# When is Self-training Effective for Parsing?

David McClosky
dmcc@cs.brown.edu

Brown Laboratory for Linguistic Information Processing (BLLIP)

Joint work with Eugene Charniak and Mark Johnson

# Outline

- What is self-training?

- Previous work

- Experimental setup
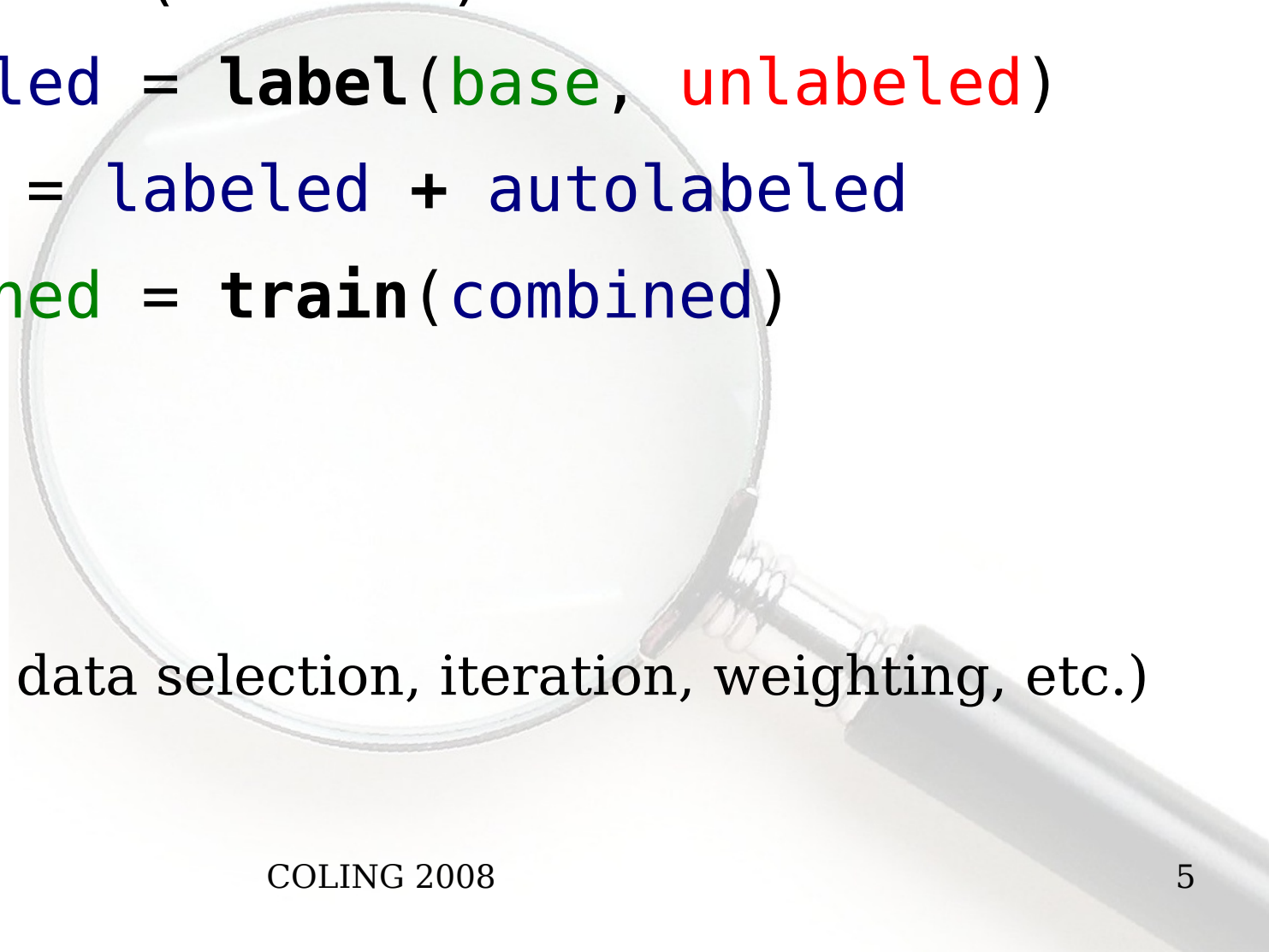
- Four hypotheses

- Conclusions

# Outline

- **What is self-training?**

- Previous work

- Experimental setup

- Four hypotheses

- Conclusions

# **Self-training Requirements**

- Labeled data

- Large amount of unlabeled data

- Statistical model:

    - `model` = **train**(`labeled data`)

    - `labels` = **label**(`model`, `unlabeled data`)

# Self-training Pseudocode

```
1  base = train(labeled)

2  autolabeled = label(base, unlabeled)

3  combined = labeled + autolabeled

4  selftrained = train(combined)
```

(Not pictured: data selection, iteration, weighting, etc.)

# Outline

- What is self-training?

- **Previous work**

- Experimental setup

- Four hypotheses

- Conclusions

# Previous Work

|  | Parser type | Seed size | Iterations | Improved? |
|---|---|---|---|---|
| **Charniak (1997)** | Generative | Large | Single | No |
| **McClosky et al. (2006)** | Gen.+Disc. | Large | Single | **Yes** |
| **Steedman et al. (2003)** | Generative | Small | Multiple | No |
| **Reichart and Rappoport (2007)** | Generative | Small | Single | **Yes** |

(large = ~40k sentences, small = <1k sentences)

**Summary of self-training for parsing experiments**

# Previous Work

| | Parser type | Seed size | Iterations | Improved? |
|---|---|---|---|---|
| **Charniak (1997)** | Generative | Large | Single | No |
| **McClosky et al. (2006)** | Gen.+Disc. | Large | Single | **Yes** |
| **Steedman et al. (2003)** | Generative | Small | Multiple | No |
| **Reichart and Rappoport (2007)** | Generative | Small | Single | **Yes** |

(large = ~40k sentences, small = <1k sentences)

**Summary of self-training for parsing experiments**

- In large seed case, generative + discriminative parser is necessary.

# Previous Work

| | Parser type | Seed size | Iterations | Improved? |
|---|---|---|---|---|
| **Charniak (1997)** | Generative | Large | Single | No |
| **McClosky et al. (2006)** | Gen.+Disc. | Large | Single | **Yes** |
| **Steedman et al. (2003)** | Generative | Small | Multiple | No |
| **Reichart and Rappoport (2007)** | Generative | Small | Single | **Yes** |

(large = ~40k sentences, small = <1k sentences)

**Summary of self-training for parsing experiments**

- In large seed case, generative + discriminative parser is necessary.

- Performing only one iteration is better than multiple iterations.

