# Novel Inference, Training and Decoding Methods over Translation Forests

**Zhifei Li**

**Center for Language and Speech Processing**

**Computer Science Department**

**Johns Hopkins University**

**Advisor: Sanjeev Khudanpur**

**Co-advisor: Jason Eisner**

# Statistical Machine Translation Pipeline

# Statistical Machine Translation Pipeline

Bilingual
  Data

# Statistical Machine Translation Pipeline

Bilingual Data → Generative Training → Translation Models

# Statistical Machine Translation Pipeline

Bilingual Data $\rightarrow$ Generative Training $\rightarrow$ Translation Models

Monolingual English

# Statistical Machine Translation Pipeline

Bilingual Data → Generative Training → Translation Models
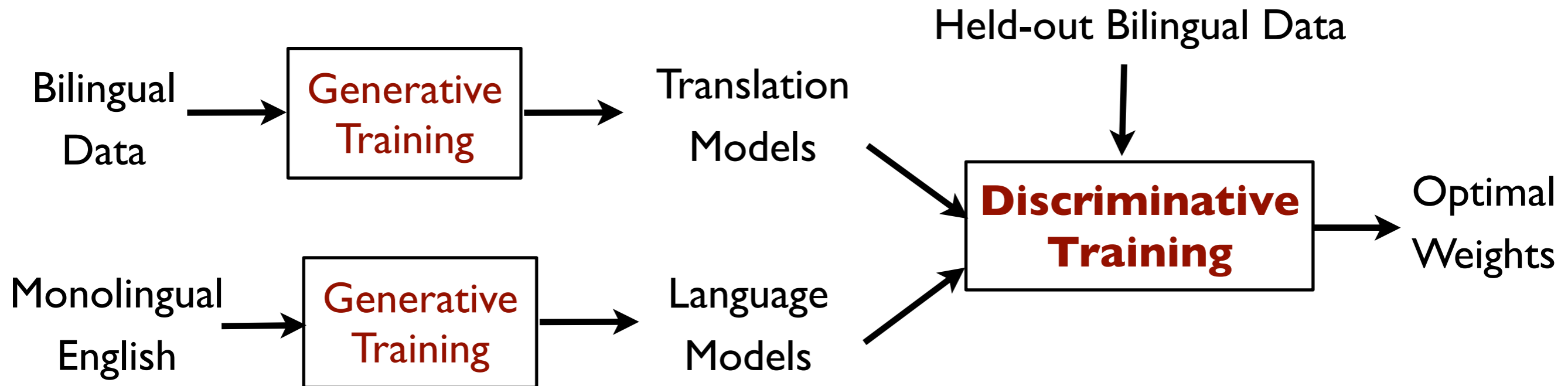
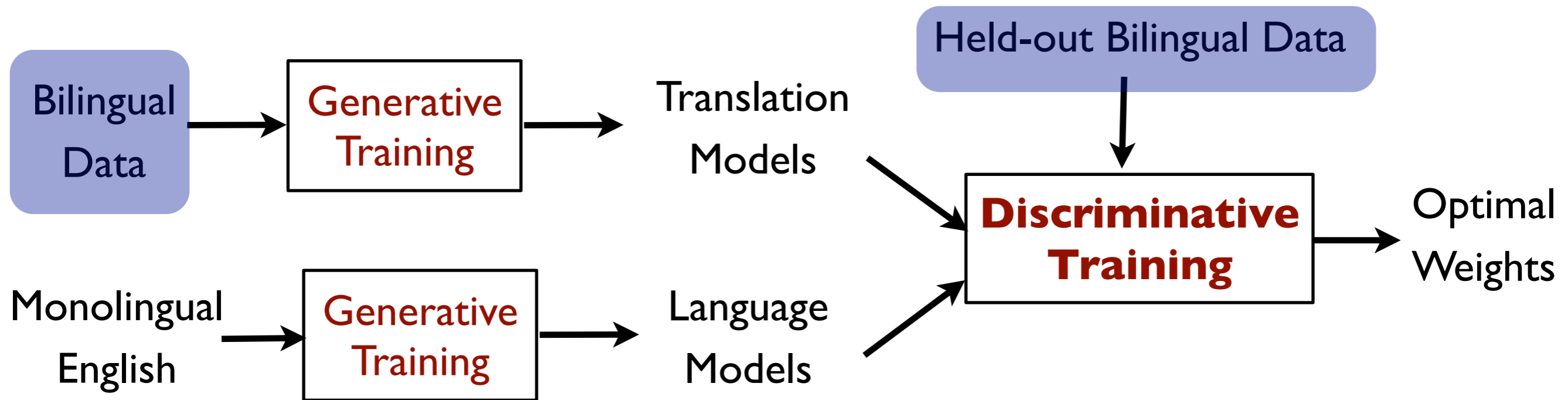Monolingual English → Generative Training → Language Models

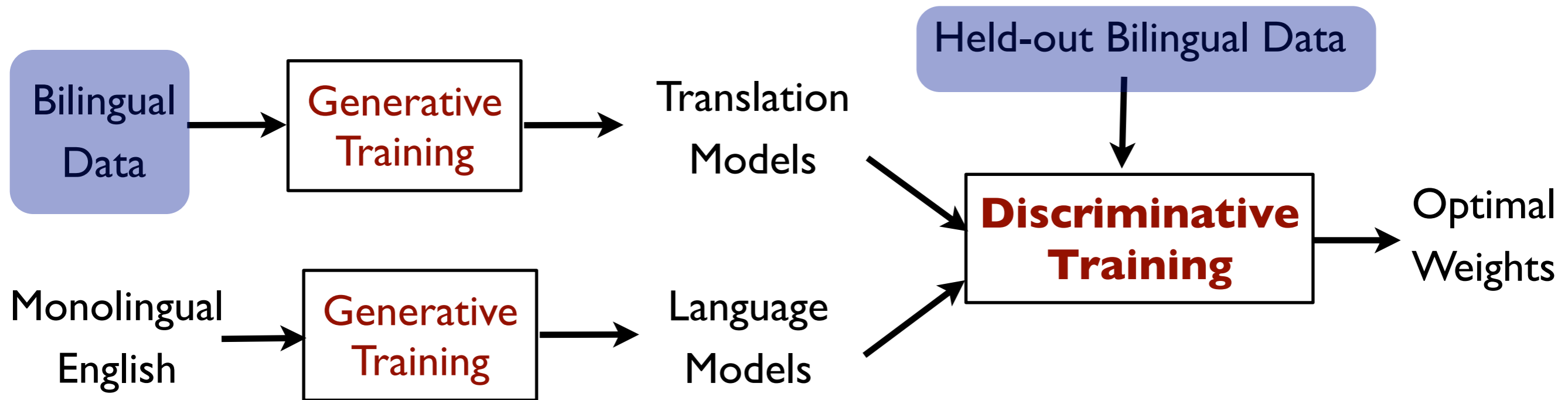# Statistical Machine Translation Pipeline
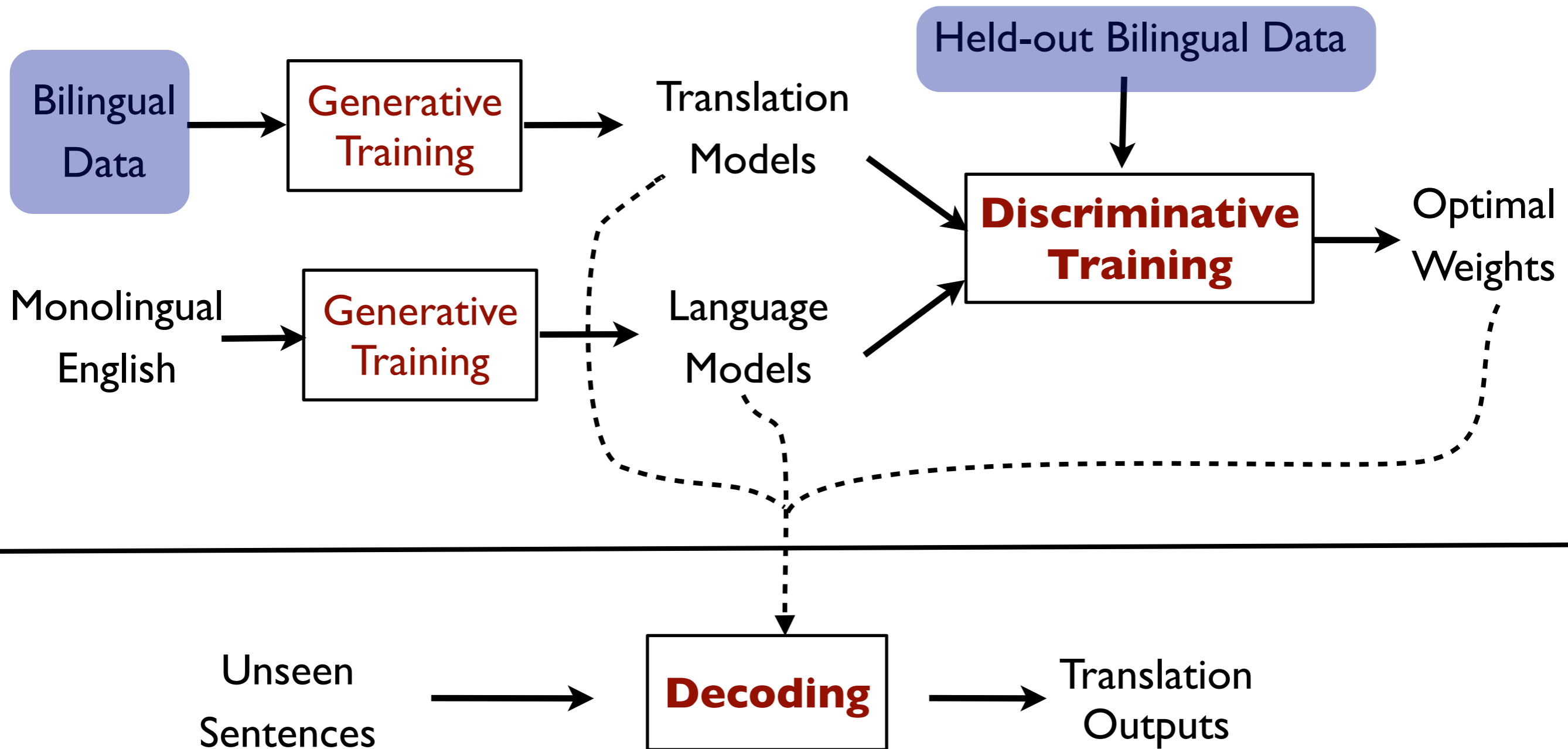
# Statistical Machine Translation Pipeline

# Statistical Machine Translation Pipeline

# Statistical Machine Translation Pipeline

# Training a Translation Model

# Training a Translation Model



垫子　上　的　猫
dianzi  shang  de  mao

# Training a Translation Model



墊子　上　的　猫
dianzi  shang  de mao

a　cat  on　the　mat

# Training a Translation Model



垫子　　上　的　猫

dianzi　shang　de　mao

a　cat　on　the　mat

# Training a Translation Model



大子　上　的　猫

dianzi shang de mao

a cat on the mat

$$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$$

# Training a Translation Model



堆子　　上　的　猫

dianzi  shang  de  mao

a  cat  on  the  mat

$X \rightarrow \langle \text{ dianzi  shang } , \text{ the  mat } \rangle$
$X \rightarrow \langle \text{ mao } , \text{ a cat } \rangle$

# Training a Translation Model



$$X \longrightarrow \langle \text{ dianzi shang , the mat } \rangle$$
$$X \longrightarrow \langle \text{ mao , a cat } \rangle$$

# Training a Translation Model



墊子　　上　的　猫

dianzi　shang de $X_0$

$X_0$　on　the　mat

$X \rightarrow \langle$ dianzi　shang , the　mat $\rangle$
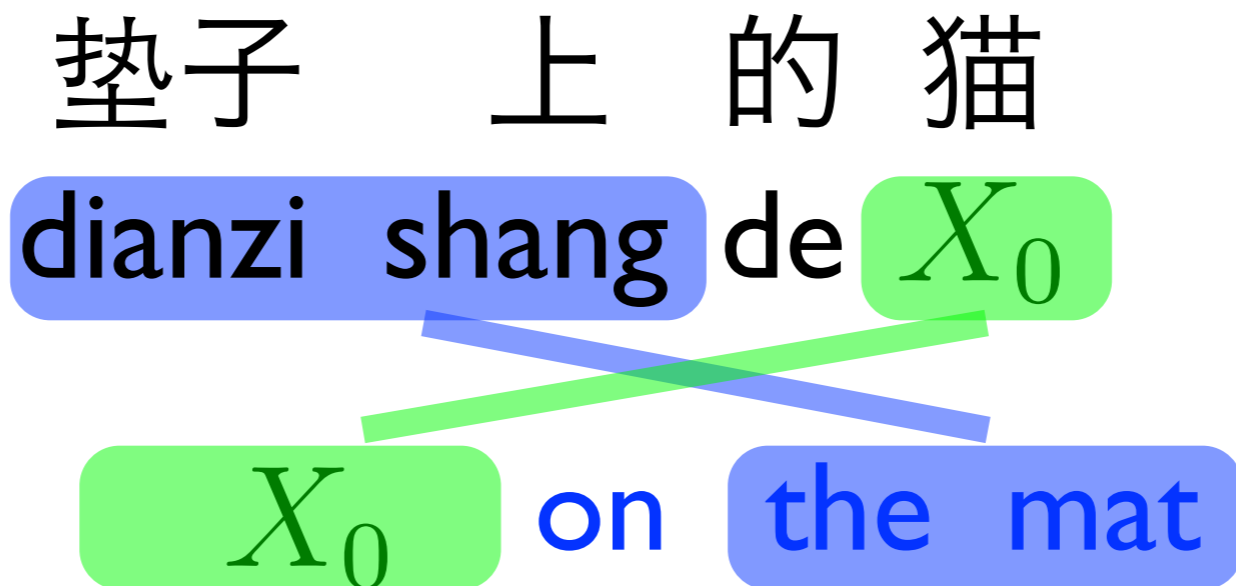$X \rightarrow \langle$ mao , 　a cat $\rangle$

# Training a Translation Model


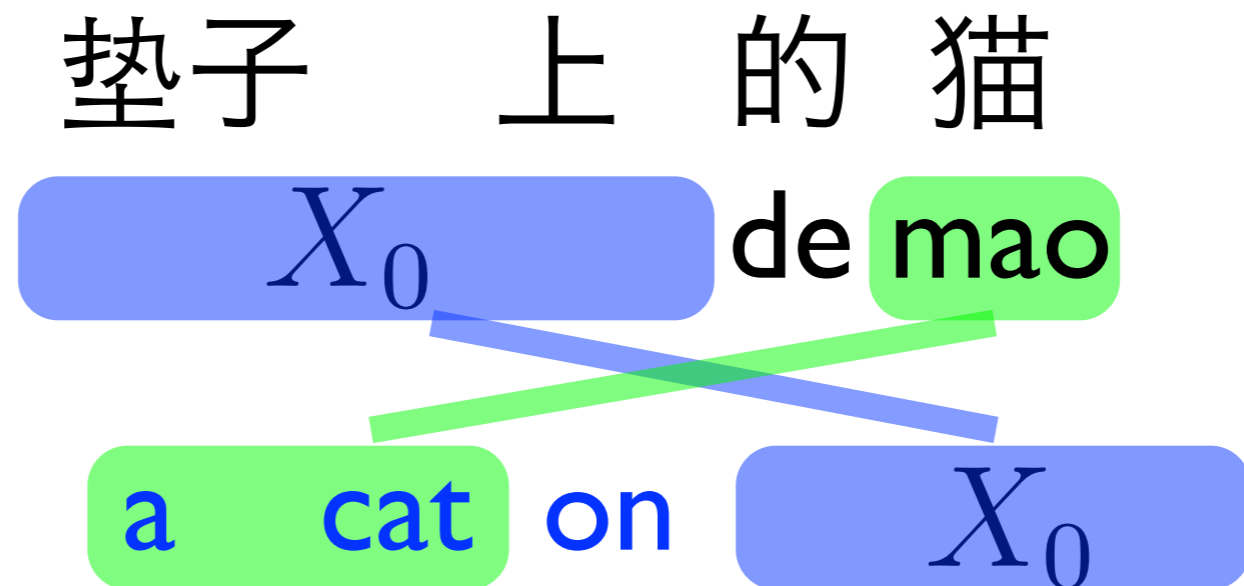
墊子　　上　的　猫

dianzi shang de $X_0$

$X_0$ on the mat

$X \longrightarrow \langle$ dianzi shang , the mat $\rangle$

$X \longrightarrow \langle$ mao , a cat $\rangle$

$X \longrightarrow \langle$ dianzi shang de $X_0$ , $X_0$ on the mat $\rangle$

# Training a Translation Model



垫子　上　的　猫

$X_0$　de mao

a　cat　on　$X_0$

$X \rightarrow \langle$ dianzi shang , the mat $\rangle$

$X \rightarrow \langle$ mao , a cat $\rangle$

$X \rightarrow \langle$ dianzi shang de $X_0$ , $X_0$ on the mat $\rangle$

# Training a Translation Model



$$X \rightarrow \langle \text{ dianzi shang }, \text{the mat} \rangle$$
$$X \rightarrow \langle \text{ mao }, \text{ a cat} \rangle$$
$$X \rightarrow \langle \text{ dianzi shang de } X_0 , X_0 \text{ on the mat} \rangle$$
$$X \rightarrow \langle X_0 \text{ de mao }, \text{a cat on } X_0 \rangle$$

