

# **Variational Decoding for Statistical Machine Translation**

**Zhifei Li, Jason Eisner, and Sanjeev Khudanpur**

**Center for Language and Speech Processing**

**Computer Science Department**

**Johns Hopkins University**

# Spurious Ambiguity

- Statistical models in MT exhibit **spurious ambiguity**
  - Many different derivations (e.g., trees or segmentations) generate the same translation string
- Regular phrase-based MT systems
  - **phrase segmentation** ambiguity
- Tree-based MT systems
  - **derivation tree** ambiguity

# Spurious Ambiguity in Phrase Segmentations

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

machine

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

machine translation software

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

machine translation software

机器

machine

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

machine translation software

机器 翻译

machine translation

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

machine translation software

机器 翻译 软件

machine translation software

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

machine translation software

- Same output:  
“**machine translation software**”
- Three different phrase  
segmentations

机器 翻译 软件

machine translation software

# Spurious Ambiguity in Phrase Segmentations

机器 翻译 软件

machine translation software

machine translation software

机器 翻译 软件

- Same output:  
“**machine translation software**”
- Three different phrase  
segmentations

machine translation software

机器 翻译 软件

machine transfer software

machine translation software

# Spurious Ambiguity in Derivation Trees

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件

$S \rightarrow (\text{机器}, \text{machine})$

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件

S->(机器, machine) S->(翻译, translation)

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件

$S \rightarrow (\text{机器}, \text{machine})$     $S \rightarrow (\text{翻译}, \text{translation})$     $S \rightarrow (\text{软件}, \text{software})$

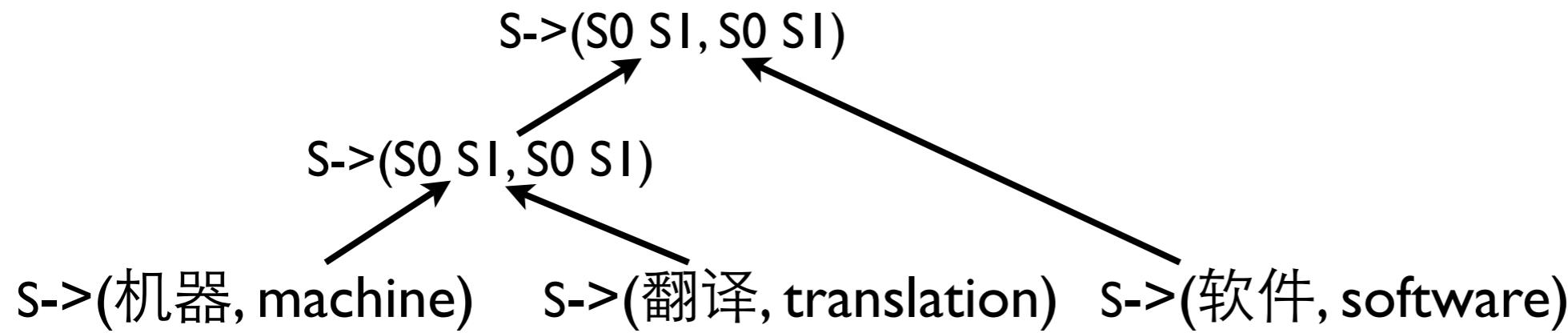
# Spurious Ambiguity in Derivation Trees

机器 翻译 软件



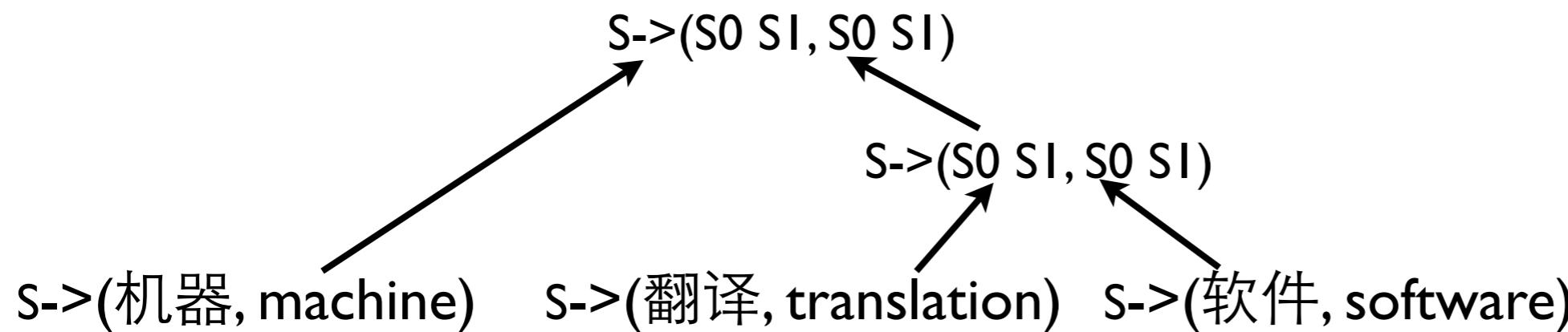
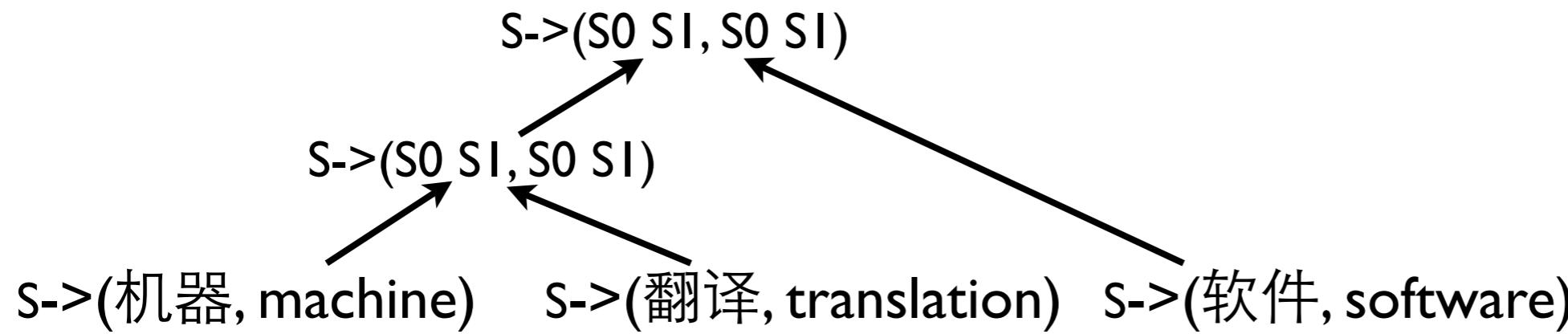
# Spurious Ambiguity in Derivation Trees

机器 翻译 软件



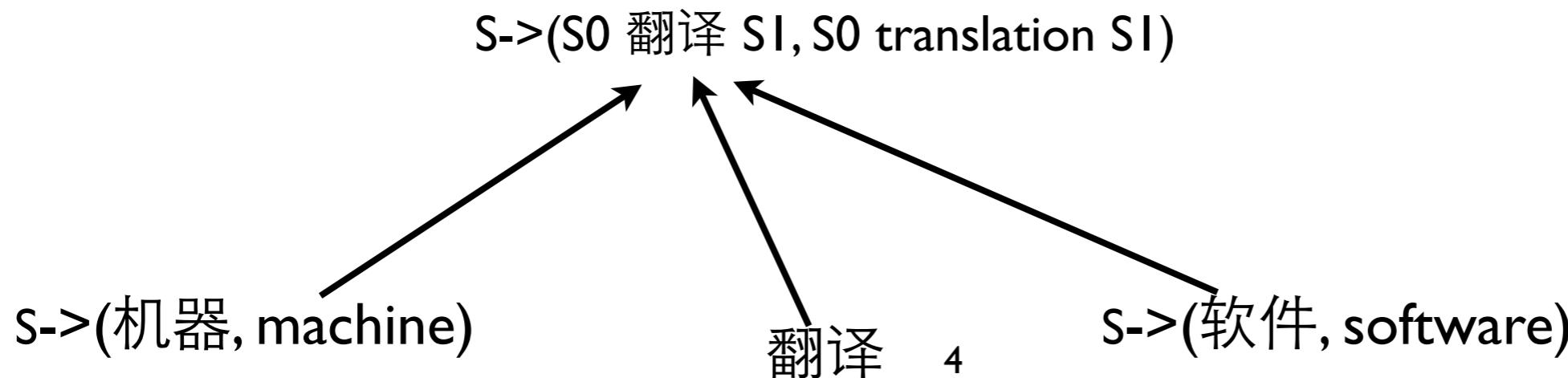
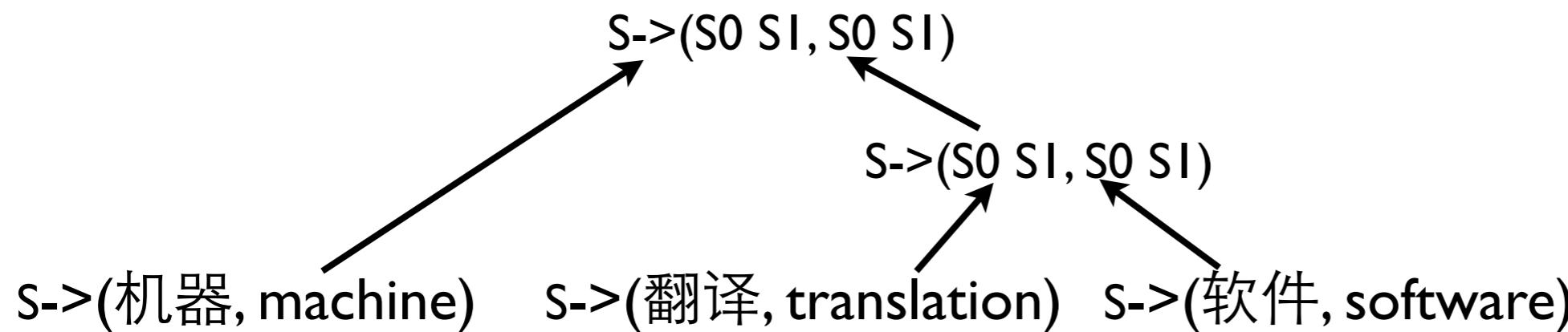
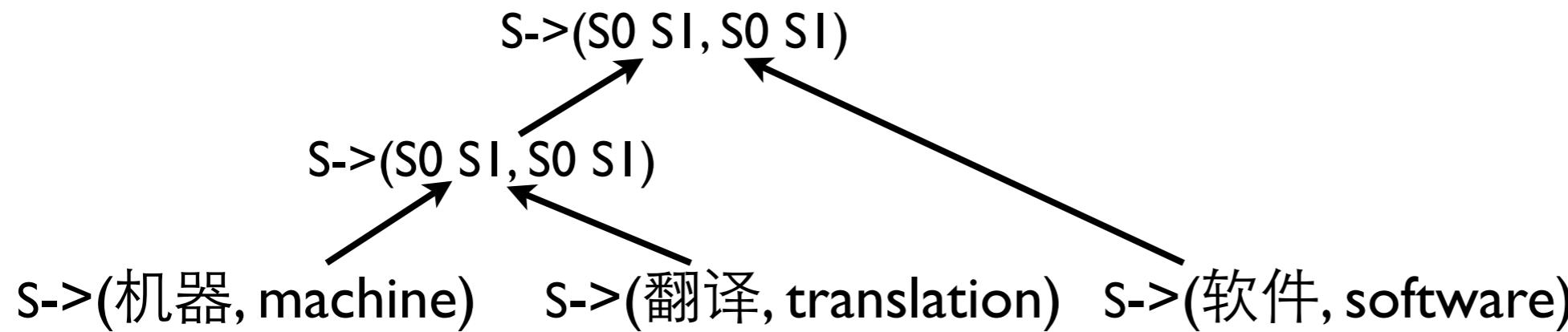
# Spurious Ambiguity in Derivation Trees

机器 翻译 软件



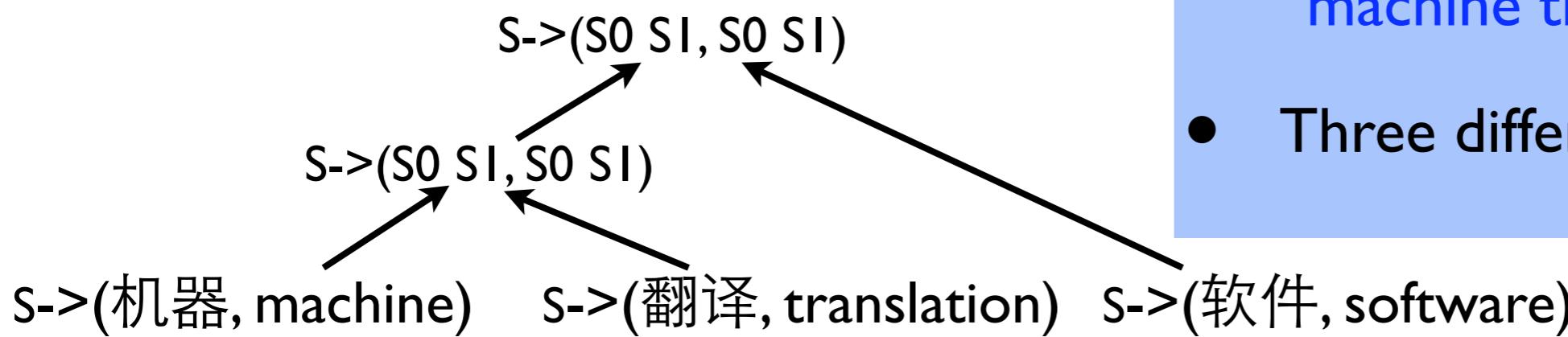
# Spurious Ambiguity in Derivation Trees

机器 翻译 软件

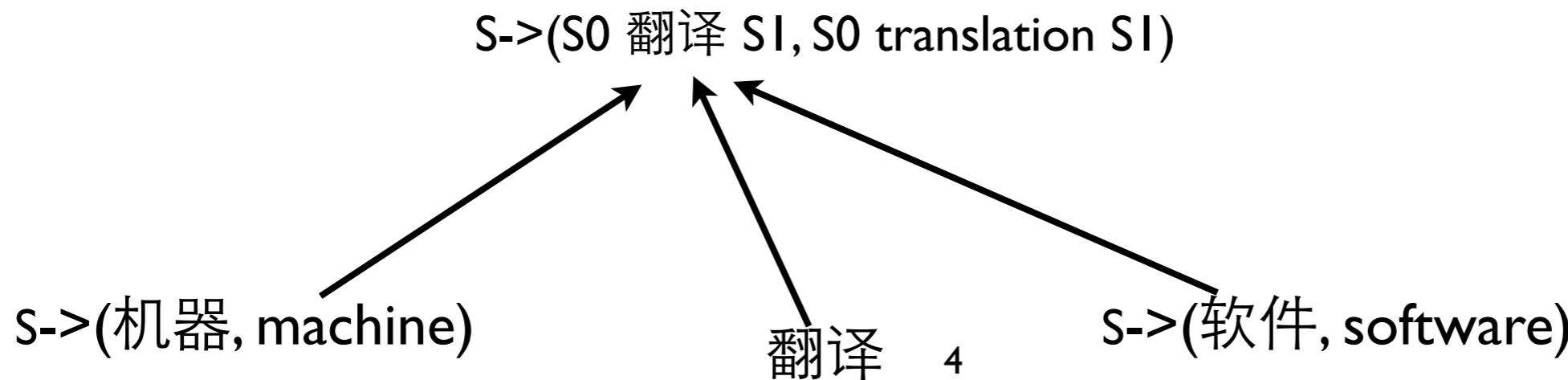
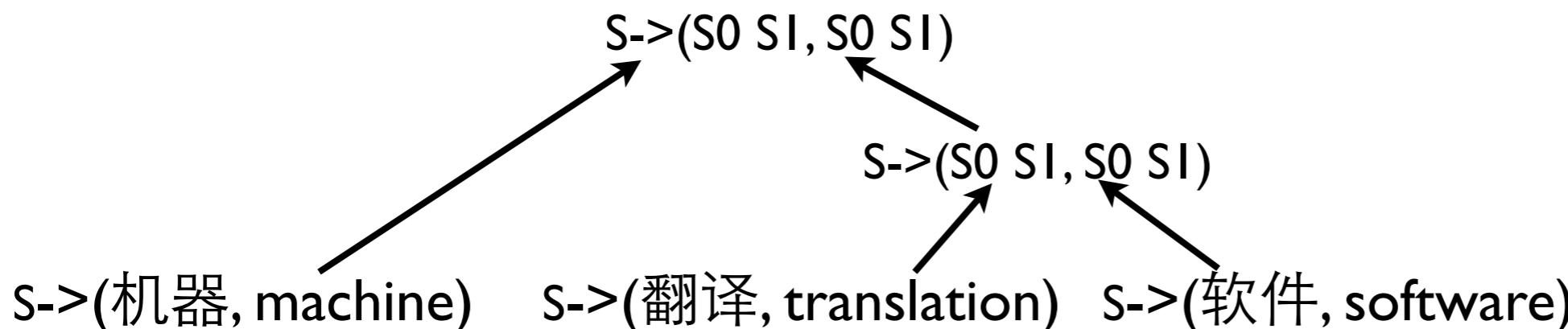


# Spurious Ambiguity in Derivation Trees

机器 翻译 软件



- Same output:  
“**machine translation software**”
- Three different derivation trees



# Maximum A Posterior (MAP) Decoding

# Maximum A Posterior (MAP) Decoding

**translation  
string**

red translation

blue translation

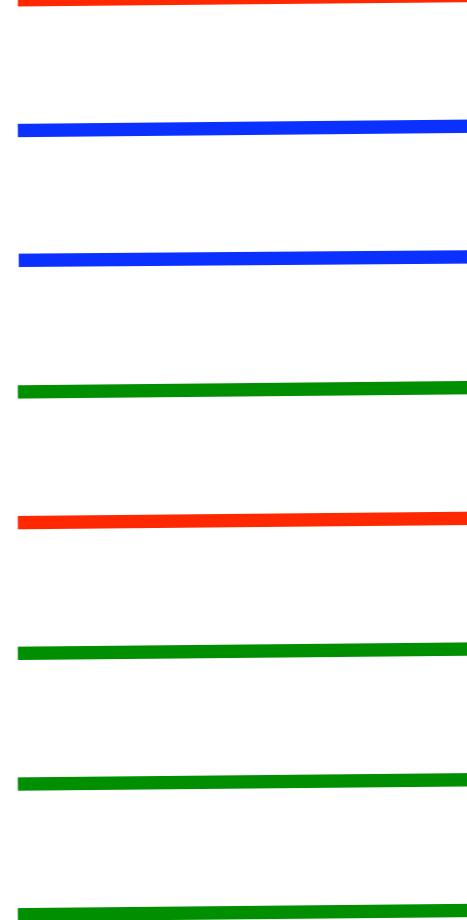
green translation

# Maximum A Posterior (MAP) Decoding

**translation  
string**

red translation

**derivation**



blue translation

green translation

# Maximum A Posterior (MAP) Decoding

**translation  
string**

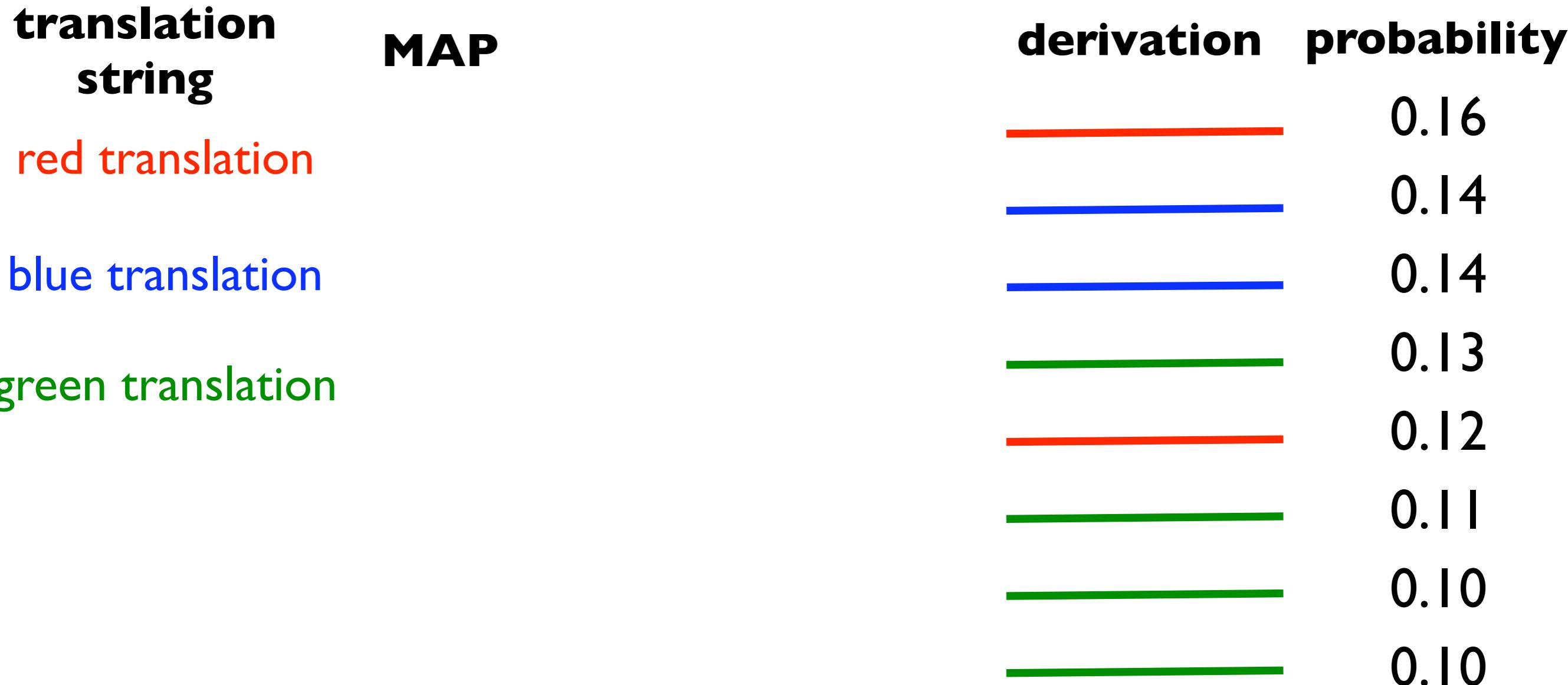
red translation

blue translation

green translation

<b>derivation</b>	<b>probability</b>
—	0.16
—	0.14
—	0.14
—	0.13
—	0.12
—	0.11
—	0.10
—	0.10

# Maximum A Posterior (MAP) Decoding



# Maximum A Posterior (MAP) Decoding

**translation  
string**

red translation

blue translation

green translation

- Exact MAP decoding

**MAP**

<b>derivation</b>	<b>probability</b>
—	0.16
—	0.14
—	0.14
—	0.13
—	0.12
—	0.11
—	0.10
—	0.10

# Maximum A Posterior (MAP) Decoding

**translation  
string**

red translation

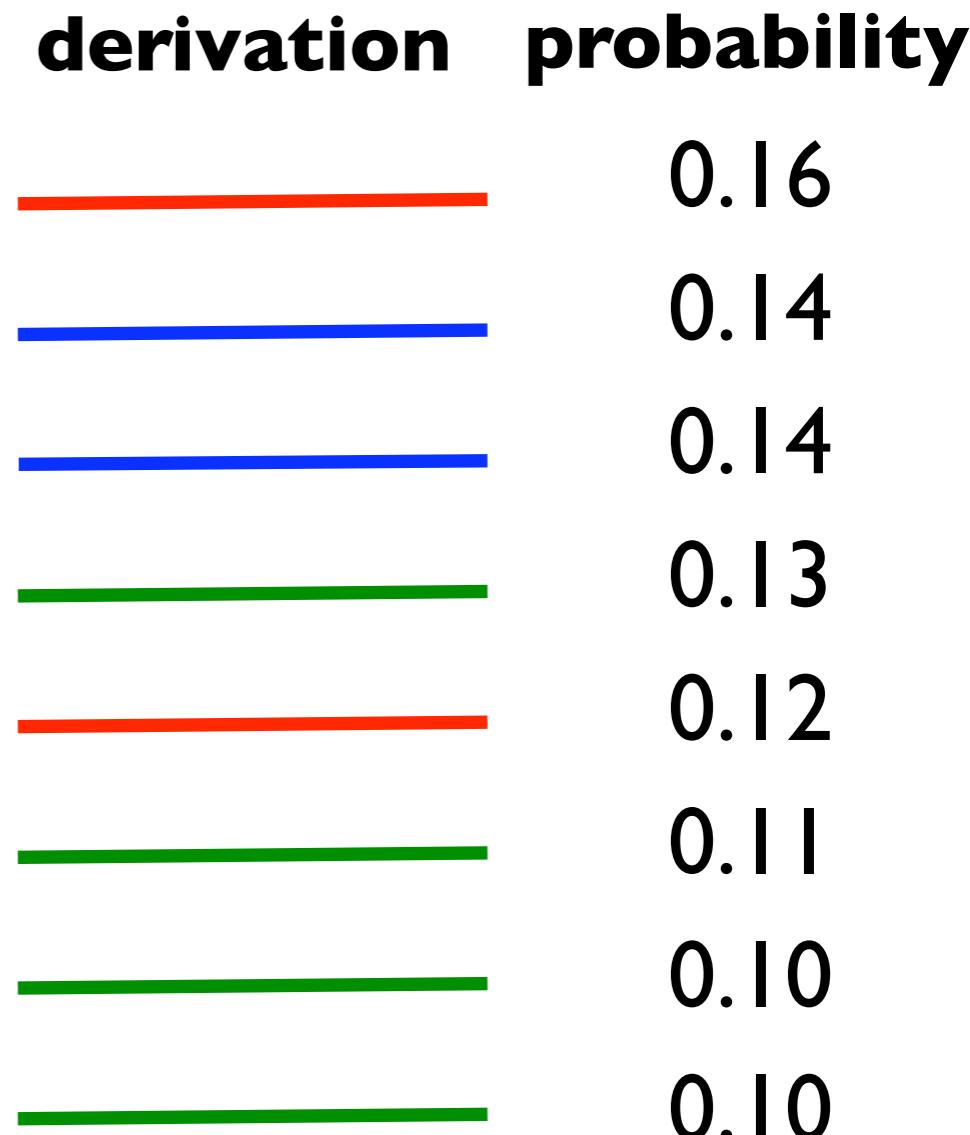
blue translation

green translation

- Exact MAP decoding

$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$



- **x**: Foreign sentence
- **y**: English sentence
- **d**: derivation

# Maximum A Posterior (MAP) Decoding

**translation  
string**

red translation

blue translation

green translation

**MAP**

<b>derivation</b>	<b>probability</b>
—	0.16
—	0.14
—	0.14
—	0.13
—	0.12
—	0.11
—	0.10
—	0.10