# Previous Analysis

| | Seed size | Predictor | |
|---|---|---|---|
| | | Length | # unknown words |
| **McClosky et al. (2006)** | Large | + | - |
| **Reichart and Rappoport (2007)** | Small | + | + |

# Previous Analysis

| | Seed size | Predictor | |
|---|---|---|---|
| | | Length | # unknown words |
| **McClosky et al. (2006)** | **Large** | + | - |
| **Reichart and Rappoport (2007)** | **Small** | + | + |

- Unknown words are a good predictor of self-training's success **only** in the small seed case.

# Outline

- What is self-training?

- Previous work

- **Experimental setup**

- Four hypotheses

- Conclusions

# Experimental Setup

- Labeled data: WSJ (Marcus et al., 1993)

- Unlabeled data: NANC (Graff, 1995)
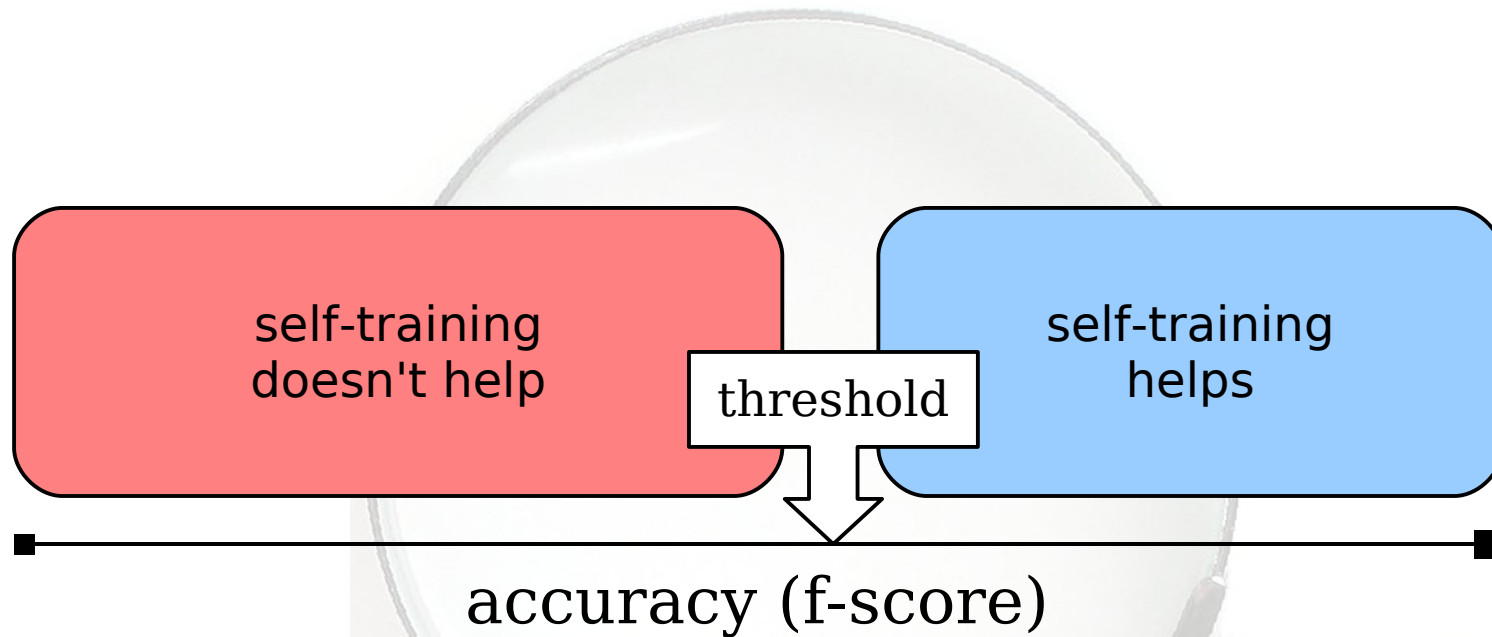
- Parser: Charniak and Johnson (2005) reranking parser

sentence → generative parser → parse

# Experimental Setup

- Labeled data: WSJ (Marcus et al., 1993)

- Unlabeled data: NANC (Graff, 1995)

- Parser: Charniak and Johnson (2005) reranking parser

sentence → $n$-best generative parser → parses

# Experimental Setup

- Labeled data: WSJ (Marcus et al., 1993)

- Unlabeled data: NANC (Graff, 1995)

- Parser: Charniak and Johnson (2005) reranking parser

sentence → *n*-best generative parser → parses → discriminative reranker → parse

# Outline

- What is self-training?

- Previous work

- Experimental setup

- **Four hypotheses**

- Conclusions

# Hypotheses for Self-training

1 Phase Transition

2 Search Errors

3 Non-generative Reranker Features

4 Bilexical Dependencies

# Phase Transition

Self-training works after a phase transition

self-training doesn't help

threshold

self-training helps

accuracy (f-score)

# Phase Transition



Parser
89.9%

Reranking Parser
91.5%

self-training
doesn't help

self-training
helps

accuracy (f-score)

sections 1, 22, 24

# Phase Transition

Parser
85.8%
10% WSJ

Reranking Parser
87.0%
10% WSJ

Parser
89.9%
100% WSJ

Reranking Parser
91.5%
100% WSJ

self-training
helps

self-training
doesn't help

self-training
helps

accuracy (f-score)

sections 1, 22, 24

➜ There is no phase transition for self-training.

# Phase Transition



Parser
85.8%
10% WSJ

Reranking Parser
87.0%
10% WSJ

Parser
89.9%
100% WSJ

Reranking Parser
91.5%
100% WSJ

self-training helps

self-training doesn't help

self-training helps

accuracy (f-score)

sections 1, 22, 24

See also: Reichart and Rappoport (2007)

➜ There is no phase transition for self-training.

# Search Errors

Self-trained models have fewer search errors

(Daniel Marcu, p.c.)