# Training a Translation Model



垫子　上　的　猫

$X_0$　de $X_1$

$X_1$ on $X_0$
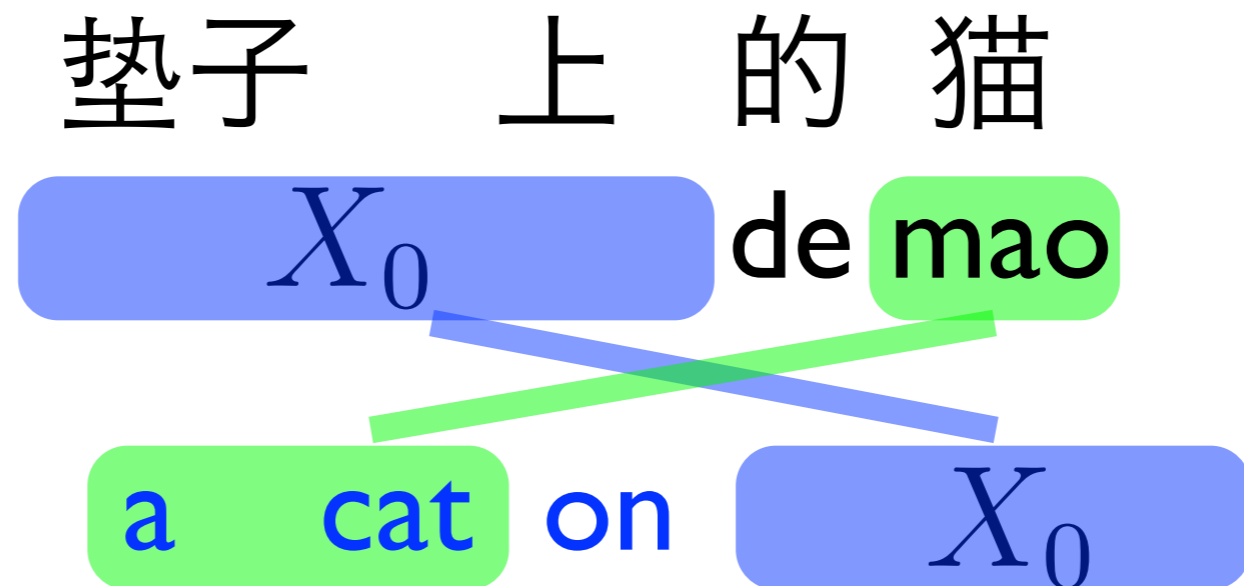
$X \rightarrow \langle$ dianzi shang , the mat $\rangle$

$X \rightarrow \langle$ mao , a cat $\rangle$

$X \rightarrow \langle$ dianzi shang de $X_0$ , $X_0$ on the mat $\rangle$

$X \rightarrow \langle X_0$ de mao , a cat on $X_0 \rangle$

# Training a Translation Model
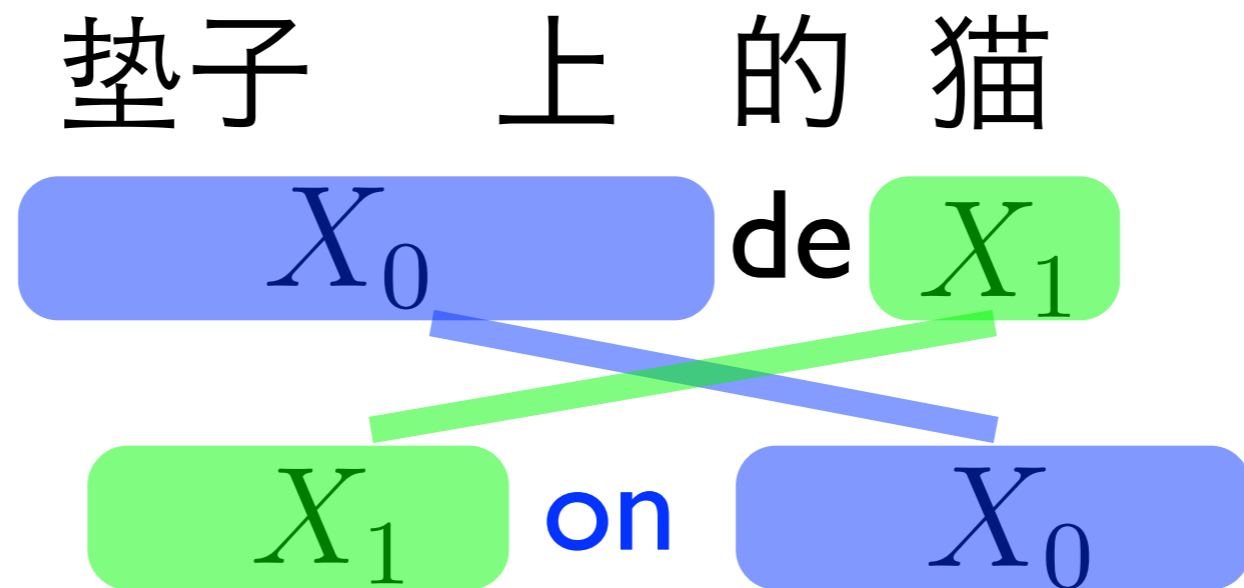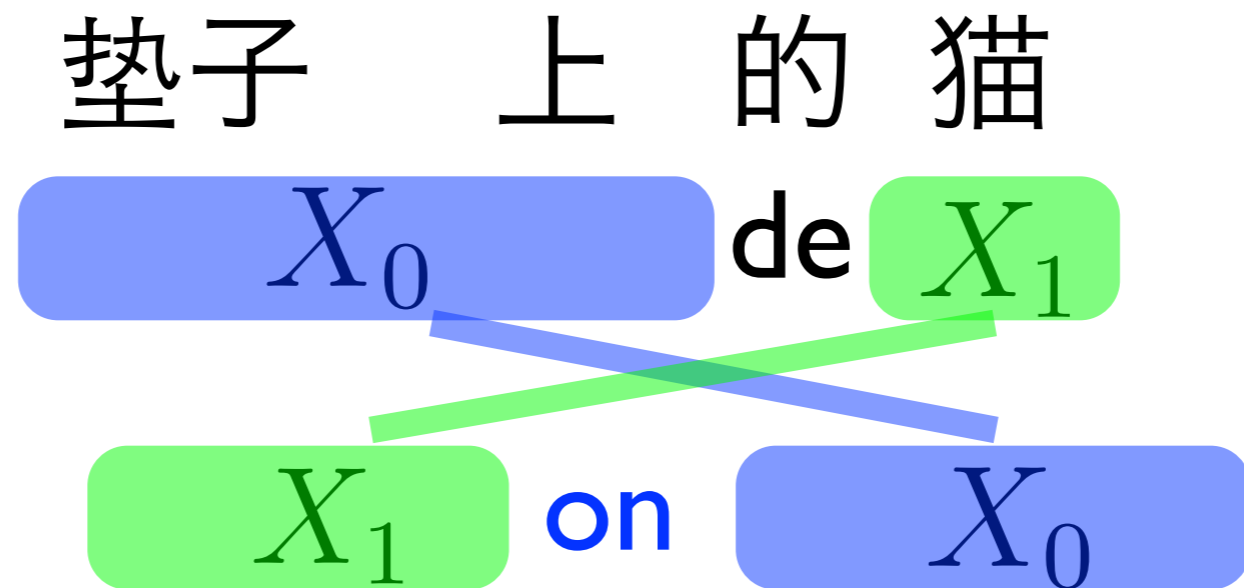
墊子 上 的 猫

$X_0$ **de** $X_1$

$X_1$ **on** $X_0$

$X \rightarrow \langle$ **dianzi shang** , the mat $\rangle$

$X \rightarrow \langle$ **mao** , a cat $\rangle$

$X \rightarrow \langle$ **dianzi shang de** $X_0$ , $X_0$ on the mat $\rangle$

$X \rightarrow \langle X_0$ **de mao** , a cat on $X_0 \rangle$

$X \rightarrow \langle X_0$ **de** $X_1$ , $X_1$ on $X_0 \rangle$

# Decoding a Test Sentence

# Decoding a Test Sentence

# Decoding a Test Sentence

墊子　上　的　狗

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi　shang de gou**

# Decoding a Test Sentence



垫子　上　的　狗

dianzi　shang de gou

the　dog　on　the　mat

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi　shang de gou**

**the　dog on the mat**

$$X \rightarrow \langle \text{ dianzi shang },\ \text{the mat } \rangle$$
$$X \rightarrow \langle \text{ gou },\ \text{the dog } \rangle$$
$$X \rightarrow \langle X_0 \text{ de } X_1 ,\ X_1 \text{ on } X_0 \rangle$$
$$S \rightarrow \langle X_0 ,\ X_0 \rangle$$

# Decoding a Test Sentence



垫子　上　的　狗

**dianzi  shang de gou**

**the    dog  on  the  mat**

| | | |
|---|---|---|
| $X$ | $\rightarrow$ | $\langle$ dianzi shang , the mat $\rangle$ |
| $X$ | $\rightarrow$ | $\langle$ gou , the dog $\rangle$ |
| $X$ | $\rightarrow$ | $\langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$ |
| $S$ | $\rightarrow$ | $\langle X_0 , X_0 \rangle$ |

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi  shang de gou**

the   dog  on  the  mat

| | | |
|---|---|---|
| $X$ | $\rightarrow$ | $\langle$ dianzi shang , the mat $\rangle$ |
| $X$ | $\rightarrow$ | $\langle$ gou , the dog $\rangle$ |
| $X$ | $\rightarrow$ | $\langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$ |
| $S$ | $\rightarrow$ | $\langle X_0 , X_0 \rangle$ |

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi  shang de gou**

<span style="color:blue">the    dog  on  the  mat</span>

| | | |
|---|---|---|
| $X$ | $\rightarrow$ | $\langle$ dianzi shang , the mat $\rangle$ |
| $X$ | $\rightarrow$ | $\langle$ gou , the dog $\rangle$ |
| $X$ | $\rightarrow$ | $\langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$ |
| $S$ | $\rightarrow$ | $\langle X_0 , X_0 \rangle$ |

dianzi  shang    de    gou

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi  shang de gou**

the　dog  on  the  mat

| | | |
|---|---|---|
| $X$ | $\rightarrow$ | $\langle$ dianzi shang , the mat $\rangle$ |
| $X$ | $\rightarrow$ | $\langle$ gou , the dog $\rangle$ |
| $X$ | $\rightarrow$ | $\langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$ |
| $S$ | $\rightarrow$ | $\langle X_0$ , $X_0 \rangle$ |

$X \rightarrow \langle$**dianzi shang,** the mat$\rangle$

|

dianzi  shang　de　gou

8

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi　shang de gou**

<span style="color:blue">**the　　dog　on　the　mat**</span>

$X \rightarrow \langle$ dianzi shang ， the mat $\rangle$

$X \rightarrow \langle$ gou ， the dog $\rangle$

$X \rightarrow \langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$

$S \rightarrow \langle X_0$ , $X_0 \rangle$

$X \rightarrow \langle$**dianzi shang**, <span style="color:blue">**the mat**</span>$\rangle$

$X \rightarrow \langle$**gou**, <span style="color:blue">**the dog**</span>$\rangle$

dianzi　shang　　de　　gou

# Decoding a Test Sentence

垫子　上　的　狗
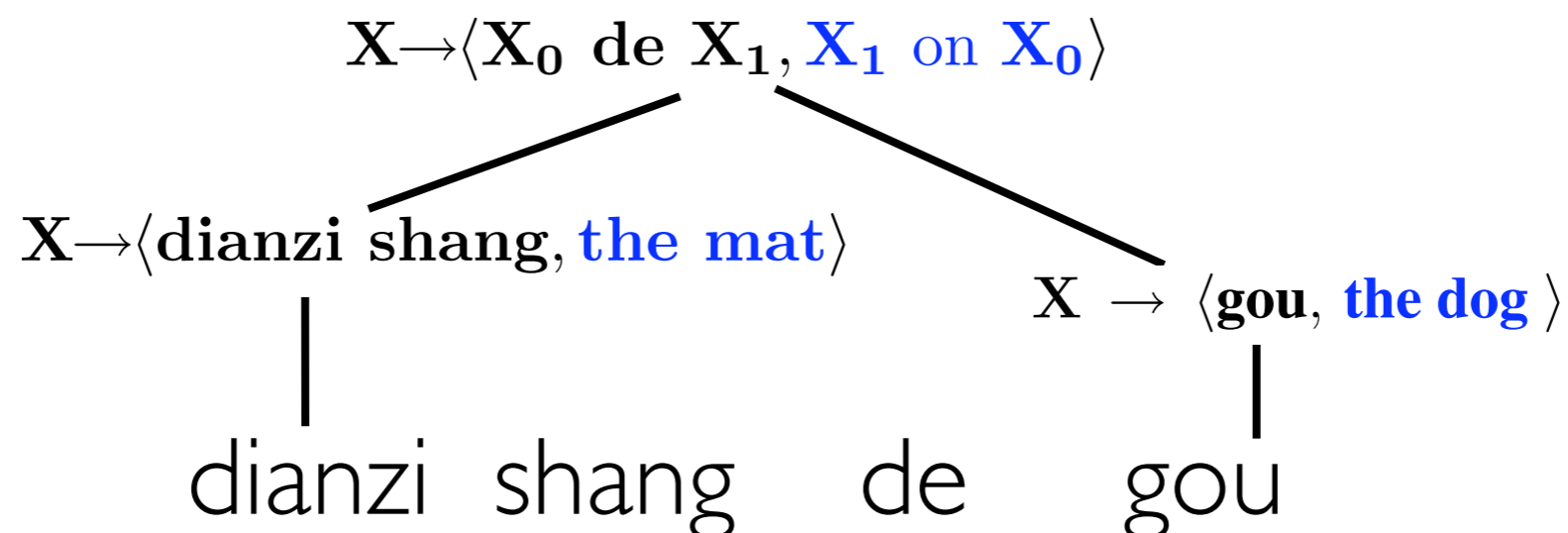
**dianzi　shang　de　gou**

<span style="color:blue">the　dog　on　the　mat</span>

$$X \rightarrow \langle \text{ dianzi shang , the mat } \rangle$$
$$X \rightarrow \langle \text{ gou , the dog } \rangle$$
$$X \rightarrow \langle X_0 \text{ de } X_1 , X_1 \text{ on } X_0 \rangle$$
$$S \rightarrow \langle X_0 , X_0 \rangle$$