- Exact MAP decoding

$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

- **x**: Foreign sentence
- **y**: English sentence
- **d**: derivation

# Maximum A Posterior (MAP) Decoding

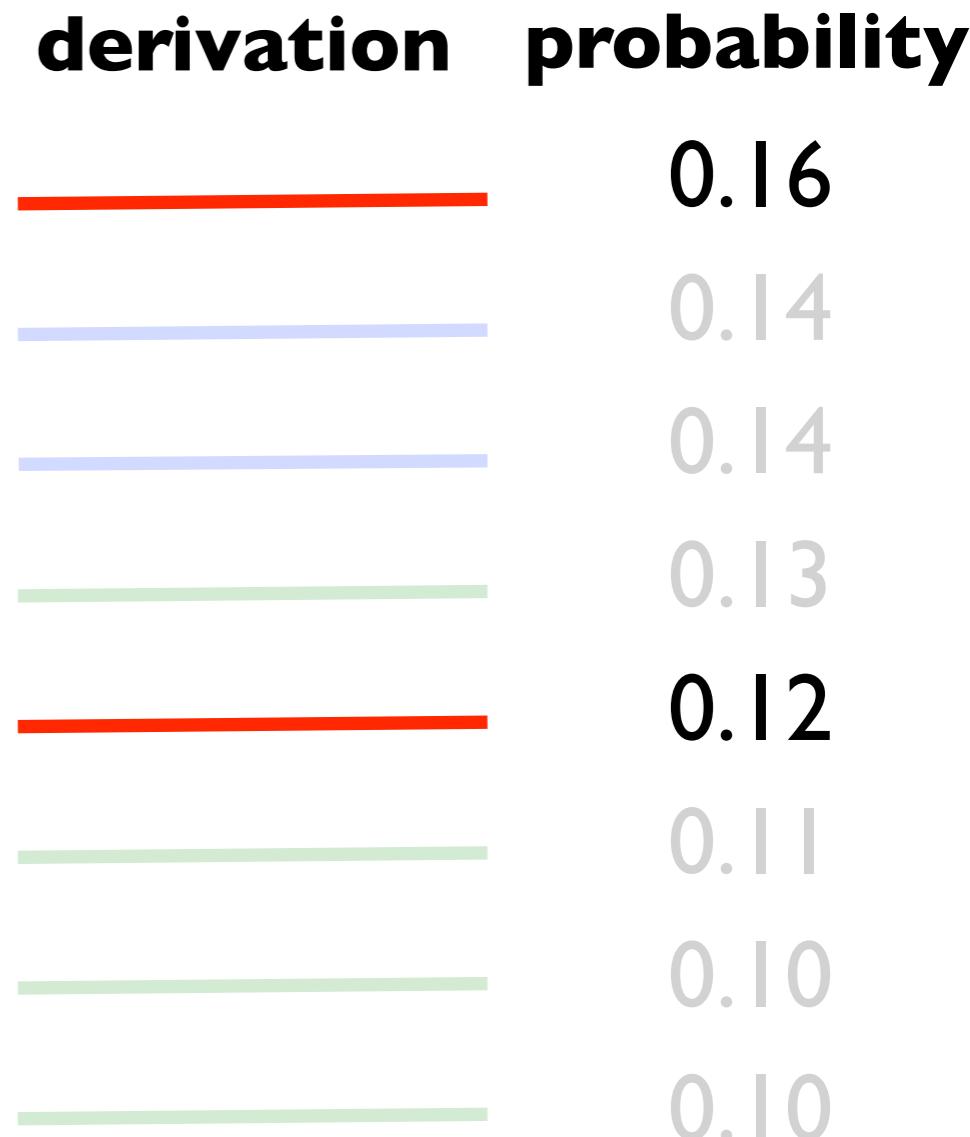
**translation  
string**

red translation

blue translation

green translation

**MAP**



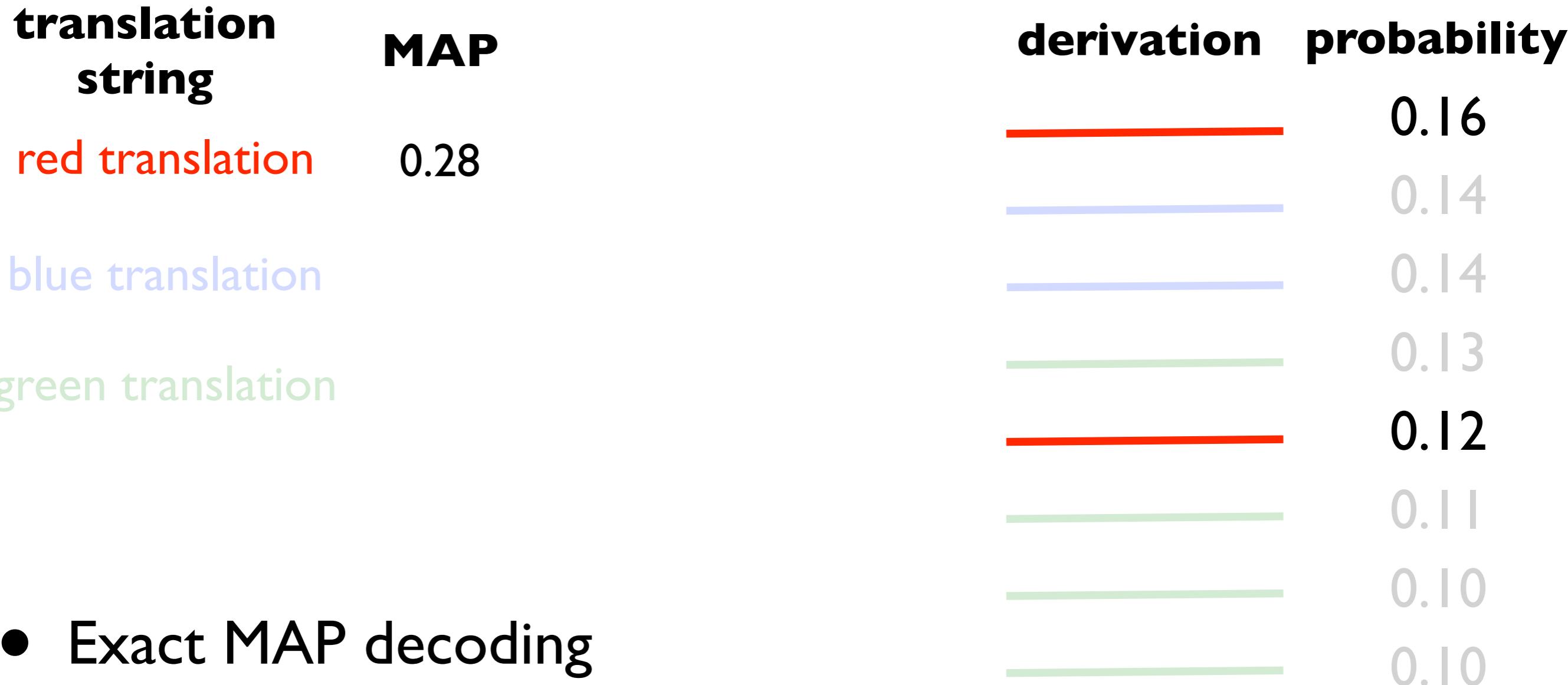
- Exact MAP decoding

$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

- **x**: Foreign sentence
- **y**: English sentence
- **d**: derivation

# Maximum A Posterior (MAP) Decoding



- Exact MAP decoding

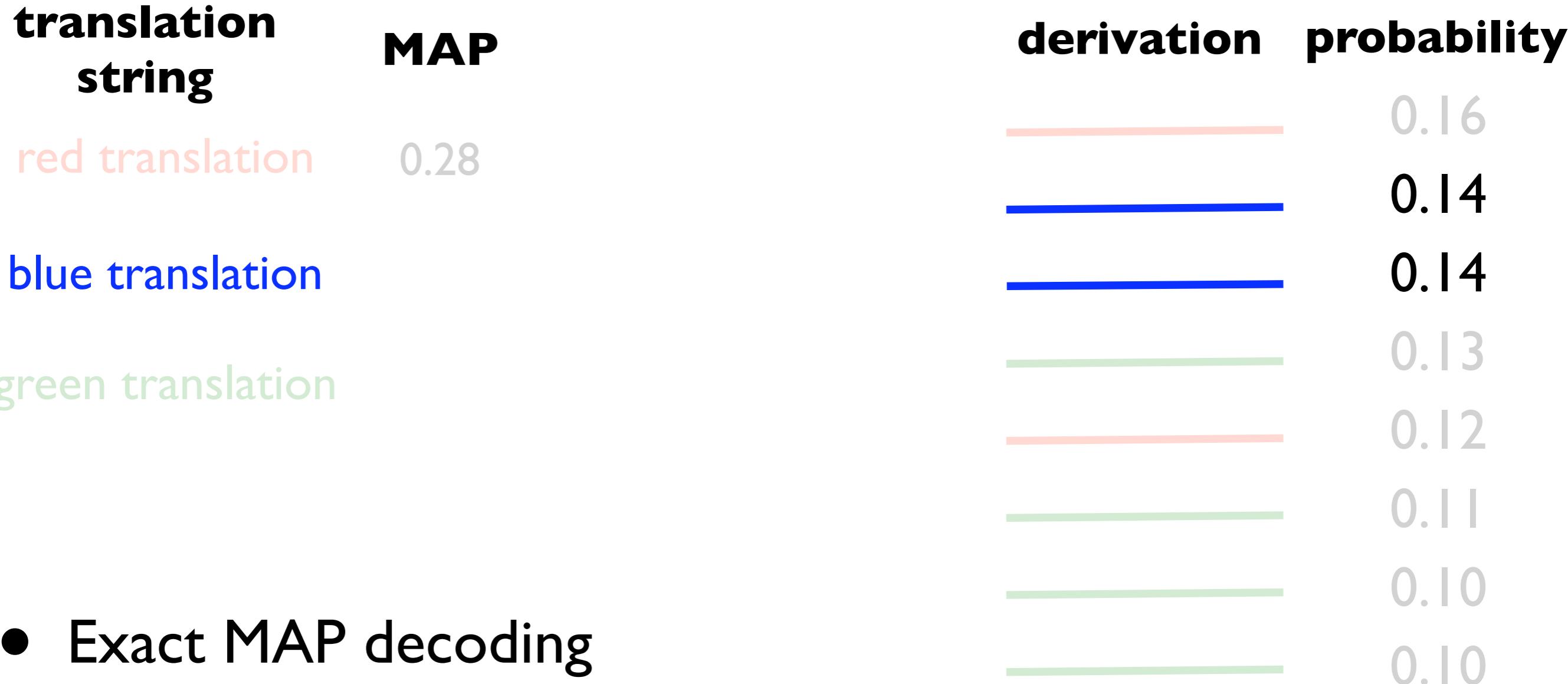
$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

6

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



- Exact MAP decoding

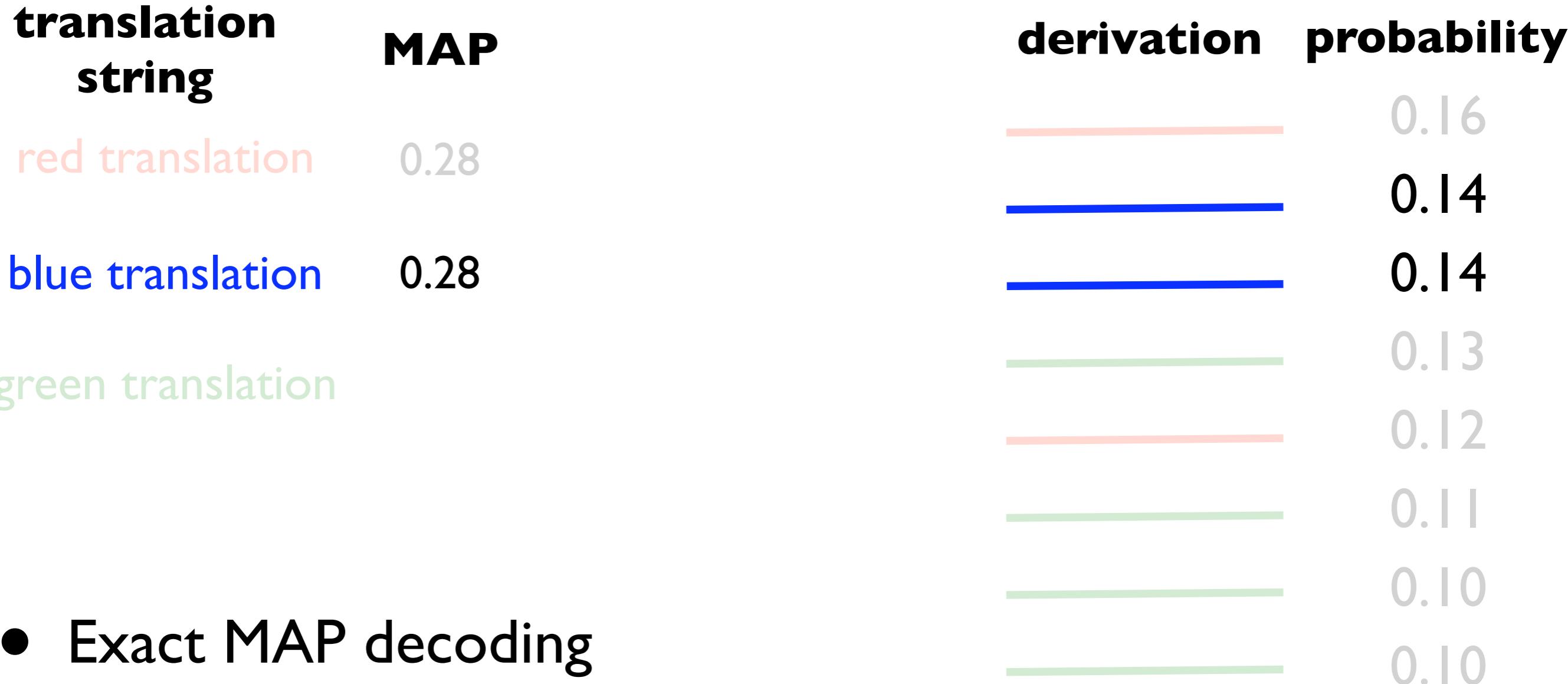
$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

7

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



- Exact MAP decoding

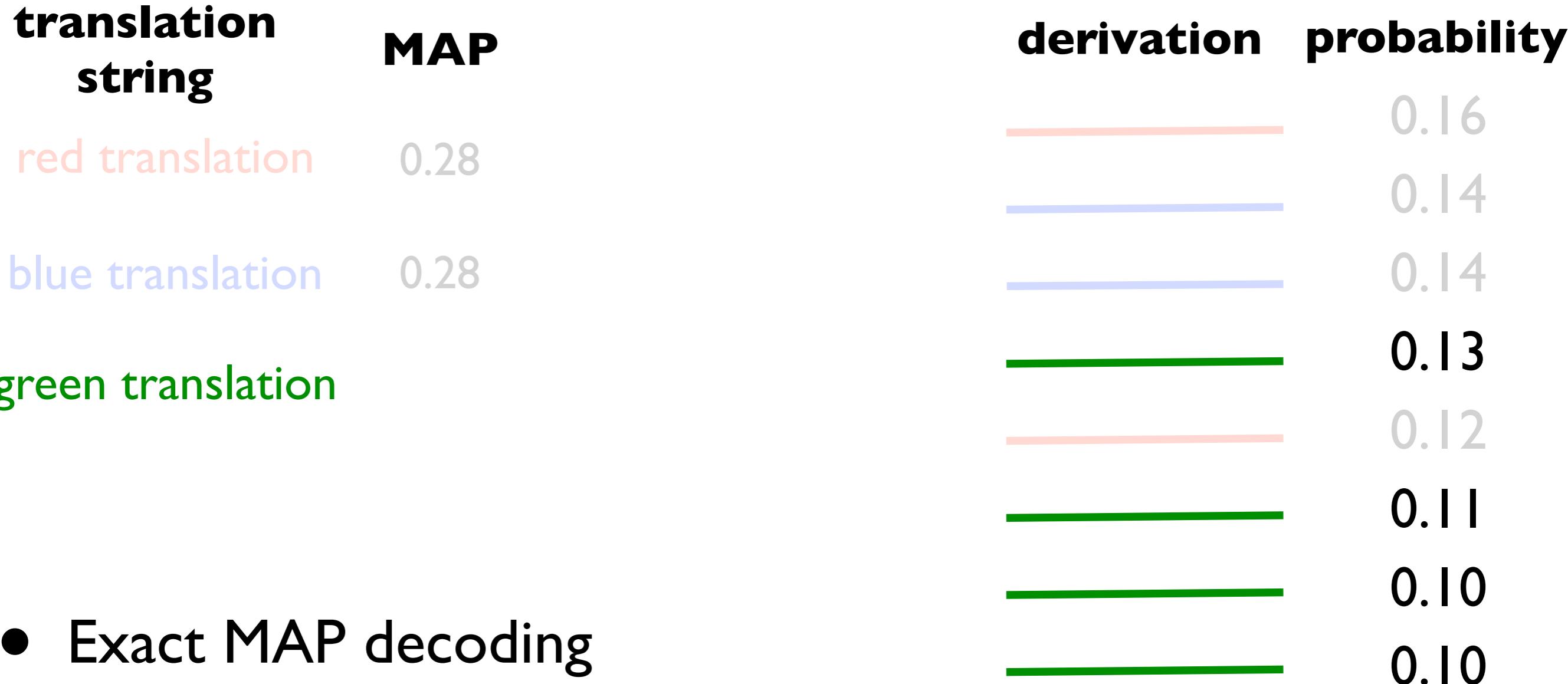
$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

7

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



- Exact MAP decoding

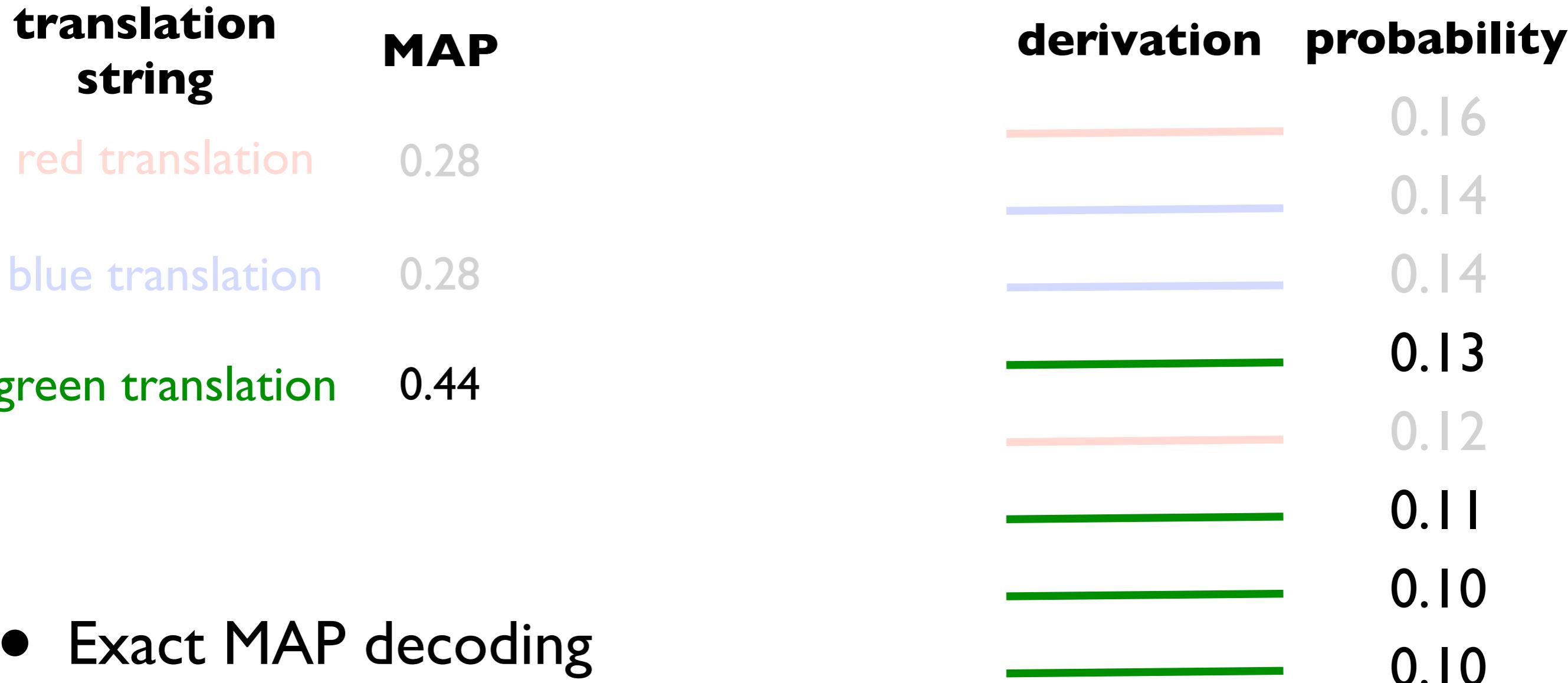
$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

8

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



- Exact MAP decoding

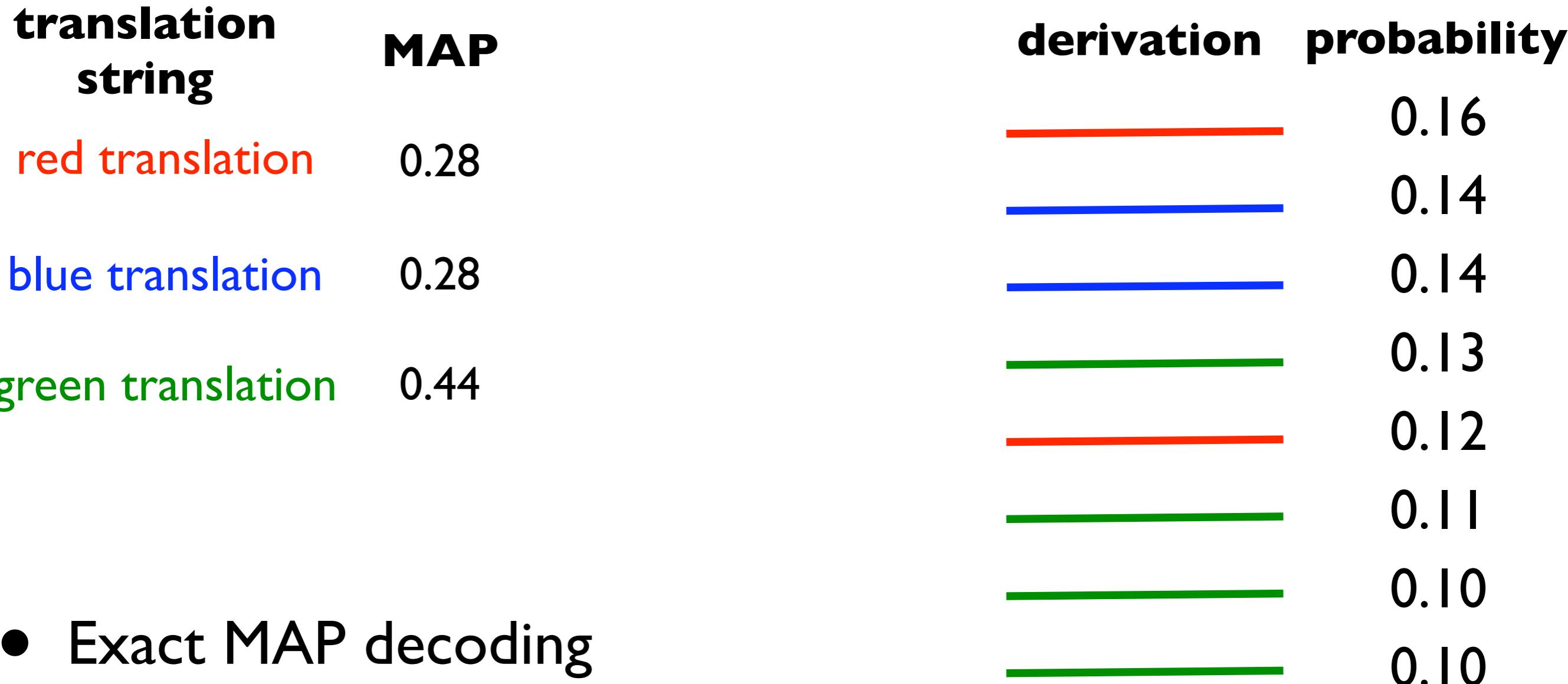
$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

8

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



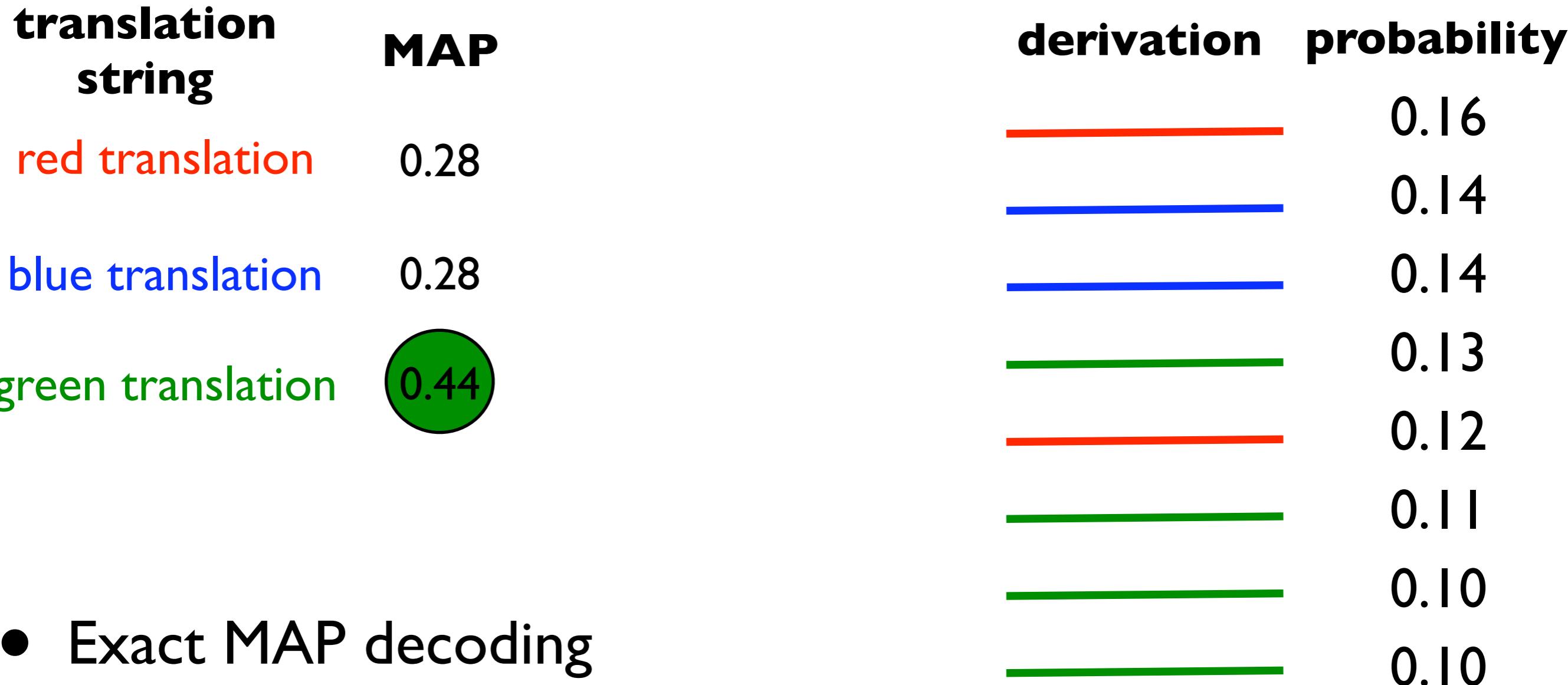
- Exact MAP decoding

$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



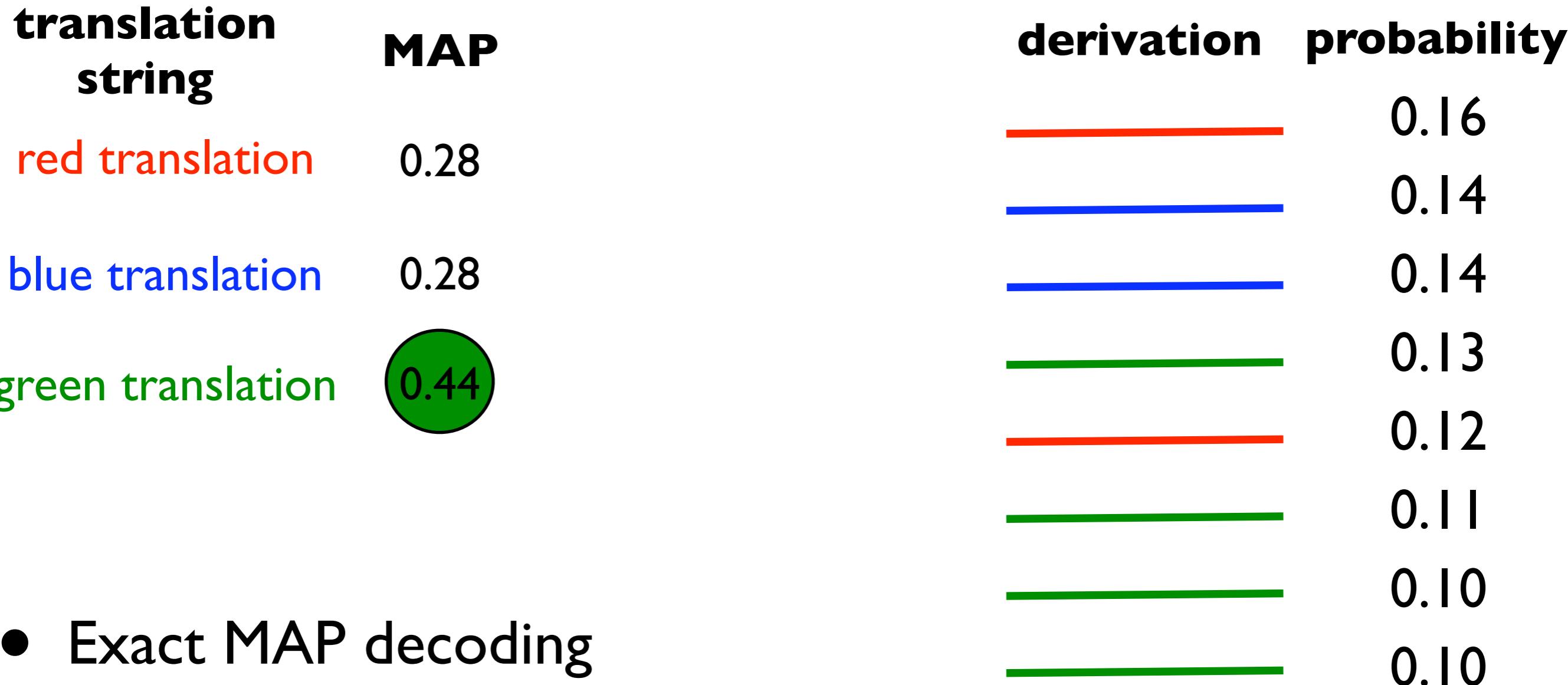
- Exact MAP decoding

$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

# Maximum A Posterior (MAP) Decoding



- Exact MAP decoding

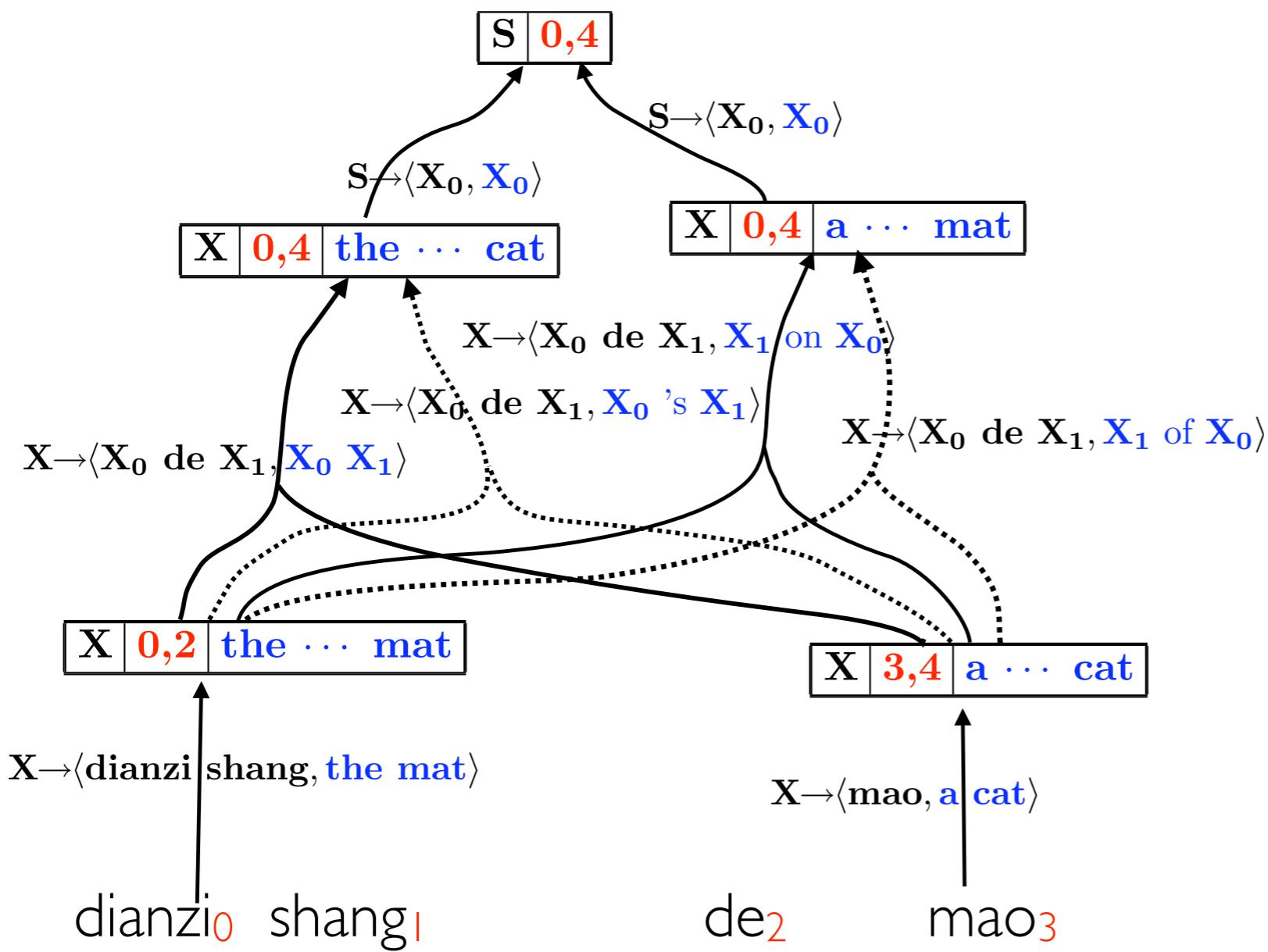
$$y^* = \arg \max_{y \in \text{Trans}(x)} p(y|x)$$