# Search Errors

🔍 Self-trained models have fewer search errors

(Daniel Marcu, p.c.)

model prefers
worse parse

model | search

errors

# Search Errors

🔍 Self-trained models have fewer search errors

(Daniel Marcu, p.c.)

model prefers
worse parse

model | search

best parse for
model not found

errors

# Some notation

**orig**  ← parses from original parser

**st**  ← parses from self-trained parser

$\text{top}_m(\mathbf{P})$  ← best parse in **P** for reranker model m

(m in **{orig,st}**)

# Comparing n-best lists

| | |
|---|---|
| Overlap of **st** and **orig** | 66.0% |
| $\text{top}_{st}(\textbf{st}) = \text{top}_{orig}(\textbf{orig})$ | 42.4% |
| $\text{top}_{st}(\textbf{st})$ in **orig** | 60.3% |
| Search errors | 2.5% |

(statistics on 5,039 sentences in sections 1, 22, 24)

- Search errors =

$\text{top}_{st}(\textbf{st})$ not in **orig** and

$\text{top}_{orig}(\textbf{st} \cup \textbf{orig}) = \text{top}_{st}(\textbf{st})$

# Decreasing Search Errors

parses — self-trained parser

parses — original parser

- Add parses from self-trained n-best list to original parser's n-best list, rescoring by original parser

parses — augmented original parser

parses — zero probability under original parser

# Evaluation

parses — **original reranking parser**

**91.5%**

parses — **self-trained reranking parser**

**92.0%**

parses — **augmented original reranking parser**

**91.7%**

(reranking parser f-score on sections 1, 22, 24)

# Evaluation

parses

original
reranking parser

**91.5%**

parses

self-trained
reranking parser

**92.0%**

➔ The original parser makes both **model** and **search** errors relative to the self-trained model.

parses

augmented original
reranking parser

**91.7%**

(reranking parser f-score on sections 1, 22, 24)

# Reranker features

Non-generative reranker features help self-training more
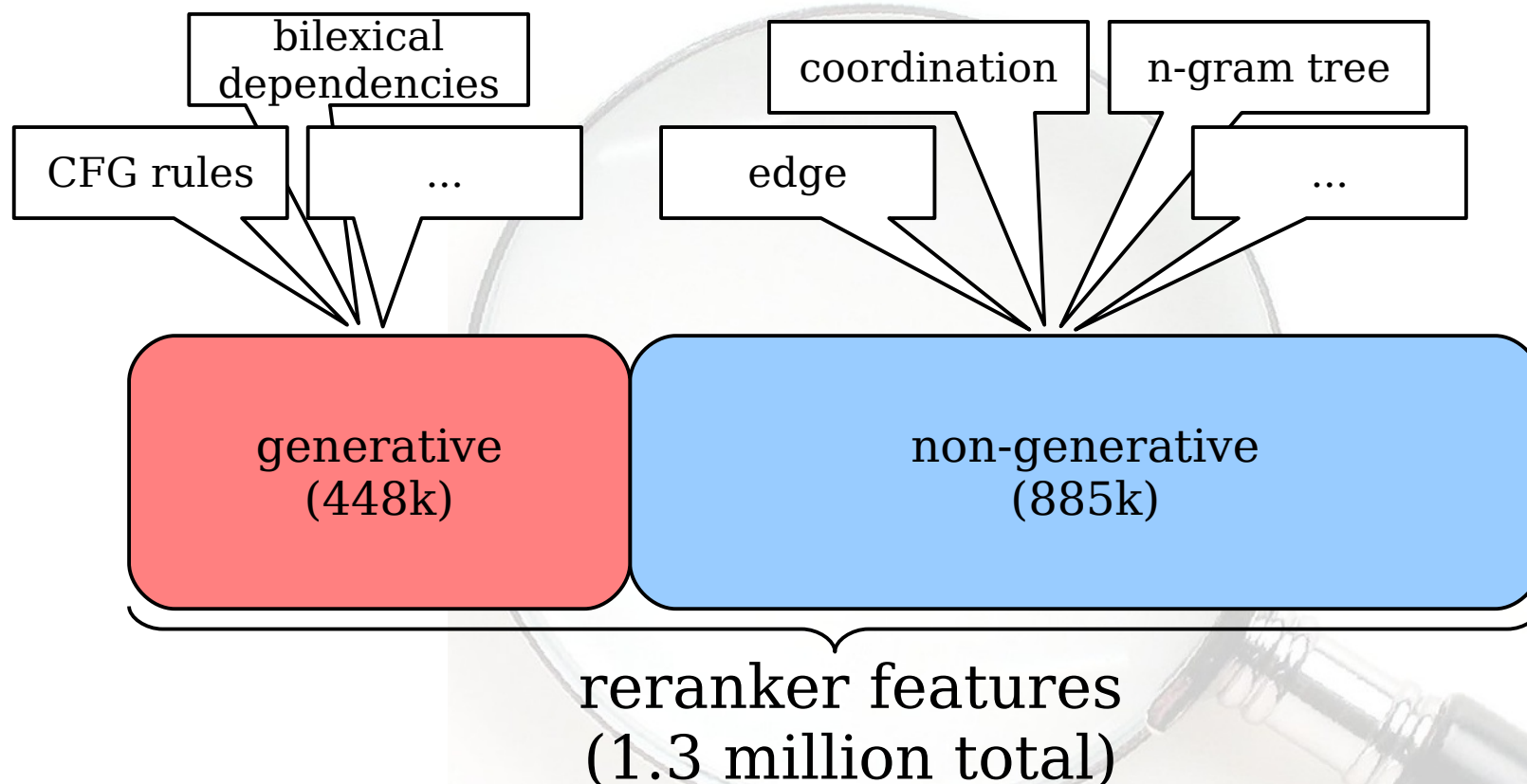
reranker features
(1.3 million total)

# Reranker features

| generative (448k) | non-generative (885k) |
|---|---|

reranker features
(1.3 million total)