$$\mathbf{X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle}$$

$$\mathbf{X \rightarrow \langle dianzi\ shang, the\ mat \rangle}$$

$$\mathbf{X \rightarrow \langle gou, the\ dog \rangle}$$

dianzi　shang　de　gou
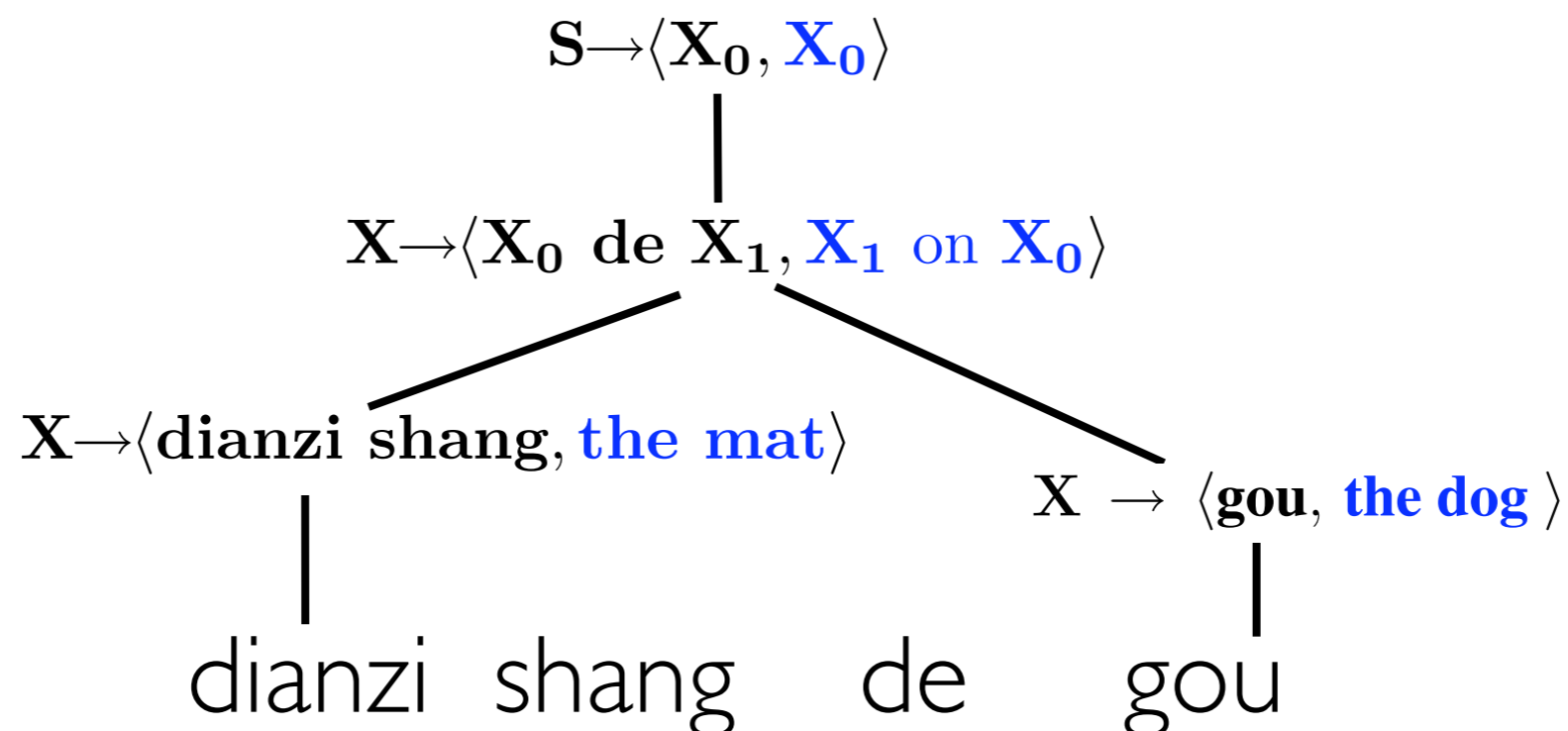
# Decoding a Test Sentence

垫子　上　的　狗

**dianzi　shang de gou**

**the　　dog　on　the　mat**

---

$X \rightarrow \langle$ dianzi shang , the mat $\rangle$

$X \rightarrow \langle$ gou , the dog $\rangle$

$X \rightarrow \langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$

$S \rightarrow \langle X_0 , X_0 \rangle$

---

$\mathbf{S} \rightarrow \langle \mathbf{X_0}, \mathbf{X_0} \rangle$

$\mathbf{X} \rightarrow \langle \mathbf{X_0}$ **de** $\mathbf{X_1}, \mathbf{X_1}$ on $\mathbf{X_0} \rangle$

$\mathbf{X} \rightarrow \langle$ **dianzi shang**, **the mat** $\rangle$

$\mathbf{X} \rightarrow \langle$ **gou**, **the dog** $\rangle$

dianzi　shang　　de　　gou

# Decoding a Test Sentence

墊子　上　的　狗

**dianzi shang de gou**

<span style="color:blue">the　dog　on　the　mat</span>
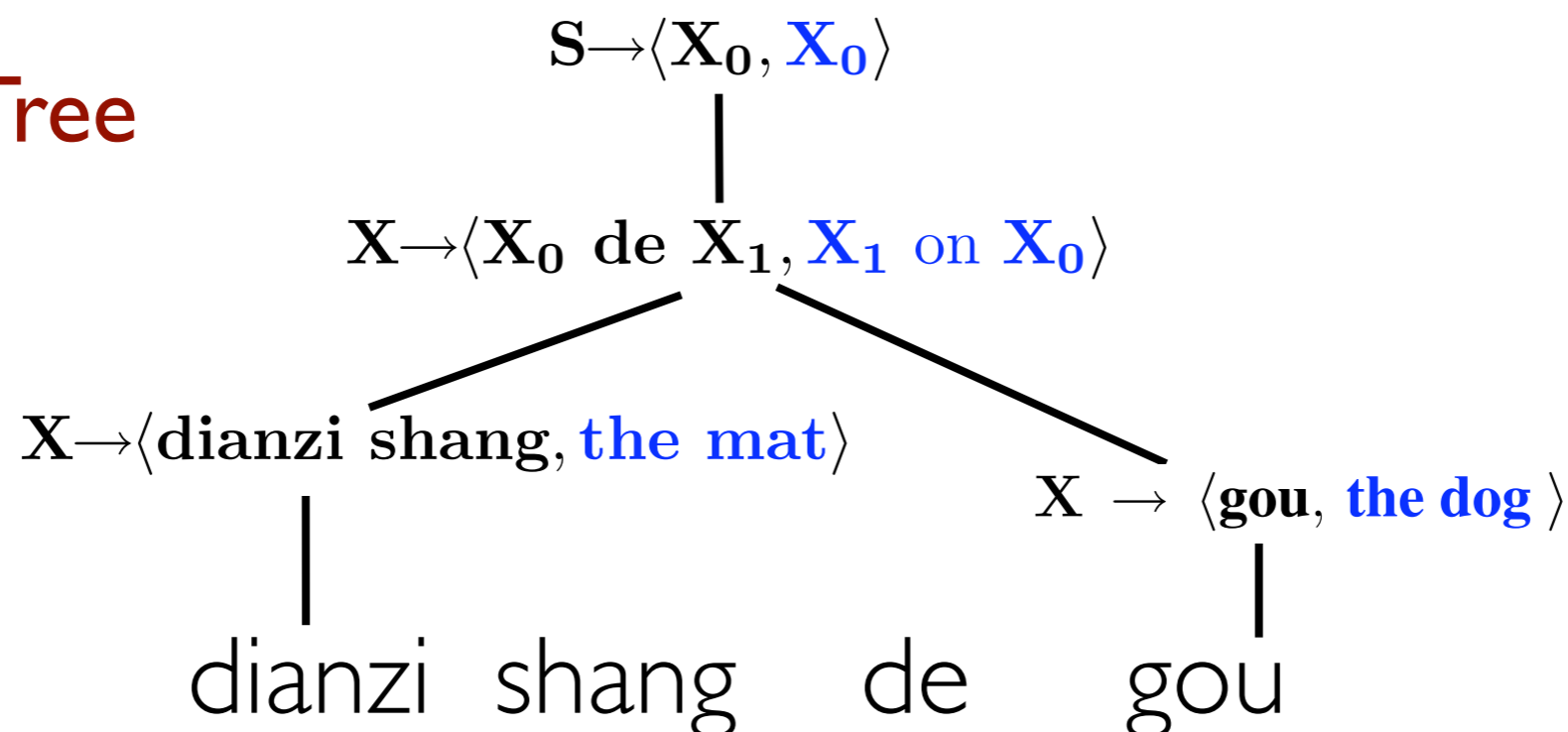
$X \rightarrow \langle$ dianzi shang , the mat $\rangle$

$X \rightarrow \langle$ gou , the dog $\rangle$

$X \rightarrow \langle X_0$ de $X_1$ , $X_1$ on $X_0 \rangle$

$S \rightarrow \langle X_0 , X_0 \rangle$

## Derivation Tree

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ gou, the dog $\rangle$

dianzi shang　de　gou

# Decoding a Test Sentence

垫子　上　的　狗

**dianzi  shang de gou**
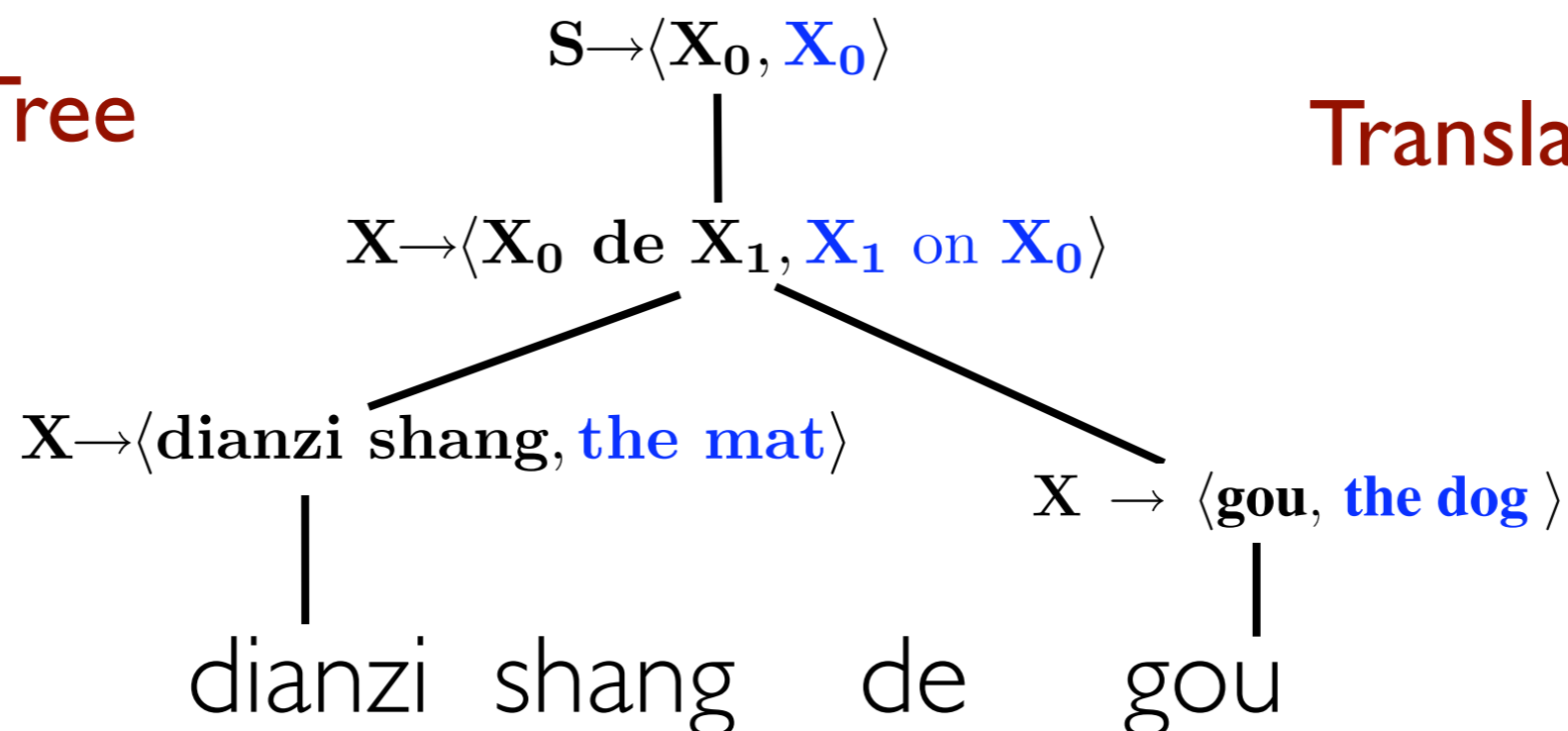
**the    dog  on  the  mat**

$$X \rightarrow \langle \text{dianzi shang} , \text{the mat} \rangle$$
$$X \rightarrow \langle \text{gou} , \text{the dog} \rangle$$
$$X \rightarrow \langle X_0 \text{ de } X_1 , X_1 \text{ on } X_0 \rangle$$
$$S \rightarrow \langle X_0 , X_0 \rangle$$

**Derivation Tree**

**Translation is easy?**

$$S \rightarrow \langle \mathbf{X_0}, \mathbf{X_0} \rangle$$

$$X \rightarrow \langle \mathbf{X_0} \text{ de } \mathbf{X_1}, \mathbf{X_1} \text{ on } \mathbf{X_0} \rangle$$

$$X \rightarrow \langle \mathbf{dianzi\ shang}, \mathbf{the\ mat} \rangle$$

$$X \rightarrow \langle \mathbf{gou}, \mathbf{the\ dog} \rangle$$

dianzi  shang   de    gou

# Translation Ambiguity

墊子　　上　的　猫
dianzi  shang de mao

a  cat  on  the  mat

# Translation Ambiguity



墊子　上　的　猫

dianzi　shang　de　mao

a　cat　on　the　mat

$X \rightarrow \langle X_0 \ de \ X_1, \ X_1 \ on \ X_0 \rangle$

9

# Translation Ambiguity

墊子　　上　的　猫

dianzi　shang　de　mao

a　cat　on　the　mat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

zhongguo　de　shoudu

capital　of　China

# Translation Ambiguity



墙子　　上　的　猫

dianzi　shang　de　mao

a　cat　on　the　mat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

zhongguo　de　shoudu

capital　of　China

# Translation Ambiguity

墊子　　上　的　猫

dianzi shang de mao

a cat on the mat

$$X \rightarrow \langle X_0 \; de \; X_1, X_1 \; on \; X_0 \rangle$$

zhongguo de shoudu

capital of China

$$X \rightarrow \langle X_0 \; de \; X_1, X_1 \; of \; X_0 \rangle$$

# Translation Ambiguity

垫子　　上　的　猫

dianzi shang de mao

a cat on the mat

$$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ on \ X_0 \rangle$$

zhongguo de shoudu

capital of China

$$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ of \ X_0 \rangle$$

wo de mao

my cat

# Translation Ambiguity

垫子　　上　的　猫

dianzi shang de mao

a cat on the mat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

zhongguo de shoudu

capital of China

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

wo de mao

my cat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$$

# Translation Ambiguity

墊子　　上　的　猫

dianzi　shang　de　mao

a　cat　on　the　mat

$X \rightarrow \langle X_0 \ \mathrm{de} \ X_1, X_1 \ \mathrm{on} \ X_0 \rangle$

zhongguo　de　shoudu

capital　of　China

$X \rightarrow \langle X_0 \ \mathrm{de} \ X_1, X_1 \ \mathrm{of} \ X_0 \rangle$

wo　de　mao

my　cat

$X \rightarrow \langle X_0 \ \mathrm{de} \ X_1, X_0 \ X_1 \rangle$

zhifei　de　mao

zhifei　's　cat

# Translation Ambiguity

墊子　　上　的　猫

dianzi shang de mao

a cat on the mat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

zhongguo de shoudu

capital of China

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$
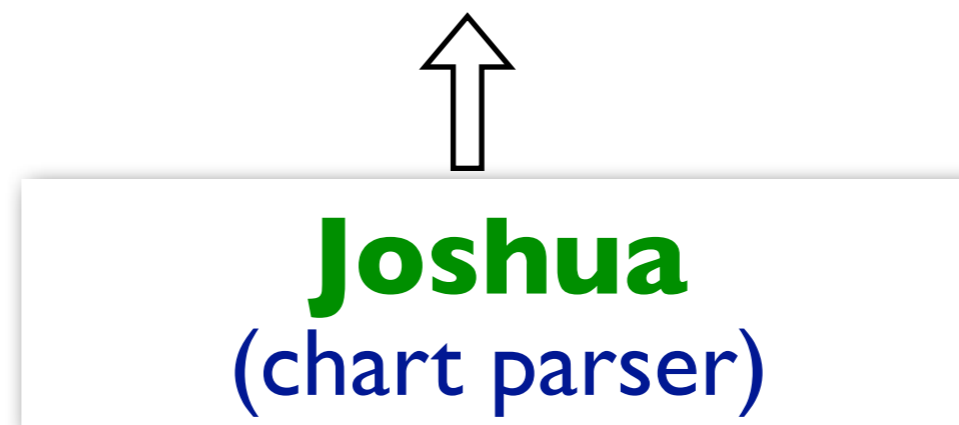
wo de mao

