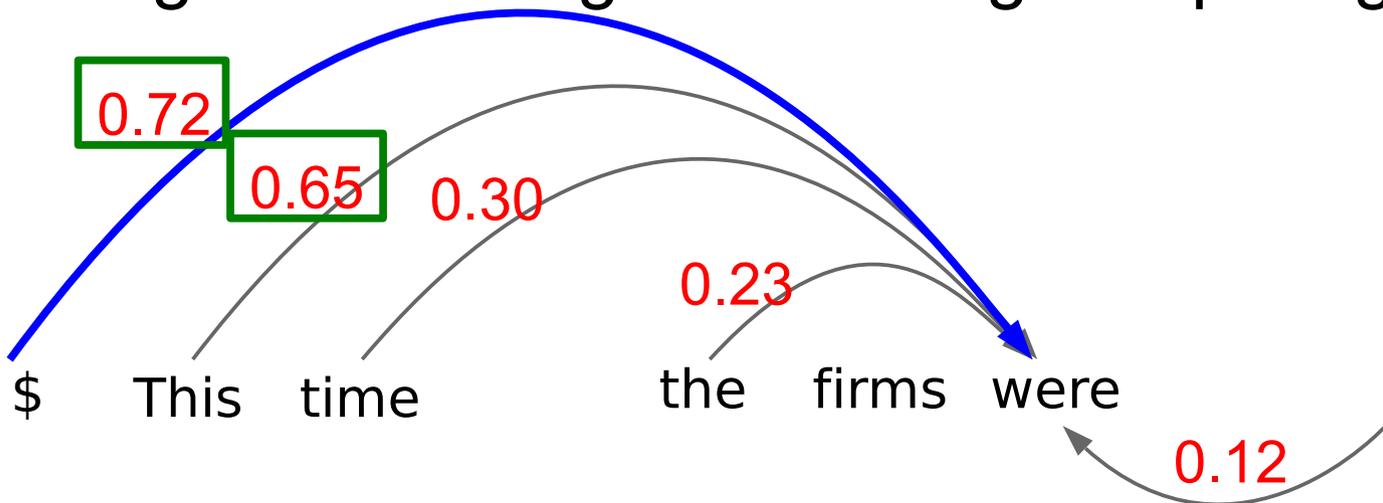
  – the    firms    : …, 0.5, 0.8, 0.85

    (scores are normalized by the sigmoid function)

  – Margins to the highest-scoring competing edge

    0.72

       0.65   0.30

               0.23

    $    This    time        the    firms    were           .

                                                    0.12

  – Index of the next feature group

# How To Train With Dynamic Features

- Training examples are not fixed in advance!

- Winners/losers from stages < k affect:

  - *Set* of edges to classify at stage k

  - The dynamic *features* of those edges at stage k

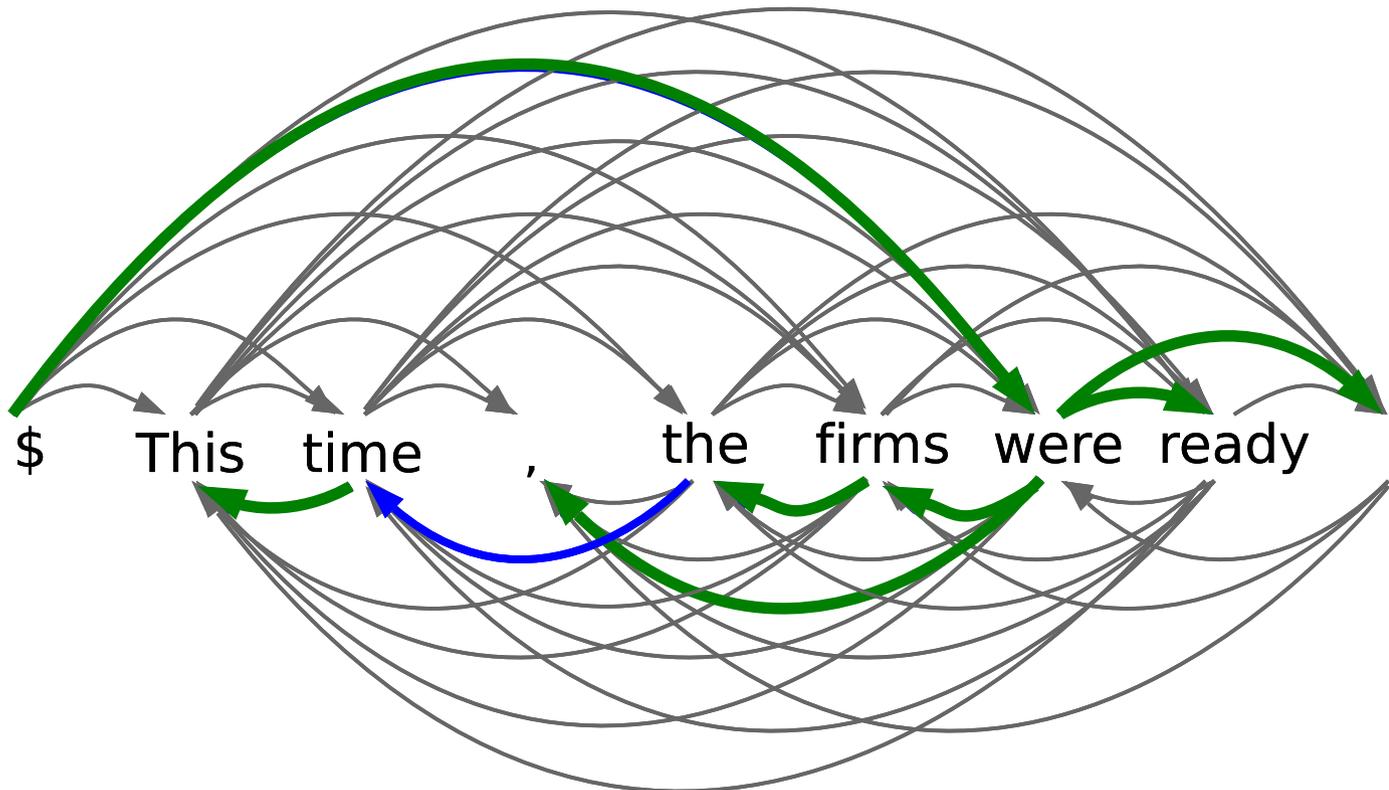- Bad decisions can cause future errors

# How To Train With Dynamic Features

- Training examples are not fixed in advance!!

