

# Outline



Vlad Eidelman



Ziyuan Wang



Adam Lopez



Jon Graehl



ThuyLinh Nguyen

- 3:20pm Parametric models: posterior regularisation. Desai
- 3:35pm Training models with rich features spaces. Vlad
- 3:50pm Decoding with complex grammars. Adam
- 4:20pm Closing remarks. Phil
- 4:25pm Finish.

# Phrase Clustering with Posterior Regularization

CLSP Summer Workshop 2010  
SMT Team  
Desai Chen  
joint work with Trevor Cohn

# Outline

- clustering problem
- EM with posterior regularization
- results and future experiments

# Phrase clustering

Phrases are defined as contiguous spans aligned with each other

i 'll bring you some now .

我 这 就 给 您 拿 一 些 。

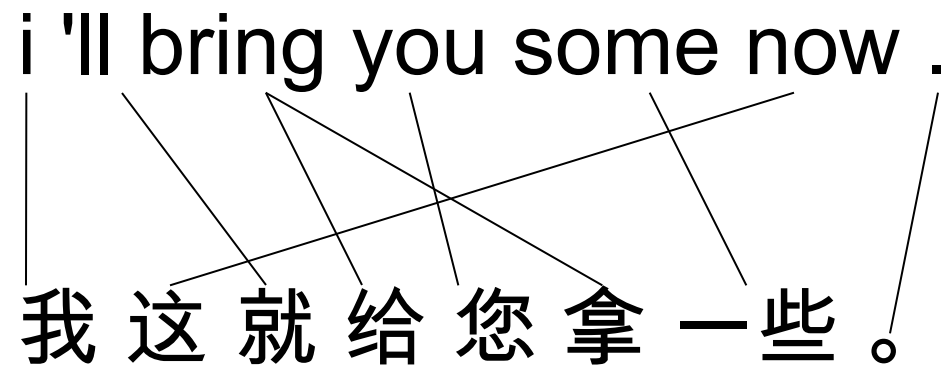
Example from btec



# Phrase clustering

Phrases are defined as contiguous spans aligned with each other

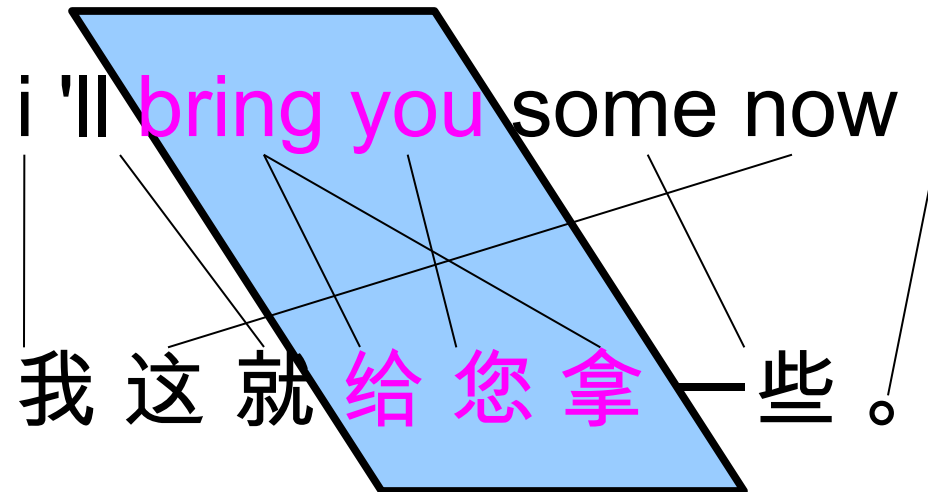
i 'll bring you some now .  
我 这 就 给 您 拿 一 些 。



Example from btec

# Phrase clustering

Phrases are defined as contiguous spans aligned with each other



Example from btec

# Phrase clustering

Phrases are defined as contiguous spans aligned with each other

i 'll bring you some now .

我 这 就 给 您 拿 一 些 。

# Phrase clustering

Contexts are words before or after the phrase:

target side context

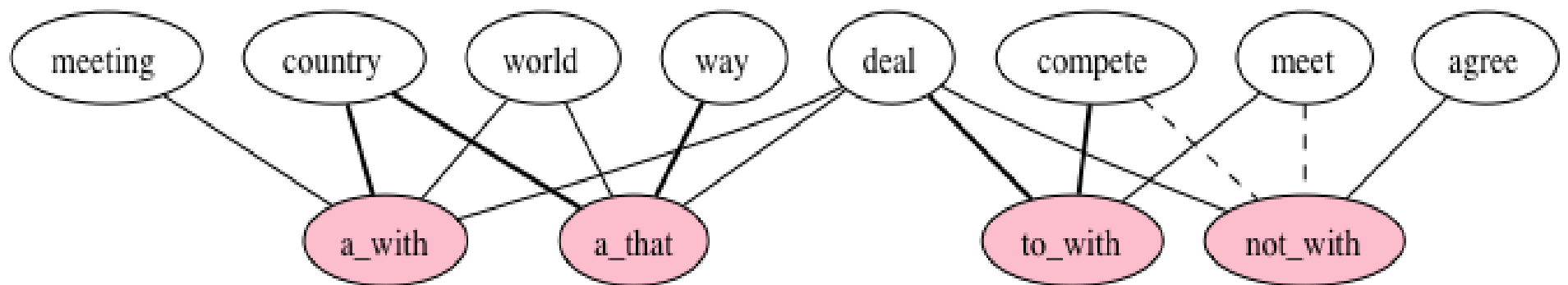
i 'll bring you some now .

我 这 就 给 您 拿 一 些 。

source side context

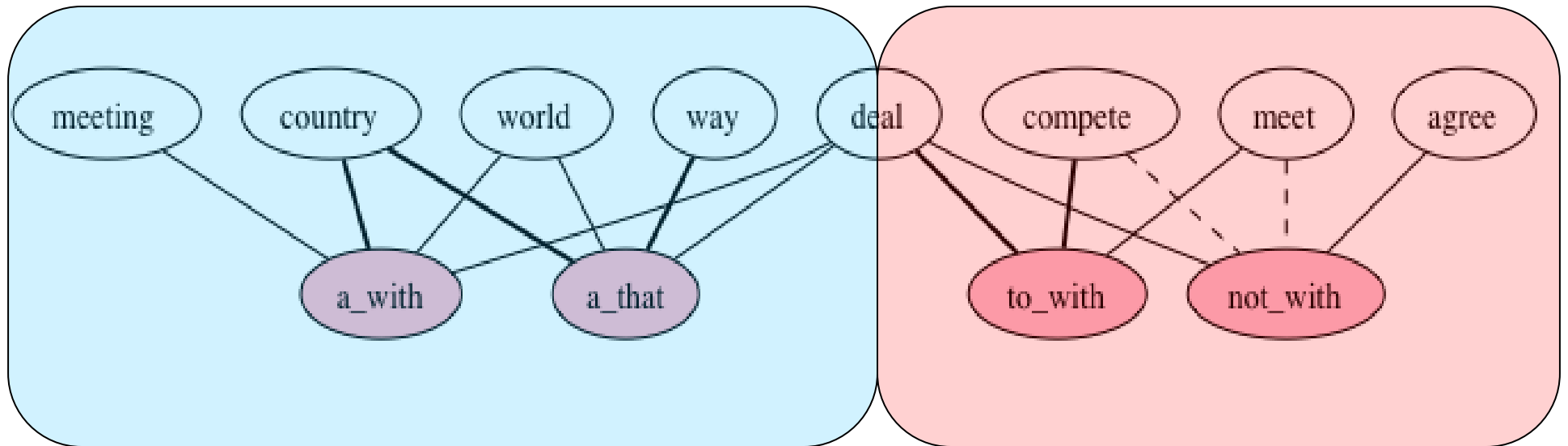
# Objective

Put all phrase-context pairs into categories



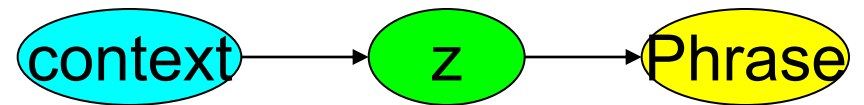
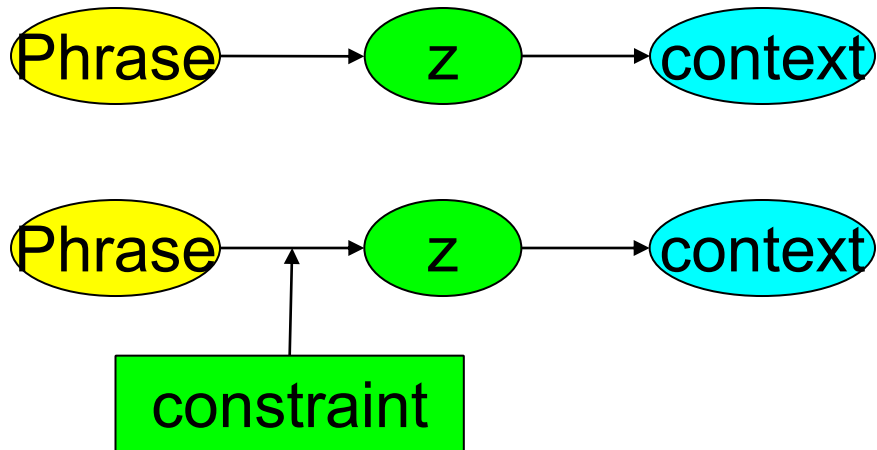
# Objective

Put all phrase-context pairs into categories



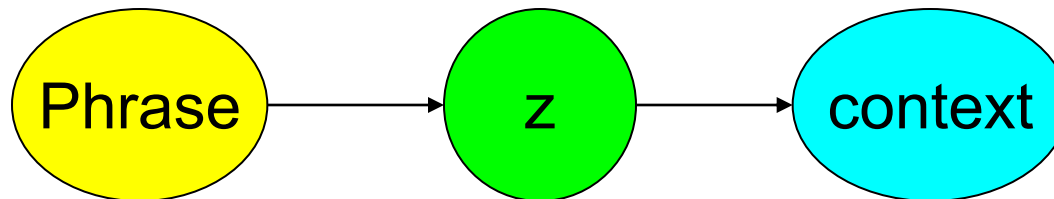
# Outline

- Where do phrases come from?
- **EM with posterior regularization**
- results and future experiment



# Expectation-Maximization

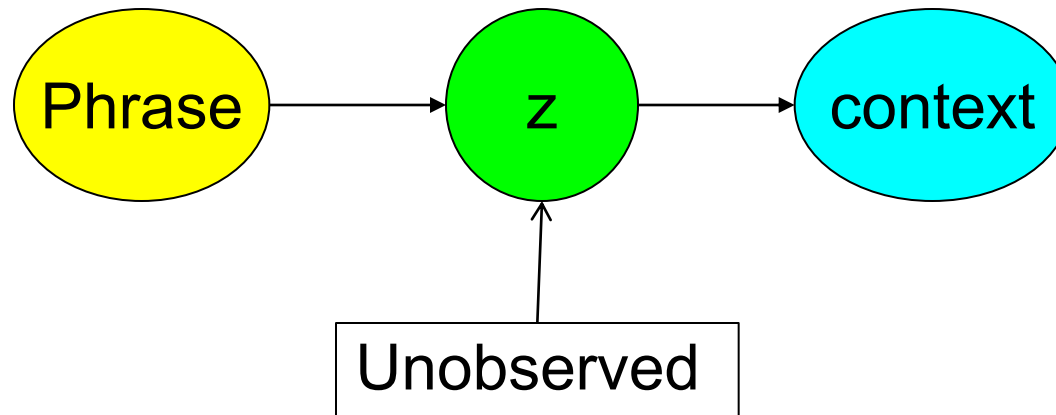
- naïve Bayes model for phrase labeling





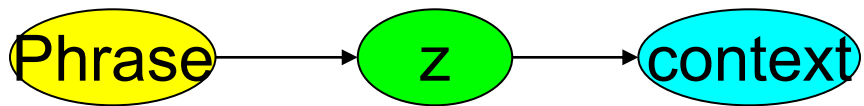
# EM clustering

- naïve Bayes model for phrase labeling



# EM clustering

- naïve Bayes model for phrase labeling

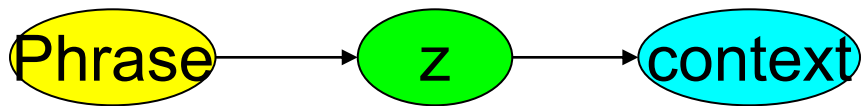


E-step

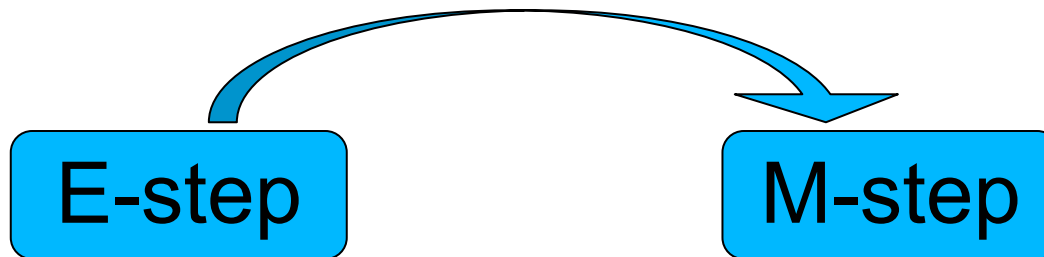
M-step

# EM clustering

- naïve Bayes model for phrase labeling

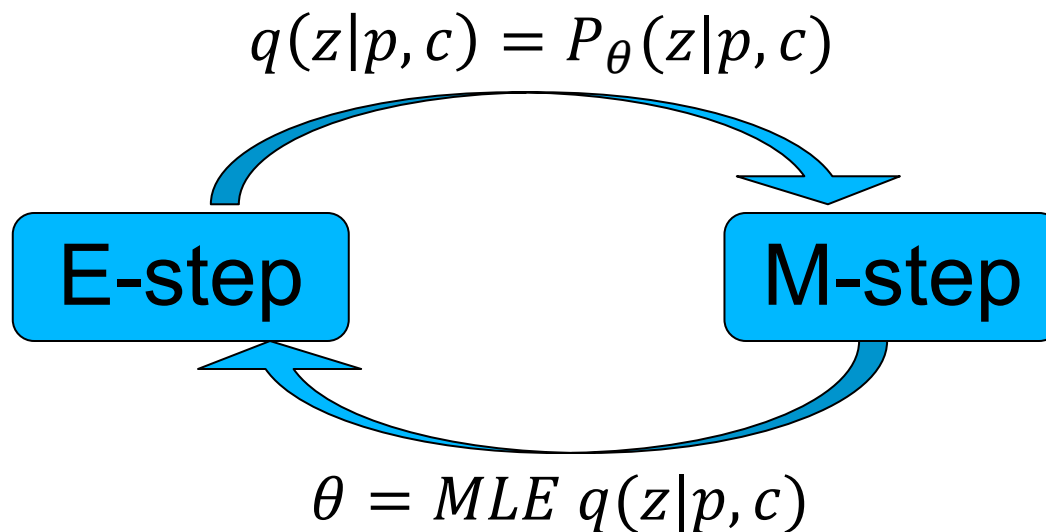
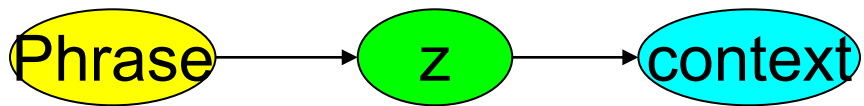


$$q(z|p, c) = P_{\theta}(z|p, c)$$



# EM clustering

- naïve Bayes model for phrase labeling

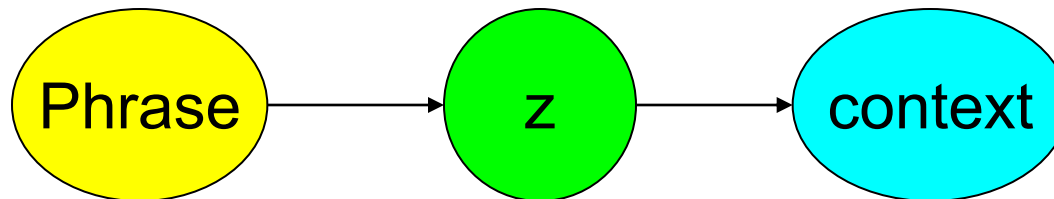


# Problem with EM

- Problem: EM uses as many categories as it wants for each phrase.
- We want to limit the number of categories associated with each phrase.

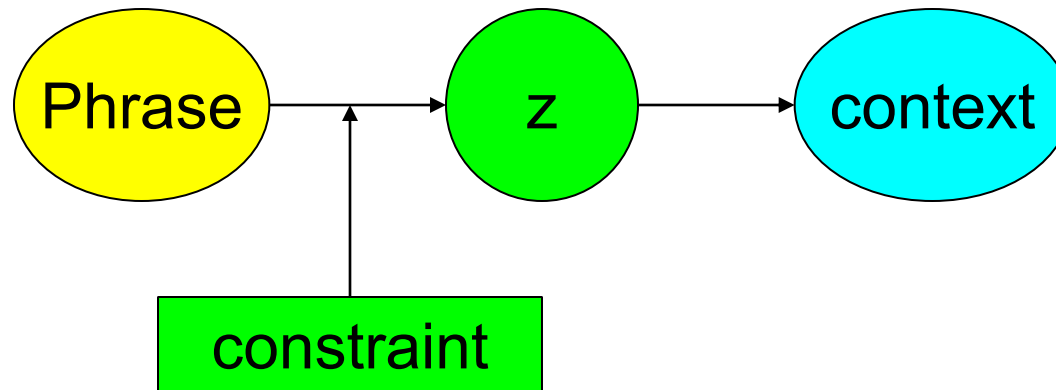
# Sparsity constraints

- Sparsity: Each phrase/context should be labeled with fewer kinds of labels.



# Sparsity constraints

- Sparsity: Each phrase/context should be labeled with fewer kinds of labels.



# Sparsity constraints

Minimize  $\sum_{p,z} \max_i P(z|p_i)$



# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

i understand there are some sightseeing bus tours here , is that right ?

there are only a few seats left in the dress circle .

well , of course there are fine restaurants .

your hotel brochure shows there are some tennis courts at your hotel .

# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

i understand there are some sightseeing bus tours here , is that right ?

there are only a few seats left in the dress circle .

well , of course there are fine restaurants .

your hotel brochure shows there are some tennis courts at your hotel .

# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

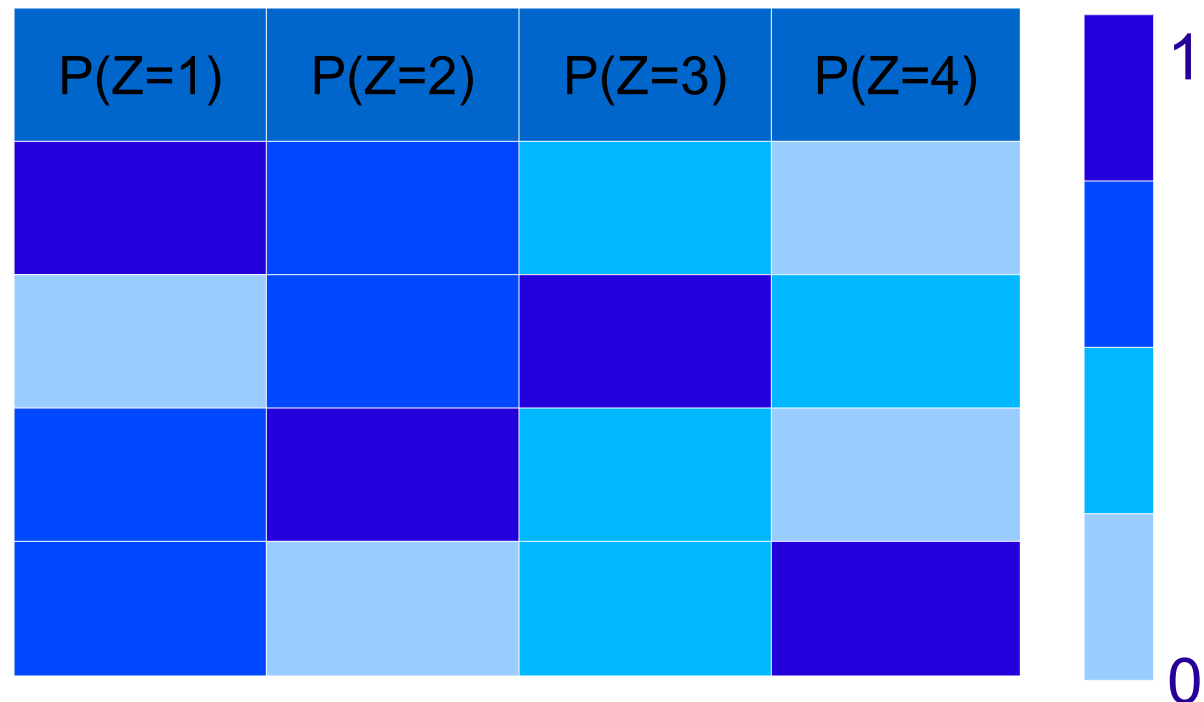
Contexts:

i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis



# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

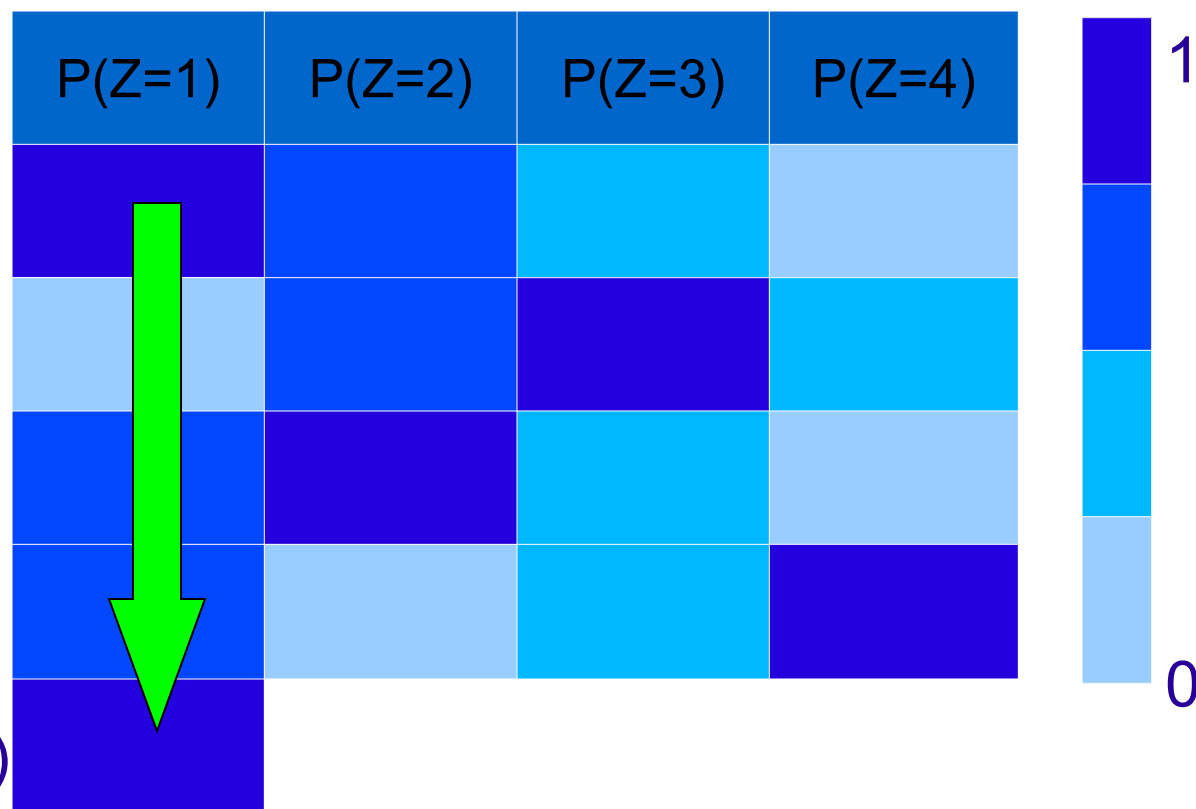
i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

max P(tag|phrase)



# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

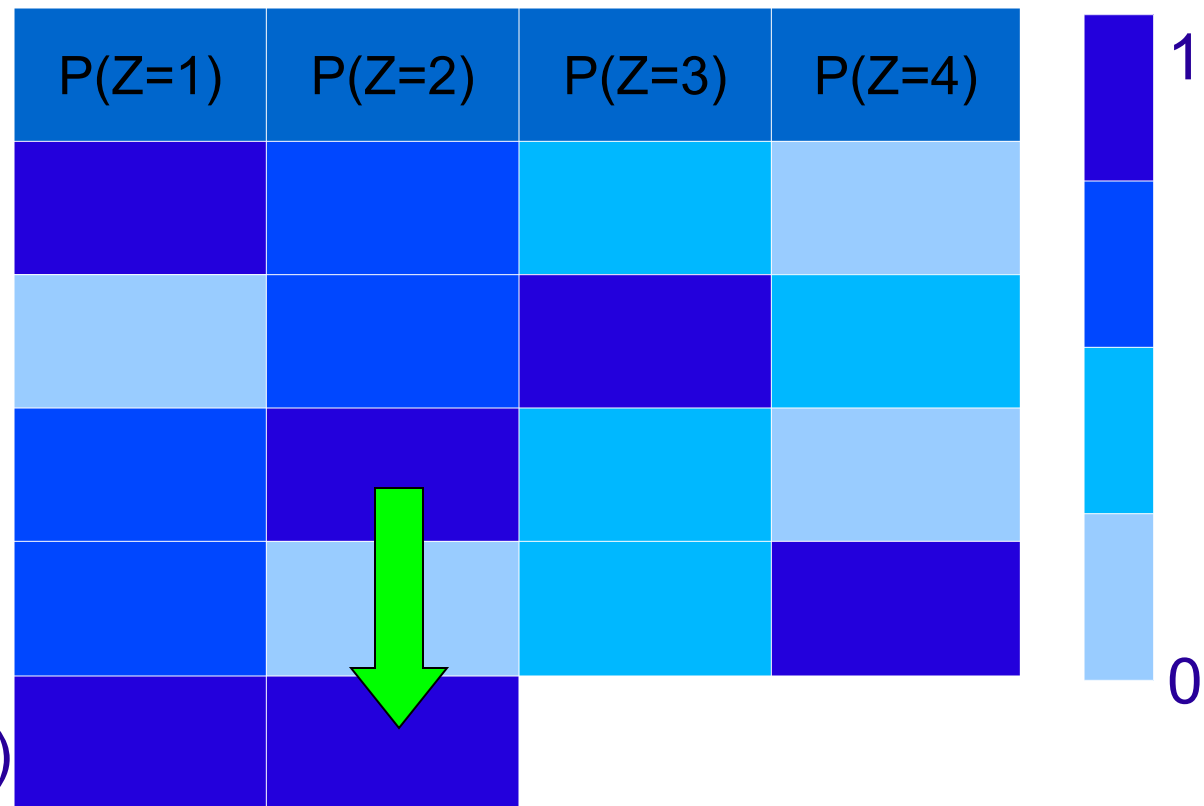
i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

max  $P(\text{tag}|\text{phrase})$



# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

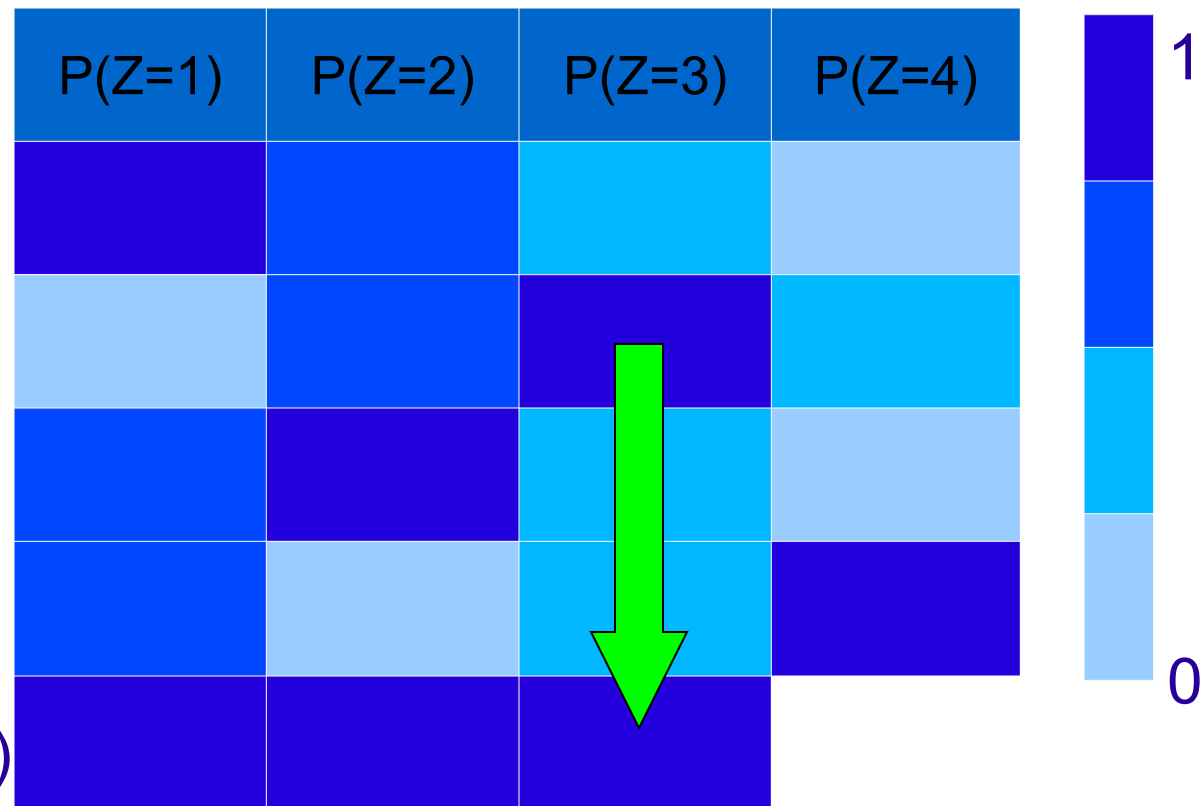
i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

max P(tag|phrase)



# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

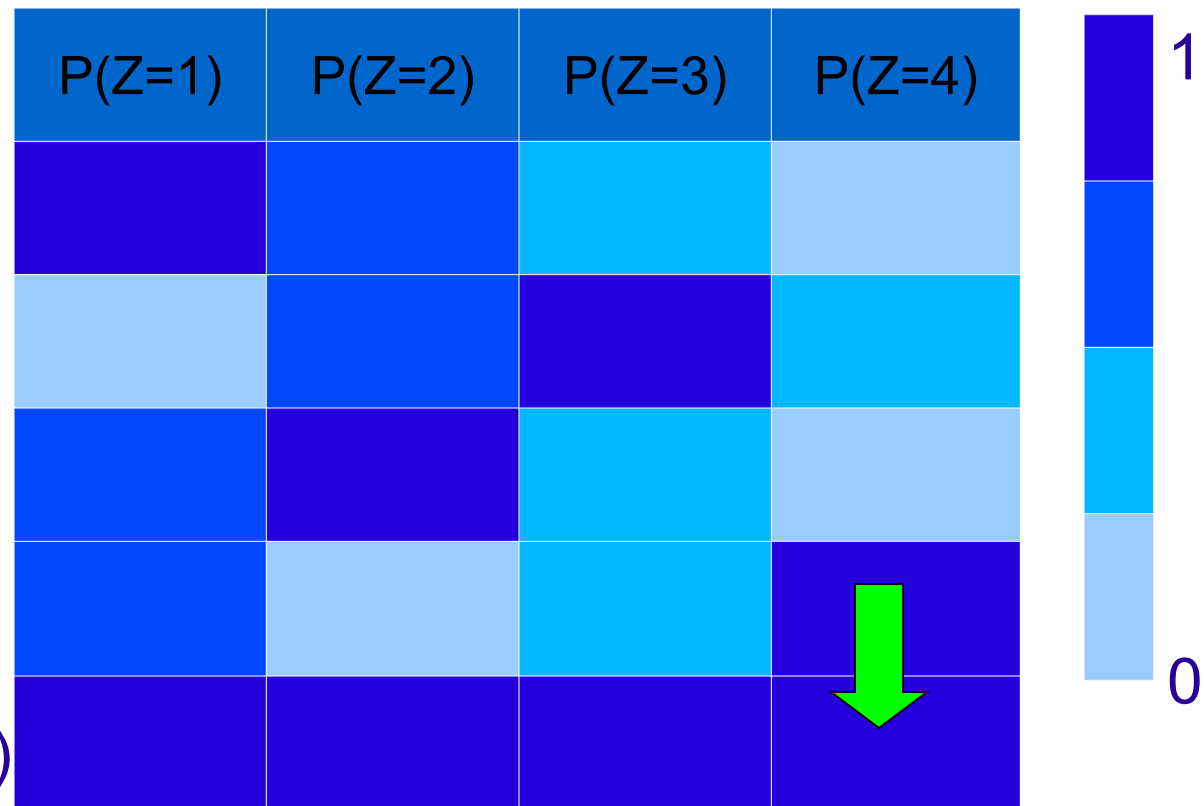
i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

max P(tag|phrase)





# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

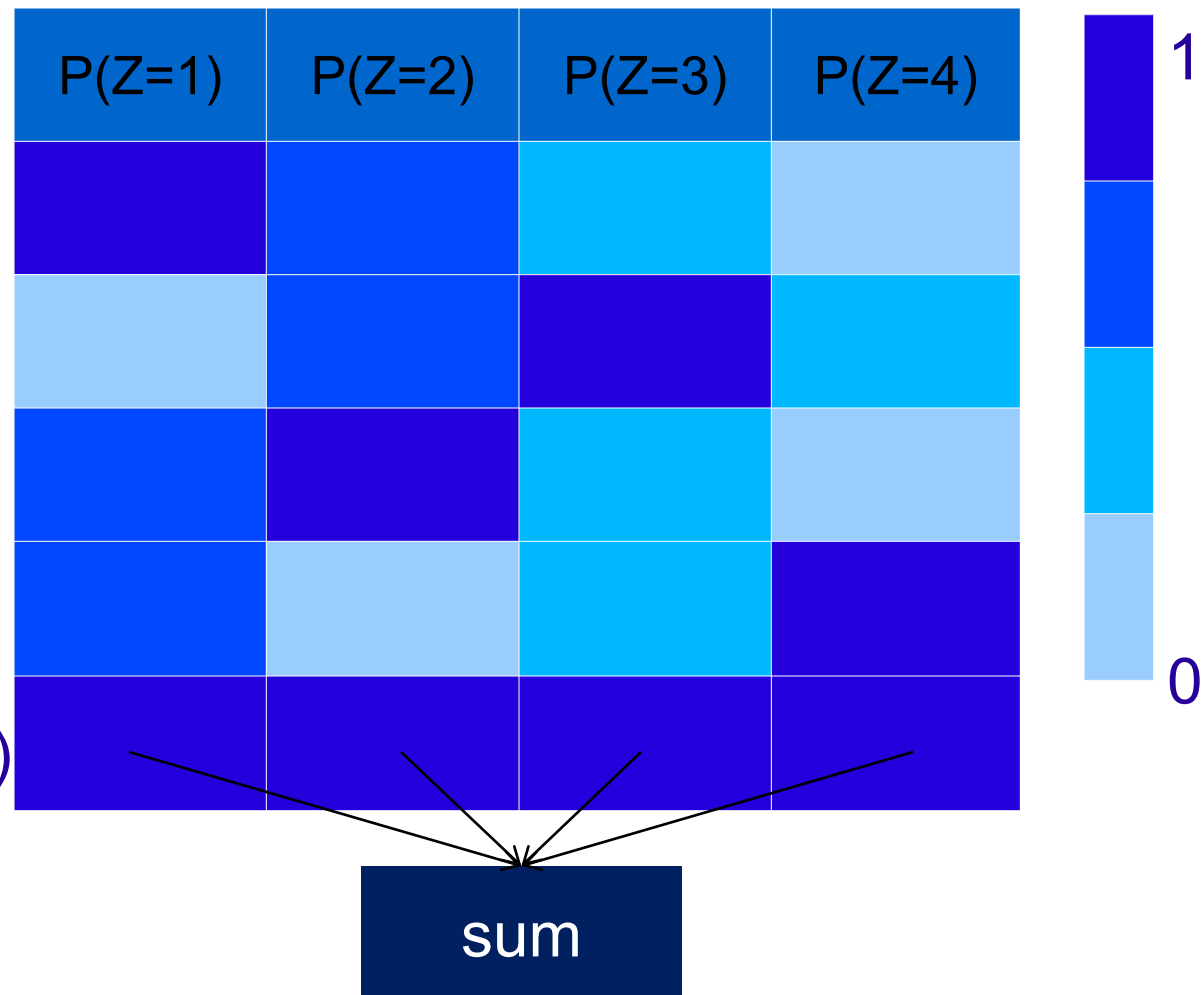
i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

$\max P(\text{tag}|\text{phrase})$



# Sparsity constraints

$$\text{Minimize } \sum_{p,z} \max_i P(z|p_i)$$

Phrase: there are

Contexts:

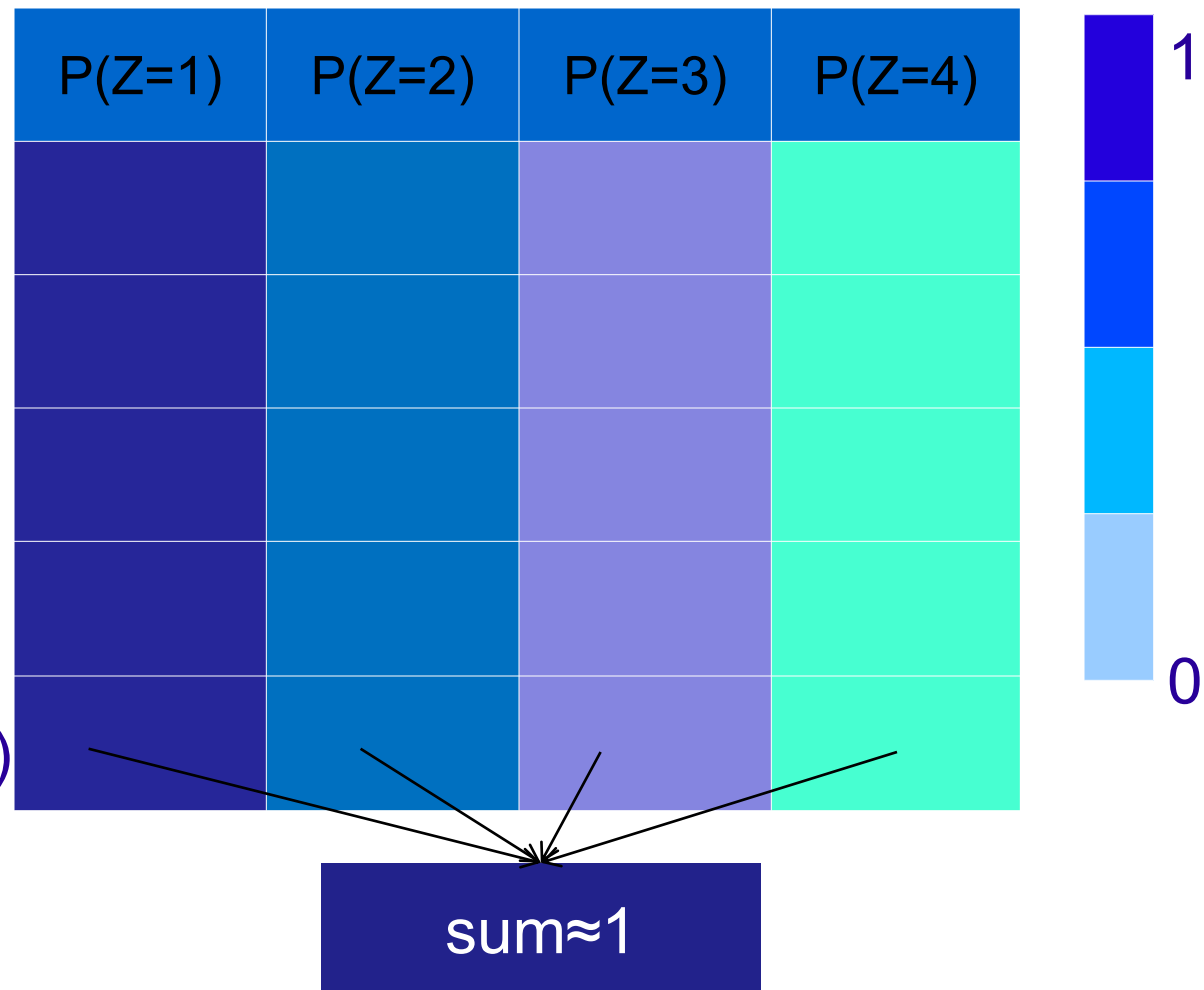
i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

$\max P(\text{tag}|\text{phrase})$

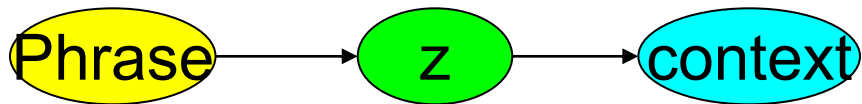


# Posterior Regularization

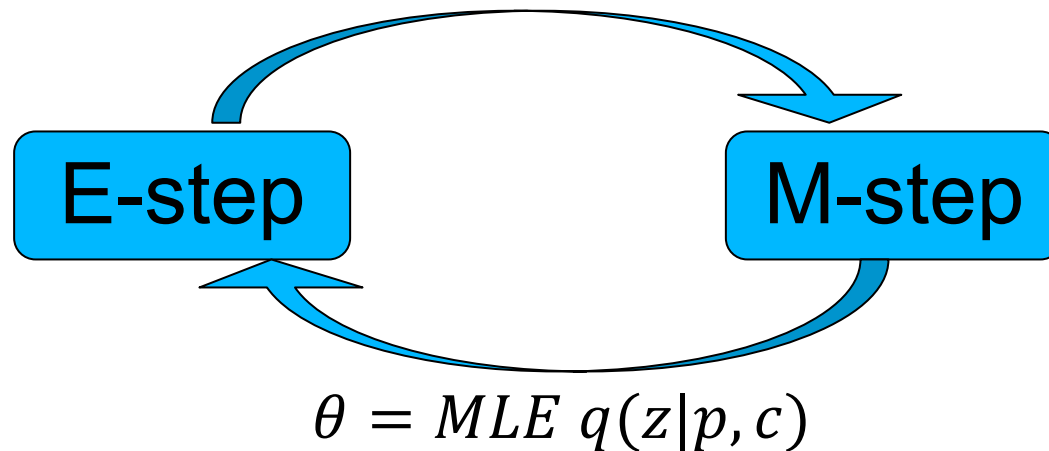
- Follows *Posterior Regularization for Structured Latent Variable Models*, Ganchev et al., 2009
- During E-step, impose constraints on the posterior  $q$  to guide the search

# Posterior Regularization

- impose constraints on the posterior  $q$

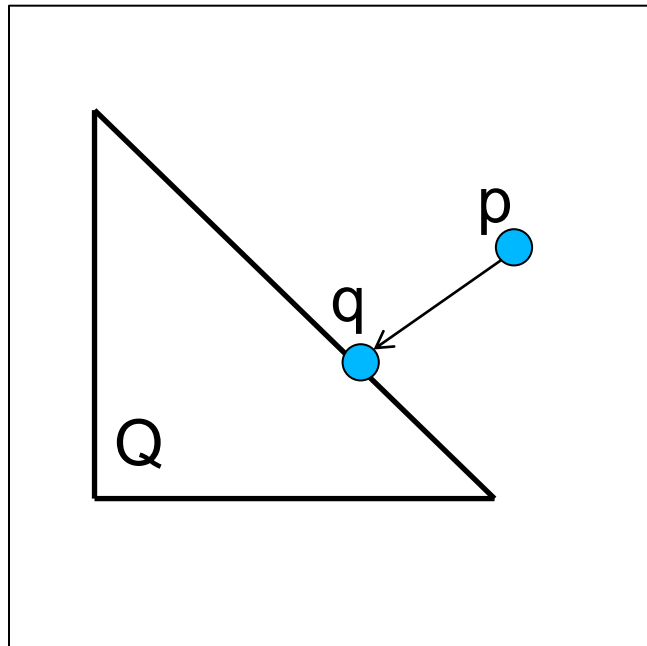
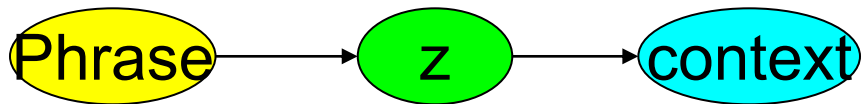


$$q(z|p, c) = \arg \min_{q \in Q} KL(q || P_{\theta})$$

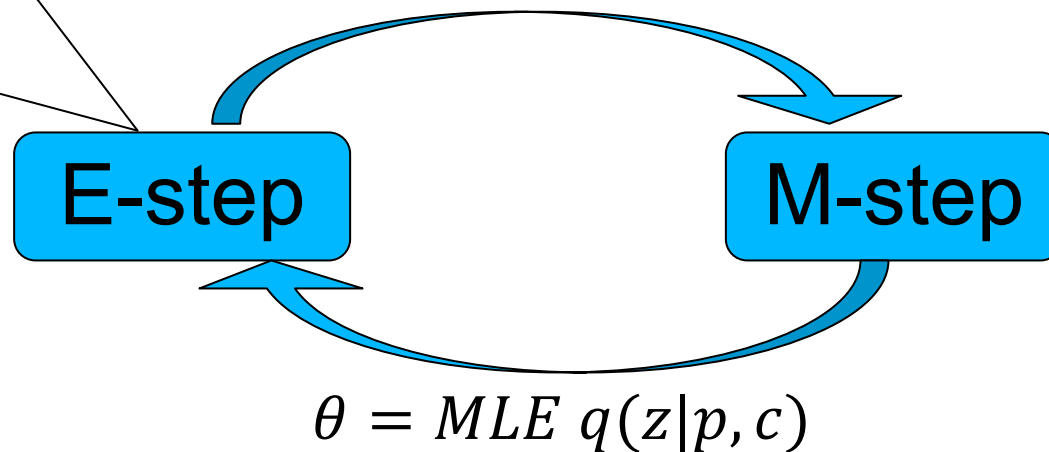


# Posterior Regularization

- impose constraints on the posterior  $q$



$$q(z|p, c) = \arg \min_{q \in Q} KL(q || P_{\theta})$$



# Sparsity constraints

Minimize  $\sum_{p,z} \max_i P(z|p_i)$

Phrase: like this

Contexts:

i understand \_ some  
sightseeing

<s> <s> \_ only a

of course \_ fine  
restaurants

brochure shows \_  
some tennis

Define feature functions:

$$\phi_{i,j}(p, z) = \begin{cases} 1 & \text{if } p = i \text{ and } z = j \\ 0 & \text{otherwise} \end{cases}$$

# Sparsity constraints

Minimize  $\sum_{p,z} \max_i P(z|p_i)$

- Soft constraint. Softness controlled by  $\sigma$ .
- During E-step, find q distribution:

$$\begin{aligned} \min_{q, c_{p,z}} & KL(q || P_{\theta}) + \sigma \sum_{p,z} c_{p,z} \\ \text{s.t.} & E_q[\phi_{p,z}] \leq c_{p,z} \end{aligned}$$

where “c”s are maximums of expectation for each word tag pair by definition.

