

Analyzing Errors of Unsupervised Learning

ACL 2008 Columbus, Ohio

June 18, 2008

Percy Liang

Dan Klein



Unsupervised grammar induction

Goal: induce hidden syntax

The man ate a tasty sandwich

Unsupervised grammar induction

Goal: induce hidden syntax

DT — NN — VBD — DT — JJ — NN
| | | | | |
The man ate a tasty sandwich

POS tagging

Unsupervised grammar induction

Goal: induce hidden syntax

DT — NN — VBD — DT — JJ — NN
| | | | | |
The man ate a tasty sandwich

DT NN VBD

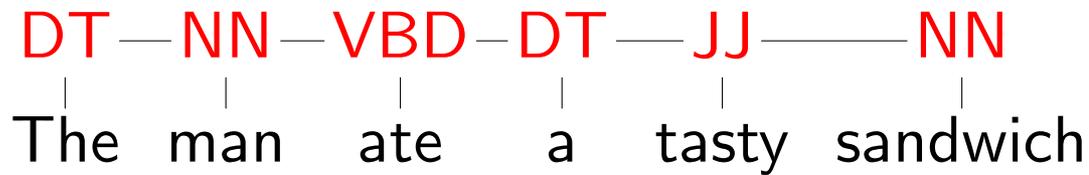
DT

JJ NN

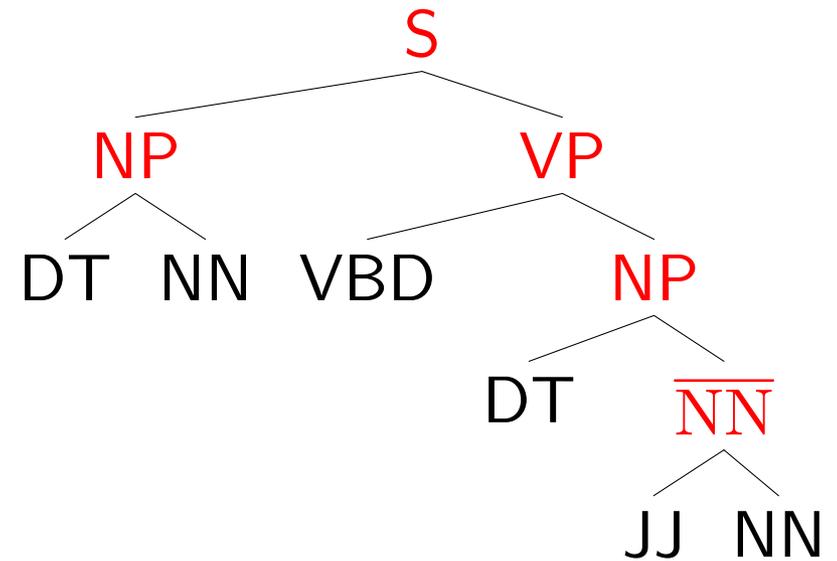
POS tagging

Unsupervised grammar induction

Goal: induce hidden syntax



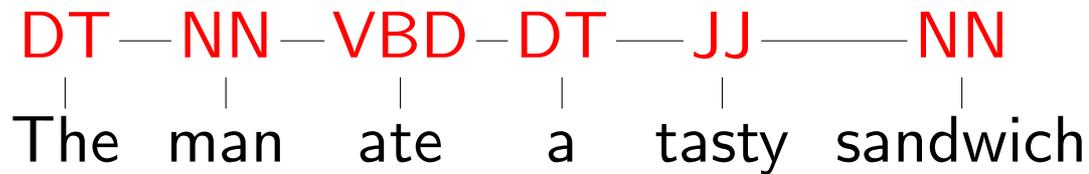
POS tagging



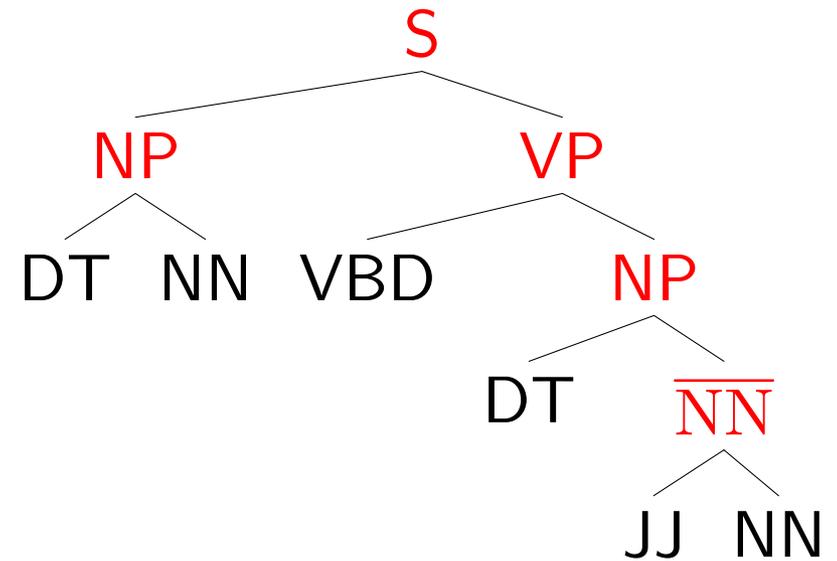
Constituency parsing

Unsupervised grammar induction

Goal: induce hidden syntax



POS tagging



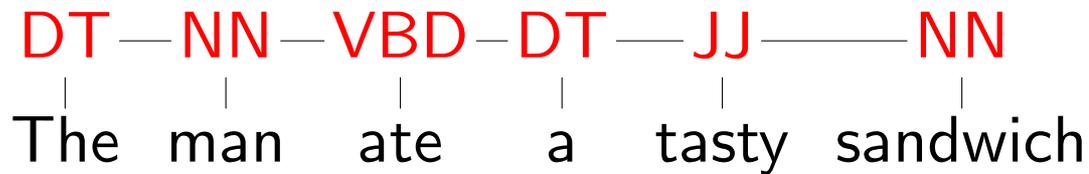
Constituency parsing

For example, on POS tagging using HMMs:

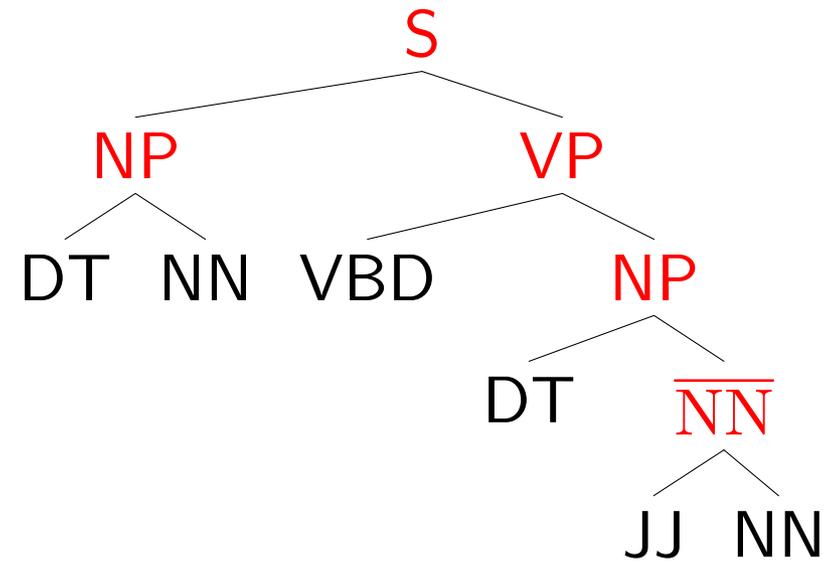
Unsupervised using EM \approx 60%

Unsupervised grammar induction

Goal: induce hidden syntax



POS tagging



Constituency parsing

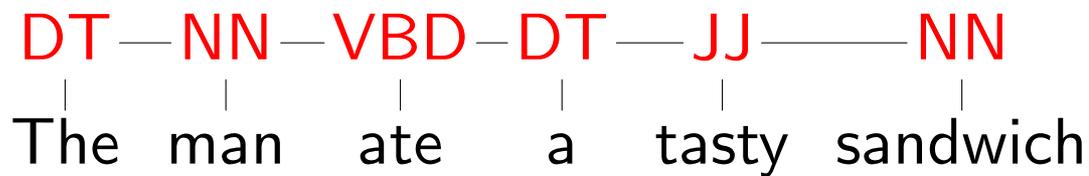
For example, on POS tagging using HMMs:

Unsupervised using EM $\approx 60\%$

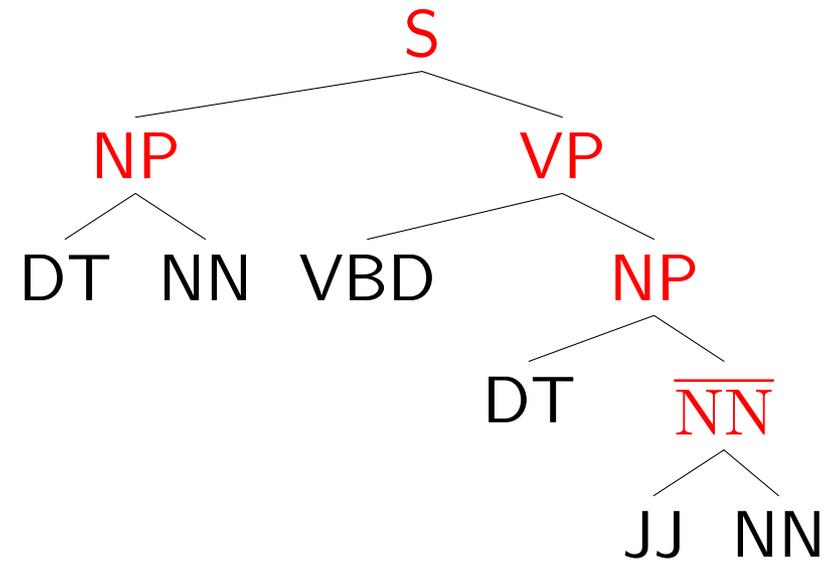
Supervised $\geq 90\%$

Unsupervised grammar induction

Goal: induce hidden syntax



POS tagging



Constituency parsing

For example, on POS tagging using HMMs:

Unsupervised using EM $\approx 60\%$

Supervised $\geq 90\%$

Why does EM fail?