$$= \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y)} p(y, d|x)$$

- $x$ : Foreign sentence
- $y$ : English sentence
- $d$ : derivation

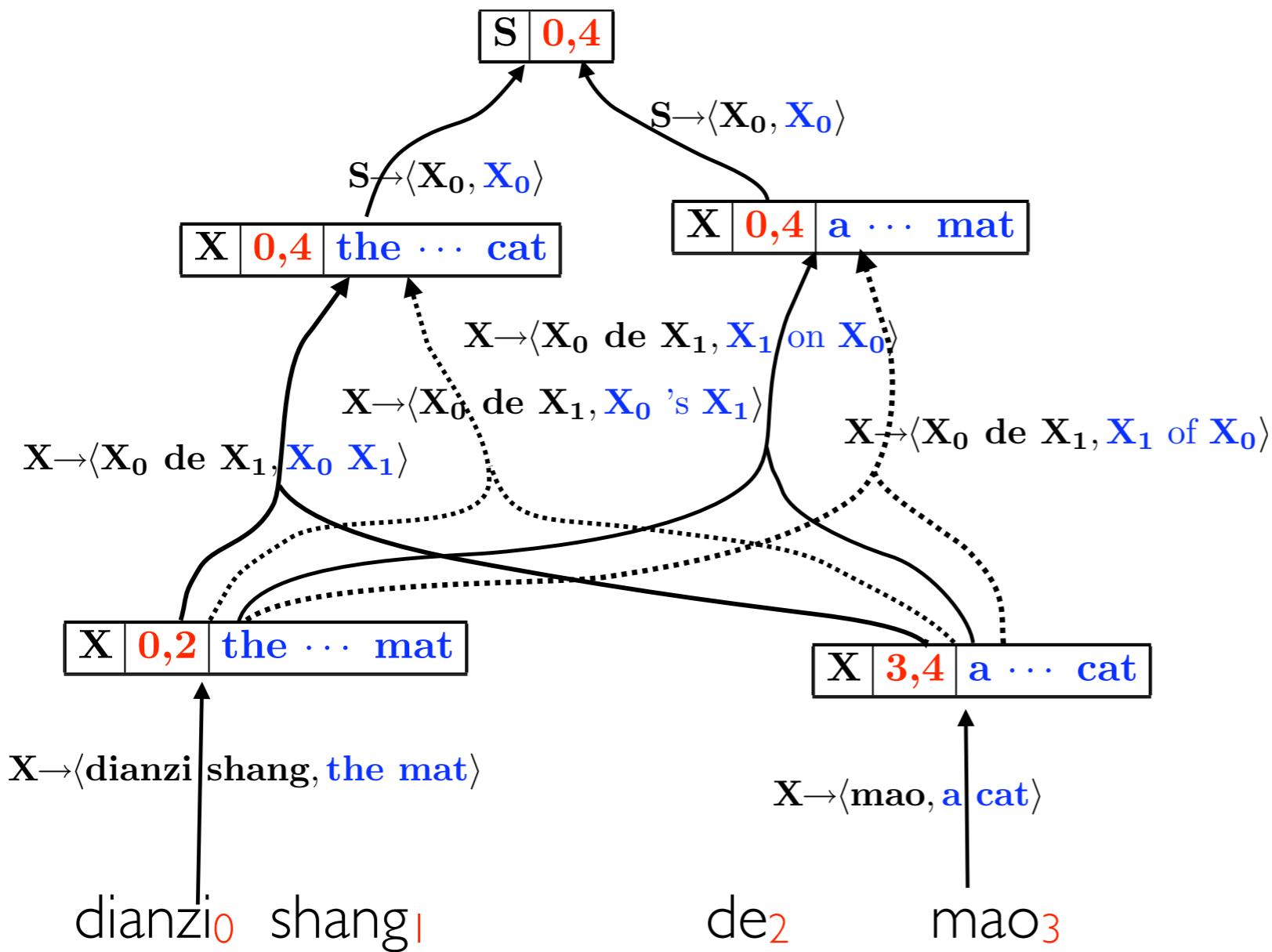
# Hypergraph as a search space

# Hypergraph as a search space

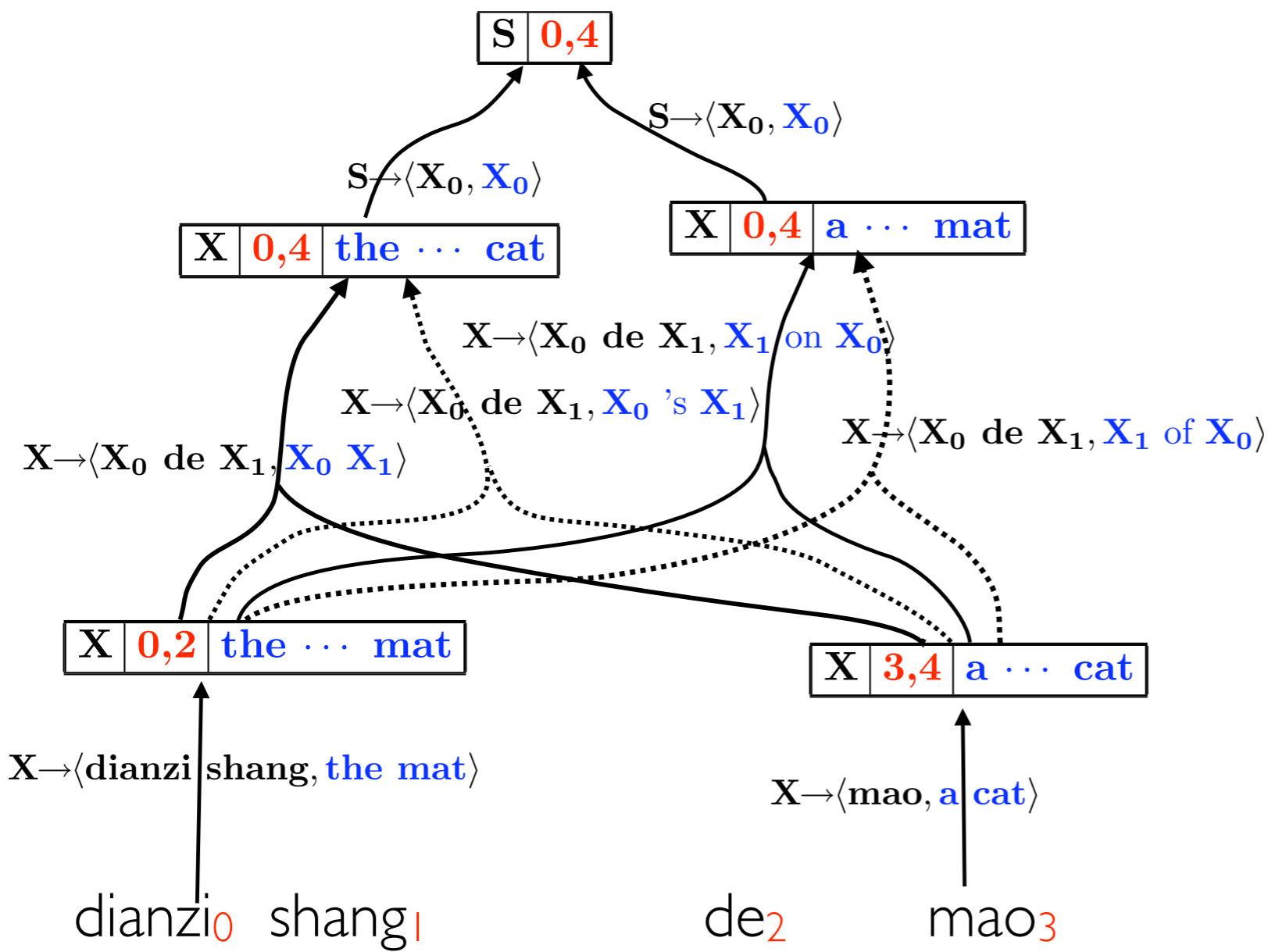


# Hypergraph as a search space

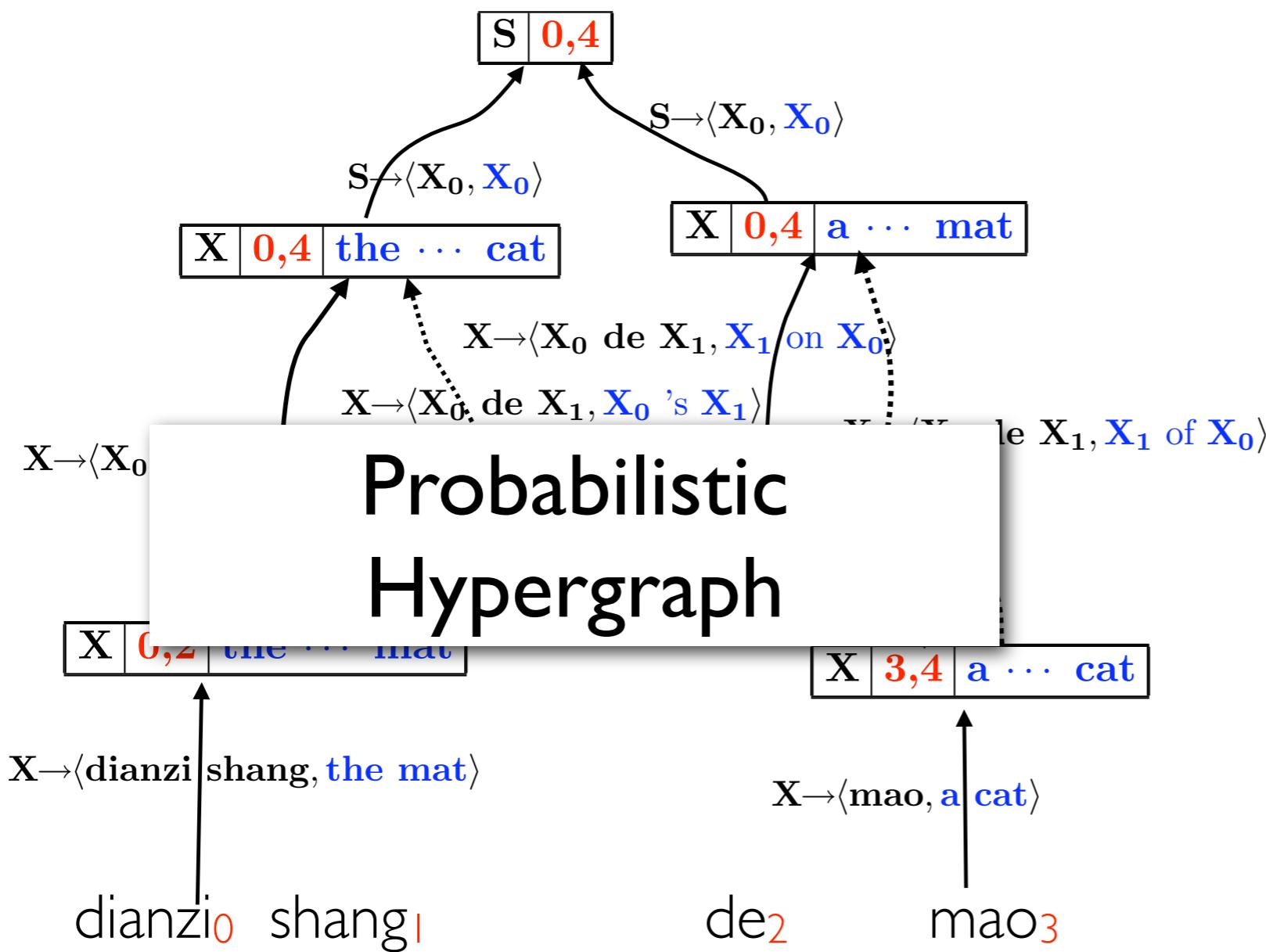
A hypergraph is a compact structure to encode exponentially many trees.



# Hypergraph as a search space

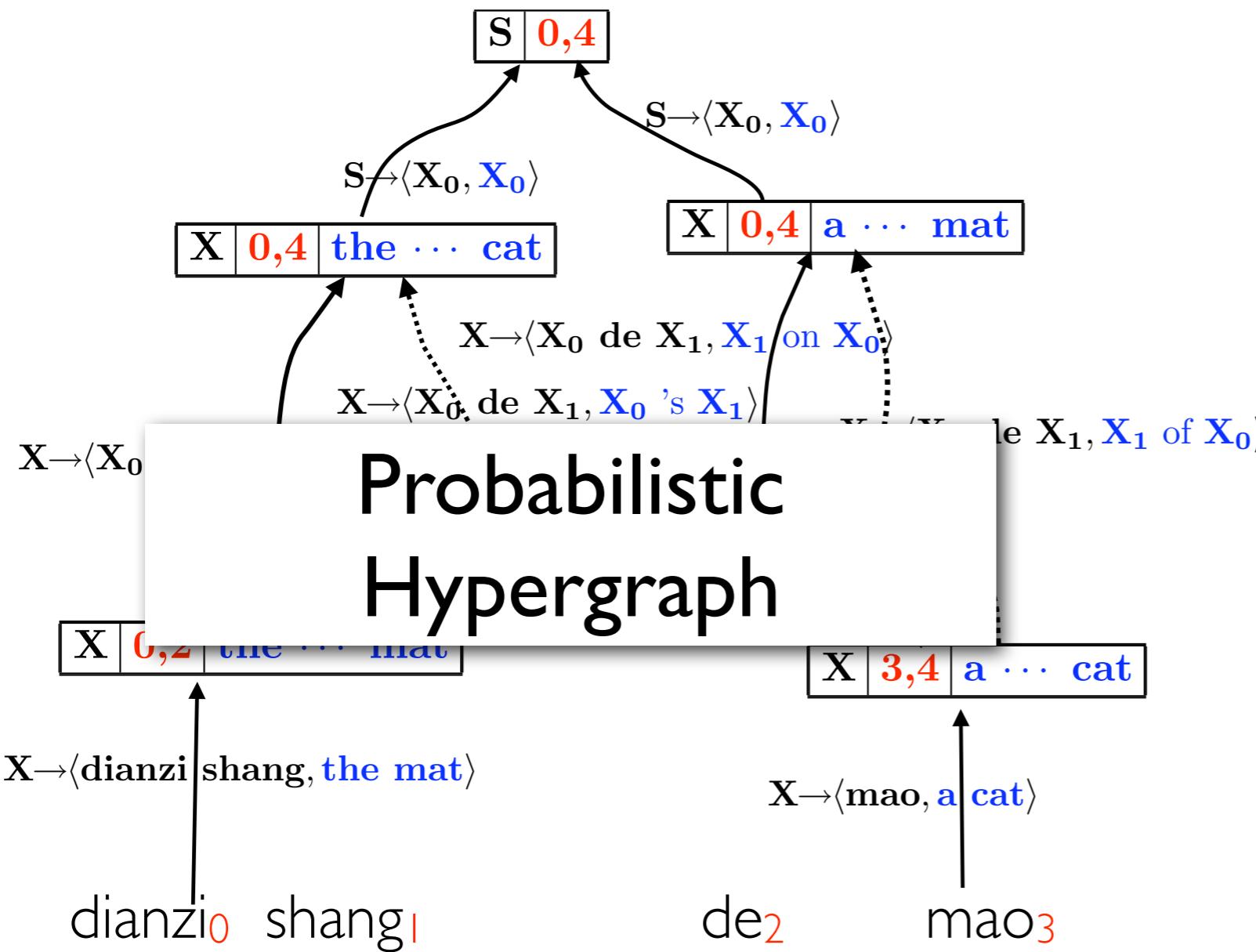


# Hypergraph as a search space



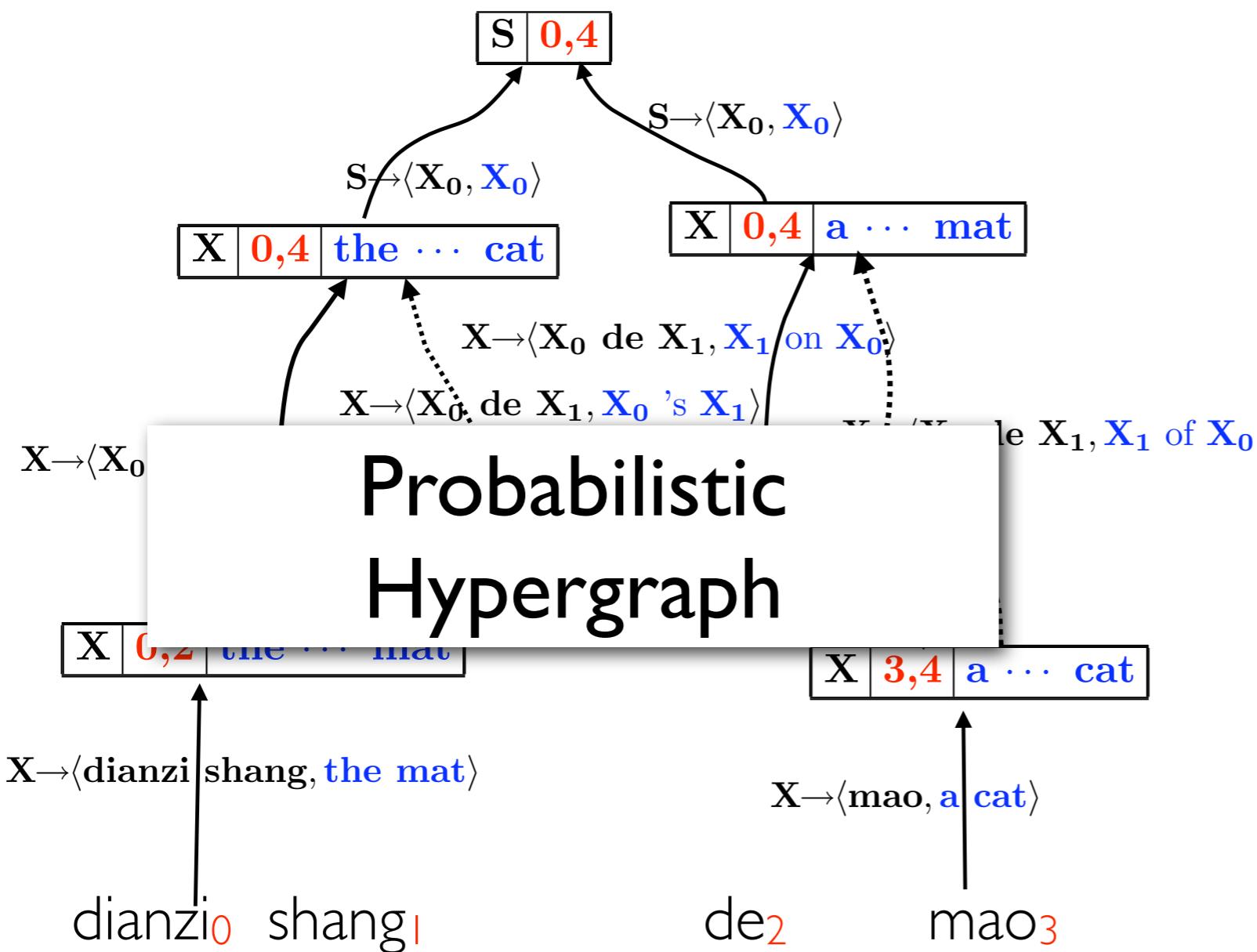
# Hypergraph as a search space

The hypergraph defines a probability distribution over **derivation trees**, i.e.  $p(y, d | x)$ ,



# Hypergraph as a search space

The hypergraph defines a probability distribution over **derivation trees**, i.e.  $p(y, d | x)$ , and also a distribution (implicit) over **strings**, i.e.  $p(y | x)$ .



# Hypergraph as a search space

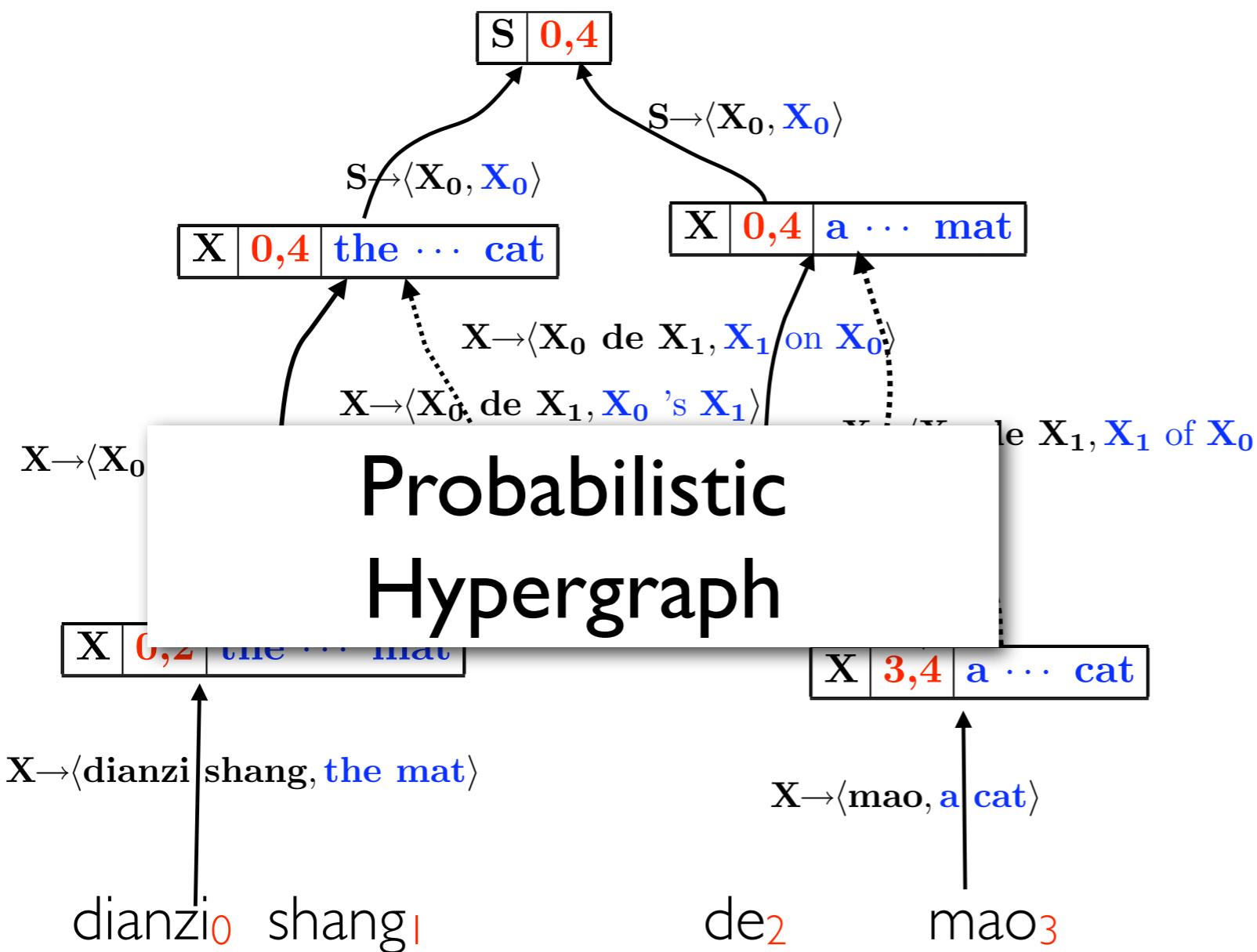
The hypergraph defines a probability distribution over **derivation trees**, i.e.  $p(y, d | x)$ ,  
and also a distribution (implicit) over **strings**, i.e.  $p(y | x)$ .

- Exact MAP decoding

$$\begin{aligned} y^* &= \arg \max_{y \in HG(x)} p(y|x) \\ &= \arg \max_{y \in HG(x)} \sum_{d \in D(x,y)} p(y, d|x) \end{aligned}$$

exponential size

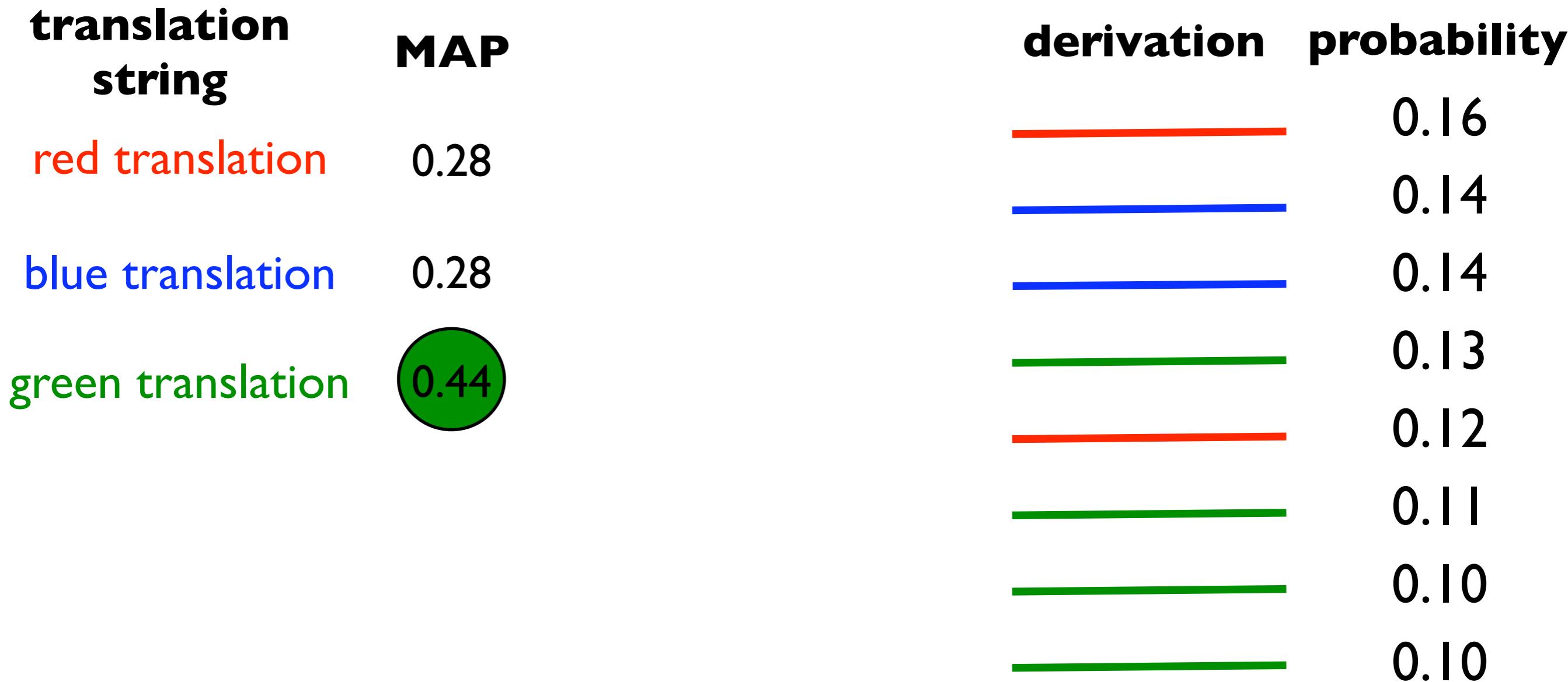
NP-hard (Sima'an 1996)



# Decoding with spurious ambiguity?

- Maximum a posterior (MAP) decoding
- Viterbi approximation
- N-best approximation (**crunching**) (May and Knight 2006)

# Viterbi Approximation



- Viterbi approximation

$$y^* = \arg \max_{y \in \text{Trans}(x)} \max_{d \in D(x,y)} p(y, d|x)$$

$$= Y(\arg \max_{d \in D(x)} p(y, d|x))$$

# Viterbi Approximation

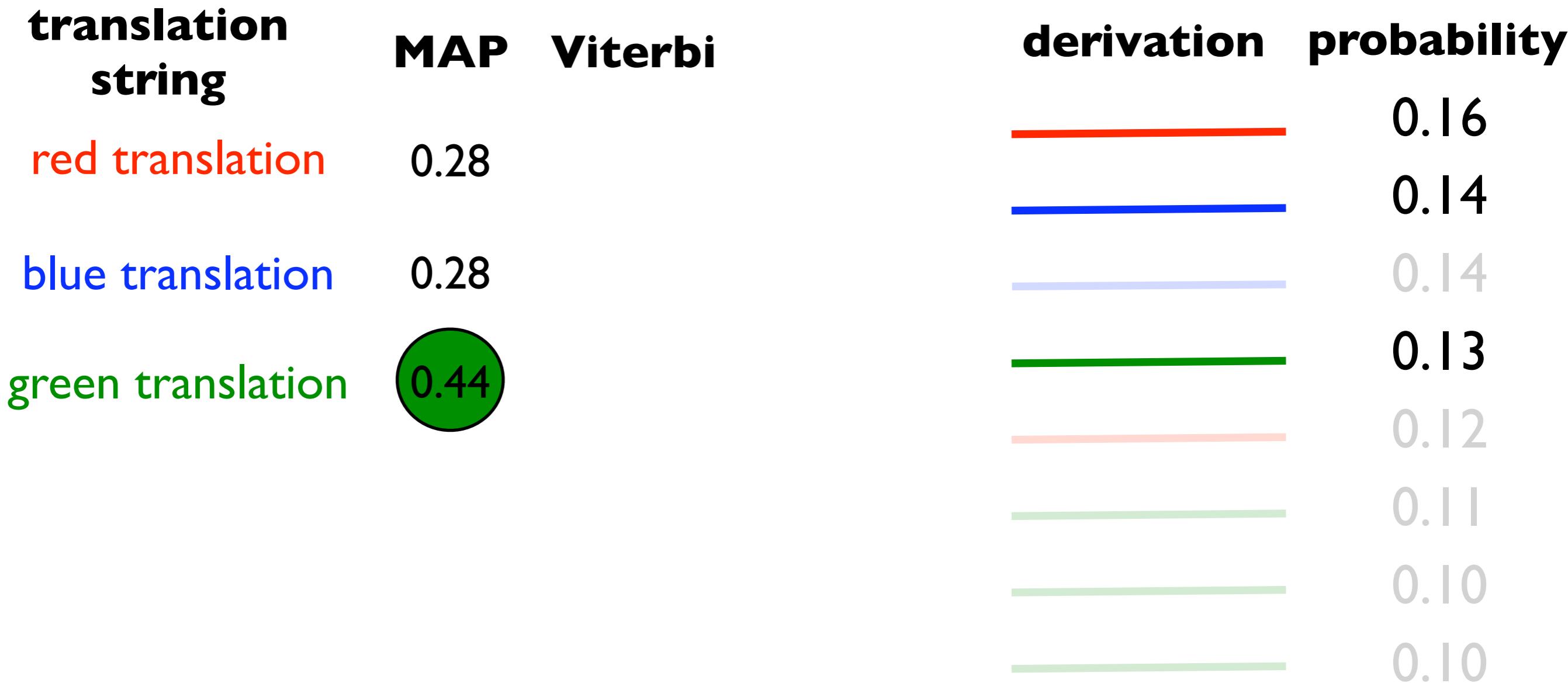
<b>translation string</b>	<b>MAP</b>	<b>Viterbi</b>	<b>derivation</b>	<b>probability</b>
red translation	0.28		—	0.16
blue translation	0.28		—	0.14
green translation	0.44		—	0.14
			—	0.13
			—	0.12
			—	0.11
			—	0.10
			—	0.10

- Viterbi approximation

$$y^* = \arg \max_{y \in \text{Trans}(x)} \max_{d \in D(x,y)} p(y, d|x)$$

$$= Y(\arg \max_{d \in D(x)} p(y, d|x))$$

# Viterbi Approximation



- Viterbi approximation

$$y^* = \arg \max_{y \in \text{Trans}(x)} \max_{d \in D(x,y)} p(y, d|x)$$

$$= Y(\arg \max_{d \in D(x)} p(y, d|x))$$

# Viterbi Approximation

translation string	MAP	Viterbi	derivation	probability
red translation	0.28	0.16	—	0.16
blue translation	0.28	0.14	—	0.14
green translation	0.44	0.13	—	0.13
			—	0.12
			—	0.11
			—	0.10
			—	0.10

- Viterbi approximation

$$y^* = \arg \max_{y \in \text{Trans}(x)} \max_{d \in D(x,y)} p(y, d|x)$$

$$= Y(\arg \max_{d \in D(x)} p(y, d|x))$$

# Viterbi Approximation

translation string	MAP	Viterbi	derivation	probability
red translation	0.28	0.16	—	0.16
blue translation	0.28	0.14	—	0.14
green translation	0.44	0.13	—	0.13
			—	0.12
			—	0.11
			—	0.10
			—	0.10

- Viterbi approximation

$$y^* = \arg \max_{y \in \text{Trans}(x)} \max_{d \in D(x,y)} p(y, d|x)$$

$$= Y(\arg \max_{d \in D(x)} p(y, d|x))$$

# N-best Approximation

<b>translation string</b>	<b>MAP</b>	<b>Viterbi</b>	<b>derivation</b>	<b>probability</b>
red translation	0.28	0.16	—	0.16
blue translation	0.28	0.14	—	0.14
green translation	0.44	0.13	—	0.13

The diagram illustrates the N-best approximation results. For each translation string (red, blue, green), two scores are provided: MAP and Viterbi. The Viterbi scores are highlighted with colored circles (red for red, blue for blue, green for green). Each score is connected by a vertical line to a horizontal bar representing its derivation. The bars are colored red, blue, and green, corresponding to the translation strings. The bars are ordered by probability from highest to lowest.

- N-best approximation (**crunching**) (May and Knight 2006)

$$y^* = \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y) \cap \text{ND}(x)} p(y, d|x)$$

# N-best Approximation

<b>translation string</b>	<b>MAP</b>	<b>Viterbi</b>	<b>4-best crunching</b>	<b>derivation</b>	<b>probability</b>
red translation	0.28	0.16			0.16
blue translation	0.28	0.14			0.14
green translation	0.44	0.13			0.13

The table shows the results of N-best approximation for three different translation strings: red, blue, and green. For each string, the MAP score is listed in the first column. The Viterbi probability is shown in a colored circle (red for red, blue for blue, green for green). The 4-best crunching results are shown as horizontal bars of decreasing length from left to right, with their corresponding probabilities listed in the last column.