- Winners/losers from stages < k affect:

  – *Set* of edges to classify at stage k

  – The dynamic *features* of those edges at stage k

- Bad decisions can cause future errors

  ## Reinforcement / Imitation Learning

- Dataset Aggregation (DAgger) (Ross et al., 2011)

  – Iterates between training and running a model

  – Learns to recover from past mistakes

# Upper Bound of Our Performance

- "Labels"

  - Gold edges always win

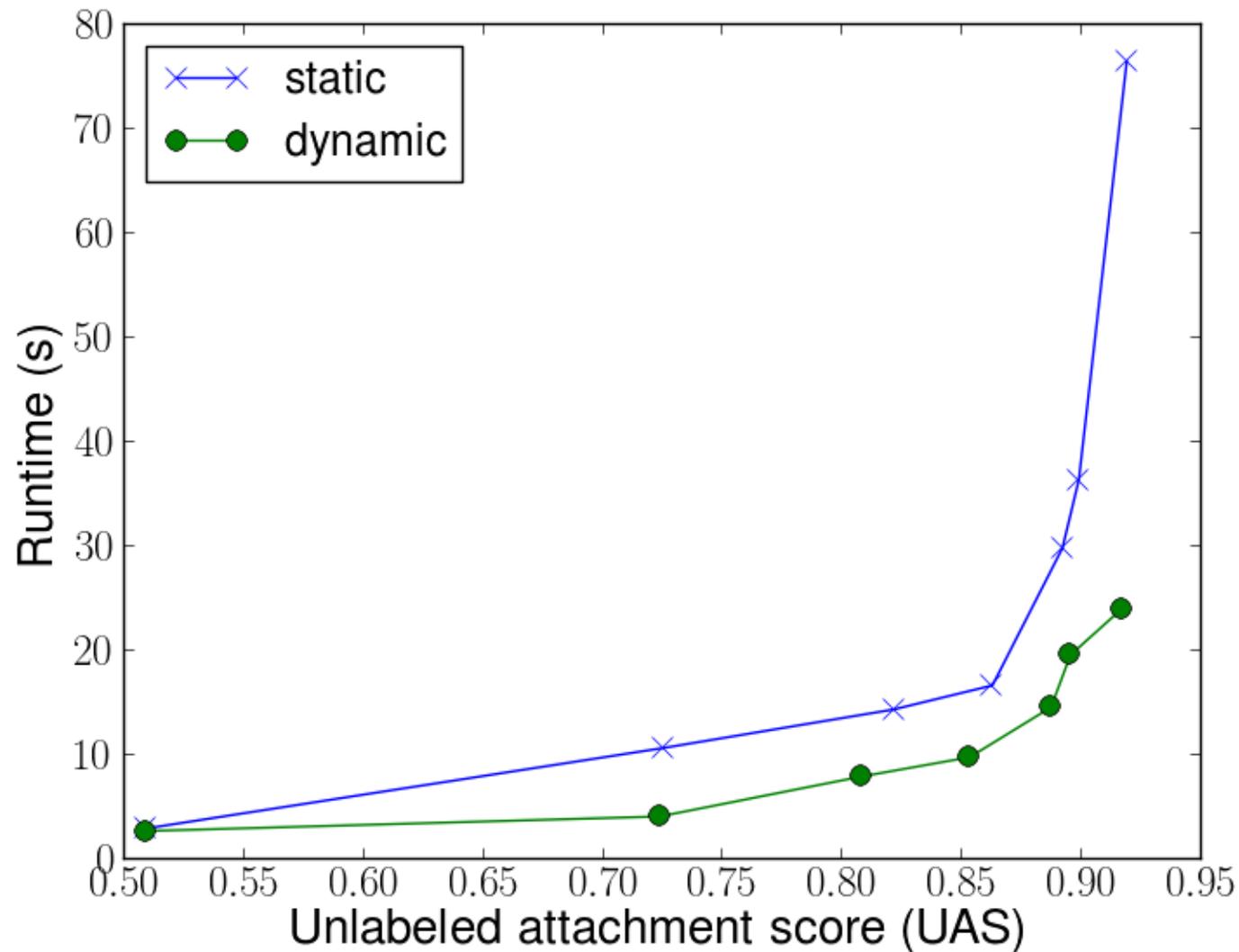  - 96.47% UAS with 2.9% first-order features

# How To Train Our Parser

1. Train parsers (non-projective, projective) *using all features*

2. Rank and group feature templates

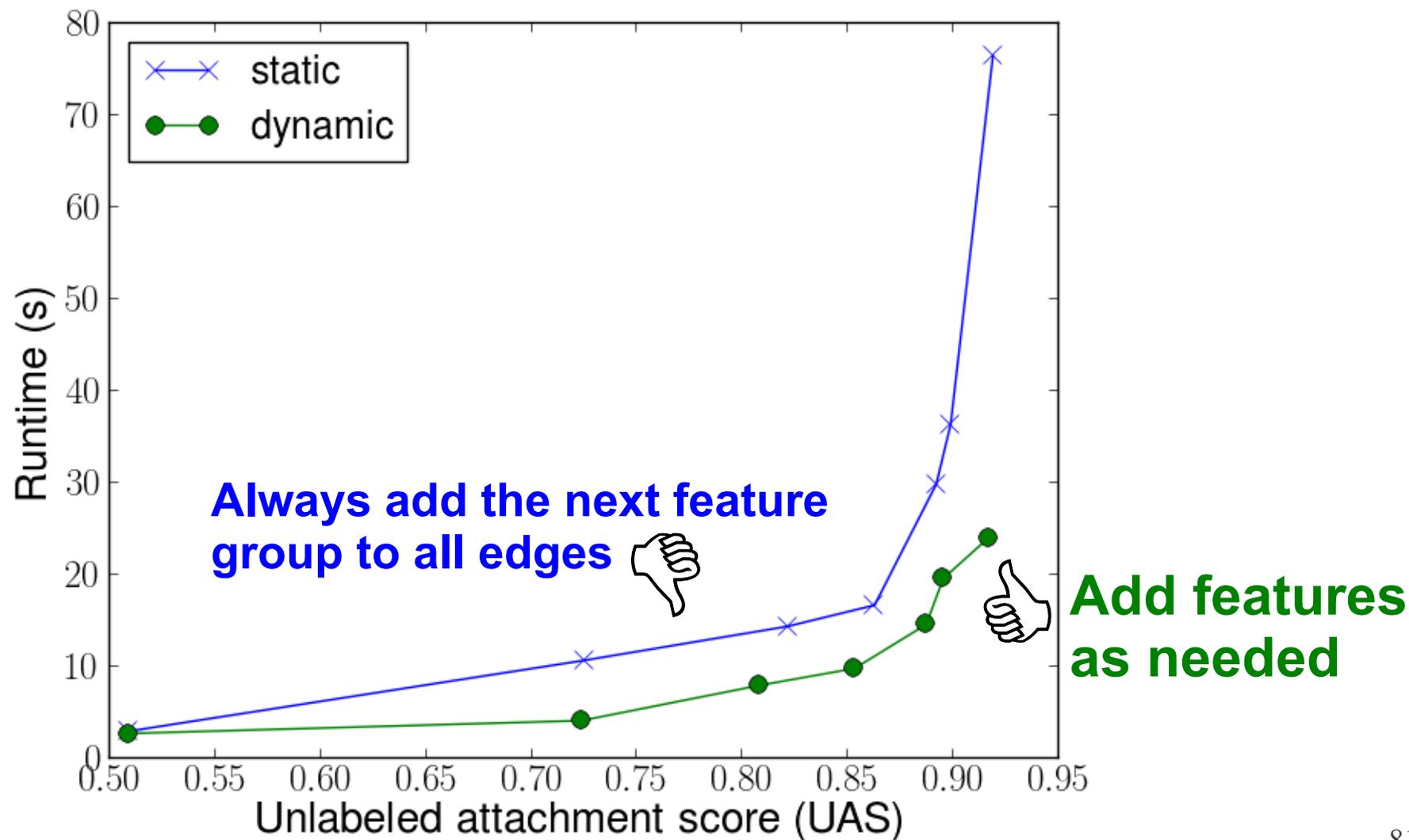3. Iteratively train a classifier to decide winners/losers

# Experiment

- Data
  - Penn Treebank: English
  - CoNLL-X: Bulgarian, Chinese, German, Japanese, Portuguese, Swedish

- Parser
  - MSTParser (McDonald et al., 2006)

- Dynamically-trained Classifier
  - LibLinear (Fan et al., 2008)