# Primitive results

- Constrained model gives clustering that's more sparse
- Clustering for a few phrases with 25 tags on BTEC ZH-EN

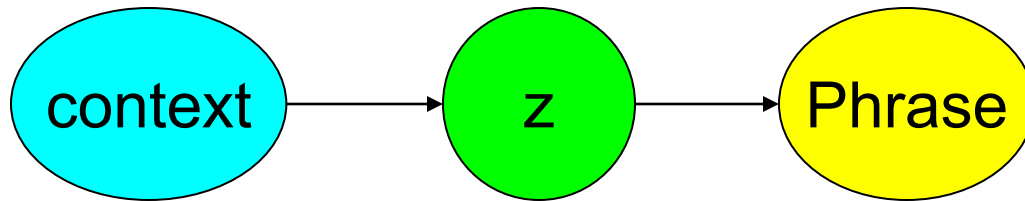
Phrase/Word	Count of the most used tag		Number of tags used	
the	1194	<b>1571</b>	11	<b>4</b>
there is	<b>53</b>	50	5	<b>4</b>
'd like	723	<b>873</b>	5	<b>2</b>



# More experiments

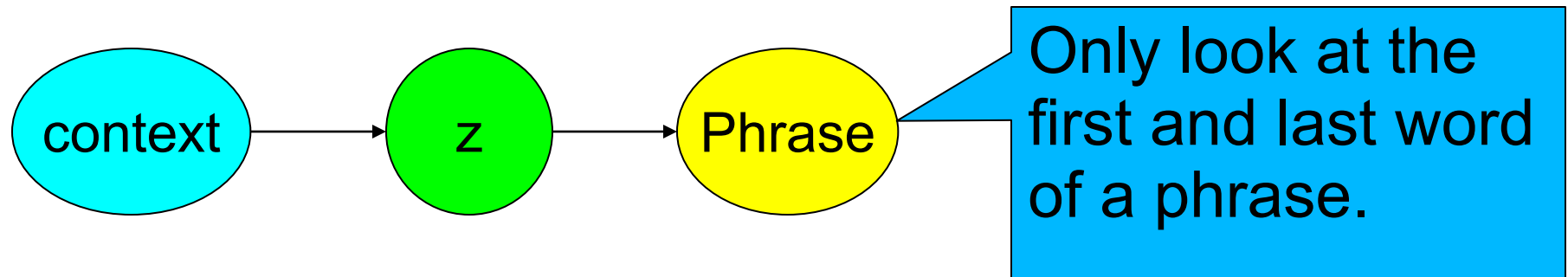
- agreement constraint: different “good” models should agree on posterior distribution
- what model to agree with: another naïve Bayes model in the reverse direction or in the other language.

# Agreement model



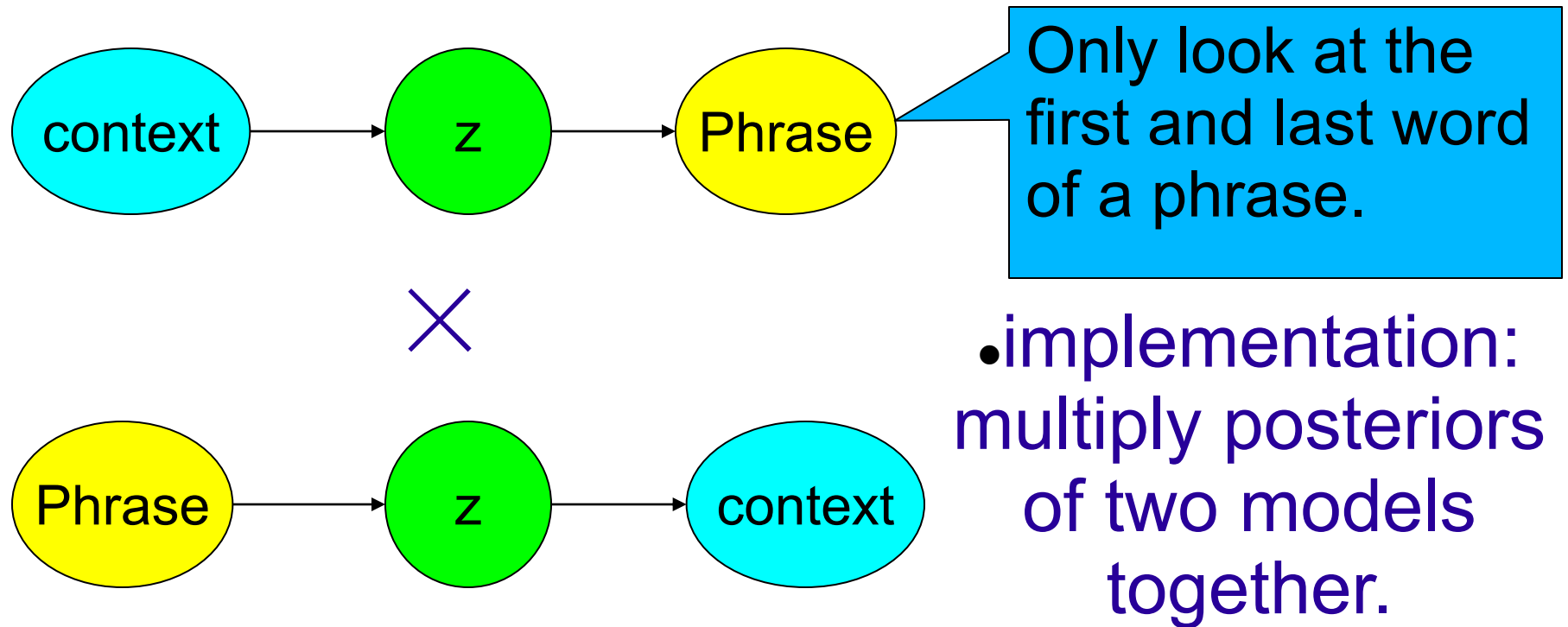
- implementation:  
multiply posteriors  
of two models  
together.

# Agreement model

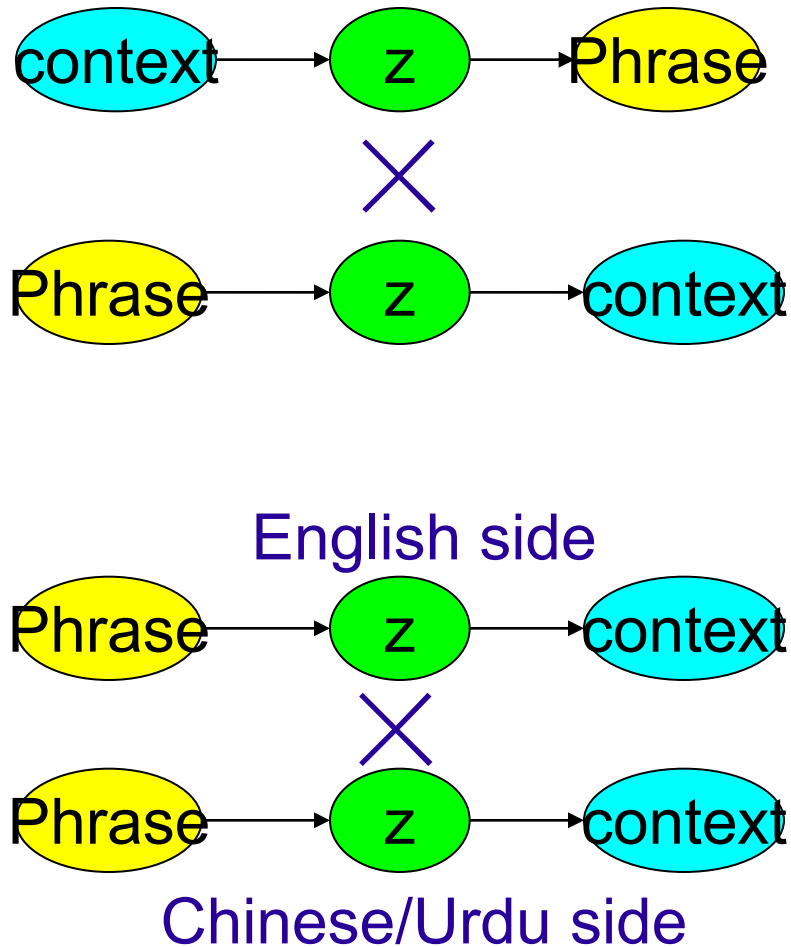


- implementation:  
multiply posteriors  
of two models  
together.

# Agreement model



# Agreement model



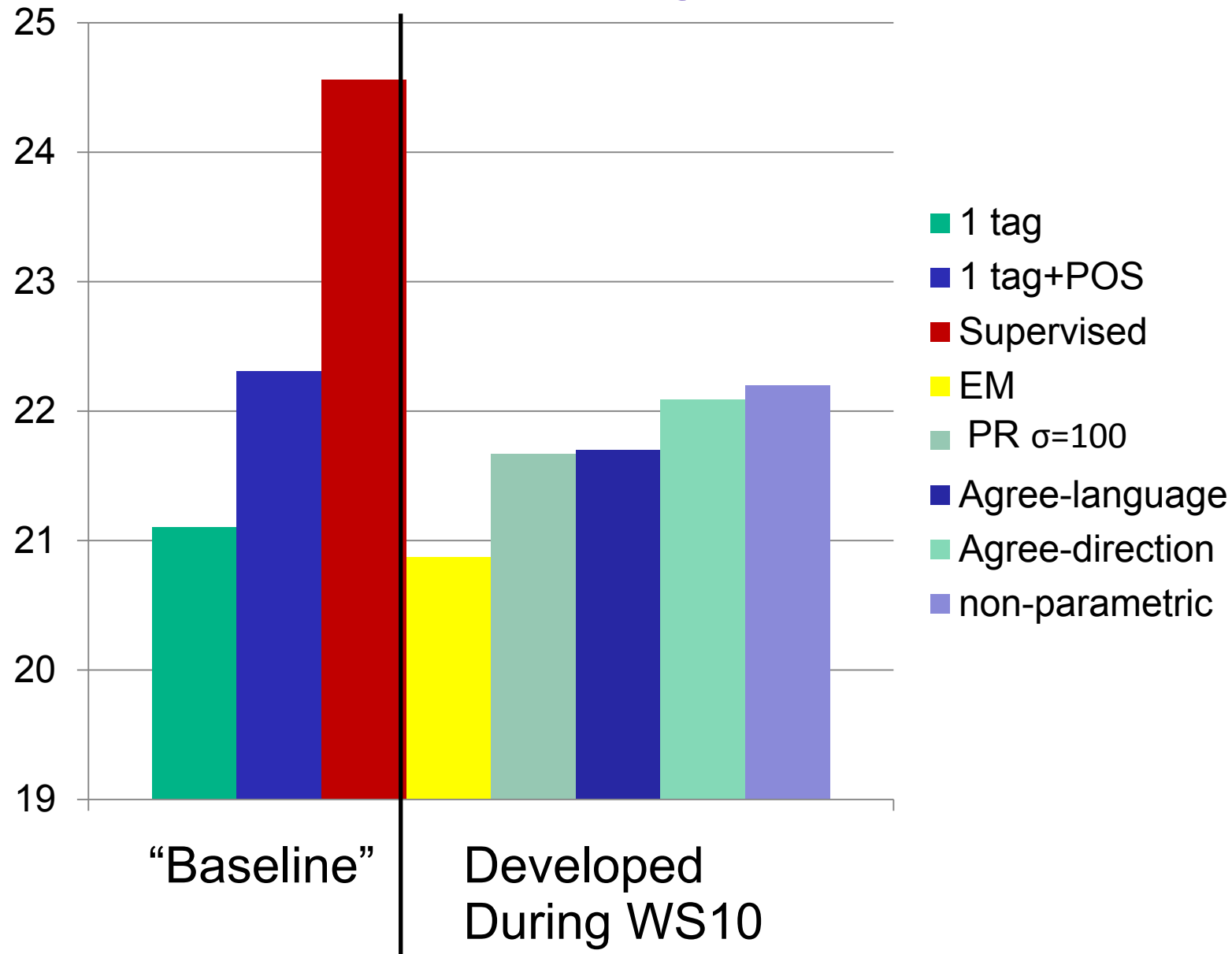
- implementation: multiply posteriors of two models together.

# Outline

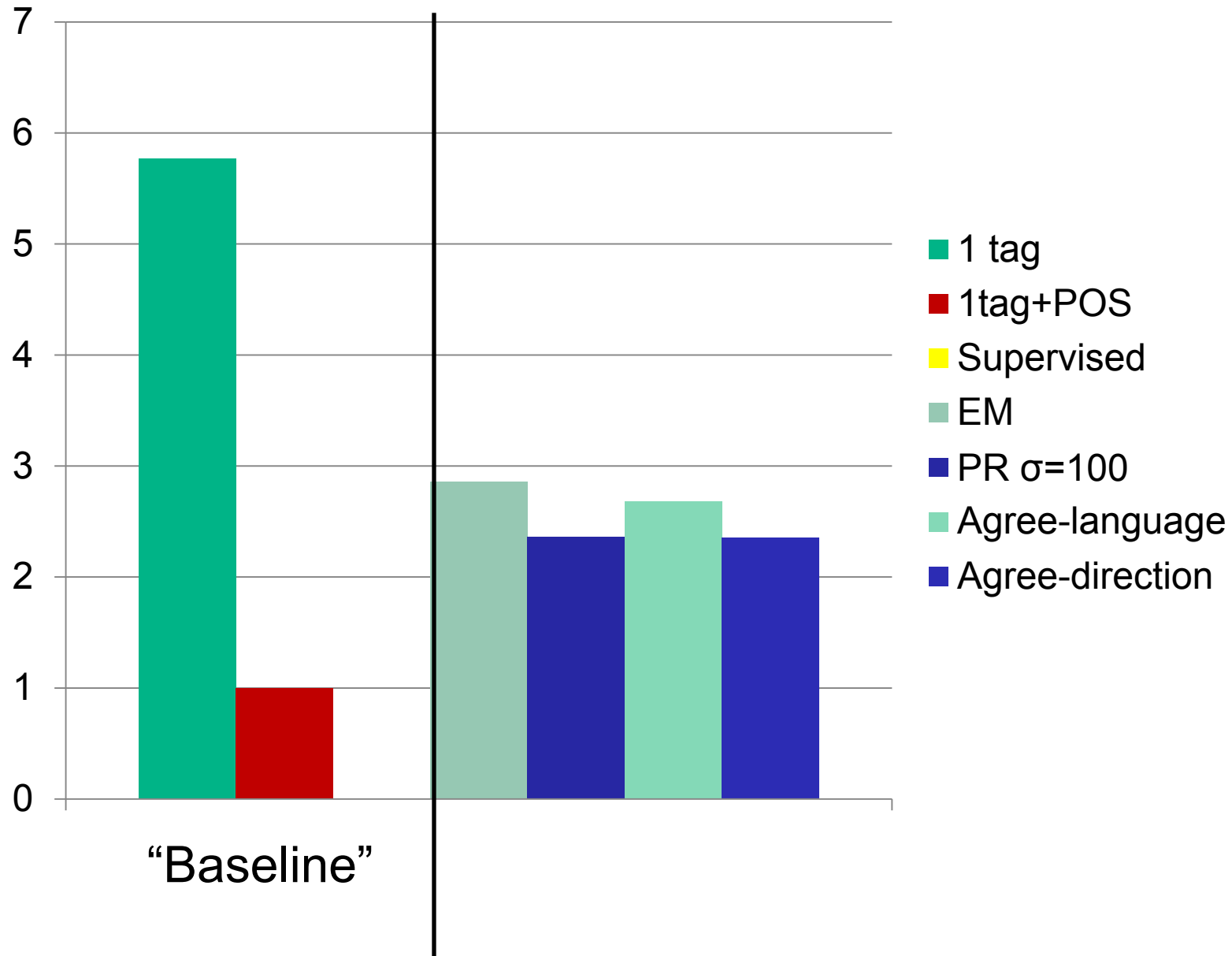
- Where do phrases come from?
- EM with posterior regularization
- **results and future experiments**

# Evaluation through the translation pipeline on Urdu-English data

BLEU score, higher is better



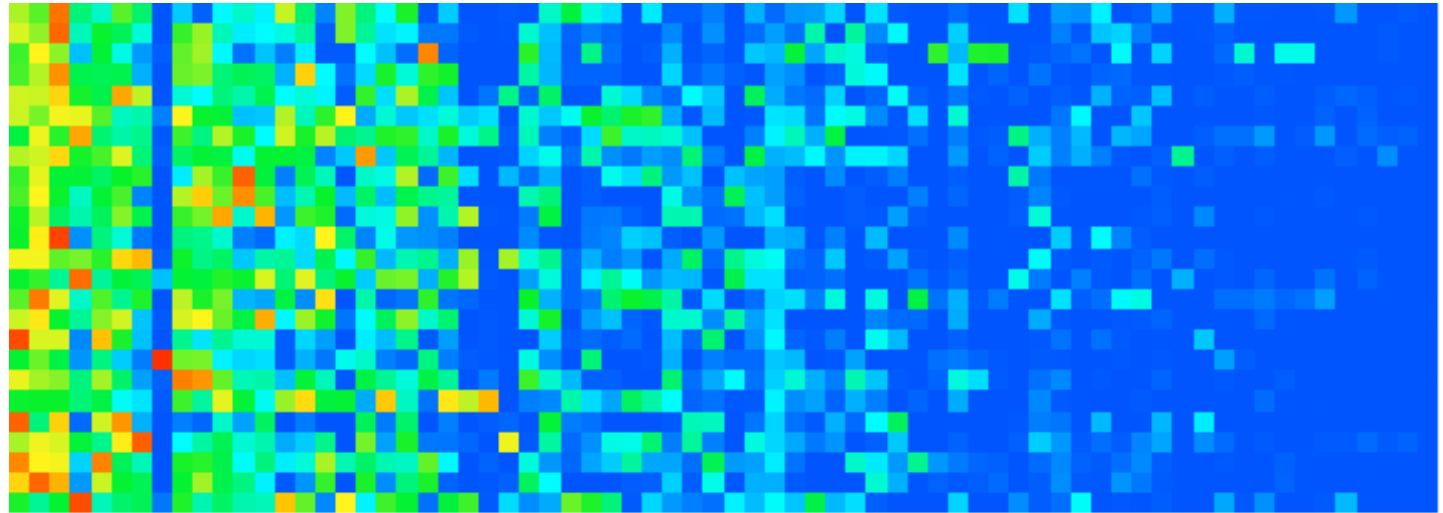
# Evaluation against supervised grammar (Conditional Entropy, lower is better)



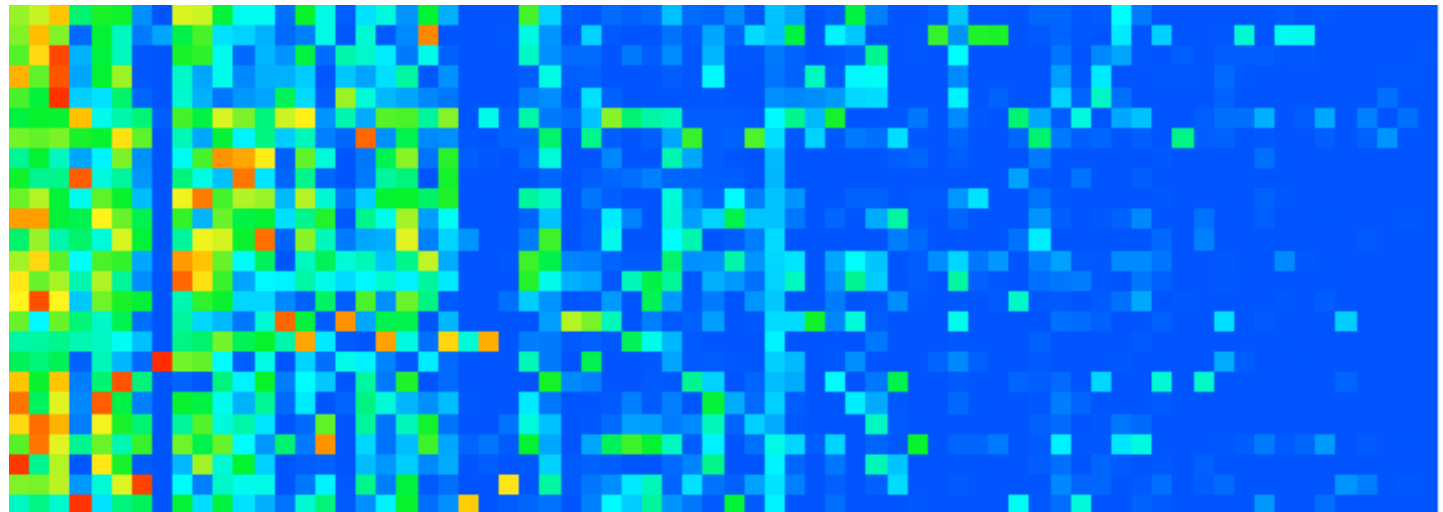


# Confusion matrix against supervised labeling

EM



Agreement  
model  
between  
languages



# Things we didn't have time to get working

- Semi-supervised training with POS tags.
- Label single-word phrases with their POS tags.

Things we didn't have time to get working

Bayesian Bayesian Bayesian

- variational Bayes inference

*Bayesian* *Bayesian* **Bayesian**

**Bayesian** **Bayesian** *Bayesian*

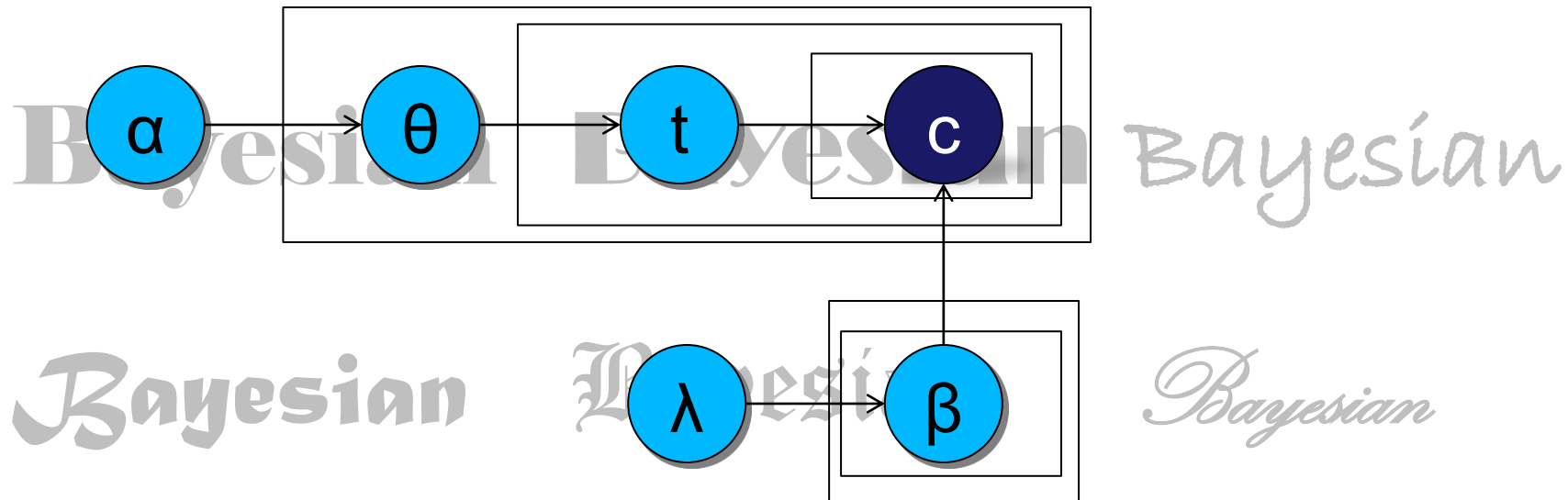
**Bayesian** *Bayesian* *Bayesian*

# Things we didn't have time to get working

Bayesian Bayesian Bayesian

- variational Bayes inference

*Bayesian Bayesian Bayesian*



# Outline

- Where do phrases come from?
- EM with posterior regularization
- results and future experiments

Thanks!

# Outline



Vlad Eidelman



Ziyuan Wang



Adam Lopez



Jon Graehl



ThuyLinh Nguyen

- 3:20pm Parametric models: posterior regularisation. Desai
- 3:35pm Training models with rich features spaces. Vlad
- 3:50pm Decoding with complex grammars. Adam
- 4:20pm Closing remarks. Phil
- 4:25pm Finish.

# Discriminative Training

Vladimir Eidelman

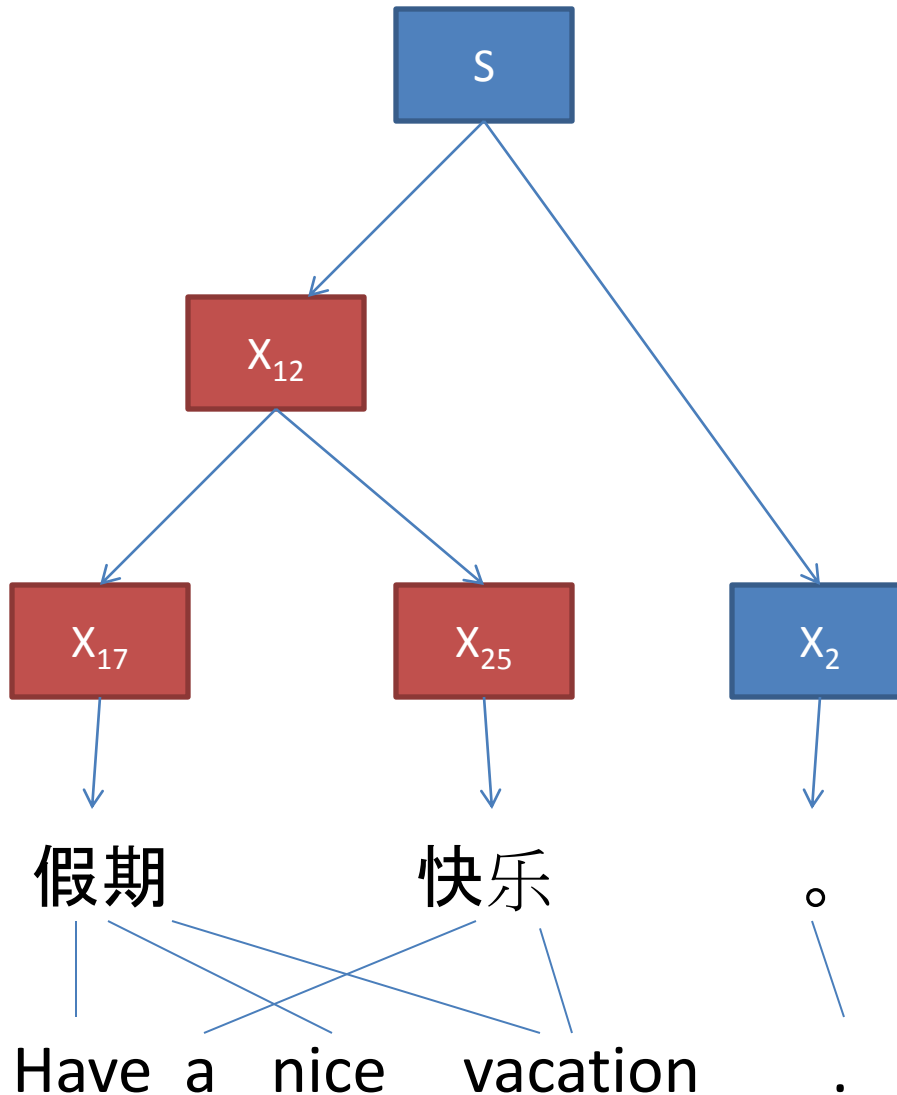
Ziyuan Wang

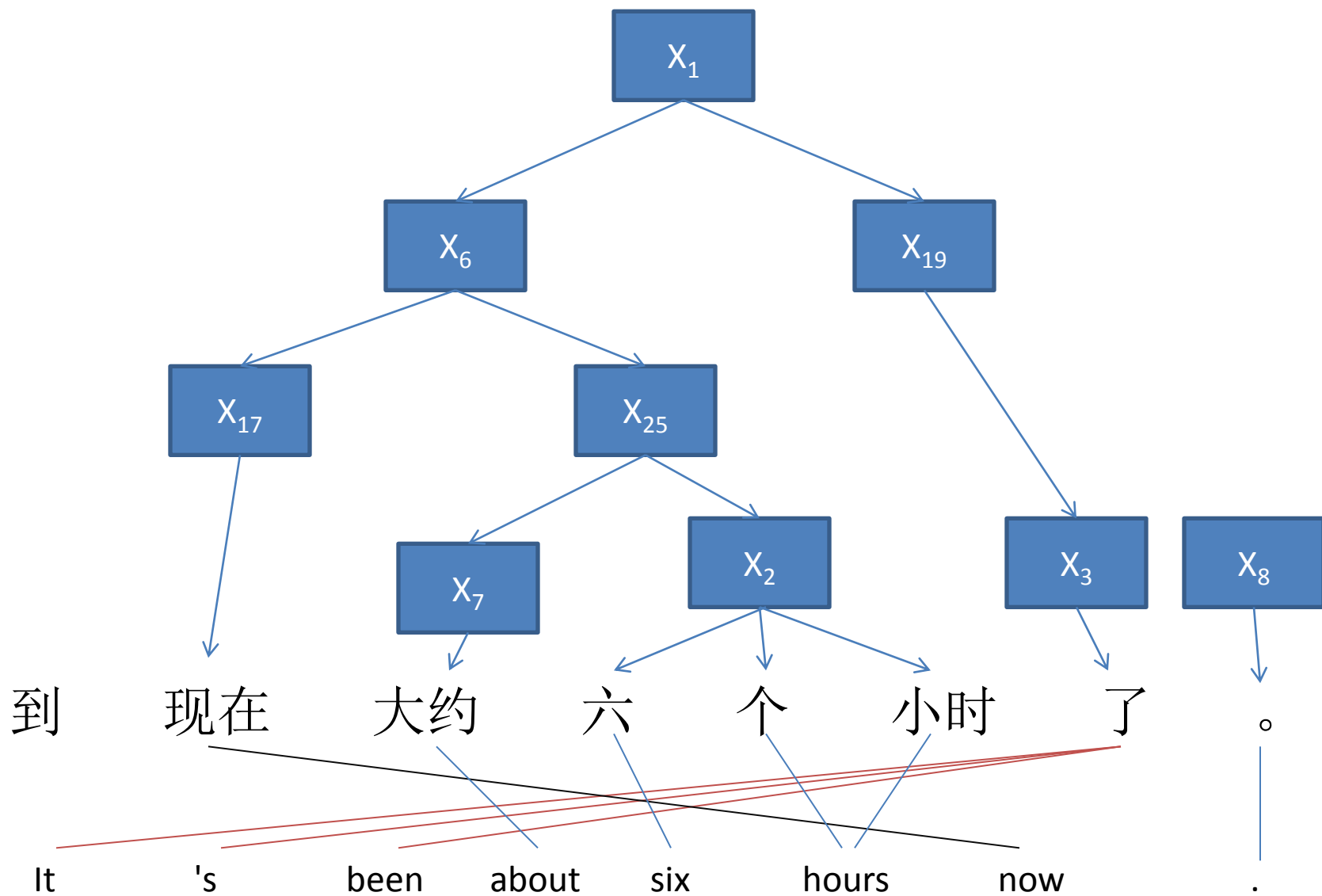
# Motivation

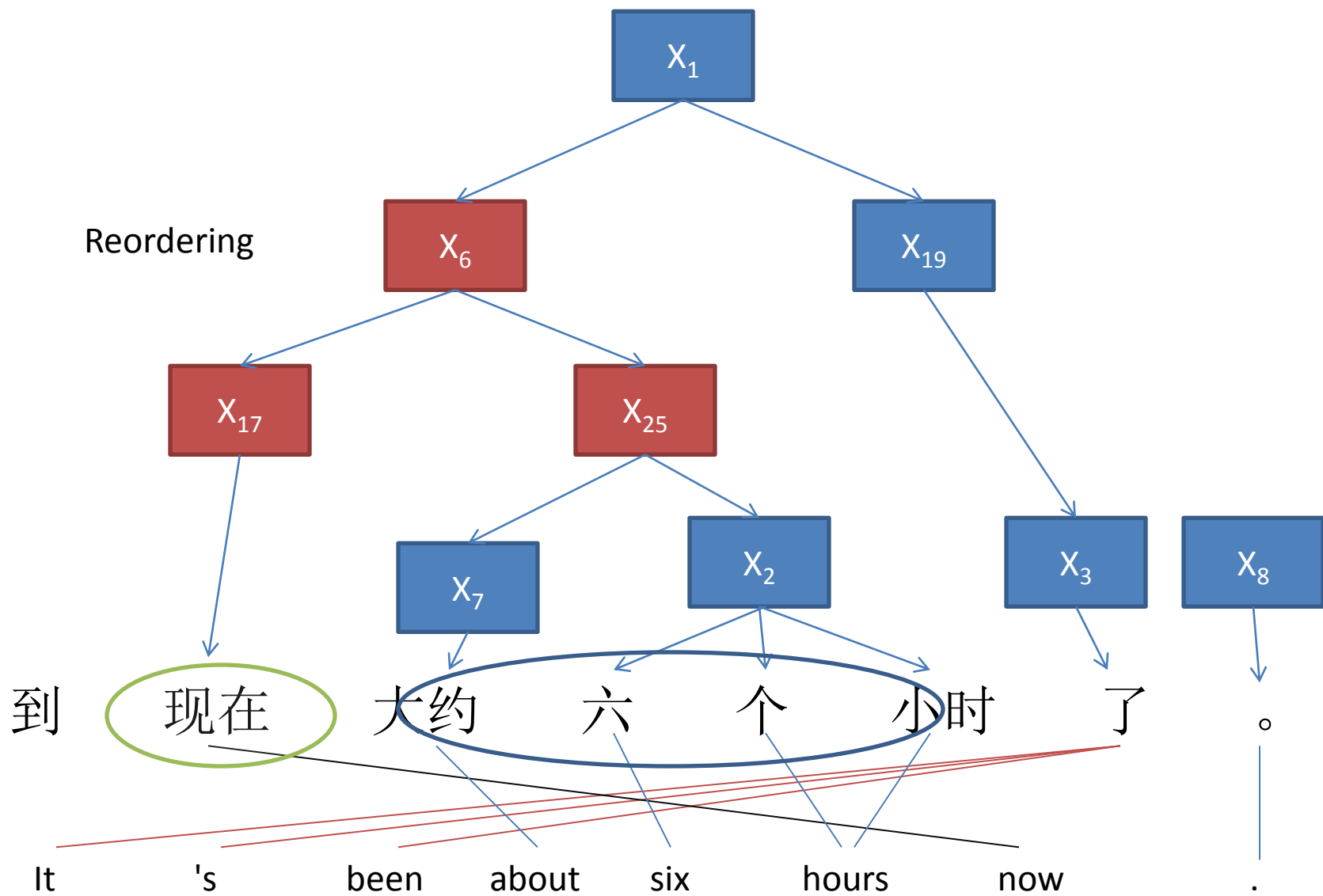
- Extract parser features from grammars
  - Source Syntax
  - Target Syntax
  - Source Context
  - Glue Features
  - OOV
  - Backoff Rule
  - Morpheme construction
- ...