Four types of errors:

Four types of errors:

Optimization error

Local optima

Four types of errors:

Optimization error

Local optima

Estimation error

Limited data

Four types of errors:

Optimization error

Local optima

Estimation error

Limited data

Approximation error

Likelihood objective $\not\Rightarrow$ accuracy

Four types of errors:

Optimization error

Local optima

Estimation error

Limited data

Approximation error

Likelihood objective $\not\Rightarrow$ accuracy

Identifiability error

Different parameter settings \rightarrow same objective

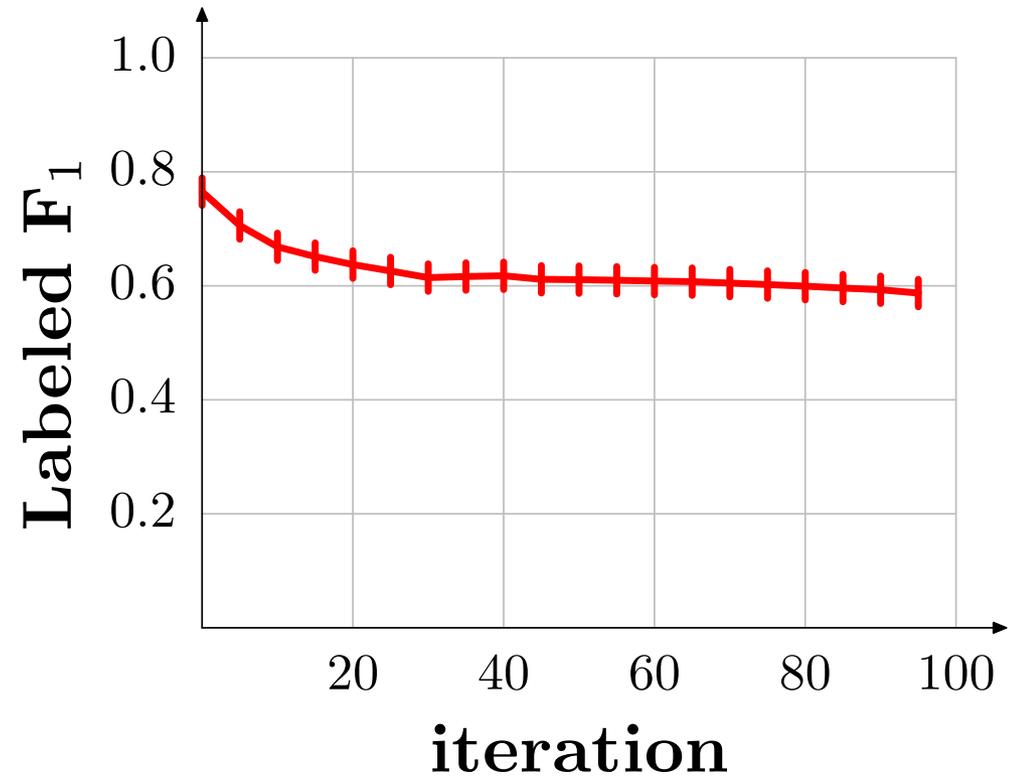
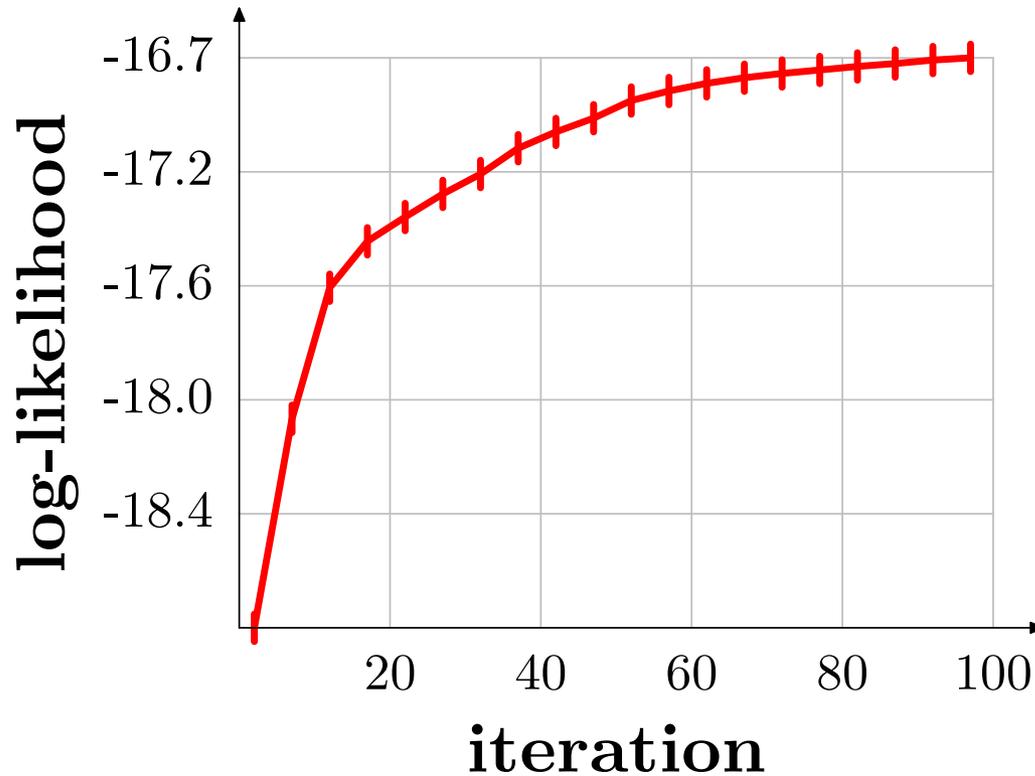
Approximation error