my cat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$$
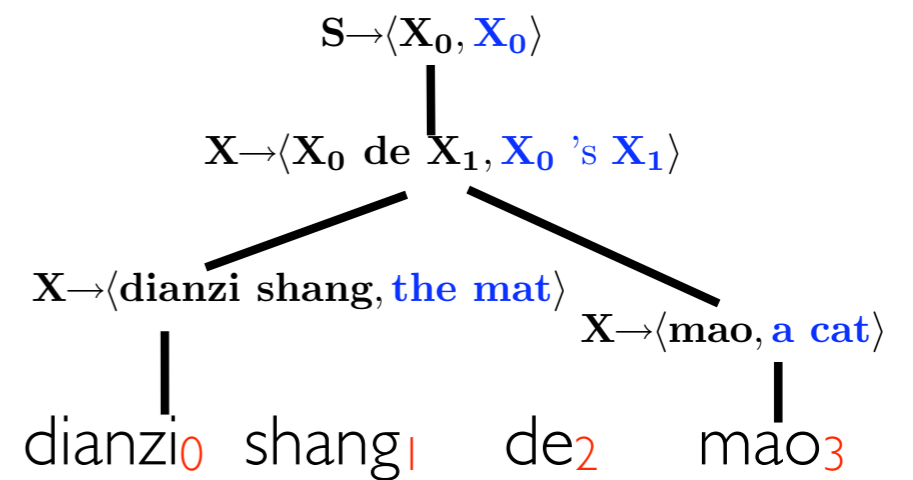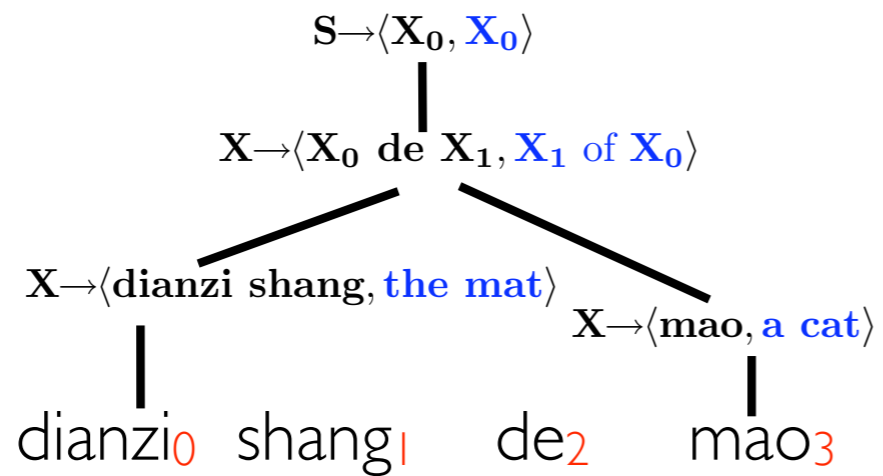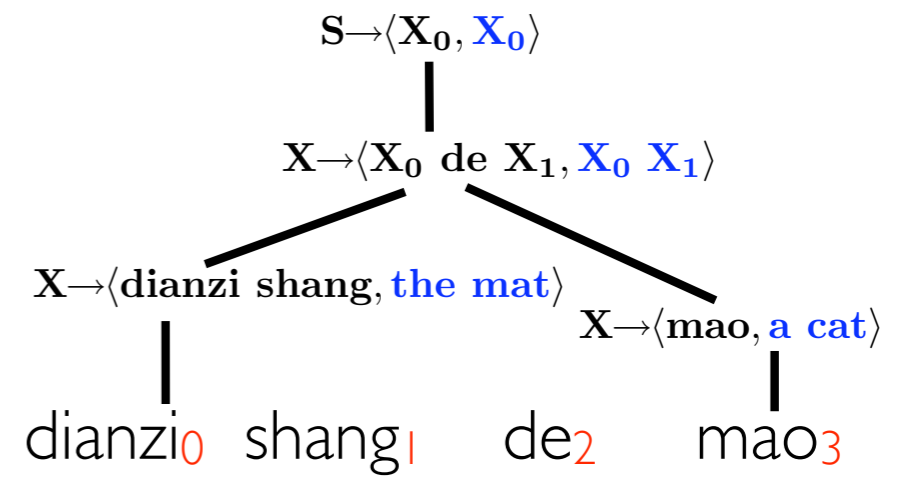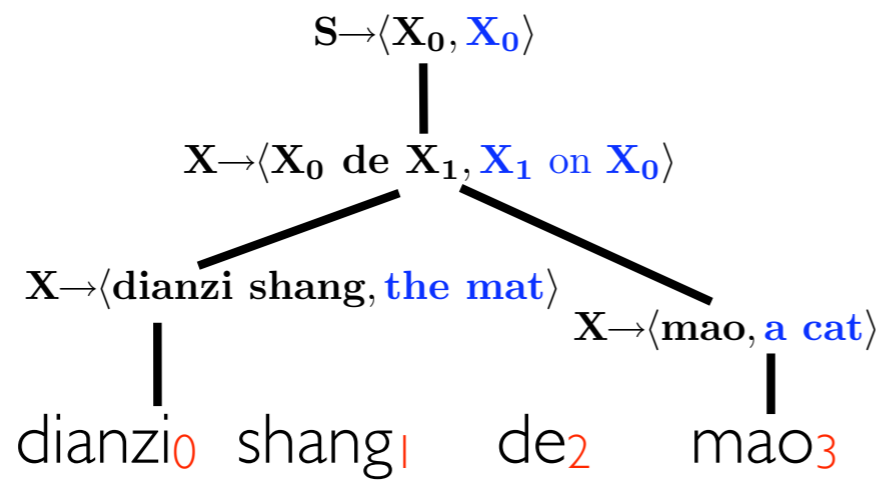
zhifei de mao

zhifei 's cat

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$$

9

dianzi  shang de mao

**Joshua**
(chart parser)

dianzi  shang de mao

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$

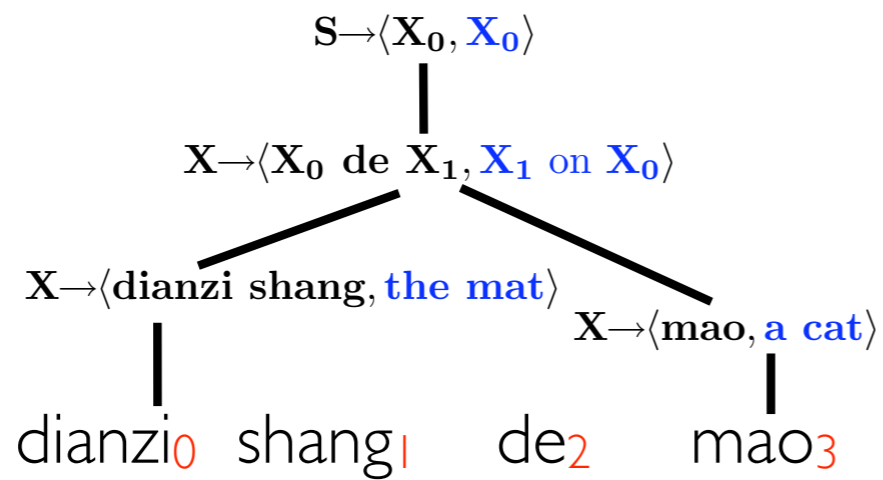$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

**Joshua**
(chart parser)

dianzi shang de mao

11

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

a cat on the mat

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

the mat a cat

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

a cat of the mat

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$
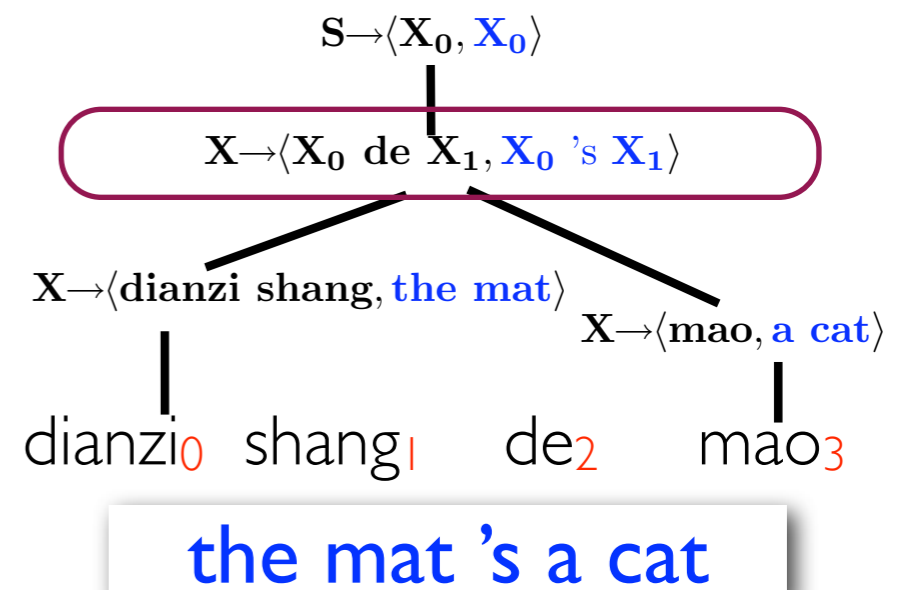
$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$ de$_2$ mao$_3$

the mat 's a cat

**Joshua**
(chart parser)

dianzi shang de mao

11

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \quad \text{de}_2 \quad \text{mao}_3$

**a cat on the mat**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \quad \text{de}_2 \quad \text{mao}_3$

**the mat a cat**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \quad \text{de}_2 \quad \text{mao}_3$

**a cat of the mat**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \quad \text{de}_2 \quad \text{mao}_3$

**the mat 's a cat**

**Joshua**
(chart parser)

dianzi shang de mao

11

hypergraph

$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \mid 0,2$

$X \mid 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$   shang$_1$                de$_2$      mao$_3$

**Joshua**
(chart parser)

dianzi  shang de mao

12

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \, X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$
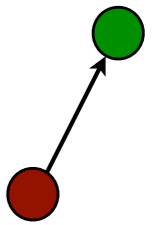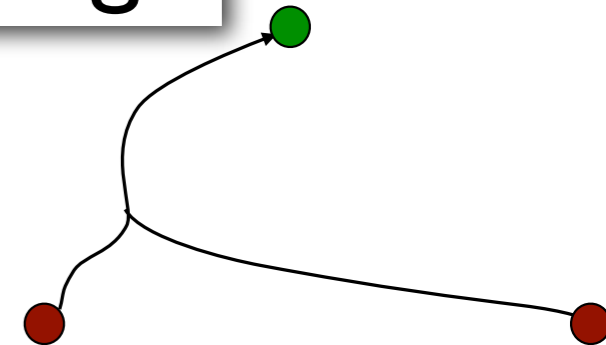
$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀  shang₁          de₂     mao₃

13

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



$S \quad 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \quad 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \quad 0,2$

$X \quad 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$ shang$_1$      de$_2$    mao$_3$

13

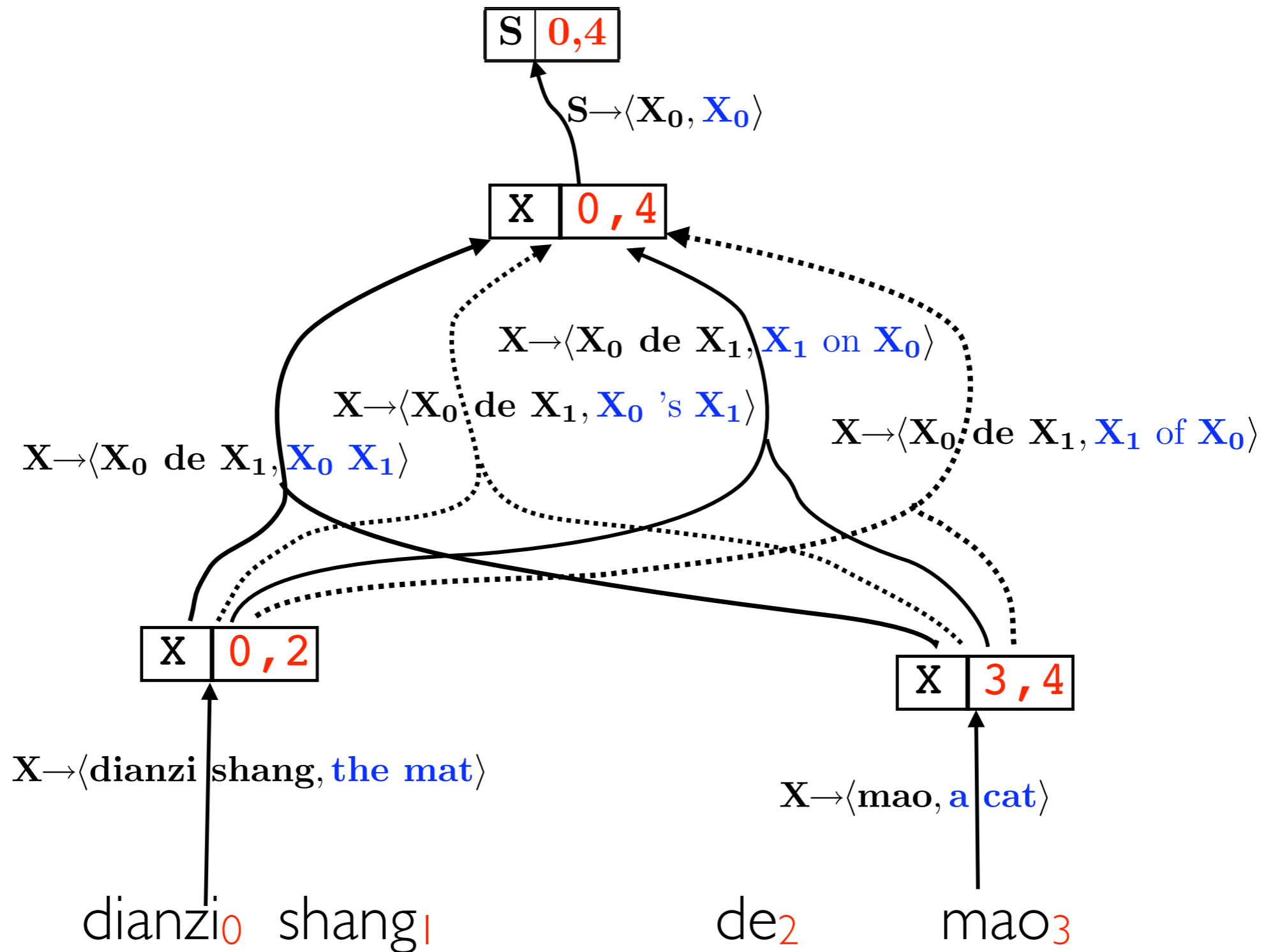# A hypergraph is a compact data structure to encode **exponentially many trees**.