# Reranker features

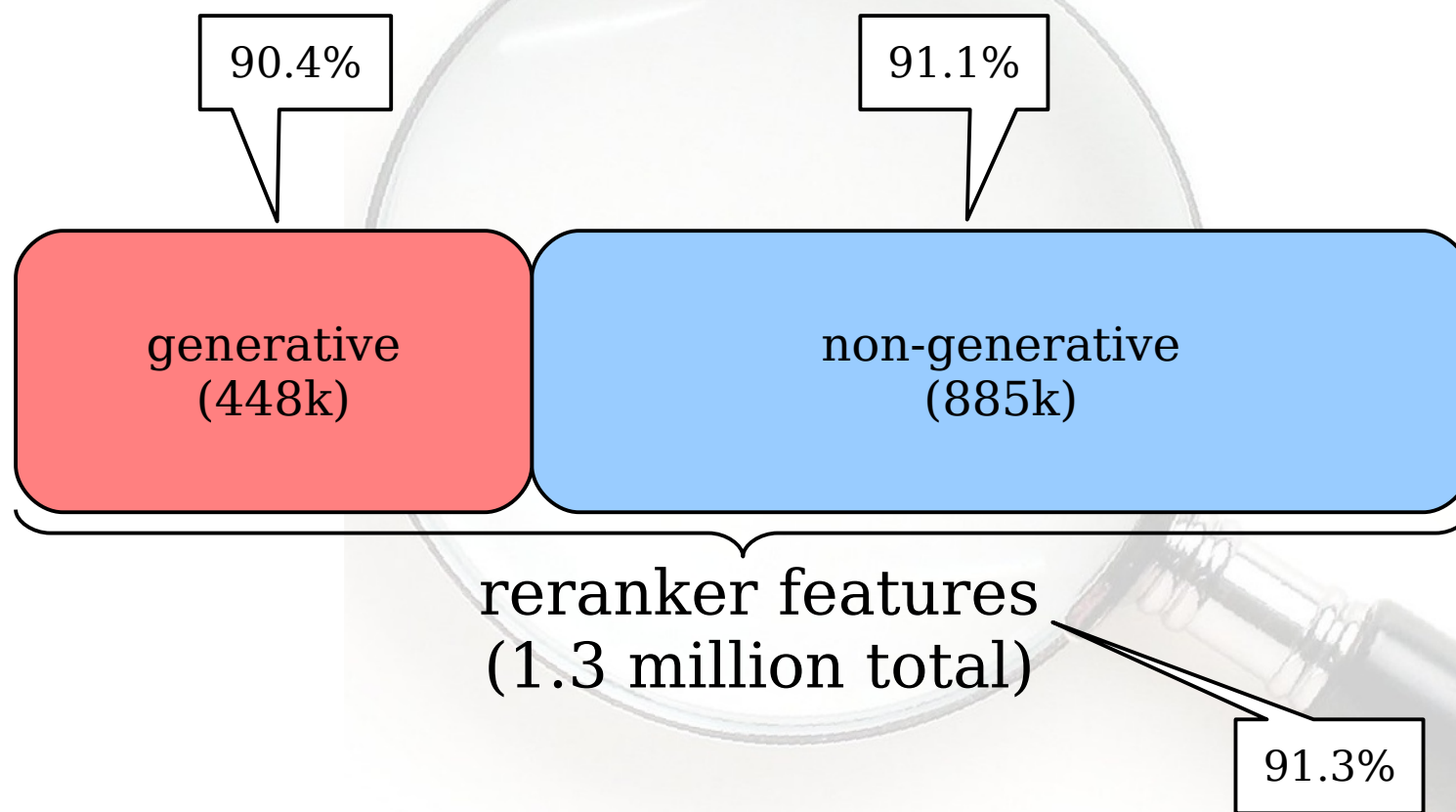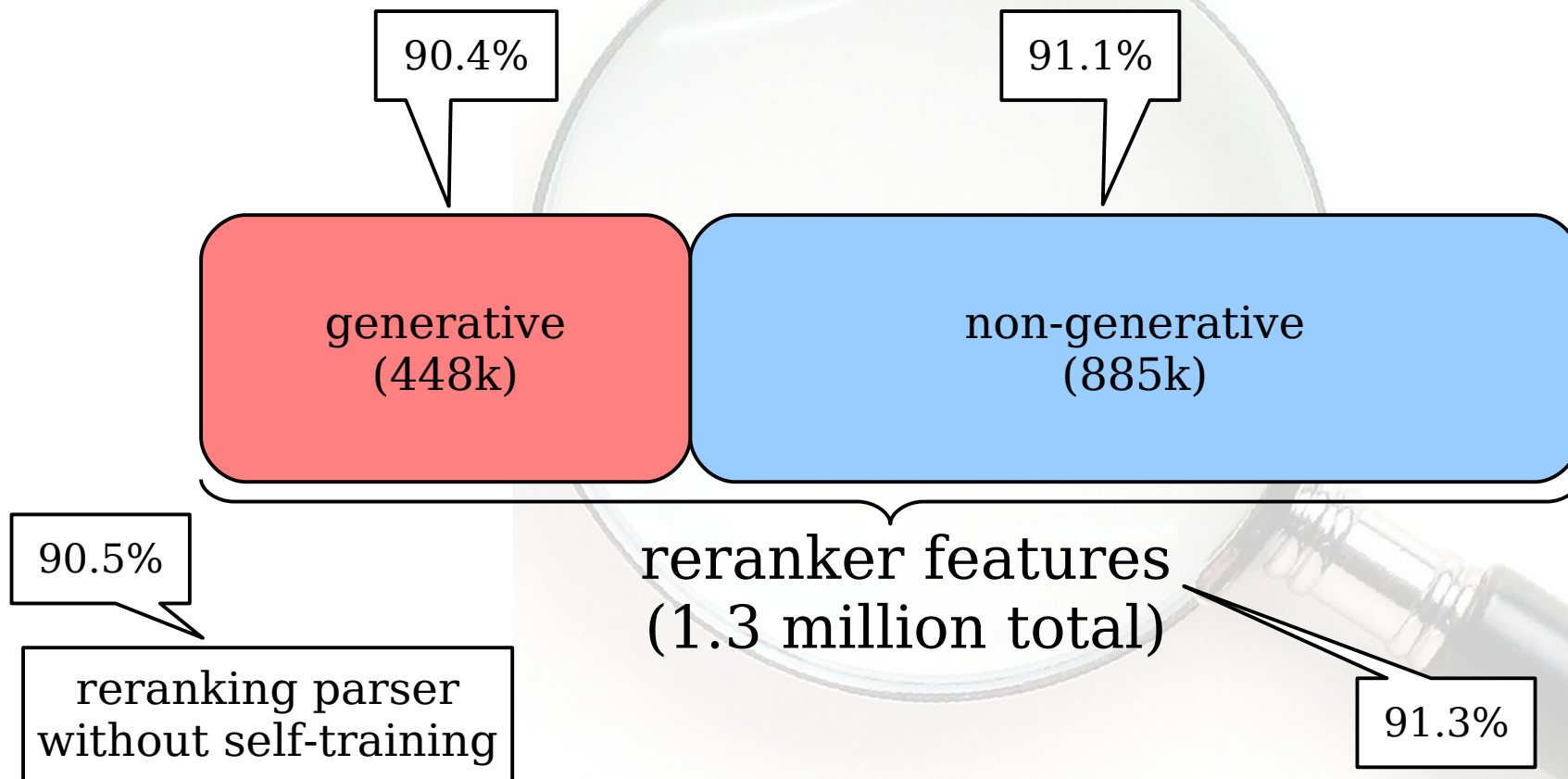Non-generative reranker features help self-training more

bilexical dependencies

CFG rules

...

coordination

n-gram tree

edge

...