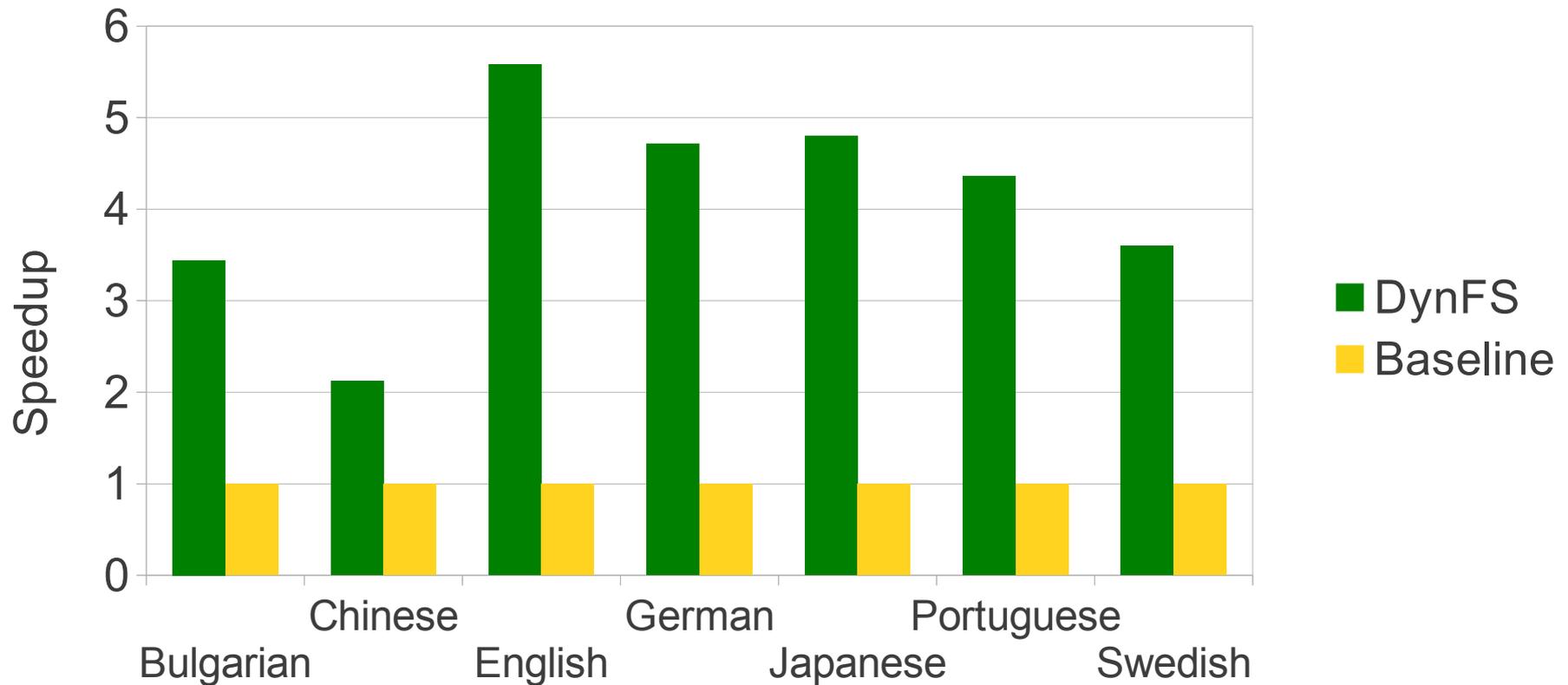# Dynamic Feature Selection Beats Static Forward Selection

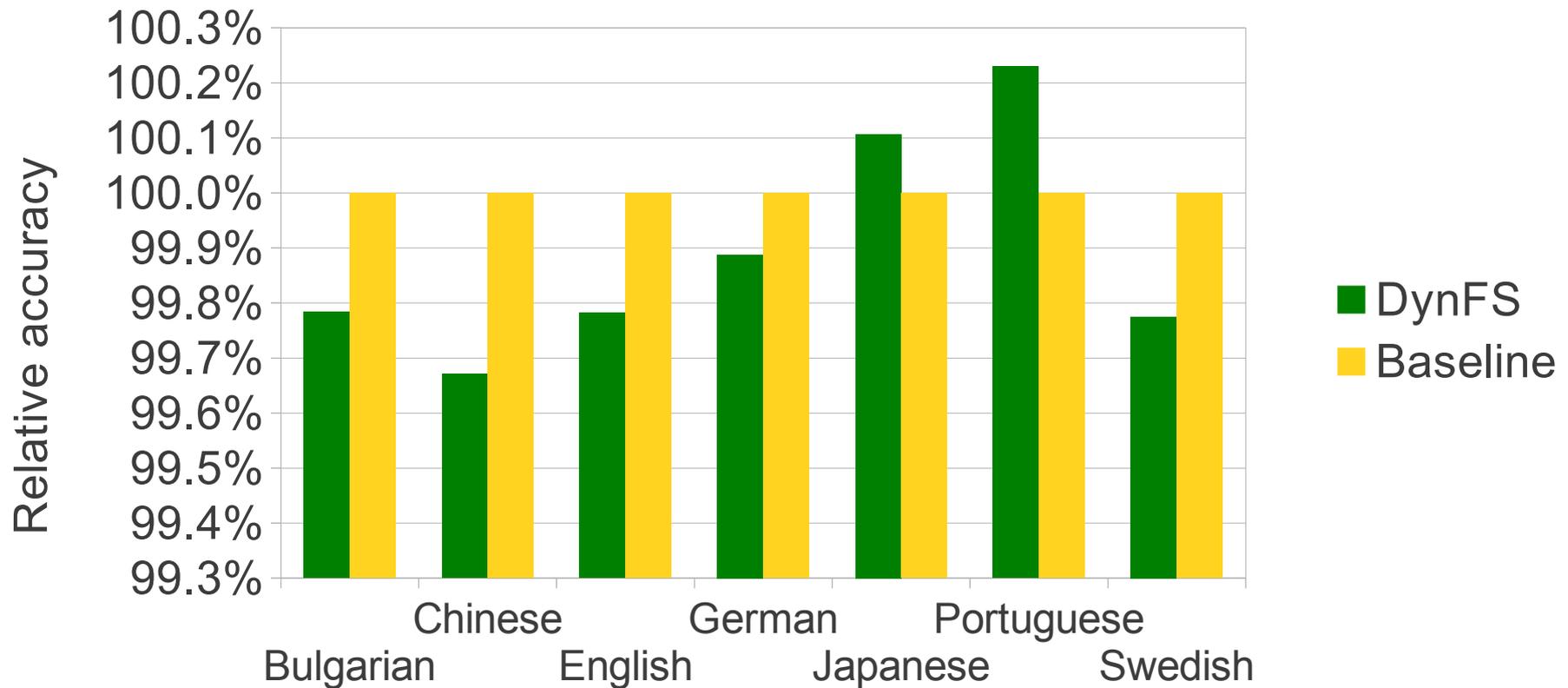# Dynamic Feature Selection Beats Static Forward Selection

# Experiment: 1st-order
# 2x to 6x speedup
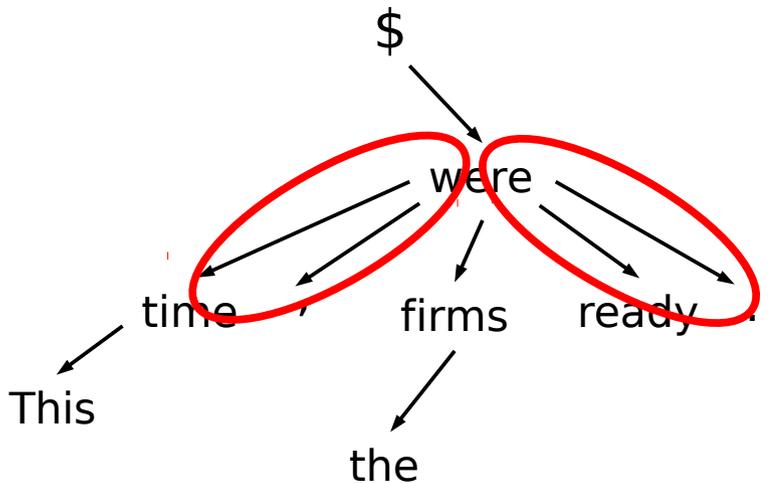
# Experiment: 1st-order
# ~0.2% loss in accuracy



$$relative\ accuracy = \frac{accuracy\ of\ the\ pruning\ parser}{accuracy\ of\ the\ full\ parser}$$

# Second-order Dependency Parsing

were  ready     .