S | 0,4

$S \rightarrow \langle X_0, X_0 \rangle$

X | 0,4

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

**node**

X | 0,2

X | 3,4

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$          de$_2$        mao$_3$

13

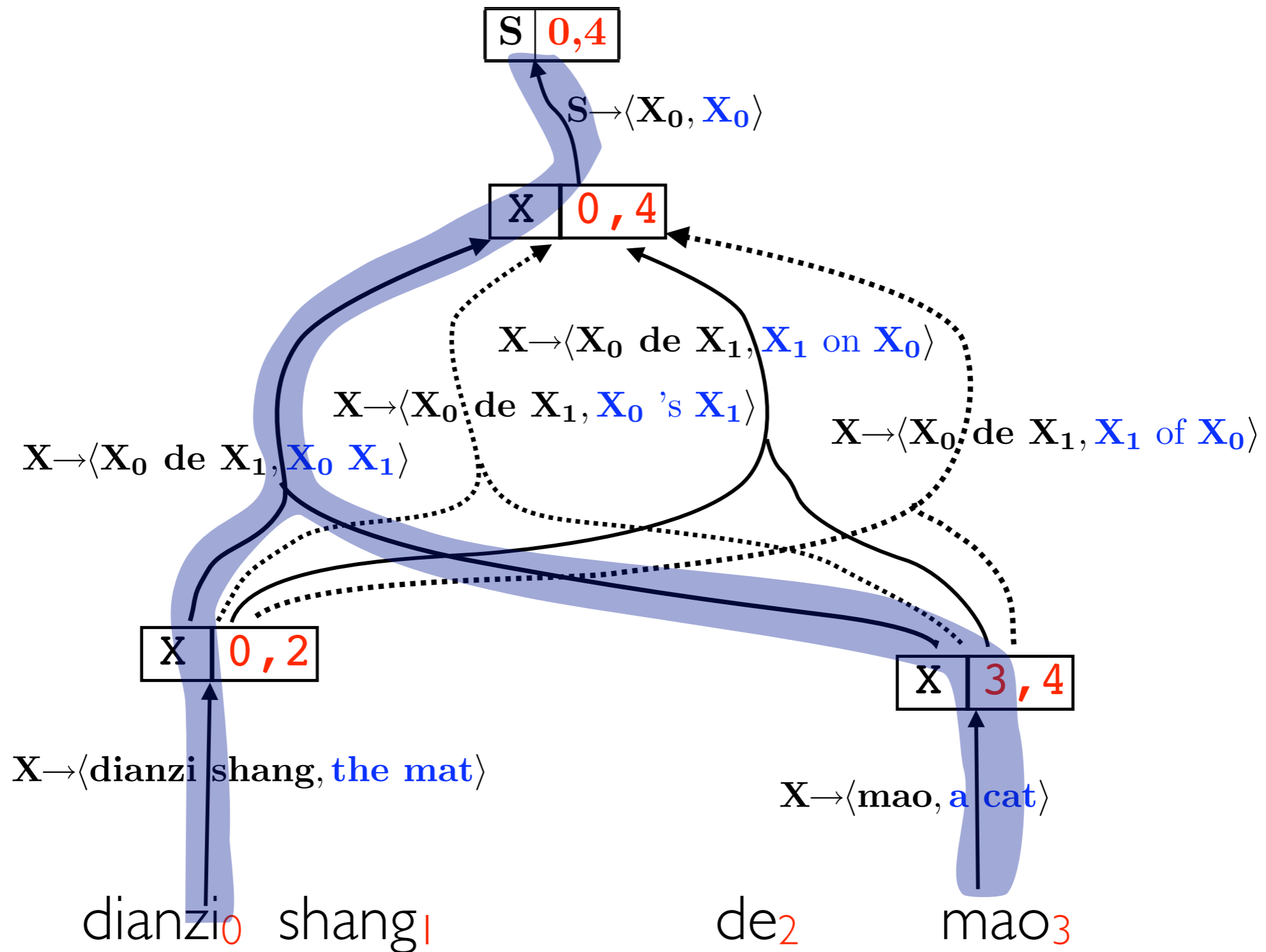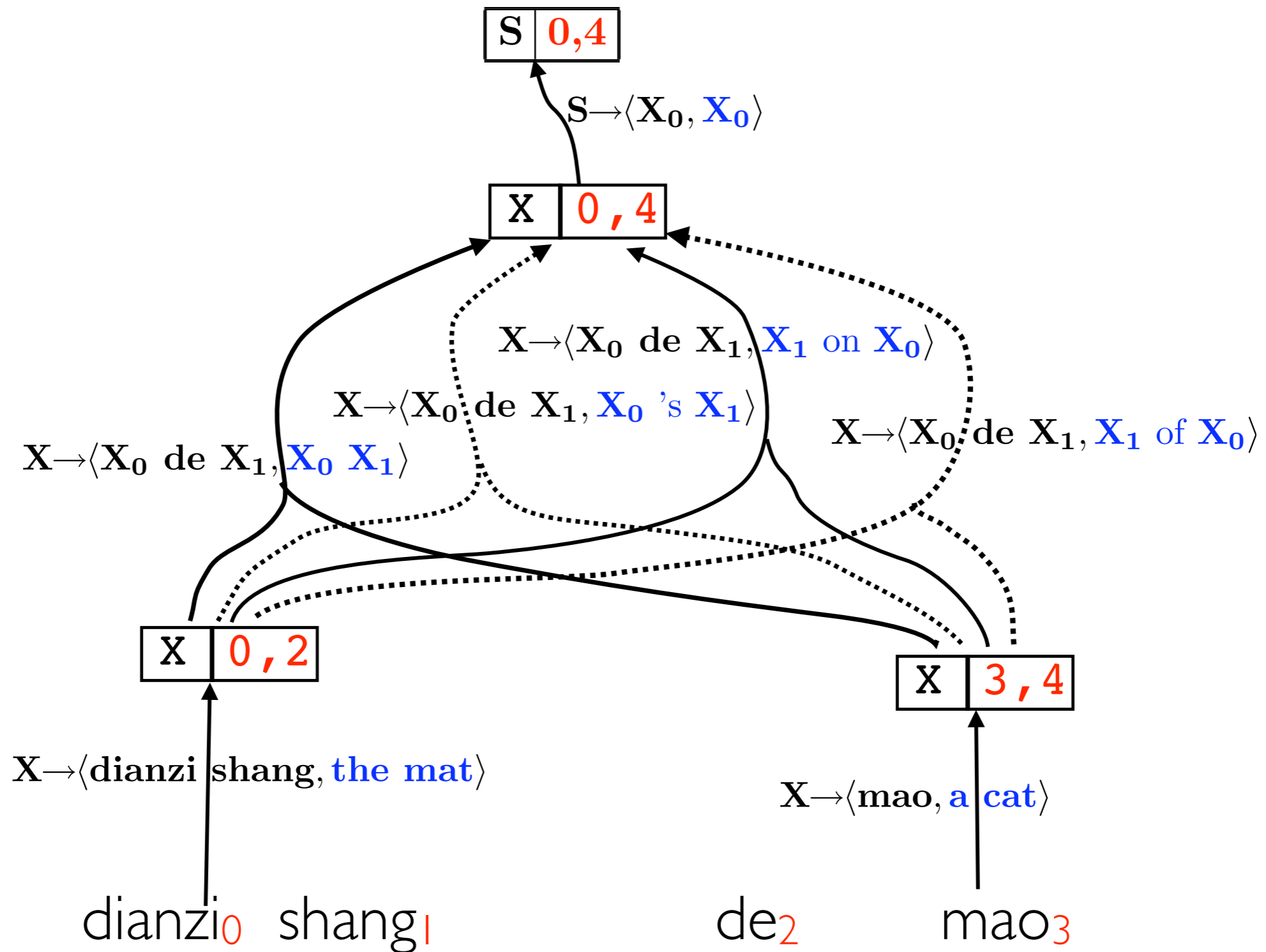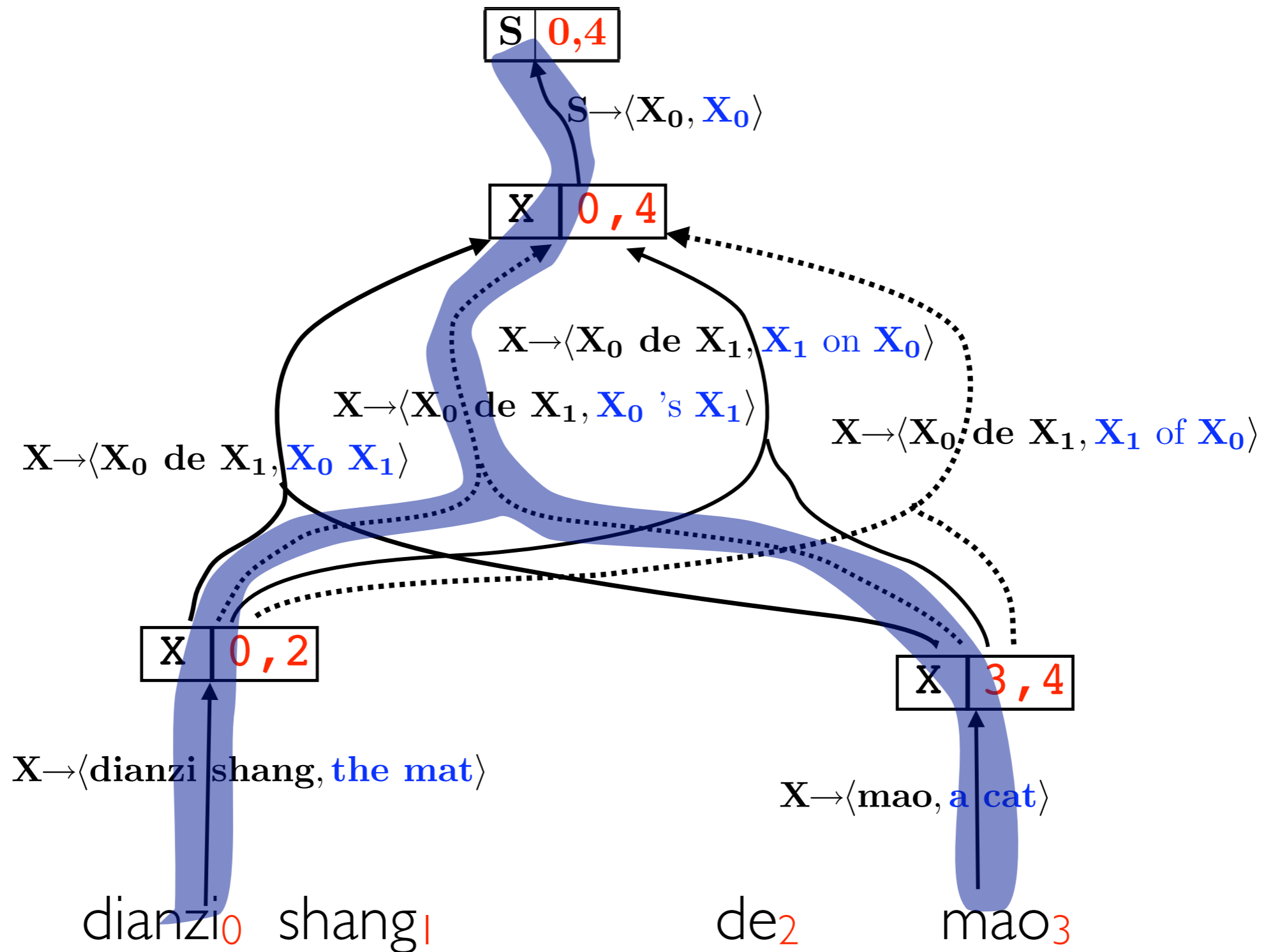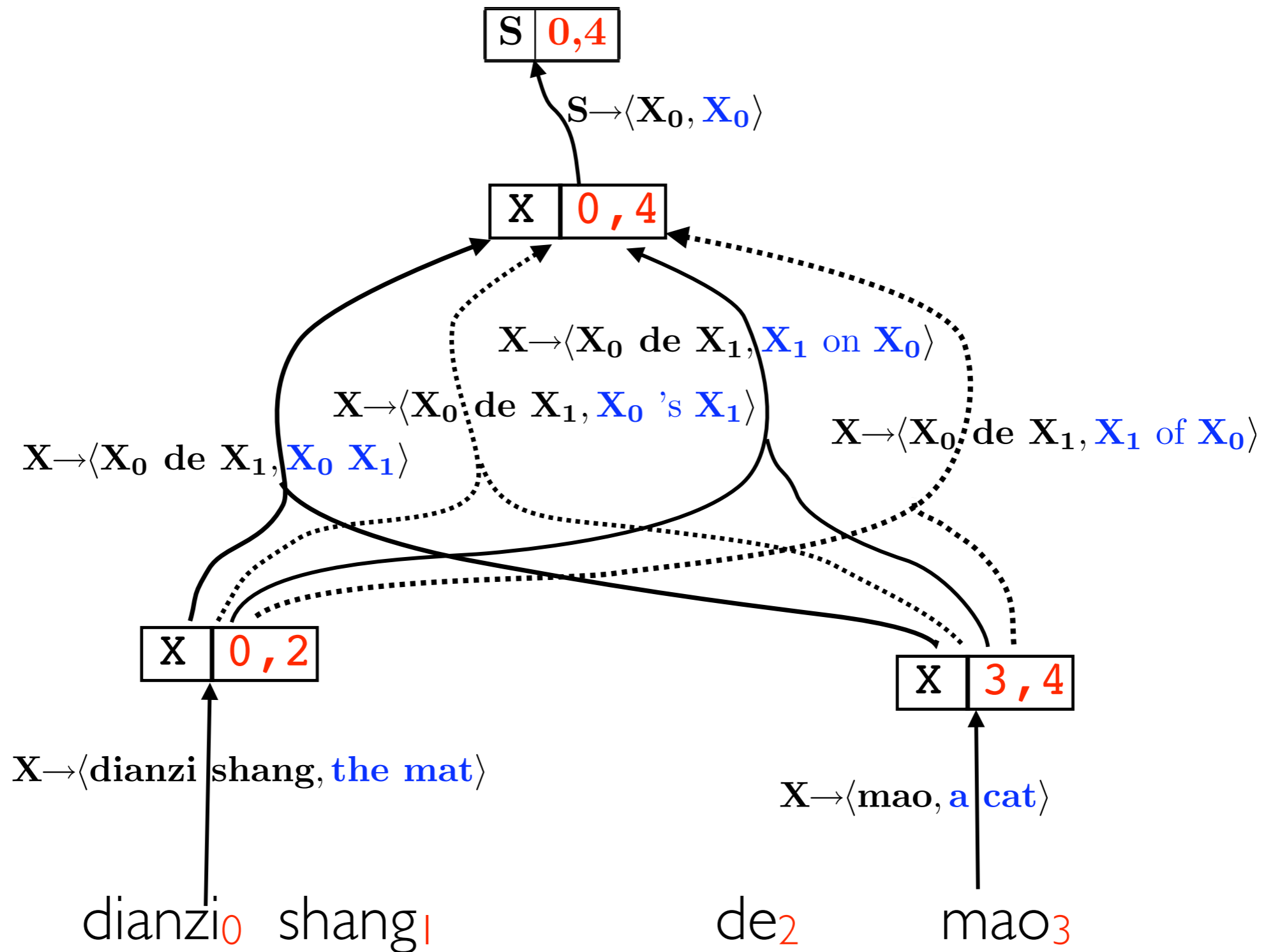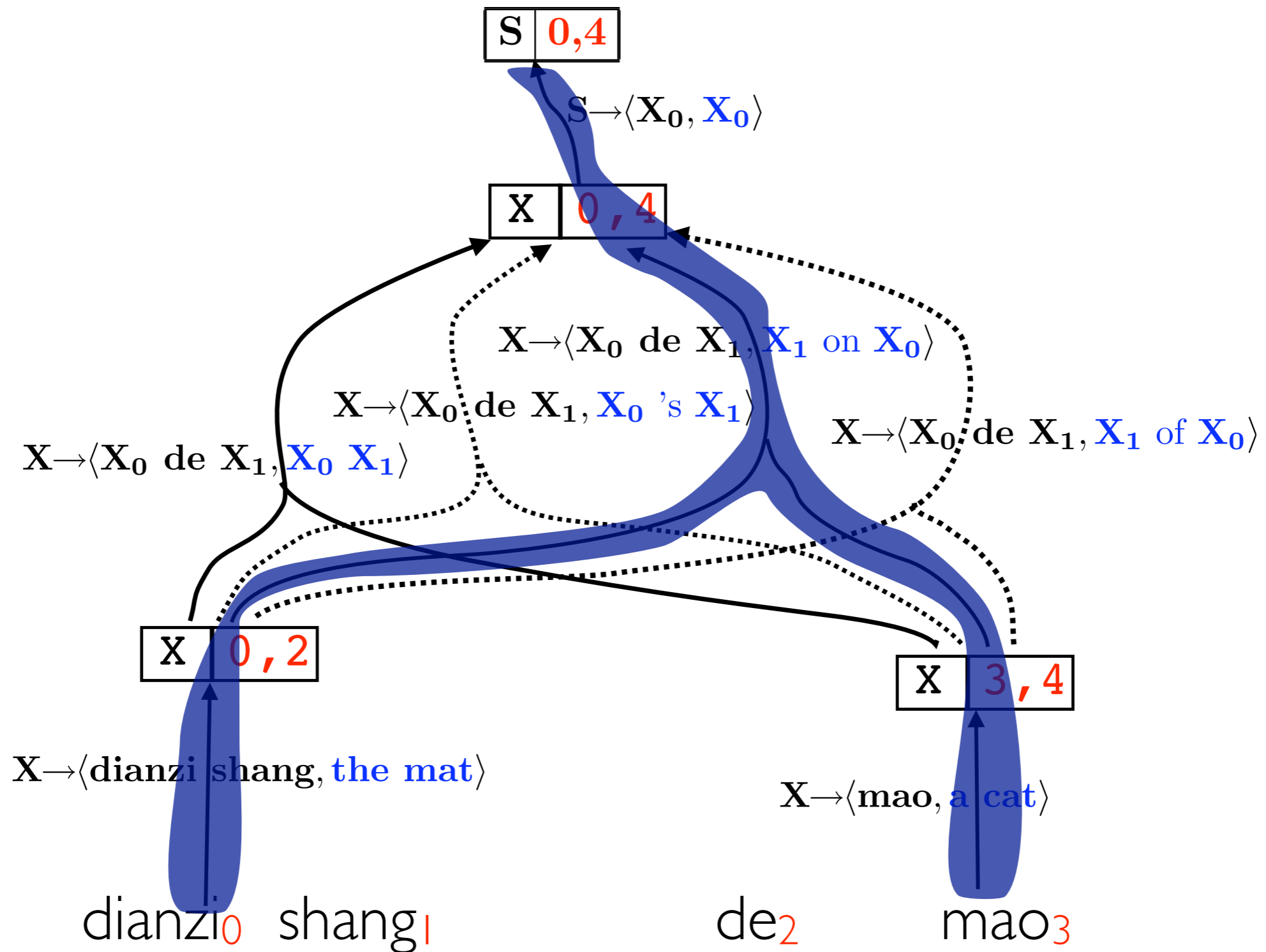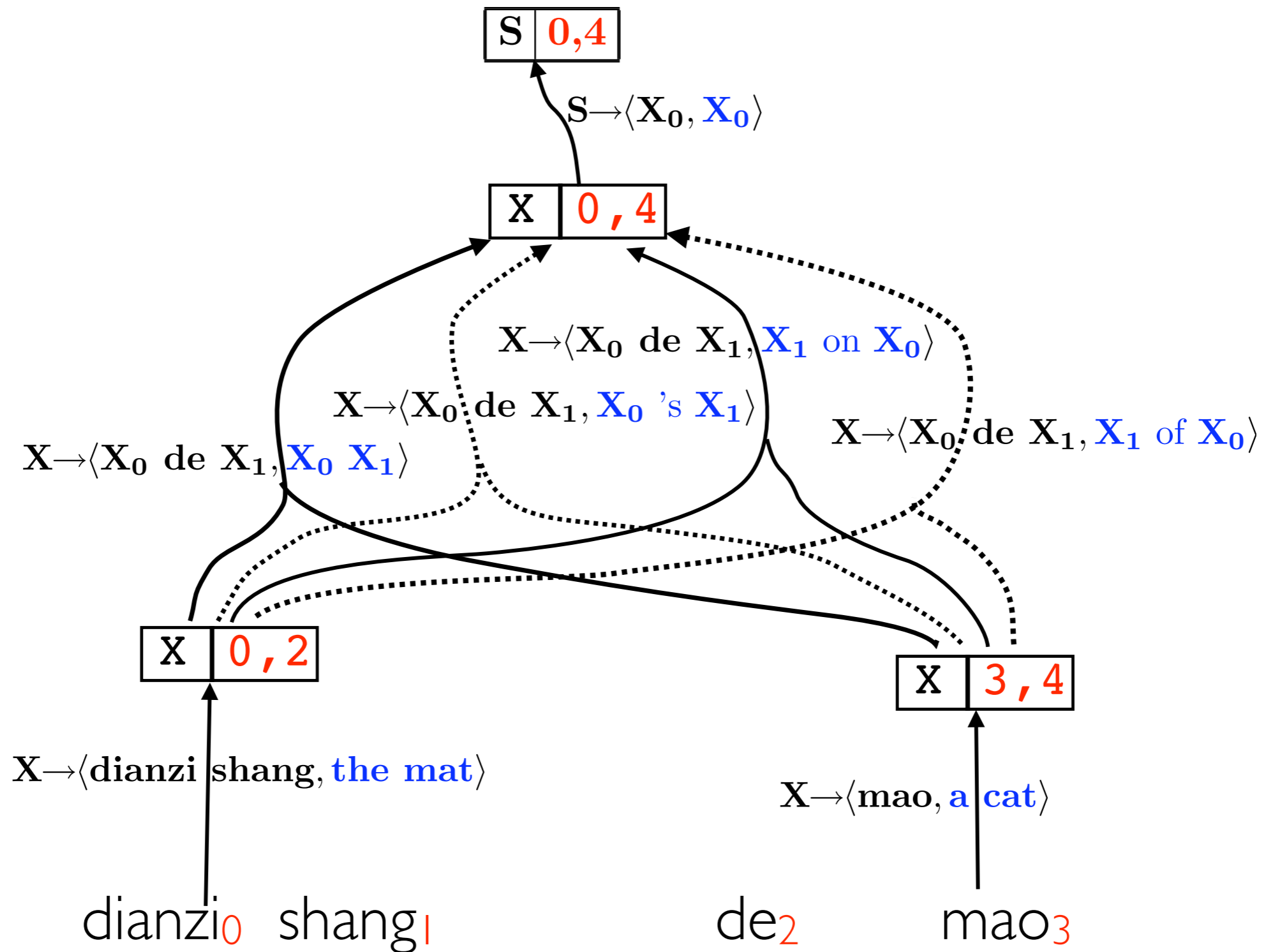# A hypergraph is a compact data structure to encode **exponentially many trees**.