**generative (448k)**

**non-generative (885k)**

reranker features
(1.3 million total)

# Reranker features

🔍 Non-generative reranker features help self-training more

(reranking parser f-score on section 24)

90.4%

91.1%

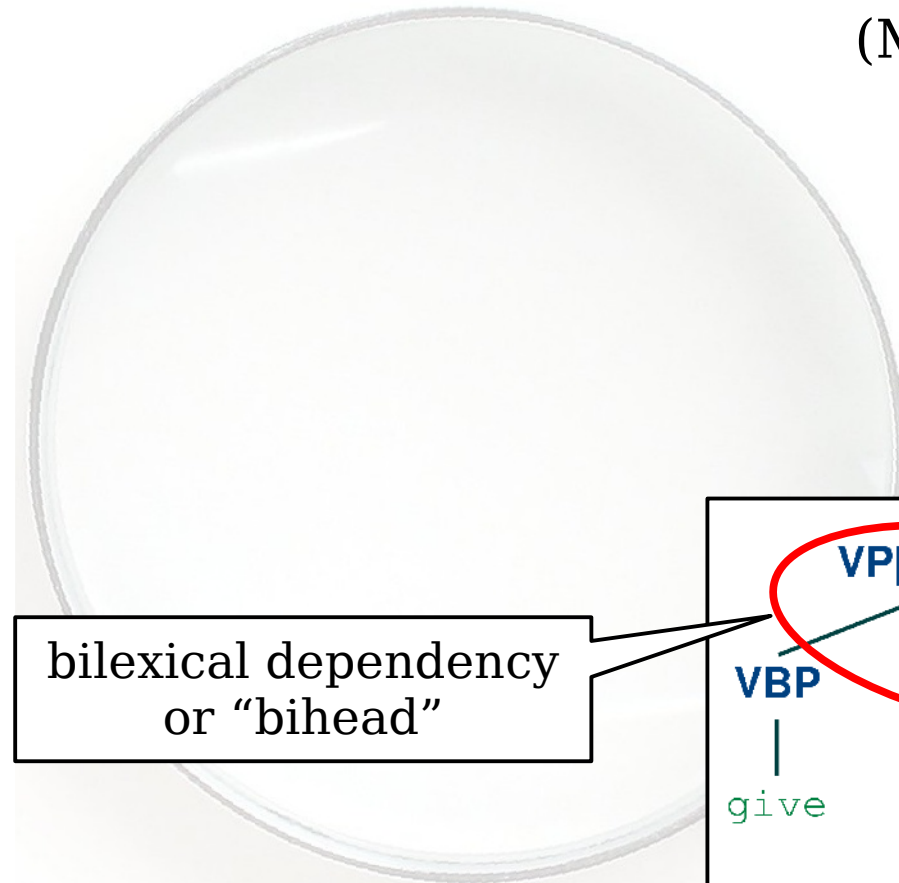generative
(448k)

non-generative
(885k)

reranker features
(1.3 million total)

91.3%

# Reranker features

Non-generative reranker features help self-training more

(reranking parser f-score on section 24)

90.4%

91.1%

generative
(448k)

non-generative
(885k)

90.5%

reranking parser
without self-training

reranker features
(1.3 million total)

91.3%

# Reranker features

➜ Non-generative reranker features are essential for self-training with a reranker.
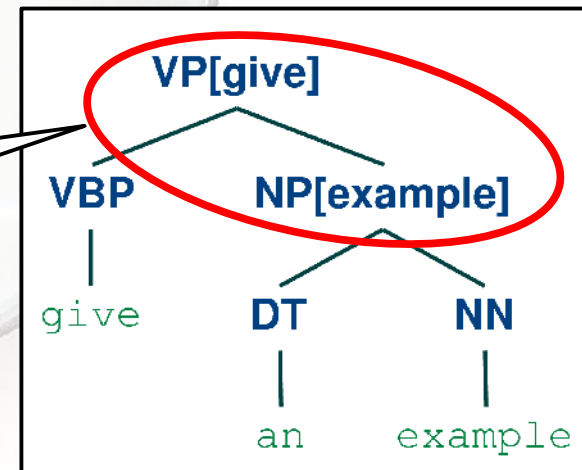
(reranking parser f-score on section 24)

90.4%

91.1%

| generative (448k) | non-generative (885k) |

reranker features
(1.3 million total)

90.5%

reranking parser
without self-training

91.3%

# Bilexical Dependencies

Self-training teaches the parser about bilexical dependencies
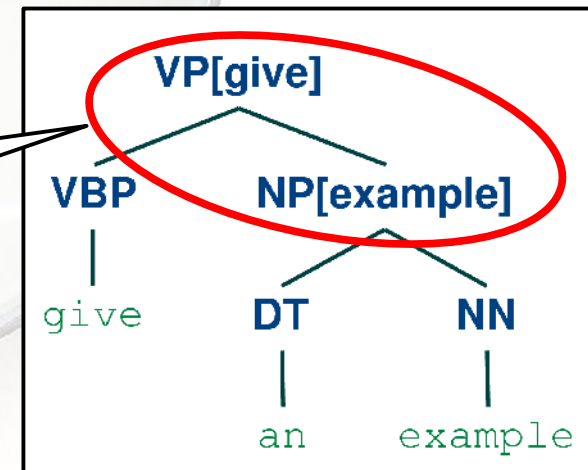
(Mitch Marcus, p.c.)

bilexical dependency or "bihead"

```
         VP[give]
        /        \
      VBP      NP[example]
       |        /     \
     give     DT      NN
              |        |
             an     example
```

# Bilexical Dependencies

Self-training teaches the parser about bilexical dependencies

(Mitch Marcus, p.c.)

Two ways to test:

- Factor analysis (as in previous work)