$

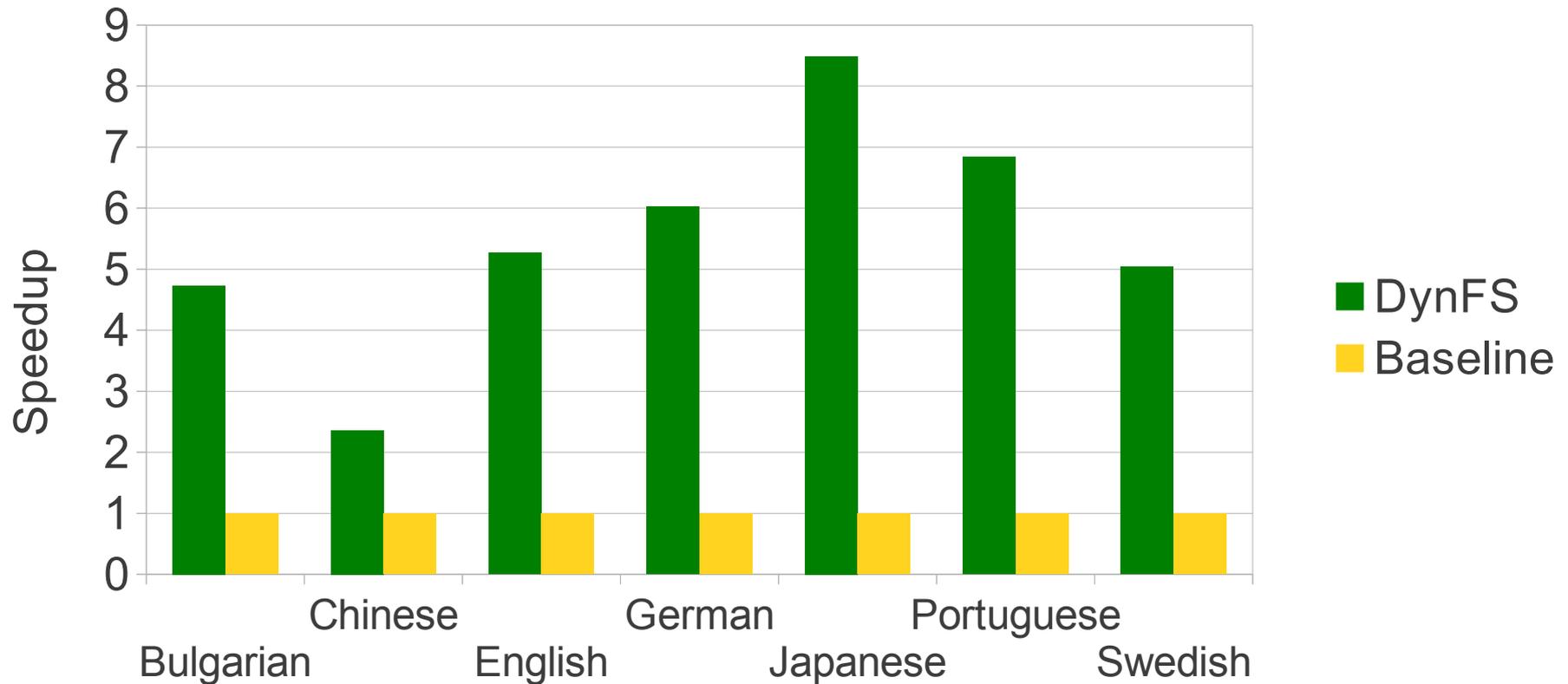were

time      firms      ready

This

the

- Features depend on the siblings as well

- First-order:
    - $O(n^2)$ substructure to score

- Second-order:
    - $O(n^3)$ substructure to score
      ~380 feature templates