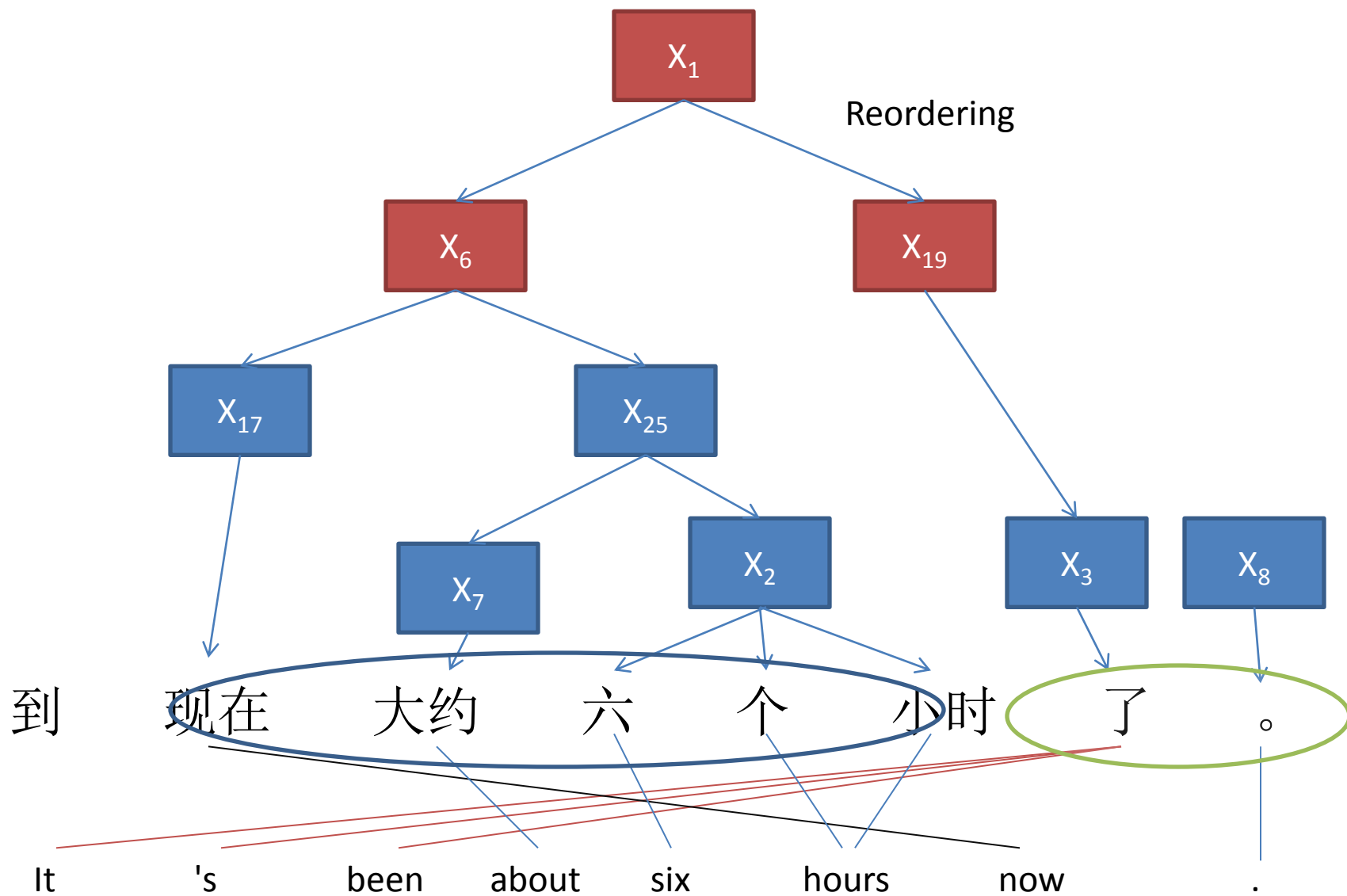


# Source Syntactic Features

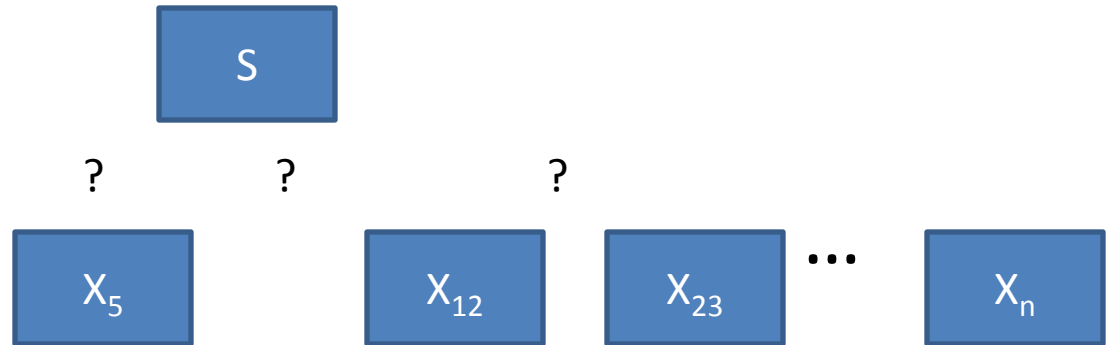




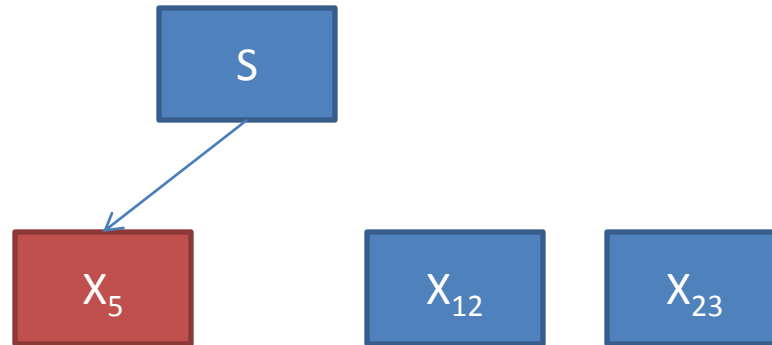




# Glue Feature

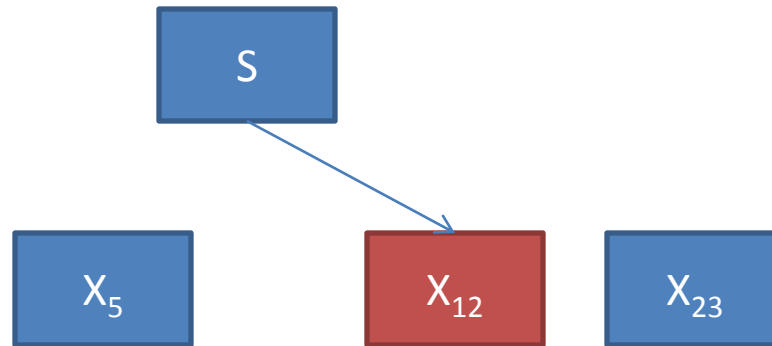


# Glue Feature



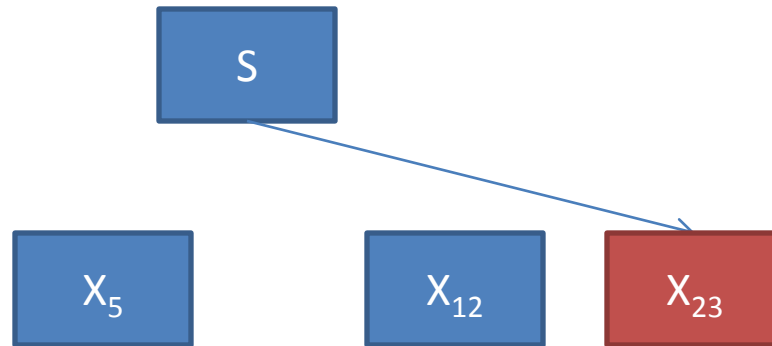
Feature:  $S_{X_5} = 1$

# Glue Feature



Feature:  $S_{X_{12}} = 1$

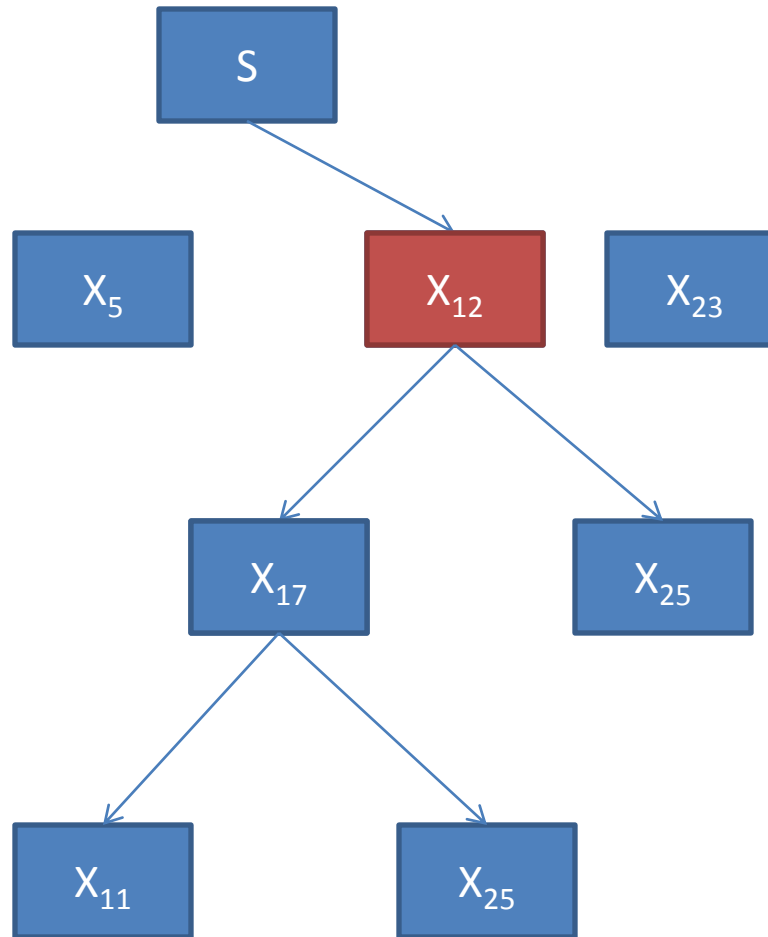
# Glue Feature



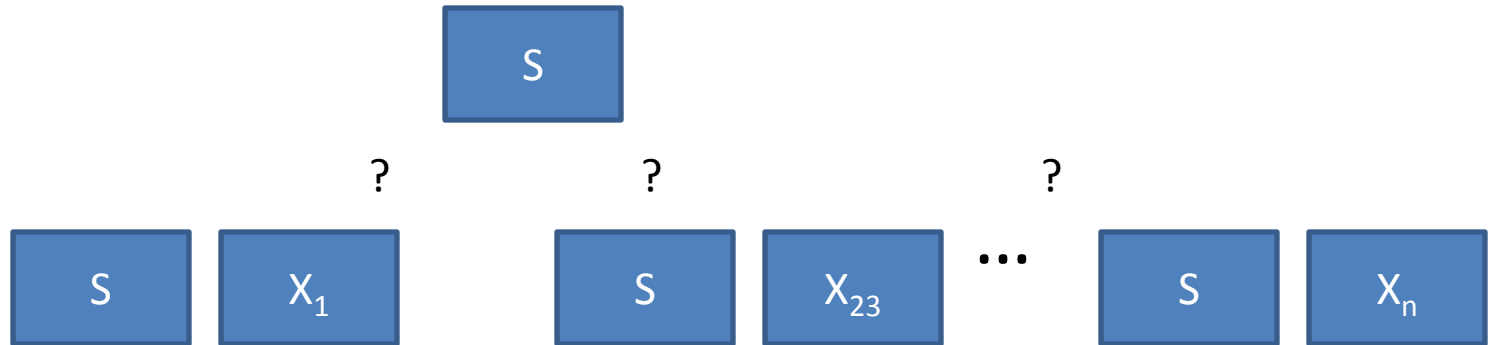
Feature:  $S_{X_{23}} = 1$



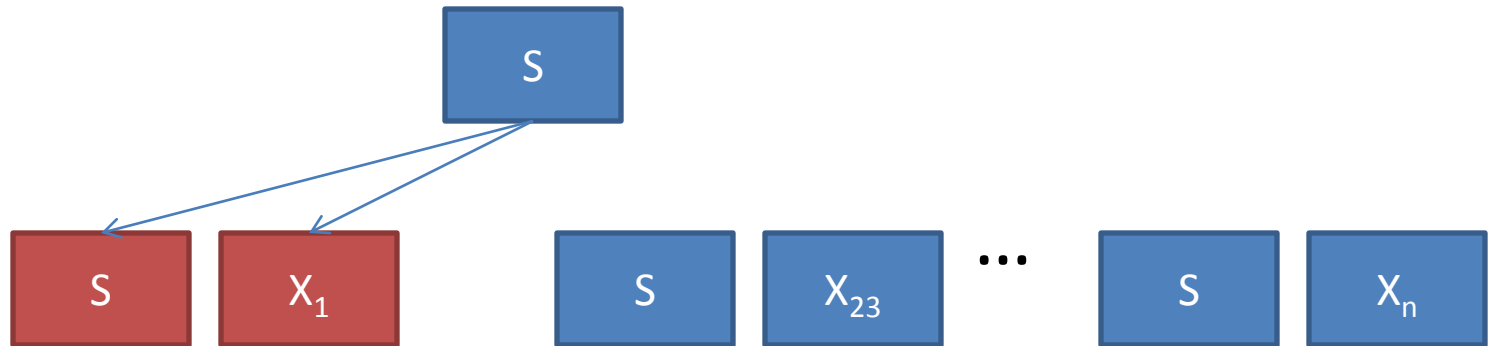
# Glue Feature



# Glue Feature

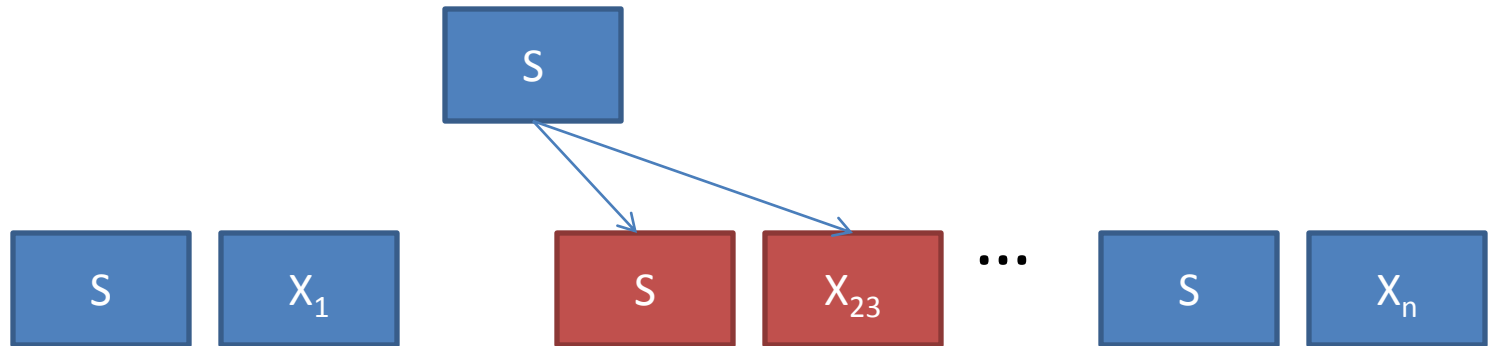


# Glue Feature



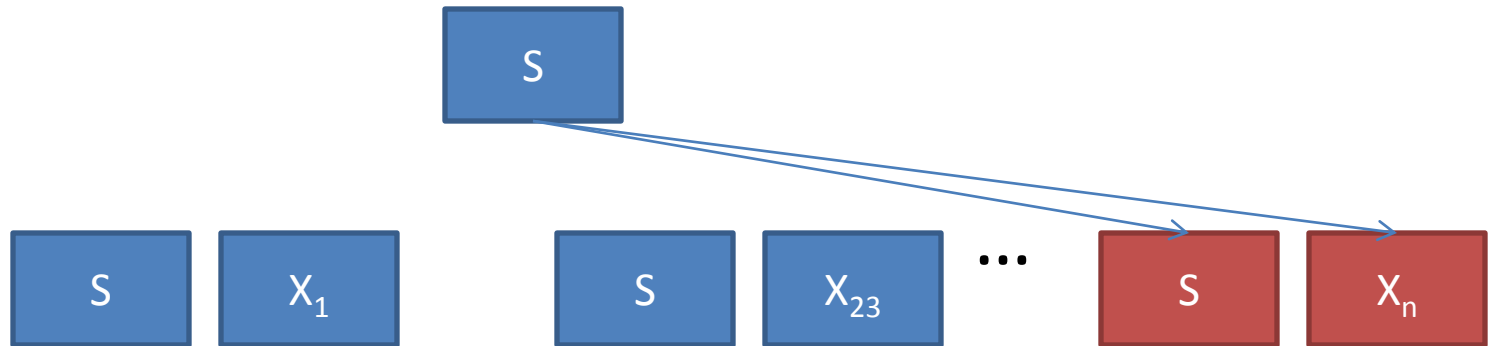
Feature:  $\text{Glue\_SX}_1 = 1$

# Glue Feature



Feature:  $\text{Glue\_SX}_{23} = 1$

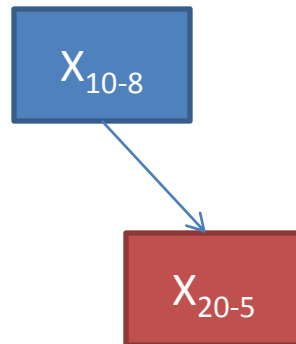
# Glue Feature



Feature:  $\text{Glue\_SX}_n = 1$

# Backoff

- In place of or combination with current dense feature



# OOV

$X_1$



Ragnarök

$X_{23}$



supercalifragilisticexpialidocious

$X_n$



6245

# OOV

$X_1$



Ragnarök

Noun?

$X_{23}$



supercalifragilisticexpialidocious

Adjective?

$X_n$



6245

Number?



# We want to...

- optimize model parameters to maximize translation quality on some metric (BLEU)
- do discriminative training so we can have features that directly help translation
- have thousands++ features

# Motivation

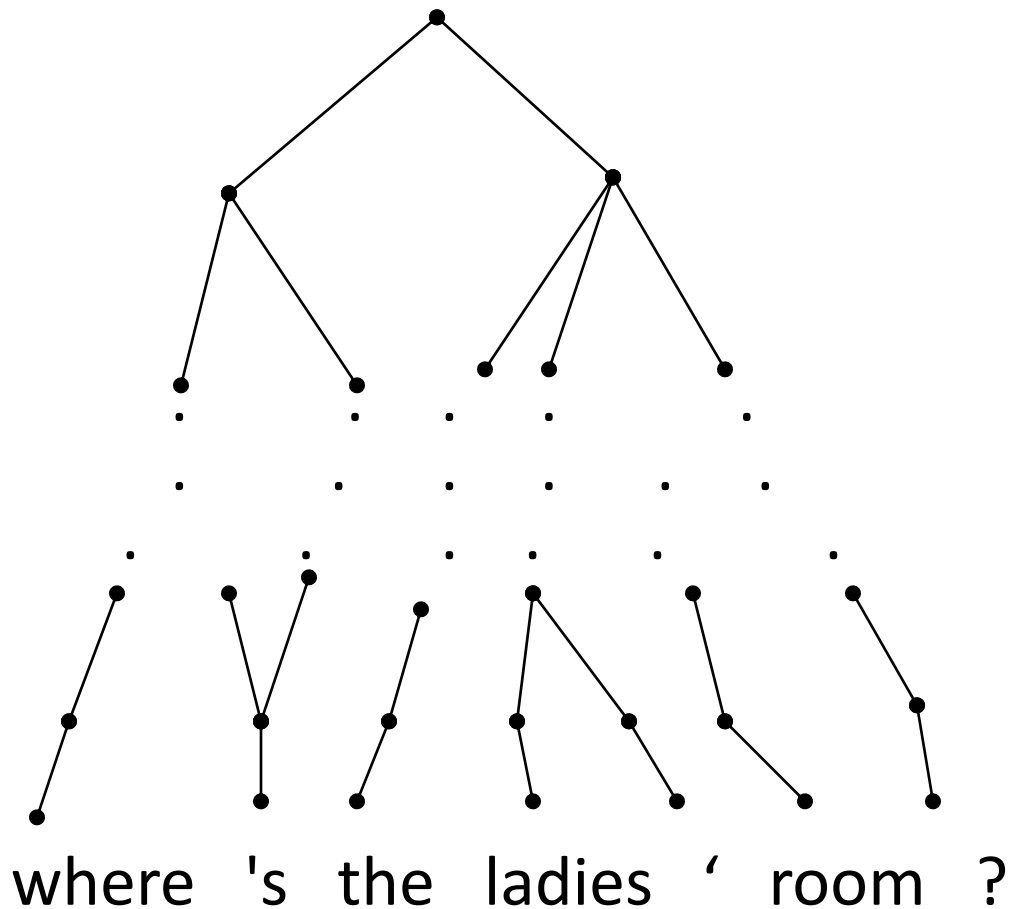
- Minimum Error Rate Training
  - Does not scale well to more than handful of features
    - $P(e)$  – Language Model
    - $P(f|e)$  – Translation Rule
    - Pass through penalty
- Alternative approaches
  - Expected BLEU training
  - MIRA
- Evaluation
  - Language invariability (parameters, iterations, etc)
  - Standardizes comparison

# Training Comparison

	MERT	MIRA	Expected BLEU
Type	1-best	Margin-based	Probabilistic
Objective	Minimize error	Minimize loss augmented score	Minimize expected error
Optimization	Line search	QP	Gradient based
Limitations	Direction of search unknown	Approximation of reference	Approximate expectation

# MIRA and Expected BLEU

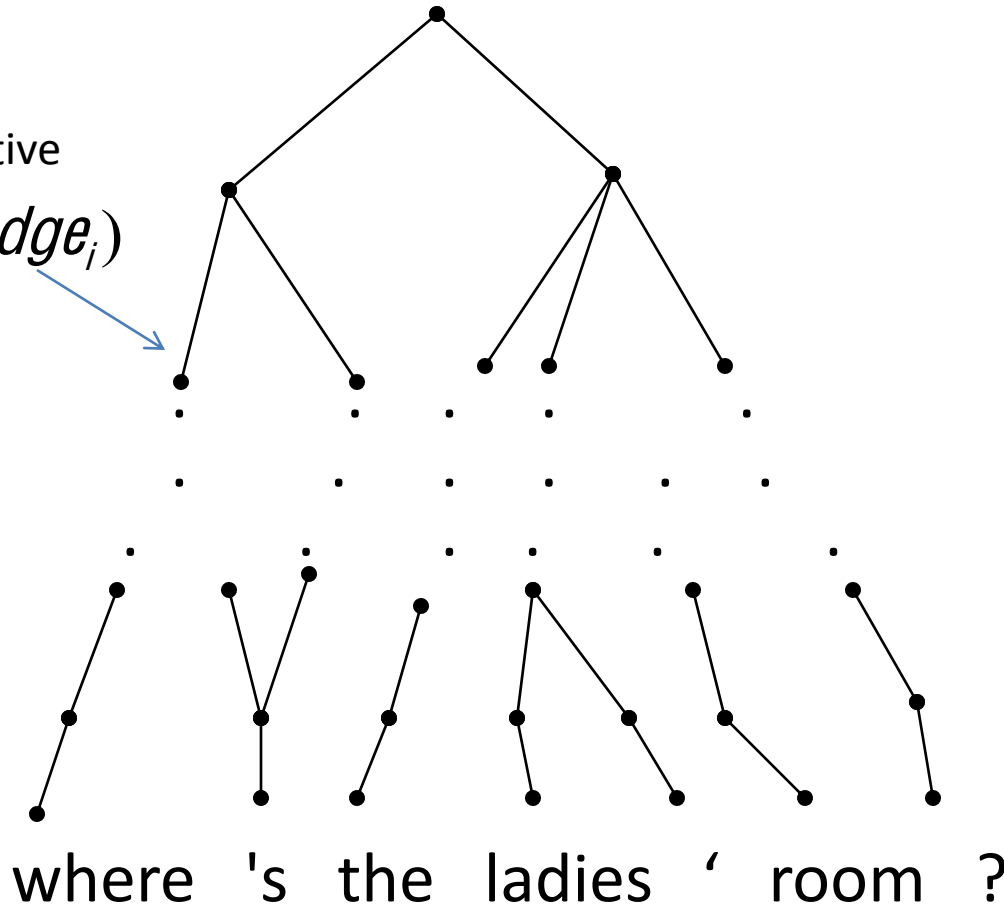
女洗手间在哪里？



# 女洗手间在哪里？

Decomposable Objective

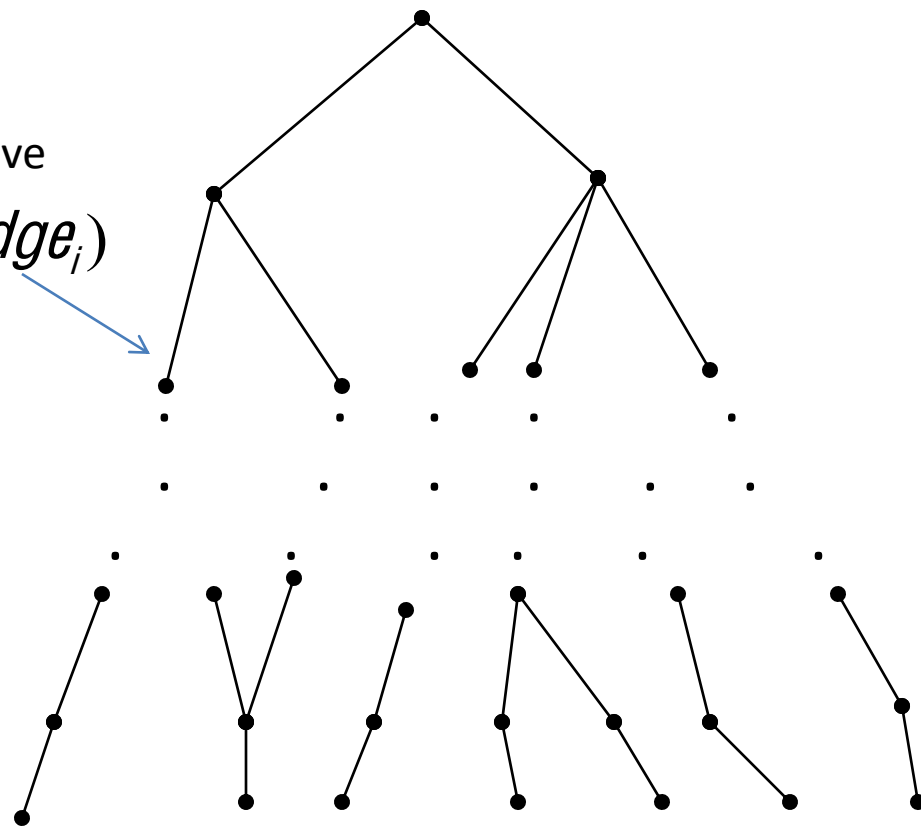
$$s(\text{edge}_i) = w f(\text{edge}_i)$$



# 女洗手间在哪里？

Decomposable Objective

$$s(\text{edge}_i) = w f(\text{edge}_i)$$



where 's the ladies ' room ?

$$\sum_{\text{edge} \in \text{derivation}} s(\text{edge}) = S(e)$$

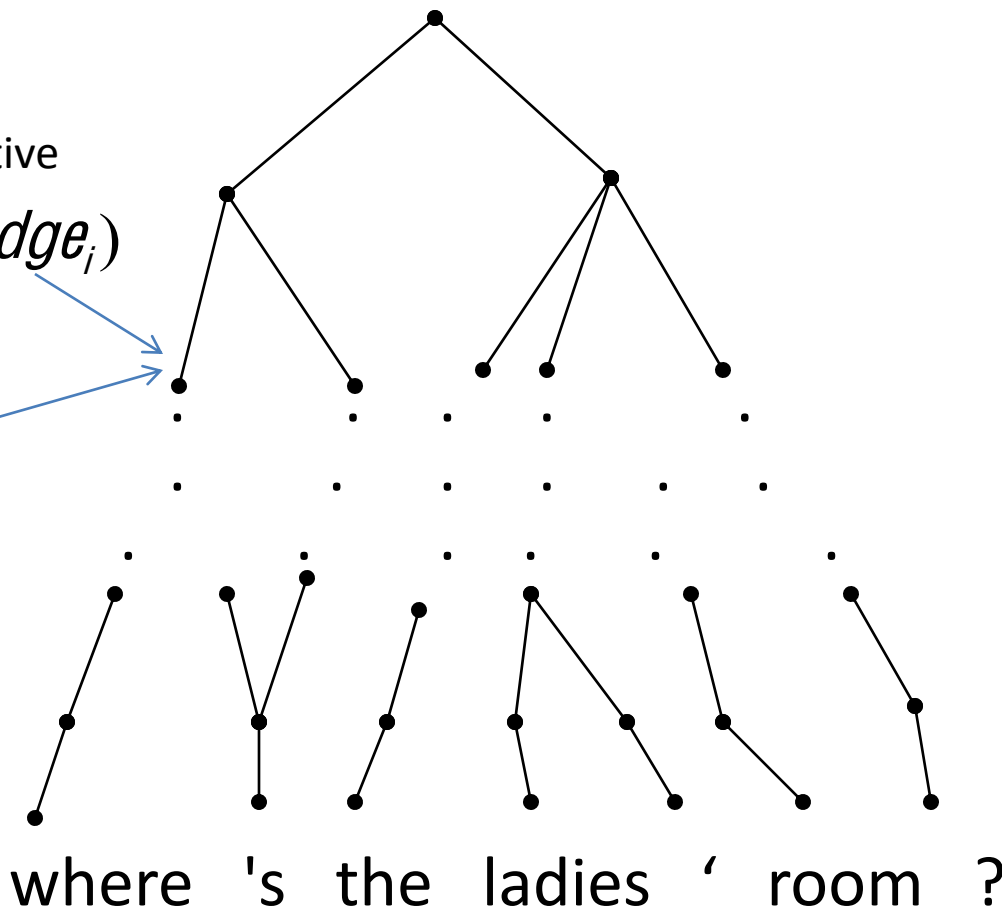
# 女洗手间在哪里？

Decomposable Objective

$$s(edge_i) = w f(edge_i)$$

Loss

*approxBLEU*



$$\sum_{edge \in derivation} s(edge) = S(e)$$

# 女洗手间在哪里？

Sentence level loss

Decomposable Objective

$$p(edge_i) = s(edge_i) = w f(edge_i)$$

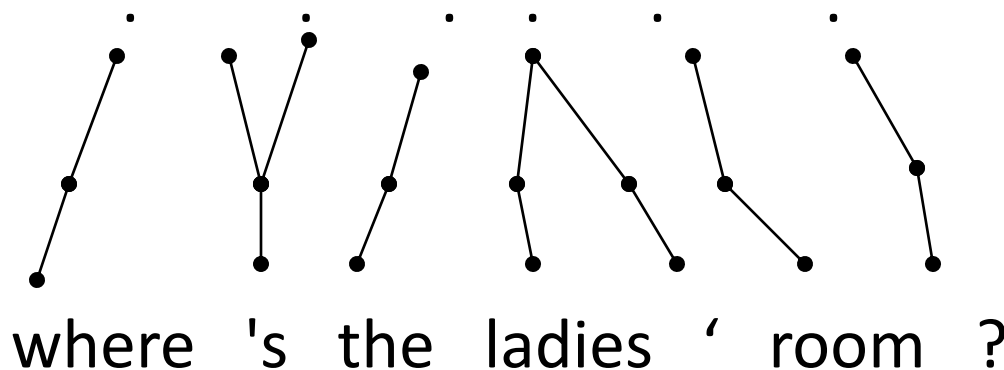
Expected BLEU  
 $E_p[f(edge_i)]$

Loss

MIRA

$$f(edge_i) = approxBLEU$$

$$f(edge_i) + p(edge_i)$$



$$\sum_{edge \in derivation} s(edge) = S(e)$$



# Margin Infused Relaxed Algorithm (MIRA)

- Online Large-Margin Learning
- Crammer and Singer (2003)
  - Multi-class classification
- Taskar (2003)
  - Extension to structured value prediction
- Watanabe (2007), Chiang (2009)
  - Application to MT

# Basic Learning Algorithm

Training data:  $D = (\mathbf{x}, \mathbf{y})$

$weight_0 = 0, total = 0, c = 0$

*for iteration*  $1 \rightarrow n$

*for*  $d = (x, y) \in D$

$weight_{c+1} = \text{update } weight_c \text{ with } d$

$total = total + weight_{c+1}$

$c = c + 1$

$weight = \frac{total}{n \times size(D)}$

# Update

$$\mathcal{S}(x, y) = w f(x, y)$$

- Learn  $w$  so that correct outputs are given higher score than incorrect ones

$$\min || w_{i+1} - w_i ||$$

- Keep the norm of the change to the weights as small as possible
- Subject to margin constraints:

$$\mathcal{S}(x, y) - \mathcal{S}(x, z) \geq \text{Loss}(y, z)$$

- Create margin between correct instance  $y$ , and incorrect instance  $z$  at least as large as the Loss of  $z$
- for all  $z$  which are possible labels of  $x$

# k-best MIRA

*Training data:  $D = (e, f)$*

*$weight_0 = 0, total = 0, c = 0$*

*for iteration  $1 \rightarrow n$*

*for  $d = (x, y) \in D$*

*Generate  $kbest(f) = \{e...e_k\}$*

*Generate margin constraints  $\forall e \in kbest(f)$*

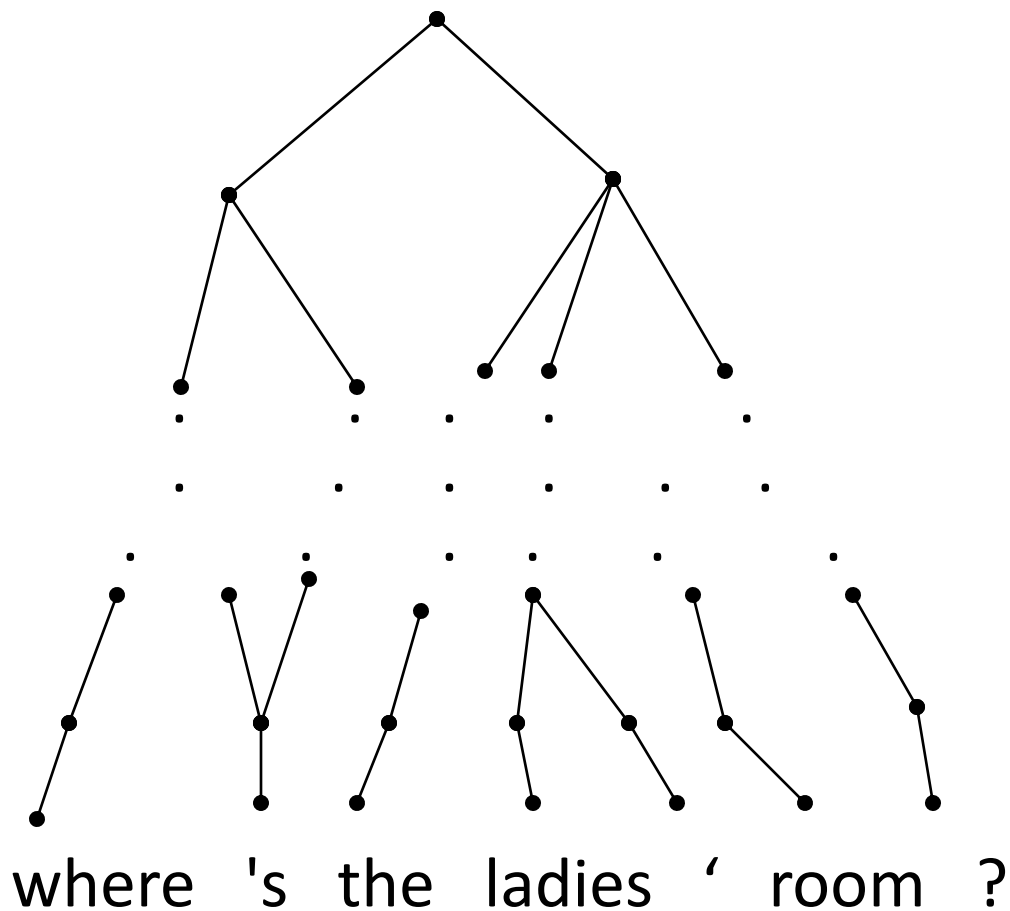
*$weight_{c+1} = \text{update } weight_c \text{ with } d$*

*$total = total + weight_{c+1}$*

*$c = c + 1$*

*$weight = \frac{total}{n \times size(D)}$*

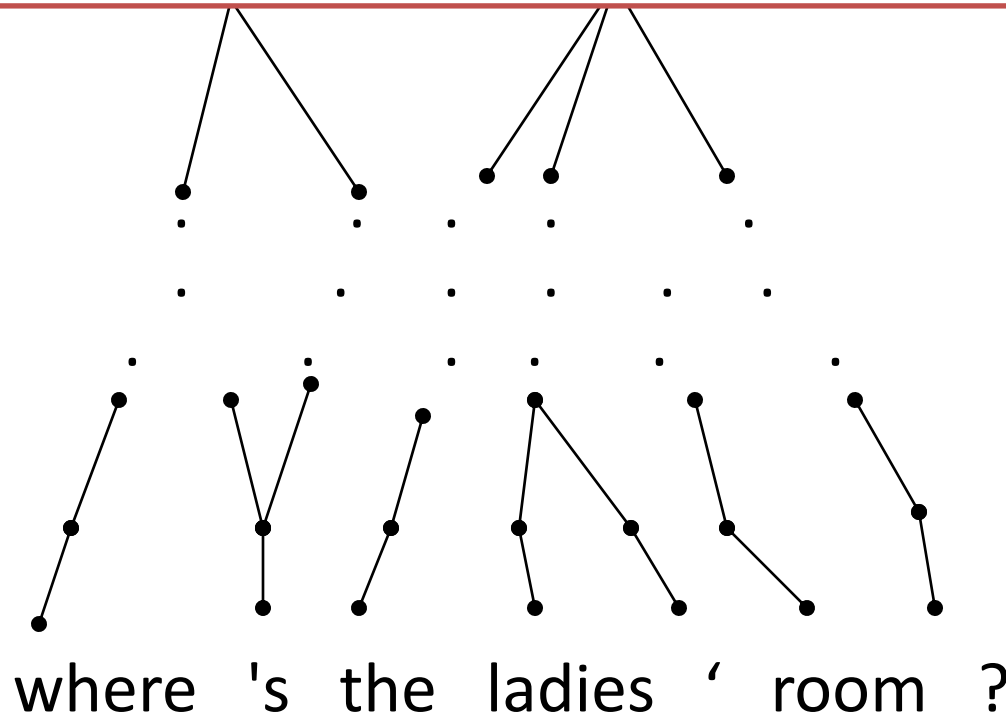
# 女洗手间在哪里？



## 女洗手间在哪里？

# Model

where 's the ladies ' room ?     LanguageModel=-6.3736...	13.661
where 's the ladies ?     LanguageModel=-5.76624...	10.8657
where 's the ladies ' ?     LanguageModel=-6.51207 ...	11.4501
where is the ladies ' room ?     LanguageModel=-7.18026 ..	14.9181
where is the ladies ?     LanguageModel=-6.5729...	11.7432



# 女士洗手间在哪里？

where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

Model+BLEU

where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| 14.9181  
where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

where 's the ladies ' room ?