- N-best approximation (**crunching**) (May and Knight 2006)

$$y^* = \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y) \cap \text{ND}(x)} p(y, d|x)$$

# N-best Approximation

<b>translation string</b>	<b>MAP</b>	<b>Viterbi</b>	<b>4-best crunching</b>	<b>derivation</b>	<b>probability</b>
red translation	0.28	0.16			0.16
blue translation	0.28	0.14			0.14
green translation	0.44	0.13			0.13
				—	0.12
				—	0.11
				—	0.10
				—	0.10

- N-best approximation (**crunching**) (May and Knight 2006)

$$y^* = \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y) \cap \text{ND}(x)} p(y, d|x)$$

# N-best Approximation

<b>translation string</b>	<b>MAP</b>	<b>Viterbi</b>	<b>4-best crunching</b>	<b>derivation</b>	<b>probability</b>
red translation	0.28	0.16	0.16	—	0.16
blue translation	0.28	0.14	0.28	—	0.14
green translation	0.44	0.13	0.13	—	0.13
				—	0.12
				—	0.11
				—	0.10
				—	0.10

- N-best approximation (**crunching**) (May and Knight 2006)

$$y^* = \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y) \cap \text{ND}(x)} p(y, d|x)$$

# N-best Approximation

translation string	MAP	Viterbi	4-best crunching	derivation	probability
red translation	0.28	0.16	0.16	—	0.16
blue translation	0.28	0.14	0.28	—	0.14
green translation	0.44	0.13	0.13	—	0.13
				—	0.12
				—	0.11
				—	0.10
				—	0.10

- N-best approximation (**crunching**) (May and Knight 2006)

$$y^* = \arg \max_{y \in \text{Trans}(x)} \sum_{d \in D(x,y) \cap \text{ND}(x)} p(y, d|x)$$

# MAP vs. Approximations

translation string	MAP	Viterbi	4-best crunching	derivation	probability
red translation	0.28	0.16	0.16	—	0.16
blue translation	0.28	0.14	0.28	—	0.14
green translation	0.44	0.13	0.13	—	0.13

Approximate derivation probabilities for the green translation:

- Red line: 0.12
- Green line: 0.11
- Green line: 0.10
- Green line: 0.10

# MAP vs. Approximations

translation string	MAP	Viterbi	4-best crunching	derivation	probability
red translation	0.28	0.16	0.16	—	0.16
blue translation	0.28	0.14	0.28	—	0.14
green translation	0.44	0.13	0.13	—	0.13

Legend: red circle = Viterbi, blue circle = 4-best crunching, green circle = MAP

- Exact MAP decoding under spurious ambiguity is **intractable**

# MAP vs. Approximations

translation string	MAP	Viterbi	4-best crunching	derivation	probability
red translation	0.28	0.16	0.16	—	0.16
blue translation	0.28	0.14	0.28	—	0.14
green translation	0.44	0.13	0.13	—	0.13

- Exact MAP decoding under spurious ambiguity is **intractable**
- Viterbi and crunching are efficient, but ignore most derivations

# MAP vs. Approximations

translation string	MAP	Viterbi	4-best crunching	derivation	probability
red translation	0.28	0.16	0.16	—	0.16
blue translation	0.28	0.14	0.28	—	0.14
green translation	0.44	0.13	0.13	—	0.13

- Exact MAP decoding under spurious ambiguity is **intractable**
- Viterbi and crunching are efficient, but ignore most derivations
- Our goal: develop an **approximation** that considers **all** the derivations **but** still allows **tractable** decoding

# Variational Decoding

# Variational Decoding

**Decoding** using **Variational** approximation

# Variational Decoding

**Decoding using Variational approximation**

**Decoding using a sentence-specific  
approximate distribution**

# Variational Decoding for MT: an Overview

# Variational Decoding for MT: an Overview

Sentence-specific decoding

# Variational Decoding for MT: an Overview

Sentence-specific decoding

Three steps:

# Variational Decoding for MT: an Overview

## Sentence-specific decoding

Three steps:

1

Generate a hypergraph

# Variational Decoding for MT: an Overview

## Sentence-specific decoding

Three steps:

1

Generate a hypergraph

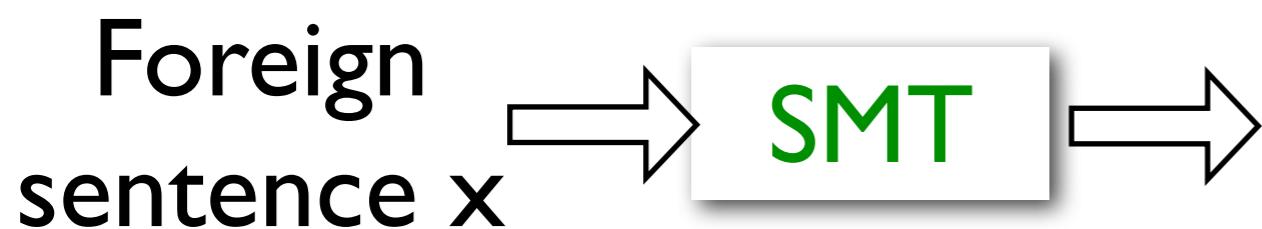
Foreign  
sentence  $x$

# Variational Decoding for MT: an Overview

## Sentence-specific decoding

Three steps:

- 1 Generate a hypergraph



# Variational Decoding for MT: an Overview

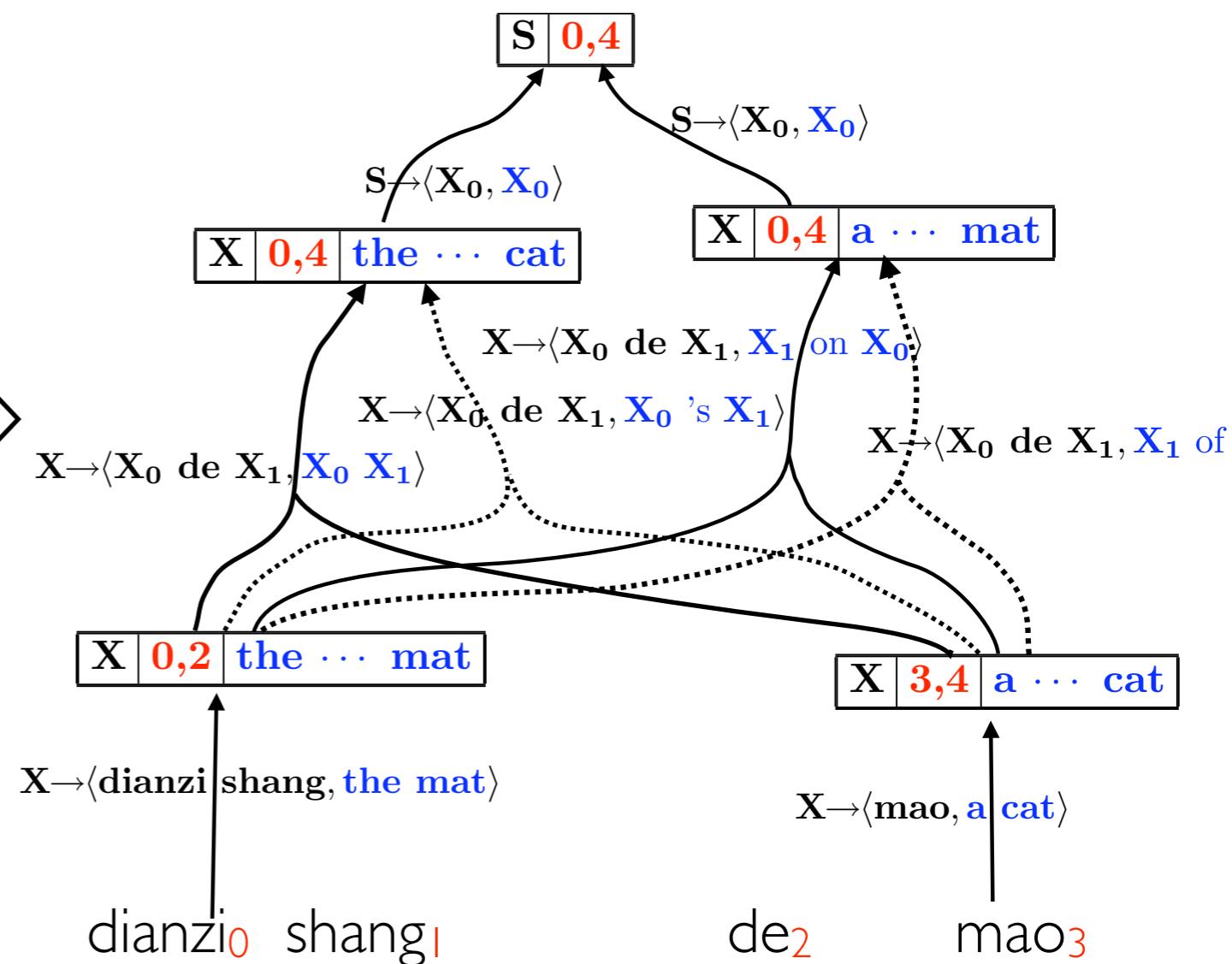
## Sentence-specific decoding

Three steps:

1

Generate a hypergraph

Foreign sentence  $x$   $\rightarrow$  SMT



# Variational Decoding for MT: an Overview

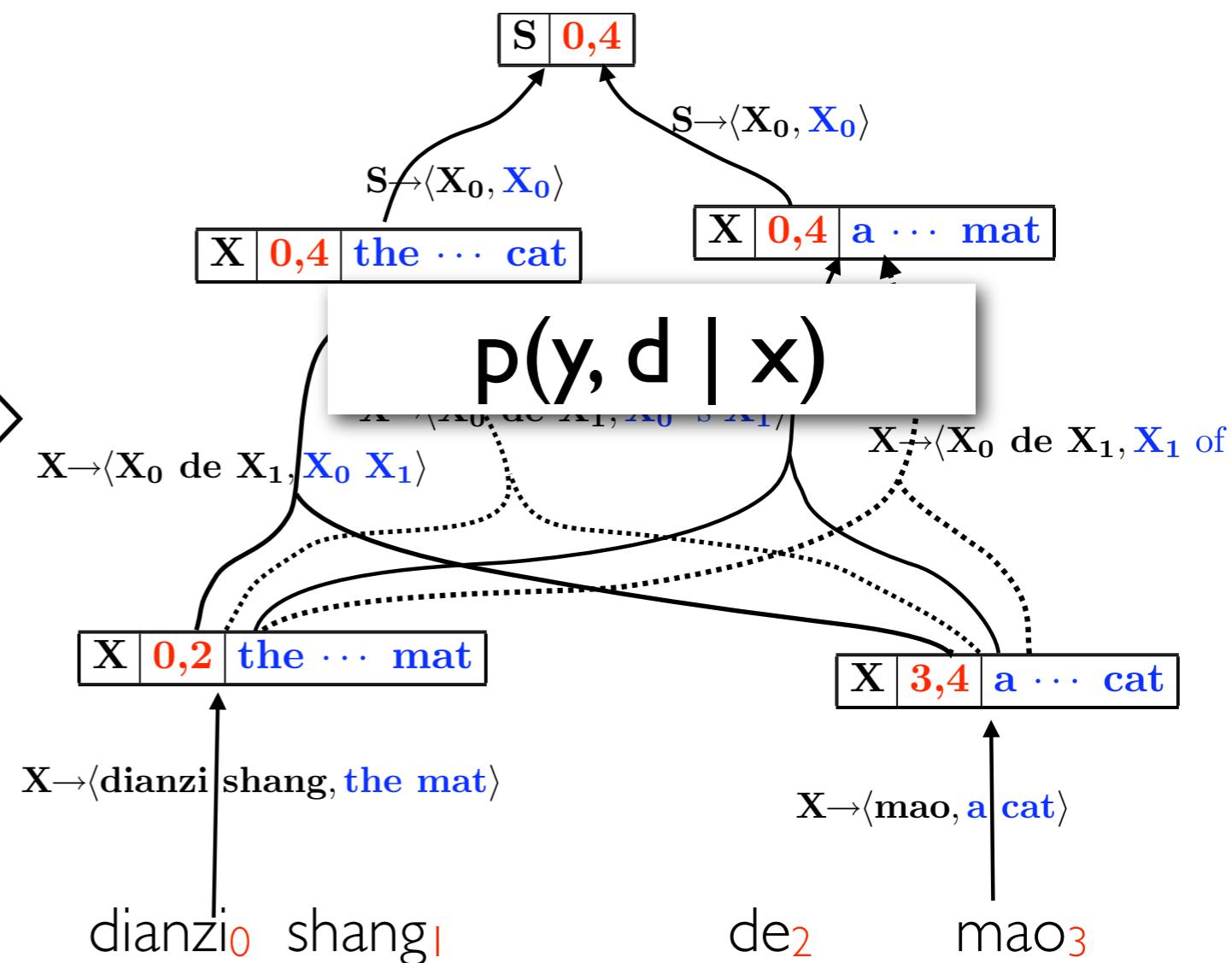
## Sentence-specific decoding

Three steps:

1

Generate a hypergraph

Foreign sentence  $x$   $\rightarrow$  SMT



# Variational Decoding for MT: an Overview

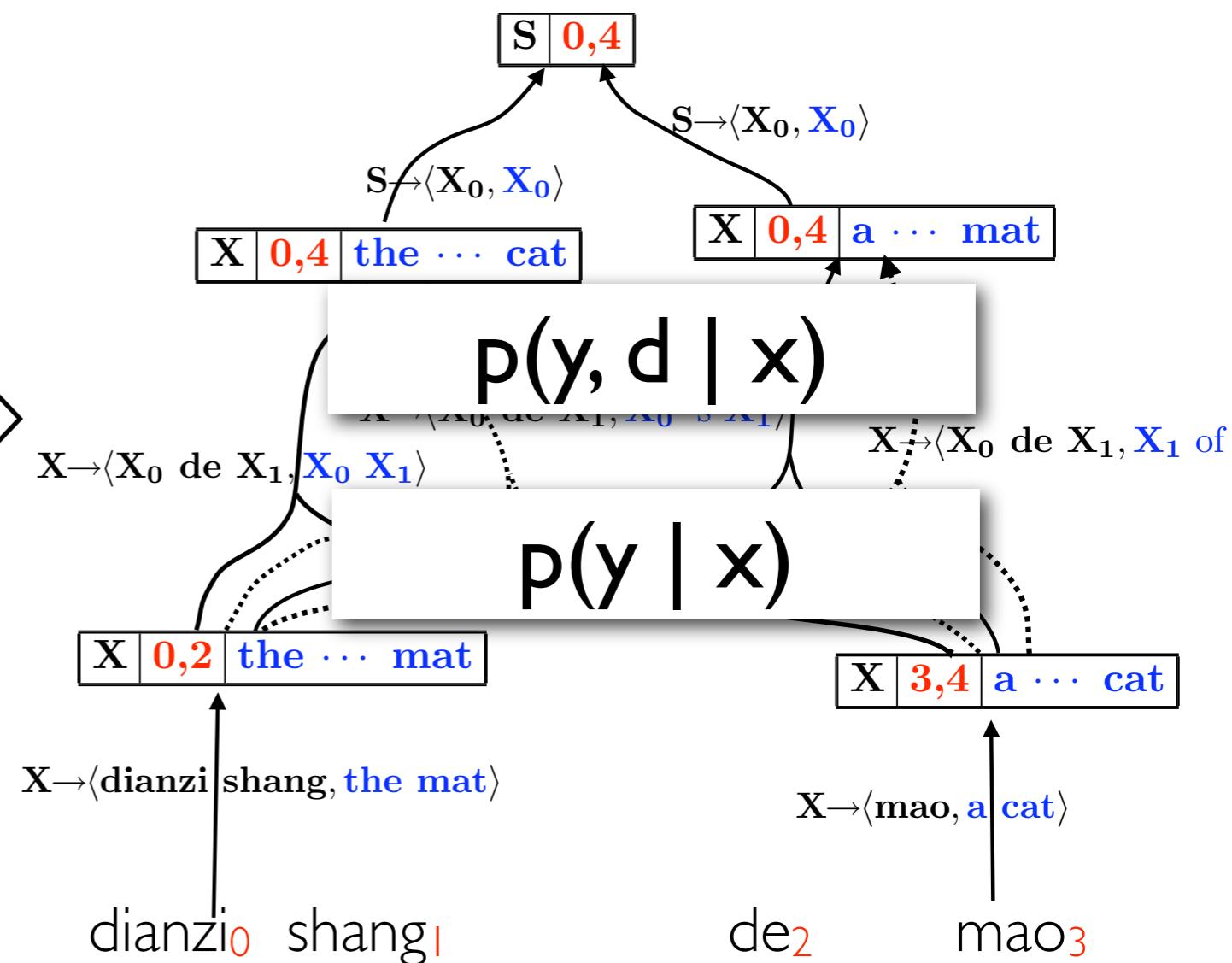
## Sentence-specific decoding

Three steps:

1

Generate a hypergraph

Foreign sentence  $x$   $\rightarrow$  SMT



# Variational Decoding for MT: an Overview

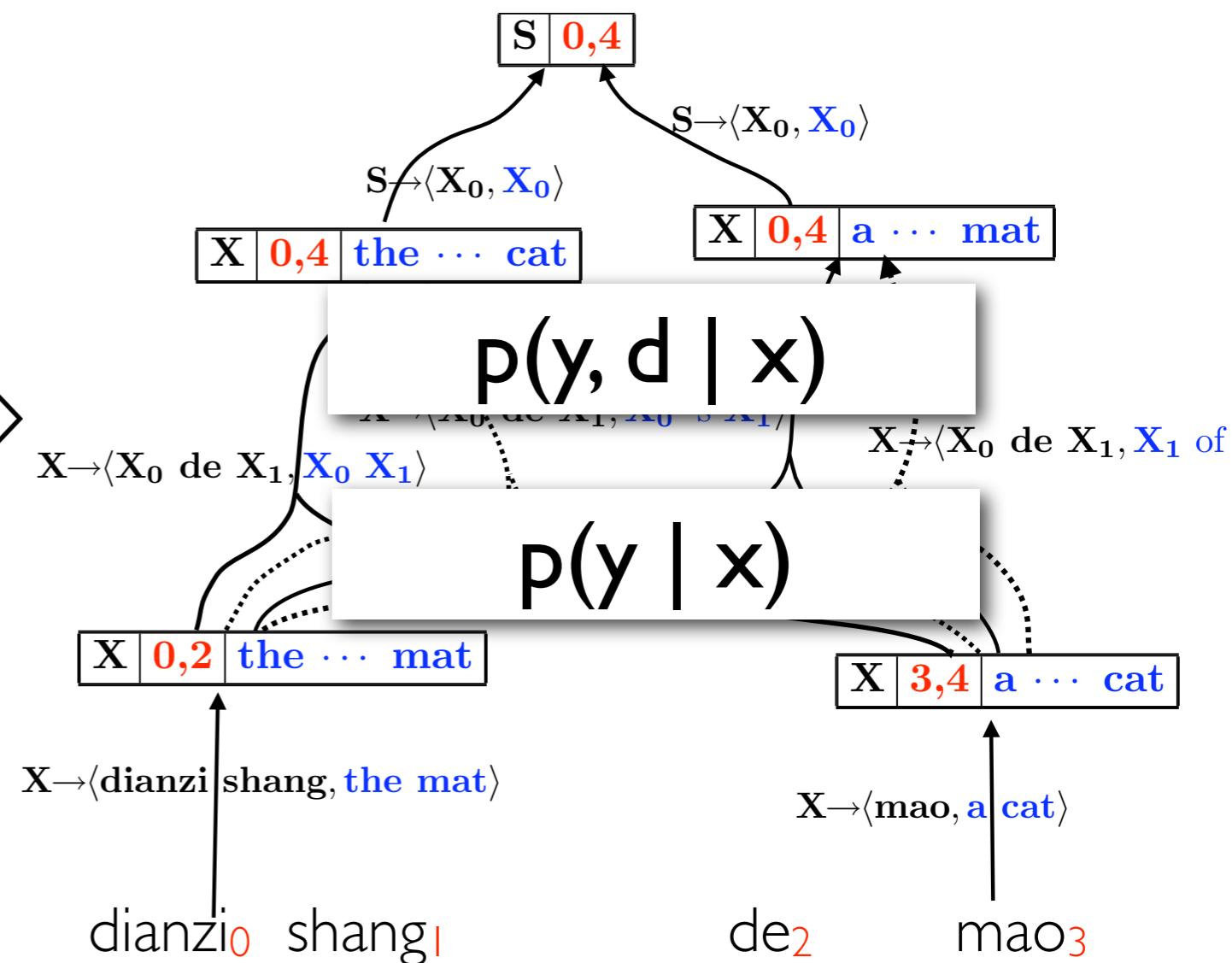
Sentence-specific decoding

Three steps:

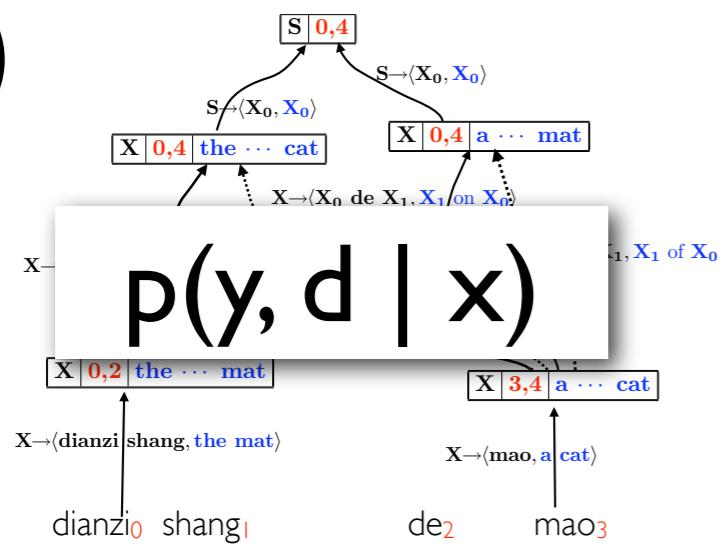
1

Generate a hypergraph

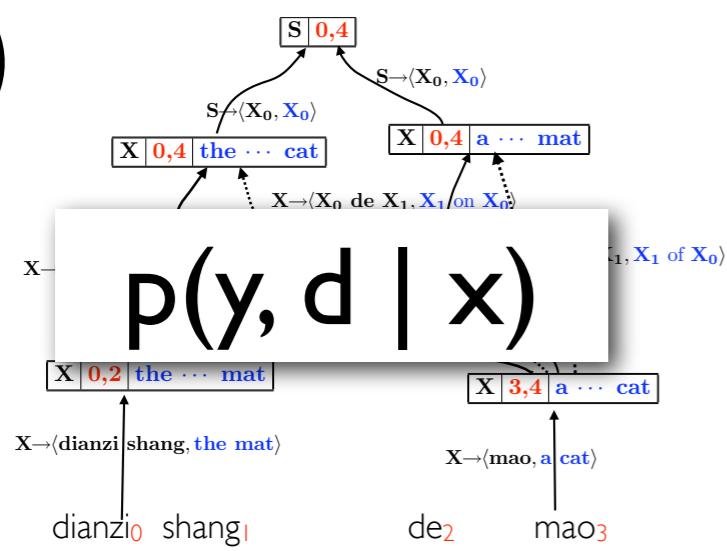
Foreign sentence  $x$   $\rightarrow$  SMT



1

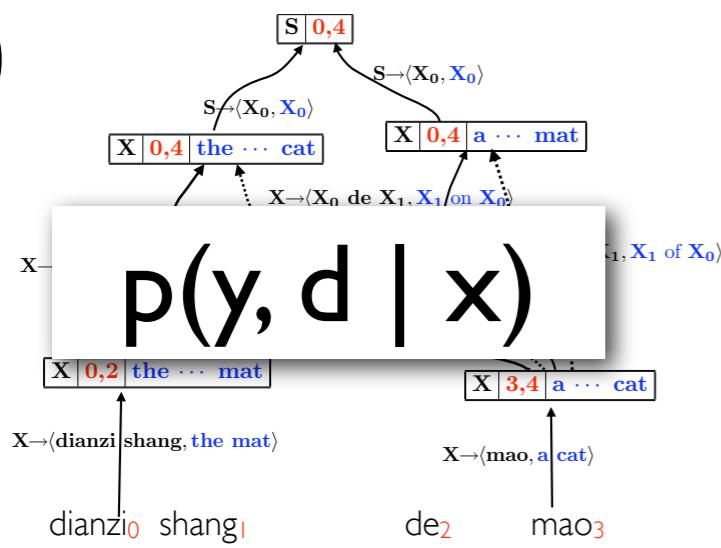


# Generate a hypergraph



## Generate a hypergraph

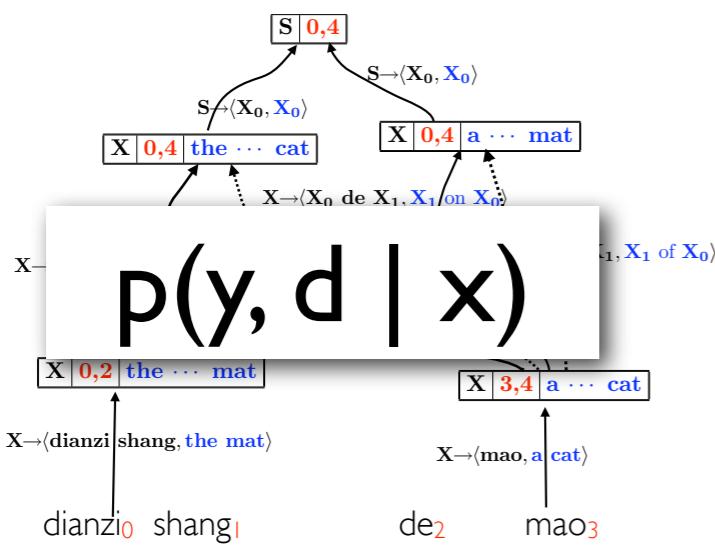
1



# Generate a hypergraph

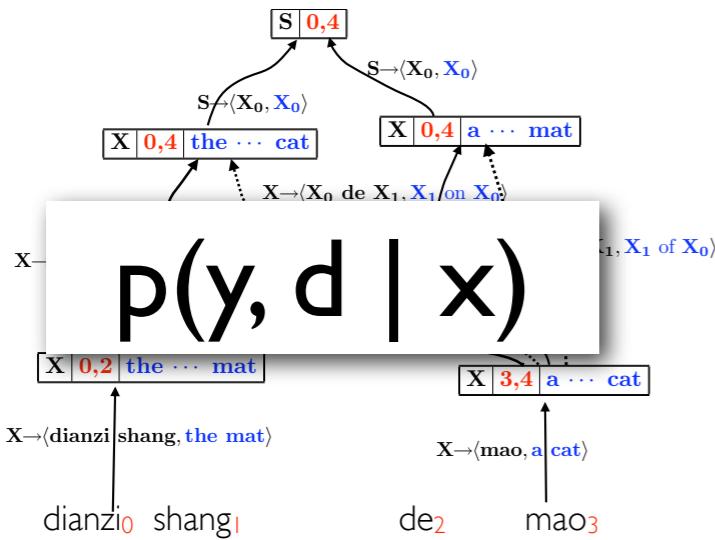
2

1

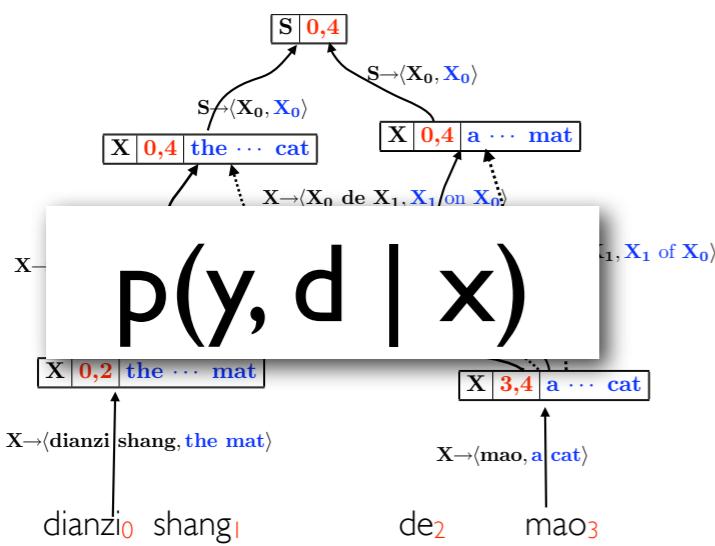


## Generate a hypergraph

2

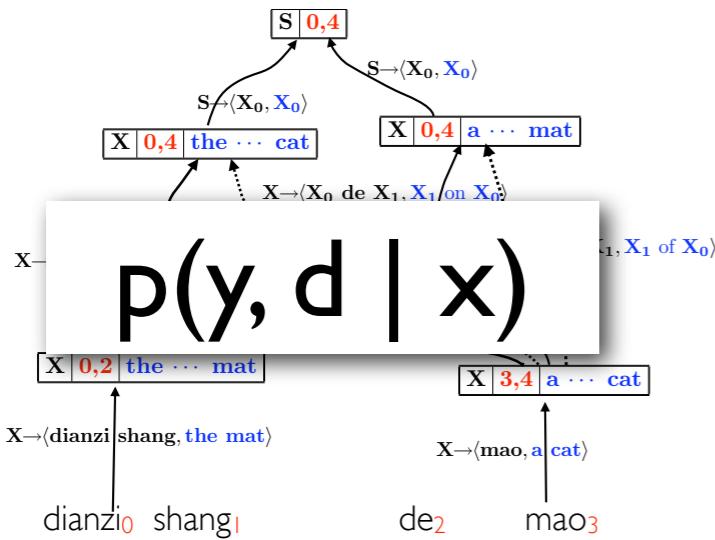


1



## Generate a hypergraph

2

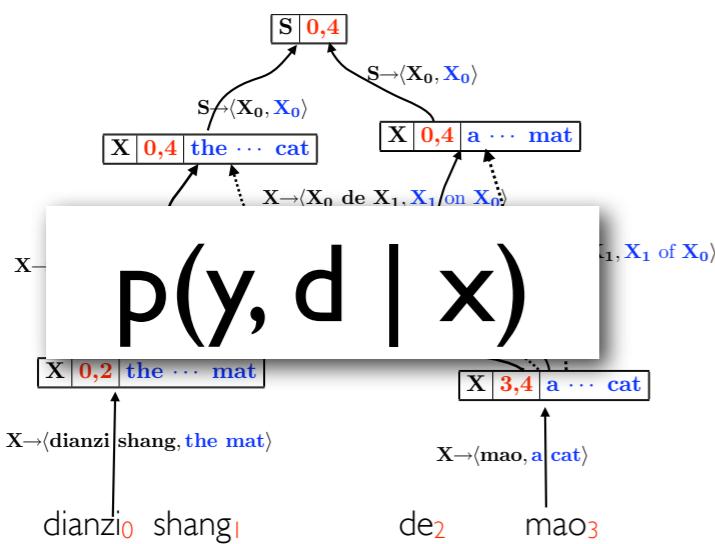


## Estimate a model from the hypergraph



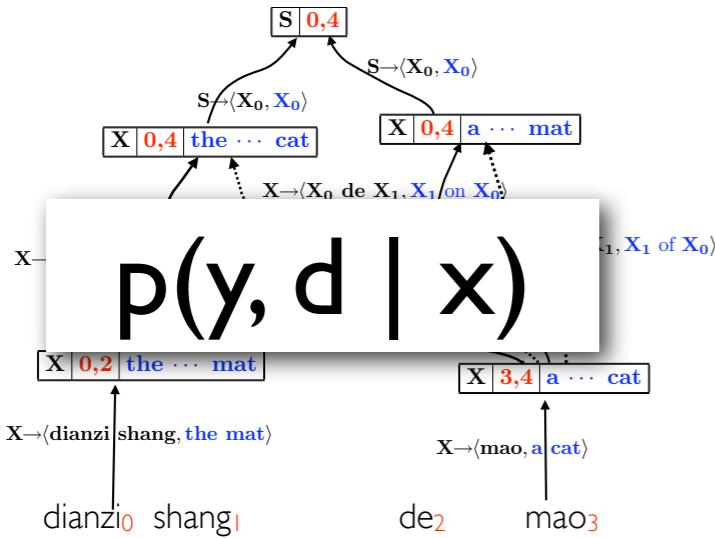
$q^*(y | x)$

1



## Generate a hypergraph

2

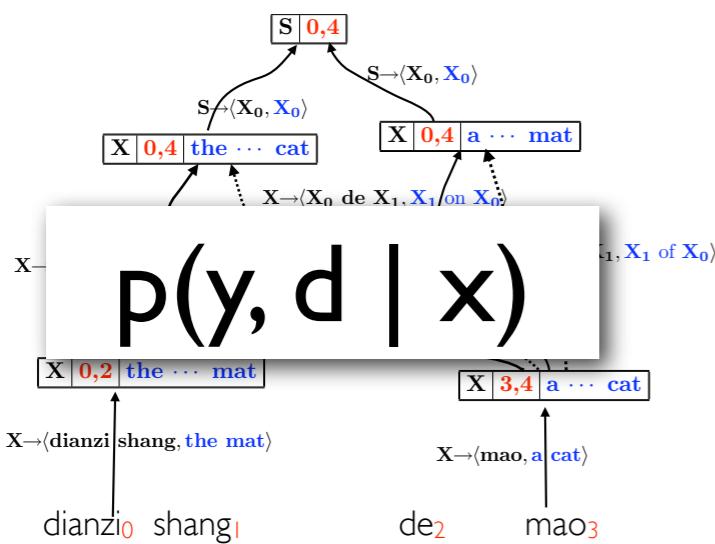


## Estimate a model from the hypergraph

$q^*$  is an n-gram model  
over output strings.