Problem: model likelihood $\not\Rightarrow$ prediction accuracy

Approximation error

Problem: model likelihood $\not\Rightarrow$ prediction accuracy

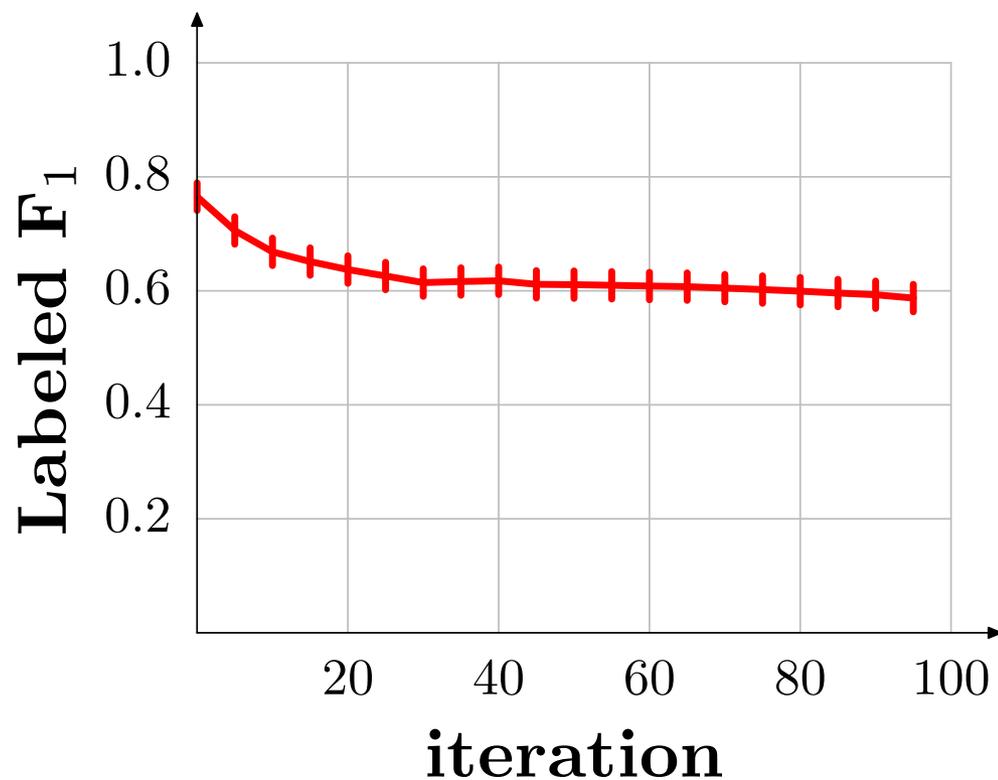
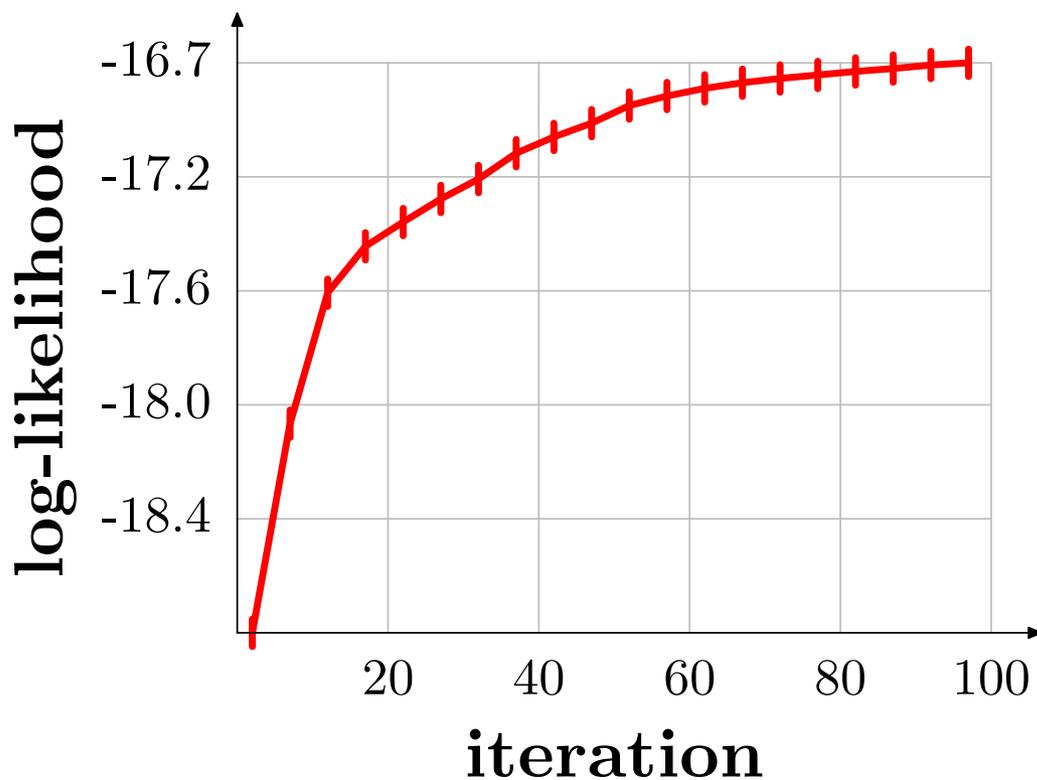
PCFG (EM starting from supervised parameter estimate):



Approximation error

Problem: model likelihood $\not\Rightarrow$ prediction accuracy

PCFG (EM starting from supervised parameter estimate):



What qualitative changes is EM making?

Migrations

For the HMM:

Truth	DT	NN	NN	RB	VBD	NNS
	The chief executive allegedly made contributions					
				↓		
Iteration 1	DT	JJ	NN	RB	VBN	NNS
	The chief executive allegedly made contributions					

Migrations

For the HMM:

Truth	DT	NN	NN	RB	VBD	NNS
	The chief executive allegedly made contributions					
				↓		
Iteration 1	DT	JJ	NN	RB	VBN	NNS
	The chief executive allegedly made contributions					

Summarize changes by a set of **migrations**:



Migrations

For the HMM:

Truth	DT	NN	NN	RB	VBD	NNS
	The chief executive allegedly made contributions					
				↓		
Iteration 1	DT	JJ	NN	RB	VBN	NNS
	The chief executive allegedly made contributions					

Summarize changes by a set of **migrations**:

NN → NN



JJ → NN

VBD → made



VBN → made

Migrations

For the HMM:

Truth	DT	NN	NN	RB	VBD	NNS
	The chief executive allegedly made contributions					
				↓		
Iteration 1	DT	JJ	NN	RB	VBN	NNS
	The chief executive allegedly made contributions					

Summarize changes by a set of **migrations**:

NN → NN



JJ → NN

VBD → made



VBN → made

What are the prominent migrations over the entire corpora?

Top HMM migrations

Iteration 1

START → ~~NN~~
NNP

Sentence-initial nouns are often proper
START Revenue/~~NN~~/~~NNP~~ rose

Top HMM migrations

Iteration 1

START → ~~NN~~
NNP

Sentence-initial nouns are often proper
*START Revenue/~~NN~~/**NNP** rose*

~~NN~~
JJ → NN

Noun adjuncts → adjectives (inconsistent gold tags)
*chief/~~NN~~/**JJ** executive/**NN** officer*

Top HMM migrations

Iteration 1

START → ~~NN~~
NNP

Sentence-initial nouns are often proper
*START Revenue/~~NN~~/**NNP** rose*

~~NN~~
JJ → NN

Noun adjuncts → adjectives (inconsistent gold tags)
*chief/~~NN~~/**JJ** executive/**NN** officer*

NNP → ~~NNP~~
NNPS

Inconsistent gold tags
*UBS Securities/~~NNP~~/**NNPS***

Top HMM migrations

Iteration 1

START → ~~NN~~
NNP

Sentence-initial nouns are often proper
START Revenue/~~NN~~/~~NNP~~ rose

~~NN~~
JJ → NN

Noun adjuncts → adjectives (inconsistent gold tags)
chief/~~NN~~/~~JJ~~ executive/~~NN~~ officer

NNP → ~~NNP~~
NNPS

Inconsistent gold tags
UBS Securities/~~NNP~~/~~NNPS~~

Iteration 2

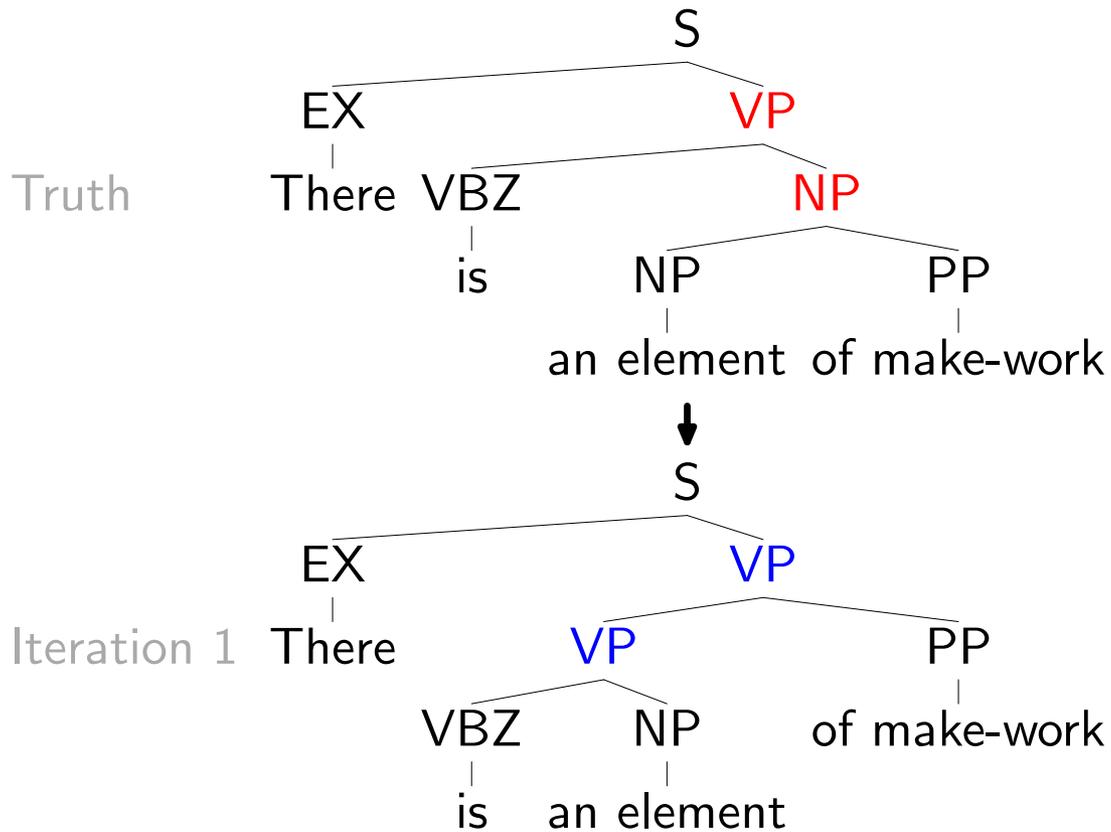
~~NN~~
JJ → NN (same as above)

START → ~~NN~~
NNP (same as above)