# 女洗手間在哪裏？

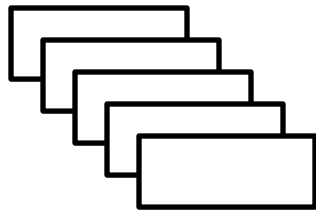
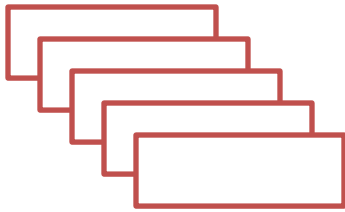
where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| 14.9181  
where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

Model - BLEU

where is the bus depot for the ladies ' room ? ||| LanguageModel=-10.7635 ||| 11.7463  
where is the bus depot for the ladies ? ||| LanguageModel=-10.1561 ||| 10.0082  
where is the bus depot for the ladies ' ? ||| LanguageModel=-10.902... ||| 10.1763  
where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
where is the bus depot for the ladies ' room . ||| LanguageModel=-11.1228 ||| 10.8613





Model	47.96
Model+ BLEU	54.08
Model- BLEU	24.10

where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

## Oracle Translation

where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| **14.9181**  
where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

where is the bus depot for the ladies ' room ? ||| LanguageModel=-10.7635 ||| 11.7463  
where is the bus depot for the ladies ? ||| LanguageModel=-10.1561 ||| 10.0082  
where is the bus depot for the ladies ' ? ||| LanguageModel=-10.902... ||| 10.1763  
where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
where is the bus depot for the ladies ' room . ||| LanguageModel=-11.1228 ||| 10.8613

where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
 where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
 where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
 where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
 where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

$\Delta f(e)$

Loss

where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| 14.9181  
 where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
 where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

where is the bus depot for the ladies ' room ? ||| LanguageModel=-10.7635 ||| 11.7463  
 where is the bus depot for the ladies ? ||| LanguageModel=-10.1561 ||| 10.0082  
 where is the bus depot for the ladies ' ? ||| LanguageModel=-10.902... ||| 10.1763  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the bus depot for the ladies ' room . ||| LanguageModel=-11.1228 ||| 10.8613

where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
 where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
 where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
 where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
 where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

$\Delta f(e)$

Loss

where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| 14.9181  
 where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
 where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

where is the bus depot for the ladies ' room ? ||| LanguageModel=-10.7635 ||| 11.7463  
 where is the bus depot for the ladies ? ||| LanguageModel=-10.1561 ||| 10.0082  
 where is the bus depot for the ladies ' ? ||| LanguageModel=-10.902... ||| 10.1763  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the bus depot for the ladies ' room . ||| LanguageModel=-11.1228 ||| 10.8613

where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
 where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
 where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
 where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
 where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

$\Delta f(e)$

Loss

where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| 14.9181  
 where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
 where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

where is the bus depot for the ladies ' room ? ||| LanguageModel=-10.7635 ||| 11.7463  
 where is the bus depot for the ladies ? ||| LanguageModel=-10.1561 ||| 10.0082  
 where is the bus depot for the ladies ' ? ||| LanguageModel=-10.902... ||| 10.1763  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the bus depot for the ladies ' room . ||| LanguageModel=-11.1228 ||| 10.8613

where 's the ladies ' room ? ||| LanguageModel=-6.3736.... ||| 13.661  
 where 's the ladies ? ||| LanguageModel=-5.76624... ||| 10.8657  
 where 's the ladies ' ? ||| LanguageModel=-6.51207 ... ||| 11.4501  
 where is the ladies ' room ? ||| LanguageModel=-7.18026 .. ||| 14.9181  
 where is the ladies ? ||| LanguageModel=-6.5729... ||| 11.7432

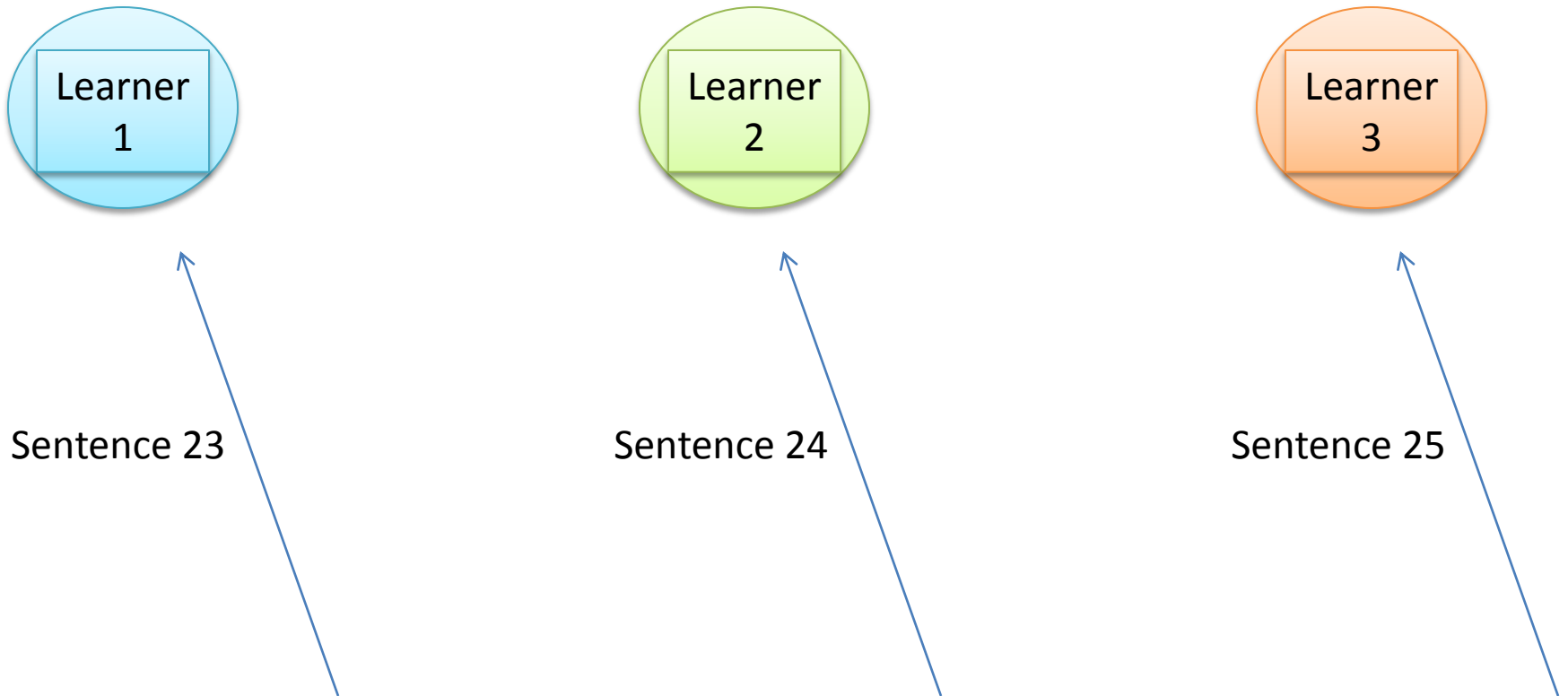
$\Delta f(e)$

Loss

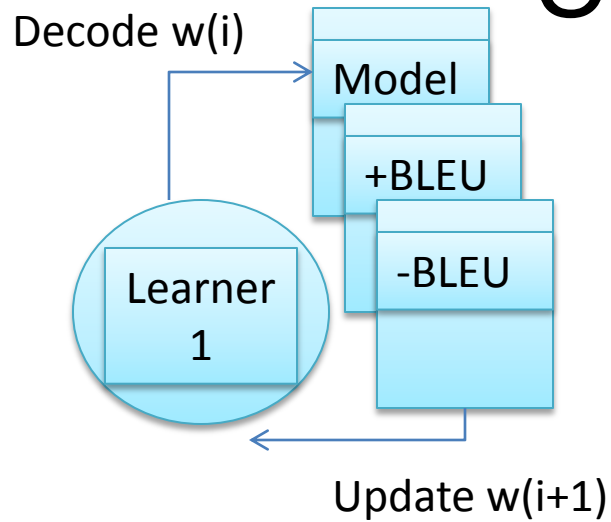
where is the ladies ' room ? ||| LanguageModel=-7.18026 ||| 14.9181  
 where is the ladies ' ? ||| LanguageModel=-7.31873 ||| 12.8778  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the ladies ? ||| LanguageModel=-6.5729 ||| 11.7432  
 where 's the ladies ? ||| LanguageModel=-5.76624 ||| 10.8657

where is the bus depot for the ladies ' room ? ||| LanguageModel=-10.7635 ||| 11.7463  
 where is the bus depot for the ladies ? ||| LanguageModel=-10.1561 ||| 10.0082  
 where is the bus depot for the ladies ' ? ||| LanguageModel=-10.902... ||| 10.1763  
 where 's the ladies ' room ? ||| LanguageModel=-6.3736 ||| 13.661  
 where is the bus depot for the ladies ' room . ||| LanguageModel=-11.1228 ||| 10.8613

# Online Updating

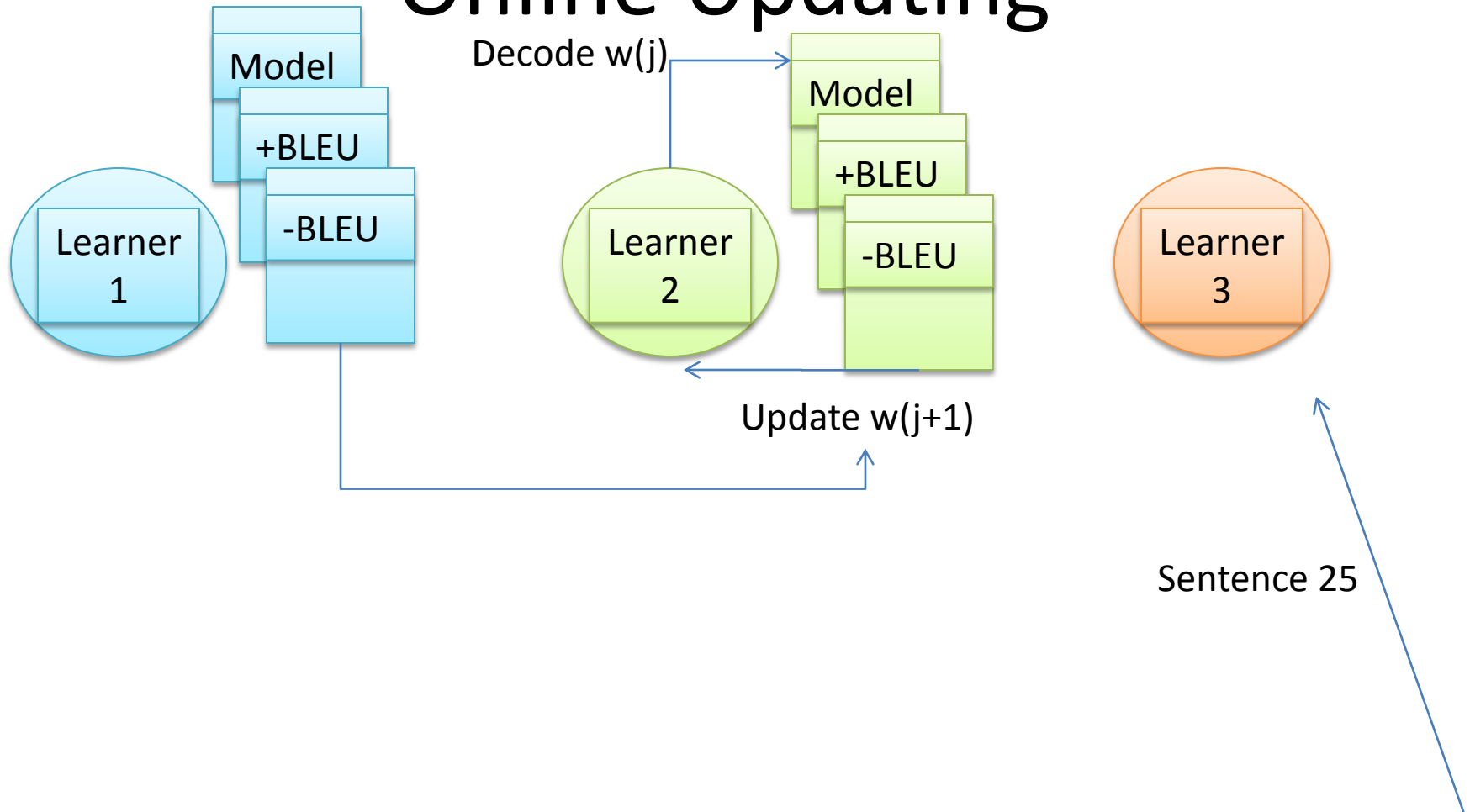


# Online Updating

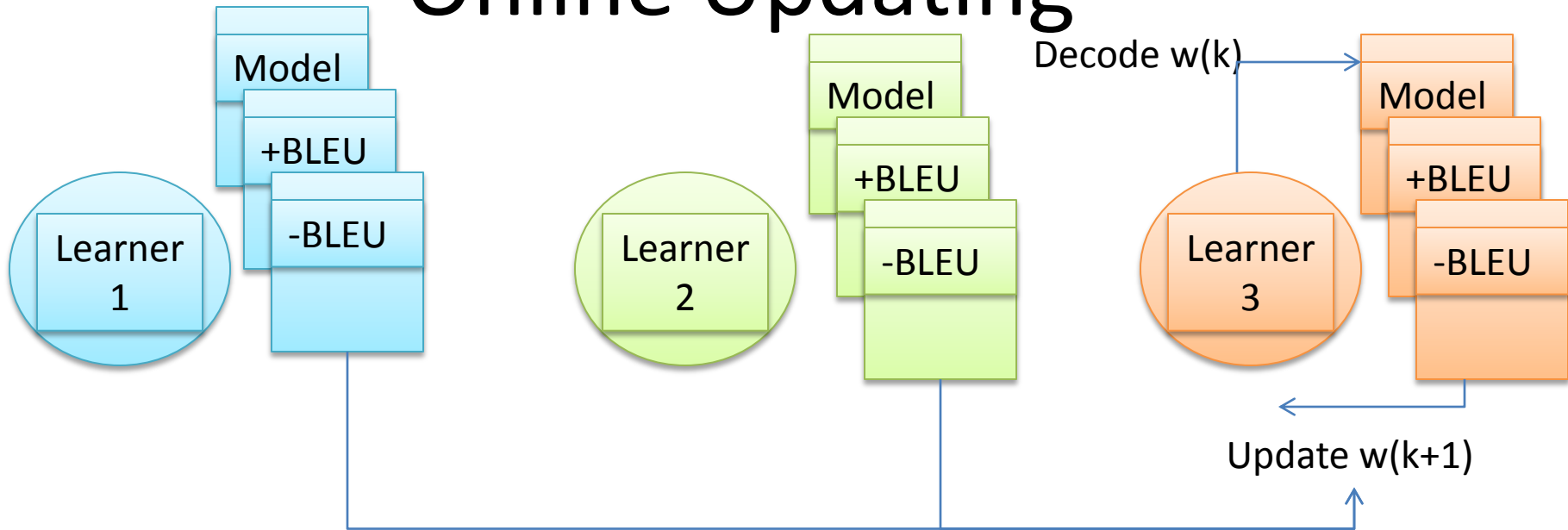




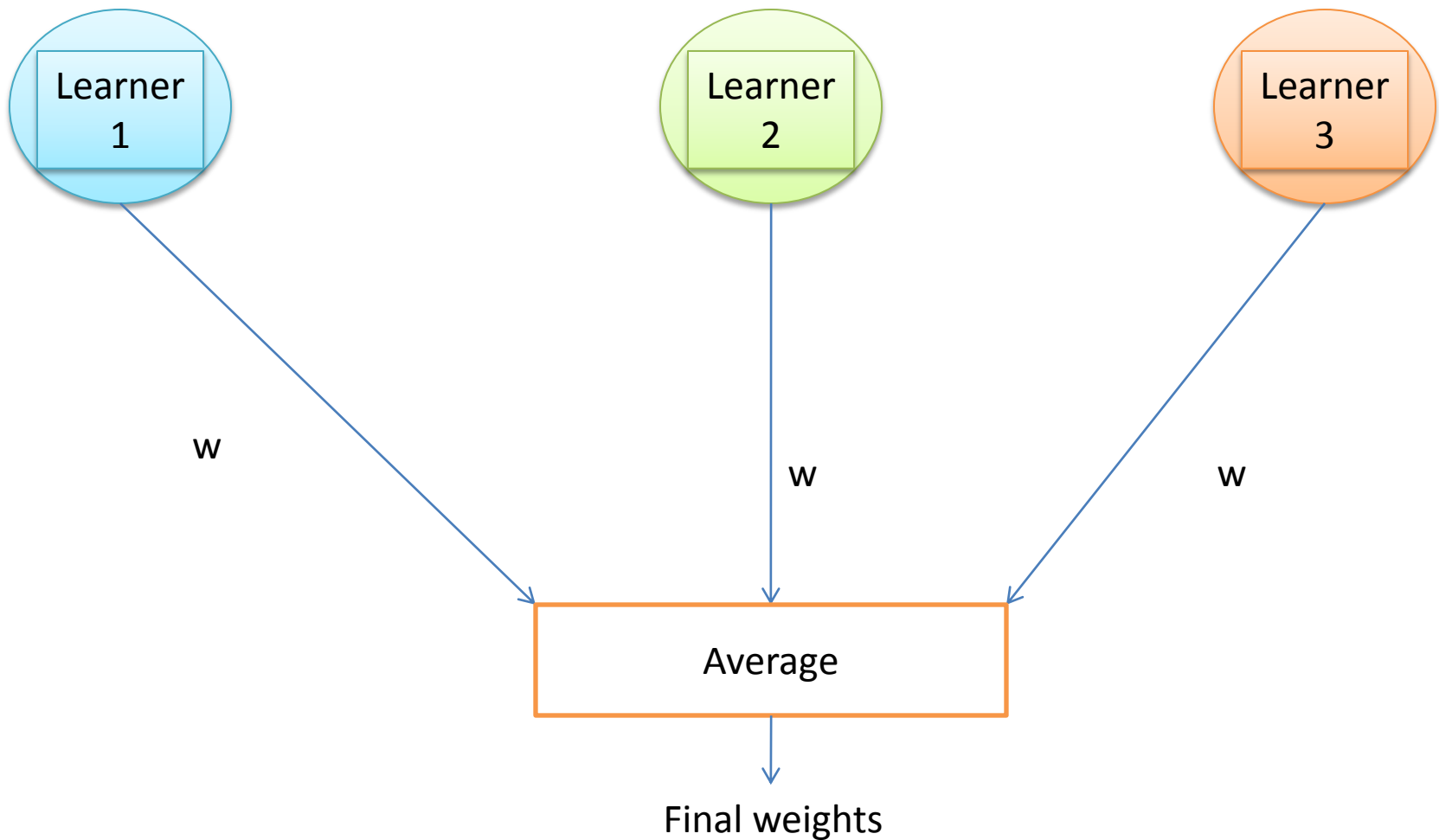
# Online Updating



# Online Updating



# Online Updating



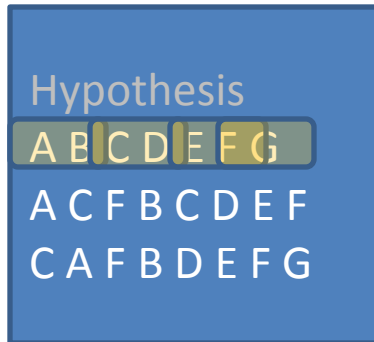
# Expected BLEU

- BLEU is just a geometric mean of ngram precisions

Hypothesis  
A B C D E F G  
A C F B C D E F  
C A F B D E F G

Reference  
A B C D E F G  
A C B D E G  
A C D E F G

# Expected BLEU



- Brevity penalty when hypothesis < reference

# Expected BLEU

- Usually perform 1-best BLEU
  - $\text{argmax}$
- Expected BLEU replaces it  $\text{argmax}$  with sum
  - Function becomes continuous w.r.t weights
- Use approximate brevity penalty
  - Replace  $\text{argmax}$  with sum
- Differentiable

# Expected BLEU

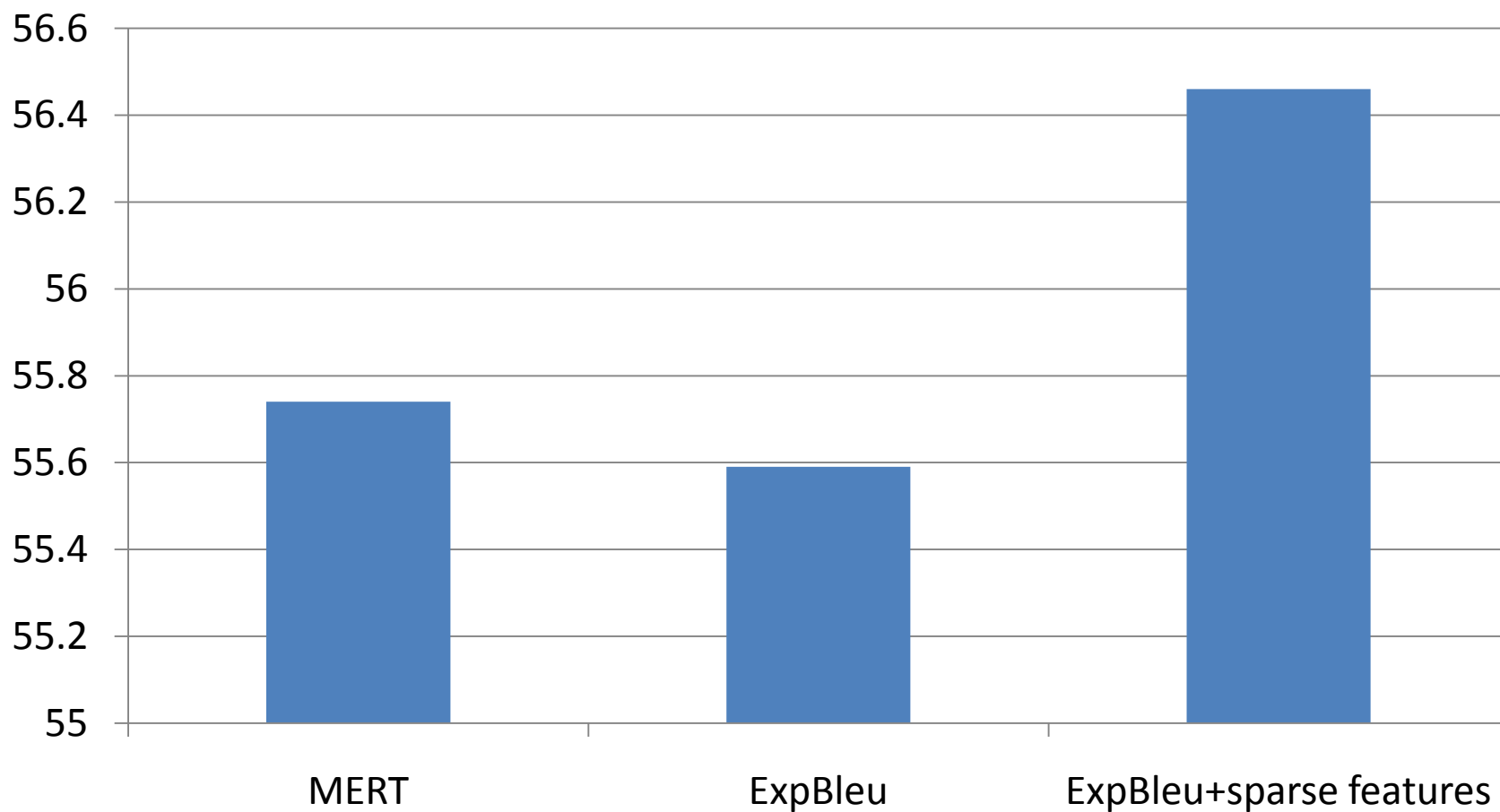
- Procedure:
  - 1) LBFGS tuning for several iterations until convergence on the hypergraph
  - 2) Re-decode the source data, generate updated hypergraph
  - 3) Repeat

# Preliminary Experiments

- Compare Expected BLEU with MIRA with equivalent grammar on same test set
- Incorporate fine-grained features
  - Source Syntax
  - Target Syntax
  - Source Context
  - Glue Features
  - OOV
  - Backoff Rule



# Preliminary Results



# Coming Soon...

## Decoding with Complex Grammars

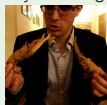
# Outline



Vlad Eidelman



Ziyuan Wang



Adam Lopez



Jon Graehl



ThuyLinh Nguyen

- 3:20pm Parametric models: posterior regularisation. Desai
- 3:35pm Training models with rich features spaces. Vlad
- 3:50pm Decoding with complex grammars. Adam
- 4:20pm Closing remarks. Phil
- 4:25pm Finish.

# Efficient Decoding for Synchronous Context-Free Grammars

Adam Lopez (Edinburgh)

Jonathan Graehl (ISI)

Chris Dyer (CMU)

with thanks to:

the whole workshop team,

Juri Ganitkevitch (JHU) & Jonny Weese (JHU)

# Efficient Decoding for Synchronous Context-Free Grammars

Adam Lopez (Edinburgh)

Jonathan Graehl (ISI)

Chris Dyer (CMU)

with thanks to:

the whole workshop team,

Juri Ganitkevitch (JHU) & Jonny Weese (JHU)



# The Story So Far

# The Story So Far

- Induce a grammar.

# The Story So Far

- Induce a grammar.
- Tune some model parameters.



# The Story So Far

- Induce a grammar.
- Tune some model parameters.
- Get a BLEU score.

# The Story So Far

- Induce a grammar.
- Tune some model parameters.
- *Decode a test set.*
- Get a BLEU score.

# The Story So Far

- Induce a grammar.
- Tune some model parameters.
  - *Decode a tuning set.*
  - *Decode a test set.*
- Get a BLEU score.

# The Story So Far

- Induce a grammar.
  - *Decode the training data.*
- Tune some model parameters.
  - *Decode a tuning set.*
  - *Decode a test set.*
- Get a BLEU score.

# The Price of Performance

- 1 Category (baseline): **20.8**
- 25 Categories: **21.7**

# The Price of Performance

- 1 Category (baseline): 3.0 sec / sentence
- 25 Categories: 52 sec / sentence

# Some Questions

# Some Questions

- Why is it so slow?



# Some Questions

- Why is it so slow?
- How can we speed it up?

# Some Questions

- Why is it so slow?
- How can we speed it up?
- What's the big idea?

# Context-Free Grammar

$$X^1 \rightarrow \textit{dianzi shang}$$

$$X^2 \rightarrow \textit{dianzi shang}$$

$$X^3 \rightarrow \textit{mao}$$

$$X^4 \rightarrow X^1 \textit{ de } X^3$$

$$X^4 \rightarrow X^2 \textit{ de } X^3$$

$$X^5 \rightarrow X^1 \textit{ de } X^3$$

$$S \rightarrow X^4$$

$$S \rightarrow X^5$$

# Context-Free Grammar

S

$$X^1 \rightarrow \textit{dianzi shang}$$

$$X^2 \rightarrow \textit{dianzi shang}$$

$$X^3 \rightarrow \textit{mao}$$

$$X^4 \rightarrow X^1 \textit{ de } X^3$$

$$X^4 \rightarrow X^2 \textit{ de } X^3$$

$$X^5 \rightarrow X^1 \textit{ de } X^3$$

$$S \rightarrow X^4$$

$$S \rightarrow X^5$$

# Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang}$

$X^2 \rightarrow \textit{dianzi shang}$

$X^3 \rightarrow \textit{mao}$

$X^4 \rightarrow X^1 \textit{ de } X^3$

$X^4 \rightarrow X^2 \textit{ de } X^3$

$X^5 \rightarrow X^1 \textit{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$

$S$   
 $\downarrow$   
 $X^4$

# Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang}$

$X^2 \rightarrow \textit{dianzi shang}$

$X^3 \rightarrow \textit{mao}$

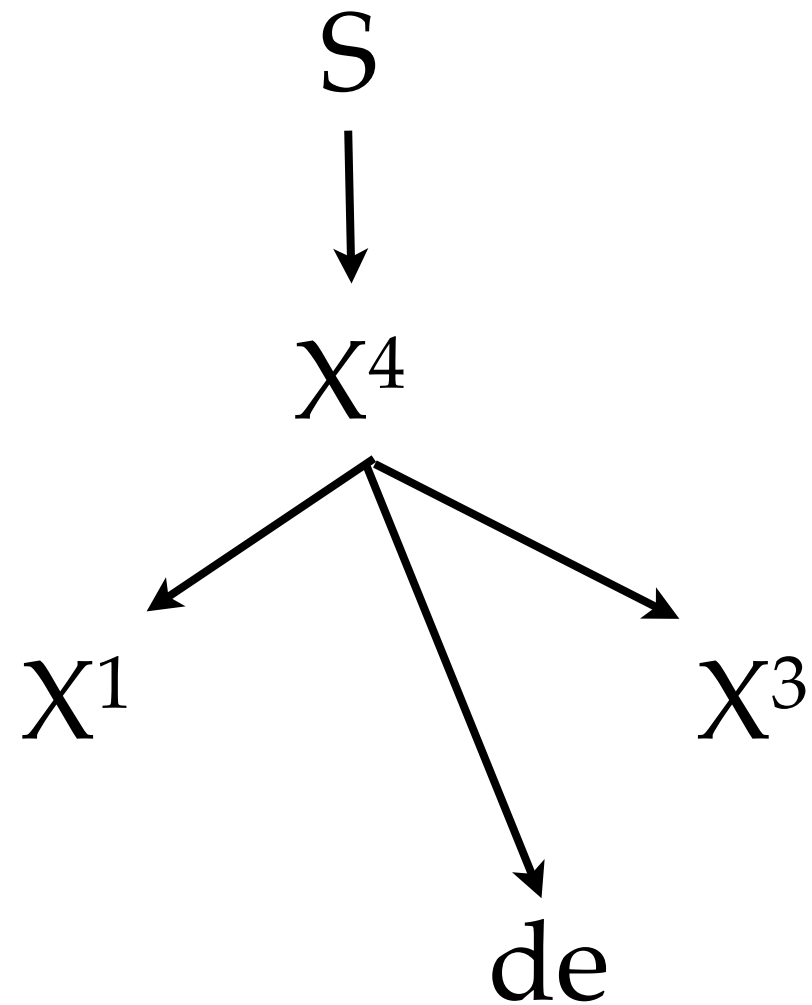
$X^4 \rightarrow X^1 \textit{ de } X^3$

$X^4 \rightarrow X^2 \textit{ de } X^3$

$X^5 \rightarrow X^1 \textit{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang}$

$X^2 \rightarrow \textit{dianzi shang}$

$X^3 \rightarrow \textit{mao}$

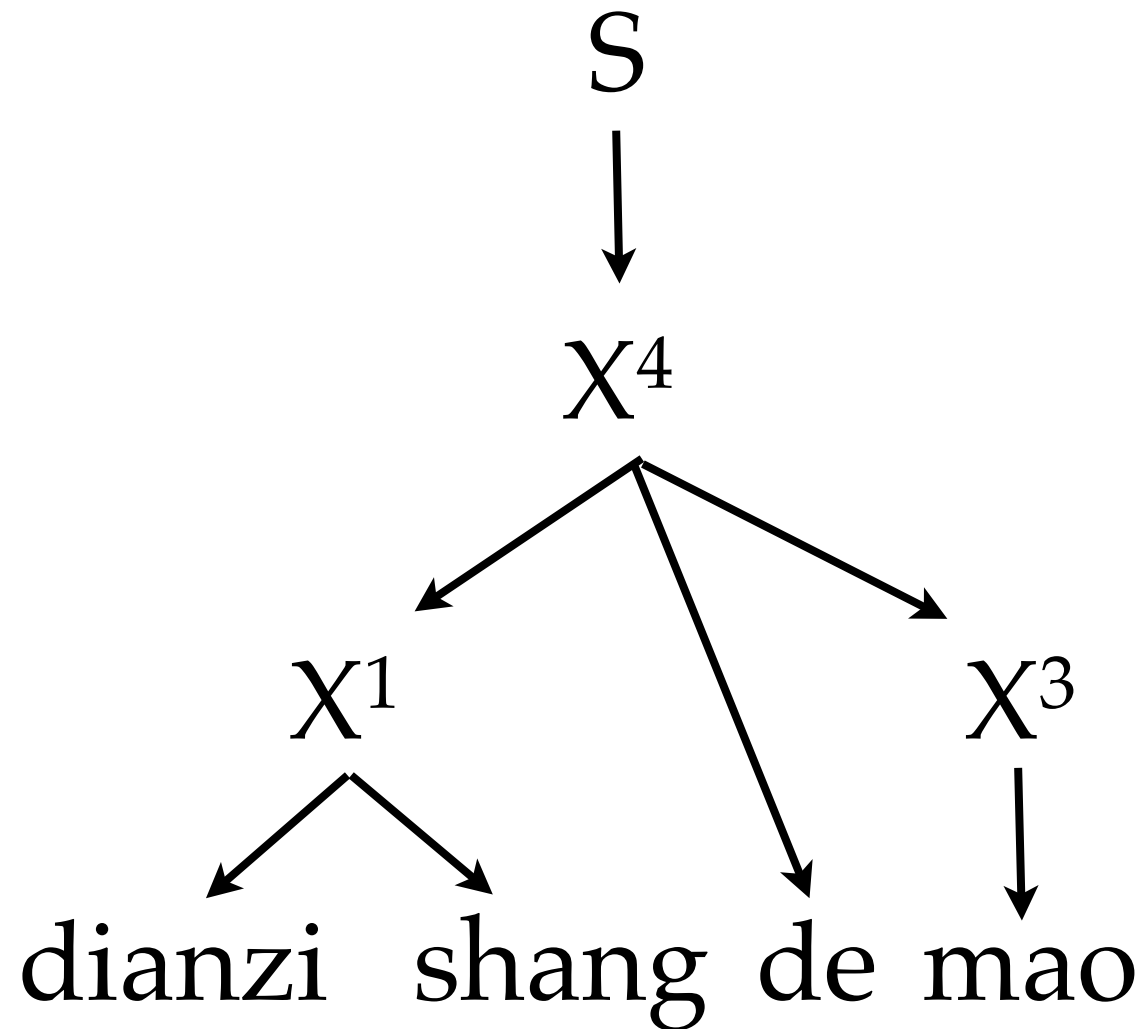
$X^4 \rightarrow X^1 \textit{ de } X^3$

$X^4 \rightarrow X^2 \textit{ de } X^3$

$X^5 \rightarrow X^1 \textit{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^2 \rightarrow \textit{dianzi shang/mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$X^4 \rightarrow X^2 \textit{ de } X^3 / X^3 \textit{ of } X^2$

$X^5 \rightarrow X^1 \textit{ de } X^3 / X^1 \textit{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$



# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$S \rightarrow X^5 / X^5$

# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$S \rightarrow X^5 / X^5$

S

S

# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$S \rightarrow X^5 / X^5$

S ..... S

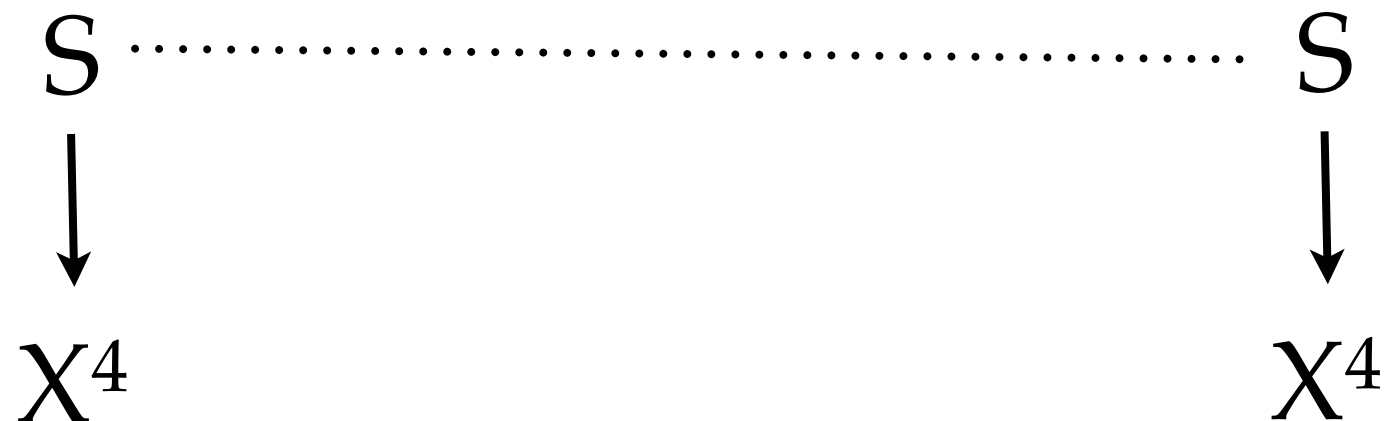
# *Synchronous* Context-Free Grammar

$X^1 \rightarrow dianzi\ shang/the\ mat$

$X^3 \rightarrow mao/the\ cat$

$X^4 \rightarrow X^1\ de\ X^3/X^3\ on\ X^1$

$S \rightarrow X^5/X^5$



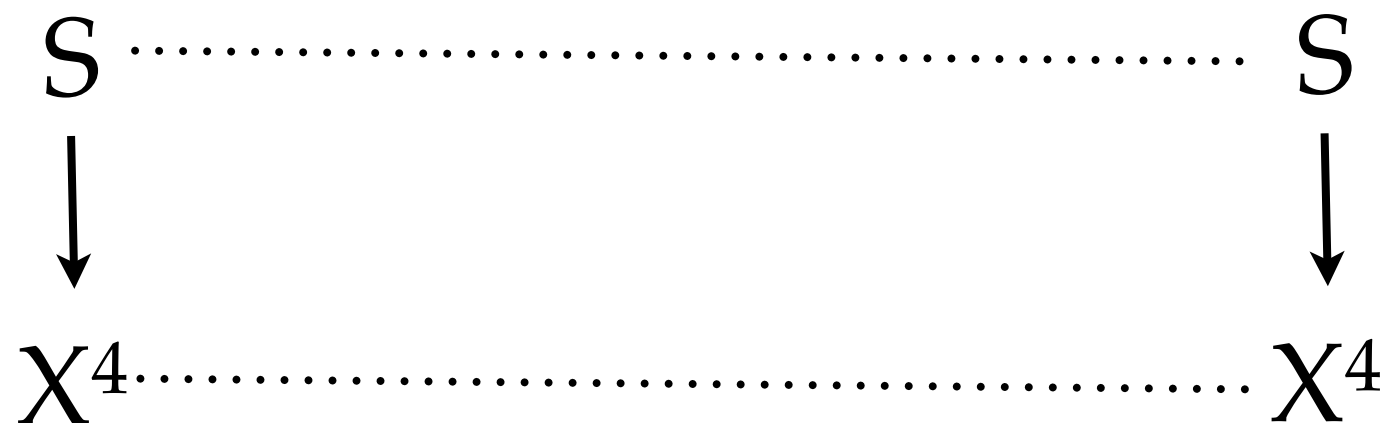
# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$S \rightarrow X^5 / X^5$



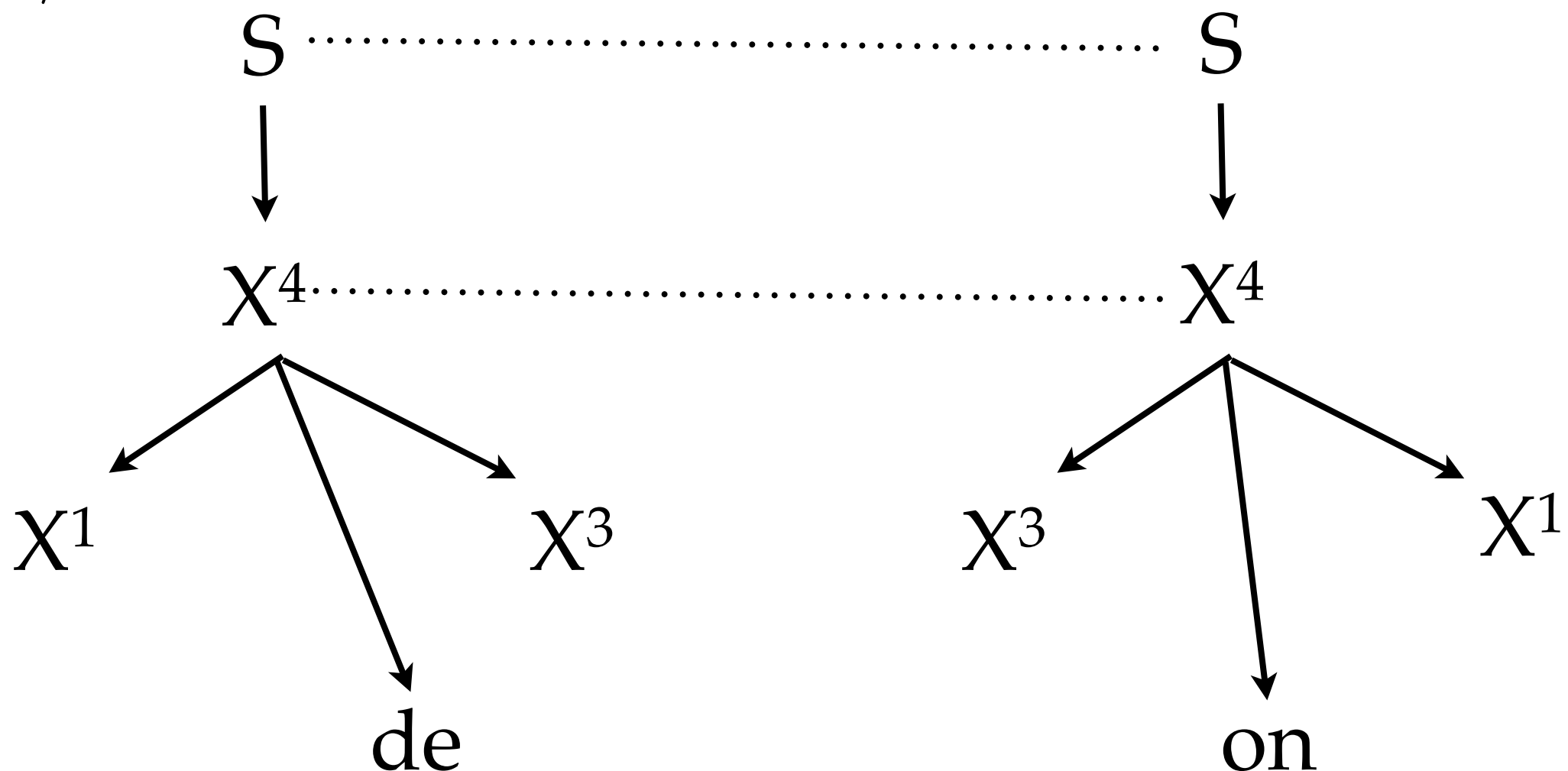
# *Synchronous* Context-Free Grammar

$X^1 \rightarrow dianzi\ shang/the\ mat$

$X^3 \rightarrow mao/the\ cat$

$X^4 \rightarrow X^1\ de\ X^3/X^3\ on\ X^1$

$S \rightarrow X^5/X^5$



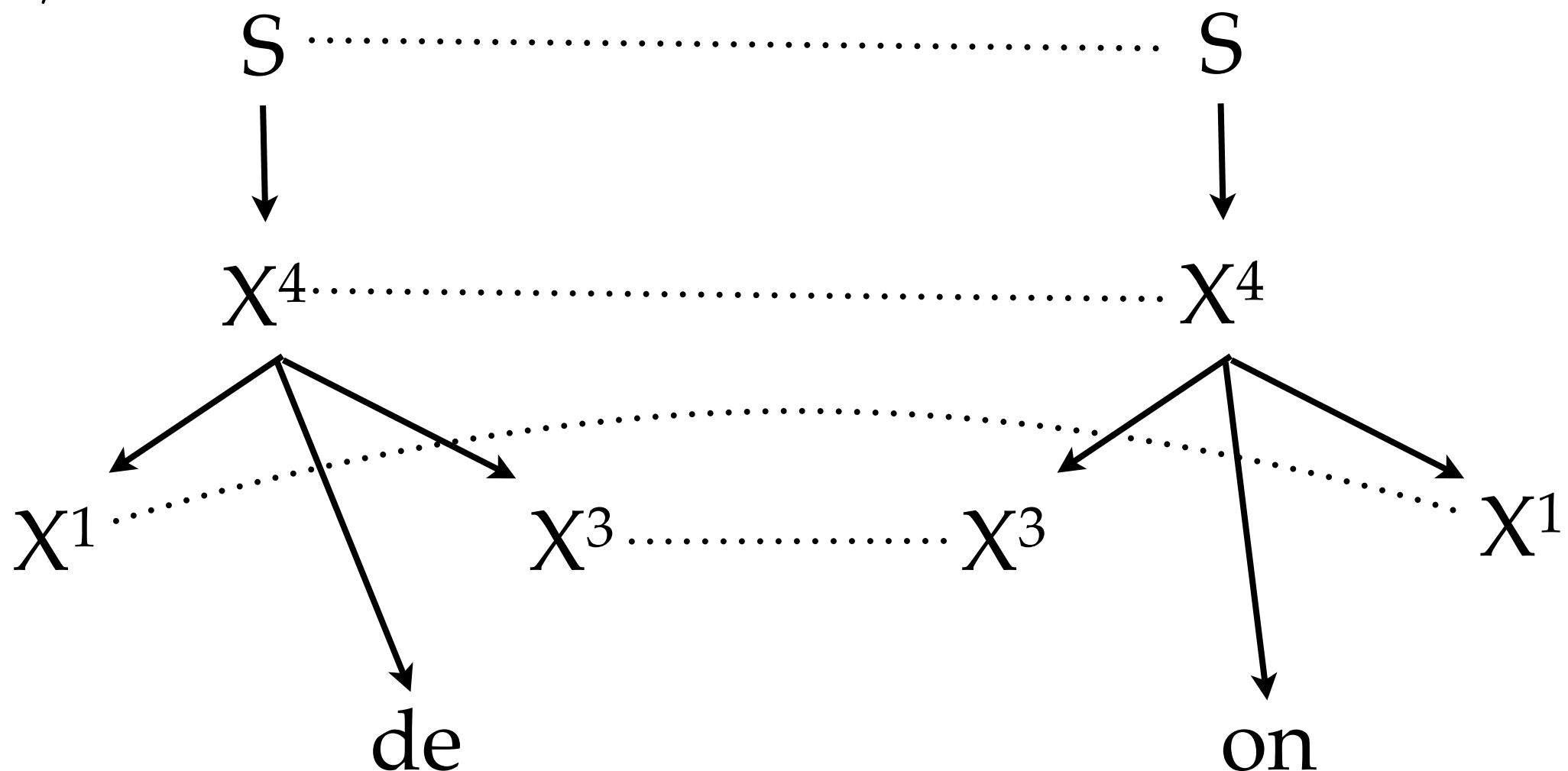
# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$S \rightarrow X^5 / X^5$