      ~96M features
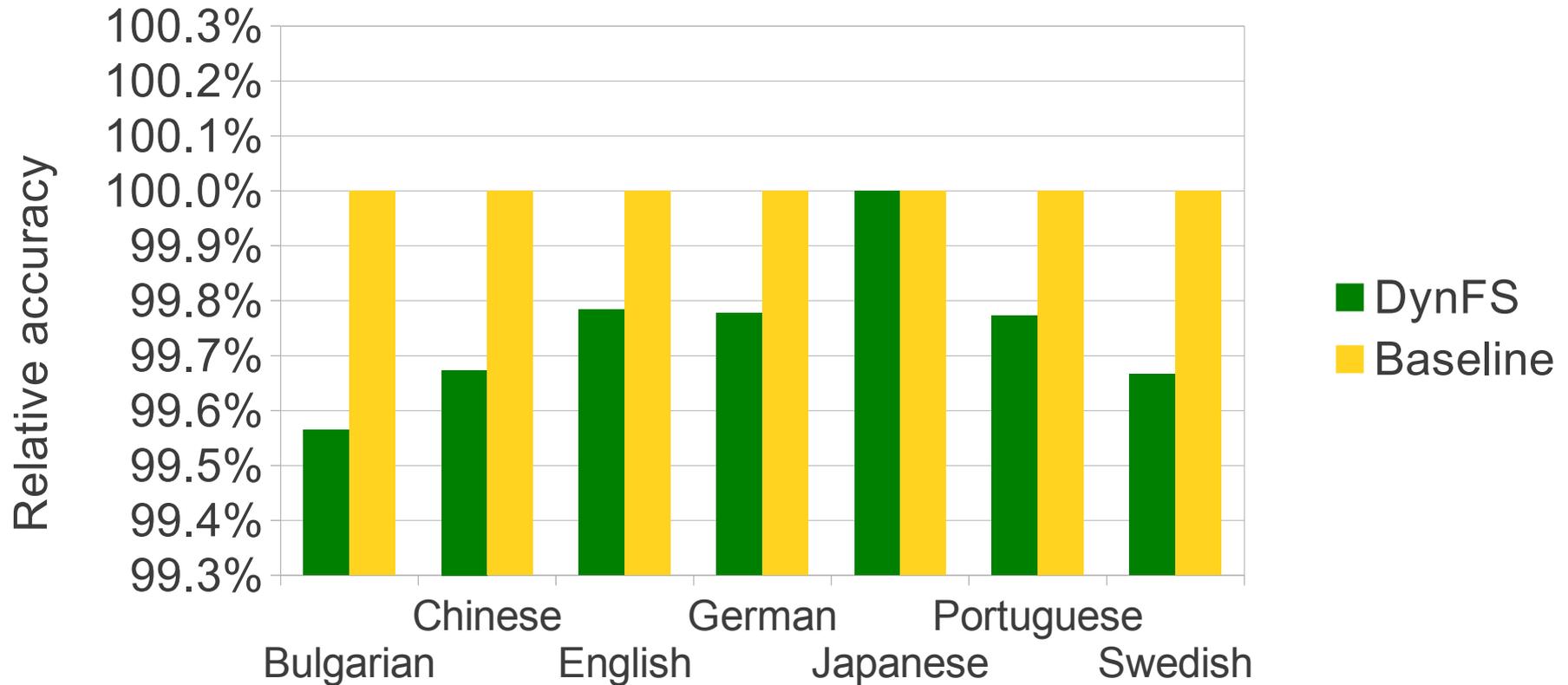
- Decoding: still $O(n^3)$
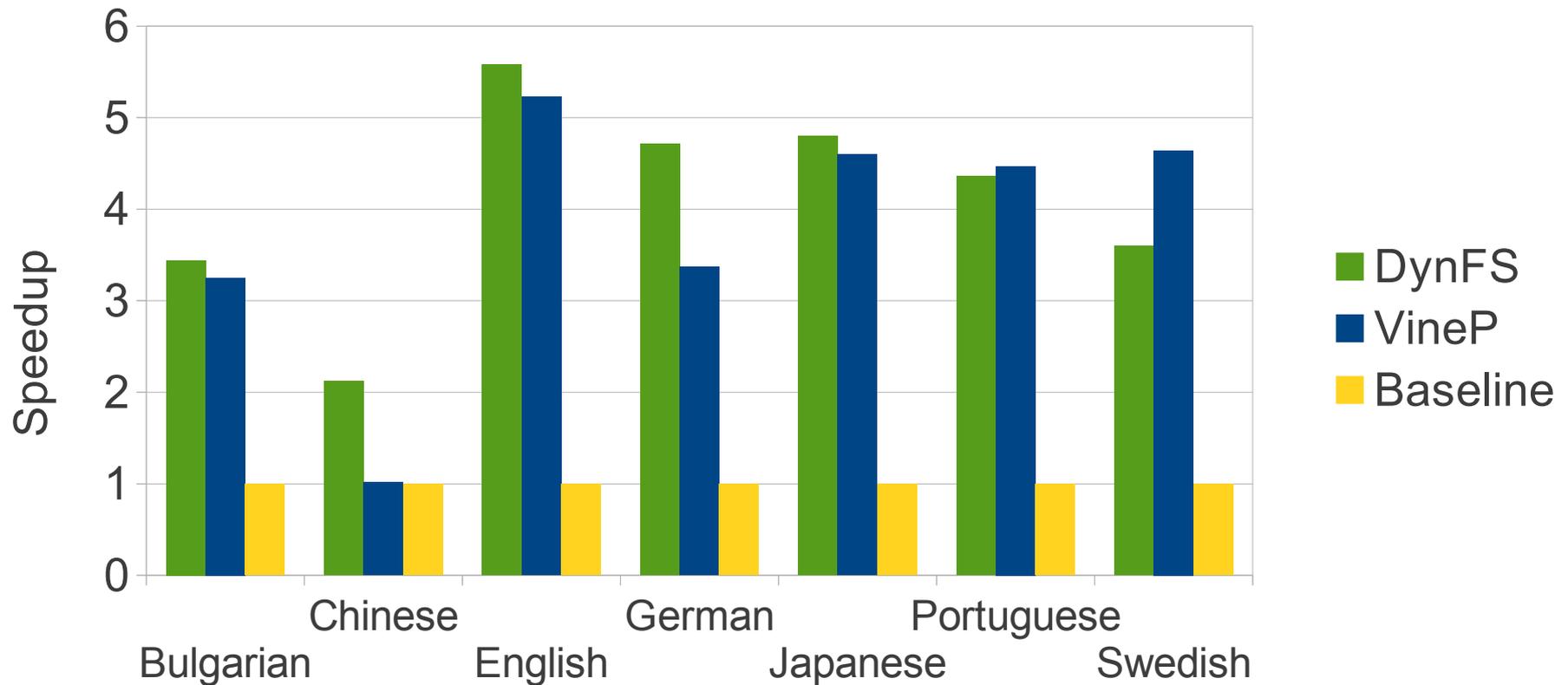
# Experiment: 2nd-order
# 2x to 8x speedup

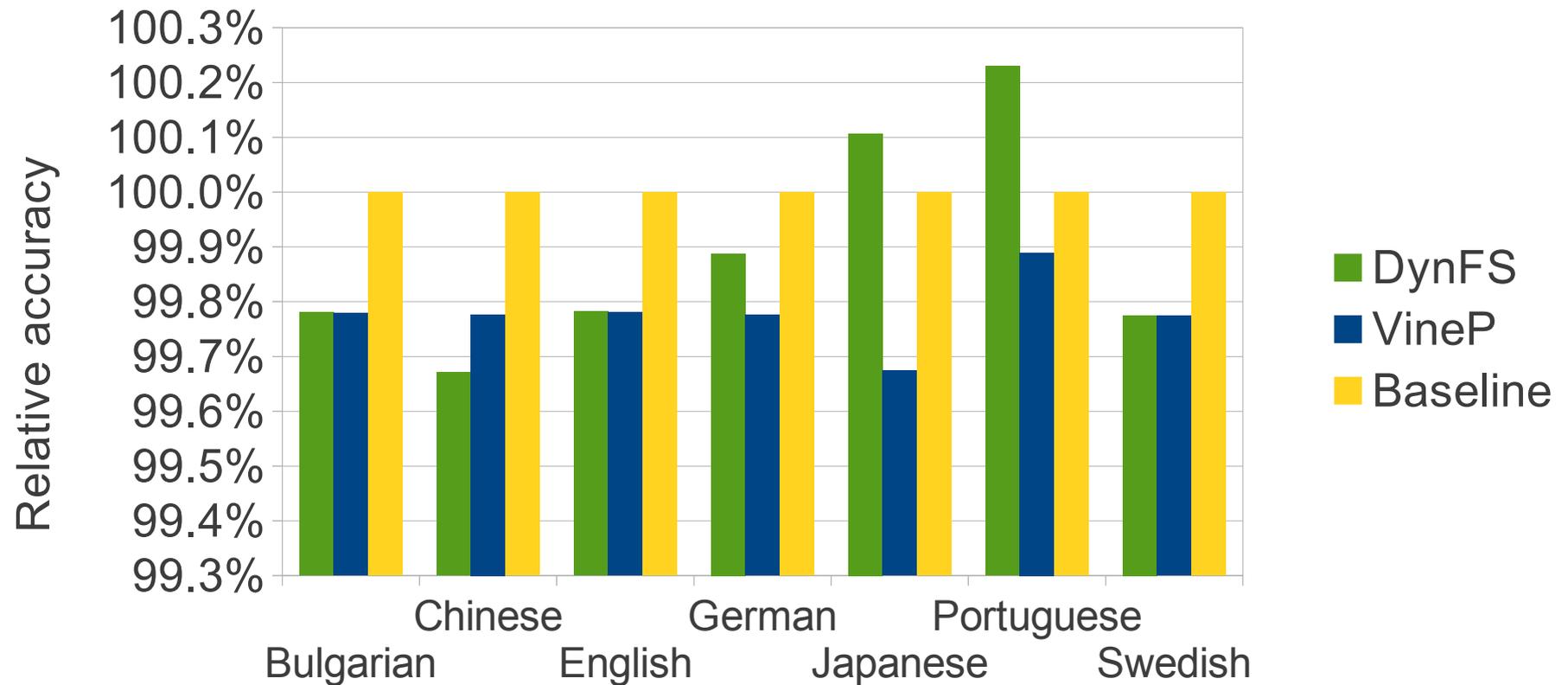# Experiment: 2nd-order
# ~0.3% loss in accuracy