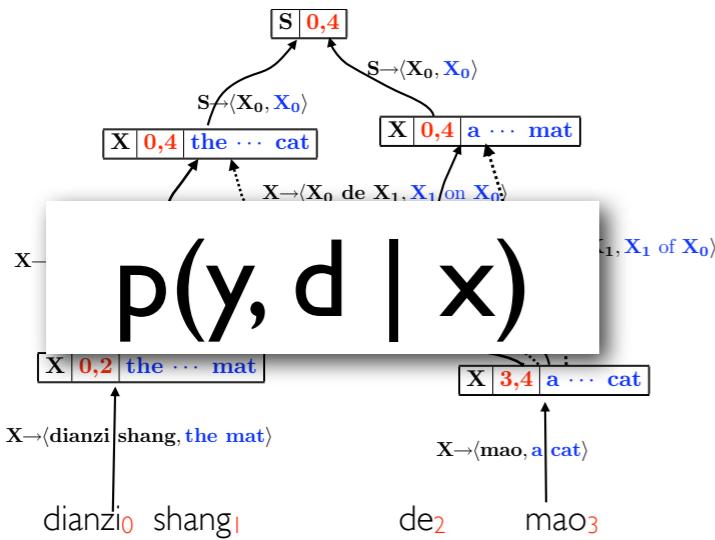
$$q^*(y | x)$$

1



## Generate a hypergraph

2



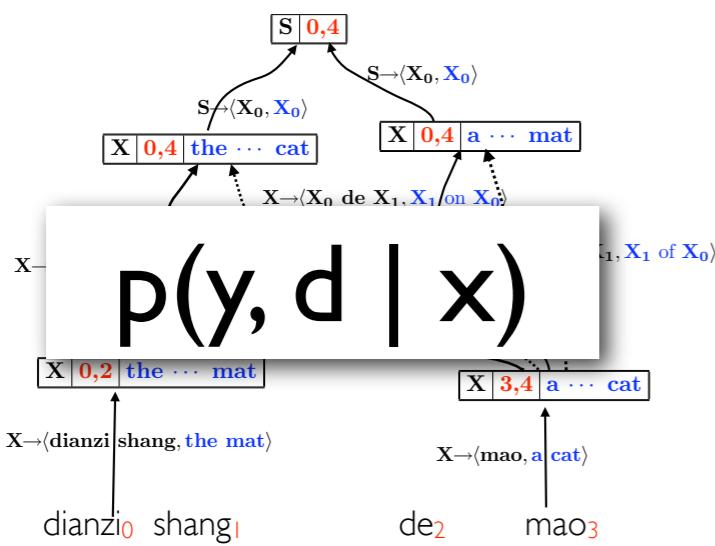
## Estimate a model from the hypergraph

$q^*$  is an n-gram model  
over output strings.

$$q^*(y | x)$$

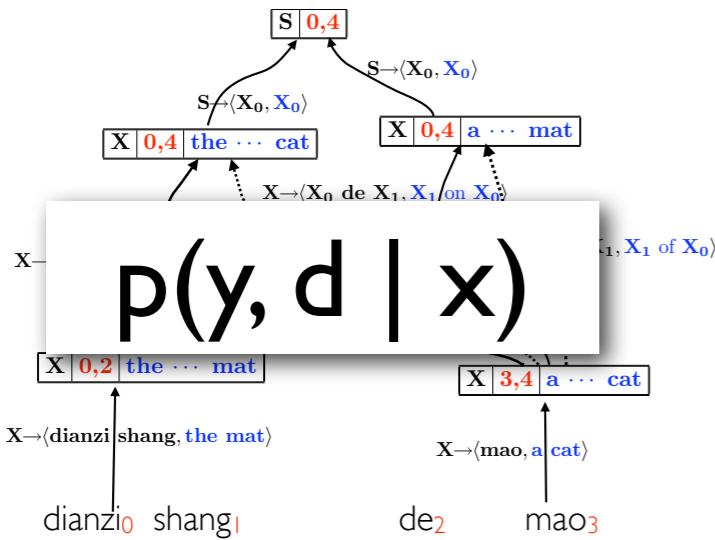
$$\approx \sum_{d \in D(x,y)} P(y,d|x)$$

1



## Generate a hypergraph

2



## Estimate a model from the hypergraph

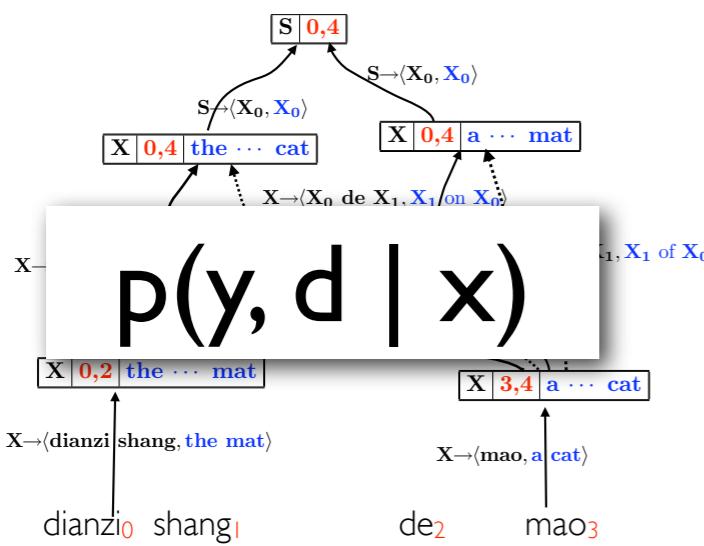
$q^*$  is an n-gram model  
over output strings.

$$q^*(y | x)$$

$$\approx \sum_{d \in D(x,y)} P(y,d|x)$$

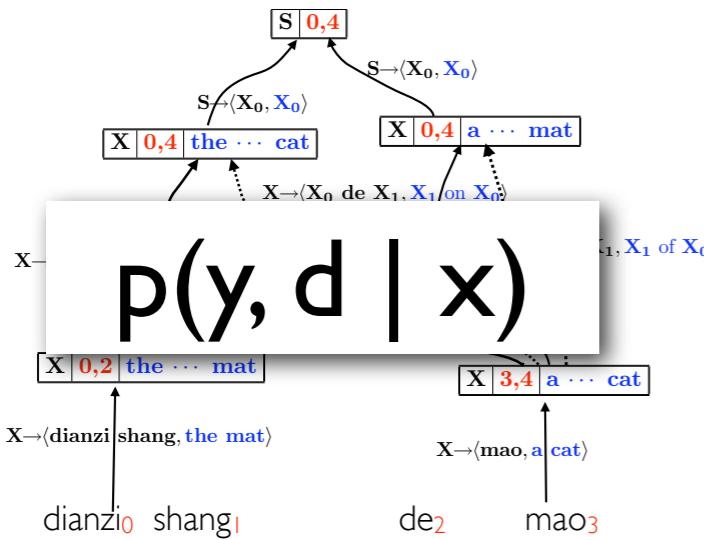
3

1



## Generate a hypergraph

2



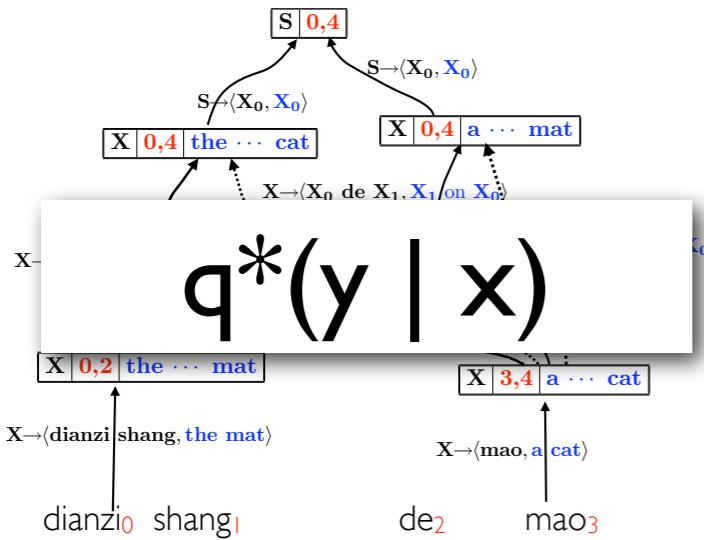
## Estimate a model from the hypergraph

$q^*$  is an n-gram model  
over output strings.

$$q^*(y | x)$$

$$\approx \sum_{d \in D(x,y)} P(y,d|x)$$

3



## Decode using $q^*$ on the hypergraph

# Variational Inference

# Variational Inference

- We want to do inference under  $p$ , but it is intractable

# Variational Inference

- We want to do inference under  $p$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

# Variational Inference

- We want to do inference under  $p$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $q^*$

# Variational Inference

- We want to do inference under  $p$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $q^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$

# Variational Inference

- We want to do inference under  $p$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $q^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$

- Then, we will use  $q^*$  as a surrogate for  $p$  in inference

# Variational Inference

- We want to do inference under  $\mathbf{p}$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $\mathbf{q}^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$

- Then, we will use  $\mathbf{q}^*$  as a surrogate for  $\mathbf{p}$  in inference

$$y^* = \arg \max_y q^*(y | x)$$

# Variational Inference

- We want to do inference under  $\mathbf{P}$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $\mathbf{q}^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$



- Then, we will use  $\mathbf{q}^*$  as a surrogate for  $\mathbf{P}$  in inference

$$y^* = \arg \max_y q^*(y | x)$$

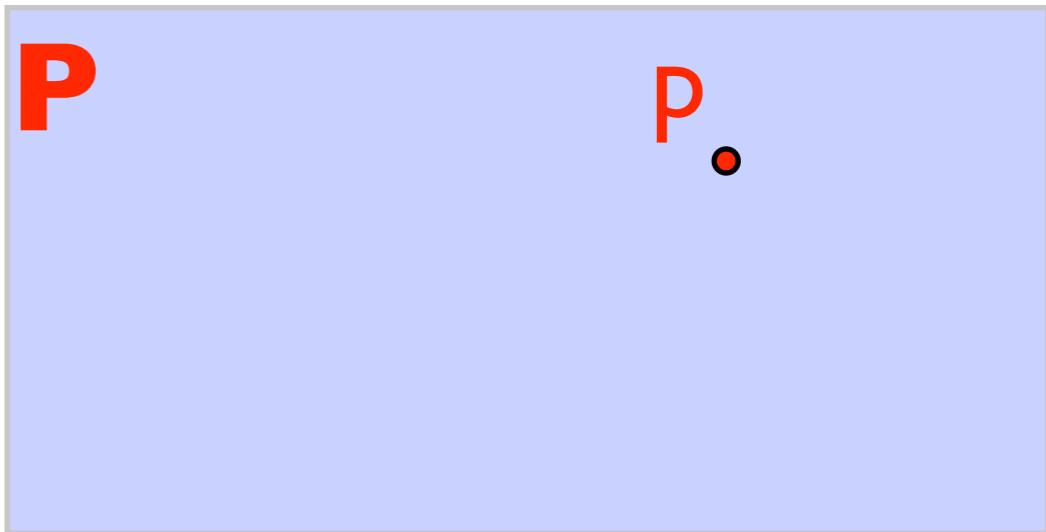
# Variational Inference

- We want to do inference under  $\textcolor{red}{P}$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $\textcolor{red}{q}^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$



- Then, we will use  $\textcolor{red}{q}^*$  as a surrogate for  $\textcolor{red}{P}$  in inference

$$y^* = \arg \max_y q^*(y | x)$$

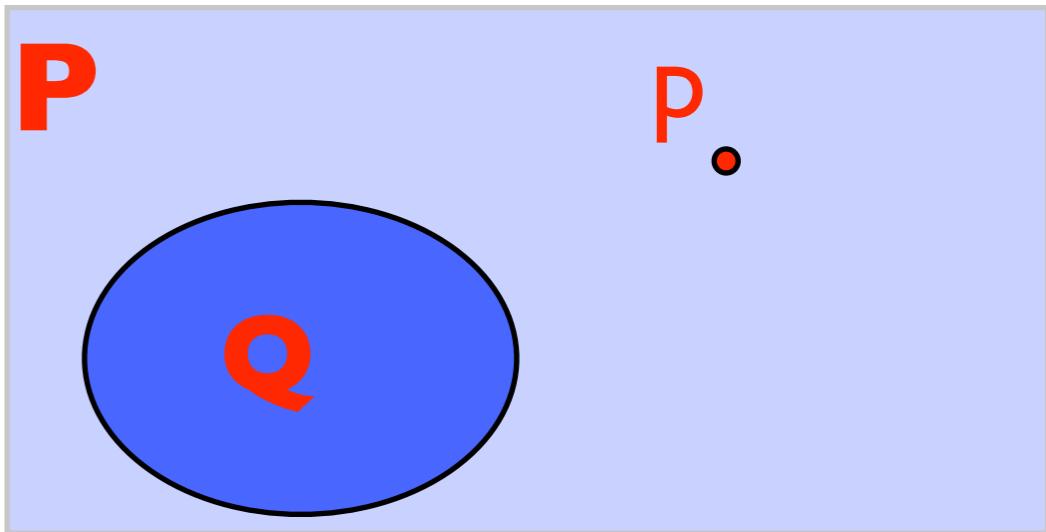
# Variational Inference

- We want to do inference under  $\textcolor{red}{P}$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $\textcolor{red}{q}^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$



- Then, we will use  $\textcolor{red}{q}^*$  as a surrogate for  $\textcolor{red}{P}$  in inference

$$y^* = \arg \max_y q^*(y | x)$$

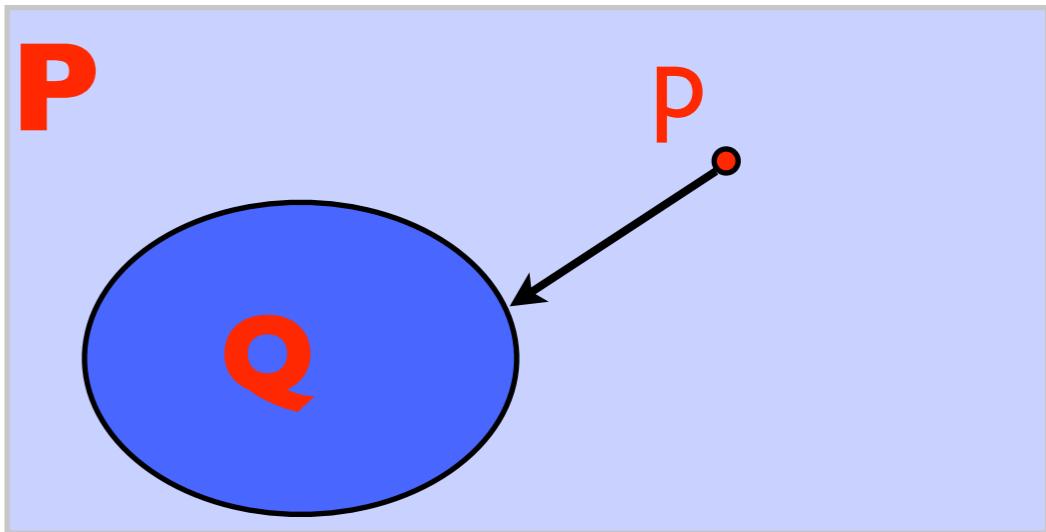
# Variational Inference

- We want to do inference under  $\mathbf{P}$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $\mathbf{q}^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$



- Then, we will use  $\mathbf{q}^*$  as a surrogate for  $\mathbf{P}$  in inference

$$y^* = \arg \max_y q^*(y | x)$$

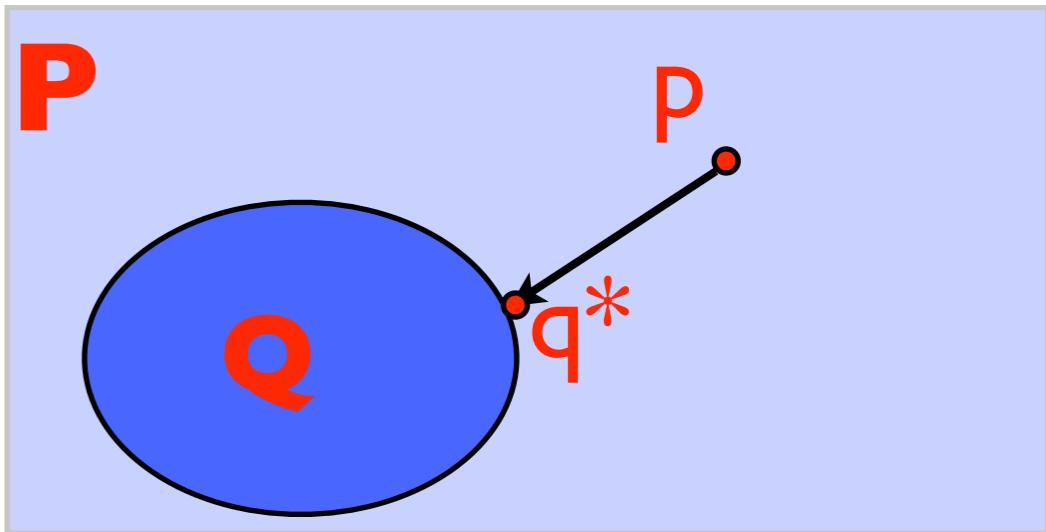
# Variational Inference

- We want to do inference under  $\mathbf{P}$ , but it is intractable

$$y^* = \arg \max_y p(y|x)$$

- Instead, we derive a simpler distribution  $\mathbf{q}^*$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q)$$



- Then, we will use  $\mathbf{q}^*$  as a surrogate for  $\mathbf{P}$  in inference

$$y^* = \arg \max_y q^*(y | x)$$

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}}$$

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \end{aligned}$$

 a family of distributions

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p \log p - p \log q) \end{aligned}$$

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log p - p \log q \end{aligned}$$

 constant

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p \log p - p \log q) \xrightarrow{\text{constant}} \\ &= \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q \end{aligned}$$

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p \log p - p \log q) \xrightarrow{\text{constant}} \\ &= \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q \end{aligned}$$

- Three questions

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p \log p - p \log q) \xrightarrow{\text{constant}} \\ &= \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q \end{aligned}$$

- Three questions

- how to parameterize  $q$ ?

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p \log p - p \log q) \xrightarrow{\text{constant}} \\ &= \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q \end{aligned}$$

- Three questions

- how to parameterize  $q$ ?
- how to estimate  $q^*$ ?

# Variational Approximation

- $q^*$ : an approximation having minimum distance to  $P$

$$\begin{aligned} q^* &= \arg \min_{q \in Q} \text{KL}(p||q) \xrightarrow{\text{a family of distributions}} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log \frac{p}{q} \\ &= \arg \min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p \log p - p \log q) \xrightarrow{\text{constant}} \\ &= \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q \end{aligned}$$

- Three questions

- how to parameterize  $q$ ?
- how to estimate  $q^*$ ?
- how to use  $q^*$  for decoding?

# Parameterization of $q \in Q$

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those *n*-grams appearing in that string

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those *n*-grams appearing in that string

3-gram model

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those *n*-grams appearing in that string

3-gram model

$y: a b c d e f$

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those *n*-grams appearing in that string

3-gram model

$y: a b c d e f$

$$q(y) = q(a) \cdot q(b|a) \cdot q(c|ab) \cdot q(d|bc) \cdot q(e|cd) \cdot q(f|de)$$

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those  $n$ -grams appearing in that string

3-gram model

$y: a b c d e f$

$$q(y) = q(a) \cdot q(b|a) \cdot q(c|ab) \cdot q(d|bc) \cdot q(e|cd) \cdot q(f|de)$$

Other ways of parameterizations are possible!

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those *n*-grams appearing in that string

3-gram model

$y: a b c d e f$

$$q(y) = q(a) \cdot q(b|a) \cdot q(c|ab) \cdot q(d|bc) \cdot q(e|cd) \cdot q(f|de)$$

# Parameterization of $q \in Q$

- Naturally, we parameterize  $q$  as an  $n$ -gram model
  - The probability of a *string* is a product of the probabilities of those  $n$ -grams appearing in that string

3-gram model

$y: a b c d e f$

$$q(y) = q(a) \cdot q(b|a) \cdot q(c|ab) \cdot q(d|bc) \cdot q(e|cd) \cdot q(f|de)$$

how to estimate these n-gram probabilities?

# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE) where  $p$  is the empirical distribution

# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE)  
where  $p$  is the empirical distribution

But in our case,  $p$  is defined **not** by a **corpus**, but by  
a **hypergraph** for a given test sentence!

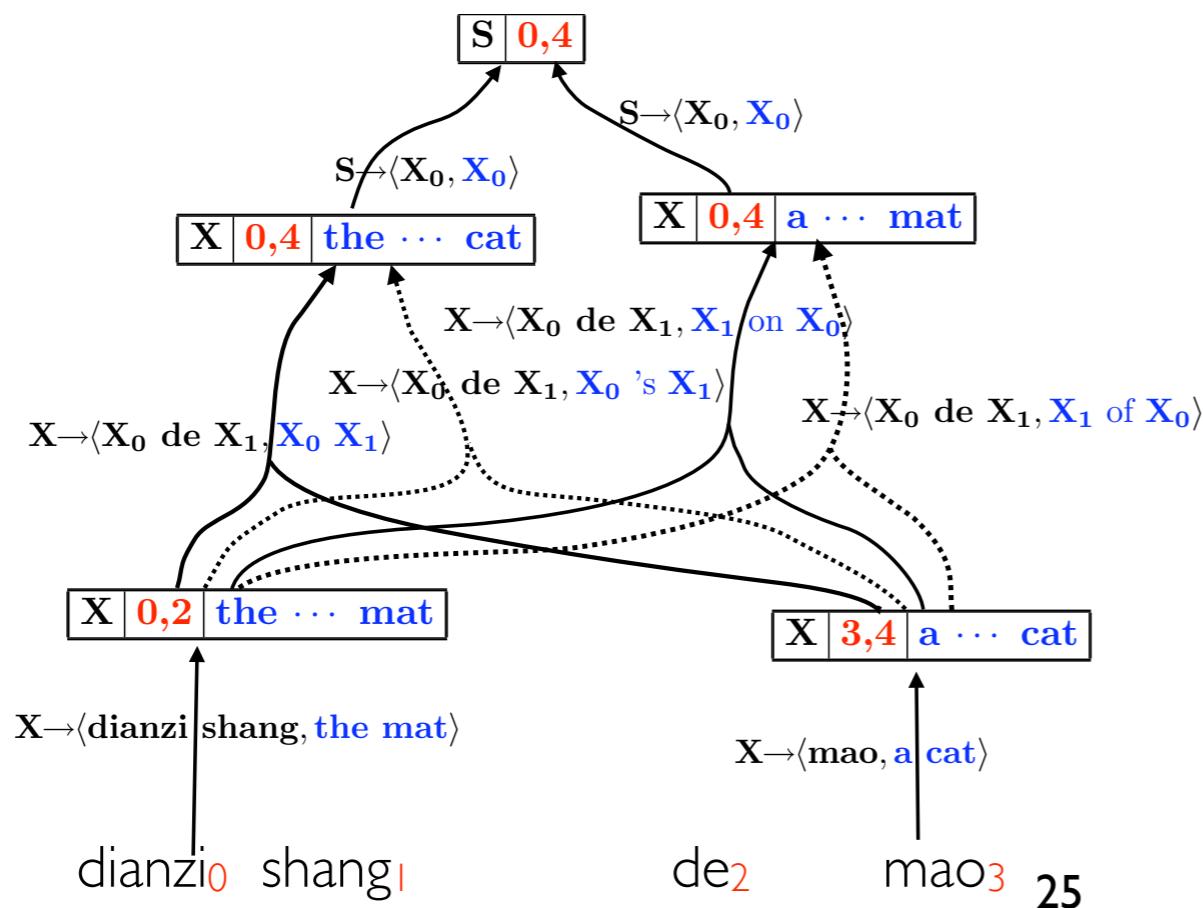
# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE)  
where  $p$  is the empirical distribution

But in our case,  $p$  is defined **not** by a **corpus**, but by a **hypergraph** for a given test sentence!



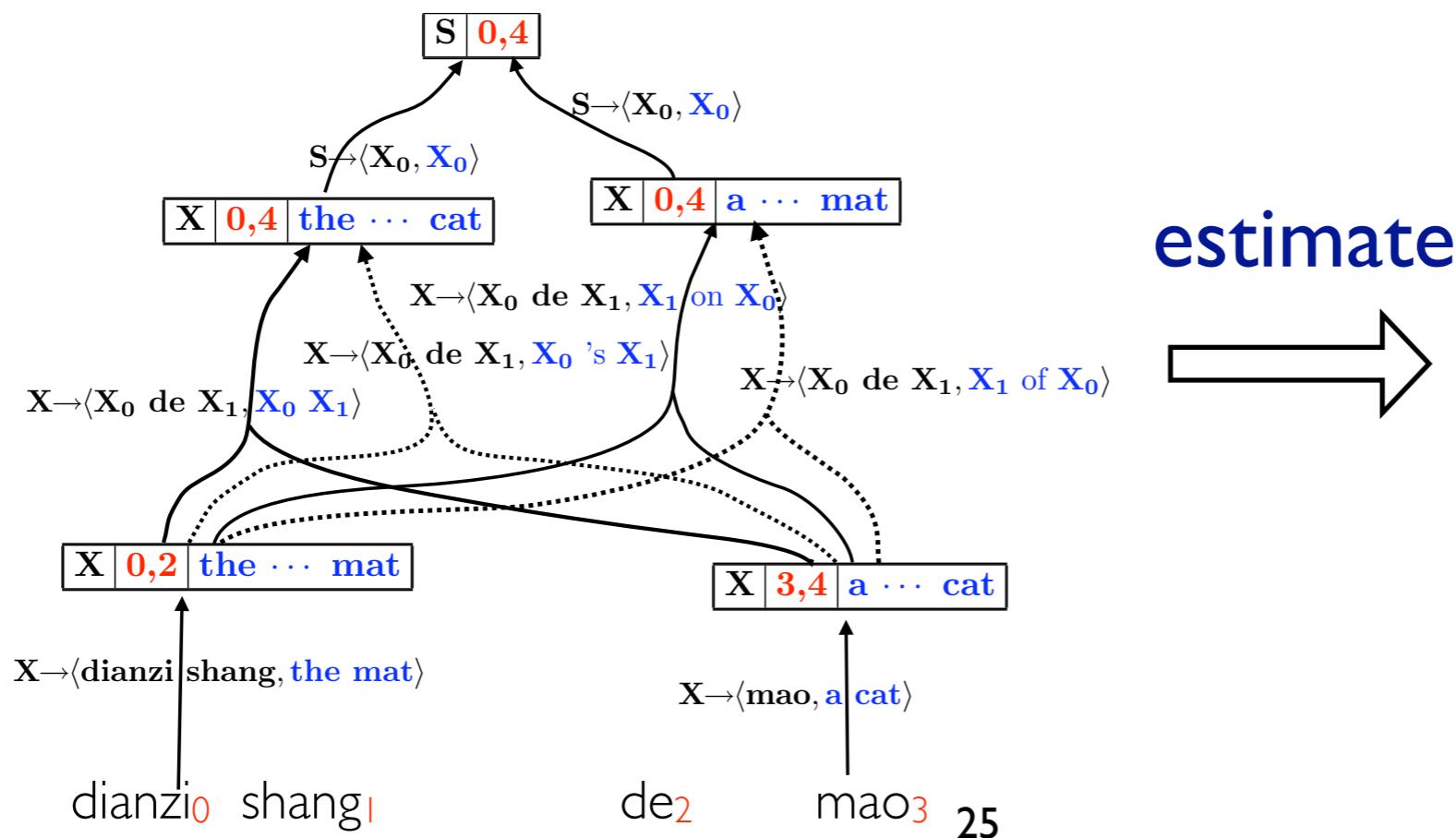
# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE)  
where  $p$  is the empirical distribution

But in our case,  $p$  is defined **not** by a **corpus**, but by a **hypergraph** for a given test sentence!



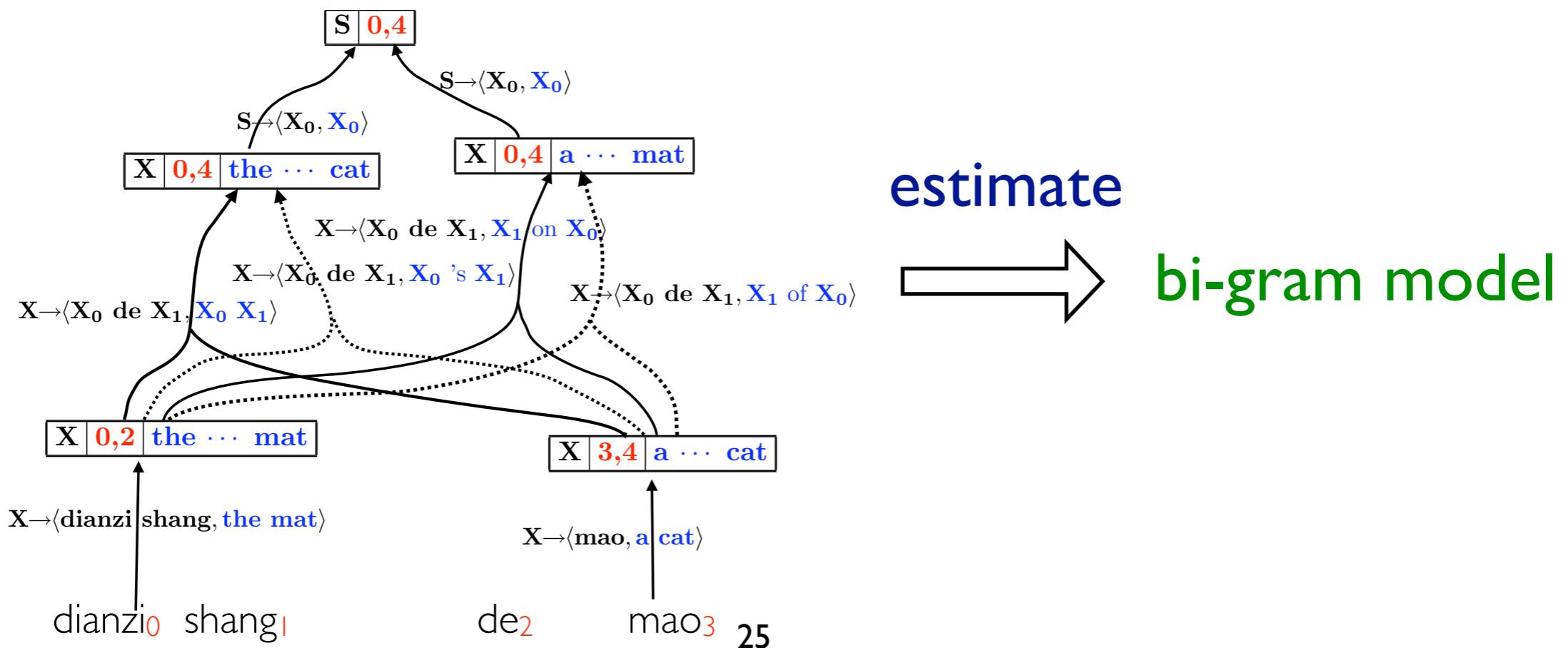
# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE) where  $p$  is the empirical distribution

But in our case,  $p$  is defined **not** by a **corpus**, but by a **hypergraph** for a given test sentence!