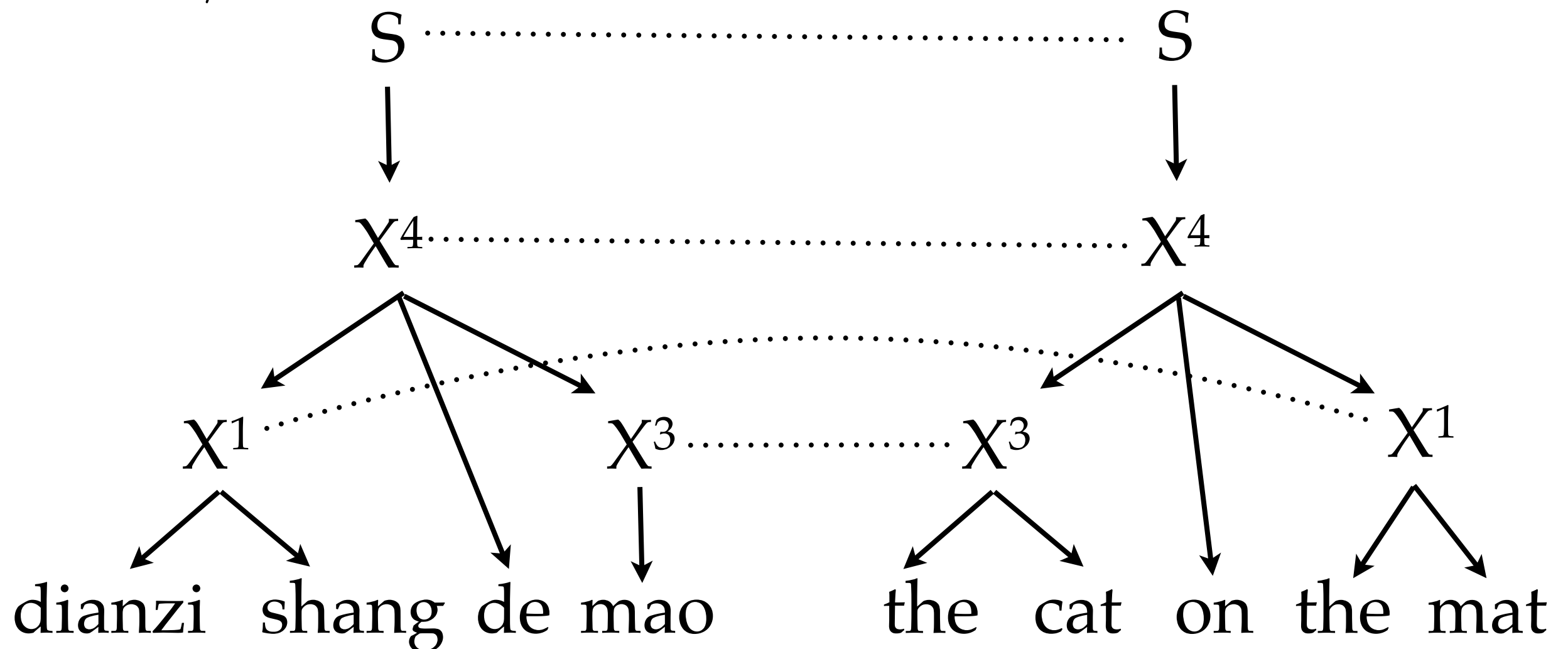
# *Synchronous* Context-Free Grammar

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$S \rightarrow X^5 / X^5$

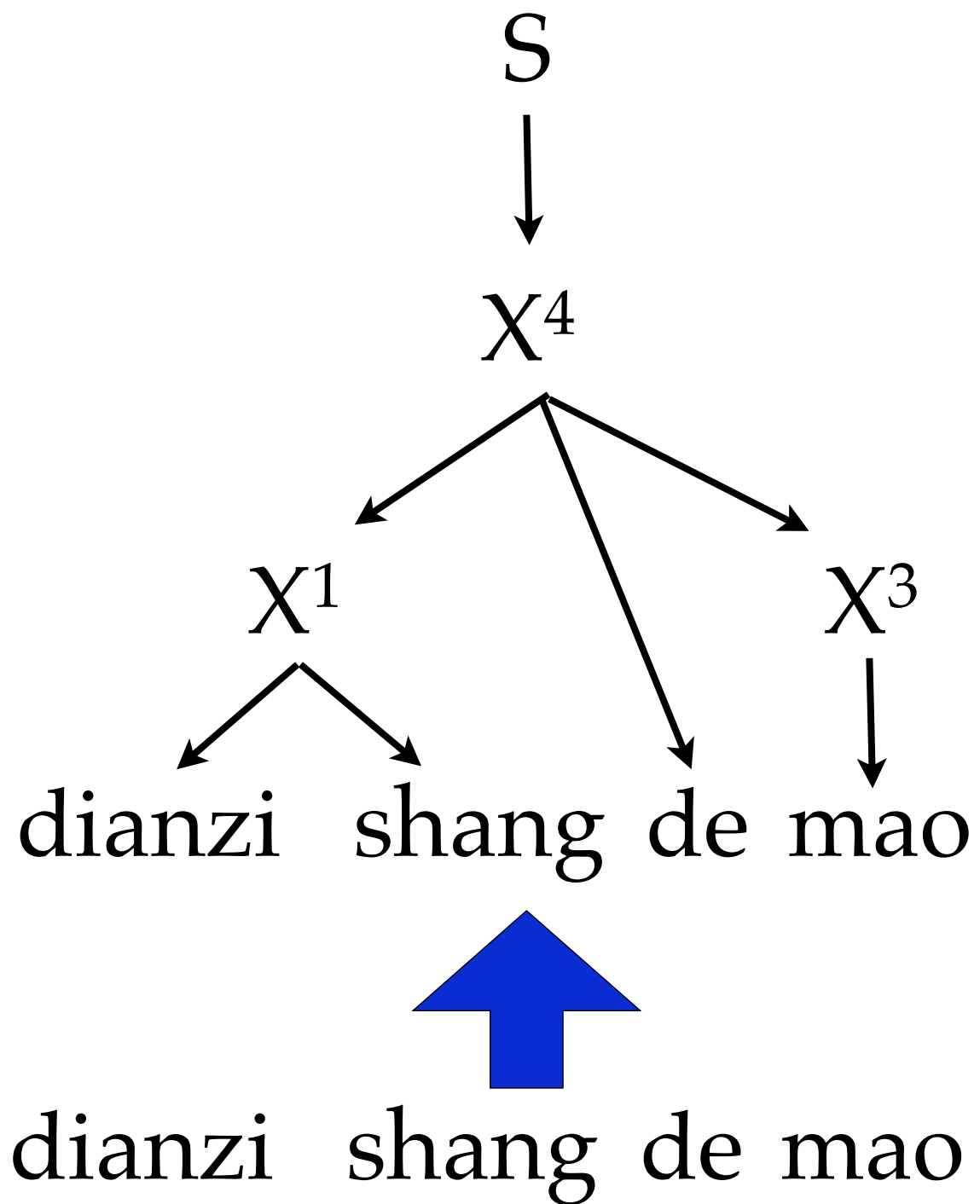




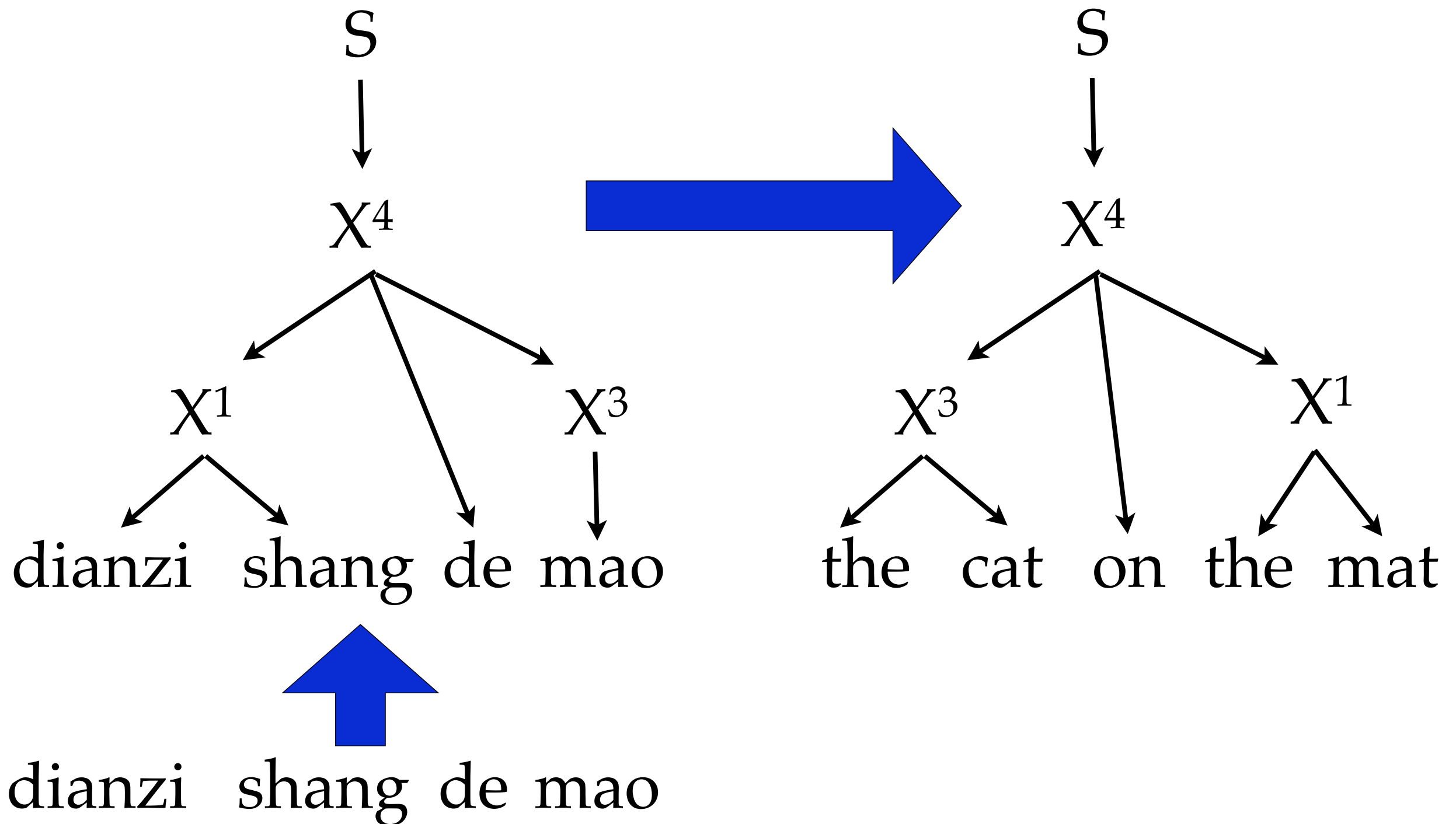
# *Synchronous* Context-Free Grammar

dianzi shang de mao

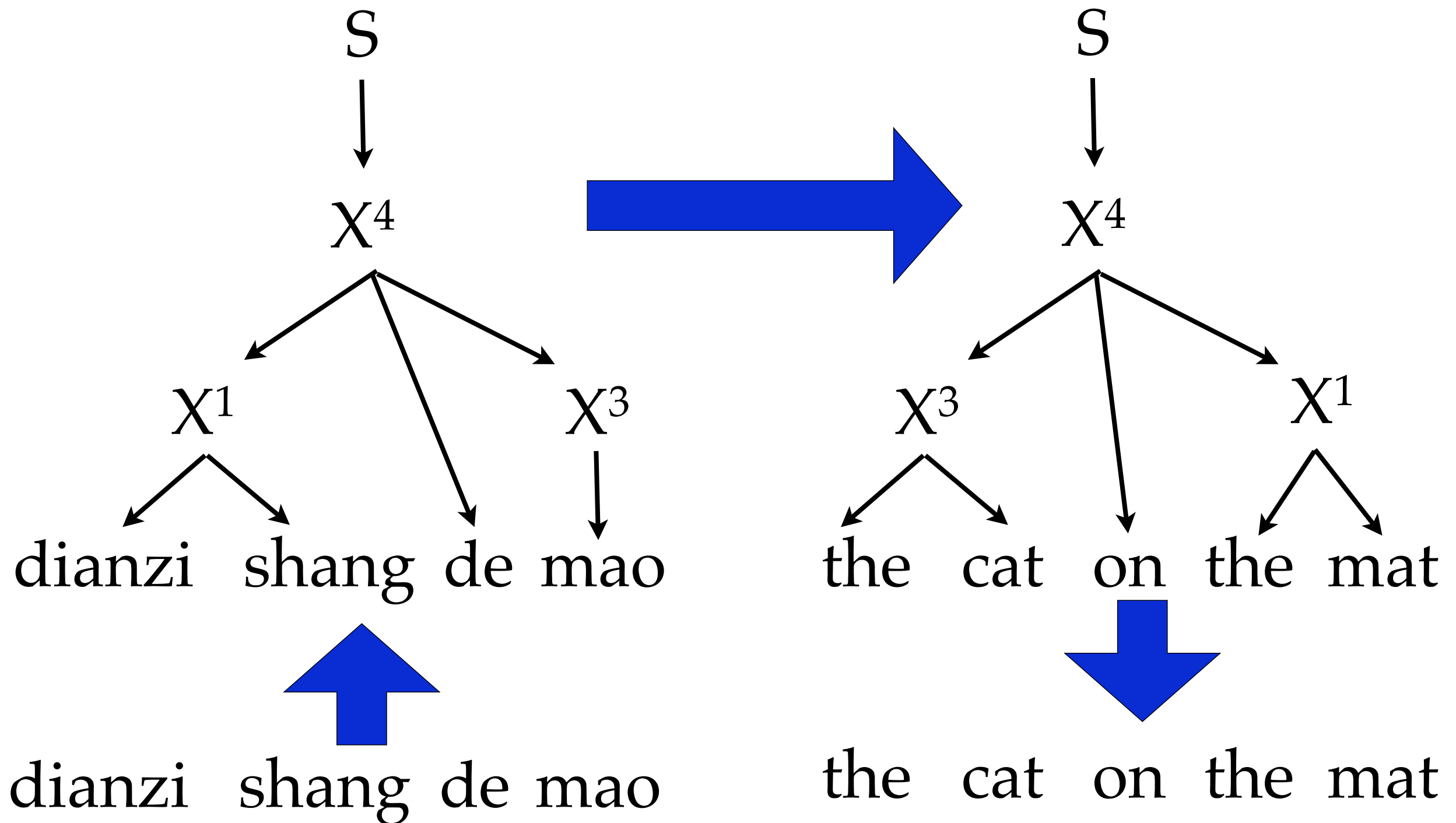
# *Synchronous* Context-Free Grammar



# *Synchronous* Context-Free Grammar



# *Synchronous* Context-Free Grammar



# Parsing

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^2 \rightarrow \textit{dianzi shang/mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$X^4 \rightarrow X^2 \textit{ de } X^3 / X^3 \textit{ of } X^2$

$X^5 \rightarrow X^1 \textit{ de } X^3 / X^1 \textit{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$

dianzi shang de mao

# Parsing

$X^1 \rightarrow \textit{dianzi shang}$

$X^2 \rightarrow \textit{dianzi shang}$

$X^3 \rightarrow \textit{mao}$

$X^4 \rightarrow X^1 \textit{ de } X^3$

$X^4 \rightarrow X^2 \textit{ de } X^3$

$X^5 \rightarrow X^1 \textit{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$

dianzi shang de mao

# Parsing

$X^1 \rightarrow \textit{dianzi shang}$

$X^2 \rightarrow \textit{dianzi shang}$

$X^3 \rightarrow \textit{mao}$

$X^4 \rightarrow X^1 \textit{ de } X^3$

$X^4 \rightarrow X^2 \textit{ de } X^3$

$X^5 \rightarrow X^1 \textit{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$

$\begin{matrix} & & 0 & & 1 & & 2 & & 3 & & 4 \\ & & \textit{dianzi} & & \textit{shang} & & \textit{de} & & \textit{mao} & & \end{matrix}$

# Parsing

$X^1 \rightarrow dianzi\ shang$

$X^2 \rightarrow dianzi\ shang$

$X^3 \rightarrow mao$

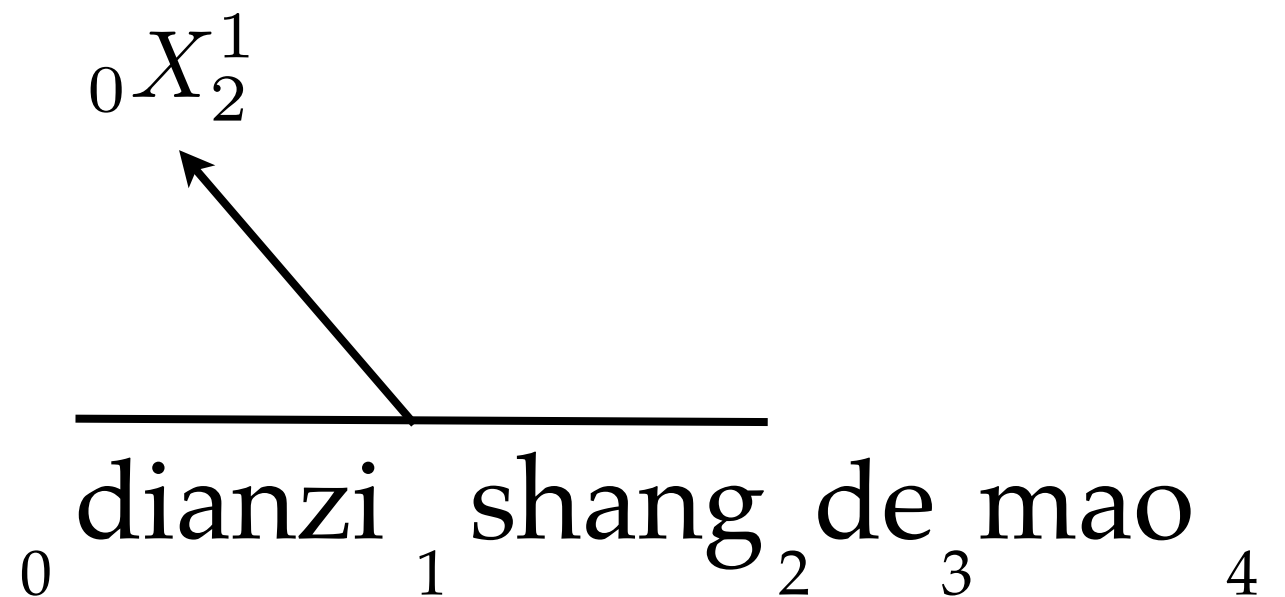
$X^4 \rightarrow X^1\ de\ X^3$

$X^4 \rightarrow X^2\ de\ X^3$

$X^5 \rightarrow X^1\ de\ X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$





# Parsing

$X^1 \rightarrow dianzi\ shang$

$X^2 \rightarrow dianzi\ shang$

$X^3 \rightarrow mao$

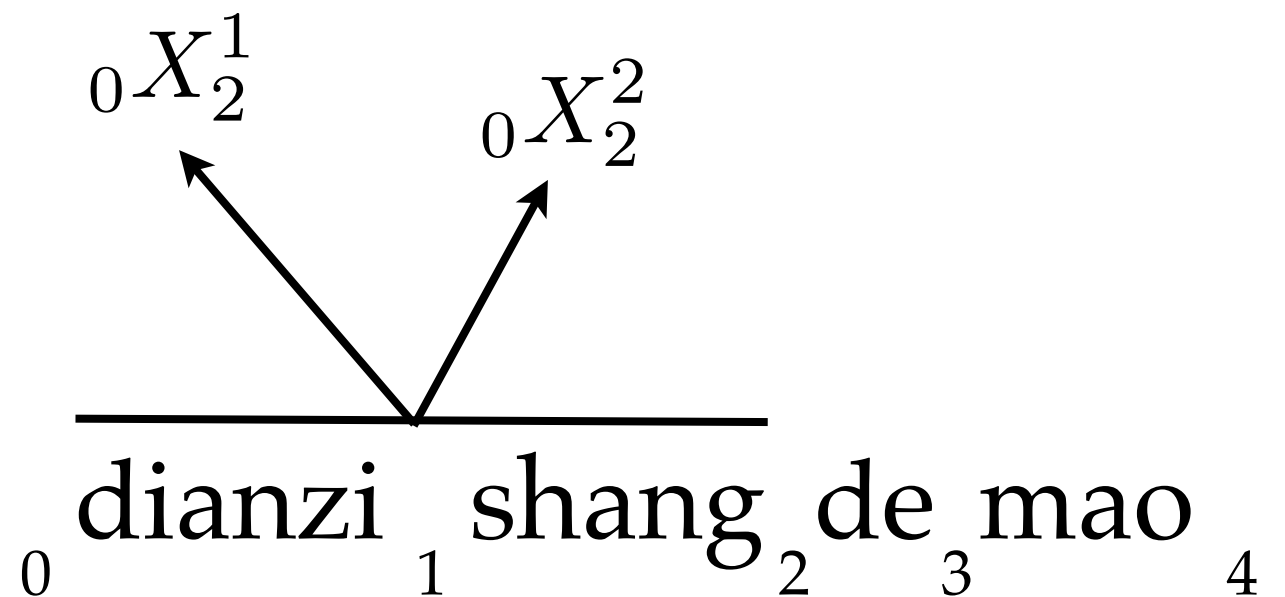
$X^4 \rightarrow X^1\ de\ X^3$

$X^4 \rightarrow X^2\ de\ X^3$

$X^5 \rightarrow X^1\ de\ X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# Parsing

$X^1 \rightarrow \text{dianzi shang}$

$X^2 \rightarrow \text{dianzi shang}$

$X^3 \rightarrow \text{mao}$

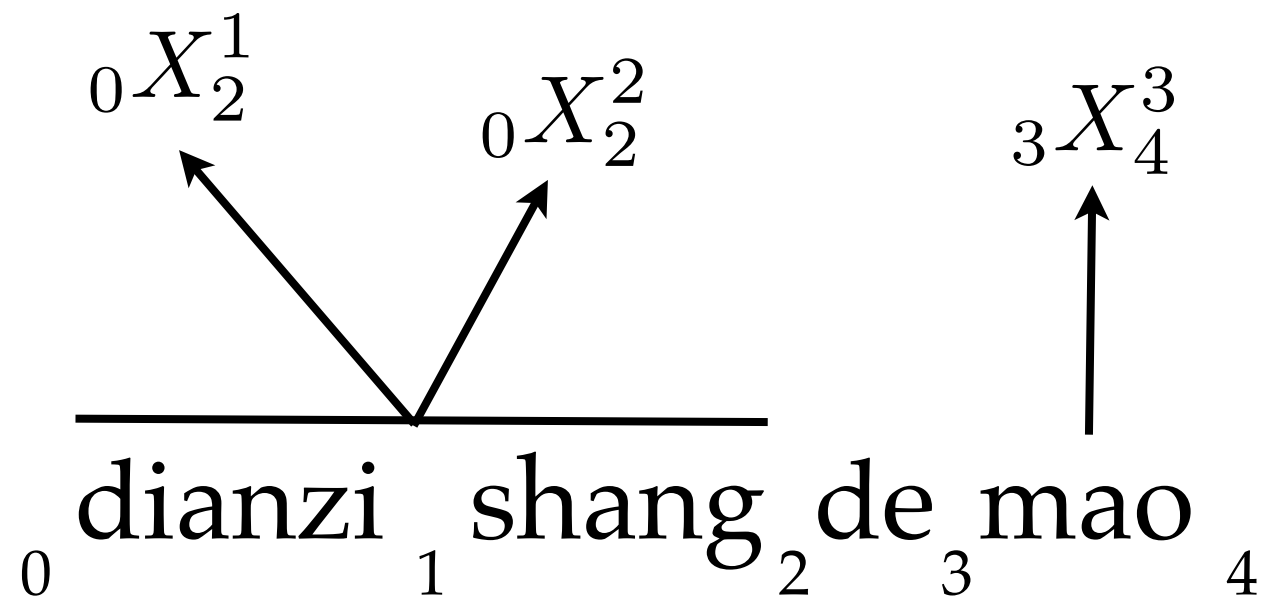
$X^4 \rightarrow X^1 \text{ de } X^3$

$X^4 \rightarrow X^2 \text{ de } X^3$

$X^5 \rightarrow X^1 \text{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# Parsing

$X^1 \rightarrow \text{dianzi shang}$

$X^2 \rightarrow \text{dianzi shang}$

$X^3 \rightarrow \text{mao}$

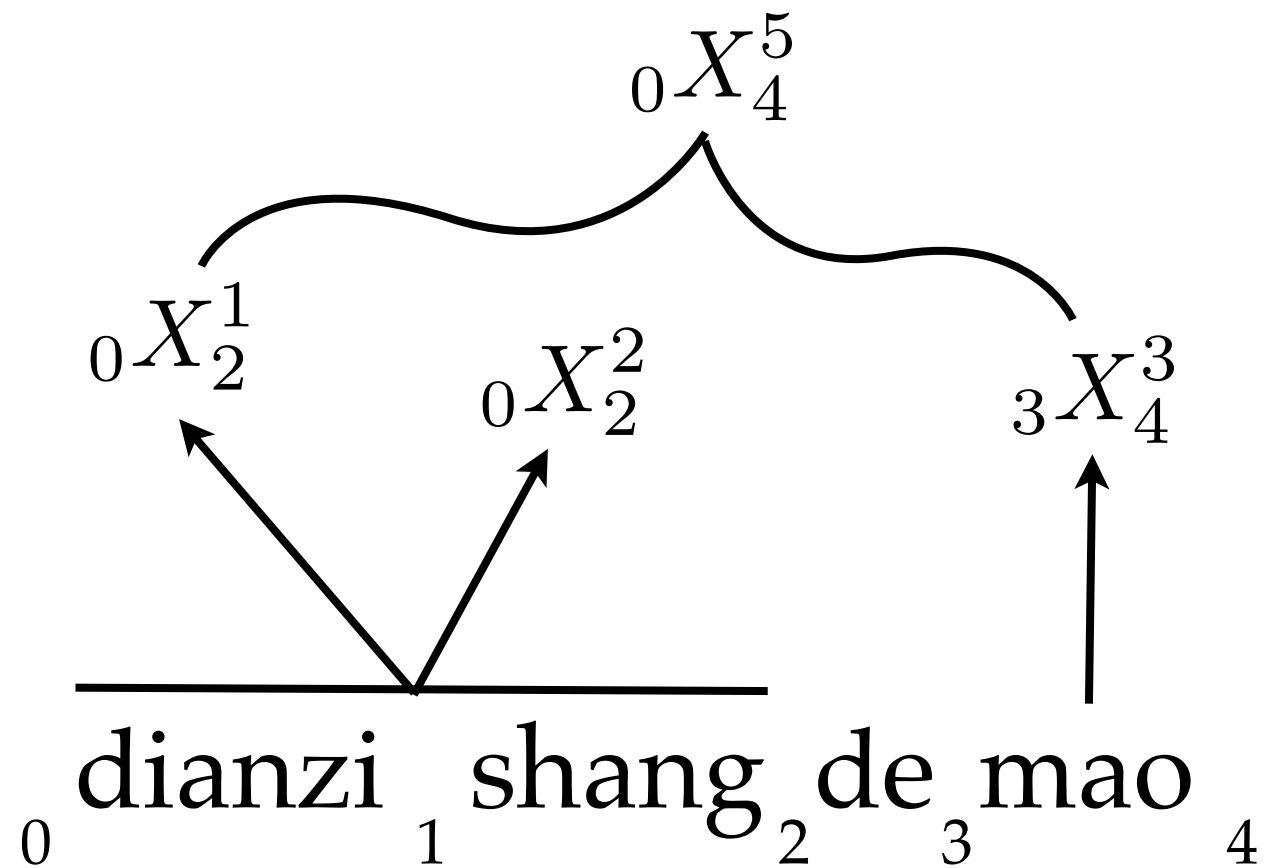
$X^4 \rightarrow X^1 \text{ de } X^3$

$X^4 \rightarrow X^2 \text{ de } X^3$

$X^5 \rightarrow X^1 \text{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# Parsing

$X^1 \rightarrow \text{dianzi shang}$

$X^2 \rightarrow \text{dianzi shang}$

$X^3 \rightarrow \text{mao}$

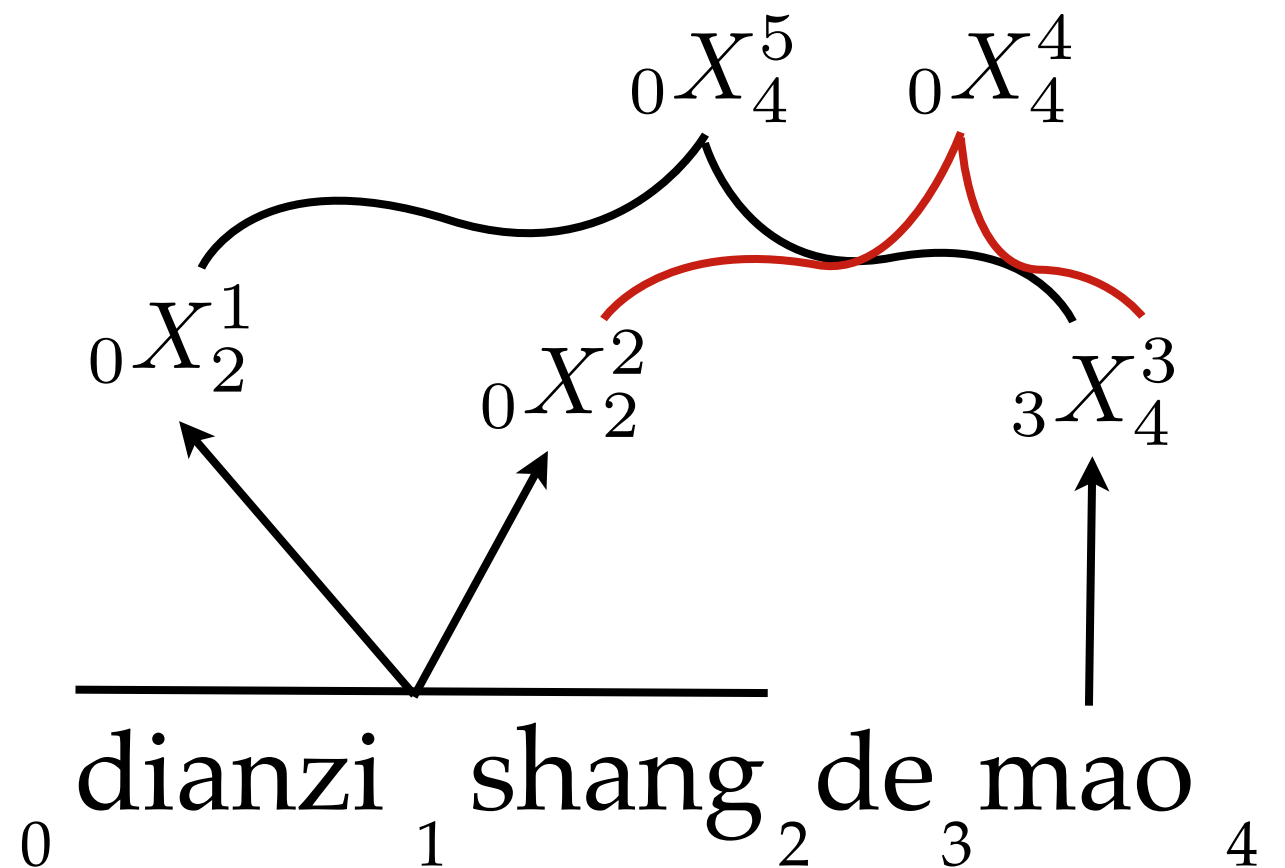
$X^4 \rightarrow X^1 \text{ de } X^3$

$X^4 \rightarrow X^2 \text{ de } X^3$

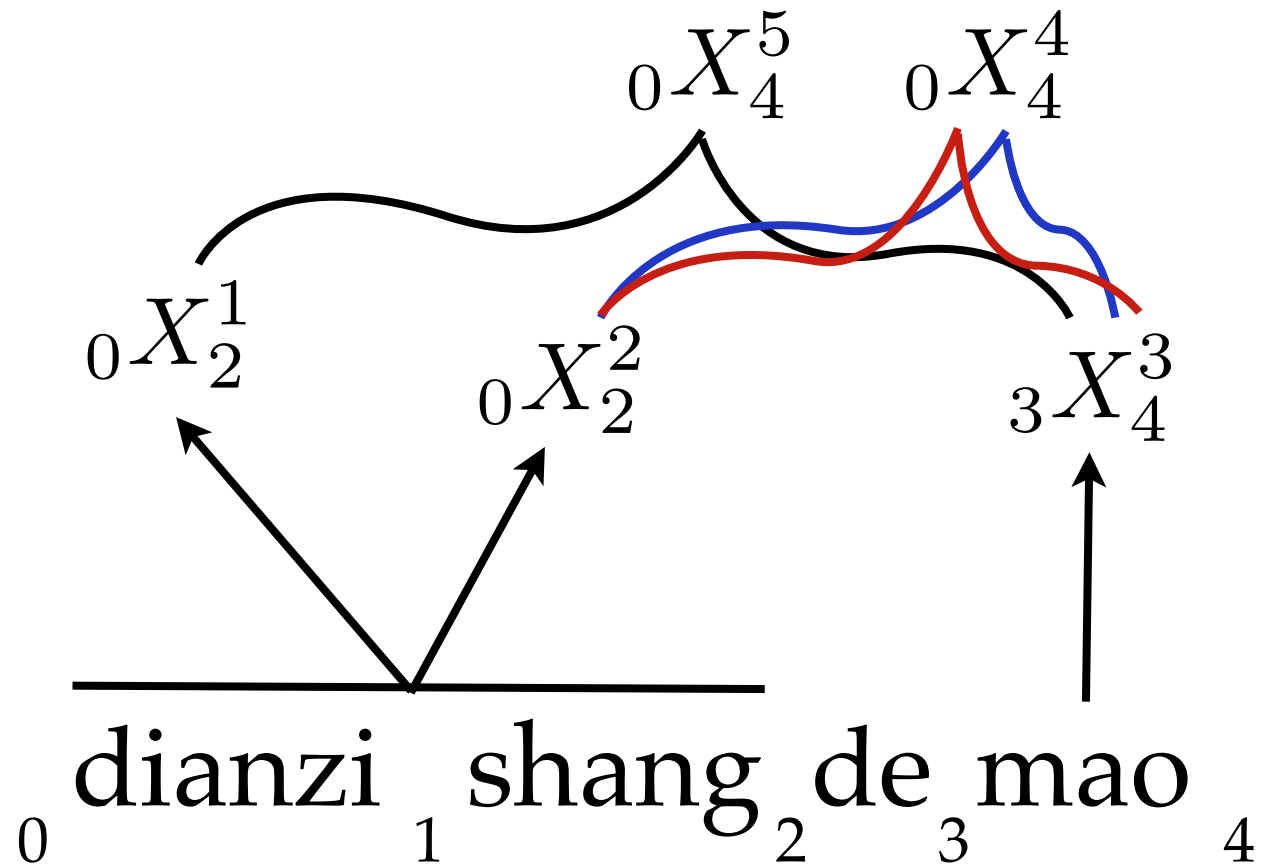
$X^5 \rightarrow X^1 \text{ de } X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# Parsing

$$X^1 \rightarrow dianzi\ shang$$
$$X^2 \rightarrow dianzi\ shang$$
$$X^3 \rightarrow mao$$
$$X^4 \rightarrow X^1 \text{ de } X^3$$
$$X^4 \rightarrow X^2 \text{ de } X^3$$
$$X^5 \rightarrow X^1 \text{ de } X^3$$
$$S \rightarrow X^4$$
$$S \rightarrow X^5$$


# Parsing

$X^1 \rightarrow dianzi\ shang$

$X^2 \rightarrow dianzi\ shang$

$X^3 \rightarrow mao$

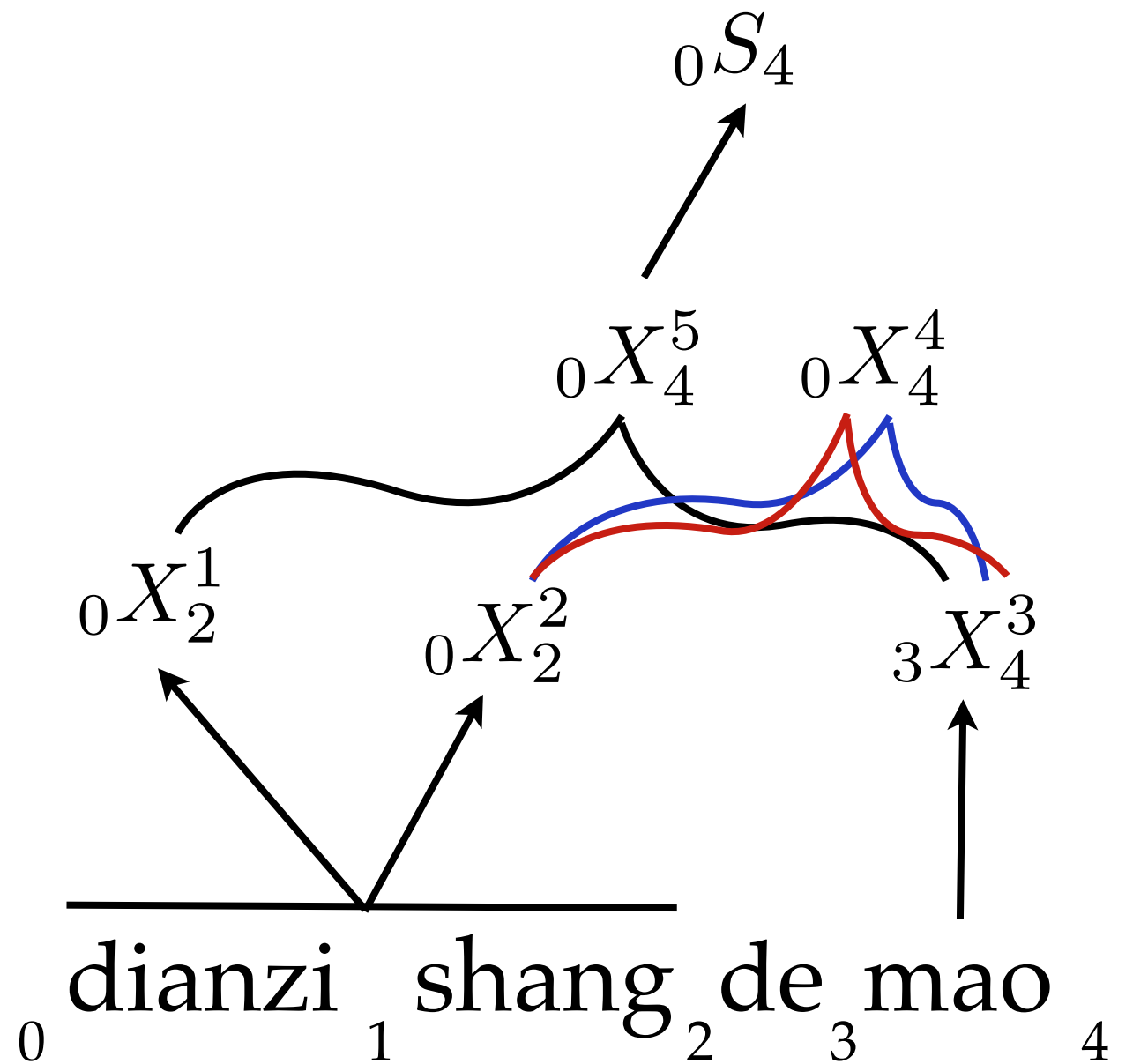
$X^4 \rightarrow X^1\ de\ X^3$

$X^4 \rightarrow X^2\ de\ X^3$

$X^5 \rightarrow X^1\ de\ X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$