- Transfer biheads distribution from self-trained model to original model
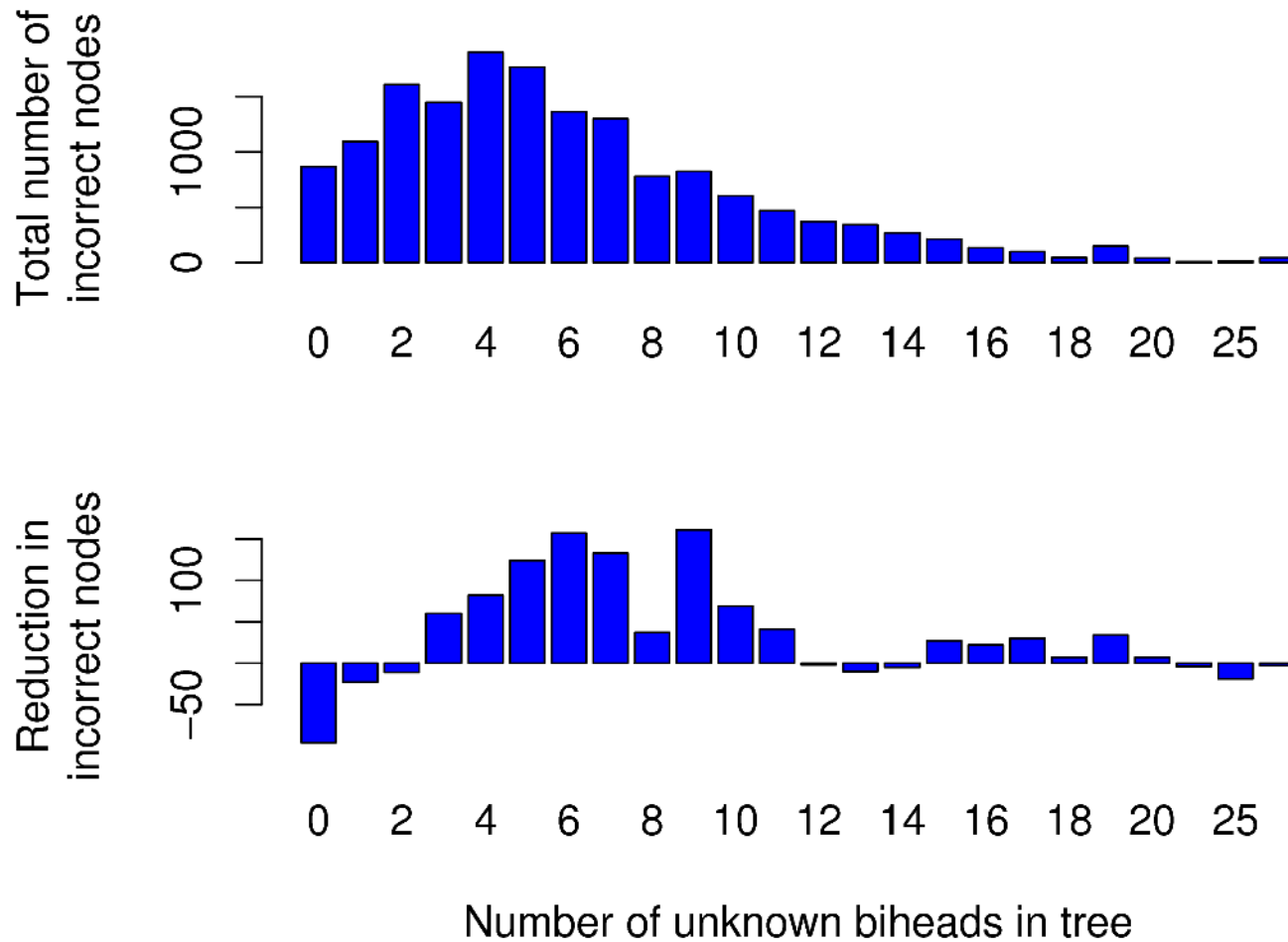
bilexical dependency or "bihead"

VP[give]

VBP | give
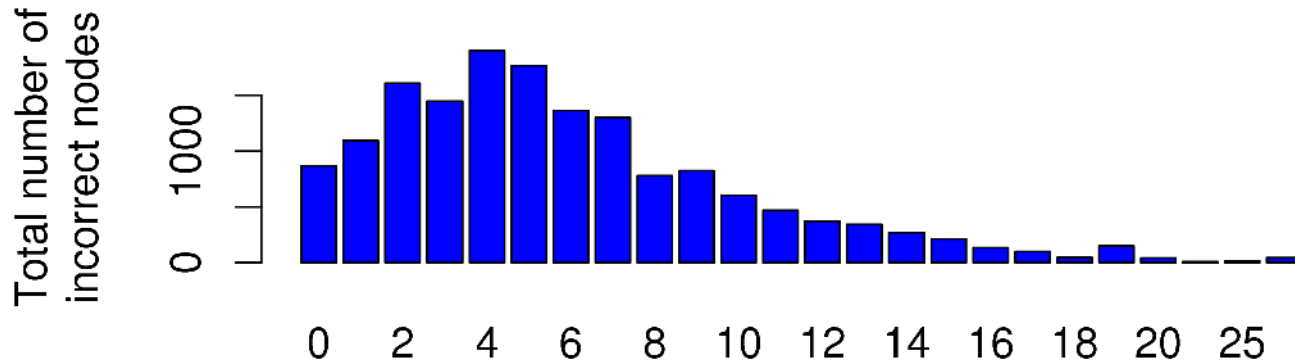
NP[example]

DT | an

NN | example

# Factor analysis: Words



Number of unknown words in tree

# Factor analysis: Words



biggest improvement when **no** unknown words!

# Factor Analysis: Biheads

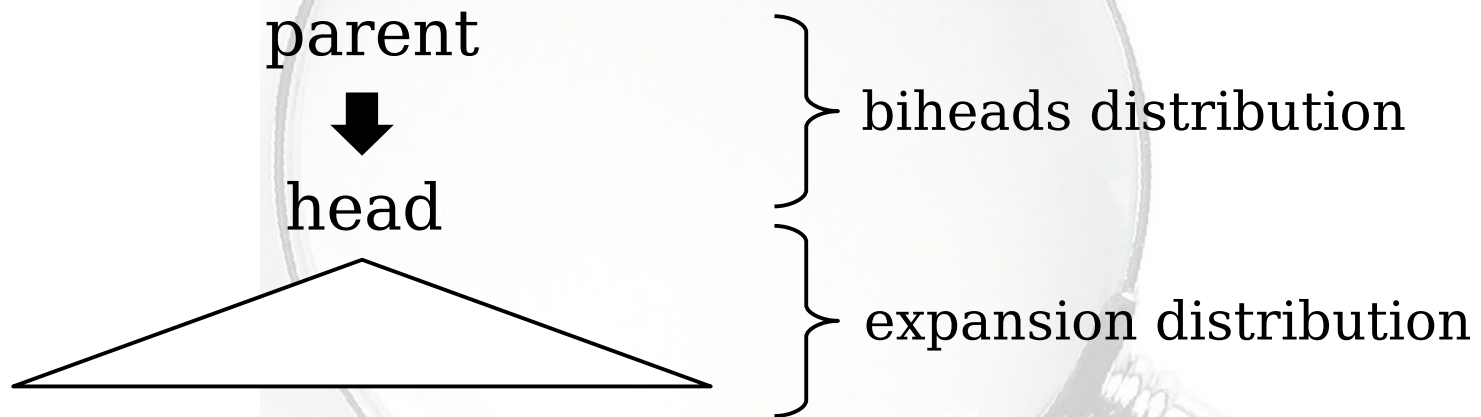# Factor Analysis: Biheads



seeing more unknown biheads helps

Number of unknown biheads in tree

# Just biheads?

- If self-trained model learns more about biheads, can we transfer that knowledge to original model?