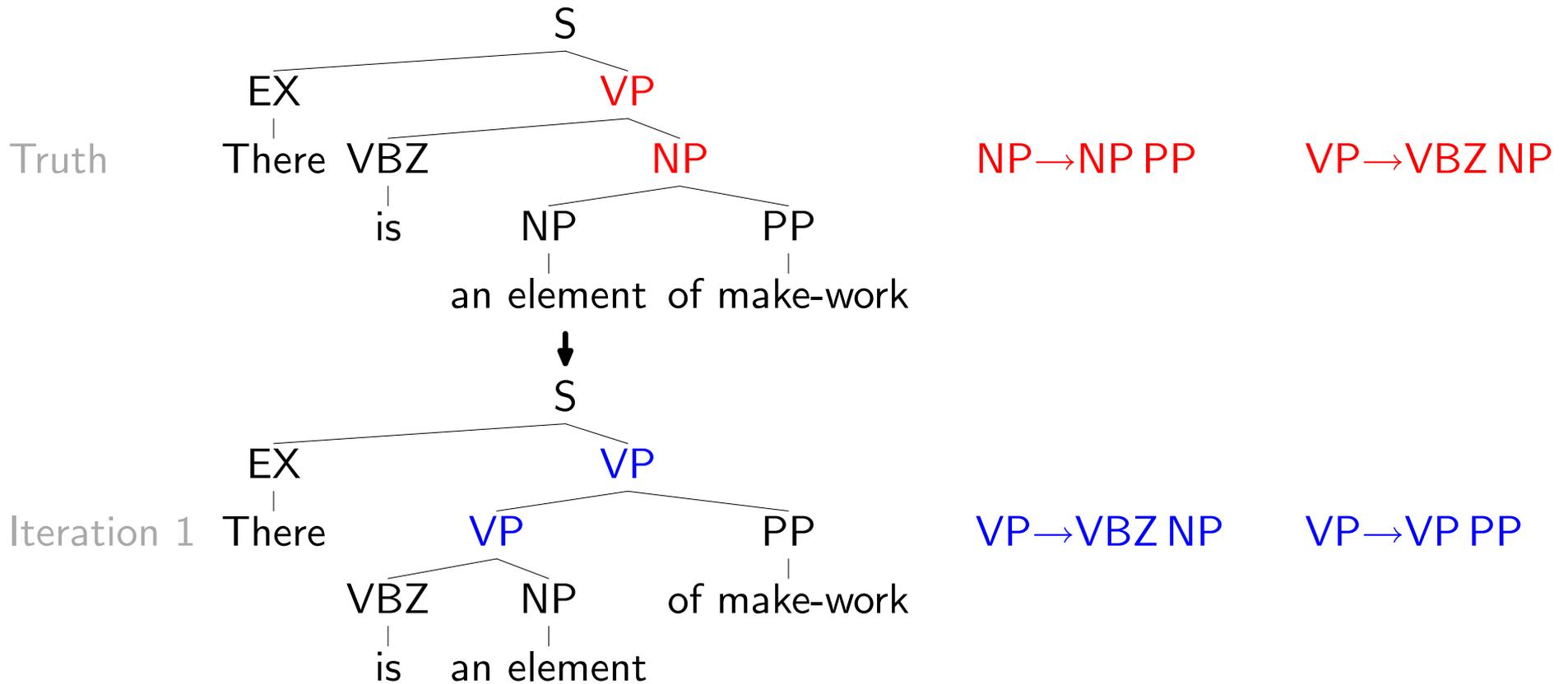
~~JJ~~
RB → TO

Inconsistent gold tags
contribute much/~~JJ~~/~~RB~~ to

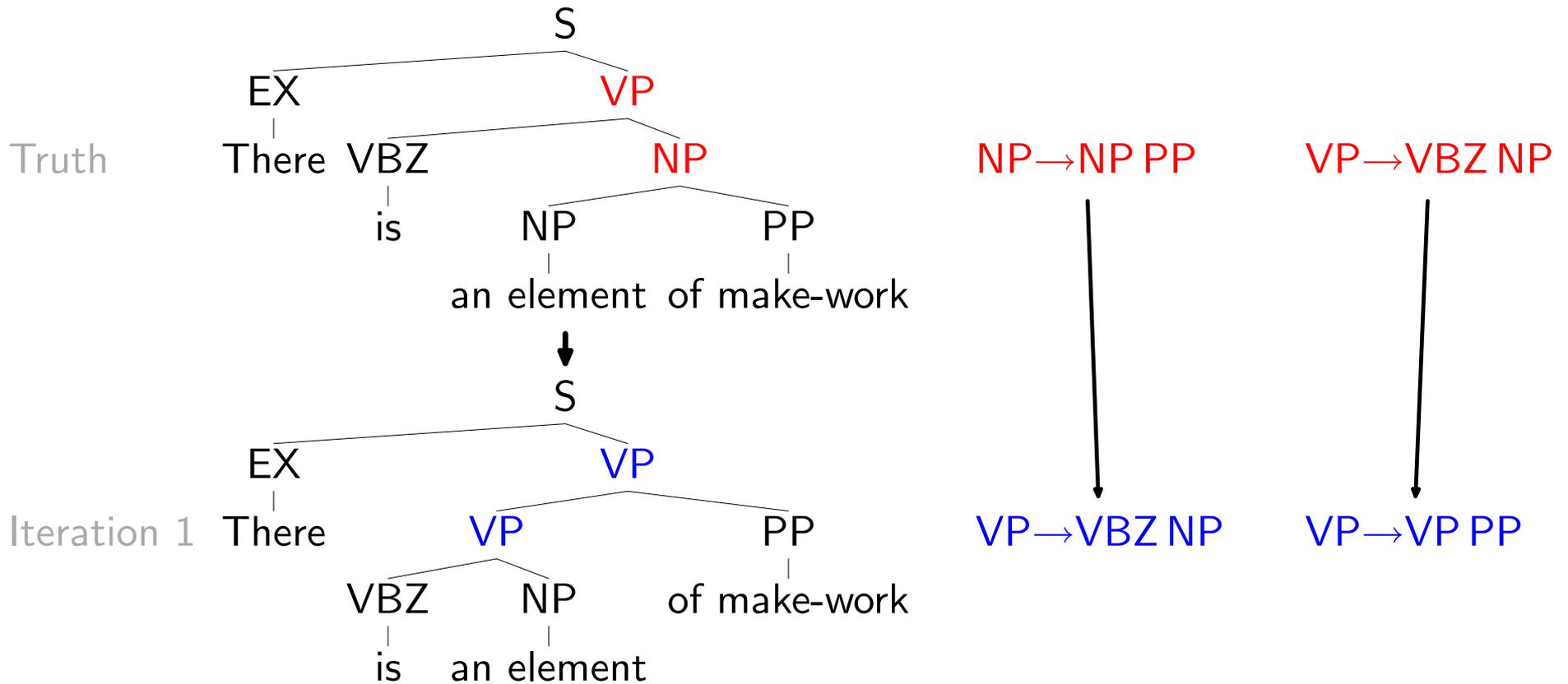
Meta-modeling for PCFGs



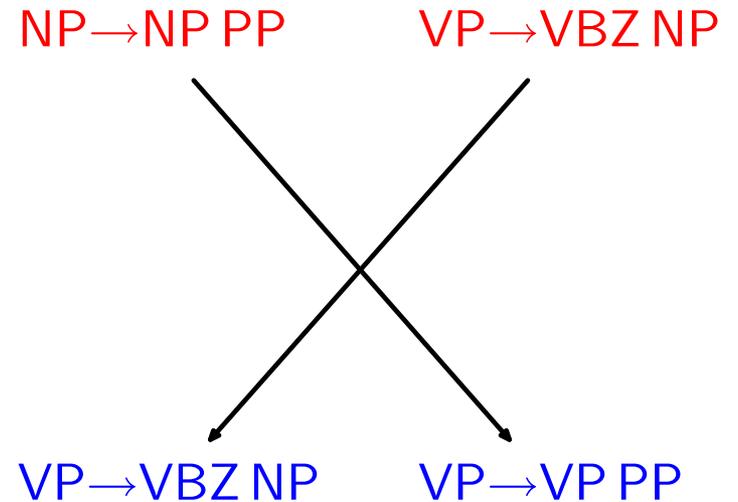
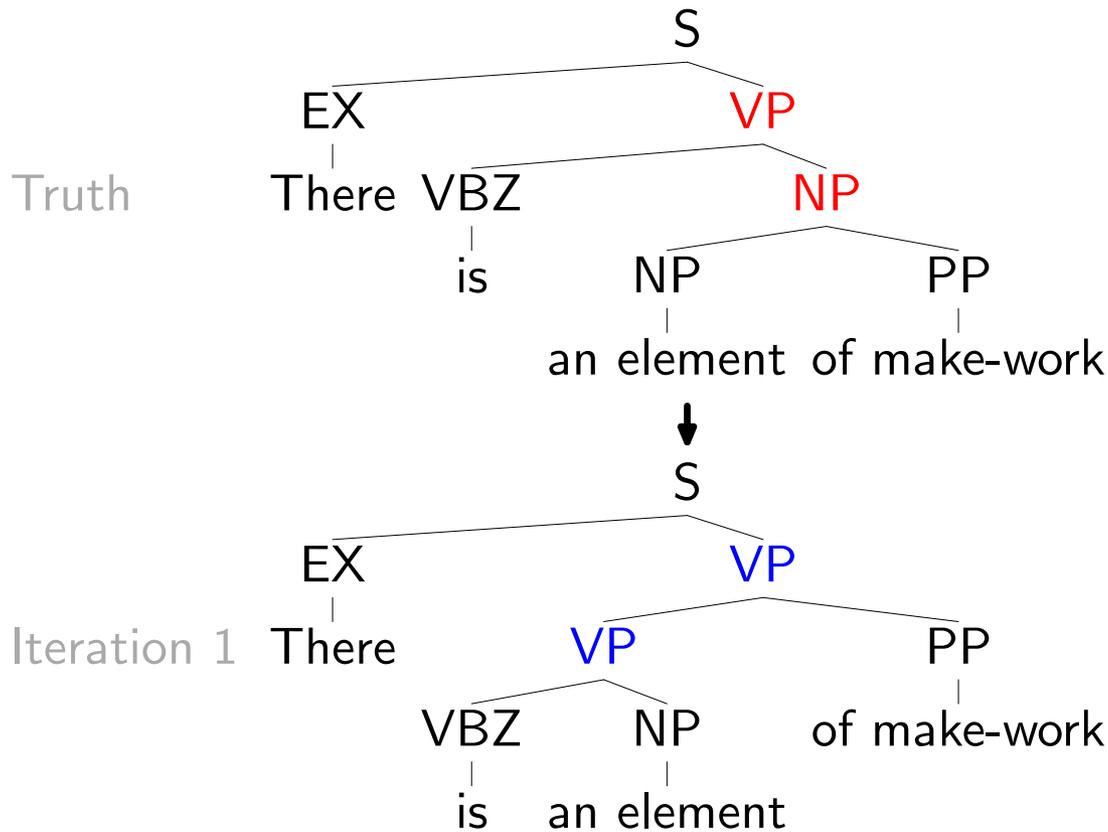
Meta-modeling for PCFGs