# Parsing

$X^1 \rightarrow dianzi\ shang$

$X^2 \rightarrow dianzi\ shang$

$X^3 \rightarrow mao$

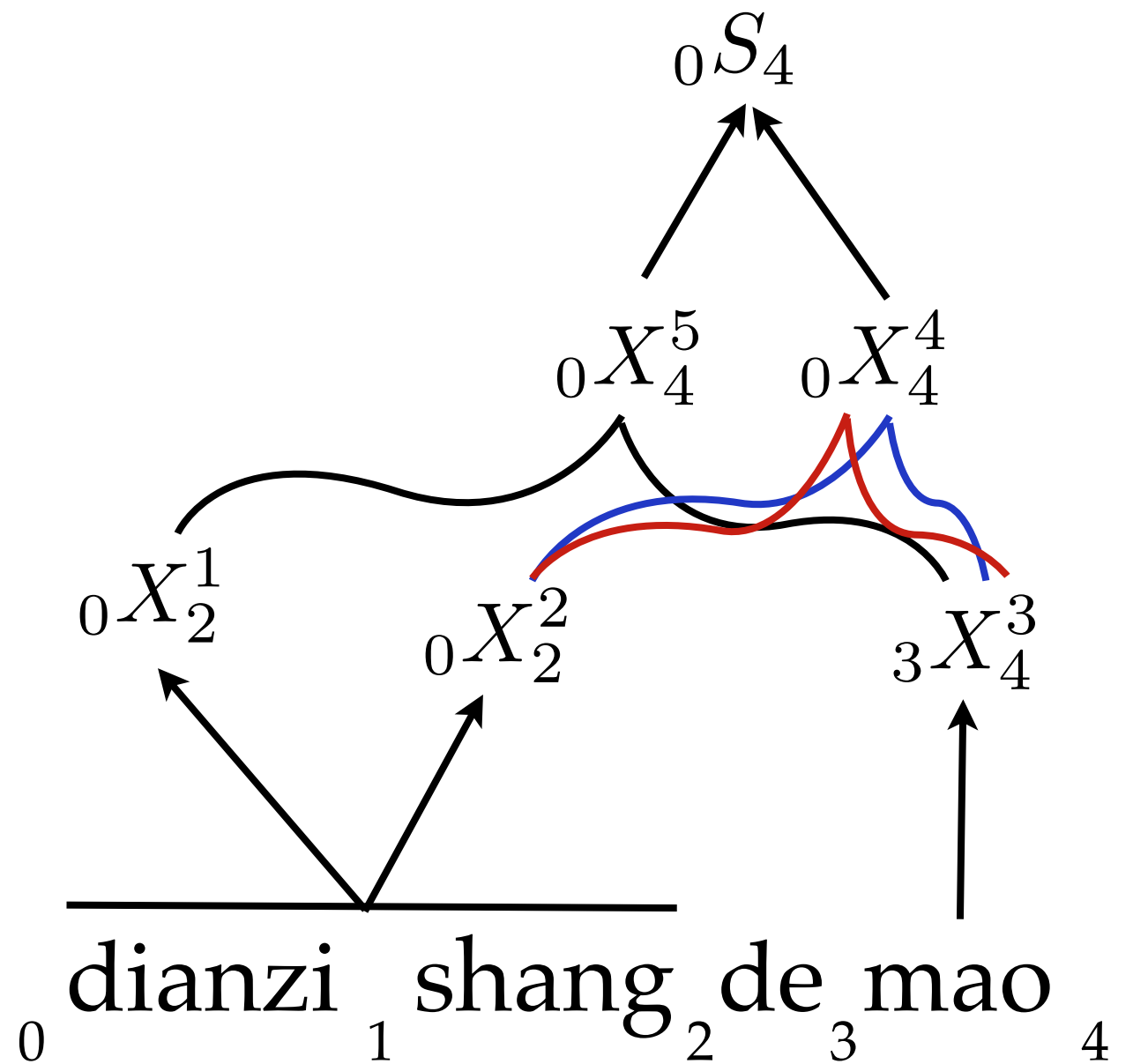
$X^4 \rightarrow X^1\ de\ X^3$

$X^4 \rightarrow X^2\ de\ X^3$

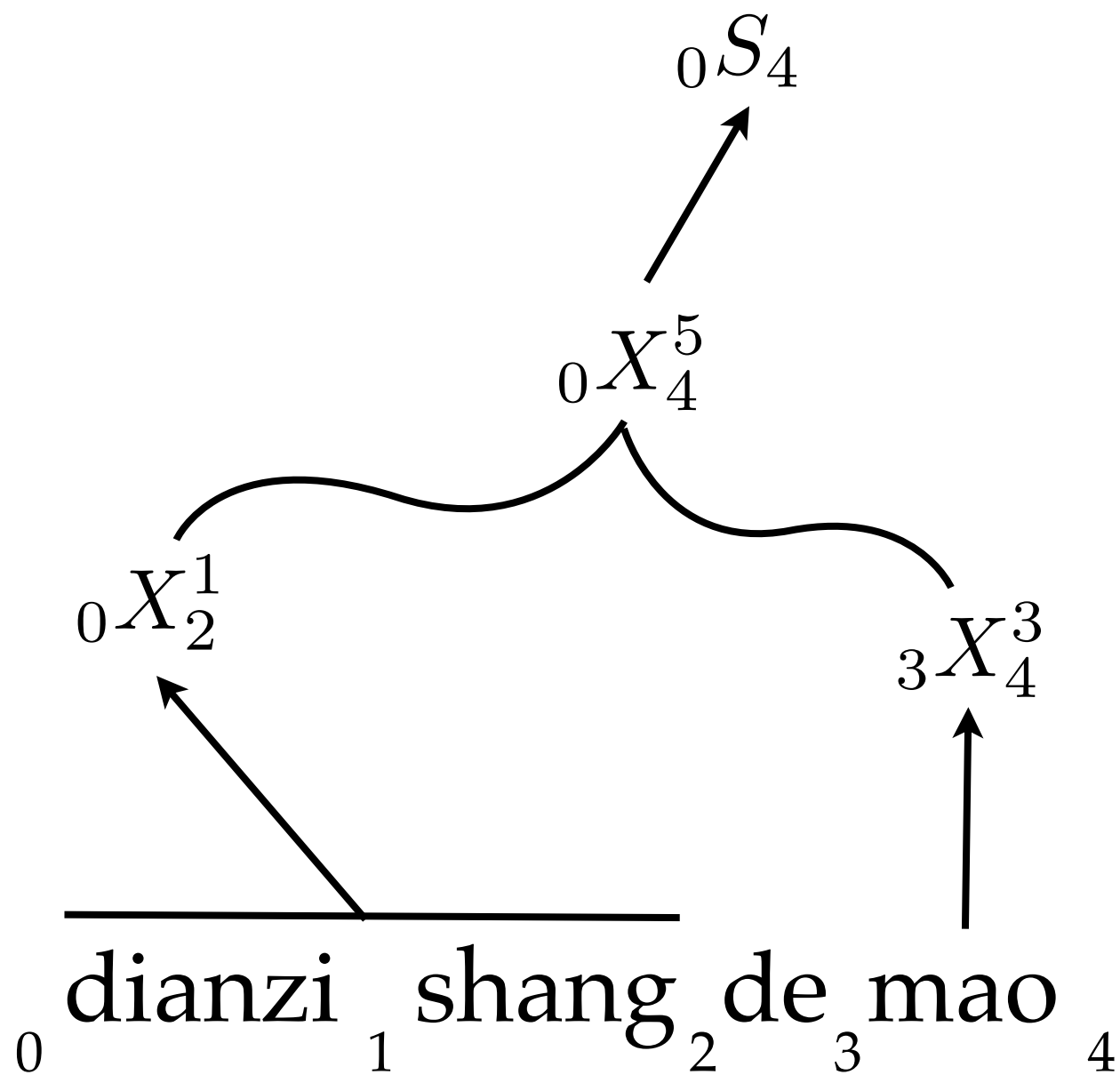
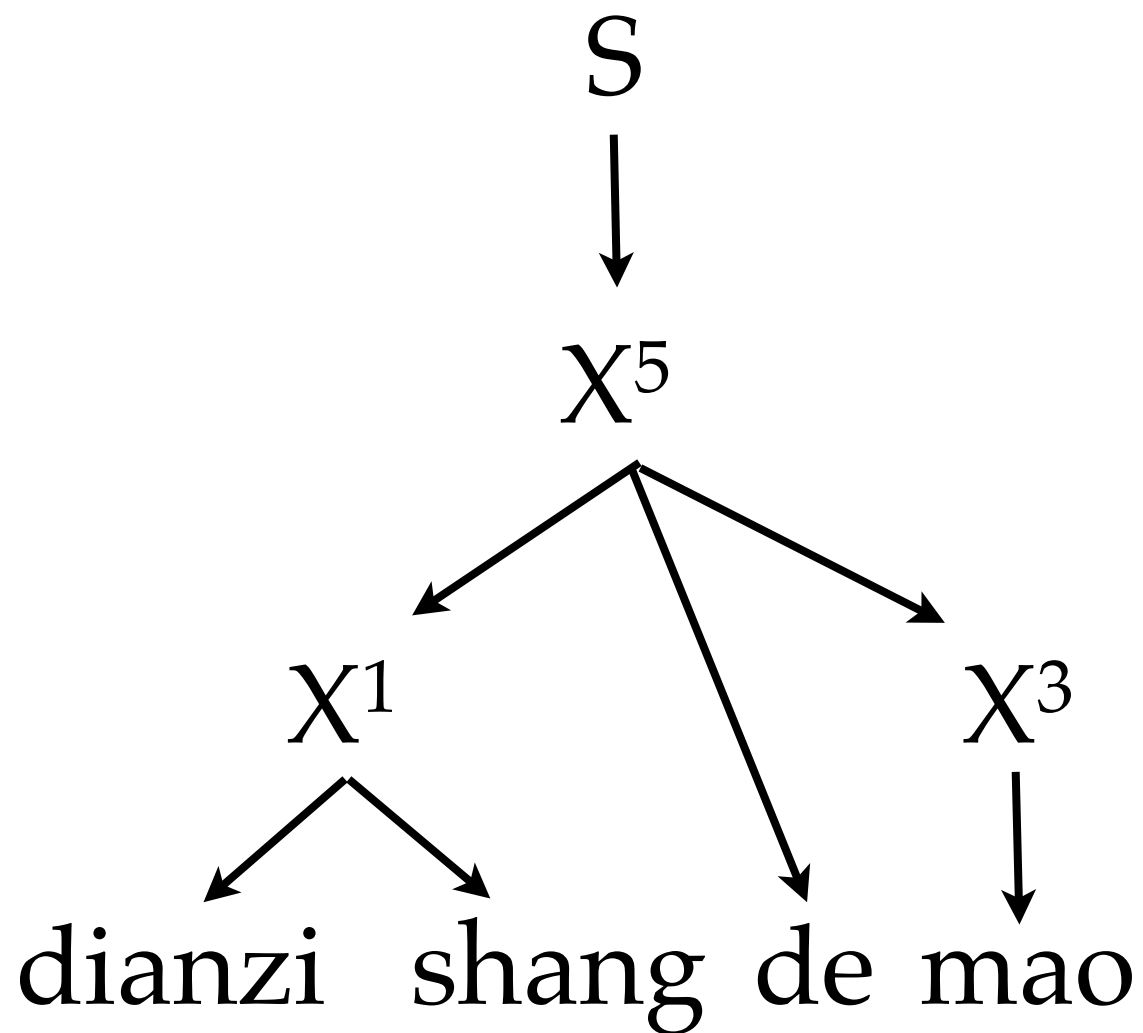
$X^5 \rightarrow X^1\ de\ X^3$

$S \rightarrow X^4$

$S \rightarrow X^5$

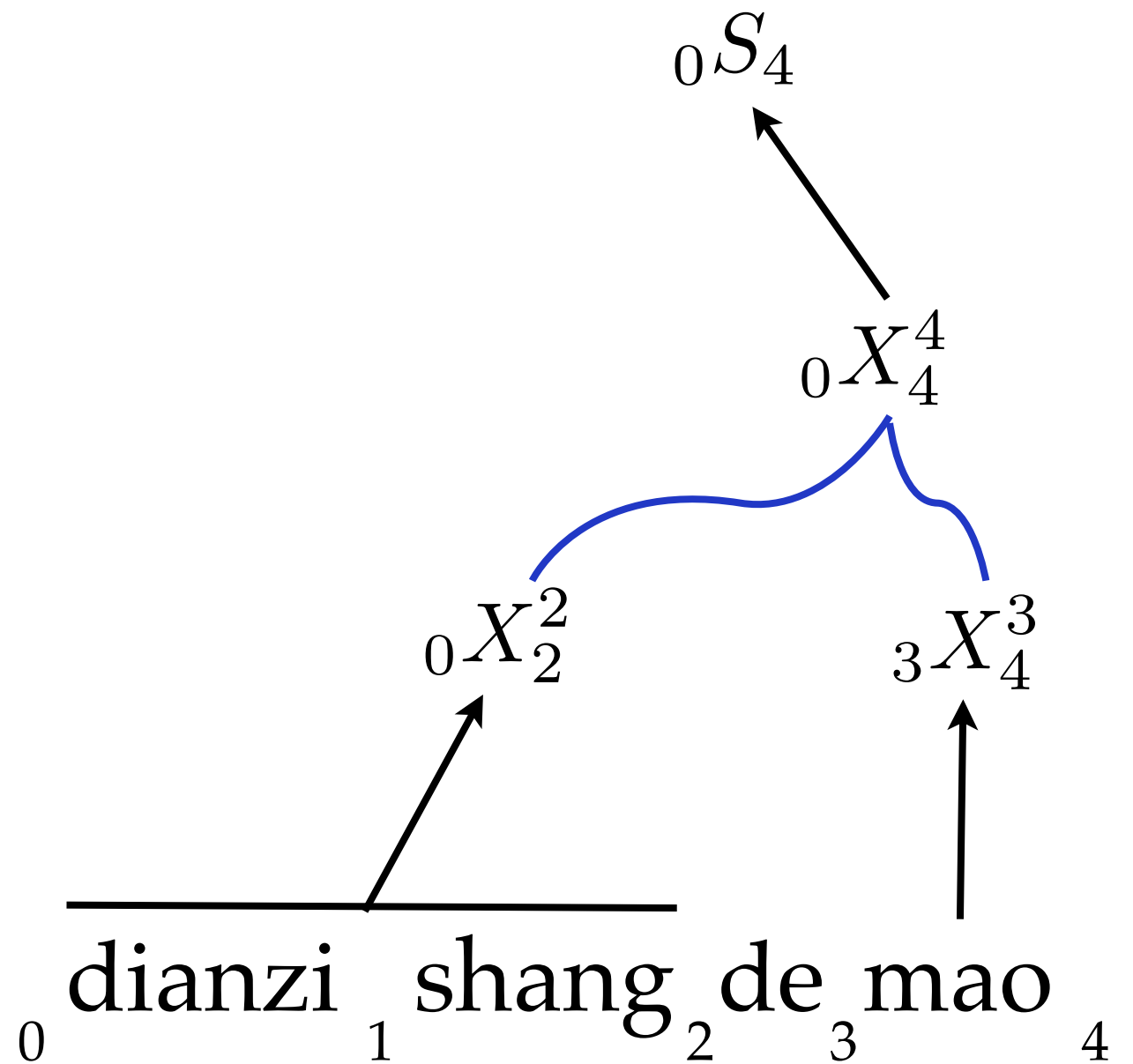
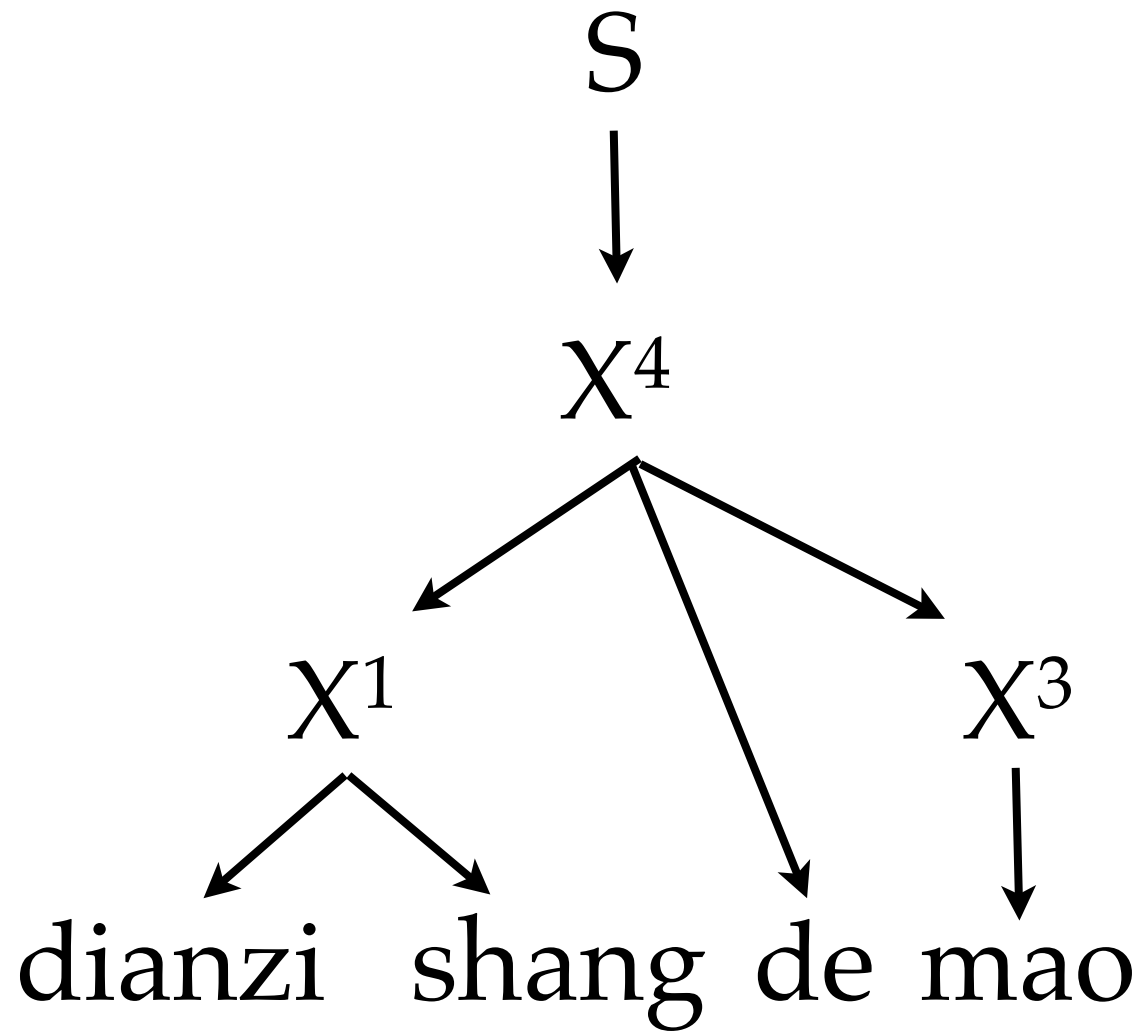


# Parsing

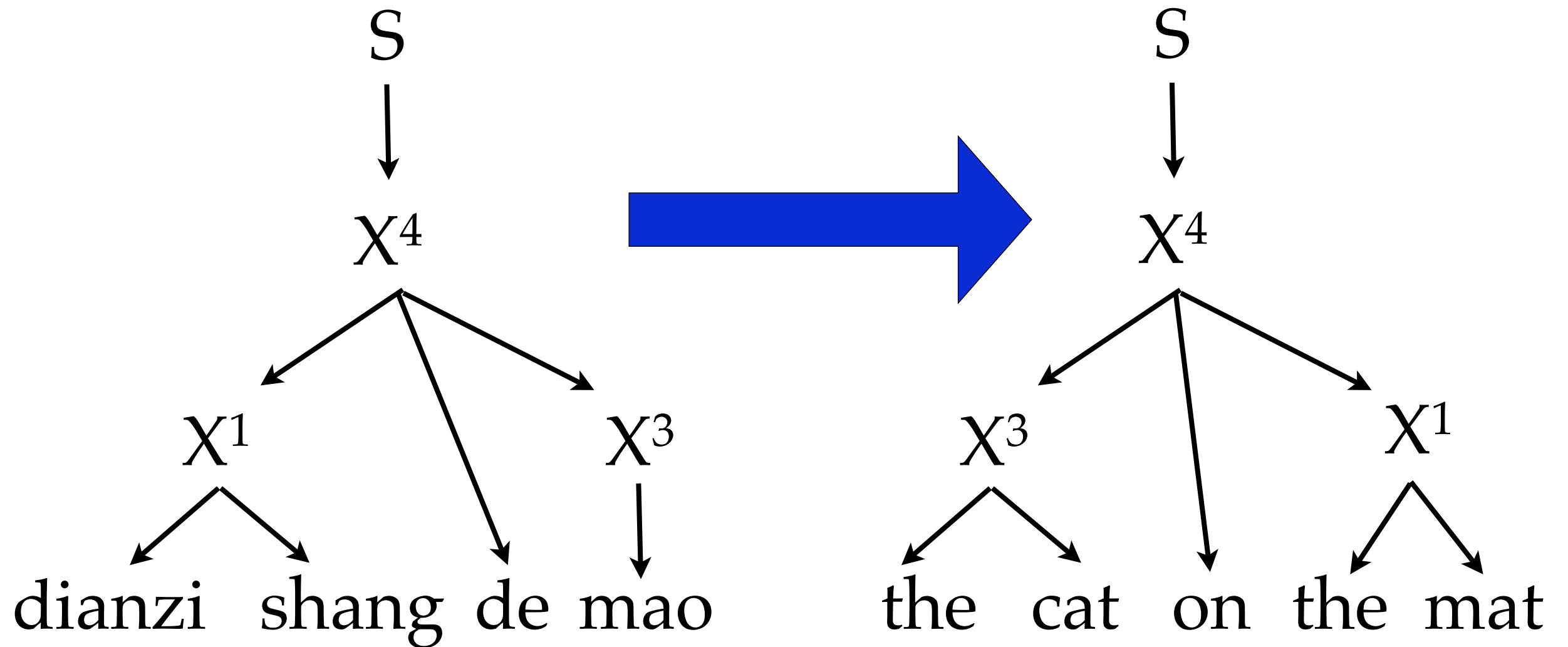




# Parsing



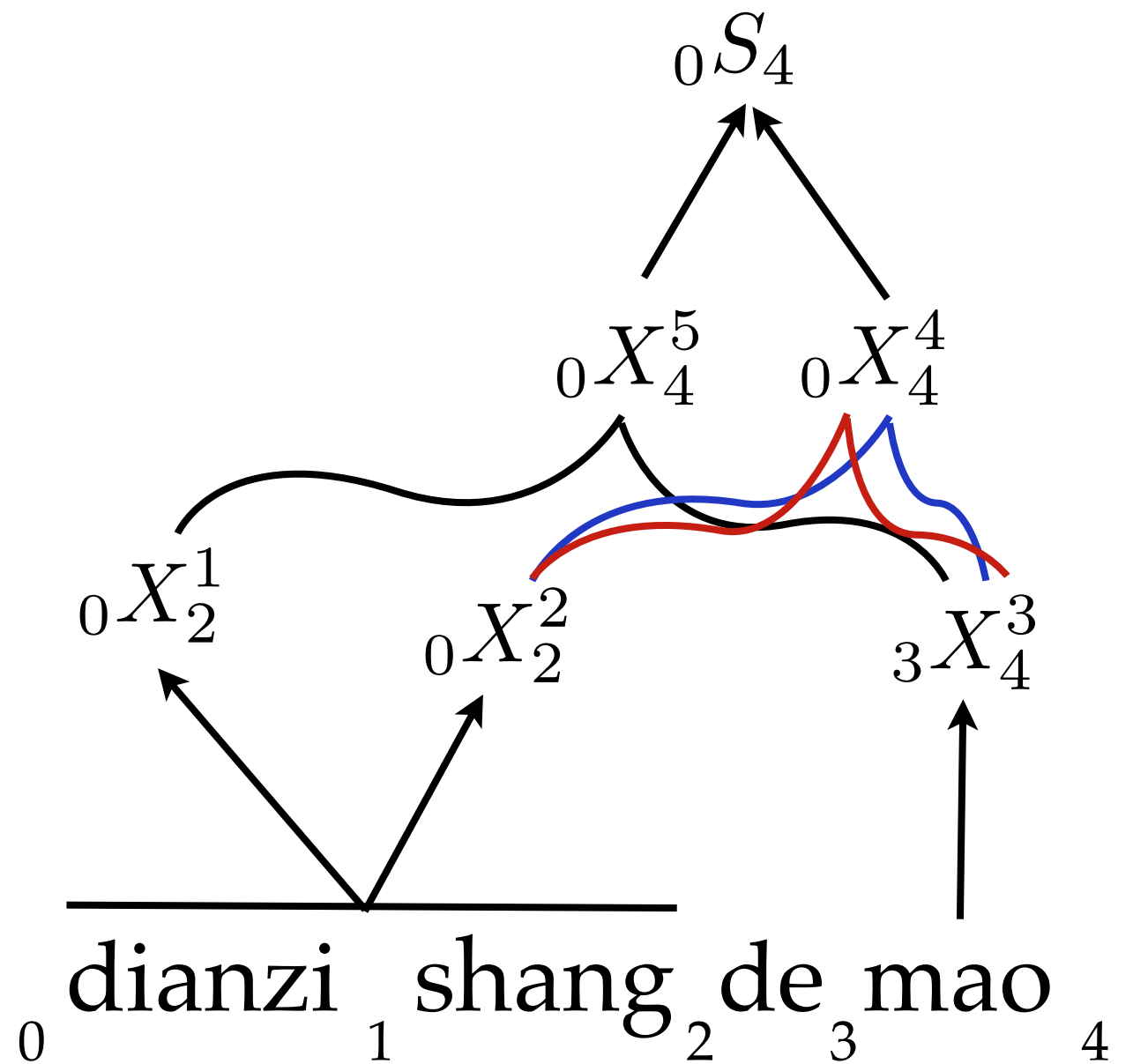
# Parsing



# Analysis

Nodes:  $O(Gn^2)$

Edges:  $O(G^3n^3)$



# Not so fast...

- Speed and memory footprint matter for both evaluation and tuning.
- What if  $G$  is really big?
- What happens when we add an  $n$ -gram language model?

# Not so fast...

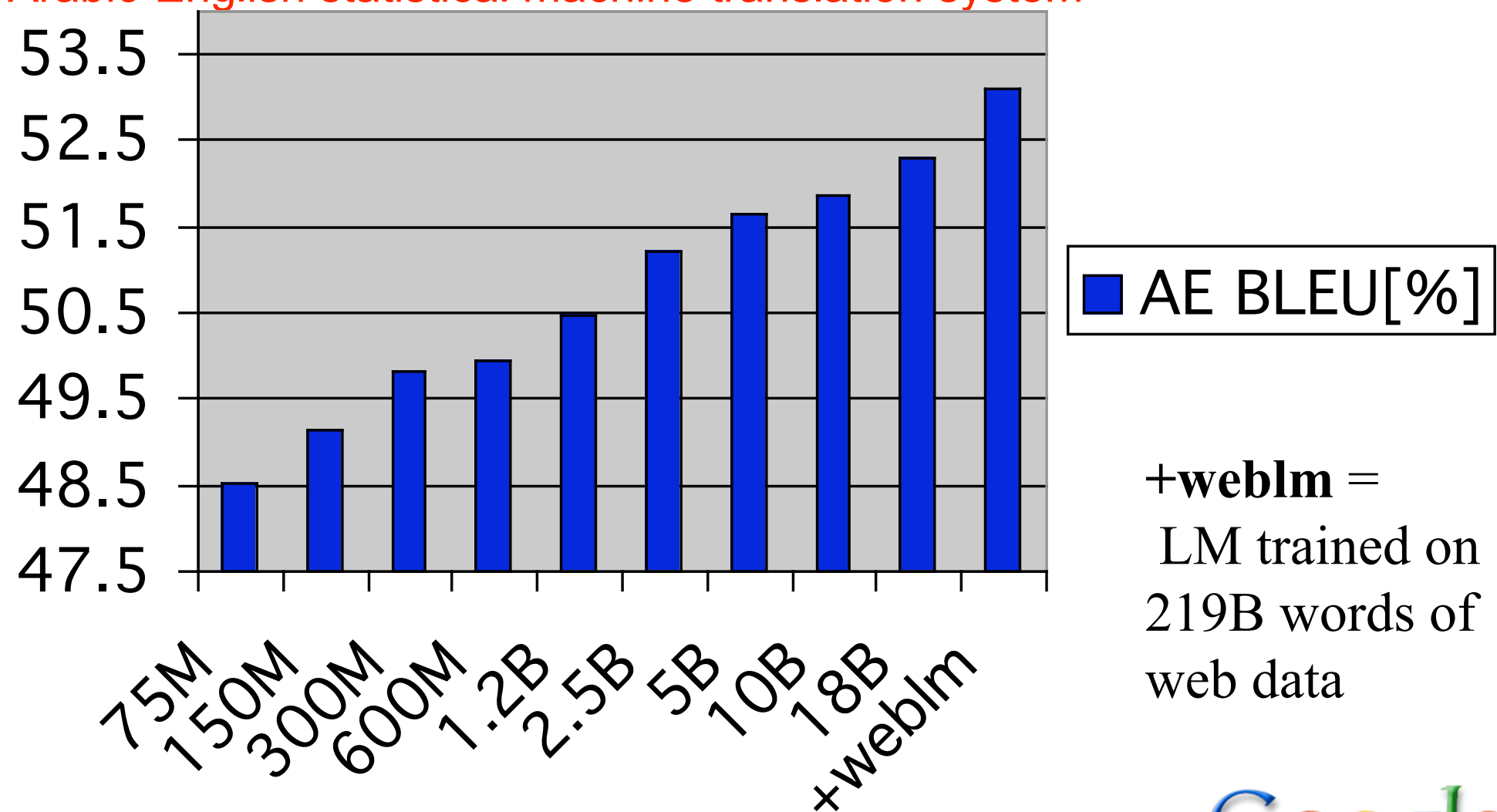
- Speed and memory footprint matter for both evaluation and tuning.
- What if  $G$  is really big?
- What happens when we add an  $n$ -gram language model?

$$\operatorname{argmax}_{\text{English}} p(\text{Urdu} | \text{English}) p(\text{English})$$

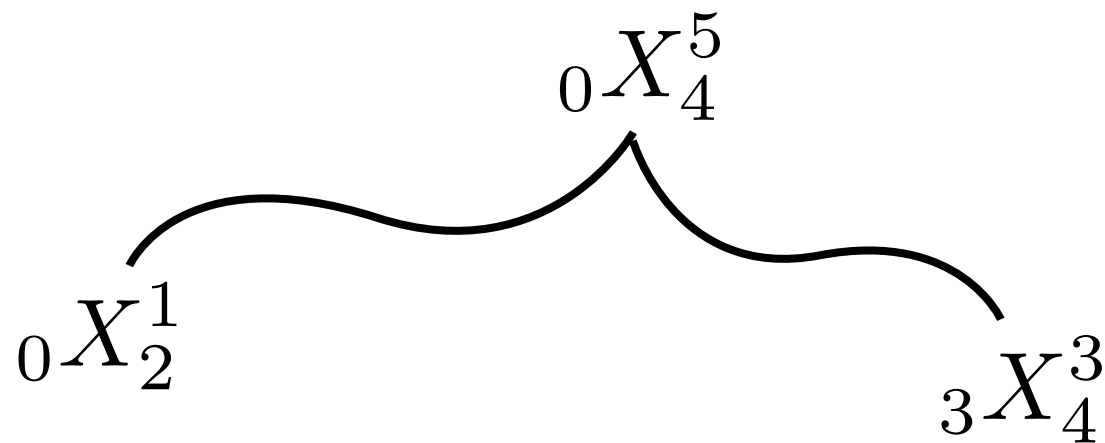
# Language Models are Important

# Language Models are Important

Impact on size of language model training data (in words) on quality of Arabic-English statistical machine translation system

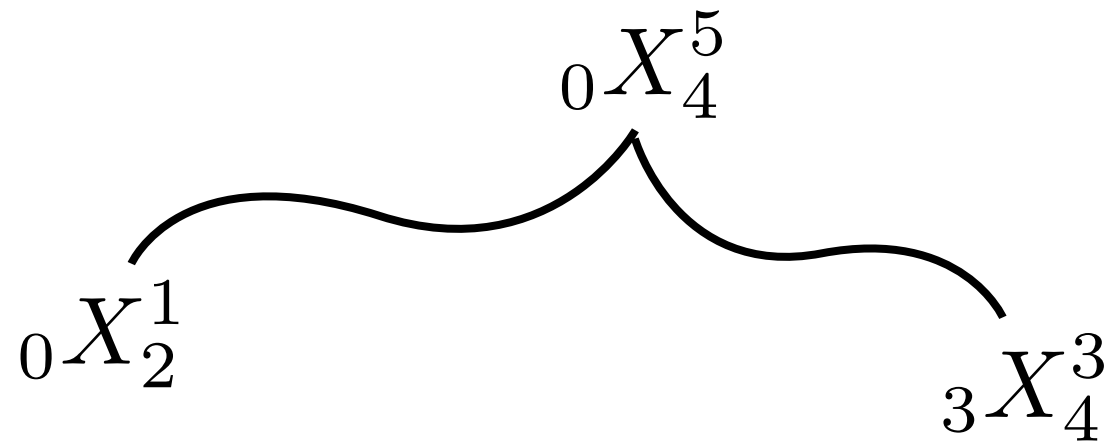


# +LM Dynamic Programming



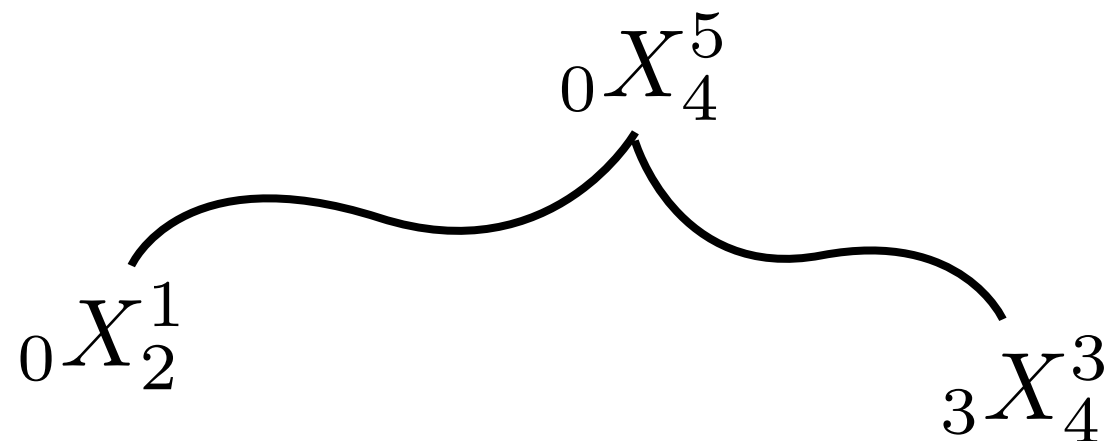


# +LM Dynamic Programming



the ... mat  
a ... mat  
mat ... mat

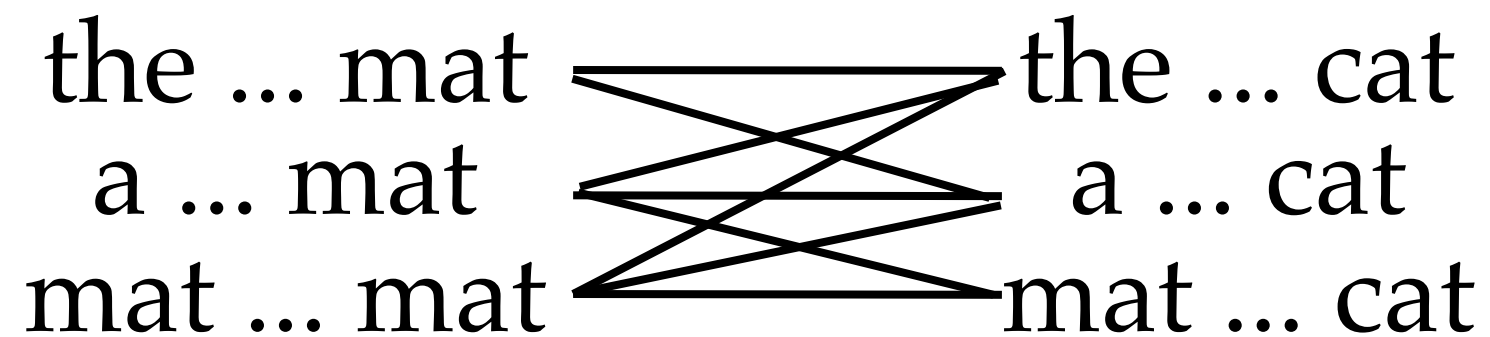
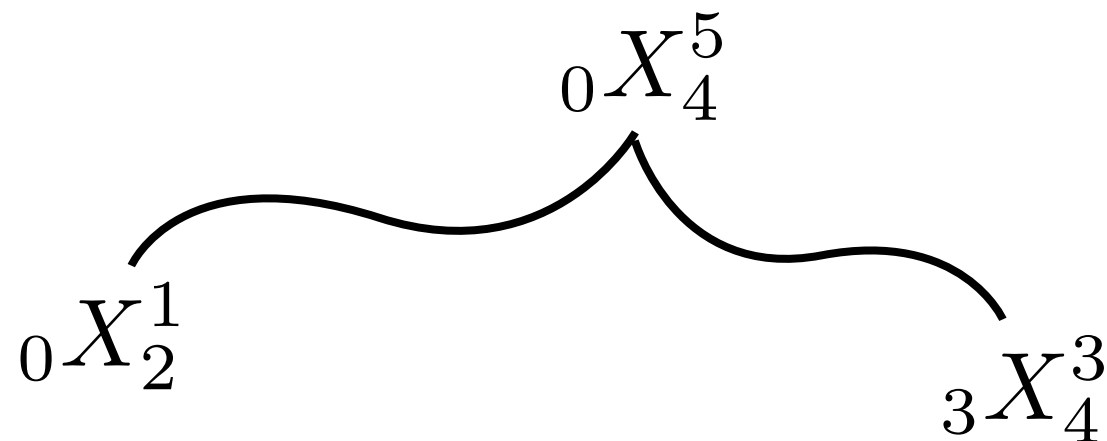
# +LM Dynamic Programming



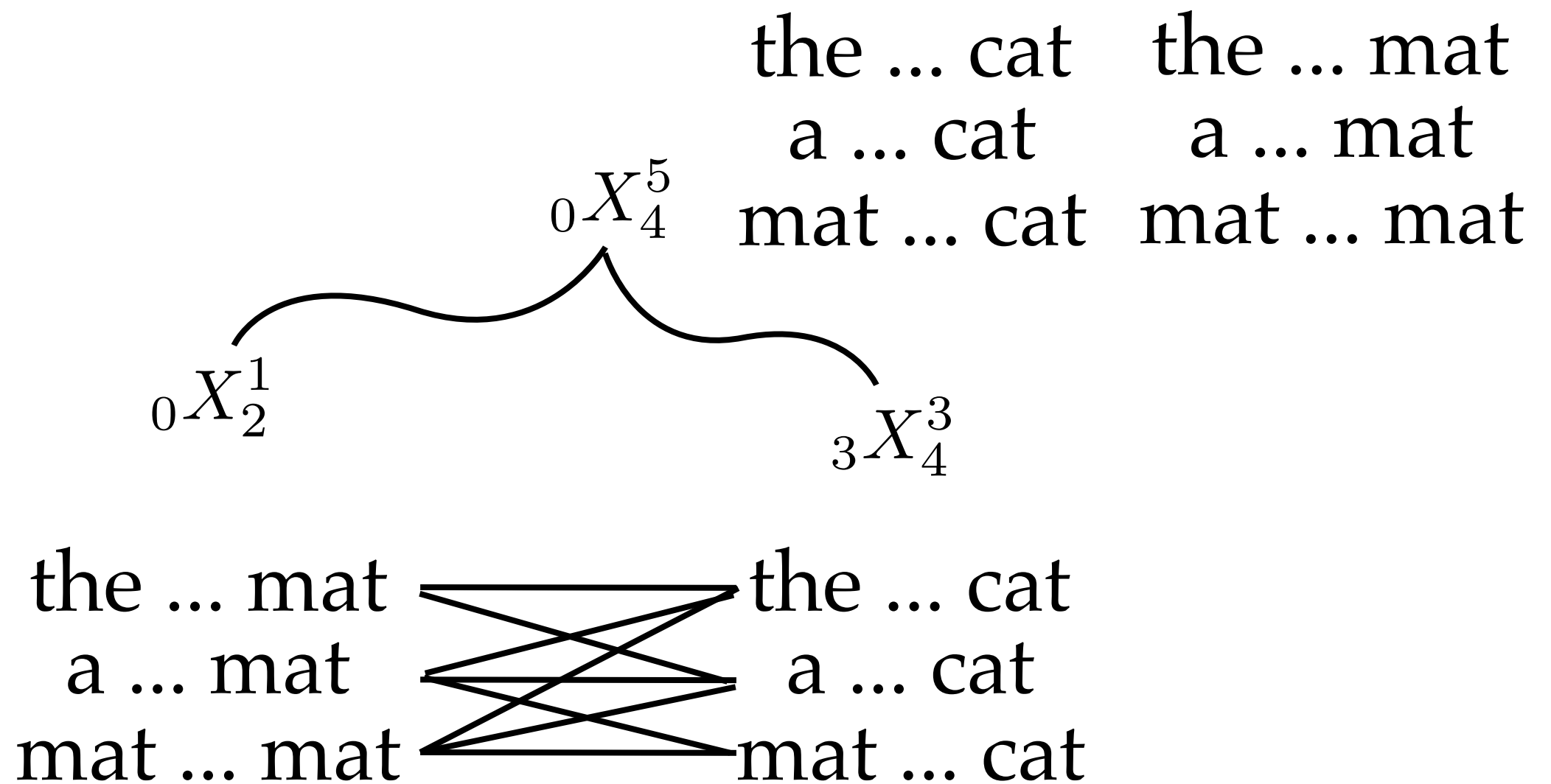
the ... mat  
a ... mat  
mat ... mat

the ... cat  
a ... cat  
mat ... cat

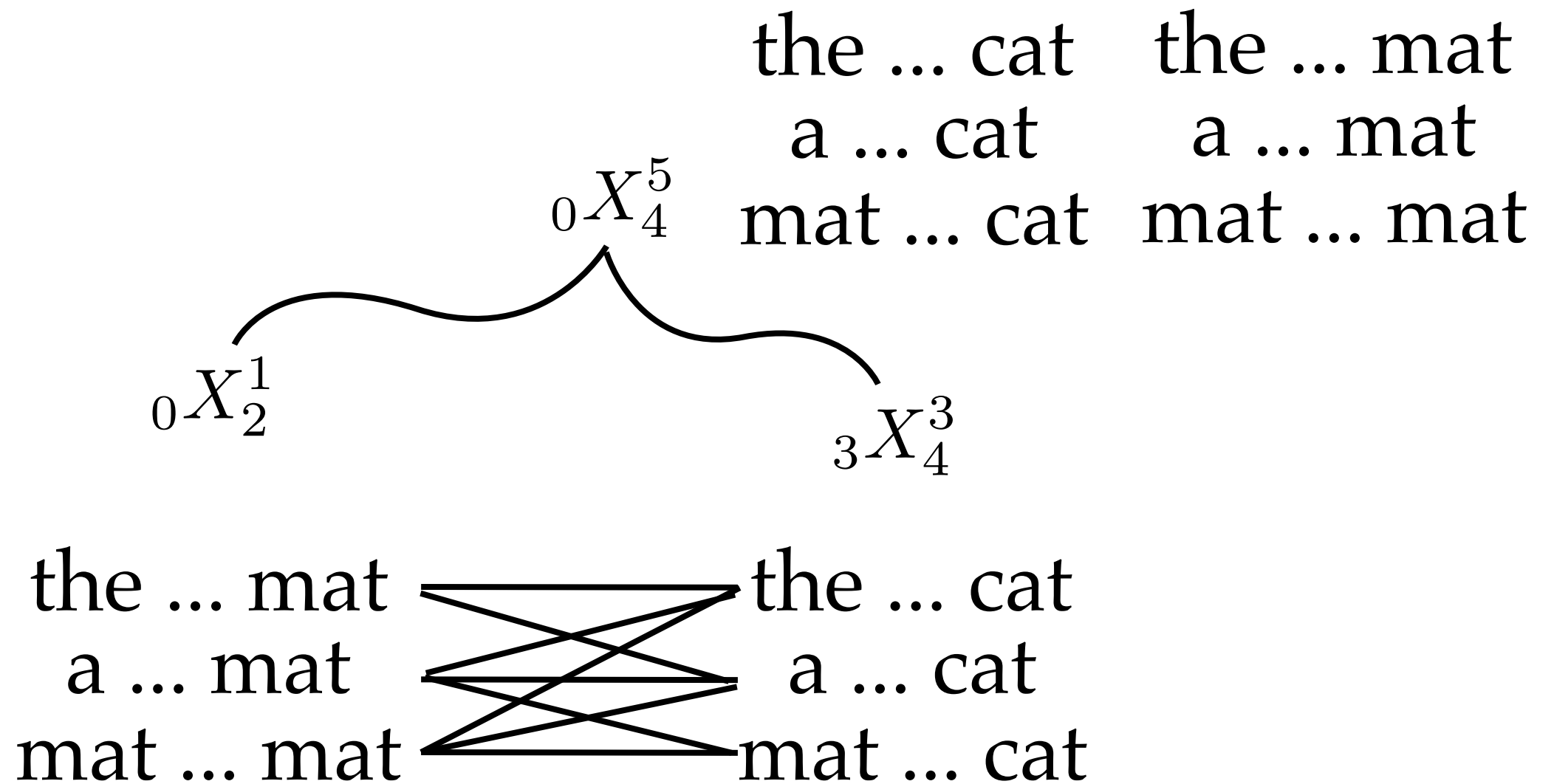
# +LM Dynamic Programming



# +LM Dynamic Programming



# +LM Dynamic Programming



Item (node) form:  $X_{i,j,q,r}$

# Cube Pruning Summary

- Parse Source
  - Result: -LM Hypergraph
- Incorporate n-grams bottom up,  
pruning +LM items along the way
  - Result: +LM Hypergraph

# Experimental Sandbox

- Urdu 25 category grammar
- cdec decoder (Dyer et al., ACL 2010)
- <http://www.cdec-decoder.org>
- <http://code.google.com/p/ws10smt>
- <http://github.com/alopez/cdec>

# cdec

- Why it's awesome:
  - Supports multiple models: linear chain CRF, SCFG, phrase-based
  - Generic hypergraph algorithms
  - Implements baseline: cube pruning