parent

head

biheads distribution

expansion distribution

**dangerously oversimplified Charniak parsing model!**

# Just biheads?

- If self-trained model learns more about biheads, can we transfer that knowledge to original model?
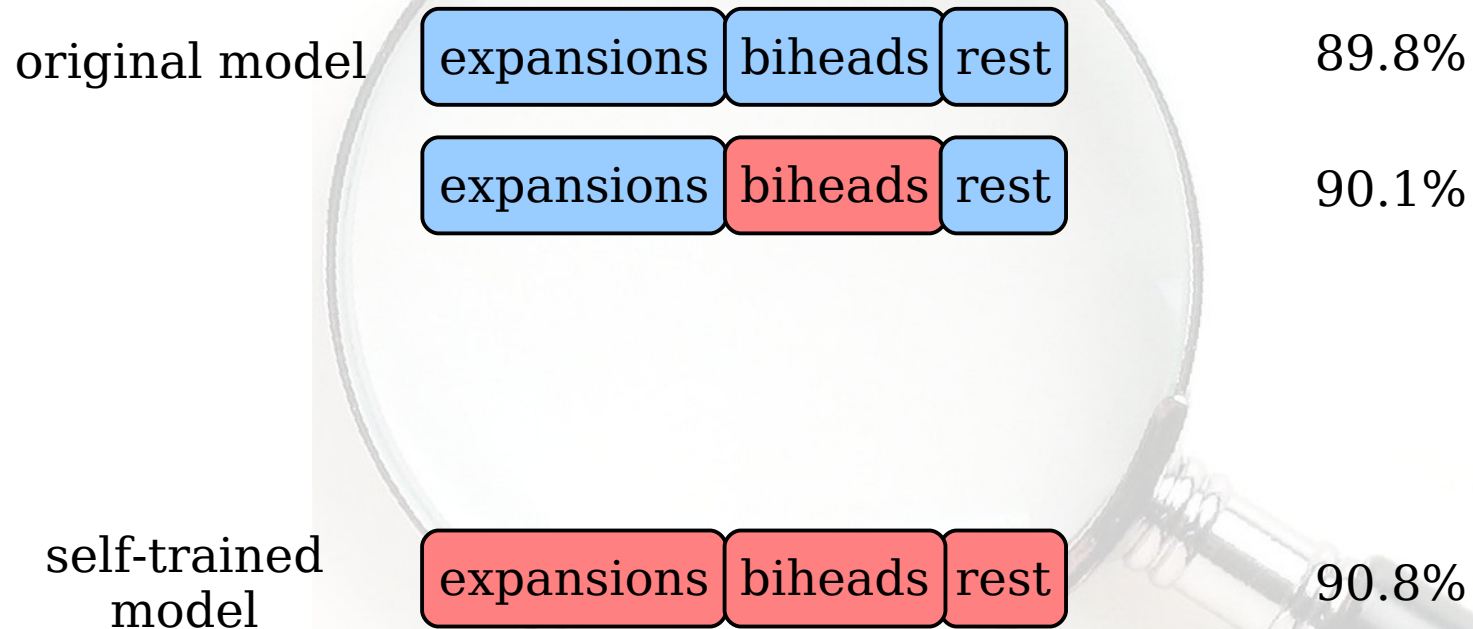
original model    | expansions | biheads | rest |    89.8%

self-trained model    | expansions | biheads | rest |    90.8%
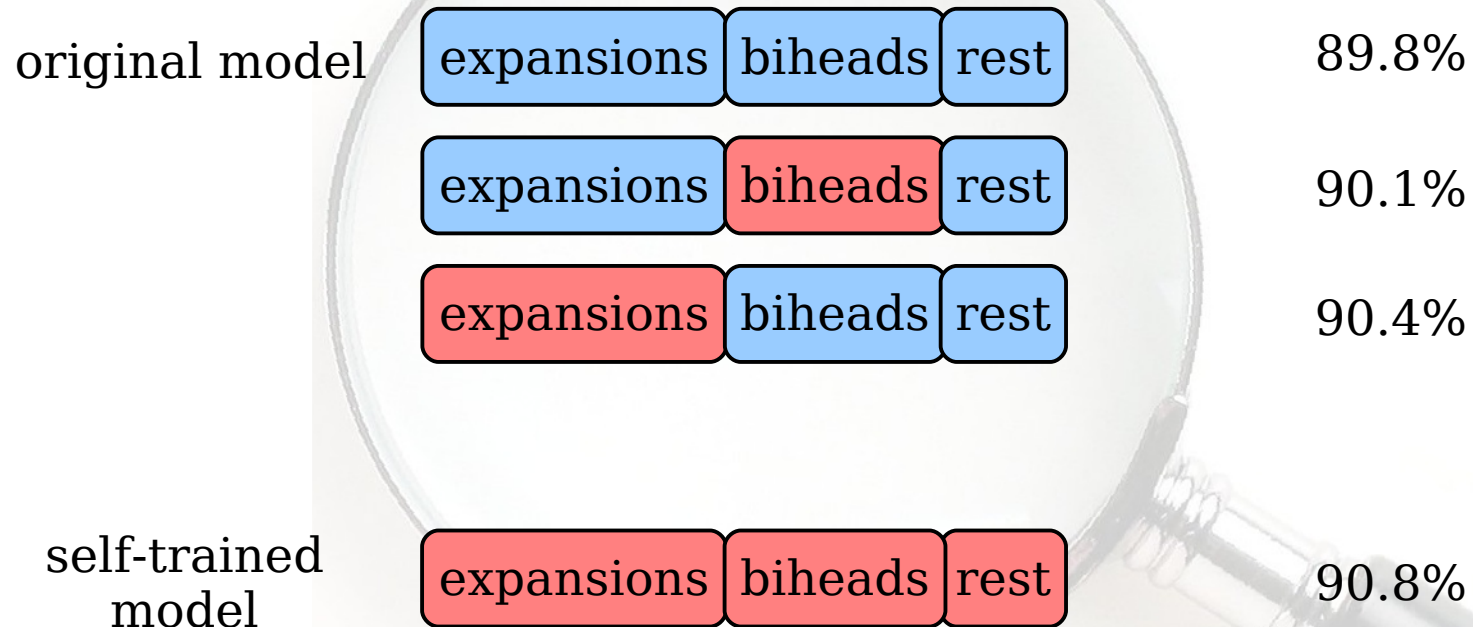
(generative parser f-score on sections 1, 22, 24)

# Just biheads?

- If self-trained model learns more about biheads, can we transfer that knowledge to original model?

original model [ expansions | biheads | rest ]   89.8%

[ expansions | biheads | rest ]   90.1%

self-trained model [ expansions | biheads | rest ]   90.8%

(generative parser f-score on sections 1, 22, 24)