# Ours vs Vine Pruning (Rush and Petrov, 2012)

- **Vine pruning**: a very fast parser that speeds up using orthogonal techniques
  - Start with short edges (*fully* scored)
  - Add long edges in if needed
- **Ours**

  - Start with all edges (*partially* scored)
  - Quickly remove unneeded edges
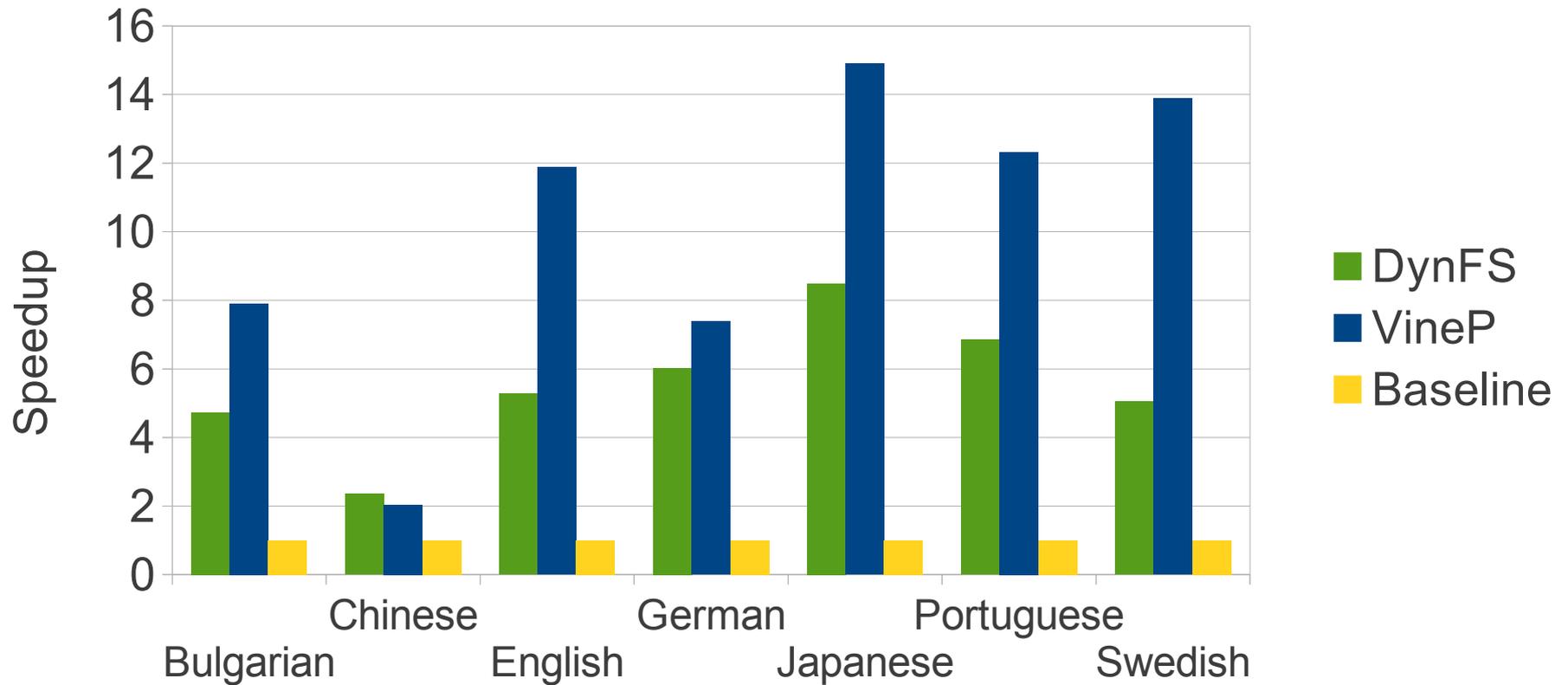- Could be combined for further speedup!

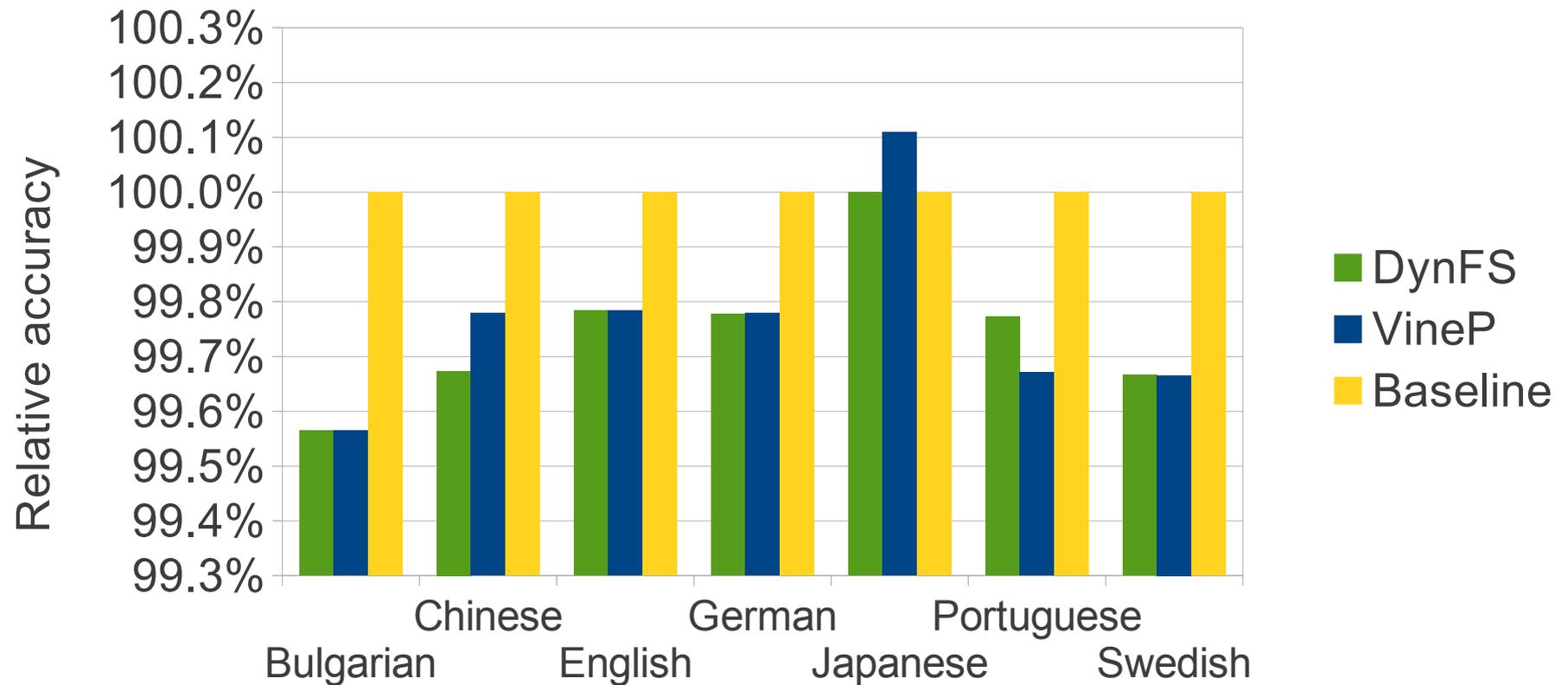# VS Vine Pruning: 1st-order comparable performance