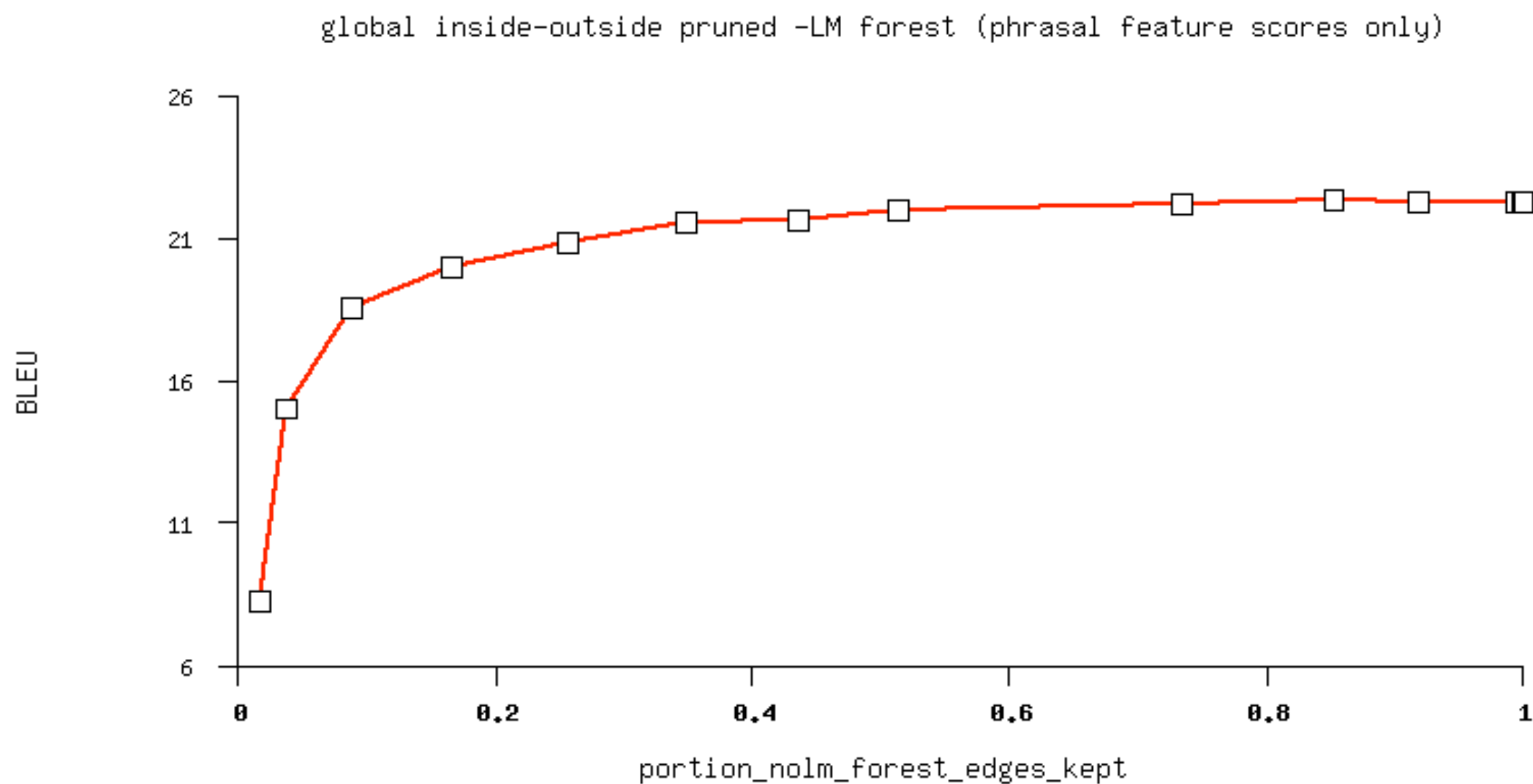
# The Size of Source Forests

- 1 Category (baseline)
  - Edges per sentence: 188,954
  - Decode time per sentence: 3.0 seconds
- 25 Categories
  - Edges per sentence: 1,242,410
  - Decode time per sentence: 52 seconds

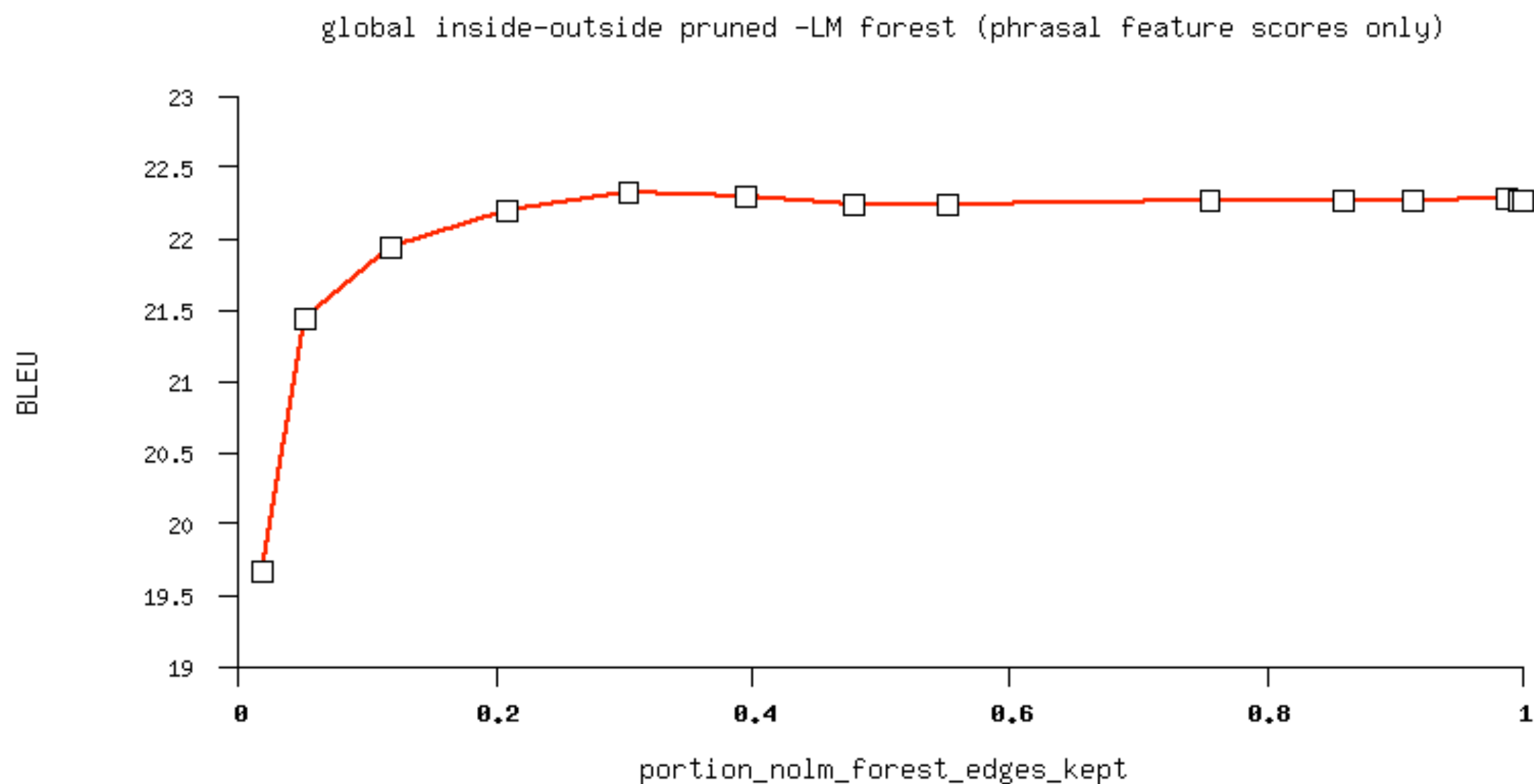
# Oracle: Does -LM Pruning Help?

- Generate *unpruned* -LM graph
- Prune using inside-outside
- Cube pruning to obtain +LM graph

# Oracle: Source Forest Pruning



# Oracle: *Tuned* Source Forest Pruning



# Coarse-to-fine Parsing

$X^1 \rightarrow dianzi\ shang/the\ mat$

$X^2 \rightarrow dianzi\ shang/mat$

$X^3 \rightarrow mao/the\ cat$

$X^4 \rightarrow X^1\ de\ X^3/X^3\ on\ X^1$

$X^4 \rightarrow X^2\ de\ X^3/X^3\ of\ X^2$

$X^5 \rightarrow X^1\ de\ X^3/X^1\ 's\ X^3$

$S \rightarrow X^4/X^4$

$S \rightarrow X^5/X^5$

# Coarse-to-fine Parsing

$X^1 \rightarrow \text{dianzi shang/the mat}$

$X^2 \rightarrow \text{dianzi shang/mat}$

$X^3 \rightarrow \text{mao/the cat}$

$X^4 \rightarrow X^1 \text{ de } X^3 / X^3 \text{ on } X^1$

$X^4 \rightarrow X^2 \text{ de } X^3 / X^3 \text{ of } X^2$

$X^5 \rightarrow X^1 \text{ de } X^3 / X^1 \text{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$

$X \rightarrow \text{dianzi shang/the ma}$

$X \rightarrow \text{dianzi shang/mat}$

$X \rightarrow \text{mao/the cat}$

$X \rightarrow X \text{ de } X / X \text{ on } X$

$X \rightarrow X \text{ de } X / X \text{ of } X$

$X \rightarrow X \text{ de } X / X \text{ 's } X$

$S \rightarrow X / X$

# Coarse-to-fine Parsing

$X^1 \rightarrow \text{dianzi shang/the mat}$

$X^2 \rightarrow \text{dianzi shang/mat}$

$X^3 \rightarrow \text{mao/the cat}$

$X^4 \rightarrow X^1 \text{ de } X^3 / X^3 \text{ on } X^1$

$X^4 \rightarrow X^2 \text{ de } X^3 / X^3 \text{ of } X^2$

$X^5 \rightarrow X^1 \text{ de } X^3 / X^1 \text{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$

$X \rightarrow \text{dianzi shang/the ma}$

$X \rightarrow \text{dianzi shang/mat}$

$X \rightarrow \text{mao/the cat}$

$X \rightarrow X \text{ de } X / X \text{ on } X$

$X \rightarrow X \text{ de } X / X \text{ of } X$

$X \rightarrow X \text{ de } X / X \text{ 's } X$

$S \rightarrow X / X$

Good news: Grammar constant shrinks

# Coarse-to-fine Parsing

$X^1 \rightarrow \text{dianzi shang/the mat}$

$X^2 \rightarrow \text{dianzi shang/mat}$

$X^3 \rightarrow \text{mao/the cat}$

$X^4 \rightarrow X^1 \text{ de } X^3 / X^3 \text{ on } X^1$

$X^4 \rightarrow X^2 \text{ de } X^3 / X^3 \text{ of } X^2$

$X^5 \rightarrow X^1 \text{ de } X^3 / X^1 \text{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$

$X \rightarrow \text{dianzi shang/the ma}$

$X \rightarrow \text{dianzi shang/mat}$

$X \rightarrow \text{mao/the cat}$

$X \rightarrow X \text{ de } X / X \text{ on } X$

$X \rightarrow X \text{ de } X / X \text{ of } X$

$X \rightarrow X \text{ de } X / X \text{ 's } X$

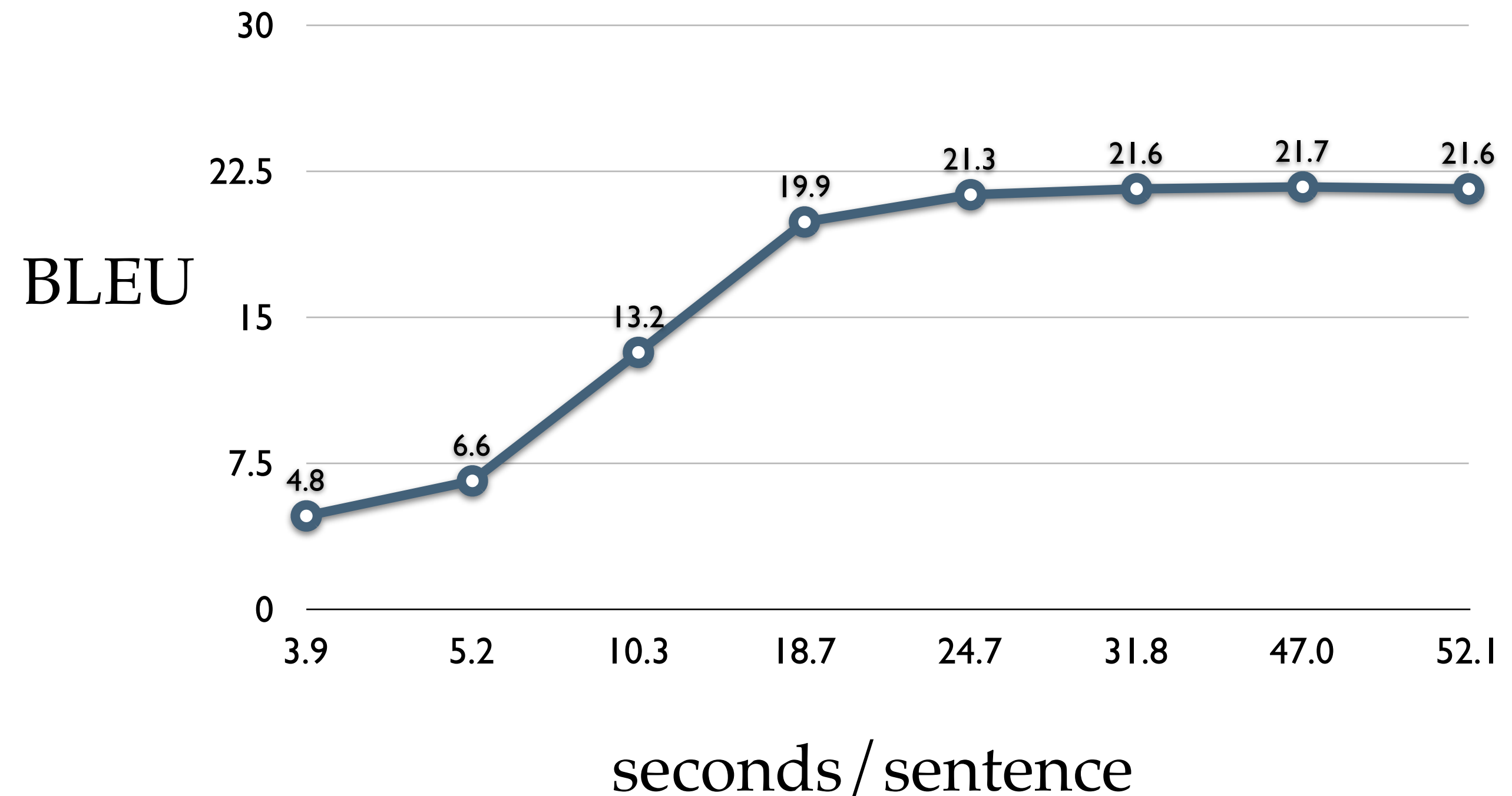
$S \rightarrow X / X$

Good news: Grammar constant shrinks

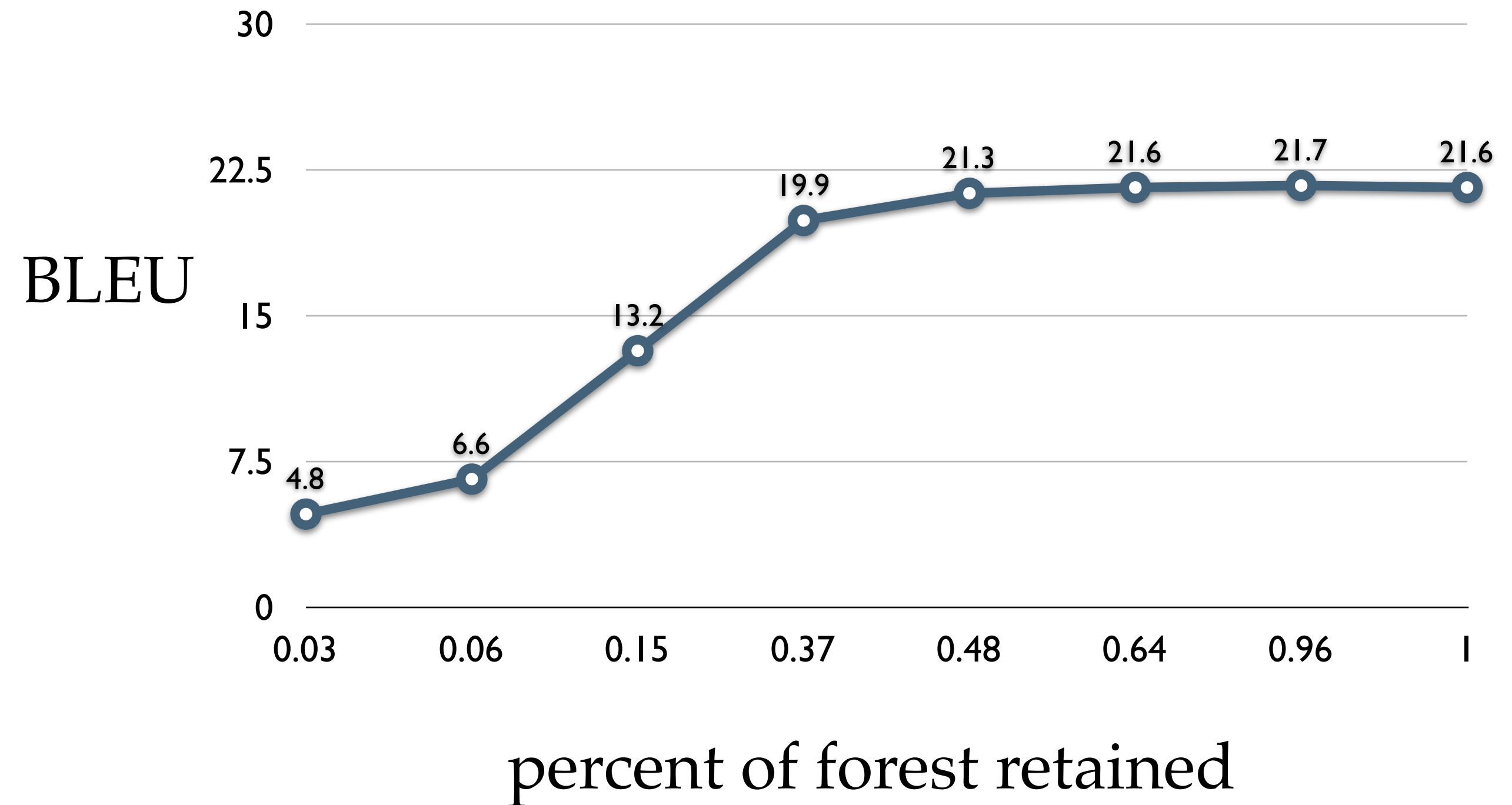
Bad news: Can sometimes prune away all parses



# Coarse-to-fine Parsing



# Coarse-to-fine Parsing



# Grammar Pruning

$X^1 \rightarrow \textit{dianzi shang/the mat}$

$X^2 \rightarrow \textit{dianzi shang/mat}$

$X^3 \rightarrow \textit{mao/the cat}$

$X^4 \rightarrow X^1 \textit{ de } X^3 / X^3 \textit{ on } X^1$

$X^4 \rightarrow X^2 \textit{ de } X^3 / X^3 \textit{ of } X^2$

$X^5 \rightarrow X^1 \textit{ de } X^3 / X^1 \textit{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$

# Grammar Pruning

$X^1 \rightarrow \text{dianzi shang/the mat}$

$X^2 \rightarrow \text{dianzi shang/mat}$

$X^3 \rightarrow \text{mao/the cat}$

$X^4 \rightarrow X^1 \text{ de } X^3 / X^3 \text{ on } X^1$

$X^4 \rightarrow X^2 \text{ de } X^3 / X^3 \text{ of } X^2$

$X^5 \rightarrow X^1 \text{ de } X^3 / X^1 \text{ 's } X^3$

$S \rightarrow X^4 / X^4$

$S \rightarrow X^5 / X^5$

$X^1 \rightarrow \text{dianzi shang/the mat}$

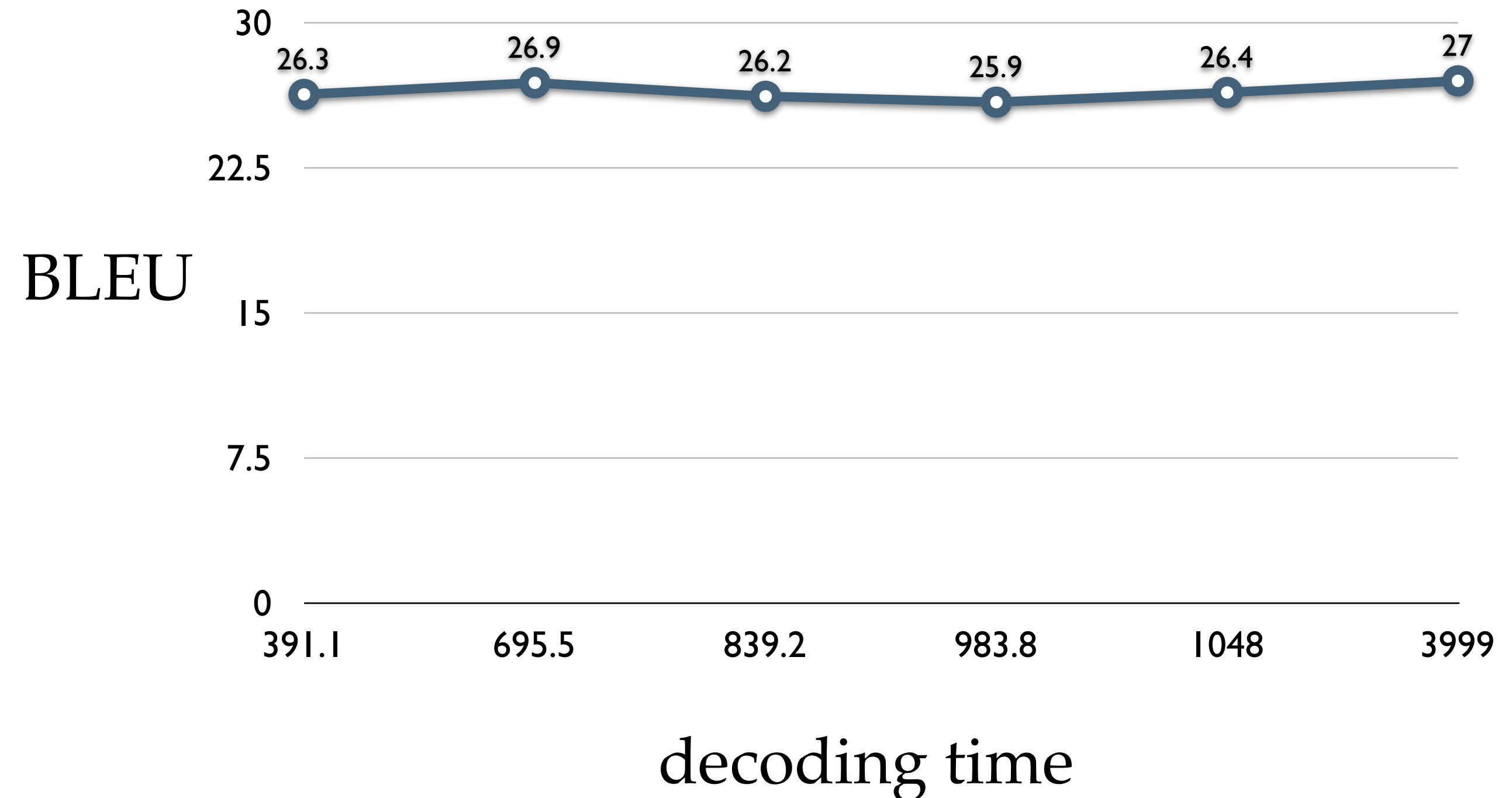
$X^3 \rightarrow \text{mao/the cat}$

$X^4 \rightarrow X^1 \text{ de } X^3 / X^3 \text{ on } X^1$

$S \rightarrow X^4 / X^4$

# Grammar pruning

(note: different data conditions)



# Current Results Summary

# Current Results Summary

- Source parse forest pruning works.

# Current Results Summary

- Source parse forest pruning works.
- Can reduce overall decoding time by 40% with *unpruned* grammars.



# Current Results Summary

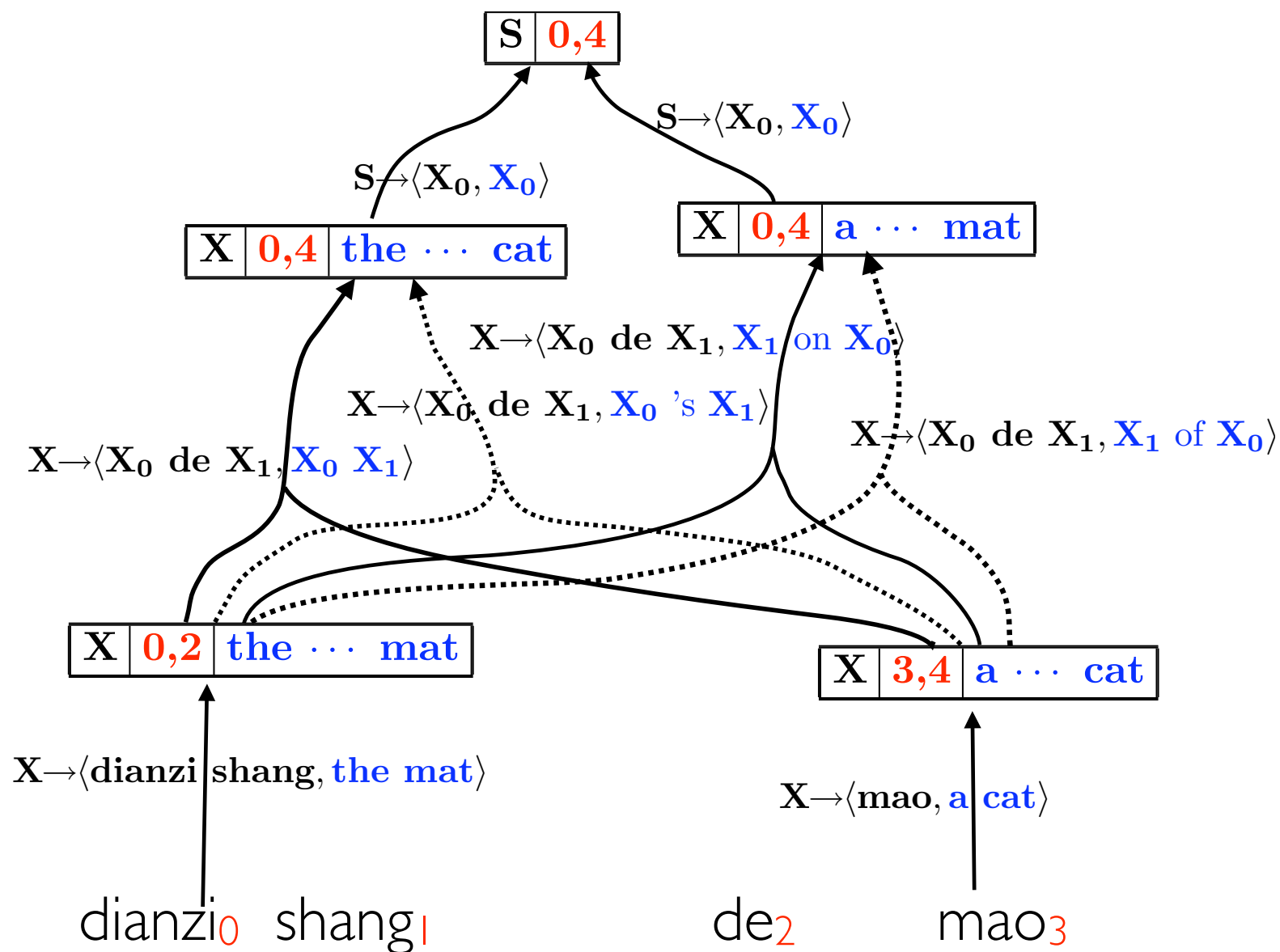
- Source parse forest pruning works.
- Can reduce overall decoding time by 40% with *unpruned* grammars.
- Can reduce decoding time by an order of magnitude with pruned grammars (*maybe* at some cost in BLEU).

# Current Results Summary

- Source parse forest pruning works.
- Can reduce overall decoding time by 40% with *unpruned* grammars.
- Can reduce decoding time by an order of magnitude with pruned grammars (*maybe* at some cost in BLEU).
- Ongoing work on more interesting algorithms...

# Parse Forests as Grammars

Input:  $\langle \text{dianzi shiang de mao}, \text{a cat on the mat} \rangle$

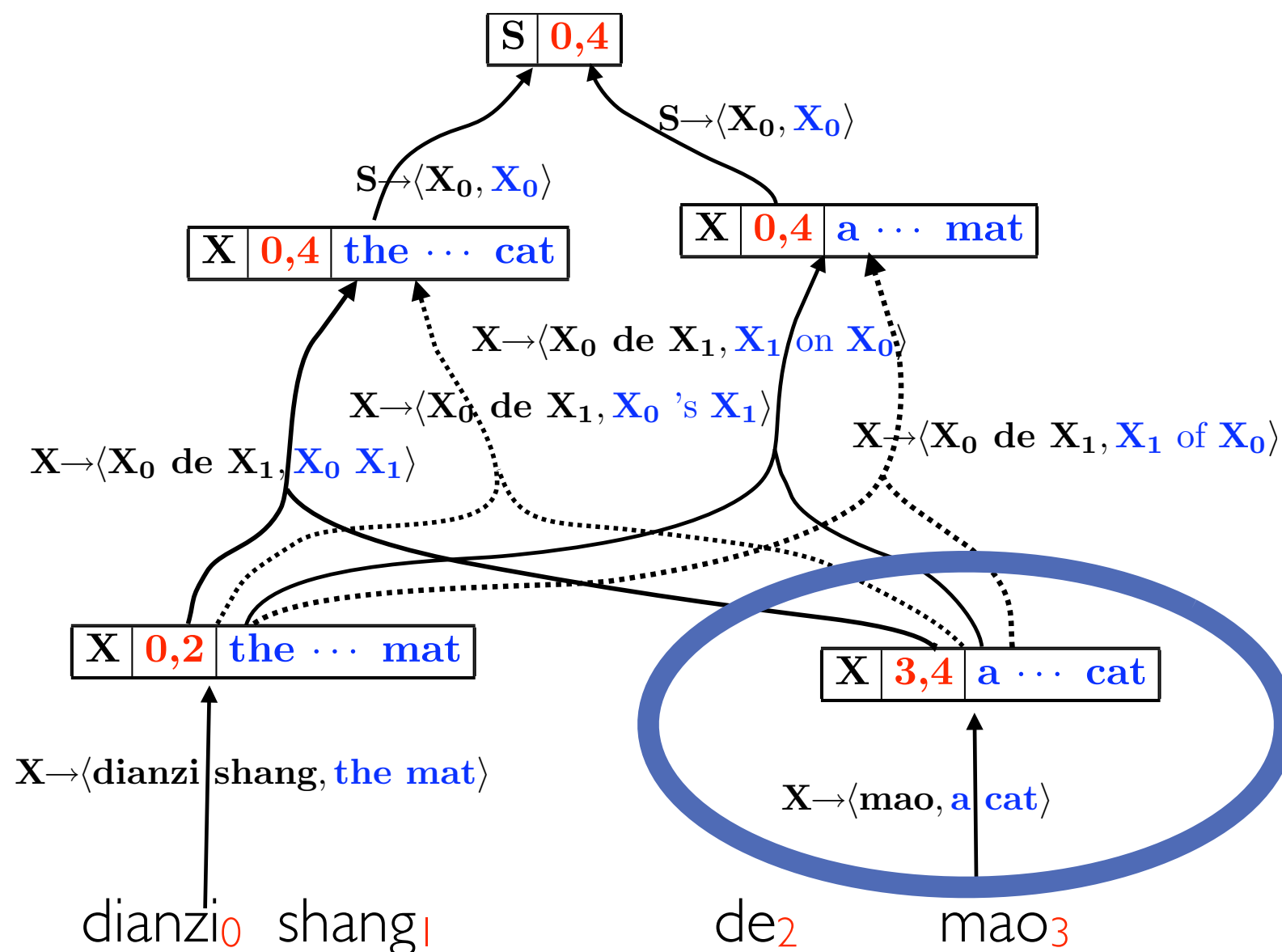


# Parse Forests as Grammars

Input:  $\langle \text{dianzi shiang de mao} , \text{a cat on the mat} \rangle$

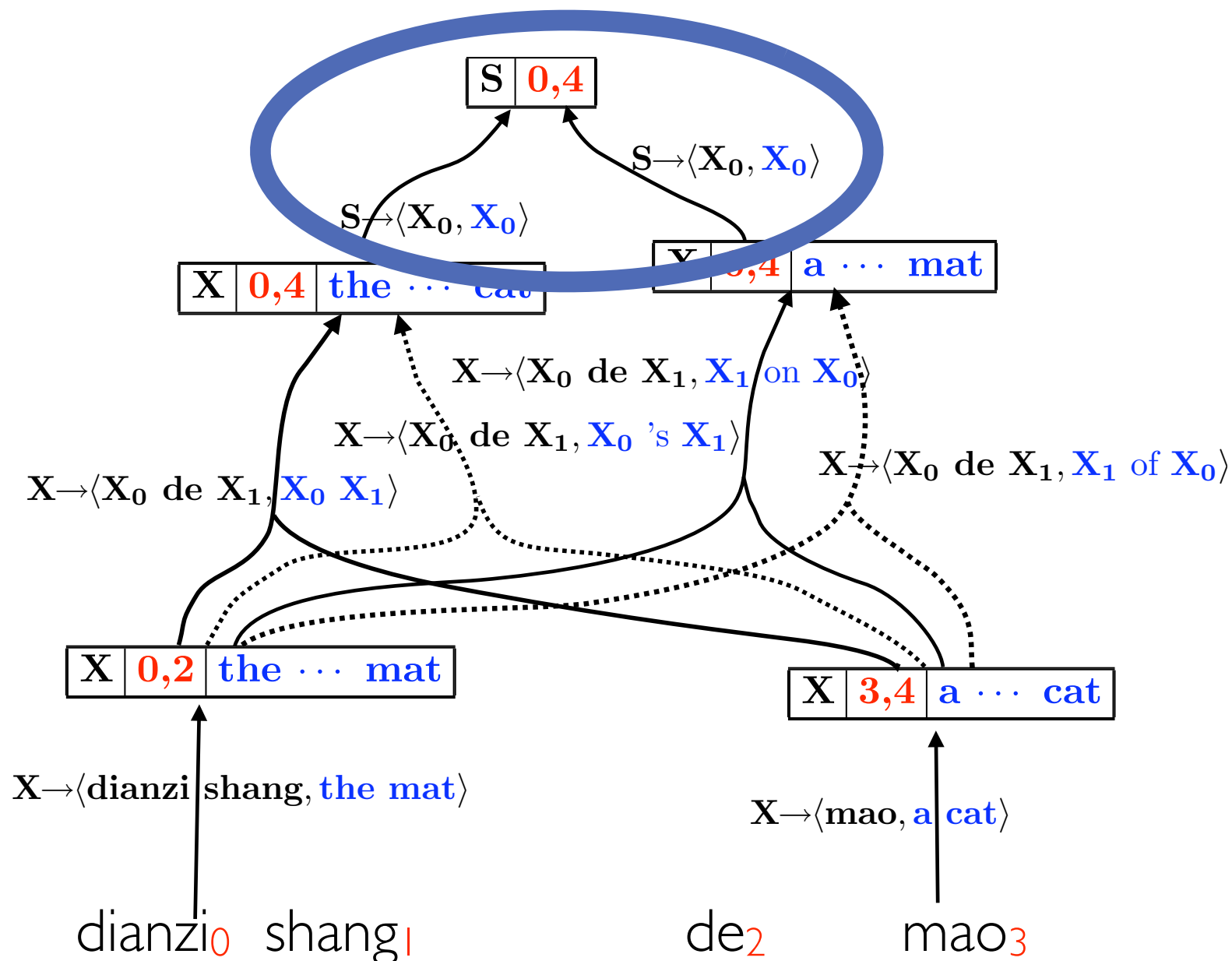
**Isomorphic CFG**

$[X34] \rightarrow \text{a cat}$



# Parse Forests as Grammars

Input:  $\langle \text{dianzi shiang de mao}, \text{a cat on the mat} \rangle$



## Isomorphic CFG

$[X34] \rightarrow \text{a cat}$

$[X02] \rightarrow \text{the mat}$

$[X04a] \rightarrow [X34] \text{ on } [X02]$

$[X04a] \rightarrow [X34] \text{ of } [X02]$

$[X04b] \rightarrow [X02] \text{'s } [X34]$

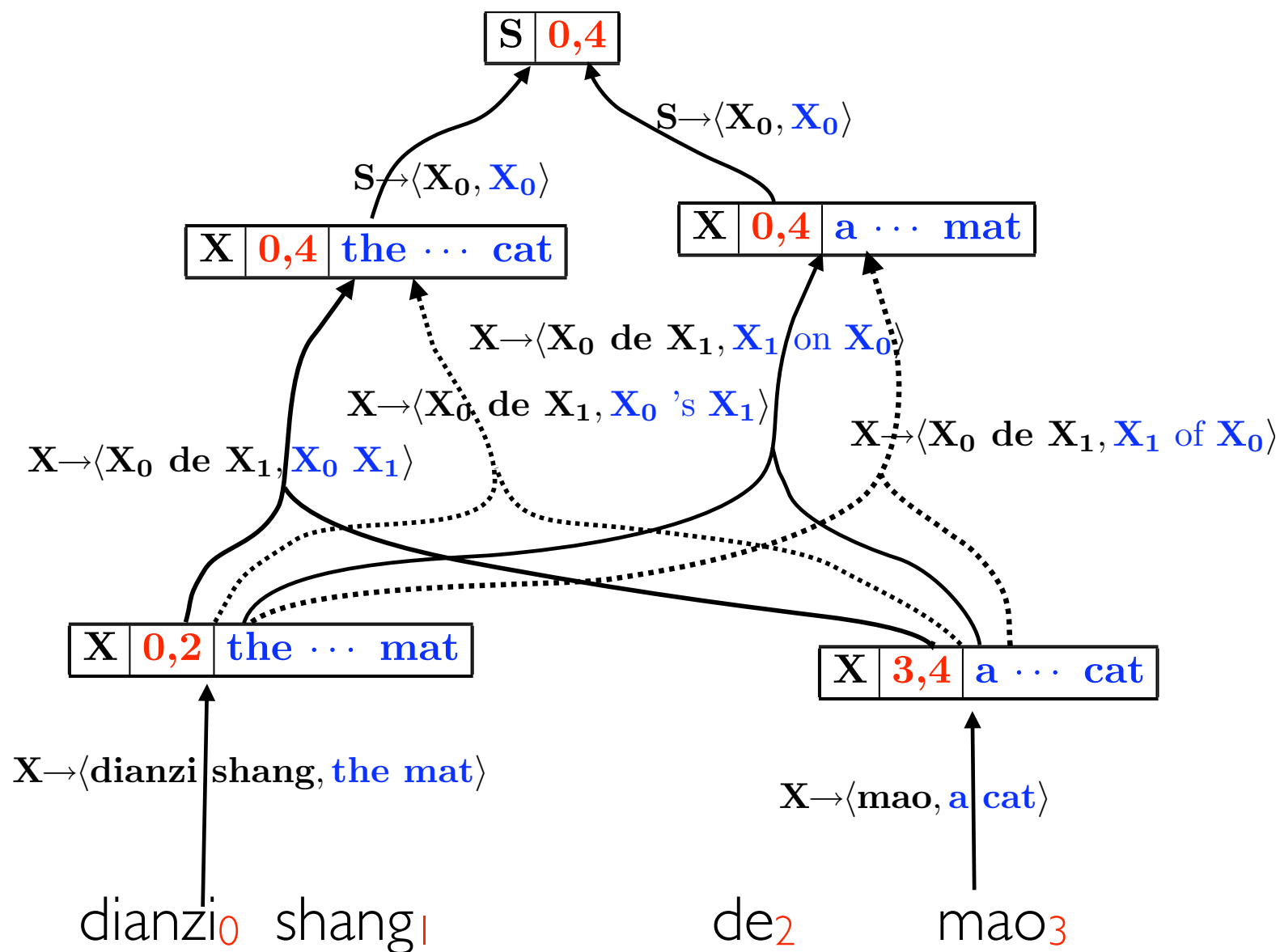
$[X04b] \rightarrow [X02] [X34]$

$[S] \rightarrow [X04a]$

$[S] \rightarrow [X04b]$

# Parse Forests as Grammars

Input:  $\langle \text{dianzi shiang de mao}, \text{a cat on the mat} \rangle$

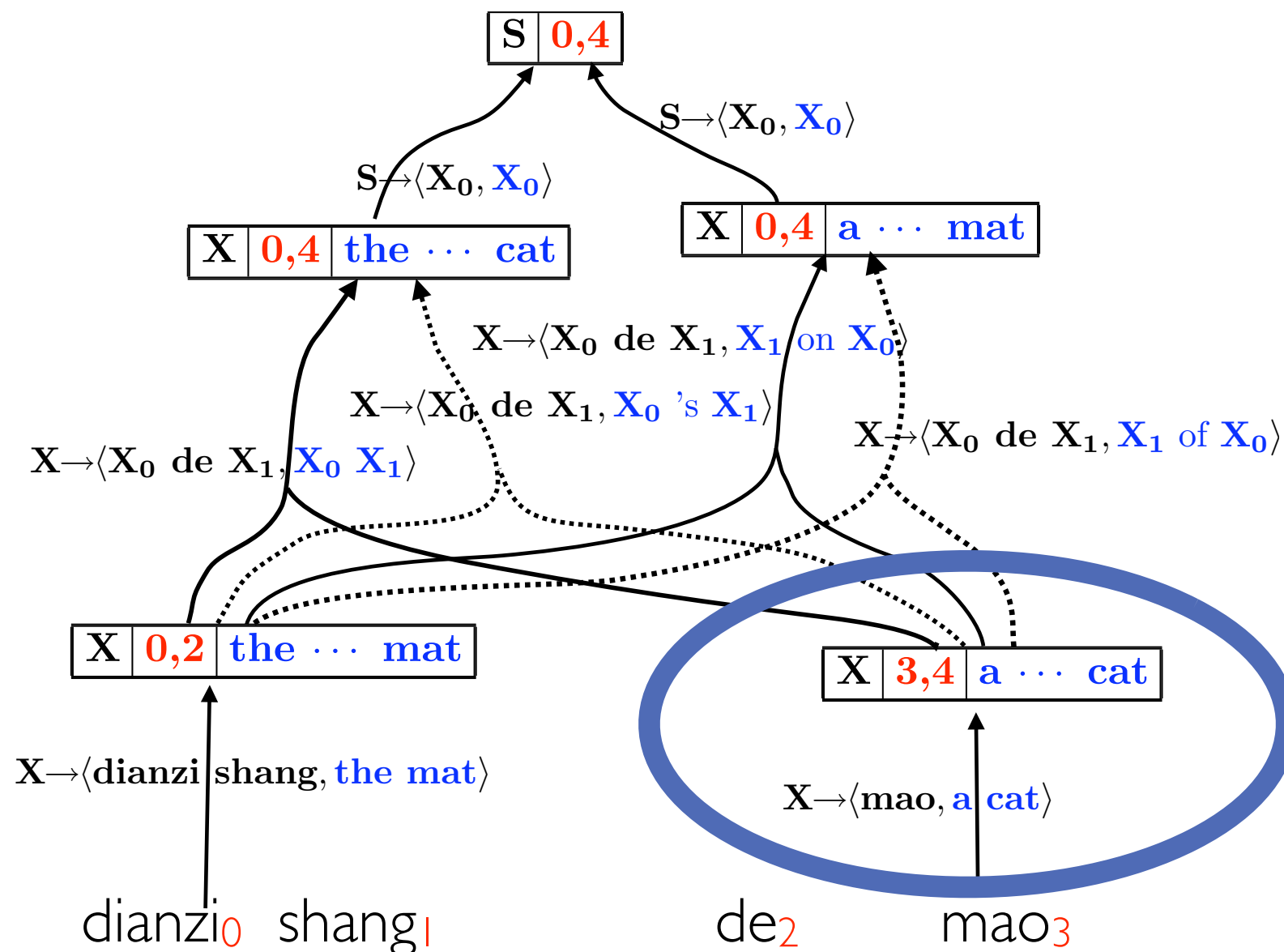


# Parse Forests as Grammars

Input:  $\langle \text{dianzi shiang de mao}, \text{a cat on the mat} \rangle$

**Isomorphic CFG**

$[X34] \rightarrow \text{a cat}$



# Translation is Intersection

- Translation by parsing is intersection
- Intersect source with grammar
- Yields a parse forest  $\rightarrow$  target grammar
- Generate with target grammar
- Intersect with target language model  
(regular language)