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

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

edge



hyperedge



hyperedge

node

$$S \mid 0,4$$

$$S \rightarrow \langle X_0, X_0 \rangle$$

$$X \mid 0,4$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \, X_1 \rangle$$

$$X \mid 0,2$$

$$X \mid 3,4$$

$$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$$

$$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$$

$\text{dianzi}_0 \quad \text{shang}_1 \qquad\qquad \text{de}_2 \qquad \text{mao}_3$

13

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

edge

hyperedge

hyperedge

node

FSA

| S | 0,4 |

$S \rightarrow \langle X_0, X_0 \rangle$

| X | 0,4 |

$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ on \ X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_0 \ 's \ X_1 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ of \ X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_0 \ X_1 \rangle$

| X | 0,2 |

| X | 3,4 |

$X \rightarrow \langle dianzi \ shang, the \ mat \rangle$

$X \rightarrow \langle mao, a \ cat \rangle$

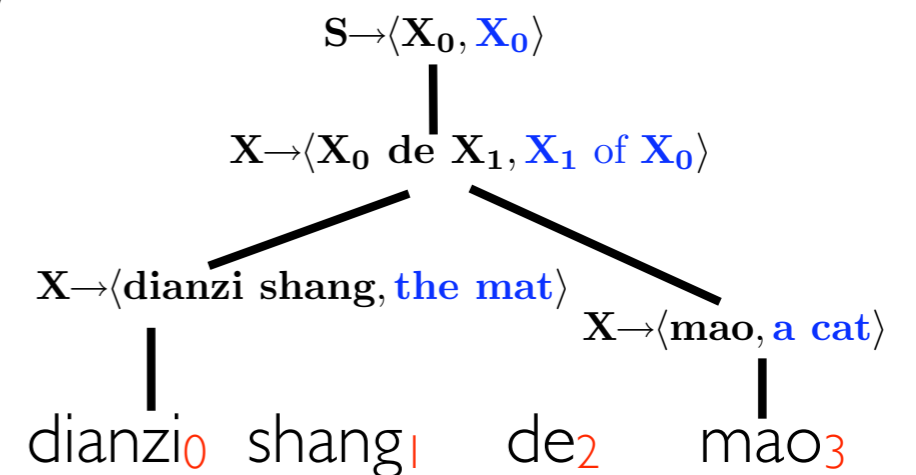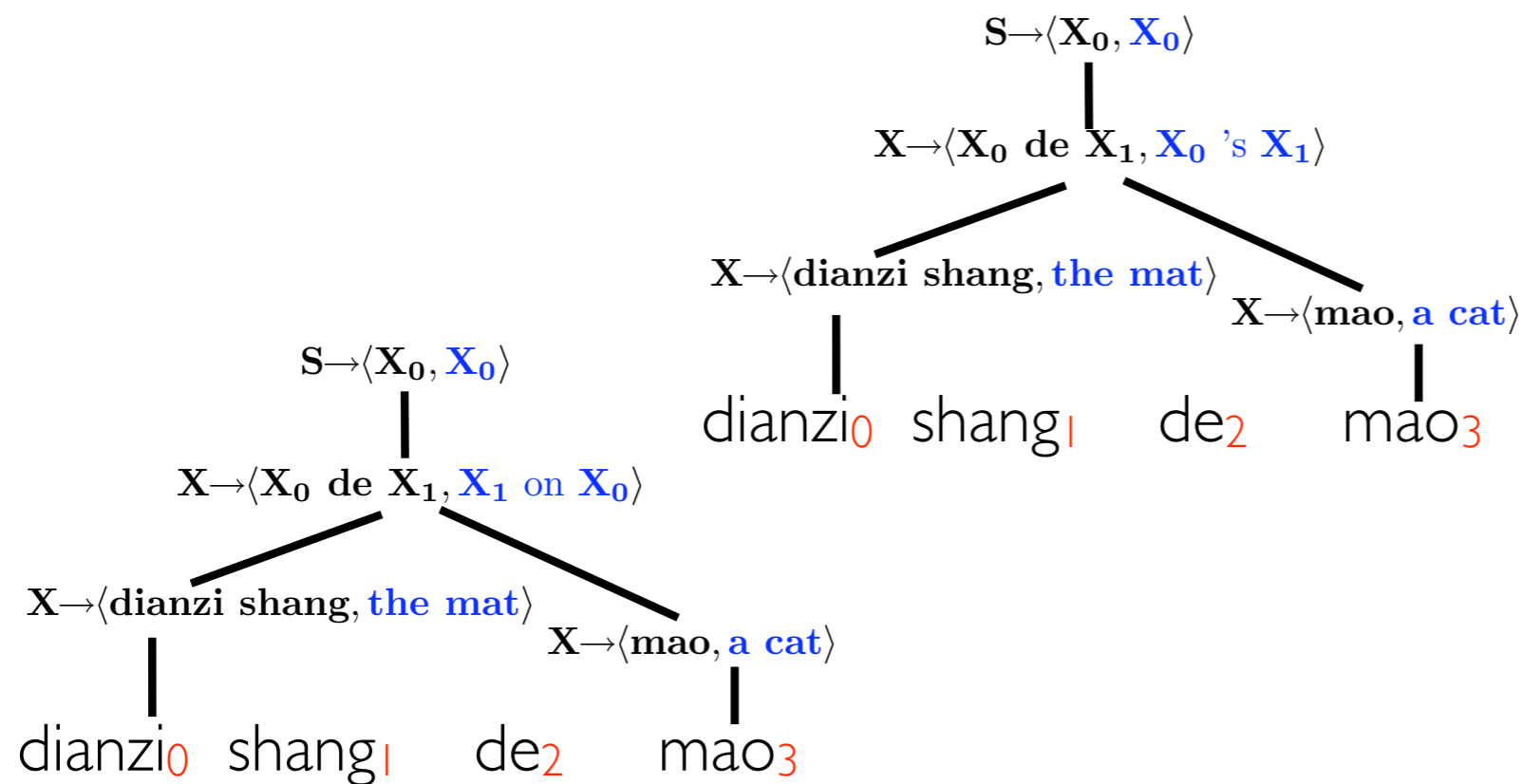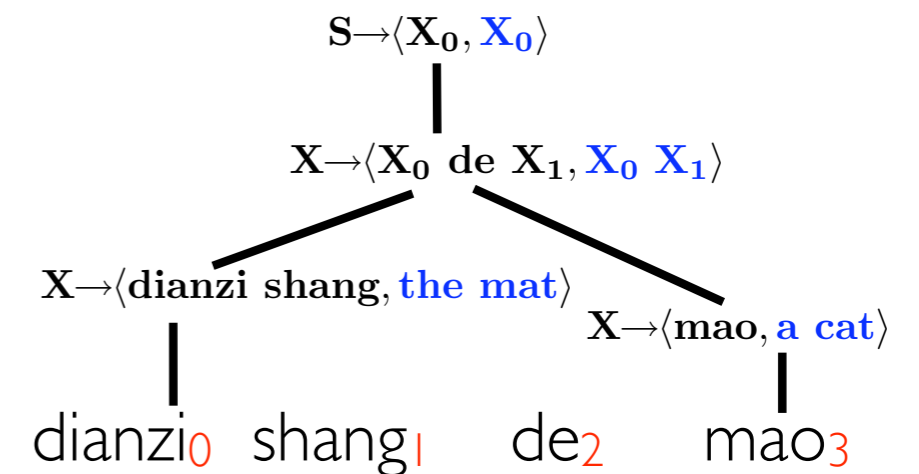$dianzi_0 \quad shang_1 \qquad de_2 \qquad mao_3$