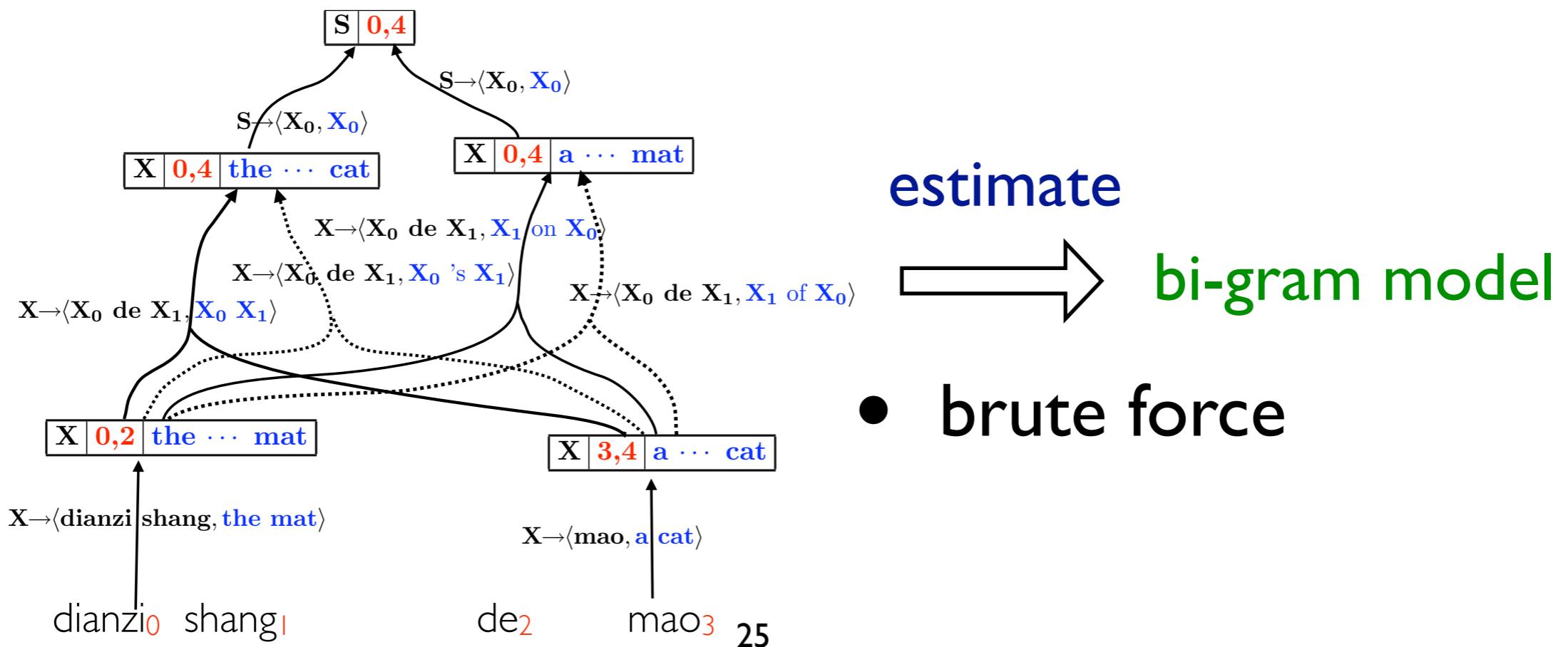
# Estimation of $q^* \in Q$

- Variational approximation

$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE) where  $p$  is the empirical distribution

But in our case,  $p$  is defined **not** by a **corpus**, but by a **hypergraph** for a given test sentence!



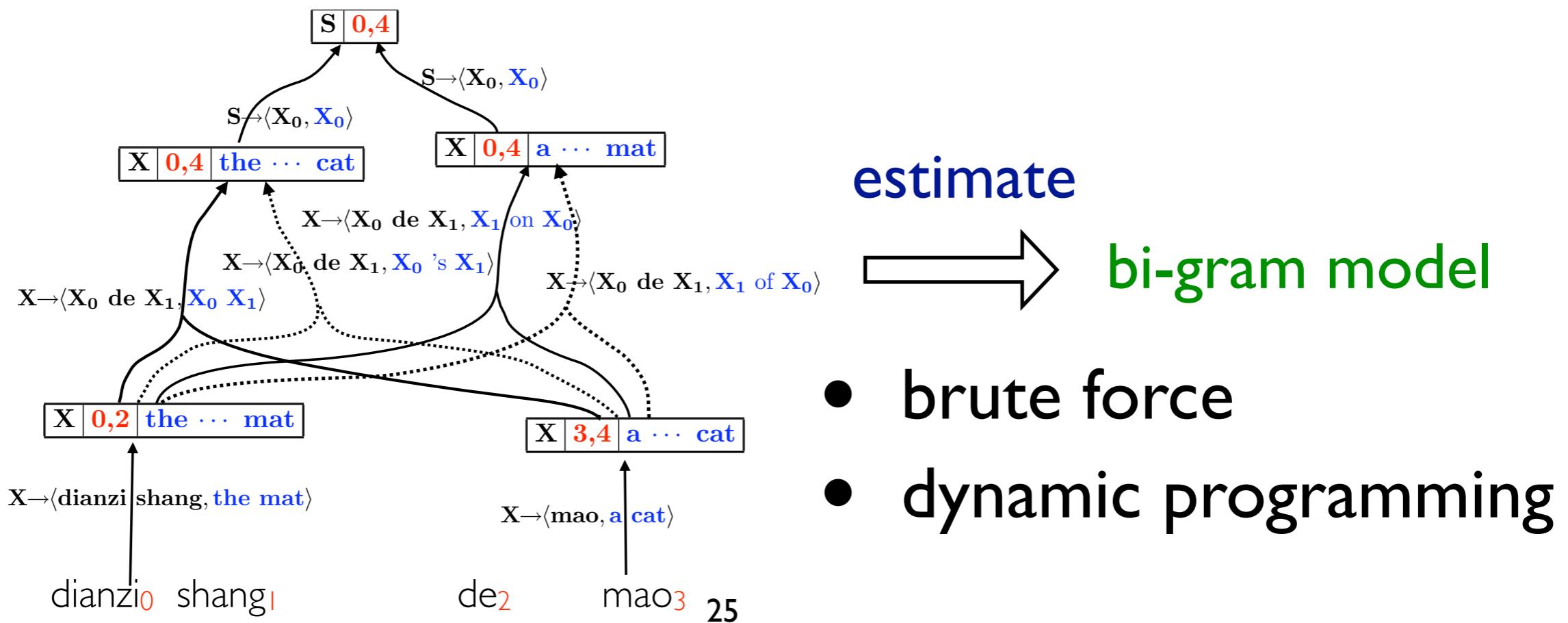
# Estimation of $q^* \in Q$

- Variational approximation

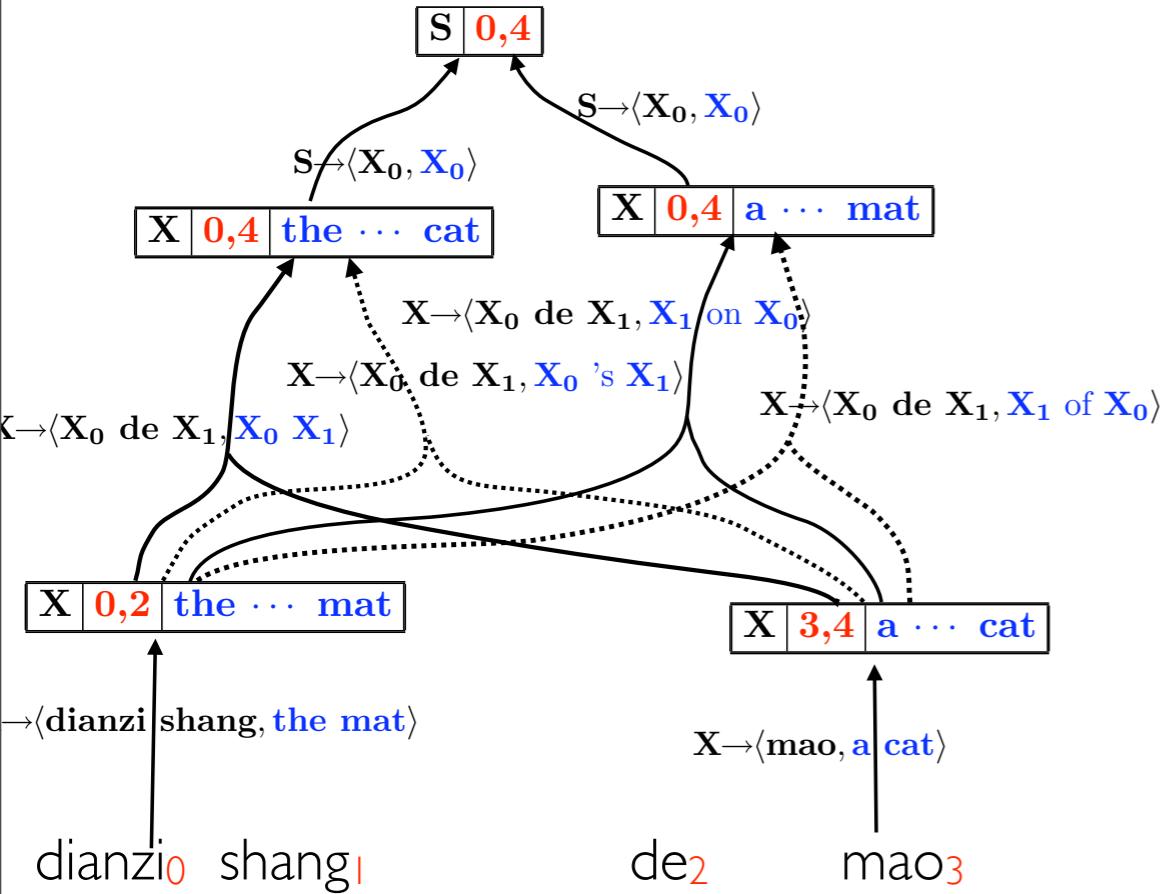
$$q^* = \arg \max_{q \in Q} \sum_{y \in \text{Trans}(x)} p \log q$$

- $q^*$  is a maximum likelihood estimate (MLE)  
where  $p$  is the empirical distribution

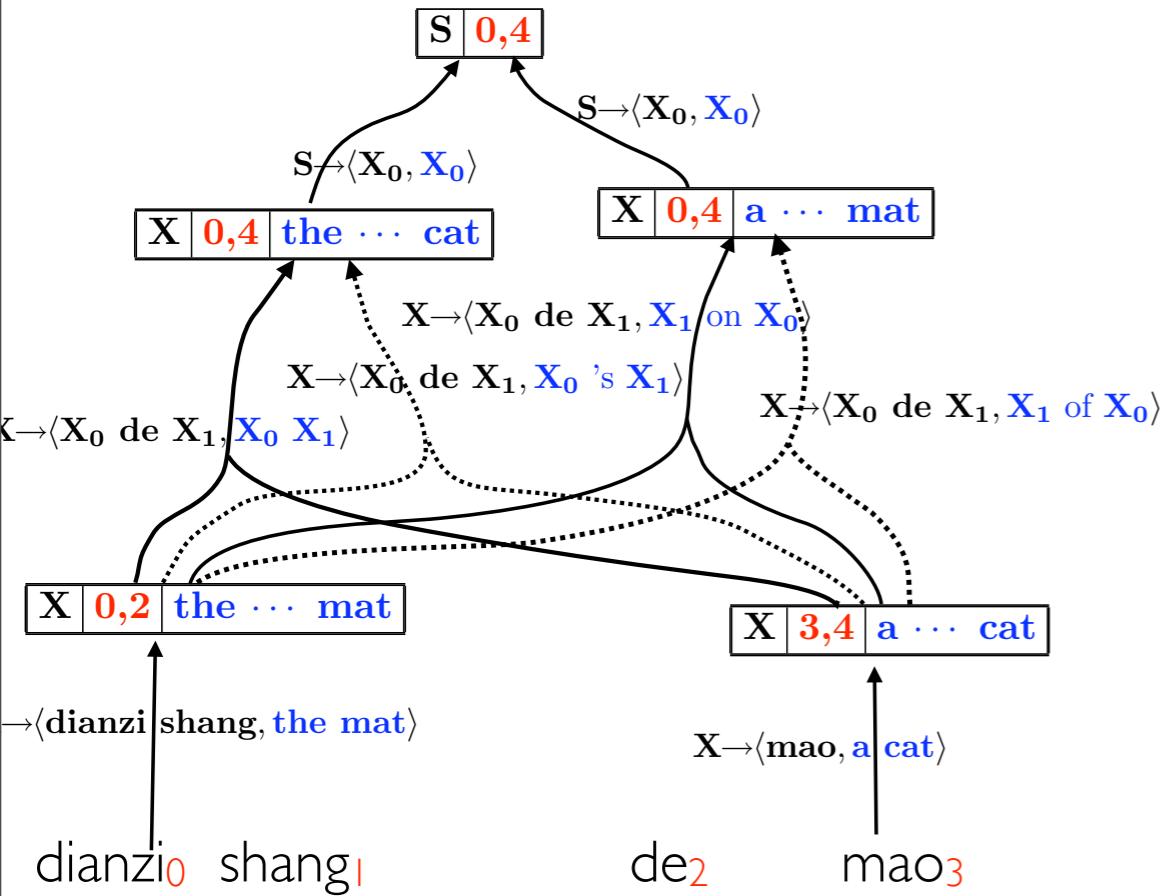
But in our case,  $p$  is defined **not** by a **corpus**, but by a **hypergraph** for a given test sentence!



# Estimating $q^*$ from a hypergraph: brute force

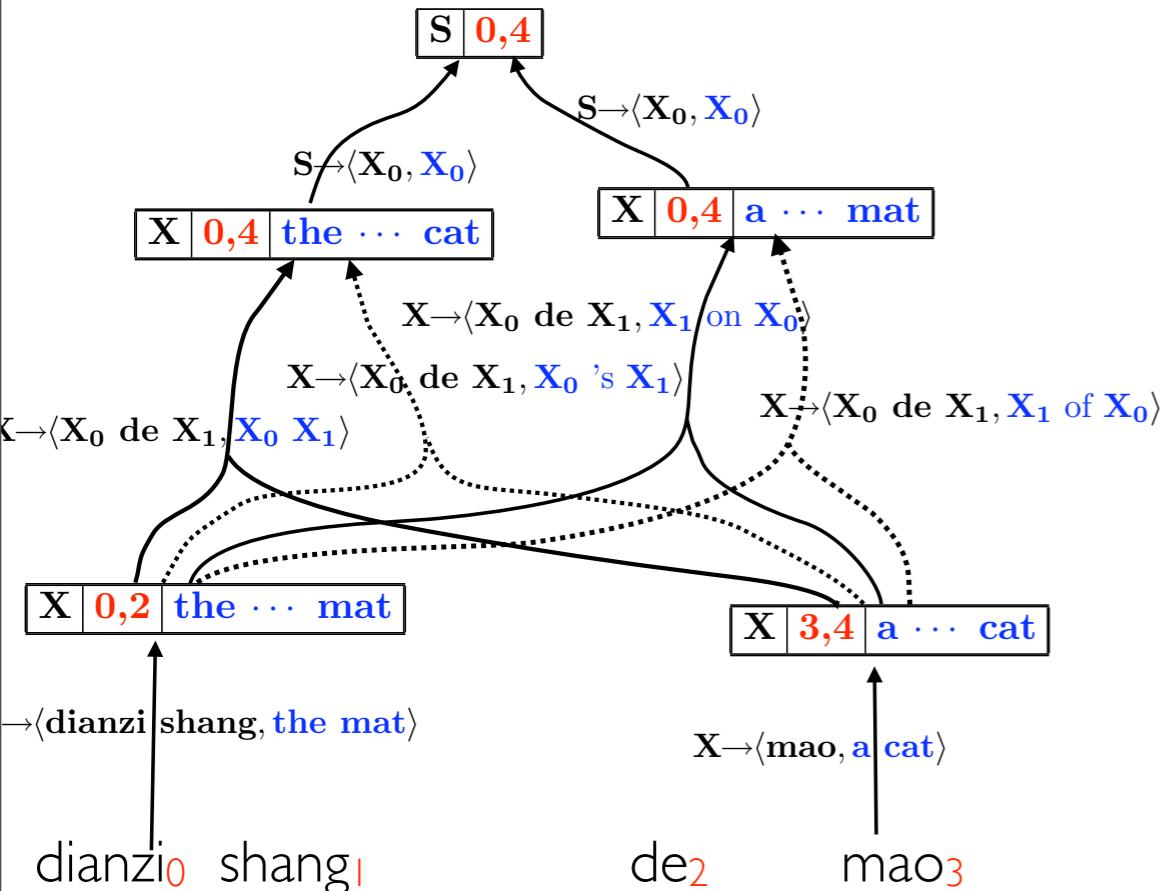


# Estimating $q^*$ from a hypergraph: brute force



Bi-gram estimation:

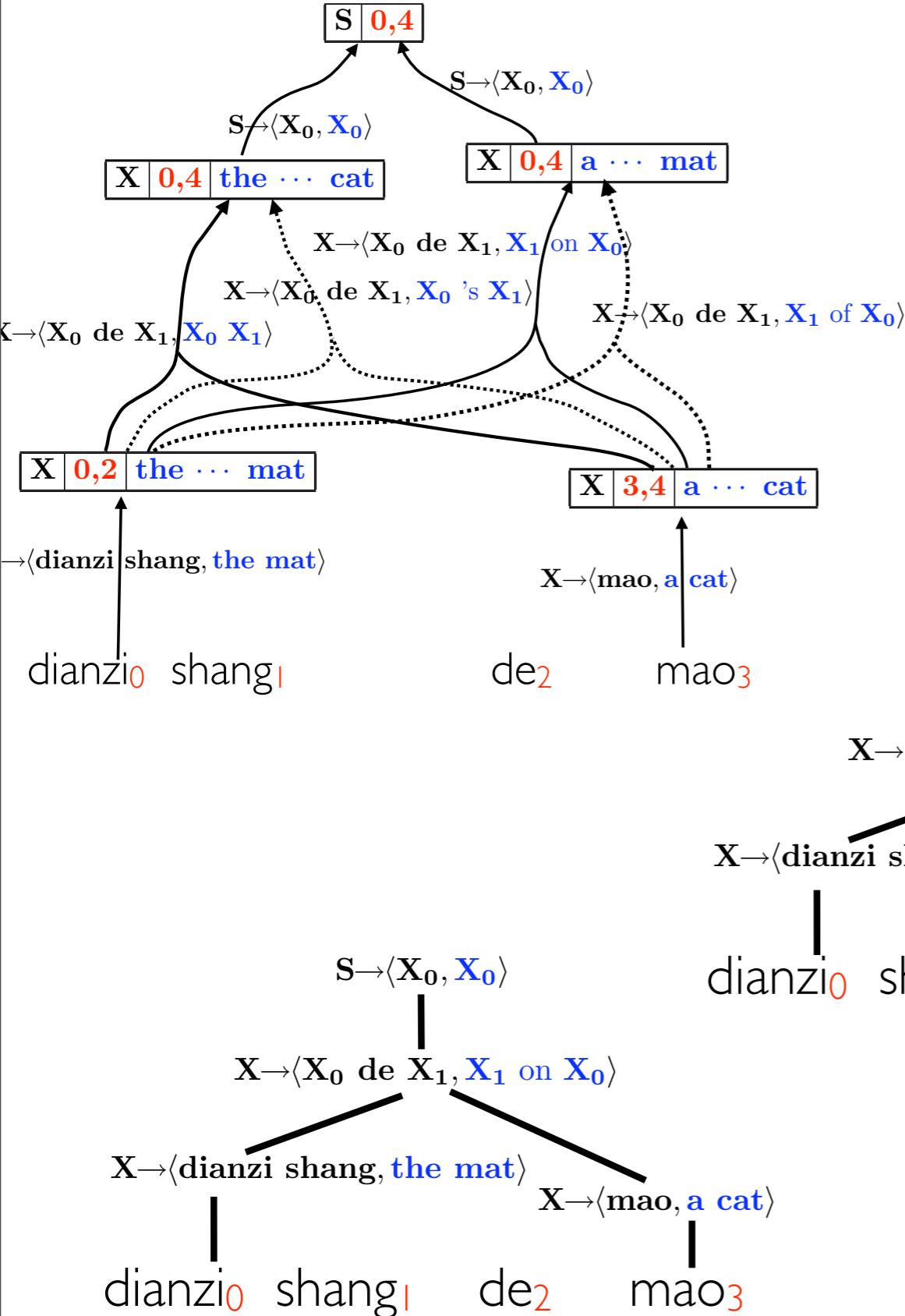
# Estimating $q^*$ from a hypergraph: brute force



Bi-gram estimation:

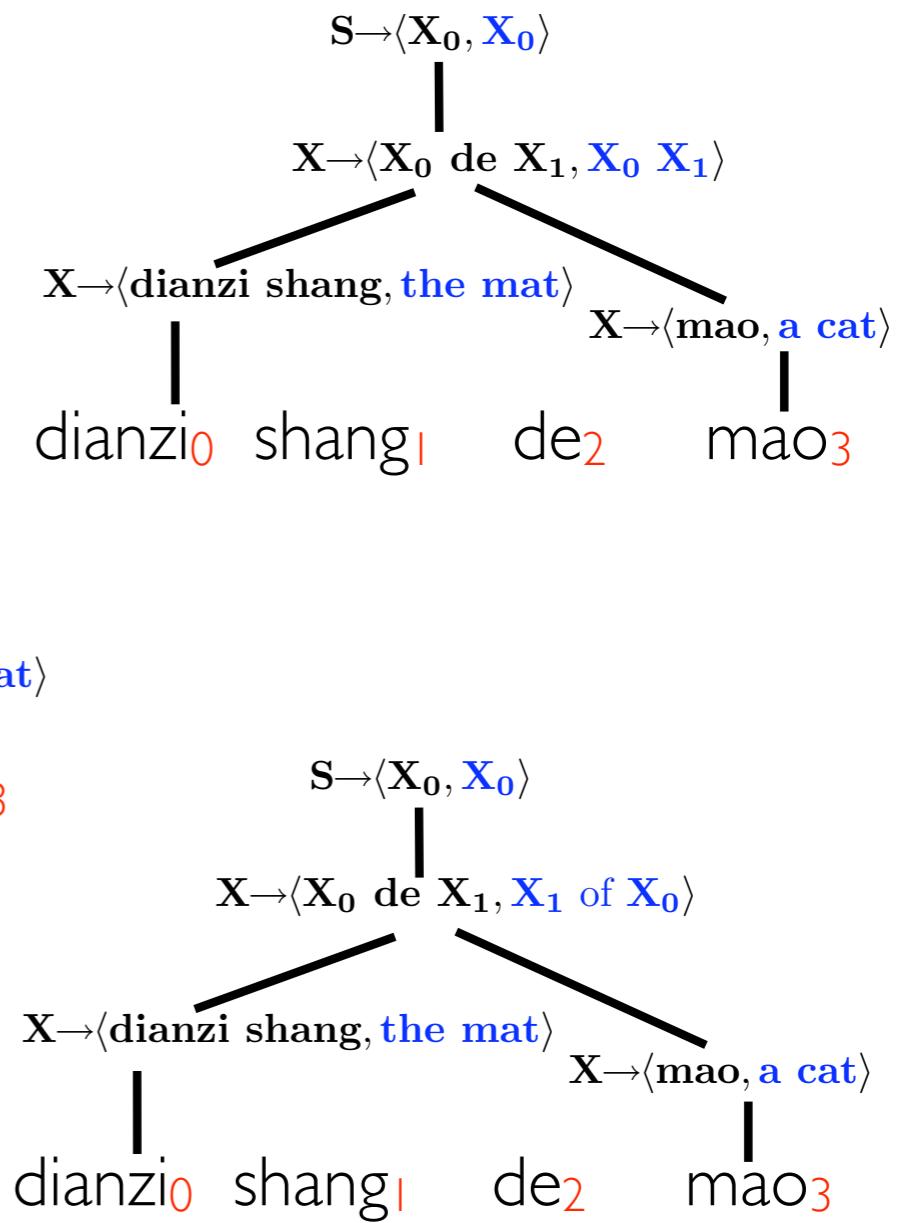
► unpack the hypergraph

# Estimating $q^*$ from a hypergraph: brute force

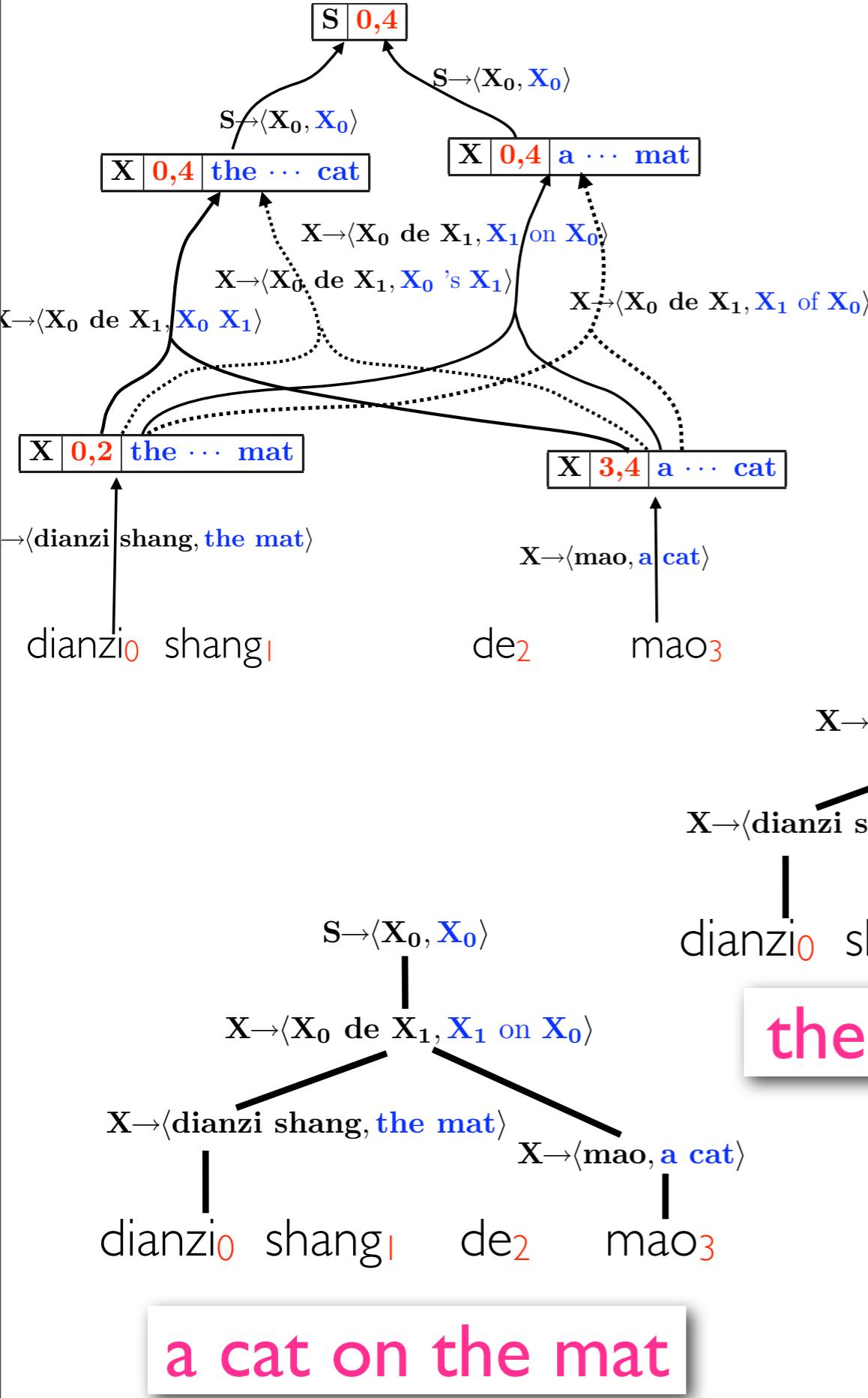


Bi-gram estimation:

► unpack the hypergraph

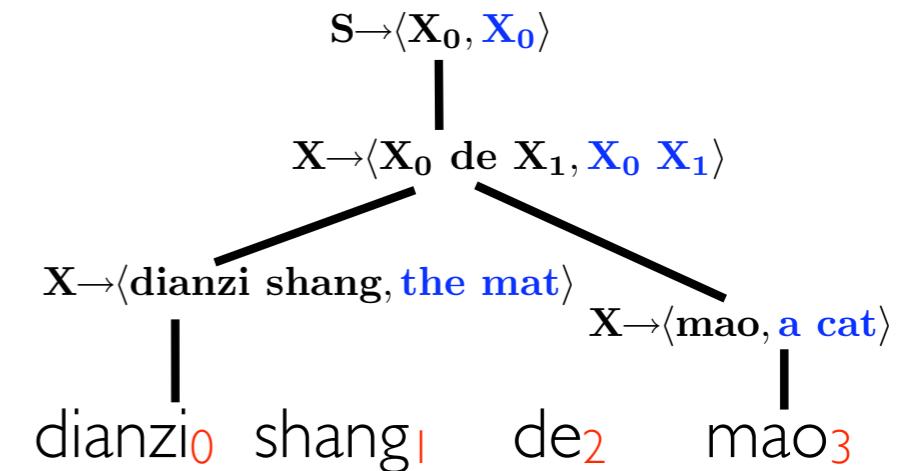


# Estimating $q^*$ from a hypergraph: brute force

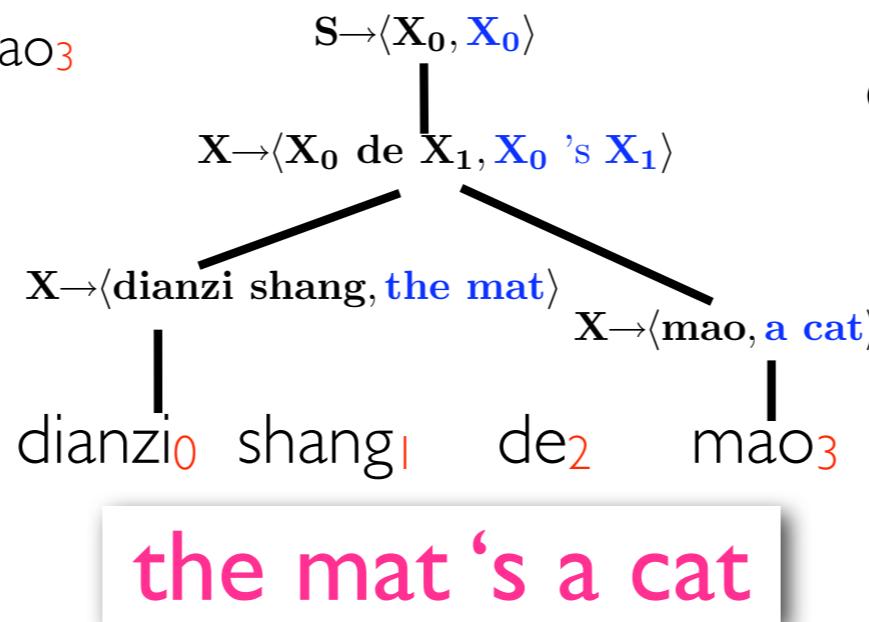


Bi-gram estimation:

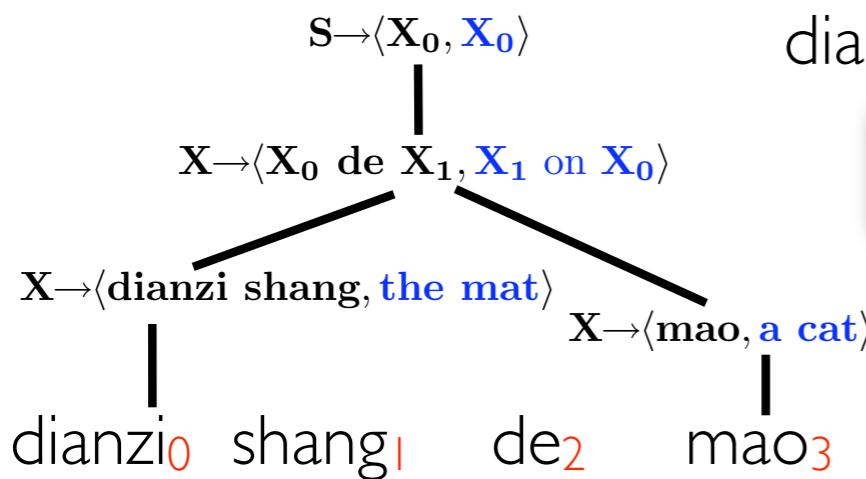
► unpack the hypergraph



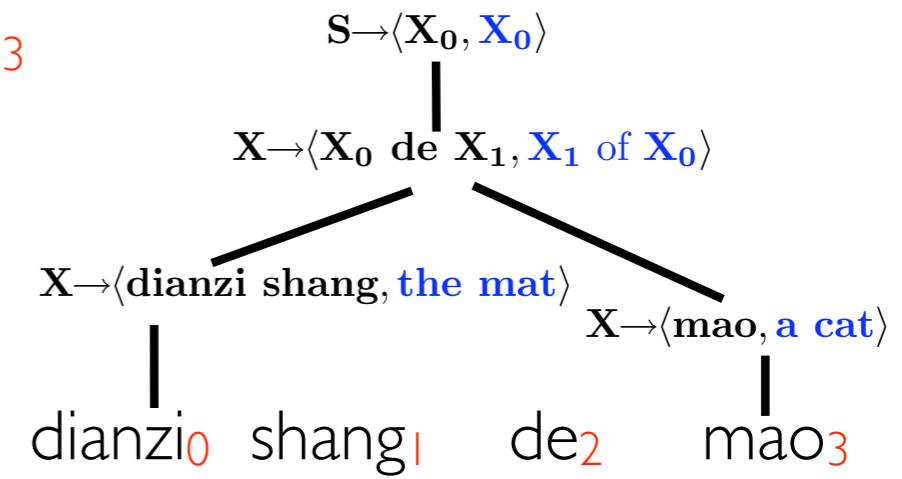
the mat a cat



the mat 's a cat

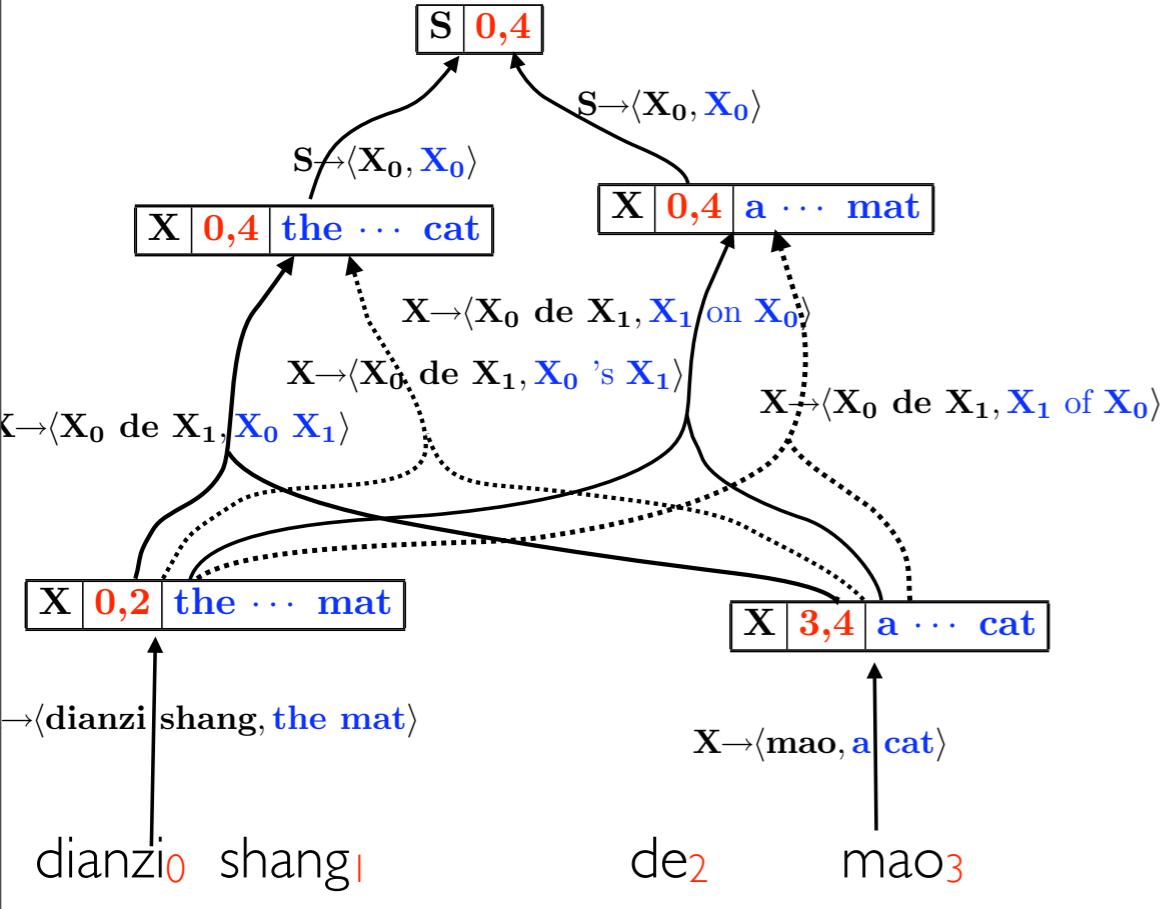


a cat on the mat



a cat of the mat

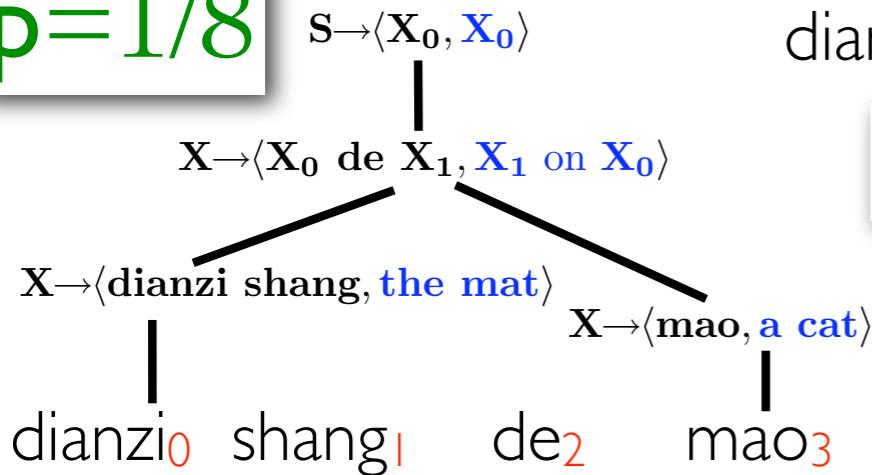
# Estimating $q^*$ from a hypergraph: brute force



Bi-gram estimation:

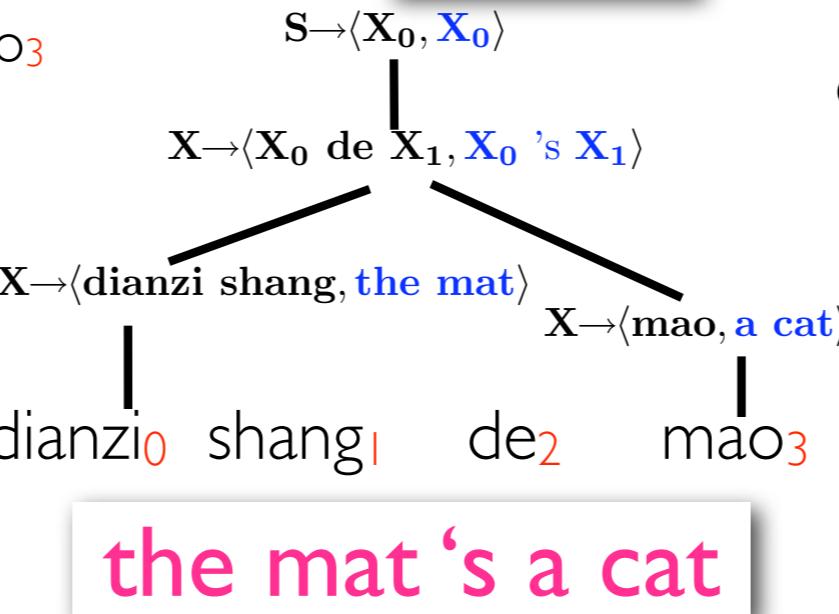
► unpack the hypergraph

$p=1/8$



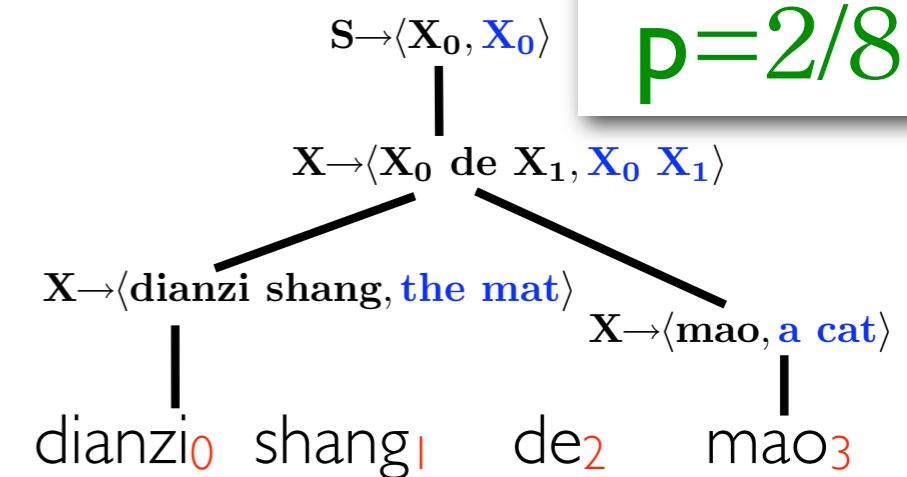
a cat on the mat

$p=3/8$



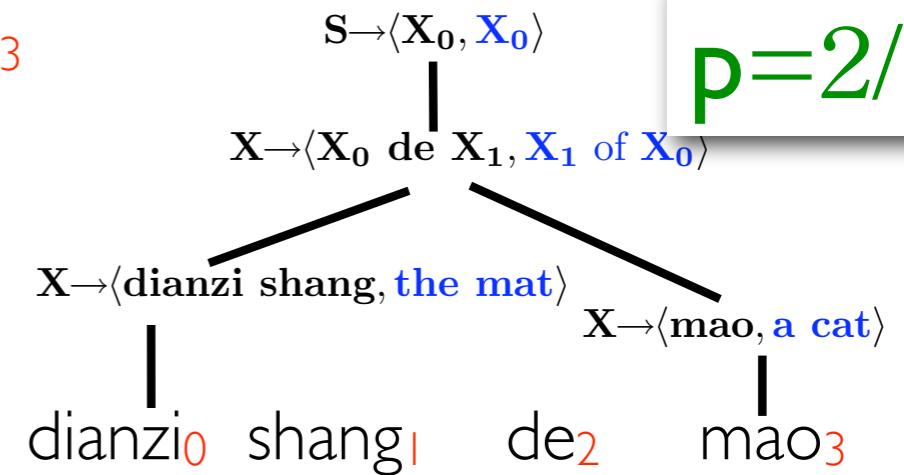
the mat 's a cat

$p=2/8$



the mat a cat

$p=2/8$



a cat of the mat

# Estimating $q^*$ from a hypergraph: brute force

a cat on the mat

$p=1/8$

the mat 's a cat

$p=3/8$

the mat a cat

$p=2/8$

a cat of the mat

$p=2/8$

# Estimating $q^*$ from a hypergraph: brute force

a cat on the mat

$p=1/8$

the mat 's a cat

$p=3/8$

the mat a cat

$p=2/8$

a cat of the mat

$p=2/8$

Bi-gram estimation:

- ▶ unpack the hypergraph

# Estimating $q^*$ from a hypergraph: brute force

a cat on the mat

$p=1/8$

the mat 's a cat

$p=3/8$

the mat a cat

$p=2/8$

a cat of the mat

$p=2/8$

Bi-gram estimation:

- ▶ unpack the hypergraph
- ▶ accumulate the **soft-count** of each bigram

# Estimating $q^*$ from a hypergraph: brute force

a cat on the mat

$p=1/8$

the mat 's a cat

$p=3/8$

the mat a cat

$p=2/8$

a cat of the mat

$p=2/8$

Bi-gram estimation:

- ▶ unpack the hypergraph
- ▶ accumulate the **soft-count** of each bigram
- ▶ normalize the counts

# Estimating $q^*$ from a hypergraph: brute force

a cat on the mat

$p=1/8$

the mat 's a cat

$p=3/8$

the mat a cat

$p=2/8$

a cat of the mat

$p=2/8$

Bi-gram estimation:

- ▶ unpack the hypergraph
- ▶ accumulate the **soft-count** of each bigram
- ▶ normalize the counts

$\Pr(\text{on} \mid \text{cat}) = 1/8$

$\Pr(\text{</s>} \mid \text{cat}) = 5/8$

$\Pr(\text{of} \mid \text{cat}) = 2/8$

# Estimating $q^*$ from a hypergraph: brute force

a cat on the mat

$p=1/8$

the mat 's a cat

$p=3/8$

the mat a cat

$p=2/8$

a cat of the mat

$p=2/8$

## Bi-gram estimation:

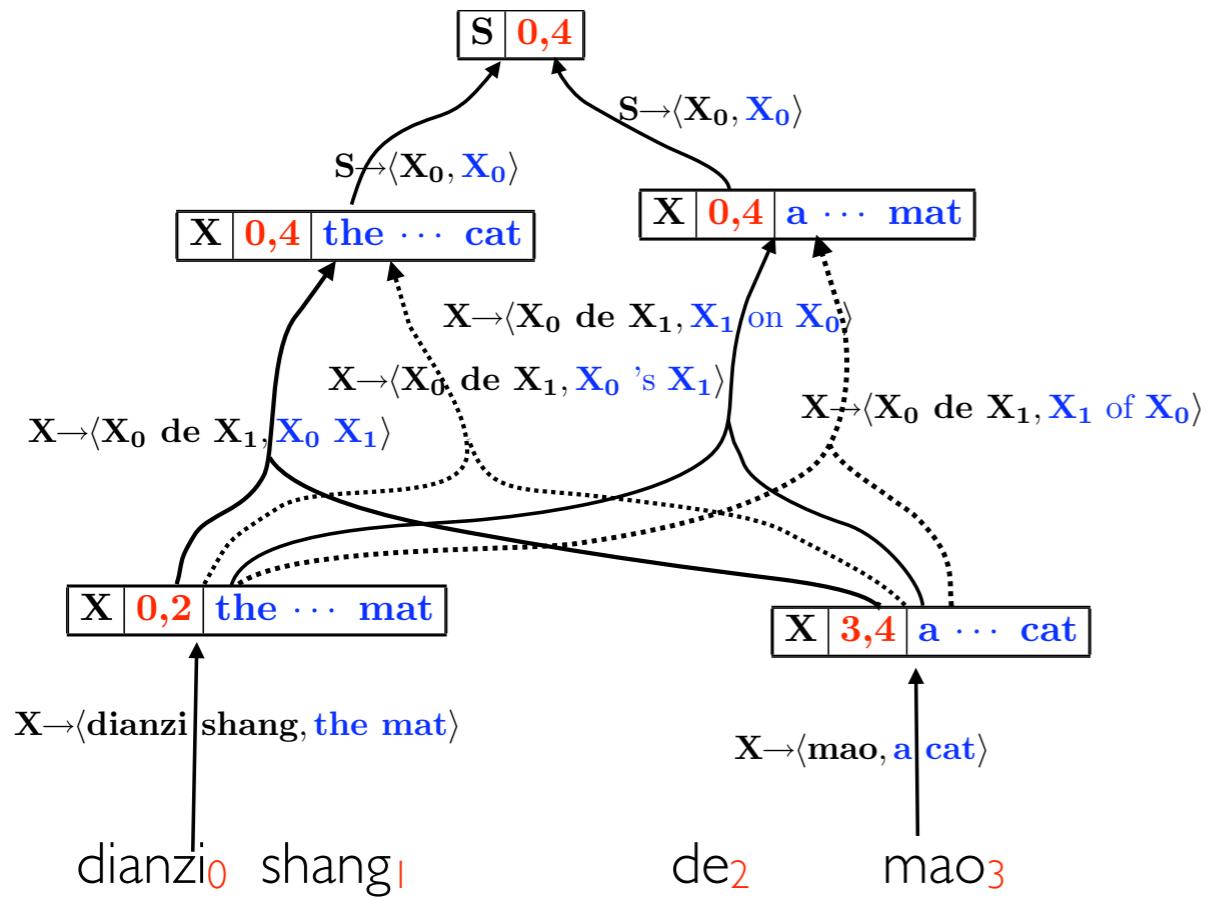
- ▶ unpack the hypergraph
- ▶ accumulate the **soft-count** of each bigram
- ▶ normalize the counts

$$\Pr(\text{on} \mid \text{cat}) = 1/8$$

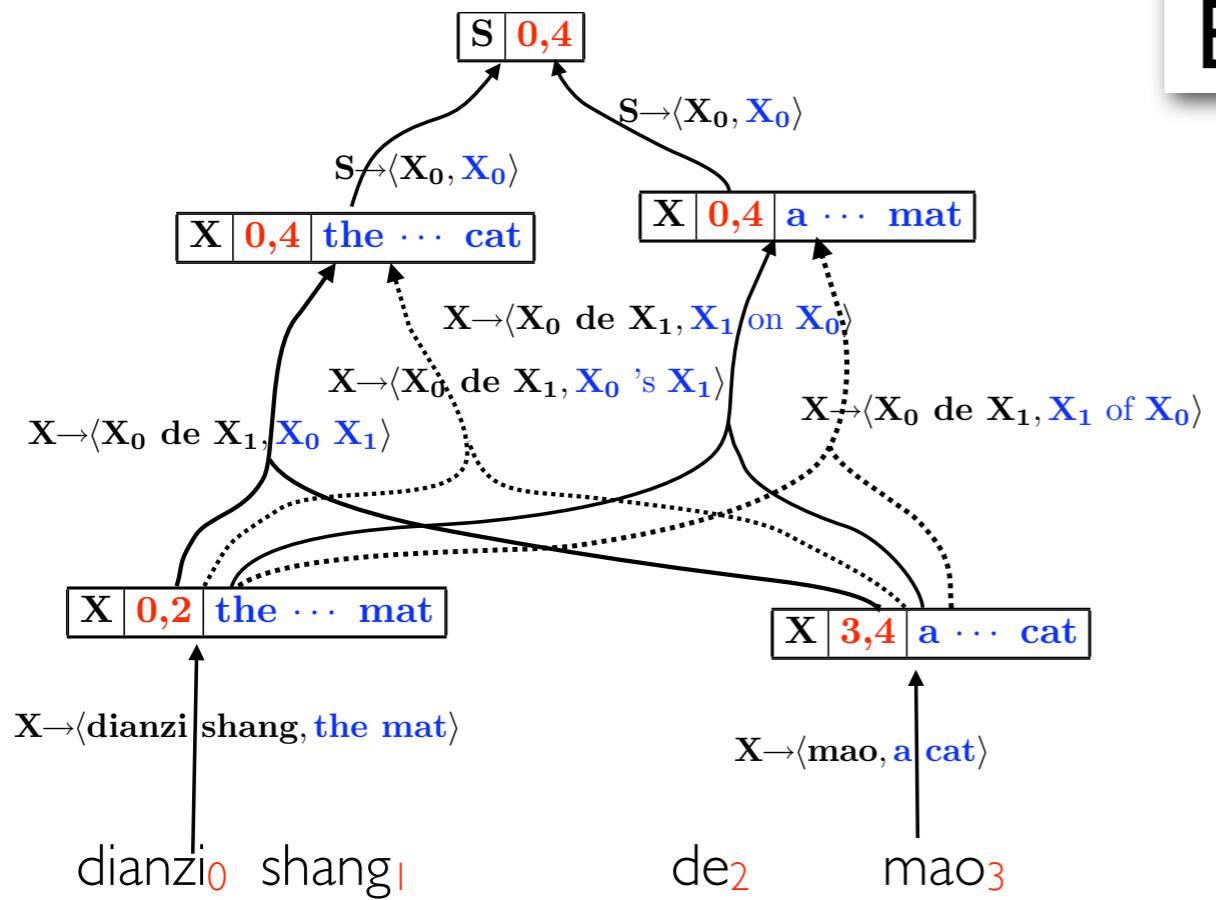
$$\Pr(\text{</s>} \mid \text{cat}) = 5/8$$

$$\Pr(\text{of} \mid \text{cat}) = 2/8$$

# Estimating $q^*$ from a hypergraph: dynamic programming

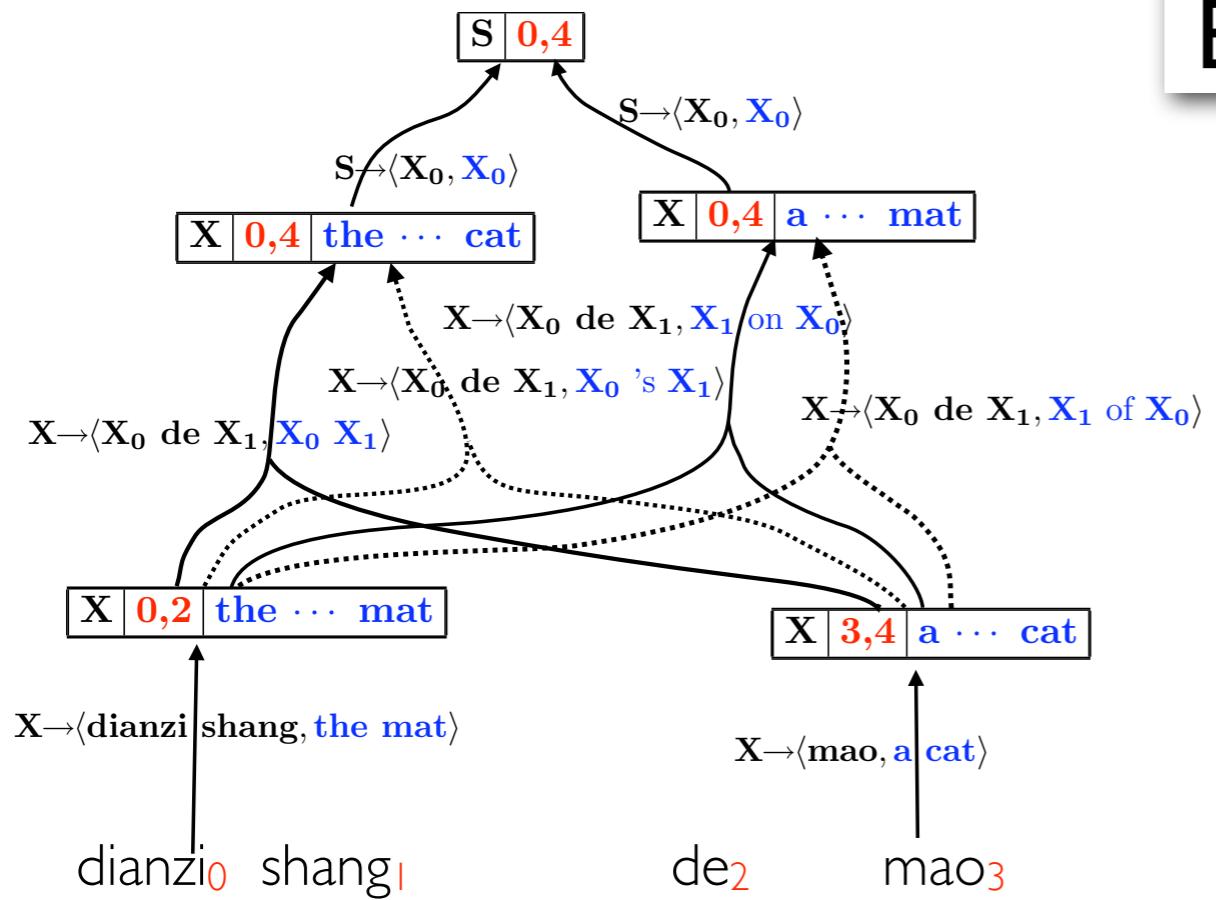


# Estimating $q^*$ from a hypergraph: dynamic programming



Bi-gram estimation:

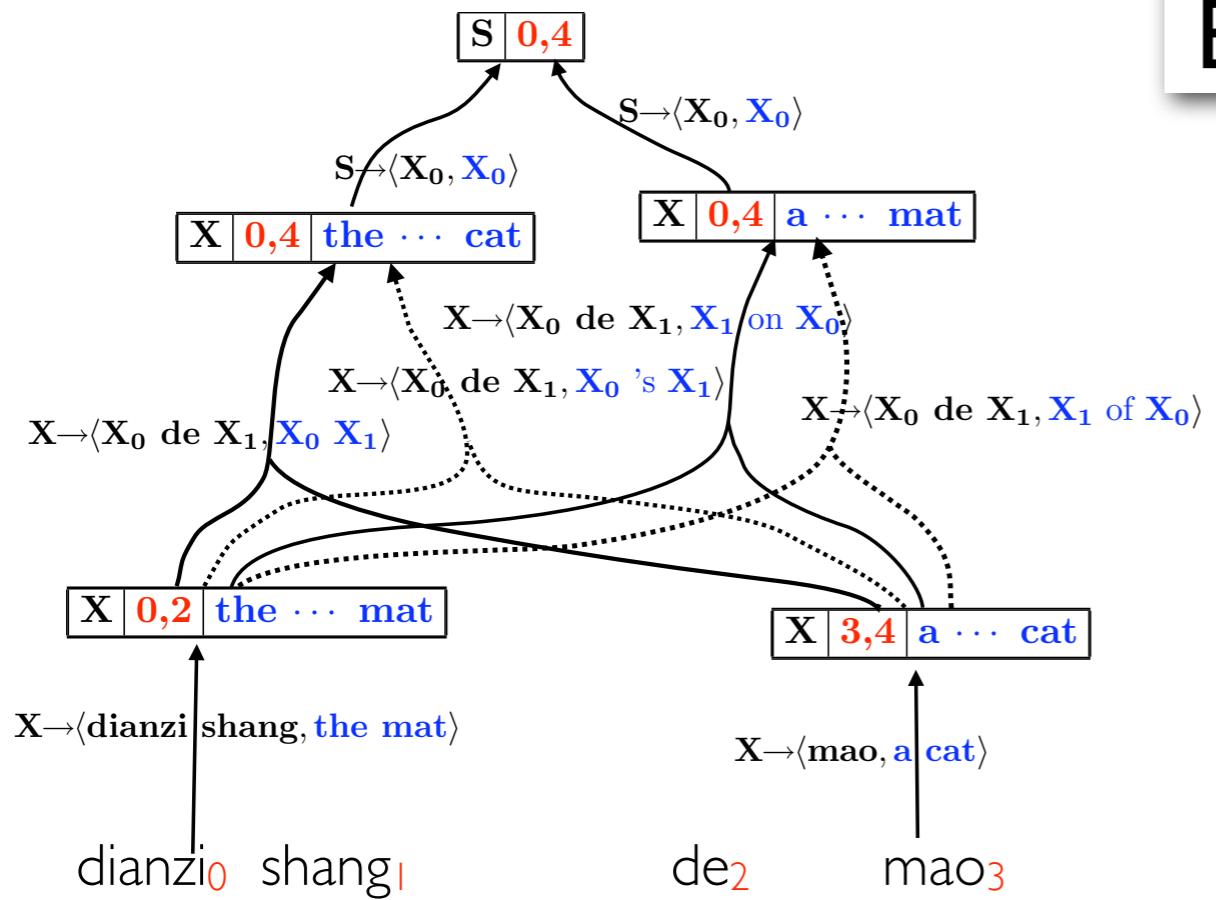
# Estimating $q^*$ from a hypergraph: dynamic programming



Bi-gram estimation:

- ▶ run inside-outside on the hypergraph

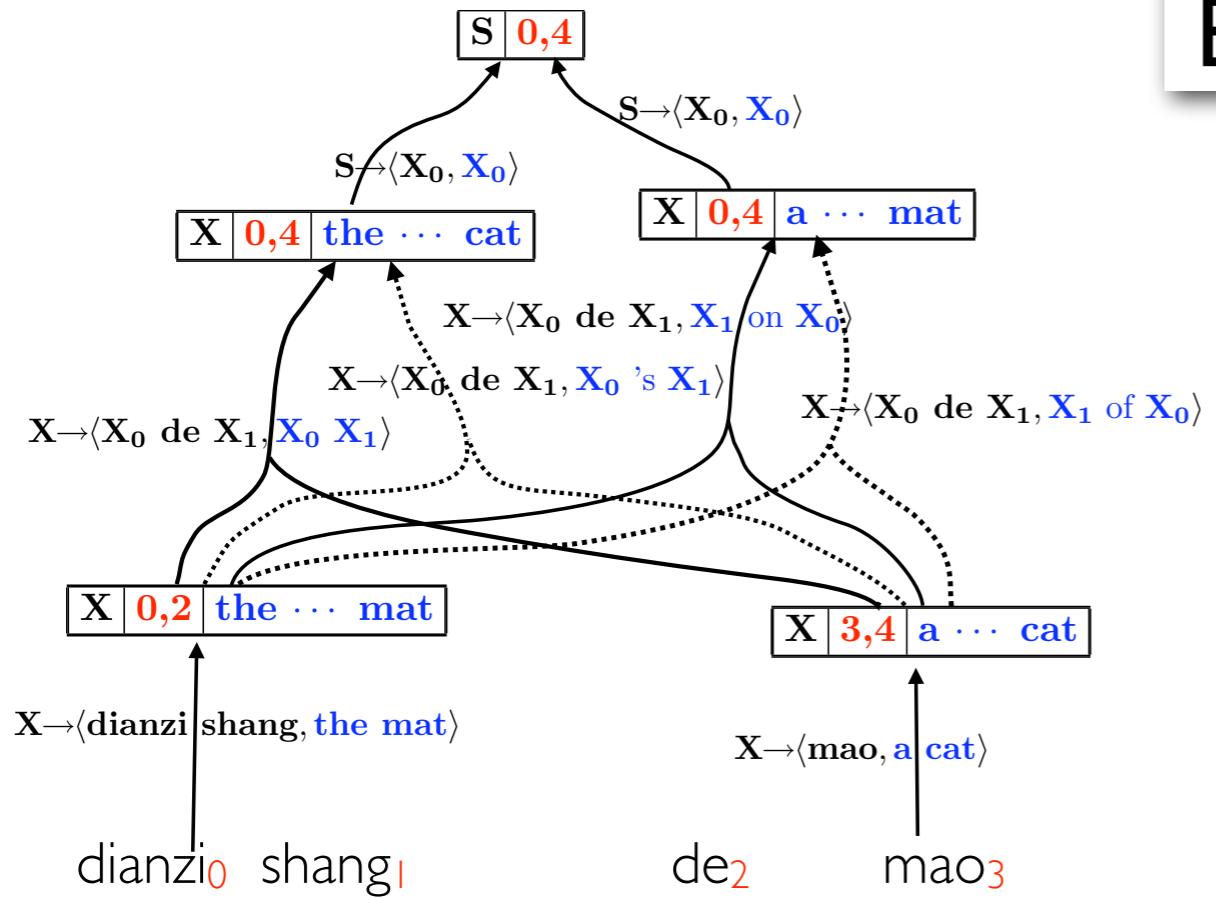
# Estimating $q^*$ from a hypergraph: dynamic programming



## Bi-gram estimation:

- ▶ run inside-outside on the hypergraph
- ▶ accumulate the **soft-count** of each bigram at each hyperedge

# Estimating $q^*$ from a hypergraph: dynamic programming



## Bi-gram estimation:

- ▶ run inside-outside on the hypergraph
- ▶ accumulate the **soft-count** of each bigram at each hyperedge
- ▶ normalize the counts

# Decoding using $q^* \in Q$

# Decoding using $q^* \in Q$

- Rescore the hypergraph  $HG(x)$

# Decoding using $q^* \in Q$

- Rescore the hypergraph  $HG(x)$

$$y^* = \arg \max_{y \in HG(x)} q^*(y|x)$$

# Decoding using $q^* \in Q$

- Rescore the hypergraph  $HG(x)$

$$y^* = \arg \max_{y \in HG(x)} q^*(y|x)$$

# Decoding using $q^* \in Q$

- Rescore the hypergraph  $HG(x)$

$$y^* = \arg \max_{y \in HG(x)} q^*(y|x)$$

$q^*$  is an n-gram model.