# Generation from Grammars

# Generation from Grammars

## ***Isomorphic CFG***

[X34] → *a cat*

[X02] → *the mat*

[X04a] → [X34] *on* [X02]

[X04a] → [X34] *of* [X02]

[X04b] → [X02] *'s* [X34]

[X04b] → [X02] [X34]

[S] → [X04a]

[S] → [X04b]

# Generation from Grammars

## ***Isomorphic CFG***

[X34]  $\rightarrow$  *a cat*

[X02]  $\rightarrow$  *the mat*

[X04a]  $\rightarrow$  [X34] *on* [X02]

[X04a]  $\rightarrow$  [X34] *of* [X02]

[X04b]  $\rightarrow$  [X02] *'s* [X34]

[X04b]  $\rightarrow$  [X02] [X34]

[S]  $\rightarrow$  [X04a]

[S]  $\rightarrow$  [X04b]

$$S_{0,4,\langle s \rangle,\langle s \rangle} \longrightarrow \bullet X_{0,4}^a$$

# Generation from Grammars

## ***Isomorphic CFG***

[X34]  $\rightarrow$  *a cat*

[X02]  $\rightarrow$  *the mat*

[X04a]  $\rightarrow$  [X34] *on* [X02]

[X04a]  $\rightarrow$  [X34] *of* [X02]


[X04b]  $\rightarrow$  [X02] *'s* [X34]

[X04b]  $\rightarrow$  [X02] [X34]

[S]  $\rightarrow$  [X04a]

[S]  $\rightarrow$  [X04b]

$X_{0,4,\langle s \rangle, \langle s \rangle}^a \rightarrow \bullet X_{3,4} \text{ on } X_{0,2}$

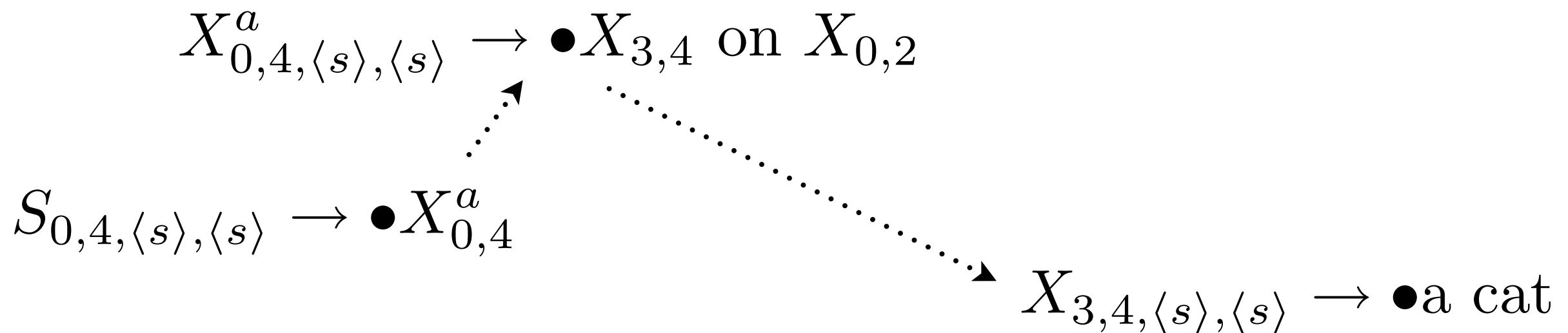


$S_{0,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet X_{0,4}^a$

# Generation from Grammars

## ***Isomorphic CFG***

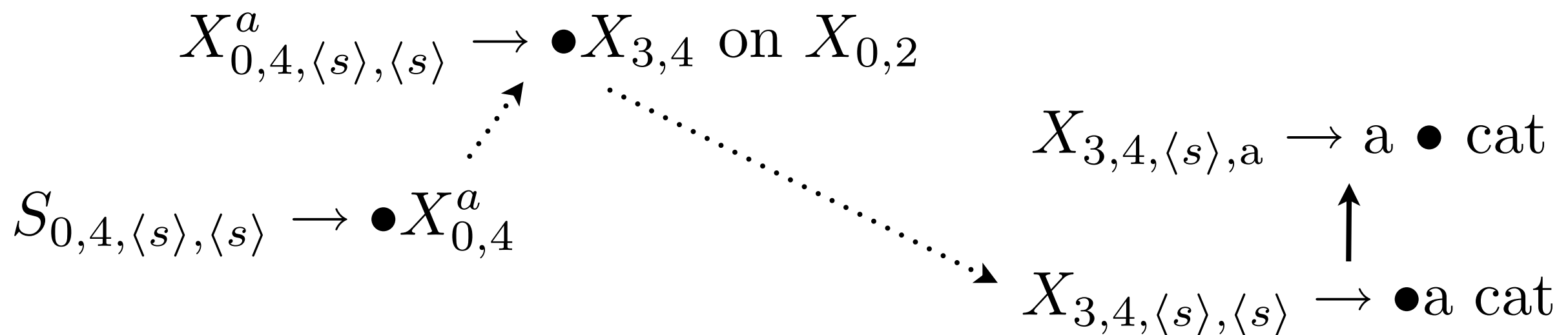
[X34]  $\rightarrow$  *a cat*  
[X02]  $\rightarrow$  *the mat*  
[X04a]  $\rightarrow$  [X34] *on* [X02]  
[X04a]  $\rightarrow$  [X34] *of* [X02]  
[X04b]  $\rightarrow$  [X02] *'s* [X34]  
[X04b]  $\rightarrow$  [X02] [X34]  
[S]  $\rightarrow$  [X04a]  
[S]  $\rightarrow$  [X04b]



# Generation from Grammars

## ***Isomorphic CFG***

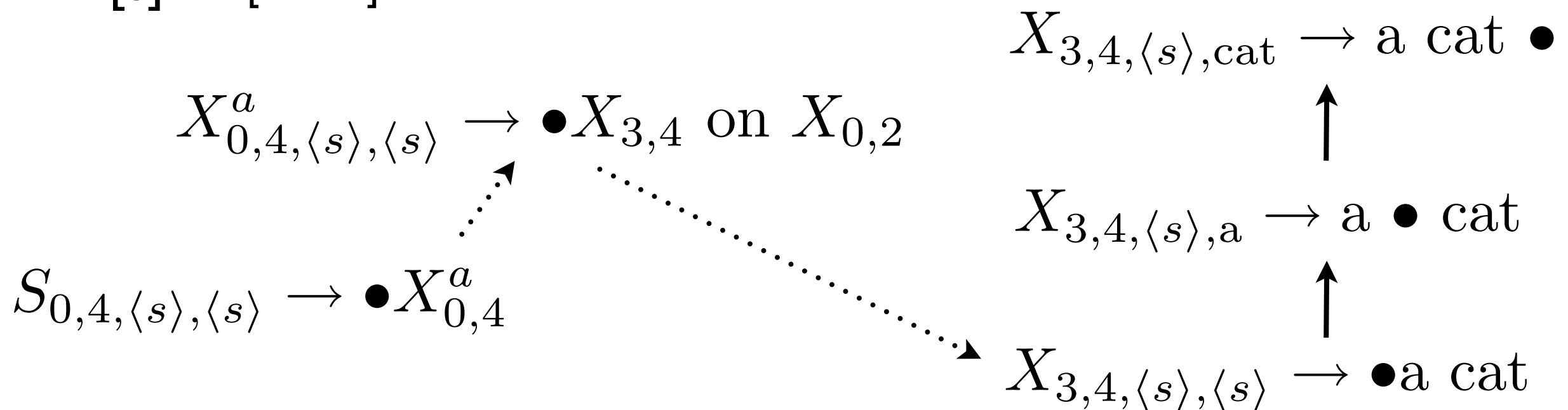
$[X34] \rightarrow a \text{ cat}$   
 $[X02] \rightarrow the \text{ mat}$   
 $[X04a] \rightarrow [X34] \text{ on } [X02]$   
 $[X04a] \rightarrow [X34] \text{ of } [X02]$   
 $[X04b] \rightarrow [X02] 's [X34]$   
 $[X04b] \rightarrow [X02] [X34]$   
 $[S] \rightarrow [X04a]$   
 $[S] \rightarrow [X04b]$



# Generation from Grammars

## ***Isomorphic CFG***

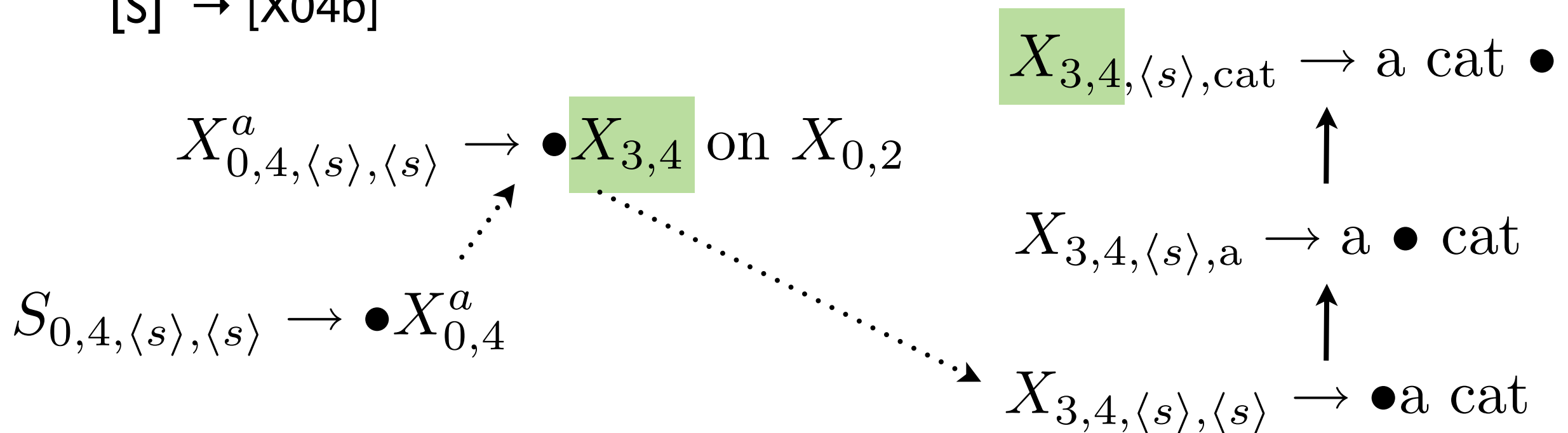
$[X34] \rightarrow a \text{ cat}$   
 $[X02] \rightarrow the \text{ mat}$   
 $[X04a] \rightarrow [X34] \text{ on } [X02]$   
 $[X04a] \rightarrow [X34] \text{ of } [X02]$   
 $[X04b] \rightarrow [X02] 's [X34]$   
 $[X04b] \rightarrow [X02] [X34]$   
 $[S] \rightarrow [X04a]$   
 $[S] \rightarrow [X04b]$



# Generation from Grammars

## ***Isomorphic CFG***

$[X34] \rightarrow a \text{ cat}$   
 $[X02] \rightarrow the \text{ mat}$   
 $[X04a] \rightarrow [X34] \text{ on } [X02]$   
 $[X04a] \rightarrow [X34] \text{ of } [X02]$   
 $[X04b] \rightarrow [X02] 's [X34]$   
 $[X04b] \rightarrow [X02] [X34]$   
 $[S] \rightarrow [X04a]$   
 $[S] \rightarrow [X04b]$

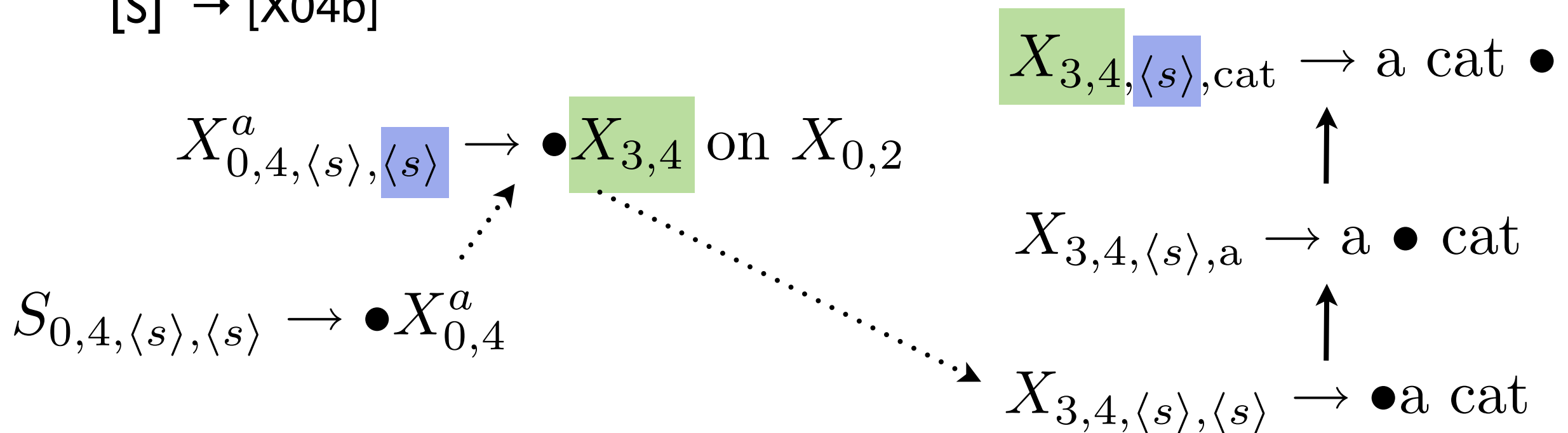




# Generation from Grammars

## ***Isomorphic CFG***

$[X34] \rightarrow a \text{ cat}$   
 $[X02] \rightarrow the \text{ mat}$   
 $[X04a] \rightarrow [X34] on [X02]$   
 $[X04a] \rightarrow [X34] of [X02]$   
 $[X04b] \rightarrow [X02] 's [X34]$   
 $[X04b] \rightarrow [X02] [X34]$   
 $[S] \rightarrow [X04a]$   
 $[S] \rightarrow [X04b]$



# Generation from Grammars

## Isomorphic CFG

[X34]  $\rightarrow$  *a cat*

[X02]  $\rightarrow$  *the mat*

[X04a]  $\rightarrow$  [X34] *on* [X02]

[X04a]  $\rightarrow$  [X34] *of* [X02]

[X04b]  $\rightarrow$  [X02] *'s* [X34]

[X04b]  $\rightarrow$  [X02] [X34]

[S]  $\rightarrow$  [X04a]

[S]  $\rightarrow$  [X04b]

$$X_{0,4,\langle s \rangle, \text{cat}}^a \rightarrow X_{3,4} \bullet \text{ on } X_{0,2}$$

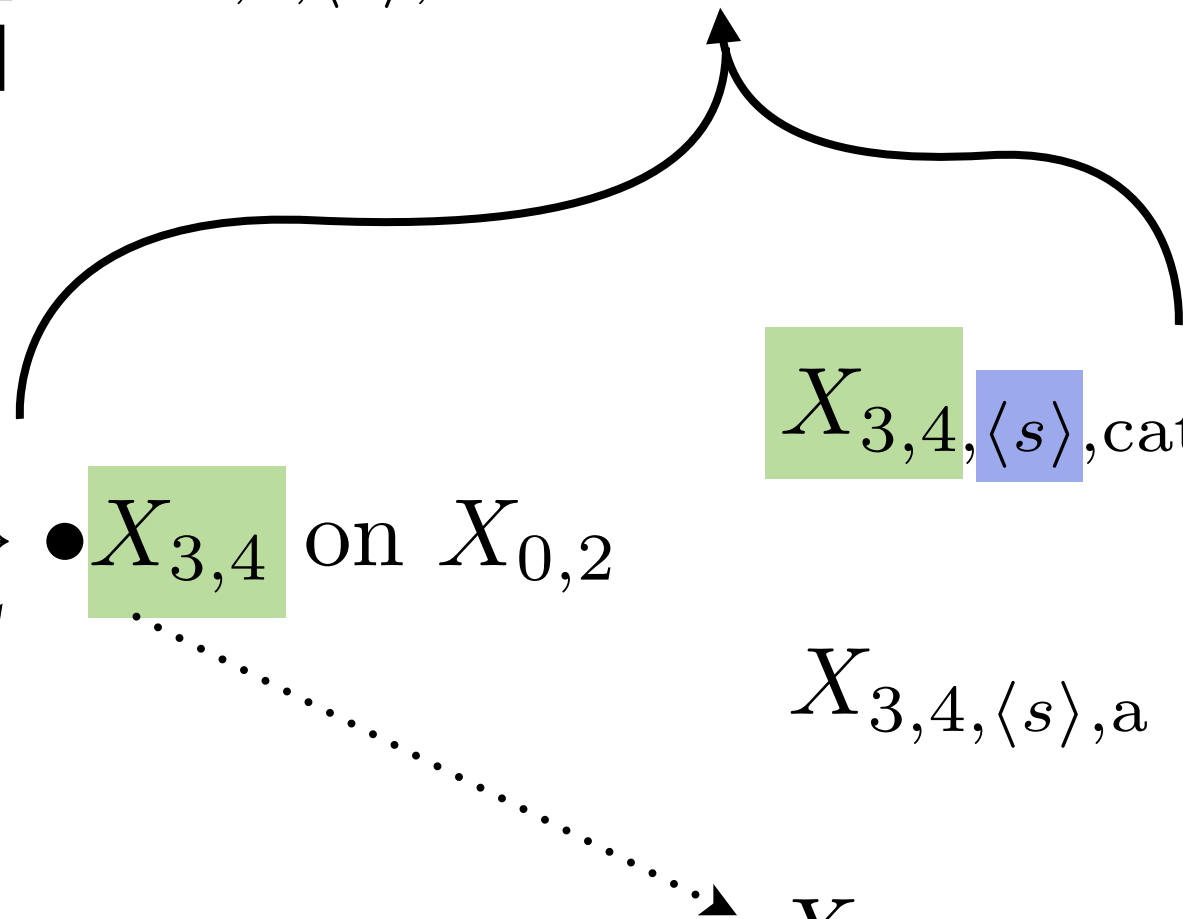
$$X_{0,4,\langle s \rangle, \langle s \rangle}^a \rightarrow \bullet X_{3,4} \text{ on } X_{0,2}$$

$$X_{3,4,\langle s \rangle, \text{cat}} \rightarrow \text{a cat} \bullet$$

$$X_{3,4,\langle s \rangle, \text{a}} \rightarrow \text{a} \bullet \text{ cat}$$

$$X_{3,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet \text{a cat}$$

$$S_{0,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet X_{0,4}^a$$



# Generation from Grammars

## Isomorphic CFG

[X34]  $\rightarrow$  *a cat*

[X02]  $\rightarrow$  *the mat*

[X04a]  $\rightarrow$  [X34] *on* [X02]

[X04a]  $\rightarrow$  [X34] *of* [X02]

[X04b]  $\rightarrow$  [X02] *'s* [X34]

[X04b]  $\rightarrow$  [X02] [X34]

[S]  $\rightarrow$  [X04a]

[S]  $\rightarrow$  [X04b]

$X_{0,4,\langle s \rangle}^a$   $\rightarrow$   $X_{3,4} \bullet$  on  $X_{0,2}$

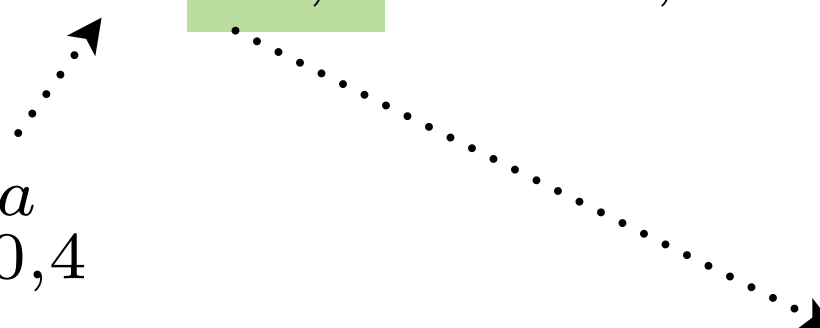
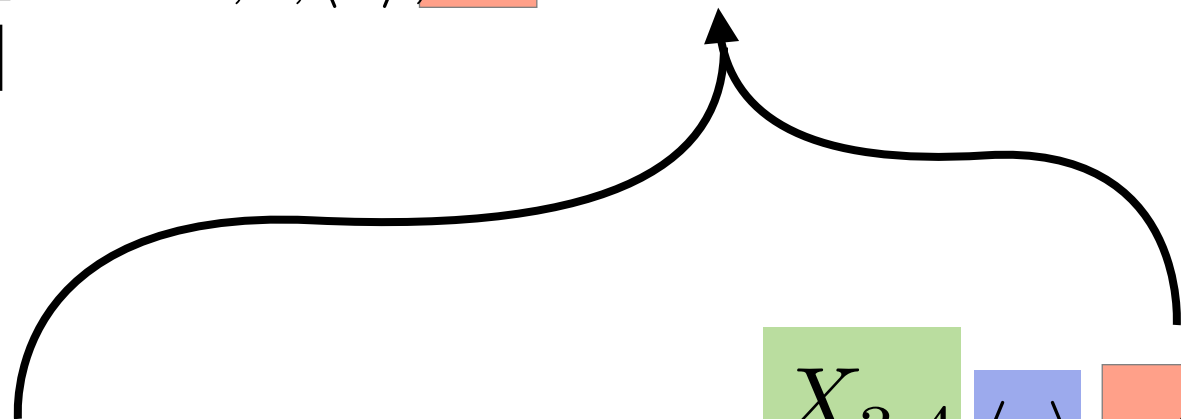
$X_{0,4,\langle s \rangle}^a$   $\rightarrow$   $\bullet X_{3,4}$  on  $X_{0,2}$

$X_{3,4,\langle s \rangle}$   $\rightarrow$  a cat  $\bullet$

$X_{3,4,\langle s \rangle,a}$   $\rightarrow$  a  $\bullet$  cat

$X_{3,4,\langle s \rangle,\langle s \rangle}$   $\rightarrow$   $\bullet$  a cat

$S_{0,4,\langle s \rangle,\langle s \rangle}$   $\rightarrow$   $\bullet X_{0,4}^a$



# Generation from Grammars

## Isomorphic CFG

[X34]  $\rightarrow$  *a cat*

[X02]  $\rightarrow$  *the mat*

[X04a]  $\rightarrow$  [X34] *on* [X02]

[X04a]  $\rightarrow$  [X34] *of* [X02]

[X04b]  $\rightarrow$  [X02] *'s* [X34]

[X04b]  $\rightarrow$  [X02] [X34]

[S]  $\rightarrow$  [X04a]

[S]  $\rightarrow$  [X04b]

$X_{0,4,\langle s \rangle, \text{on}} \rightarrow X_{3,4} \text{ on} \bullet X_{0,2}$

$X_{0,4,\langle s \rangle, \text{cat}}^a \rightarrow X_{3,4} \bullet \text{on} X_{0,2}$

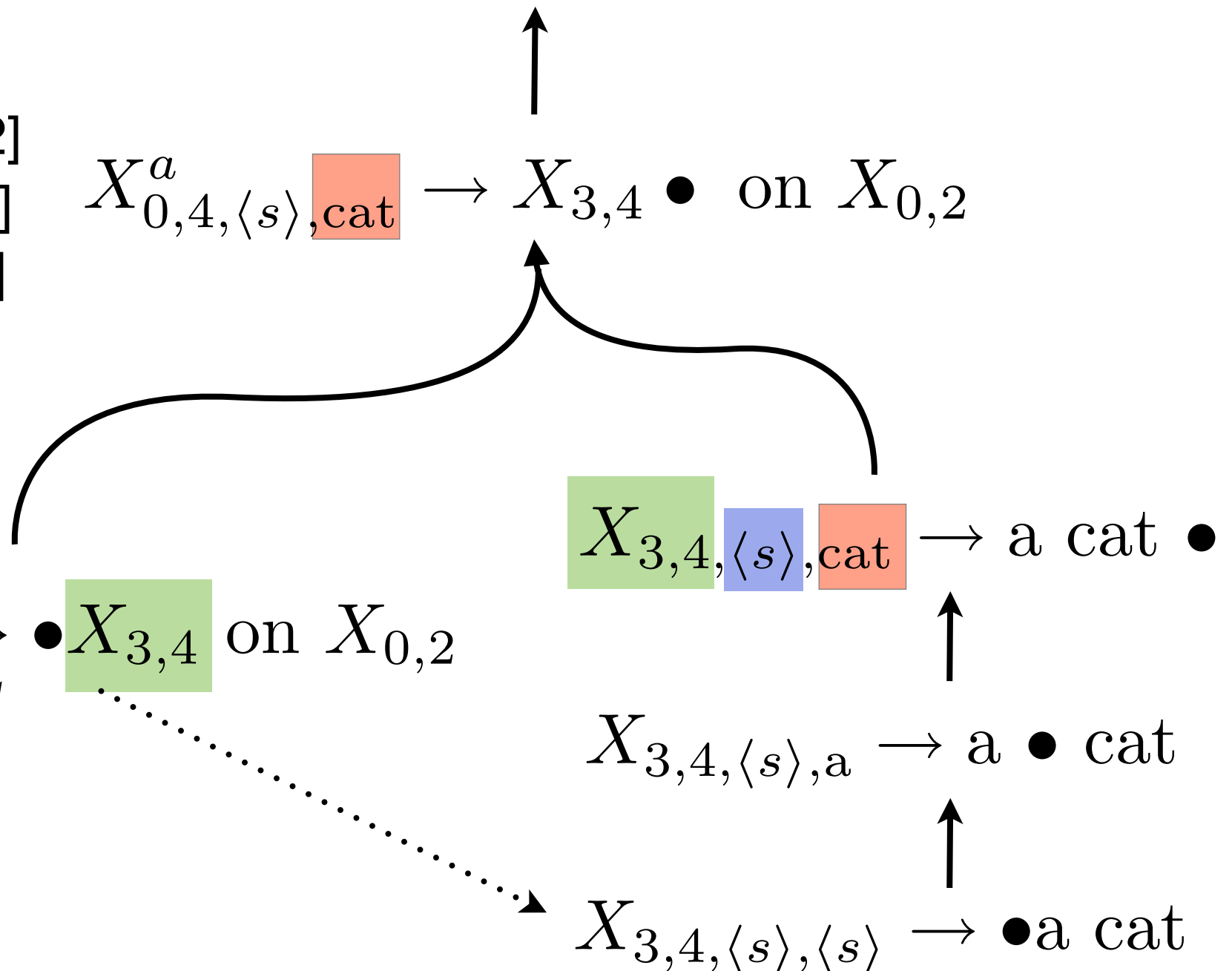
$X_{0,4,\langle s \rangle, \langle s \rangle}^a \rightarrow \bullet X_{3,4} \text{ on} X_{0,2}$

$X_{3,4,\langle s \rangle, \text{cat}} \rightarrow \text{a cat} \bullet$

$X_{3,4,\langle s \rangle, \text{a}} \rightarrow \text{a} \bullet \text{cat}$

$X_{3,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet \text{a cat}$

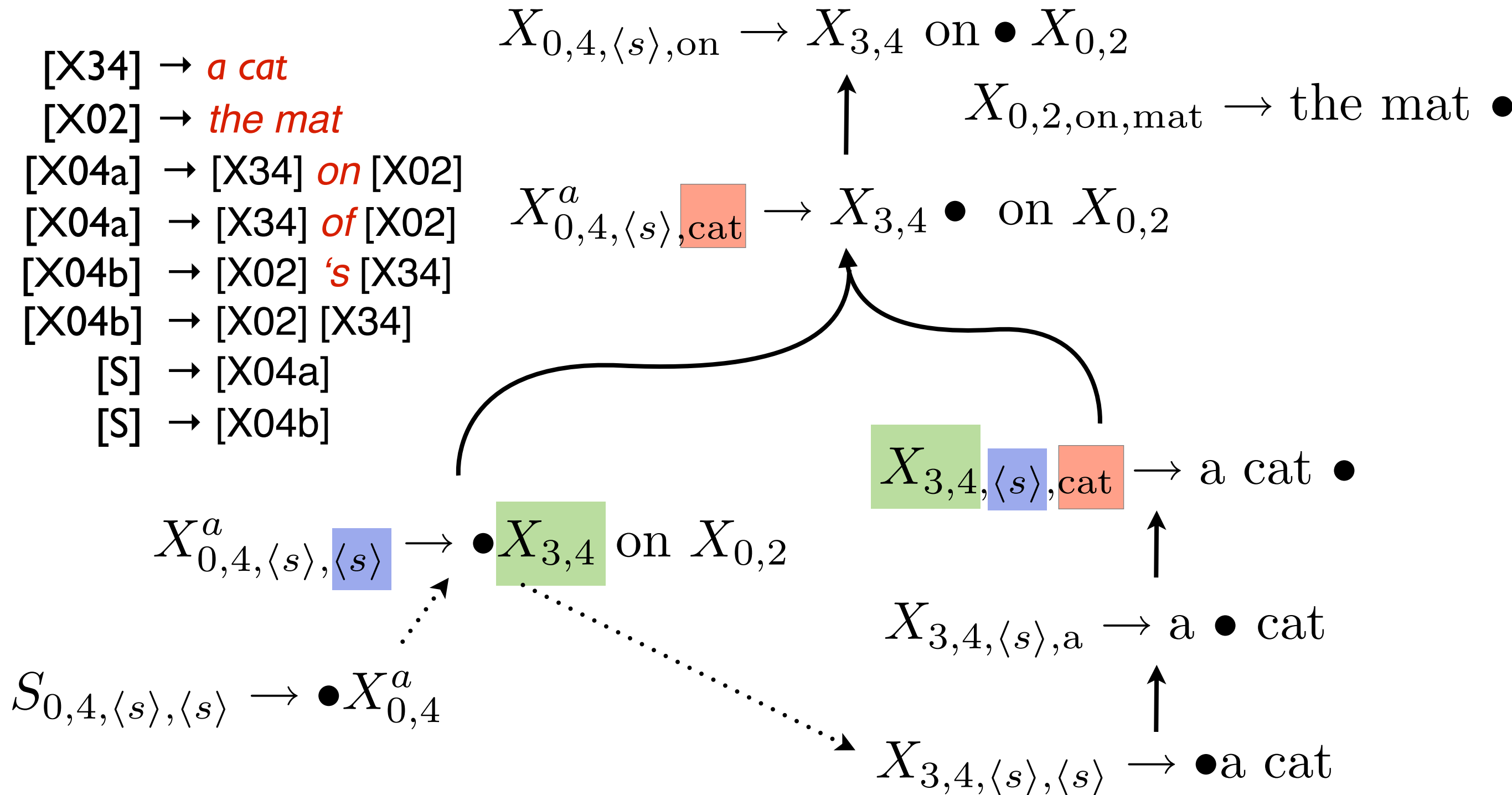
$S_{0,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet X_{0,4}^a$



# Generation from Grammars

## Isomorphic CFG

$[X34] \rightarrow a \text{ cat}$   
 $[X02] \rightarrow the \text{ mat}$   
 $[X04a] \rightarrow [X34] \text{ on } [X02]$   
 $[X04a] \rightarrow [X34] \text{ of } [X02]$   
 $[X04b] \rightarrow [X02] 's [X34]$   
 $[X04b] \rightarrow [X02] [X34]$   
 $[S] \rightarrow [X04a]$   
 $[S] \rightarrow [X04b]$



# Generation from Grammars

# Isomorphic CFG

**[X34]**  $\rightarrow$  *a cat*

**[X02]** → *the mat*

**[X04a] → [X34] *on* [X02]**

**[X04a]** → **[X34]** *of* **[X02]**

**[X04b] → [X02] 's [X34]**

$$[X04b] \rightarrow [X02] [X34]$$

**[S] → [X04a]**

**[S] → [X04b]**

$$X_{0,4,\langle s \rangle, \text{mat}} \rightarrow X_{3,4} \text{ on } X_{0,2} \bullet$$

$$X_{0,4,\langle s \rangle,\text{on}} \longrightarrow X_{3,4} \text{ on } \bullet X_{0,2}$$

$$X_{0,2,\text{on},\text{mat}} \rightarrow \text{the mat } \bullet$$

$$X_{0,4,\langle s \rangle, \text{cat}}^a \rightarrow X_{3,4} \bullet \text{ on } X_{0,2}$$

$$X_{0,4,\langle s \rangle, \langle s \rangle}^a \rightarrow \bullet X_{3,4} \text{ on } X_{0,2}$$

$$X_{3,4,\langle s \rangle, \text{cat}} \rightarrow \text{a cat} \bullet$$

$$X_{3,4,\langle s \rangle, \mathbf{a}} \rightarrow \mathbf{a} \bullet \text{cat}$$

$$X_{3,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet \text{a cat}$$

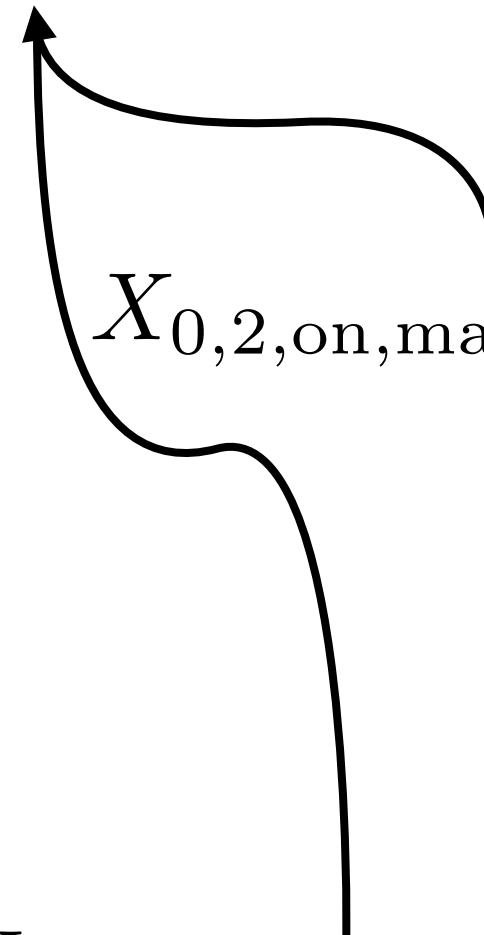
$$S_{0,4,\langle s \rangle, \langle s \rangle} \rightarrow \bullet X_{0,4}^a$$

# Generation from Grammars

$X_{0,4,\langle s \rangle, \text{mat}} \rightarrow X_{3,4} \text{ on } X_{0,2}$

$X_{0,2, \text{on}, \text{mat}} \rightarrow \text{the mat}$

$X_{3,4,\langle s \rangle, \text{cat}} \rightarrow \text{a cat}$



# Generation from Grammars

$$X_{0,4,\langle s \rangle, \text{mat}} \rightarrow X_{3,4} \text{ on } X_{0,2}$$

## ***Isomorphic CFG***

[X34]  $\rightarrow$  *a cat*

[X02]  $\rightarrow$  *the mat*

[X04a]  $\rightarrow$  [X34] *on* [X02]

[X04a]  $\rightarrow$  [X34] *of* [X02]

[X04b]  $\rightarrow$  [X02] *'s* [X34]

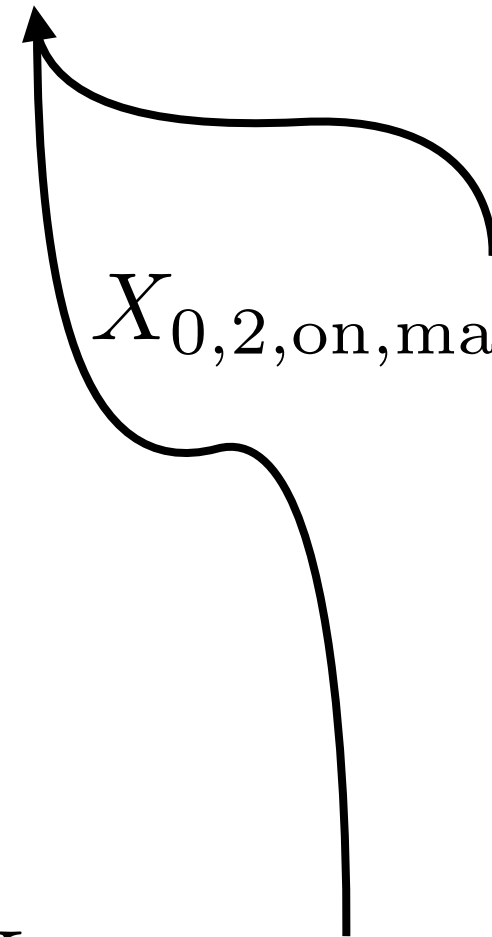
[X04b]  $\rightarrow$  [X02] [X34]

[S]  $\rightarrow$  [X04a]

[S]  $\rightarrow$  [X04b]

$X_{0,2,\text{on},\text{mat}} \rightarrow$  the mat

$X_{3,4,\langle s \rangle, \text{cat}} \rightarrow$  a cat





# Theoretical Outcomes

- Works for arbitrary grammars, not just binary grammars
- Asymptotically faster than cube pruning (“hook trick”, Liang et al. 2006).
- Produces lots of admissible heuristics ( $A^*$ )
- No cube pruning: everything is monotonic.

# Conclusions

- Faster algorithms are needed to make induced grammars practical.
- Workshop made significant progress towards this goal.
- More improvements are underway.

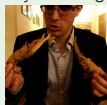
# Outline



Vlad Eidelman



Ziyuan Wang



Adam Lopez



Jon Graehl



ThuyLinh Nguyen

- 3:20pm Parametric models: posterior regularisation. Desai
- 3:35pm Training models with rich features spaces. Vlad
- 3:50pm Decoding with complex grammars. Adam
- 4:20pm Closing remarks. Phil
- 4:25pm Finish.

# Statistical machine translation: state-of-the-art

Urdu → English

اس حملہ کے بعد بڑی تعداد میں مقامی باشندوں نے علاقوں کو خالی کر دیا ہے .



In this attack a large number of local residents has should vacate areas.

- Current state-of-the-art translation models struggle with language pairs which exhibit large differences in structure.

# Statistical machine translation: our unsupervised grammars

Urdu → English

اس حملہ کے بعد بڑی تعداد میں مقامی باشندوں نے علاقوں کو خالی کر دیا ہے .



After this attack, a large number of local residents have to vacate the areas.

- In this workshop we've made some small steps towards better translations for difficult language pairs.

# Statistical machine translation: our unsupervised grammars

Urdu → English

اس حملہ کے بعد بڑی تعداد میں مقامی باشندوں نے علاقوں کو خالی کر دیا ہے .



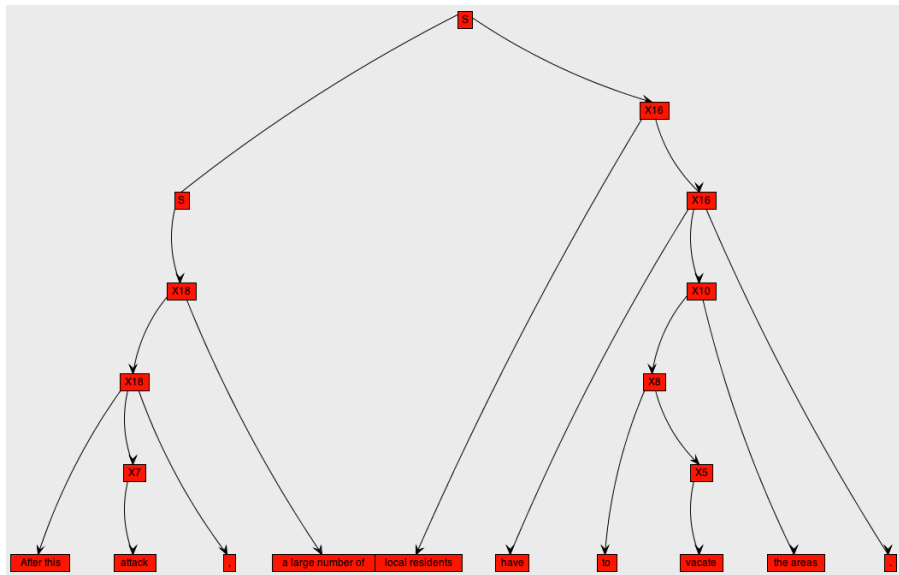
After this attack, a large number of local residents have to vacate the areas.

- In this workshop we've made some small steps towards better translations for difficult language pairs.

Google Translate:

*\*After the attack a number of local residents has blank areas.*

# Induced Translation Structure



# What we've achieved:

- The first unsupervised labelled SCFG induction algorithms:
  - ▶ by clustering translation phrases which occur in the same context we can learn which phrases are substituteable,
  - ▶ we have implemented parametric and non-parametric Bayesian clustering algorithms and shown positive results on real translation tasks.
- Improved SCFG decoders that efficiently decode grammars with many labels:
  - ▶ we have created faster search algorithms tuned for syntactic grammars.
- Discriminative training regimes to leverage features extracted from these grammars:
  - ▶ we've implemented two large scale discriminative algorithms for training our models.



# What we've achieved:

- The first unsupervised labelled SCFG induction algorithms:
  - ▶ by clustering translation phrases which occur in the same context we can learn which phrases are substituteable,
  - ▶ we have implemented parametric and non-parametric Bayesian clustering algorithms and shown positive results on real translation tasks.
- Improved SCFG decoders that efficiently decode grammars with many labels:
  - ▶ we have created faster search algorithms tuned for syntactic grammars.
- Discriminative training regimes to leverage features extracted from these grammars:
  - ▶ we've implemented two large scale discriminative algorithms for training our models.

Thank you.