13

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

edge



hyperedge



| S | 0,4 |

$S \rightarrow \langle X_0, X_0 \rangle$

hyperedge

| X | 0,4 |

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

node

| X | 0,2 |

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

| X | 3,4 |

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

FSA

Packed Forest

$\text{dianzi}_0 \quad \text{shang}_1 \qquad \text{de}_2 \qquad \text{mao}_3$

13

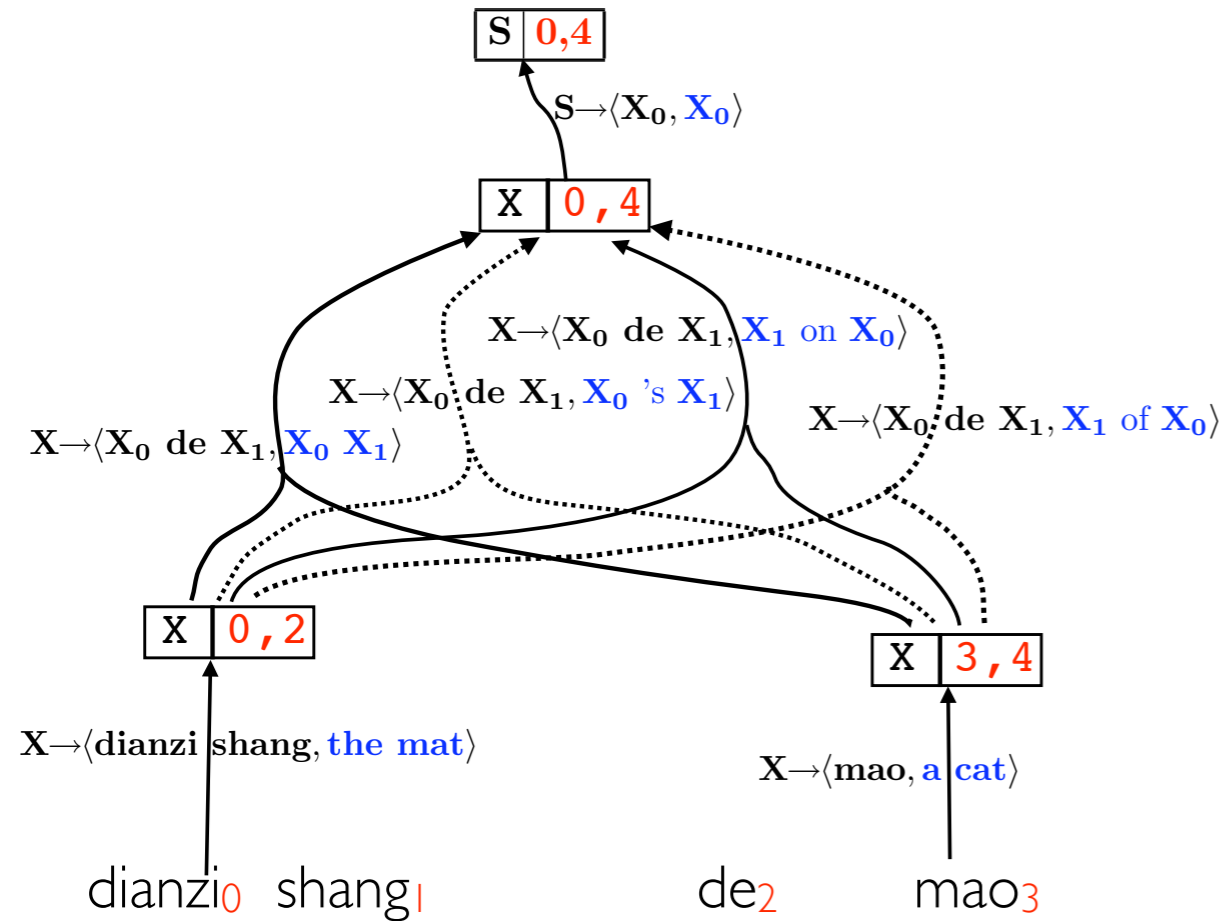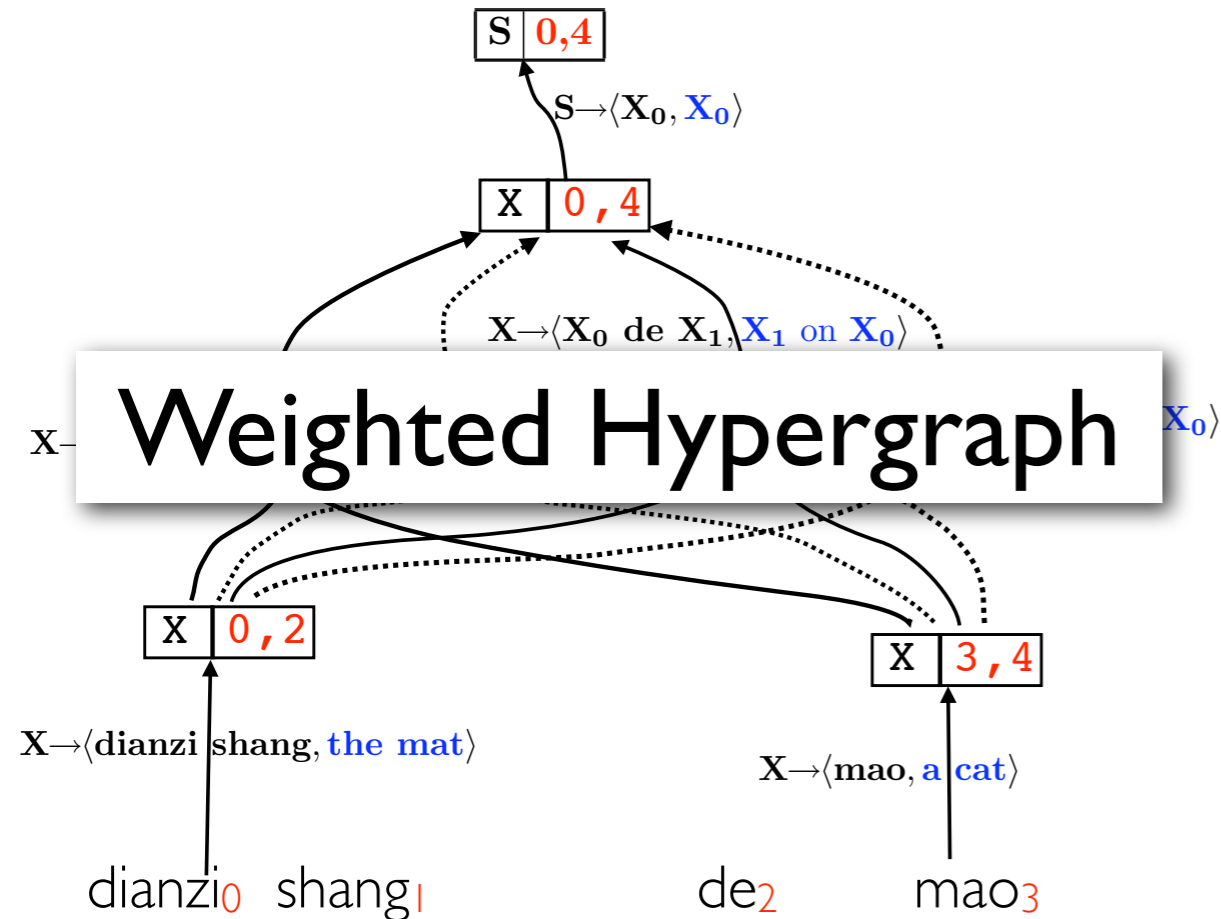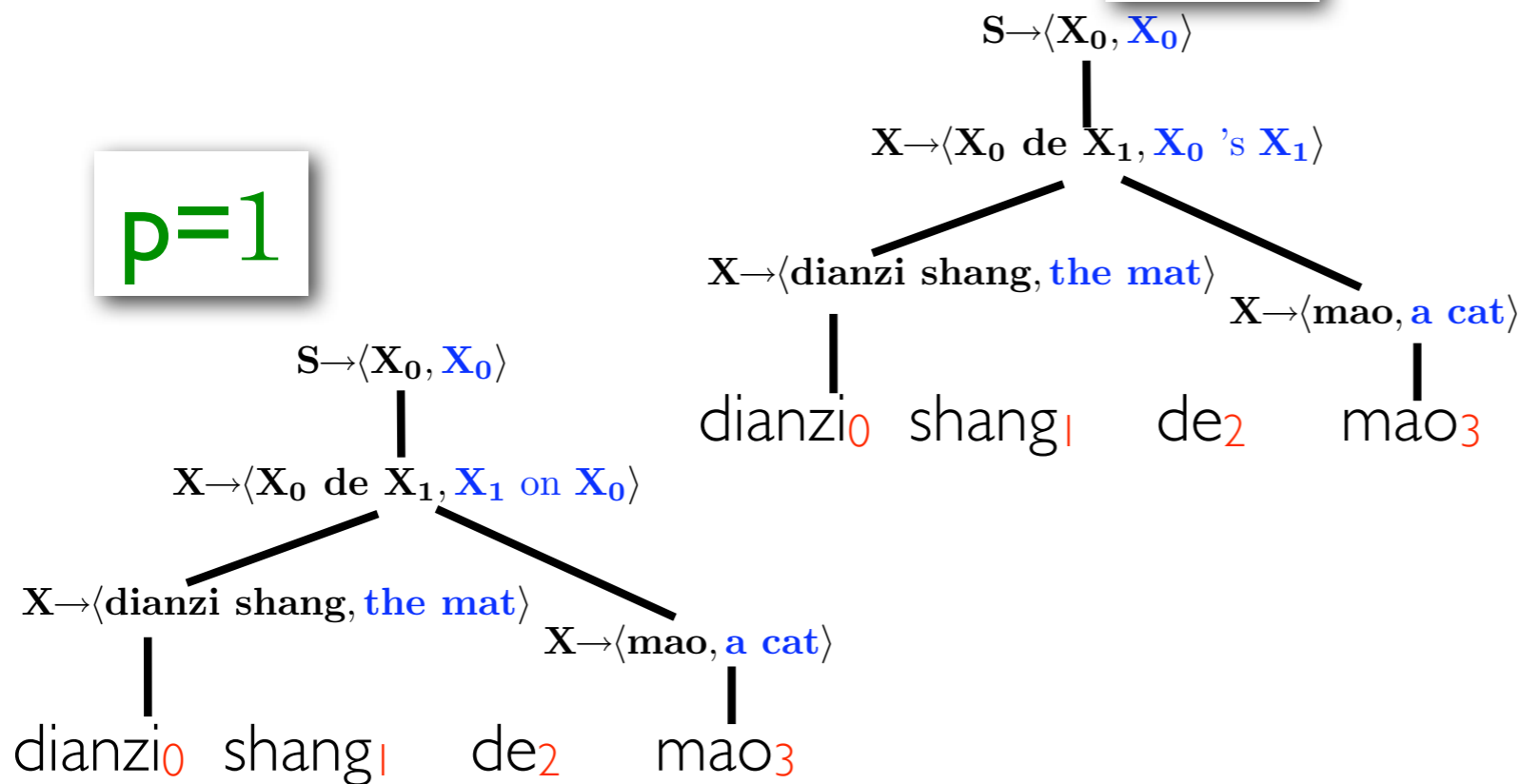# A hypergraph is a compact data structure to encode **exponentially many trees**.



| S | 0,4 |

$S \rightarrow \langle X_0, X_0 \rangle$

| X | 0,4 |

$X \rightarrow \langle X_0 \ \text{de} \ X_1, X_1 \ \text{on} \ X_0 \rangle$

$X \rightarrow \langle X_0 \ \text{de} \ X_1, X_0 \ \text{'s} \ X_1 \rangle$

$X \rightarrow \langle X_0 \ \text{de} \ X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \ \text{de} \ X_1, X_1 \ \text{of} \ X_0 \rangle$

| X | 0,2 |

| X | 3,4 |

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$        de$_2$    mao$_3$

14

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



$S \;|\; 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \;|\; 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \; X_1 \rangle$

$X \;|\; 0,2$

$X \;|\; 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \qquad \text{de}_2 \qquad \text{mao}_3$

14

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



$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \mid 0,2$

$X \mid 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$        de$_2$    mao$_3$

15

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



S | 0,4

$S \rightarrow \langle X_0, X_0 \rangle$

X | 0,4

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

X | 0,2

X | 3,4

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$      de$_2$      mao$_3$

15

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



$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \mid 0,2$

$X \mid 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$           de$_2$     mao$_3$

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



$$S \rightarrow \langle X_0, X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$$

$$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$$

$$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$$

dianzi$_0$  shang$_1$      de$_2$    mao$_3$

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



$$S \mid 0,4$$

$$S \rightarrow \langle X_0, X_0 \rangle$$

$$X \mid 0,4$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \mid 0,2$$

$$X \mid 3,4$$

$$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$$

$$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$$

dianzi$_0$  shang$_1$          de$_2$     mao$_3$

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



| S | 0,4 |

$S \rightarrow \langle X_0, X_0 \rangle$

| X | 0,4 |

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

| X | 0,2 |

| X | 3,4 |

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$          de$_2$     mao$_3$

17

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

**Structure sharing**



$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \mid 0,2$

$X \mid 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \qquad \text{de}_2 \quad \text{mao}_3$

17

# Why Hypergraphs?

- Contains a much larger hypothesis space than a k-best list

- General compact data structure

  - special cases include

    - finite state machine (e.g., lattice),

    - and/or graph

    - packed forest

  - can be used for speech, parsing, tree-based MT systems, and many more

S | 0,4

$S \rightarrow \langle X_0, X_0 \rangle$

X | 0,4

$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ on \ X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_0 \text{'s} \ X_1 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ of \ X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_0 \ X_1 \rangle$

X | 0,2

X | 3,4

$X \rightarrow \langle dianzi \ shang, the \ mat \rangle$

$X \rightarrow \langle mao, a \ cat \rangle$

dianzi₀  shang₁     de₂     mao₃

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle dianzi \ shang, the \ mat \rangle$

$X \rightarrow \langle mao, a \ cat \rangle$

dianzi₀  shang₁     de₂     mao₃

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_0 \text{'s} \ X_1 \rangle$

$X \rightarrow \langle dianzi \ shang, the \ mat \rangle$

$X \rightarrow \langle mao, a \ cat \rangle$

dianzi₀  shang₁     de₂     mao₃

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ on \ X_0 \rangle$

$X \rightarrow \langle dianzi \ shang, the \ mat \rangle$

$X \rightarrow \langle mao, a \ cat \rangle$

dianzi₀  shang₁     de₂     mao₃

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \ de \ X_1, X_1 \ of \ X_0 \rangle$

$X \rightarrow \langle dianzi \ shang, the \ mat \rangle$

$X \rightarrow \langle mao, a \ cat \rangle$

dianzi₀  shang₁     de₂     mao₃

**Weighted Hypergraph**

S | 0,4

S→⟨X₀, X₀⟩

X | 0,4

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨ ... , X₀⟩

X | 0,2

X→⟨dianzi shang, the mat⟩

X | 3,4

X→⟨mao, a cat⟩

dianzi₀  shang₁        de₂    mao₃

**p=1**

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨dianzi shang, the mat⟩   X→⟨mao, a cat⟩

dianzi₀ shang₁    de₂   mao₃

**p=3**

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ 's X₁⟩

X→⟨dianzi shang, the mat⟩   X→⟨mao, a cat⟩

dianzi₀ shang₁    de₂   mao₃

**p=2**

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ X₁⟩

X→⟨dianzi shang, the mat⟩   X→⟨mao, a cat⟩

dianzi₀ shang₁    de₂   mao₃

**p=2**

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ of X₀⟩

X→⟨dianzi shang, the mat⟩   X→⟨mao, a cat⟩

dianzi₀ shang₁    de₂   mao₃

19

**S** | **0,4**

$S \rightarrow \langle X_0, X_0 \rangle$

**X** | **0,4**

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

## Weighted Hypergraph

$X_0 \rangle$

**X** | **0,2**

**X** | **3,4**

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$          de$_2$    mao$_3$

## Linear model:

$$p(d \mid x) = \theta \cdot \Phi(d, x)$$

**p=2**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$    de$_2$    mao$_3$

**p=3**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$    de$_2$    mao$_3$

**p=1**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$    de$_2$    mao$_3$

**p=2**

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$    de$_2$    mao$_3$

19

Linear model:

foreign input

$$p(d \mid x) = \theta \cdot \Phi(d, x)$$

19

**S** **0,4**

S→⟨**X₀**, **X₀**⟩

**X** **0,4**

X→⟨**X₀ de X₁**, **X₁ on X₀**⟩

Weighted Hypergraph

X→ **X₀**

**X** **0,2**

**X** **3,4**

X→⟨**dianzi shang**, **the mat**⟩

X→⟨**mao**, **a cat**⟩

dianzi₀ shang₁     de₂     mao₃

Linear model:

derivation    foreign input

$$p(d \mid x) = \theta \cdot \Phi(d, x)$$

**p=2**

S→⟨**X₀**, **X₀**⟩

X→⟨**X₀ de X₁**, **X₀ X₁**⟩

X→⟨**dianzi shang**, **the mat**⟩

X→⟨**mao**, **a cat**⟩

dianzi₀ shang₁     de₂     mao₃

**p=3**

S→⟨**X₀**, **X₀**⟩

X→⟨**X₀ de X₁**, **X₀ 's X₁**⟩

X→⟨**dianzi shang**, **the mat**⟩

X→⟨**mao**, **a cat**⟩

dianzi₀ shang₁     de₂     mao₃

**p=1**

S→⟨**X₀**, **X₀**⟩

X→⟨**X₀ de X₁**, **X₁ on X₀**⟩

X→⟨**dianzi shang**, **the mat**⟩

X→⟨**mao**, **a cat**⟩

dianzi₀ shang₁     de₂     mao₃

**p=2**

S→⟨**X₀**, **X₀**⟩

X→⟨**X₀ de X₁**, **X₁ of X₀**⟩

X→⟨**dianzi shang**, **the mat**⟩

X→⟨**mao**, **a cat**⟩

dianzi₀ shang₁     de₂     mao₃

19

Linear model:

$$p(d \mid x) = \theta \cdot \Phi(d, x)$$

derivation

foreign input

features

Weighted Hypergraph

S 0,4

S→⟨X₀, X₀⟩

X 0,4

X→⟨X₀ de X₁, X₁ on X₀⟩

X₀⟩

X 0,2

X 3,4

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁      de₂      mao₃

p=2

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ X₁⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁      de₂      mao₃

p=3

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ 's X₁⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁      de₂      mao₃

p=1

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁      de₂      mao₃

p=2

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ of X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁      de₂      mao₃

19

**S** `0,4`

$S \rightarrow \langle X_0, X_0 \rangle$

**X** `0,4`

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

## Weighted Hypergraph

$X_0 \rangle$

**X** `0,2`

**X** `3,4`

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$dianzi_0 \quad shang_1 \qquad de_2 \qquad mao_3$

## Linear model:

derivation

foreign input

$$p(d \mid x) = \theta \cdot \Phi(d, x)$$

weights

features

p=2

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$dianzi_0 \quad shang_1 \qquad de_2 \qquad mao_3$

p=3

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$dianzi_0 \quad shang_1 \qquad de_2 \qquad mao_3$

p=1

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$dianzi_0 \quad shang_1 \qquad de_2 \qquad mao_3$

p=2

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$dianzi_0 \quad shang_1 \qquad de_2 \qquad mao_3$

Log-linear model:

$$p(d \mid x) = \frac{e^{\theta \cdot \Phi(d,x)}}{Z(x)}$$

Z=2+1+3+2=8

Probabilistic Hypergraph

S | 0,4

S→⟨X₀, X₀⟩

X | 0,4

X→⟨X₁ of X₀⟩

X | 0,2

X→⟨dianzi shang, the mat⟩

X | 3,4

X→⟨mao, a cat⟩

dianzi₀  shang₁       de₂    mao₃

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ X₁⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ 's X₁⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ of X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