# Decoding using $q^* \in Q$

- Rescore the hypergraph  $HG(x)$

$$y^* = \arg \max_{y \in HG(x)} q^*(y|x)$$

$q^*$  is an n-gram model.

- have efficient dynamic programming algorithms
  - score the hypergraph using an n-gram model

# Decoding using $q^* \in Q$

- Rescore the hypergraph  $HG(x)$

$$y^* = \arg \max_{y \in HG(x)} q^*(y|x)$$

$q^*$  is an n-gram model.

- have efficient dynamic programming algorithms
  - score the hypergraph using an n-gram model

John already told you how to do this 😊

# KL divergences under different variational models

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q) = H(p, q) - H(p)$$

# KL divergences under different variational models

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q) = H(p, q) - H(p)$$

Measure bits/word	$\overline{H}(p)$	$\overline{\text{KL}}(p  \cdot)$			
		$q_1^*$	$q_2^*$	$q_3^*$	$q_4^*$
MT'04	1.36	0.97	0.32	0.21	0.17
MT'05	1.37	0.94	0.32	0.21	0.17

# KL divergences under different variational models

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q) = H(p, q) - H(p)$$

Measure bits/word	$\bar{H}(p)$	$\overline{\text{KL}}(p  \cdot)$			
		$q_1^*$	$q_2^*$	$q_3^*$	$q_4^*$
MT'04	1.36	0.97	0.32	0.21	0.17
MT'05	1.37	0.94	0.32	0.21	0.17

- The larger the order **n** is, the smaller the KL divergence is!
- The reduction of KL divergence happens mostly when switching from unigram to bigram

# KL divergences under different variational models

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q) = H(p, q) - H(p)$$

Measure bits/word	$\bar{H}(p)$	$\overline{\text{KL}}(p  \cdot)$			
		$q_1^*$	$q_2^*$	$q_3^*$	$q_4^*$
MT'04	1.36	0.97	0.32	0.21	0.17
MT'05	1.37	0.94	0.32	0.21	0.17

# KL divergences under different variational models

$$q^* = \arg \min_{q \in Q} \text{KL}(p||q) = H(p, q) - H(p)$$

Measure bits/word	$\bar{H}(p)$	$\overline{\text{KL}}(p  \cdot)$			
		$q_1^*$	$q_2^*$	$q_3^*$	$q_4^*$
MT'04	1.36	0.97	0.32	0.21	0.17
MT'05	1.37	0.94	0.32	0.21	0.17

How to compute them on a **hypergraph**?

see (Li and Eisner, EMNLP'09)

# BLEU scores when using a single variational n-gram model

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
1gram	25.9	24.5
2gram	<b>36.1</b>	<b>33.4</b>
3gram	36.0	33.1
4gram	35.8	32.9

# BLEU scores when using a single variational n-gram model

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
1gram	25.9	24.5
2gram	<b>36.1</b>	<b>33.4</b>
3gram	36.0	33.1
4gram	35.8	32.9

- unigram performs very badly

# BLEU scores when using a single variational n-gram model

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
1gram	25.9	24.5
2gram	<b>36.1</b>	<b>33.4</b>
3gram	36.0	33.1
4gram	35.8	32.9

- unigram performs very badly
- bigram achieves best BLEU scores

# BLEU scores when using a single variational n-gram model

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
1gram	25.9	24.5
2gram	<b>36.1</b>	<b>33.4</b>
3gram	36.0	33.1
4gram	35.8	32.9

- unigram performs very badly
- bigram achieves best BLEU scores ???

# BLEU scores when using a single variational n-gram model

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
1gram	25.9	24.5
2gram	<b>36.1</b>	<b>33.4</b>
3gram	36.0	33.1
4gram	35.8	32.9

- unigram performs very badly
- bigram achieves best BLEU scores

???

modeling error in  $P$



BLEU cares about both low- and high-order  
n-gram matches

BLEU cares about both low- and high-order n-gram matches

- Interpolating variational n-gram model for different n

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

BLEU cares about both low- and high-order n-gram matches

- Interpolating variational n-gram model for different n

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

Viterbi and variational are different ways in approximating P

BLEU cares about both low- and high-order n-gram matches

- Interpolating variational n-gram model for different n

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

Viterbi and variational are different ways in approximating P

$$y^* = \arg \max_{y \in \text{HG}(x)} \left( \sum_n \theta_n \cdot \log q_n^*(y \mid x) + \theta_v \cdot \log p_{\text{Viterbi}}(y \mid x) \right)$$

BLEU cares about both low- and high-order n-gram matches

- Interpolating variational n-gram model for different n

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

Viterbi and variational are different ways in approximating P

$$y^* = \arg \max_{y \in \text{HG}(x)} \left( \sum_n \theta_n \cdot \log q_n^*(y \mid x) + \theta_v \cdot \boxed{\log p_{\text{Viterbi}}(y \mid x)} \right)$$

# Minimum Bayes Risk (MBR) decoding?

(Tromble et al. 2008)

(Denero et al. 2009)

# Minimum Risk Decoding

- Maximum A Posterior (MAP) decoding
  - find the most **probable** translation string

$$y^* = \arg \max_{y \in \text{HG}(x)} p(y|x)$$

- Minimum risk decoding
  - find the **consensus** translation string

$$y^* = \arg \min_{y \in \text{HG}(x)} \text{Risk}(y)$$

$$\text{Risk}(y) = \sum_{y'} L(y, y') p(y' | x)$$

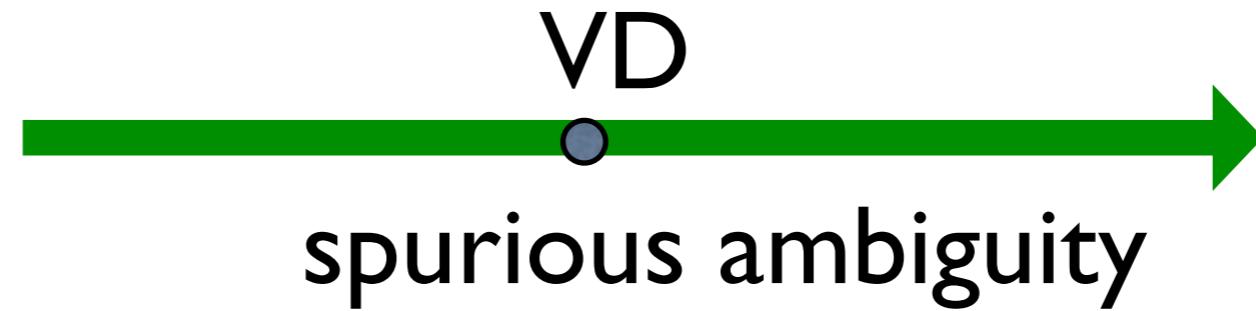
# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)

# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)

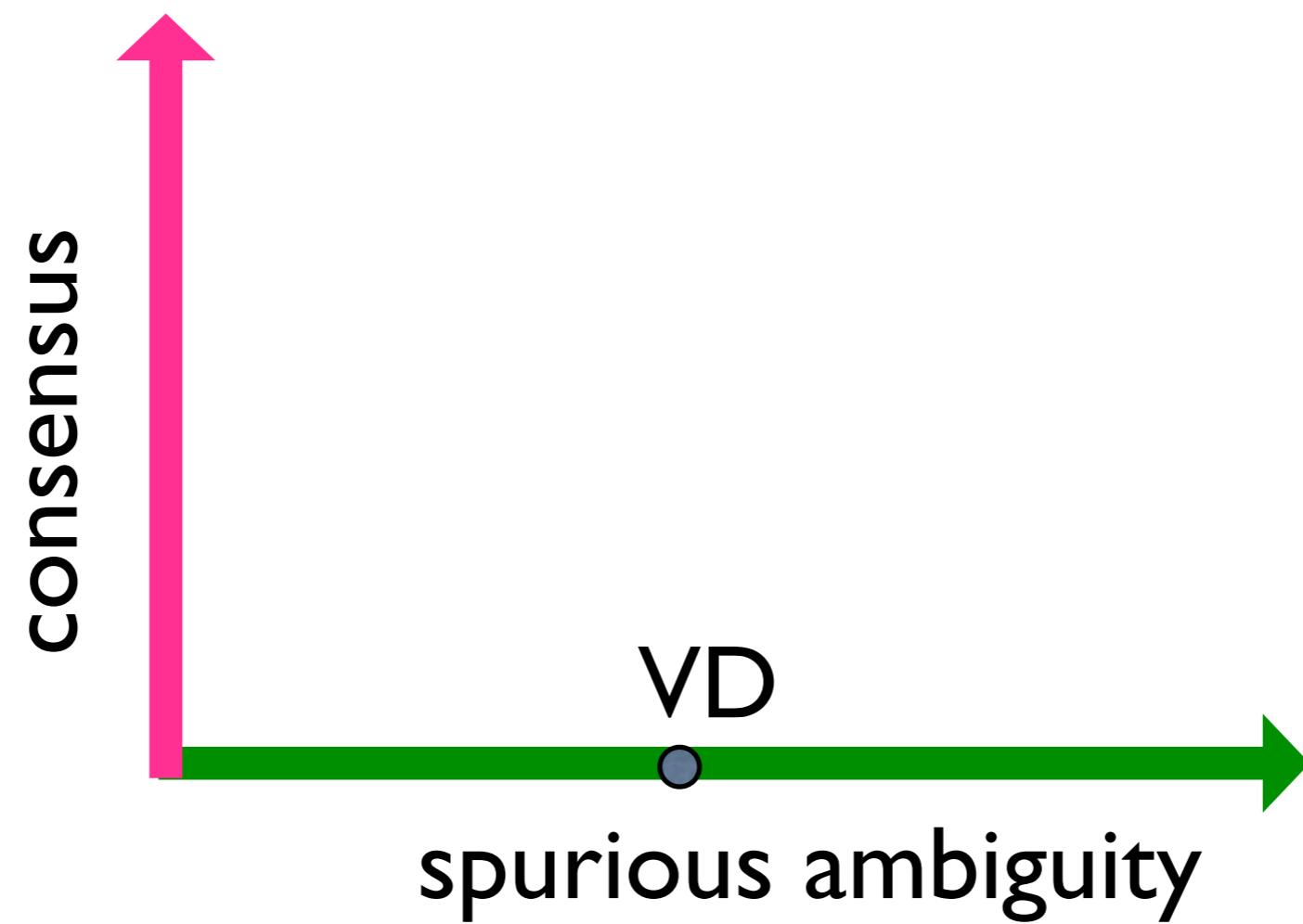


spurious ambiguity

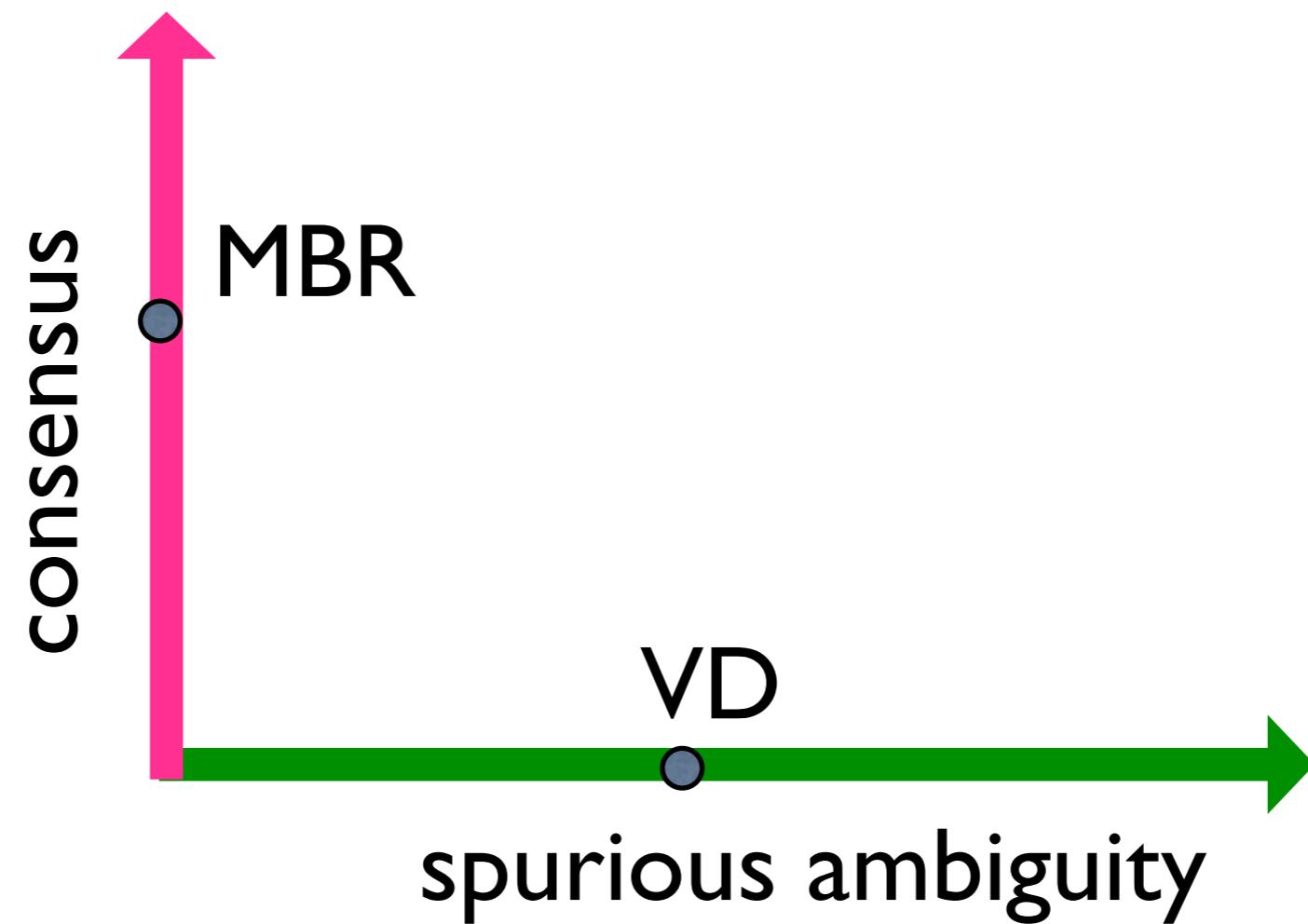
# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)



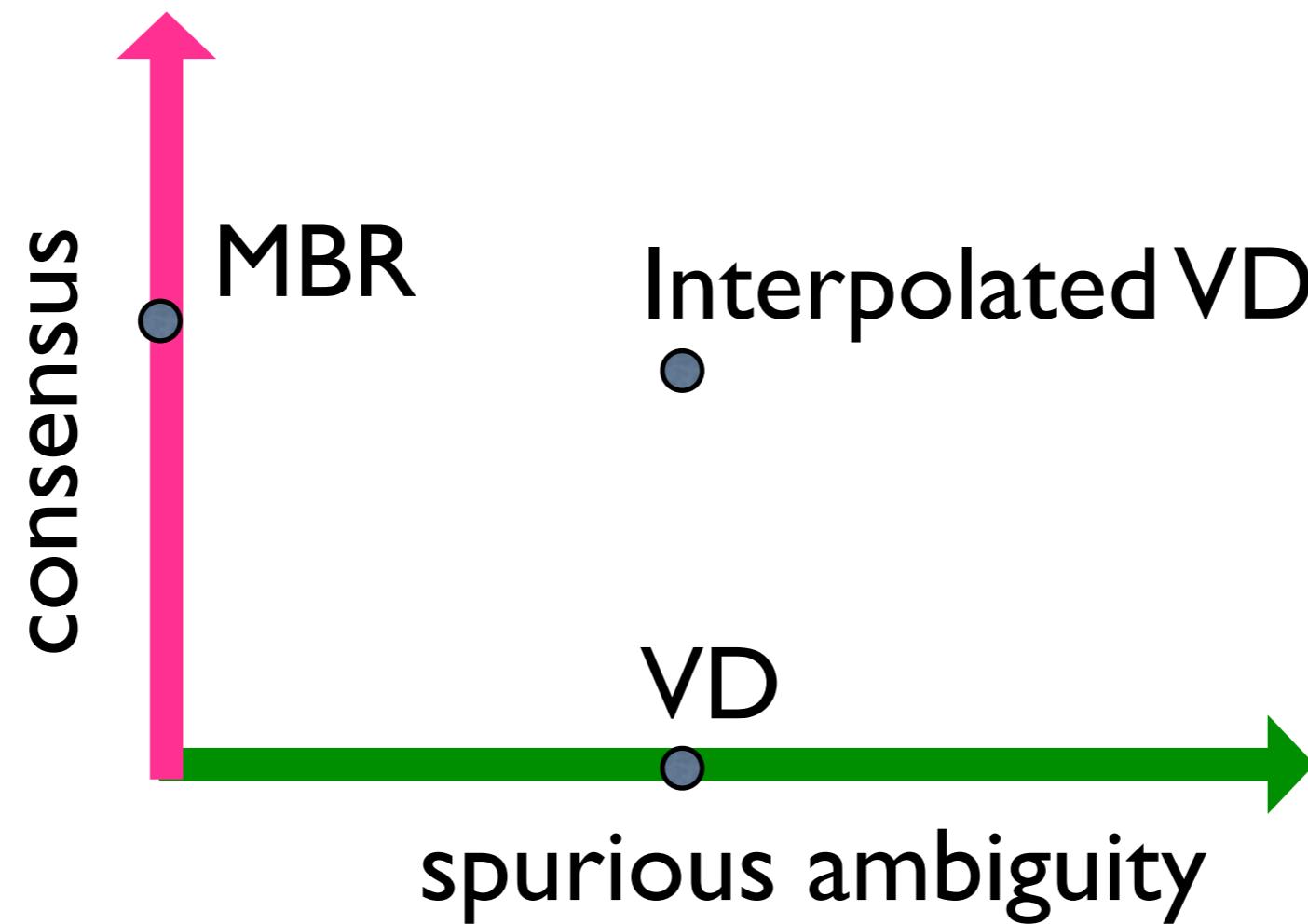
# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)



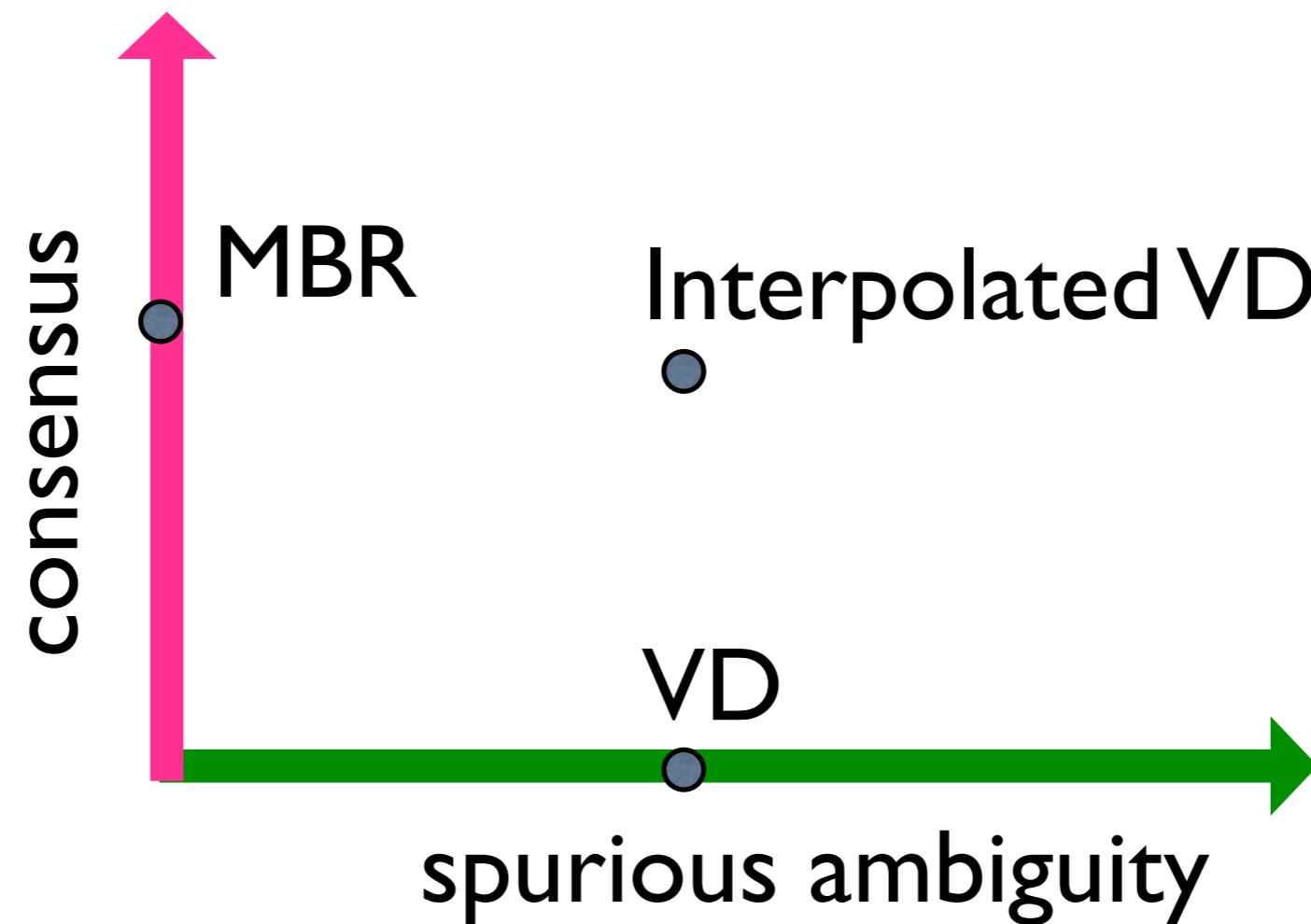
# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)



# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)



# Variational Decoding(VD) vs. MBR (Tromble et al. 2008)



Both BLEU metric and our variational distributions happen to use n-gram dependencies.

- Variational decoding with interpolation

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

$$q_n(y \mid x) = \prod_{w \in W_n} q(r(w) \mid h(w), x)^{c_w(y)}$$

$$q(r(w) \mid h(w), x) = \frac{\sum_{y'} c_w(y') p(y' \mid x)}{\sum_{y'} c_{h(w)}(y') p(y' \mid x)}$$

- Minimum risk decoding (Tromble et al. 2008)

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot g_n(y \mid x)$$

$$g_n(y \mid x) = \sum_{w \in W_n} g(w \mid x) c_w(y)$$

$$g(w \mid x) = \sum_{y'} \delta_w(y') p(y' \mid x)$$

- Variational decoding with interpolation

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

decision rule

$$q_n(y \mid x) = \prod_{w \in W_n} q(r(w) \mid h(w), x)^{c_w(y)}$$

$$q(r(w) \mid h(w), x) = \frac{\sum_{y'} c_w(y') p(y' \mid x)}{\sum_{y'} c_{h(w)}(y') p(y' \mid x)}$$

- Minimum risk decoding (Tromble et al. 2008)

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot g_n(y \mid x)$$

decision rule

$$g_n(y \mid x) = \sum_{w \in W_n} g(w \mid x) c_w(y)$$

$$g(w \mid x) = \sum_{y'} \delta_w(y') p(y' \mid x)$$

- Variational decoding with interpolation

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

decision rule

$$q_n(y \mid x) = \prod_{w \in W_n} q(r(w) \mid h(w), x)^{c_w(y)}$$

n-gram model

$$q(r(w) \mid h(w), x) = \frac{\sum_{y'} c_w(y') p(y' \mid x)}{\sum_{y'} c_{h(w)}(y') p(y' \mid x)}$$

- Minimum risk decoding (Tromble et al. 2008)

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot g_n(y \mid x)$$

decision rule

$$g_n(y \mid x) = \sum_{w \in W_n} g(w \mid x) c_w(y)$$

n-gram model

$$g(w \mid x) = \sum_{y'} \delta_w(y') p(y' \mid x)$$

- Variational decoding with interpolation

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

decision rule

$$q_n(y \mid x) = \prod_{w \in W_n} q(r(w) \mid h(w), x)^{c_w(y)}$$

n-gram model

$$q(r(w) \mid h(w), x) = \frac{\sum_{y'} c_w(y') p(y' \mid x)}{\sum_{y'} c_{h(w)}(y') p(y' \mid x)}$$

n-gram probability

- Minimum risk decoding (Tromble et al. 2008)

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot g_n(y \mid x)$$

decision rule

$$g_n(y \mid x) = \sum_{w \in W_n} g(w \mid x) c_w(y)$$

n-gram model

$$g(w \mid x) = \sum_{y'} \delta_w(y') p(y' \mid x)$$

n-gram probability

- Variational decoding with interpolation

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y \mid x)$$

$$q_n(y \mid x) = \prod_{w \in W_n} q(r(w) \mid h(w), x)^{c_w(y)}$$

$$q(r(w) \mid h(w), x) = \frac{\sum_{y'} c_w(y') p(y' \mid x)}{\sum_{y'} c_{h(w)}(y') p(y' \mid x)}$$

- Minimum risk decoding (Tromble et al. 2008)

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot g_n(y \mid x)$$

$$g_n(y \mid x) = \sum_{w \in W_n} g(w \mid x) c_w(y)$$

$$g(w \mid x) = \sum_{y'} \delta_w(y') p(y' \mid x)$$

- Variational decoding with interpolation

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot \log q_n^*(y | x)$$

$$q_n(y | x) = \prod_{w \in W_n} q(r(w) | h(w), x)^{c_w(y)}$$

$$q(r(w) | h(w), x) = \frac{\sum_{y'} c_w(y') p(y' | x)}{\sum_{y'} c_{h(w)}(y') p(y' | x)}$$

- Minimum risk decoding (Tromble et al. 2008)

$$y^* = \arg \max_{y \in \text{HG}(x)} \sum_n \theta_n \cdot g_n(y | x)$$

$$g_n(y | x) = \sum_{w \in W_n} g(w | x) c_w(y)$$

non-probabilistic

$$g(w | x) = \sum_{y'} \delta_w(y') p(y' | x)$$

very expensive to compute

# BLEU Results on Chinese-English NIST MT Tasks

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
MBR ( $K=1000$ )	35.8	32.7
Crunching ( $N=10000$ )	35.7	32.8
Crunching+MBR ( $N=10000$ )	35.8	32.7
Variational (1to4gram+wp+vt)	<b>36.6</b>	<b>33.5</b>

# BLEU Results on Chinese-English NIST MT Tasks

Decoding scheme	MT'04	MT'05
Viterbi	35.4	32.6
MBR ( $K=1000$ )	35.8	32.7
Crunching ( $N=10000$ )	35.7	32.8
Crunching+MBR ( $N=10000$ )	35.8	32.7
Variational (1to4gram+wp+vt)	<b>36.6</b>	<b>33.5</b>

- variational decoding improves over Viterbi, MBR, and crunching

# Conclusions

- Exact MAP decoding with spurious ambiguity is intractable
- Viterbi or N-best approximations are efficient, but ignore most derivations
- We developed a variational approximation, which considers all derivations but still allows tractable decoding
- Our variational decoding improves a state of the art baseline

# Future directions

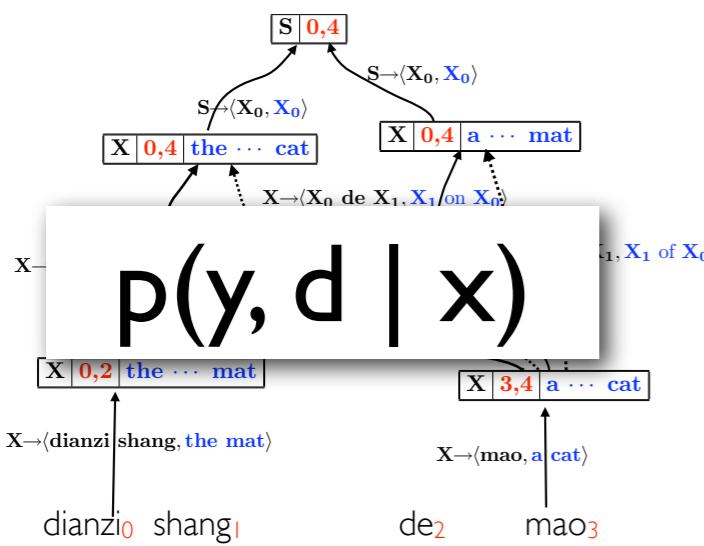
- The MT pipeline is full of intractable problems
  - variational approximation is a principled way to tackle these problems
- Decoding with spurious ambiguity is a common problem in many other NLP applications
  - Models with latent variables
  - Data oriented parsing (DOP)
  - Hidden Markov Models (HMM)
  - .....

Thank you!  
谢谢！

A circular photograph of a Joshua tree in a desert environment. The tree has a thick, textured trunk and several branches topped with clusters of sharp, light-green spines. The background shows a vast, flat desert plain under a clear blue sky with a few wispy clouds.

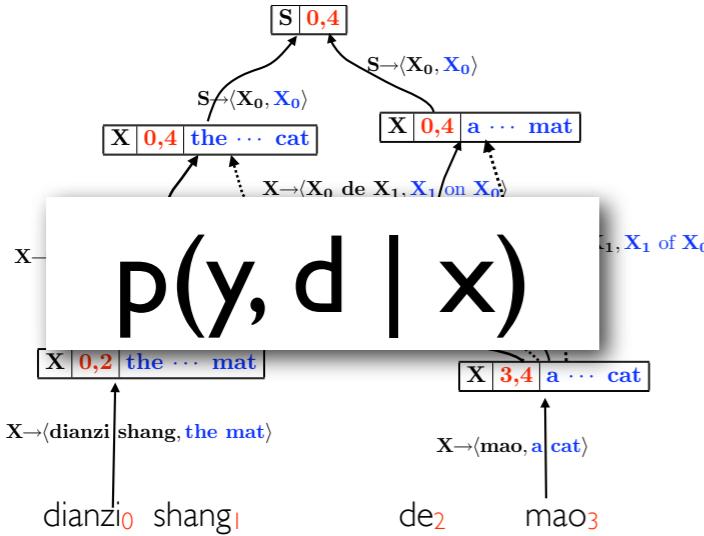
joshua

1



## Generate a hypergraph

2



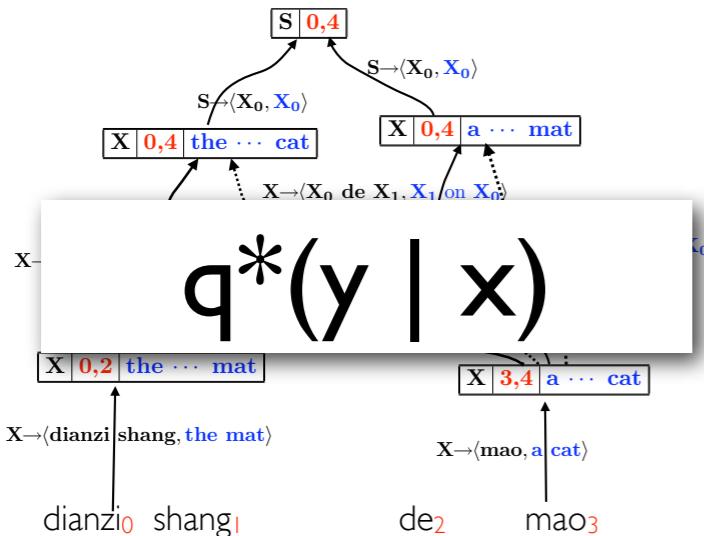
## Estimate a model from the hypergraph

$q^*$  is an n-gram model  
over output strings.

$$q^*(y | x)$$

$$\approx \sum_{d \in D(x,y)} P(y,d|x)$$

3



## Decode using $q^*$ on the hypergraph