Meta-modeling for PCFGs

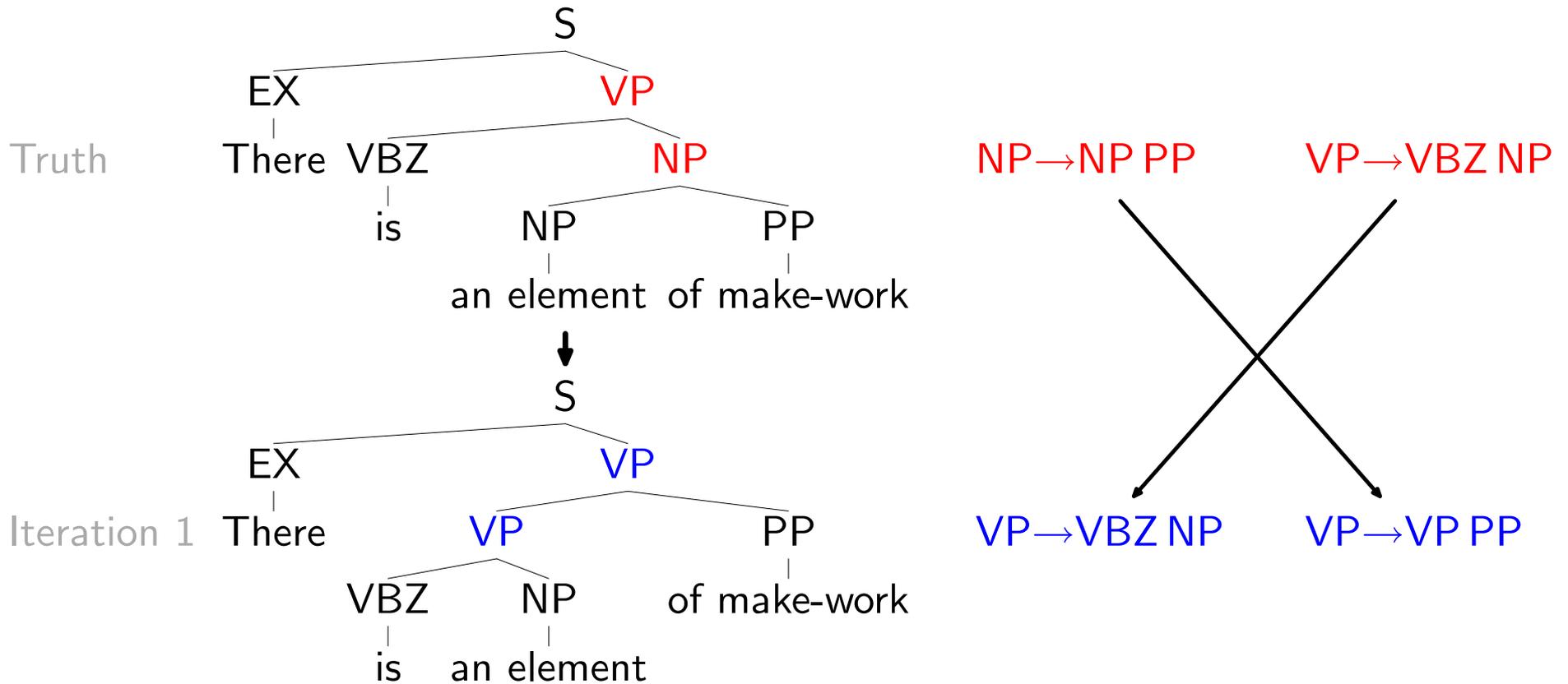


Meta-modeling for PCFGs



Migrations less clear due to uncertainty in tree structure...

Meta-modeling for PCFGs

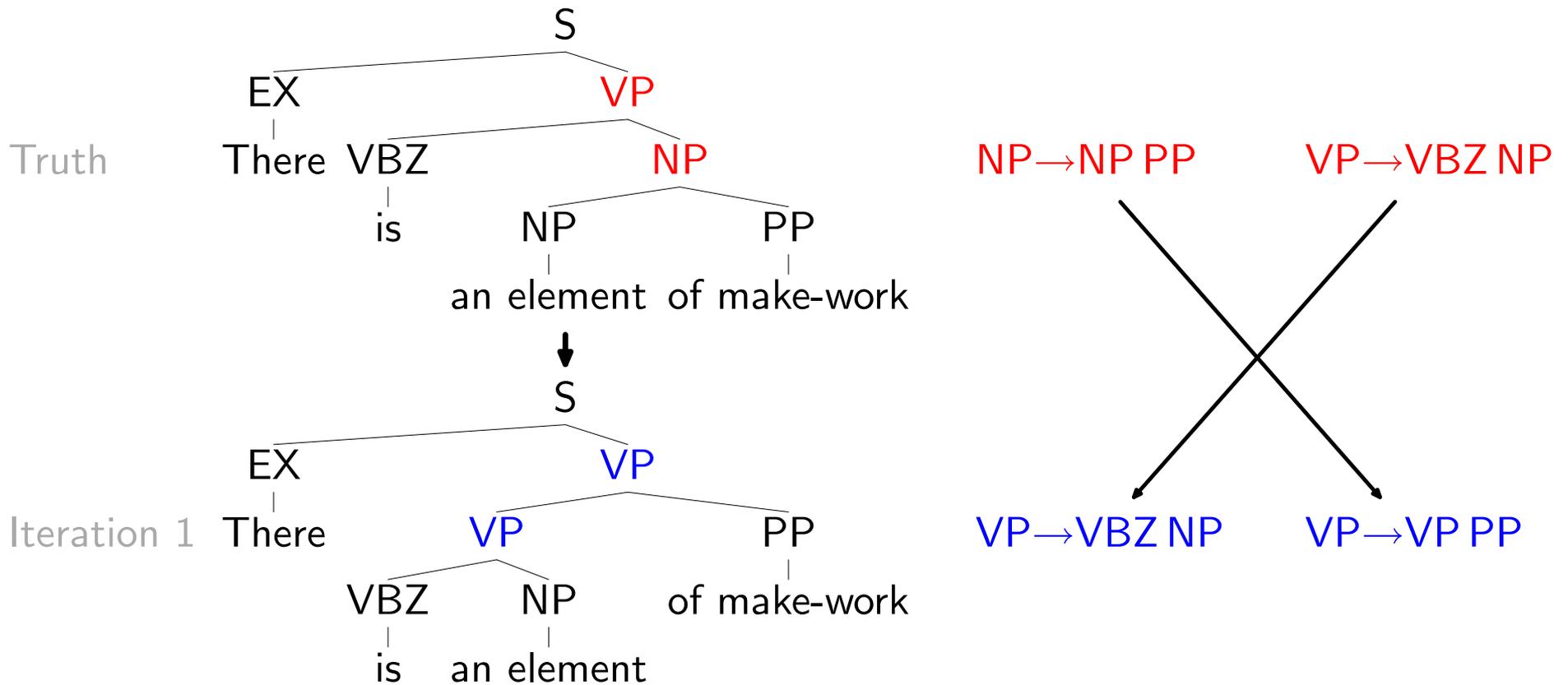


Migrations less clear due to uncertainty in tree structure...

Our approach: use a meta-model

- Migrations are hidden alignments to be learned
- Fit using EM

Meta-modeling for PCFGs



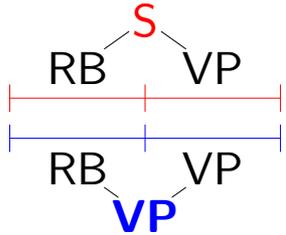
Migrations less clear due to uncertainty in tree structure...