# VS Vine Pruning: 1st-order

# VS Vine Pruning: 2nd-order

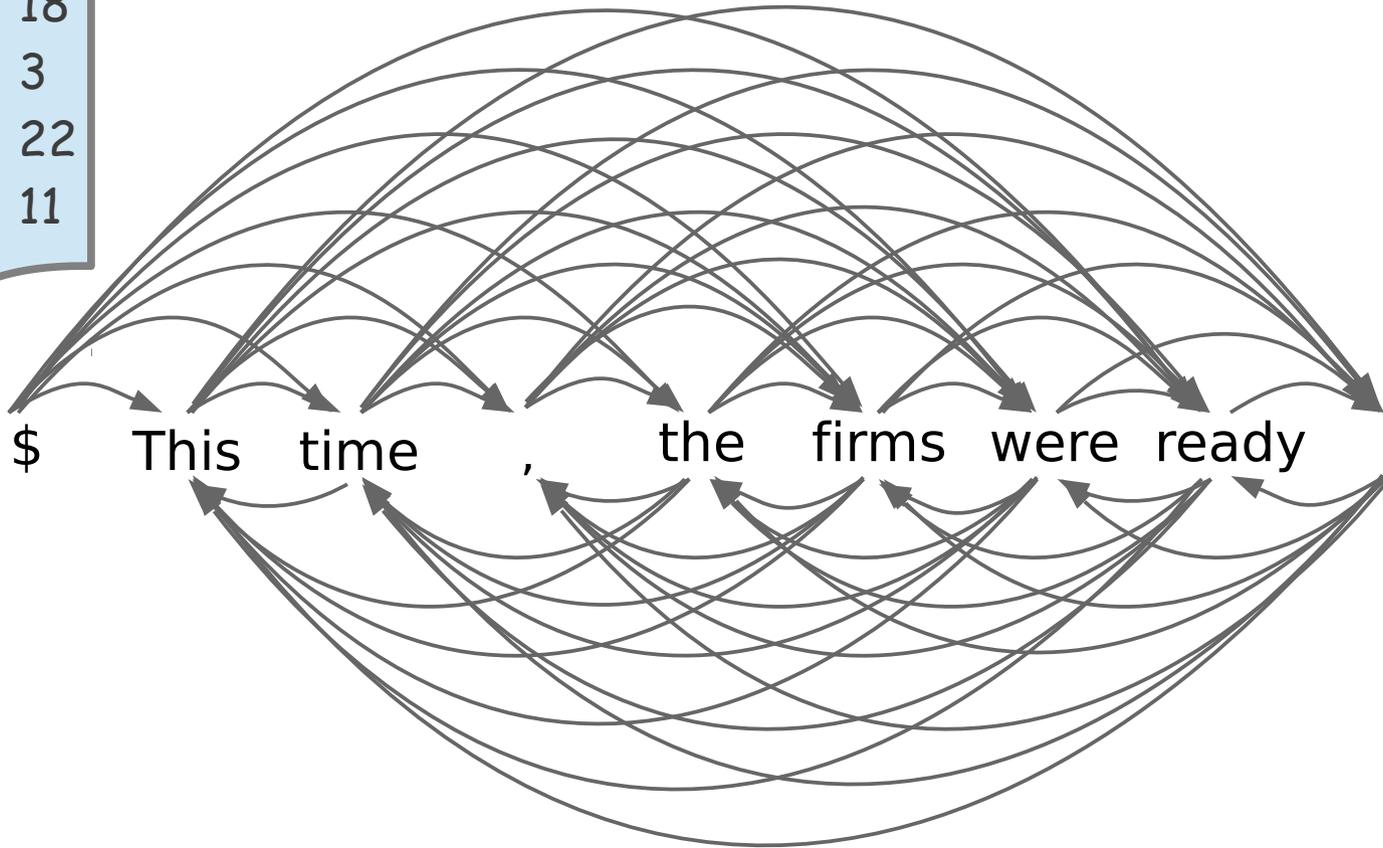# VS Vine Pruning: 2nd-order

# Conclusion

- **Feature computation** is expensive in structured prediction

- Commitment should be made **dynamically**

- Early commitment to edges reduce both searching and scoring time

- Can be used in other feature-rich models for structured prediction

# Backup Slides

# Static dictionary pruning (Rush and Petrov, 2012)

VB → CD:      18
VB ← CD:       3
NN → VBG:   22
NN ← VBG:   11
...

$    This    time        ,        the    firms    were    ready    .

# Reinforcement Learning 101

- Markov Decision Process (MDP)

  - State: all the information helping us to make decisions

  - Action: things we choose to do

  - Reward: criteria for evaluating actions

  - Policy: the "brain" that makes the decision

- Goal

  - Maximize the expected future reward

# Policy Learning

- Markov Decision Process (MDP)

$$\pi\ (\ \text{the}\quad \text{firms} + \text{context}) = \text{add} / \text{lock}$$

  - reward = accuracy + $\lambda \cdot$ speed

- Reinforcement learning

  - Delayed reward

  - Long time to converge

- Imitation learning

  - Mimic the oracle

  - Reduced to supervised classification problem

# Imitation Learning

- Oracle

  - (near) optimal performance

  - generate target action in any given state

$\pi$ ( the  firms + context) = *lock*
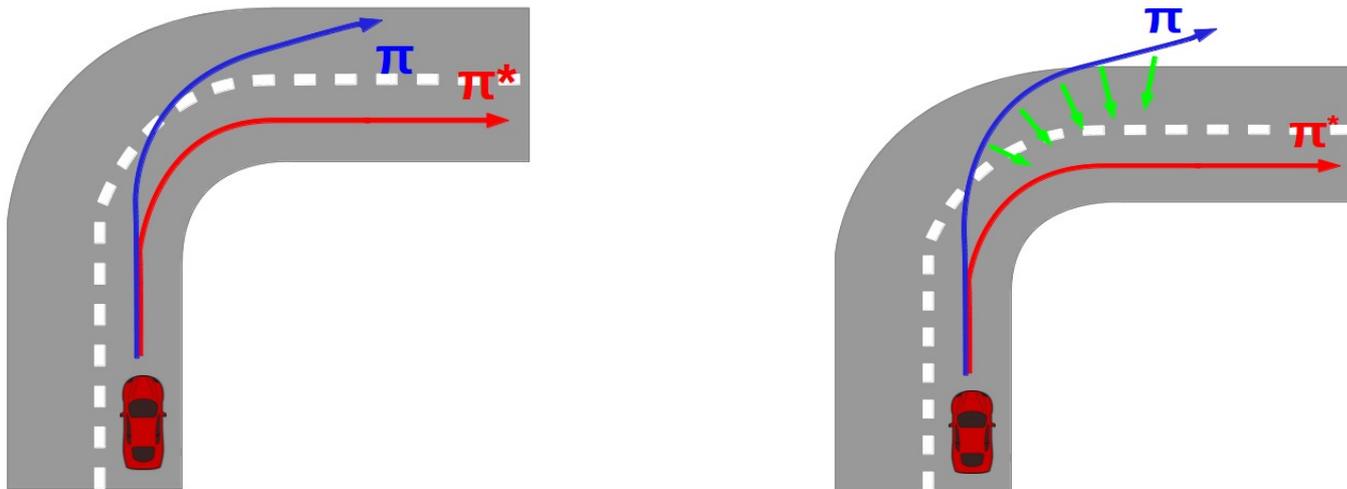
$\pi$ (time  ,  the + context) = *add*

...