# Just biheads?

- If self-trained model learns more about biheads, can we transfer that knowledge to original model?
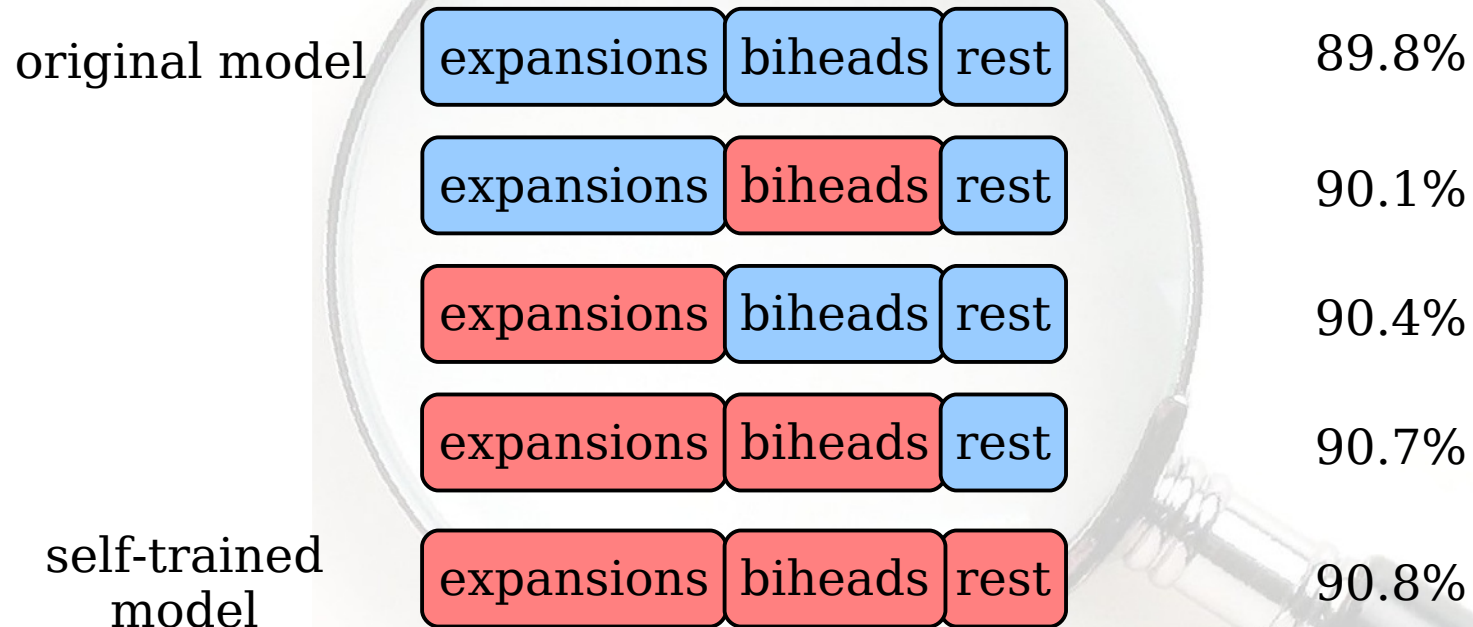
| | | |
|---|---|---|
| original model | expansions biheads rest | 89.8% |
| | expansions **biheads** rest | 90.1% |
| | **expansions** biheads rest | 90.4% |
| self-trained model | **expansions biheads rest** | 90.8% |

(generative parser f-score on sections 1, 22, 24)

# Just biheads?

- If self-trained model learns more about biheads, can we transfer that knowledge to original model?

| original model | expansions | biheads | rest | 89.8% |
| | expansions | biheads | rest | 90.1% |
| | expansions | biheads | rest | 90.4% |
| | expansions | biheads | rest | 90.7% |
| self-trained model | expansions | biheads | rest | 90.8% |

(generative parser f-score on sections 1, 22, 24)

# Just biheads?

➔ Self-training improves biheads distribution, but **also** expansions distribution.

| | | | | |
|---|---|---|---|---|
| original model | expansions | biheads | rest | 89.8% |
| | expansions | biheads | rest | 90.1% |
| | expansions | biheads | rest | 90.4% |
| | expansions | biheads | rest | 90.7% |
| self-trained model | expansions | biheads | rest | 90.8% |

(generative parser f-score on sections 1, 22, 24)

# Outline

- What is self-training?

- Previous work

- Experimental setup

- Four hypotheses

- **Conclusions**

COLING 2008

# Hypothesis Results

1 **Phase Transition**

Disproved.

2 **Search Errors**

Reducing search errors helps but model errors remain.

3 **Non-generative Reranker Features**

Reranker features must be different from generative parser.

4 **Bilexical Dependencies**

Biheads correlated with self-training improvements.
Self-training helps all parts of the parser, not just biheads.

# Two different cases?

| | Predictor | | | Need reranker? |
|---|---|---|---|---|
| | Length | unk. words | unk. biheads | |
| **Large seed** | + | - | + | + |
| **Small seed** | + | + | ? | - |

# A Possible Connection

- **Small seed case:** Vocabulary is sparse so unknown **words** may be resolved in unlabeled text.

- **Large seed case:** Learn new **head information** in unlabeled text.  Non-generative features in the reranker needed to handle more complex constructions.

# Future Work

- "Bilexical dependencies" experiments with small seed size.

- Different parsers (Collins, LTAG, discriminative, ...) with and without rerankers.

- More hypotheses...
  (Audience participation?)

# Acknowledgments

**Questions?**

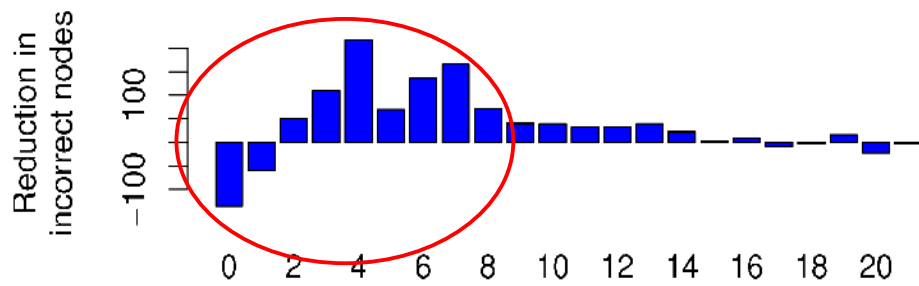`http://bllip.cs.brown.edu/selftraining/`

# Extra Slides

# Bigrams and Biheads

# Bigrams and Biheads



seeing more unknown bigrams/biheads helps