Our approach: use a meta-model

- Migrations are hidden alignments to be learned
- Fit using EM (convex, similar to IBM model 1)

Top PCFG migrations learned by meta-model

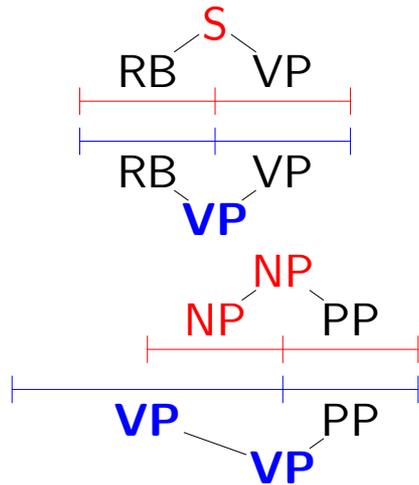
Iteration 1



Sentential adverbs \rightarrow VP adverbs

Top PCFG migrations learned by meta-model

Iteration 1

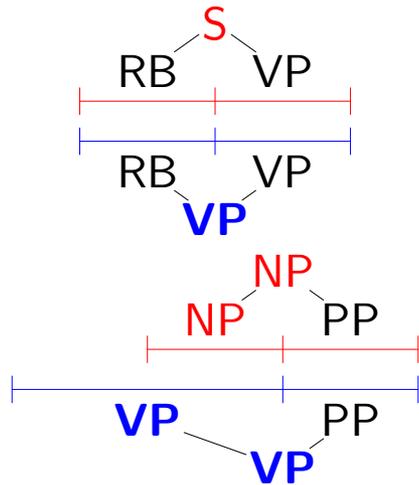


Sentential adverbs → VP adverbs

PPs raised from NPs to verbal level

Top PCFG migrations learned by meta-model

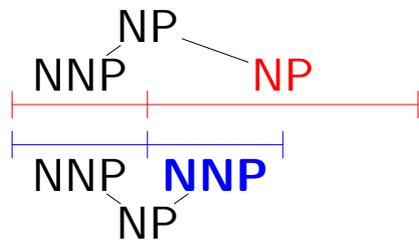
Iteration 1



Sentential adverbs → VP adverbs

PPs raised from NPs to verbal level

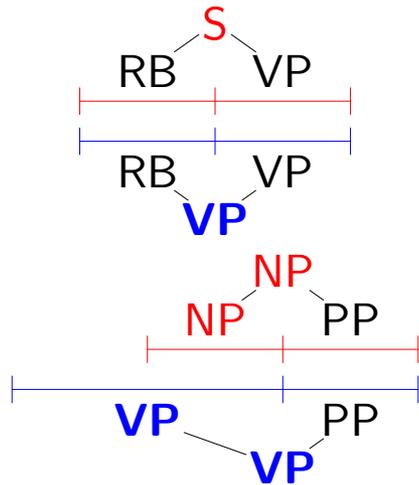
Iteration 2



Right-branching → left-branching structures

Top PCFG migrations learned by meta-model

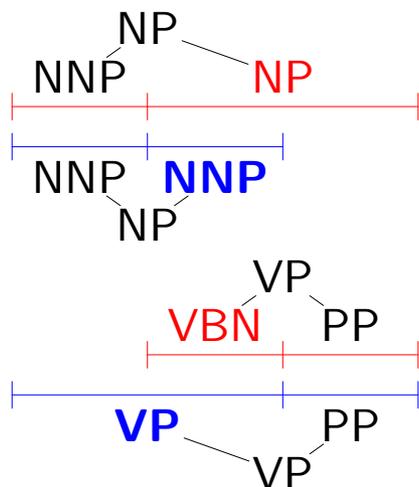
Iteration 1



Sentential adverbs → VP adverbs

PPs raised from NPs to verbal level

Iteration 2



Right-branching → left-branching structures

PP raised to higher VP

Meta-modeling summary

- Meta-model: a diagnostic tool to analyze errors systematically

Meta-modeling summary

- Meta-model: a diagnostic tool to analyze errors systematically
- General phenomenon: **regularization** of syntactic structure

Meta-modeling summary

- Meta-model: a diagnostic tool to analyze errors systematically
- General phenomenon: **regularization** of syntactic structure

✓ Approximation error
Identifiability error
Estimation error
Optimization error

Identifiability error

\mathbf{x} : input sentence

\mathbf{y} : hidden output

$p_{\theta}(\mathbf{x}, \mathbf{y})$: joint distribution with parameters θ

Identifiability error

\mathbf{x} : input sentence

\mathbf{y} : hidden output

$p_{\theta}(\mathbf{x}, \mathbf{y})$: joint distribution with parameters θ

Non-identifiability:



Identifiability error

\mathbf{x} : input sentence

\mathbf{y} : hidden output

$p_{\theta}(\mathbf{x}, \mathbf{y})$: joint distribution with parameters θ

Non-identifiability:

Learning is indifferent...

$$p_{\theta_1}(\mathbf{x}) = p_{\theta_2}(\mathbf{x})$$
$$\boxed{\theta_1} \quad ? \quad \boxed{\theta_2}$$

Identifiability error

\mathbf{x} : input sentence

\mathbf{y} : hidden output

$p_{\theta}(\mathbf{x}, \mathbf{y})$: joint distribution with parameters θ

Non-identifiability:

Learning is indifferent...

$$p_{\theta_1}(\mathbf{x}) = p_{\theta_2}(\mathbf{x})$$

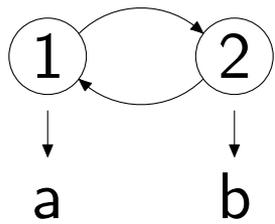
$$\boxed{\theta_1} \quad ? \quad \boxed{\theta_2}$$

but matters to prediction (bad!)

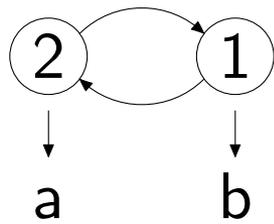
$$p_{\theta_1}(\mathbf{y} | \mathbf{x}) \neq p_{\theta_2}(\mathbf{y} | \mathbf{x})$$

Examples of non-identifiability

- Label symmetries



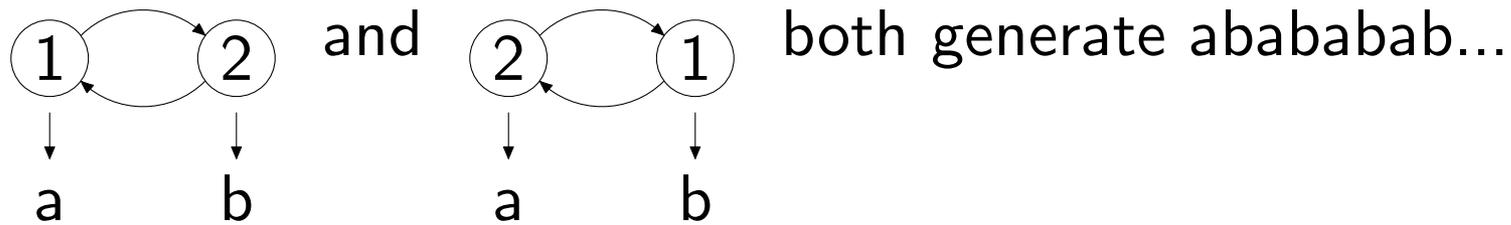
and



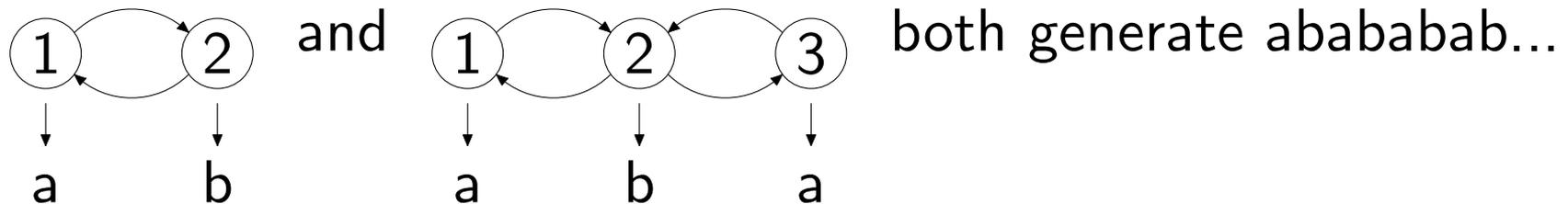
both generate abababab...

Examples of non-identifiability

- Label symmetries

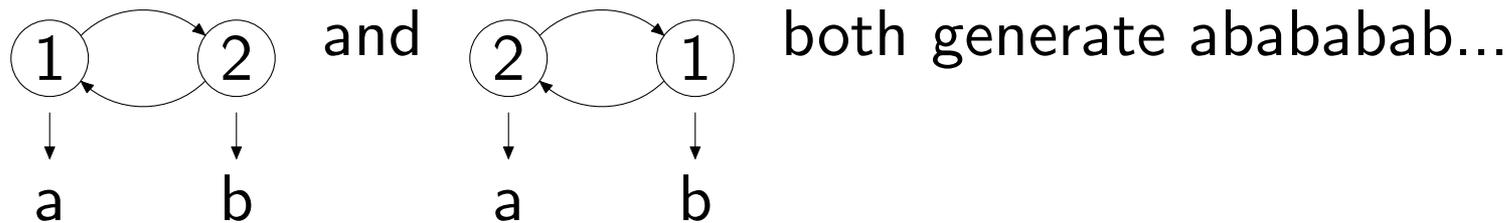


- K -state HMM (if true distribution is $< K$ -state HMM)