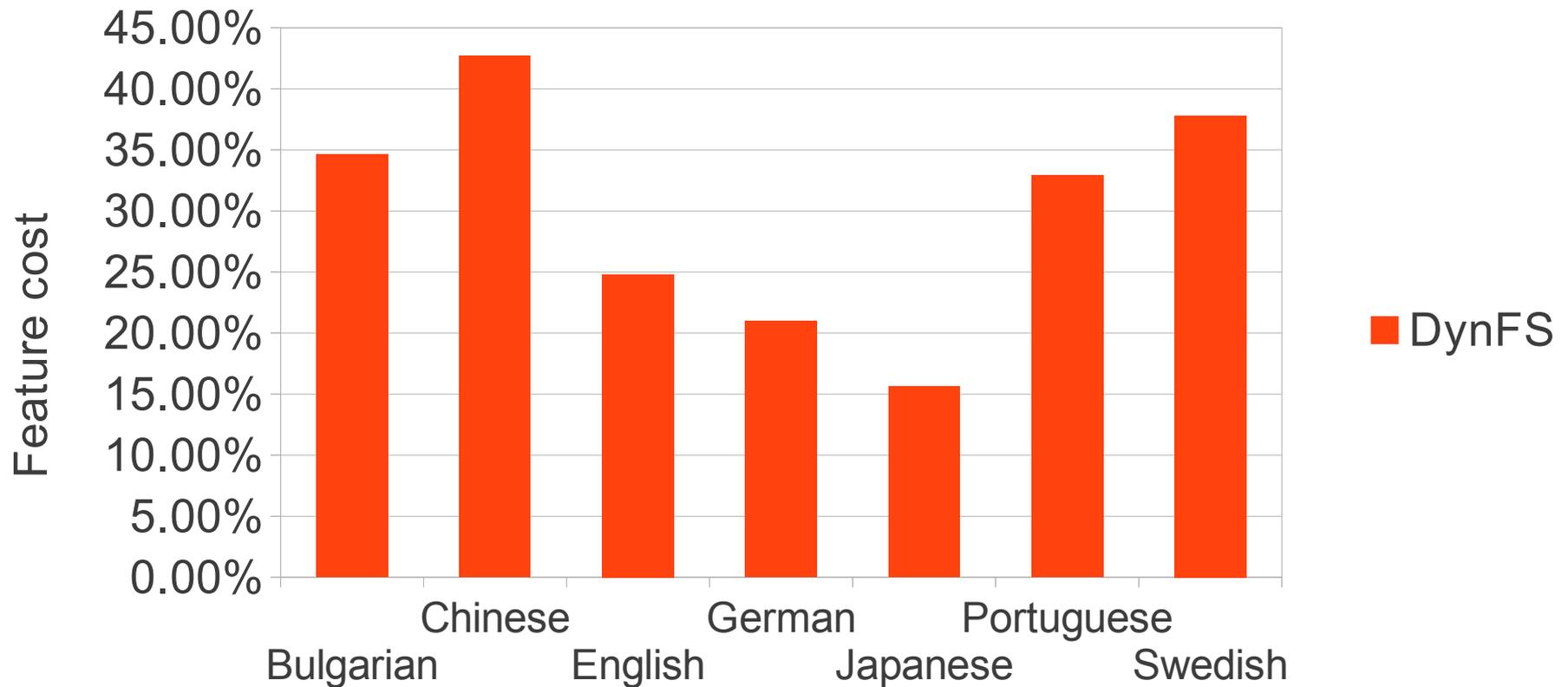$$\{\psi(s), \pi^*(s)\}$$

Binary classifier

# Dataset Aggregation (DAgger)

- Collect data from the oracle only
  - Different distribution at training and test time
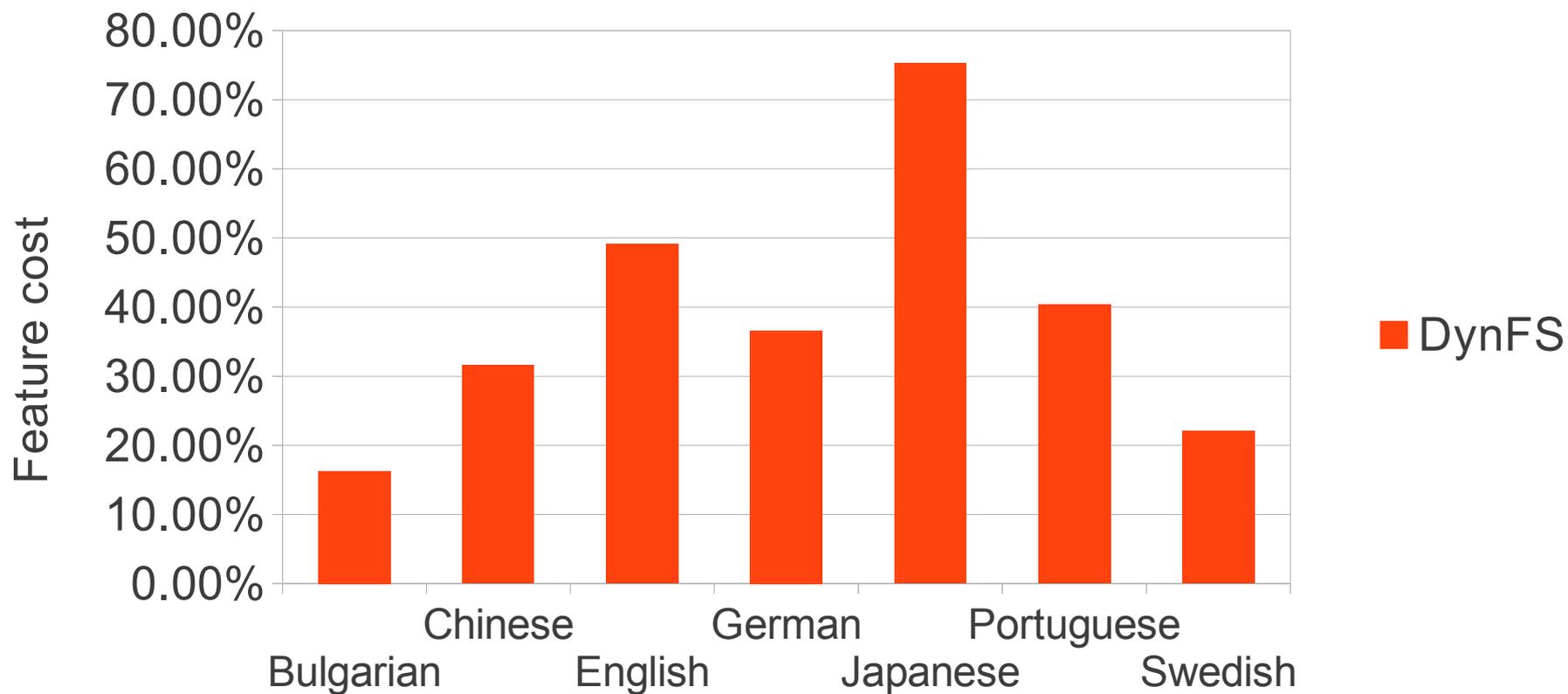- Iterative policy training



- Correct the learner's mistake
- Obtain a policy performs well under its own policy distribution

# Experiment (1st-order)



$$cost = \frac{\text{\# feature templates used}}{\text{total \# feature templates on the statically pruned graph}}$$

# Experiment (2nd-order)

# Second-order Parsing

.

# Second-order Parsing