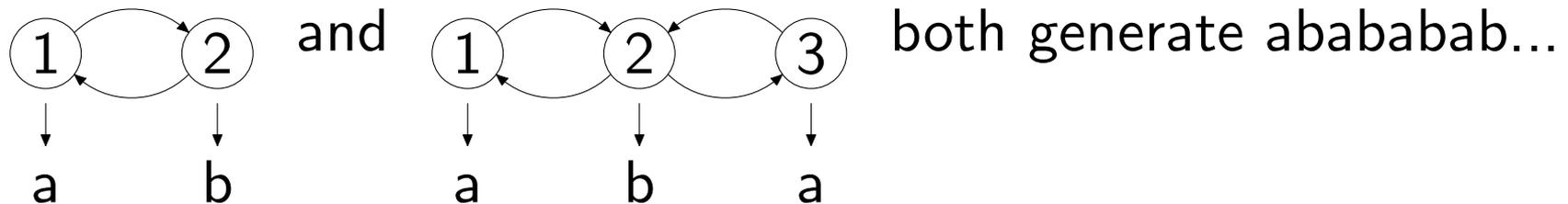


Examples of non-identifiability

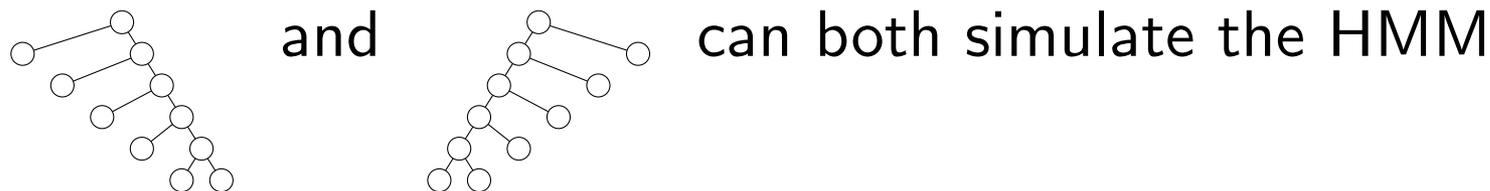
- Label symmetries



- K -state HMM (if true distribution is $< K$ -state HMM)

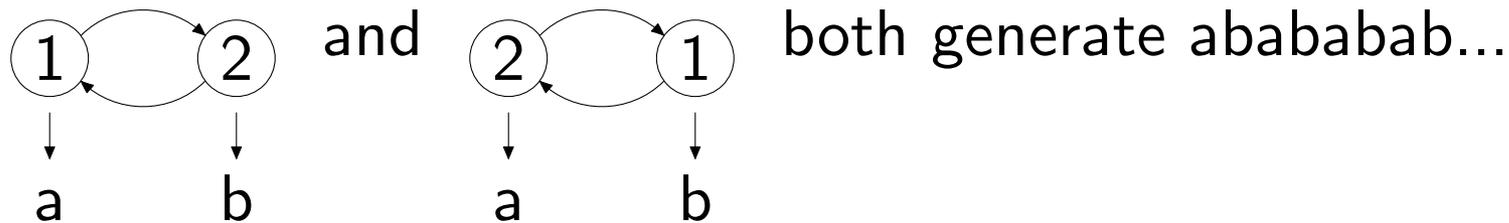


- PCFG (if true distribution is HMM)

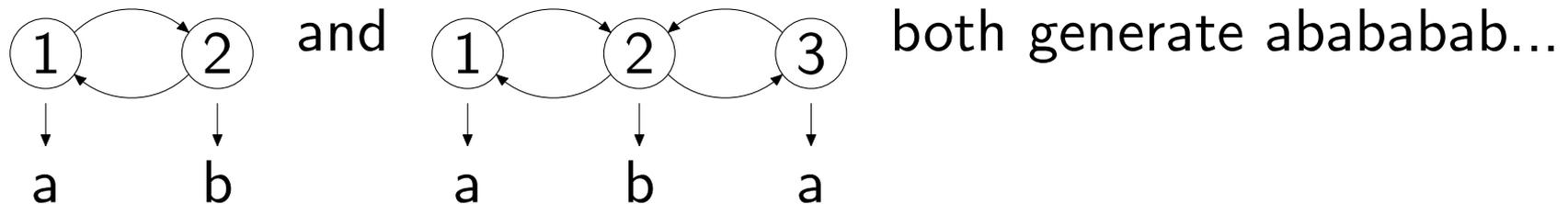


Examples of non-identifiability

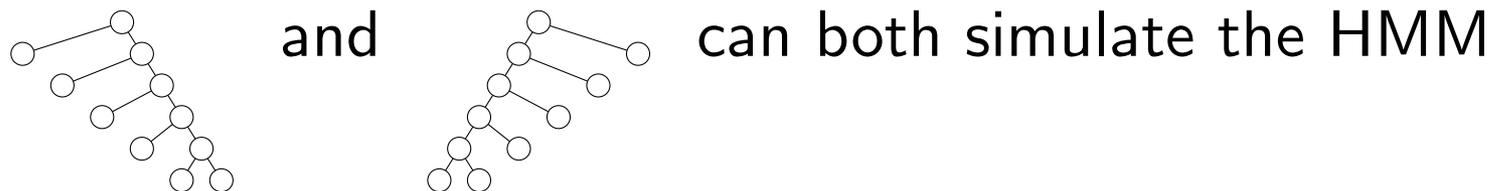
- Label symmetries



- K -state HMM (if true distribution is $< K$ -state HMM)



- PCFG (if true distribution is HMM)



Real data is complex, so last two are not an issue

Identifiability and distance

Given θ_1 and θ_2 , how to measure distance between them?

Want distance $\left(\begin{array}{c} \text{Diagram 1} \\ \text{Diagram 2} \end{array} \right) = 0$

The diagram shows two directed graphs side-by-side, enclosed in large parentheses. The left graph has two nodes, 1 and 2, in a row. Node 1 has a downward arrow to the label 'a'. Node 2 has a downward arrow to the label 'b'. There are directed edges from 1 to 2 and from 2 to 1. The right graph has two nodes, 2 and 1, in a row. Node 2 has a downward arrow to the label 'a'. Node 1 has a downward arrow to the label 'b'. There are directed edges from 2 to 1 and from 1 to 2. A comma is placed between the two graphs. To the right of the closing parenthesis is an equals sign followed by a zero.

Identifiability and distance

Given θ_1 and θ_2 , how to measure distance between them?

Want distance $\left(\begin{array}{cc} \text{1} & \text{2} & \text{2} & \text{1} \\ \downarrow & \downarrow & \downarrow & \downarrow \\ \text{a} & \text{b} & \text{a} & \text{b} \end{array} \right) = 0$

- Computing **label-permutation invariant distance** is NP-hard
- We use bipartite matching to find lower and upper bounds

Identifiability and distance

Given θ_1 and θ_2 , how to measure distance between them?

Want distance $\left(\begin{array}{cc} \text{1} & \text{2} & \text{2} & \text{1} \\ \downarrow & \downarrow & \downarrow & \downarrow \\ \text{a} & \text{b} & \text{a} & \text{b} \end{array} \right) = 0$

- Computing **label-permutation invariant distance** is NP-hard
- We use bipartite matching to find lower and upper bounds

✓ Approximation error
✓ Identifiability error
Estimation error
Optimization error

Estimation and optimization errors

Experiment setup:

- Take some parameters θ^* (say, supervised estimate on real data)
- Use θ^* to generate synthetic data
- Can we recover θ^* using EM?

Estimation and optimization errors

Experiment setup:

- Take some parameters θ^* (say, supervised estimate on real data)
- Use θ^* to generate synthetic data
- Can we recover θ^* using EM?



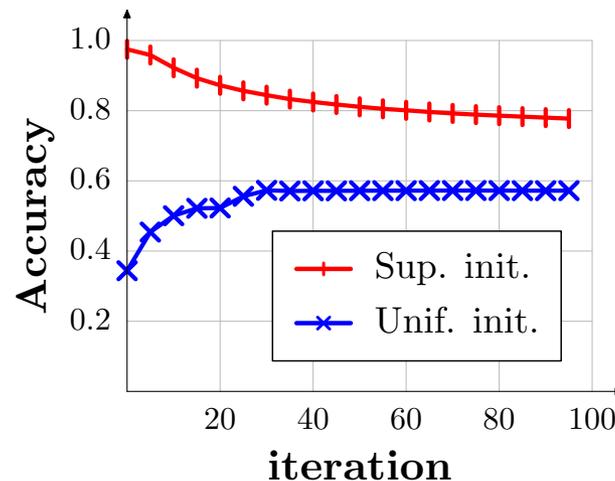
Estimation and optimization errors

Experiment setup:

- Take some parameters θ^* (say, supervised estimate on real data)
- Use θ^* to generate synthetic data
- Can we recover θ^* using EM? **No?**



HMM on 5K examples:



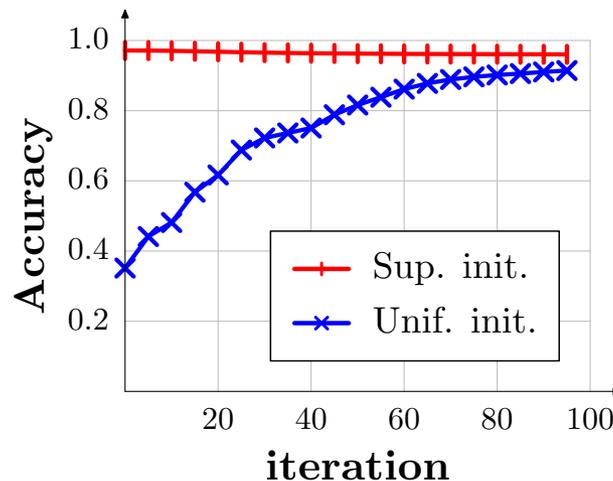
Estimation and optimization errors

Experiment setup:

- Take some parameters θ^* (say, supervised estimate on real data)
- Use θ^* to generate synthetic data
- Can we recover θ^* using EM? **Yes!**



HMM on 500K examples:

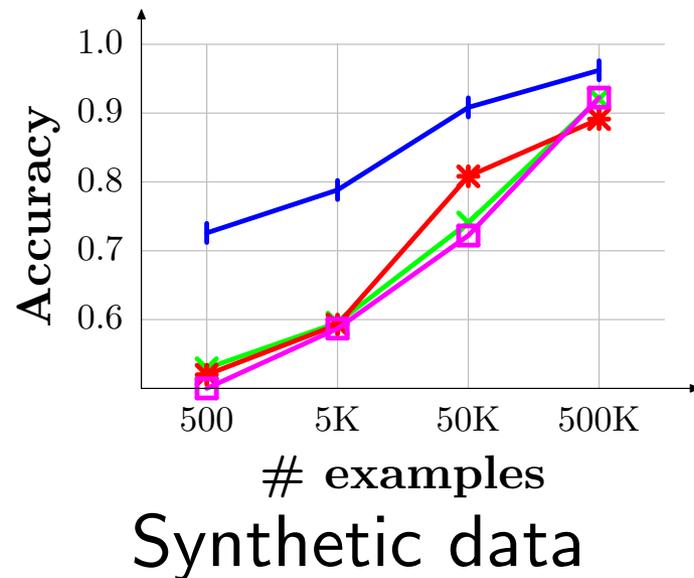


Optimization error decreases with more data

On HMM model (similar for PCFG and a dependency model):

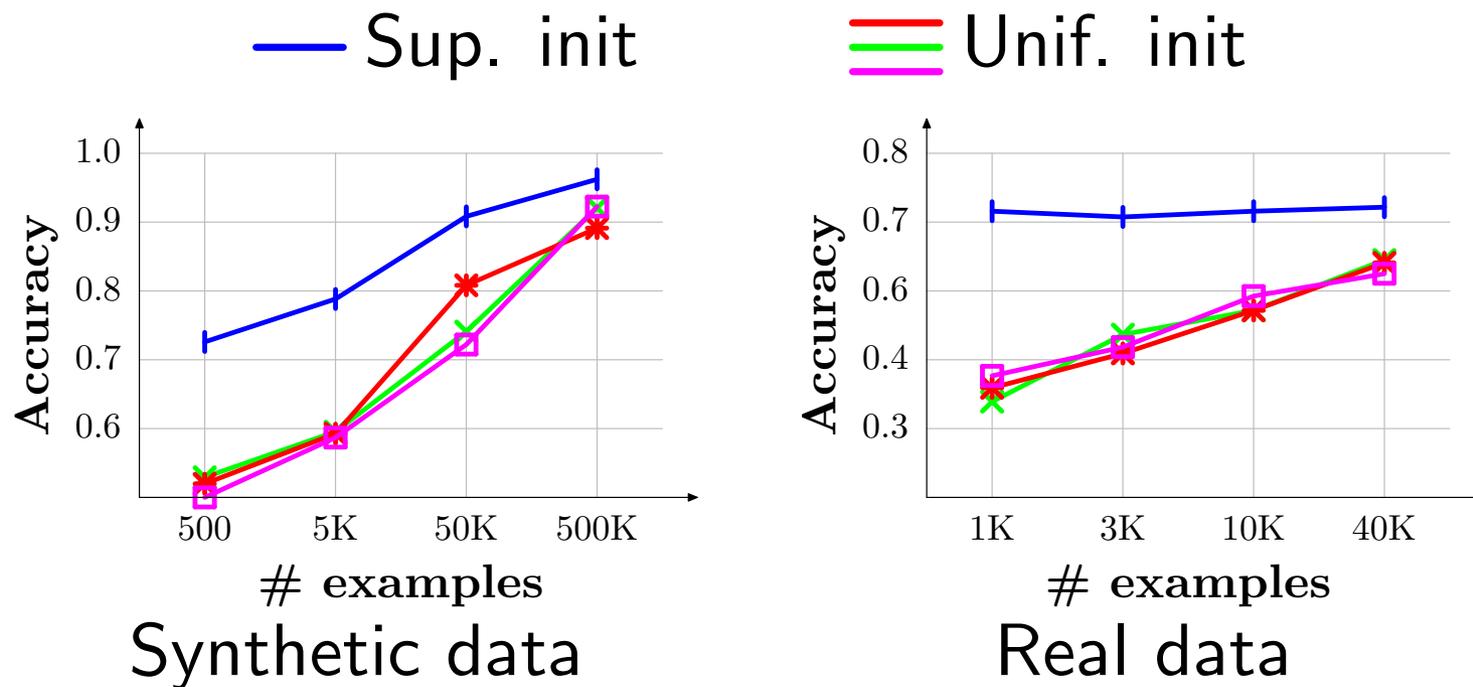
— Sup. init

— Unif. init



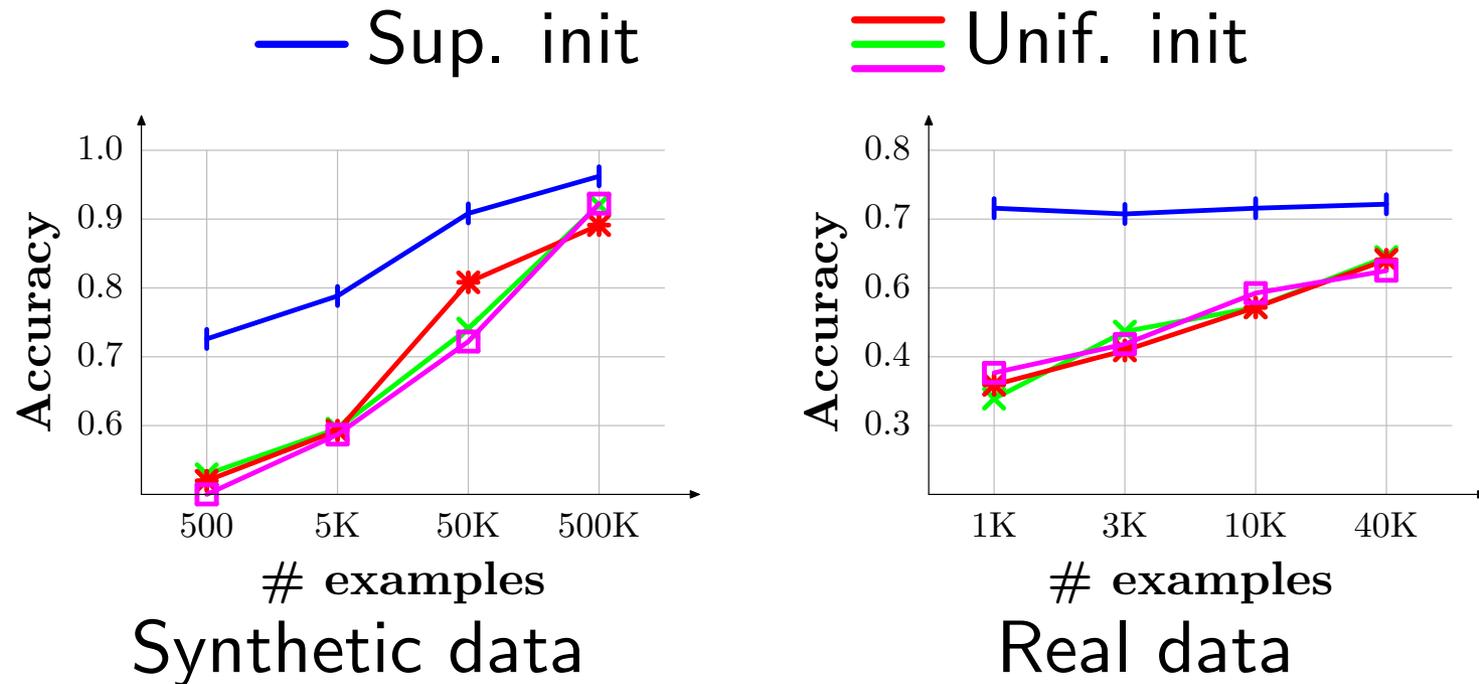
Optimization error decreases with more data

On HMM model (similar for PCFG and a dependency model):



Optimization error decreases with more data

On HMM model (similar for PCFG and a dependency model):



Why does this phenomenon happen?

- **Intuition:** with more data, EM can pick up the salient patterns more easily
- Was also shown for mixture of Gaussians [Srebro, 2006]

Summary

- ✓ Approximation error

 - Meta-model: tool for systematic error analysis

Summary

- ✓ Approximation error

 - Meta-model: tool for systematic error analysis

- ✓ Identifiability error

 - Distance robust to label symmetries

Summary

- ✓ Approximation error
 - Meta-model: tool for systematic error analysis
- ✓ Identifiability error
 - Distance robust to label symmetries
- ✓ Estimation error
 - Decreases with more data

Summary

- ✓ Approximation error
 - Meta-model: tool for systematic error analysis
- ✓ Identifiability error
 - Distance robust to label symmetries
- ✓ Estimation error
 - Decreases with more data
- ✓ Optimization error
 - Decreases with more data!