20

S | 0,4

S→⟨X₀, X₀⟩

$$p(d \mid x) = \frac{e^{\theta \cdot \Phi(d,x)}}{Z(x)}$$

X | 0,4

Probabilistic Hypergraph

**Log-linear model:**

X₁ of X₀

X→⟨X

Z=2+1+3+2=8

X | 0,2

X | 3,4

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁          de₂     mao₃

p=2/8

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ X₁⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

p=3/8

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₀ 's X₁⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

p=1/8

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

p=2/8

S→⟨X₀, X₀⟩

X→⟨X₀ de X₁, X₁ of X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁    de₂    mao₃

Probabilistic Hypergraph

S | 0,4

$S \to \langle X_0, X_0 \rangle$

X | 0,4

$X \to \langle X$ ... $X_1$ of $X_0 \rangle$

X | 0,2

$X \to \langle dianzi\ shang, the\ mat \rangle$

dianzi₀ shang₁

X | 3,4

$X \to \langle mao, a\ cat \rangle$

de₂ mao₃

p=2/8

$S \to \langle X_0, X_0 \rangle$

$X \to \langle X_0\ de\ X_1, X_0\ X_1 \rangle$

$X \to \langle dianzi\ shang, the\ mat \rangle$

$X \to \langle mao, a\ cat \rangle$

dianzi₀ shang₁ de₂ mao₃

p=3/8

$S \to \langle X_0, X_0 \rangle$

$X \to \langle X_0\ de\ X_1, X_0\ 's\ X_1 \rangle$

$X \to \langle dianzi\ shang, the\ mat \rangle$

$X \to \langle mao, a\ cat \rangle$

dianzi₀ shang₁ de₂ mao₃

p=1/8

$S \to \langle X_0, X_0 \rangle$

$X \to \langle X_0\ de\ X_1, X_1\ on\ X_0 \rangle$

$X \to \langle dianzi\ shang, the\ mat \rangle$

$X \to \langle mao, a\ cat \rangle$

dianzi₀ shang₁ de₂ mao₃

p=2/8

$S \to \langle X_0, X_0 \rangle$

$X \to \langle X_0\ de\ X_1, X_1\ of\ X_0 \rangle$

$X \to \langle dianzi\ shang, the\ mat \rangle$

$X \to \langle mao, a\ cat \rangle$

dianzi₀ shang₁ de₂ mao₃

21

The hypergraph defines a probability distribution over **trees**!

Probabilistic Hypergraph

$S \rightarrow \langle X_0, X_0 \rangle$

$X$ 0,4

$X \rightarrow \langle X$ ... $X_1$ of $X_0 \rangle$

$X$ 0,2

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X$ 3,4

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi₀ shang₁

de₂ mao₃

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi₀ shang₁    de₂    mao₃

p=3/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_0$ 's $X_1 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi₀ shang₁    de₂    mao₃

p=1/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_1$ on $X_0 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi₀ shang₁    de₂    mao₃

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_1$ of $X_0 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi₀ shang₁    de₂    mao₃

21

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

Probabilistic Hypergraph

S→⟨X₀, X₀⟩

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁

de₂ mao₃

p=2/8

S→⟨X₀, X₀⟩
X→⟨X₀ de X₁, X₀ X₁⟩
X→⟨dianzi shang, the mat⟩
X→⟨mao, a cat⟩
dianzi₀ shang₁  de₂  mao₃

p=3/8

S→⟨X₀, X₀⟩
X→⟨X₀ de X₁, X₀ 's X₁⟩
X→⟨dianzi shang, the mat⟩
X→⟨mao, a cat⟩
dianzi₀ shang₁  de₂  mao₃

p=1/8

S→⟨X₀, X₀⟩
X→⟨X₀ de X₁, X₁ on X₀⟩
X→⟨dianzi shang, the mat⟩
X→⟨mao, a cat⟩
dianzi₀ shang₁  de₂  mao₃

p=2/8

S→⟨X₀, X₀⟩
X→⟨X₀ de X₁, X₁ of X₀⟩
X→⟨dianzi shang, the mat⟩
X→⟨mao, a cat⟩
dianzi₀ shang₁  de₂  mao₃

21

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4$

X→⟨X ... $X_1$ of $X_0$⟩

Probabilistic Hypergraph

$X \mid 0,2$

$X \mid 3,4$

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁        de₂     mao₃

S | 0,4

$S \rightarrow \langle X_0, X_0 \rangle$

X | 0,4

Probabilistic Hypergraph

$X \rightarrow \langle X$ ... $X_1$ of $X_0 \rangle$

X | 0,2

X | 3,4

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi$_0$  shang$_1$          de$_2$      mao$_3$

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

Which translation do we present to a user?          Decoding

22

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

Which translation do we present to a user?   Decoding

How do we set the parameters Θ?   Training

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

Which translation do we present to a user?  Decoding

How do we set the parameters Θ?  Training

What atomic operations do we need to perform?  Atomic Inference

22

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

| **training** (e.g., mert) | **decoding** (e.g., mbr) |
|---|---|
| **atomic inference operations** (e.g., finding one-best, k-best or expectation, inference can be *exact* or *approximate*) | |

Which translation do we present to a user?    Decoding

How do we set the parameters Θ?    Training

What atomic operations do we need to perform?    Atomic Inference

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

| training (e.g., mert) | decoding (e.g., mbr) |
|---|---|
| **atomic inference operations** (e.g., finding one-best, k-best or expectation, inference can be *exact* or *approximate*) | |

Which translation do we present to a user?     Decoding

How do we set the parameters Θ?     Training

What atomic operations do we need to perform?     Atomic Inference

Why are the problems difficult?

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

| **training**<br>(e.g., mert) | **decoding**<br>(e.g., mbr) |
|---|---|
| **atomic inference operations**<br>(e.g., finding one-best, k-best or expectation,<br>inference can be *exact* or *approximate*) | |

Which translation do we present to a user?    Decoding

How do we set the parameters Θ?    Training

What atomic operations do we need to perform? Atomic Inference

Why are the problems difficult?

- brute-force will be too slow as there are exponentially many trees, so require sophisticated dynamic programs

The hypergraph defines a probability distribution over **trees**!

the distribution is parameterized by Θ

| **training** (e.g., mert) | **decoding** (e.g., mbr) |
| --- | --- |
| **atomic inference operations** (e.g., finding one-best, k-best or expectation, inference can be *exact* or *approximate*) | |

Which translation do we present to a user?    Decoding

How do we set the parameters Θ?    Training

What atomic operations do we need to perform? Atomic Inference

Why are the problems difficult?

- brute-force will be too slow as there are exponentially many trees, so require sophisticated dynamic programs

- sometimes intractable, require approximations

# Inference, Training and Decoding on Hypergraphs

- **Atomic Inference**
  - finding one-best derivations

| Graph | Topological | Best-first | | |
|---|---|---|---|---|
| | | no heuristic | with heuristic | with hierarchy |
| FSA | Viterbi | Dijkstra | $A^*$ | HA$^*$ |
| Hypergraph | CYK | Knuth | Klein and Manning | Generalized $A^*$ |

  - finding k-best derivations
  - computing expectations (e.g., of features)

- **Training**
  - Perceptron, conditional random field (CRF), minimum error rate training (MERT), minimum risk, and MIRA

- **Decoding**
  - Viterbi decoding, maximum a posterior (MAP) decoding, and minimum Bayes risk (MBR) decoding

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training

  ‣ minimum imputed risk

  ‣ contrastive language model estimation

- Variational Decoding

- First- and Second-order Expectation Semirings

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training
  - ▸ minimum imputed risk
  - ▸ contrastive language model estimation

- Variational Decoding

- First- and Second-order Expectation Semirings

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training
  - ▸ minimum imputed risk
  - ▸ contrastive language model estimation

  } main focus

- Variational Decoding

- First- and Second-order Expectation Semirings

# Training Setup

- Each **training example** consists of
  - a foreign sentence (from which a hypergraph is generated to represent many possible translations)
  - a reference translation

    x: dianzi  shang de mao

    y: a  cat  on  the  mat

- **Training**
  - adjust the parameters Θ so that the reference translation is preferred by the model

$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X$

$X \mid 0,2$

$X \mid 3,4$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \qquad \text{de}_2 \quad \text{mao}_3$

# Supervised: Minimum Empirical Risk

# Supervised: Minimum Empirical Risk

- Minimum Empirical Risk Training

$$\theta^* = \arg\min_{\theta} \sum_{x,y} \tilde{p}(x,y) \mathrm{L}(\delta_\theta(x), y)$$

- Minimum Empirical Risk Training

$$\theta^* = \arg\min_\theta \sum_{x,y} \boxed{\tilde{p}(x,y)} \mathrm{L}(\delta_\theta(x), y)$$

empirical
distribution

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_\theta \sum_{x,y} \tilde{p}(x,y) \mathrm{L}(\delta_\theta(x), y)$$

empirical
distribution

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_\theta \sum_{x,y} \boxed{\tilde{p}(x,y)} \boxed{\mathrm{L}(\delta_\theta(x), y)}$$

empirical
distribution

X

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_\theta \sum_{x,y} \boxed{\tilde{p}(x,y)} \boxed{\mathrm{L}(\delta_\theta(x), y)}$$

empirical
distribution

$$x \longrightarrow \boxed{\delta_\theta}$$

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_\theta \sum_{x,y} \tilde{p}(x,y) \mathrm{L}(\delta_\theta(x), y)$$

empirical distribution

MT decoder

$$x \longrightarrow \boxed{\delta_\theta}$$

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_{\theta} \sum_{x,y} \boxed{\tilde{p}(x,y)}\boxed{\mathrm{L}(\delta_\theta(x), y)}$$

empirical distribution

MT decoder

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_{\theta} \sum_{x,y} \boxed{\tilde{p}(x,y)} \boxed{\mathrm{L}(\delta_\theta(x), y)}$$

empirical distribution

MT decoder

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

MT output

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_{\theta} \sum_{x,y} \tilde{p}(x,y) \, \mathrm{L}(\delta_{\theta}(x), y)$$

MT decoder

empirical distribution

$$x \longrightarrow \boxed{\delta_{\theta}} \longrightarrow \delta_{\theta}(x) \longleftrightarrow y$$

loss

MT output

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_{\theta} \sum_{x,y} \boxed{\tilde{p}(x,y)} \boxed{\mathrm{L}(\delta_\theta(x), y)}$$

empirical distribution

MT decoder

negated BLEU

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

MT output

26

# Supervised: Minimum Empirical Risk

- ## Minimum Empirical Risk Training

$$\theta^* = \arg\min_{\theta} \sum_{x,y} \tilde{p}(x,y) \mathrm{L}(\delta_{\theta}(x), y)$$

empirical distribution

MT decoder

negated BLEU

$$x \longrightarrow \boxed{\delta_{\theta}} \longrightarrow \delta_{\theta}(x) \longleftrightarrow \mathrm{loss} \quad y$$

MT output

- ## Uniform Empirical Distribution

$$\theta^* = \arg\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} \mathrm{L}(\delta_{\theta}(x_i), \tilde{y}_i)$$

# Supervised: Minimum Empirical Risk

- ## Minimum Empirical Risk Training

$$\theta^* = \arg\min_\theta \sum_{x,y} \tilde{p}(x,y) \mathrm{L}(\delta_\theta(x), y)$$

empirical distribution

MT decoder

negated BLEU

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \overset{\text{loss}}{\longleftrightarrow} y$$

MT output

- ## Uniform Empirical Distribution

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)$$

- MERT
- CRF
- Peceptron

# Supervised: Minimum Empirical Risk

- **Minimum Empirical Risk Training**

$$\theta^* = \arg\min_\theta \sum_{x,y} \tilde{p}(x,y) \, \mathrm{L}(\delta_\theta(x), y)$$

empirical distribution

MT decoder

negated BLEU

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \longleftrightarrow y$$

loss

MT output

- **Uniform Empirical Distribution**

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)$$

- MERT
- CRF
- Peceptron

**What if the input x is missing?**

# Unsupervised: Minimum Imputed Risk

- **Minimum Empirical Risk Training**

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftrightarrow{\text{loss}} y$$

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} L(\delta_\theta(x_i), \tilde{y}_i)$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \overset{\text{loss}}{\longleftrightarrow} y$$

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)$$

- Minimum Imputed Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training



$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$

- Minimum Imputed Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftrightarrow{\text{loss}} y$$

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$

- Minimum Imputed Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftrightarrow{\text{loss}} y$$

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$

- Minimum Imputed Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$$\tilde{y}_i$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$

$$\mathsf{x} \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftrightarrow{\text{loss}} \mathsf{y}$$

- Minimum Imputed Risk Training
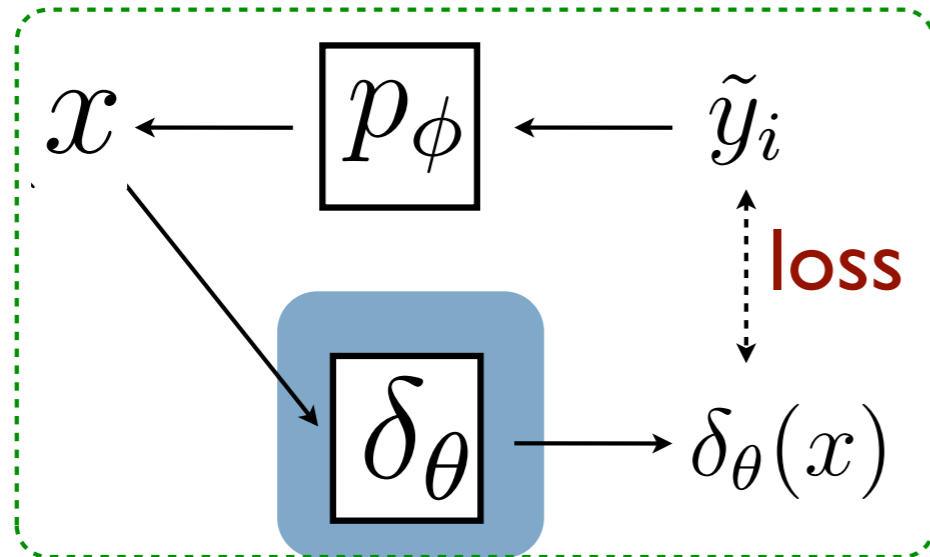
$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$p_\phi$: reverse model

$$\boxed{p_\phi} \longleftarrow \tilde{y}_i$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$

- Minimum Imputed Risk Training
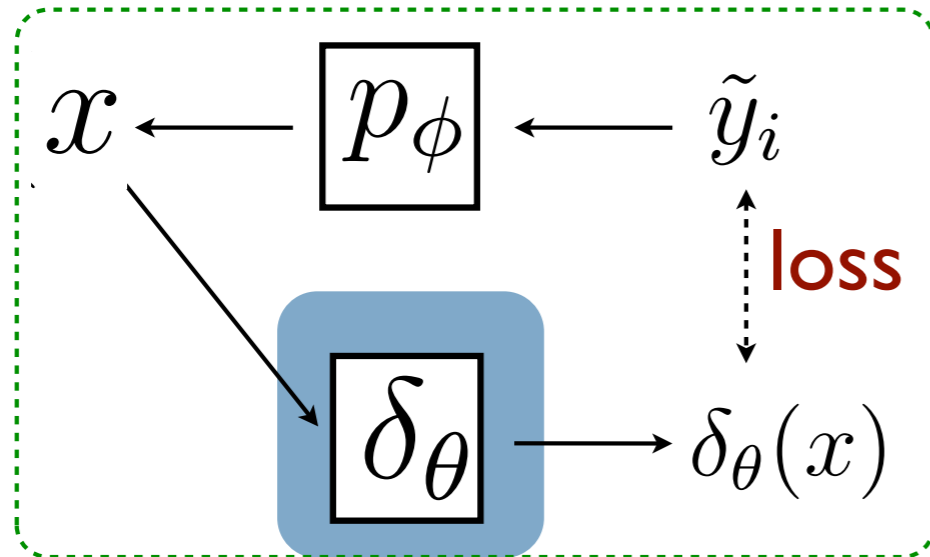
$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$p_\phi$: reverse model

$x$: imputed input

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

27

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$

- Minimum Imputed Risk Training
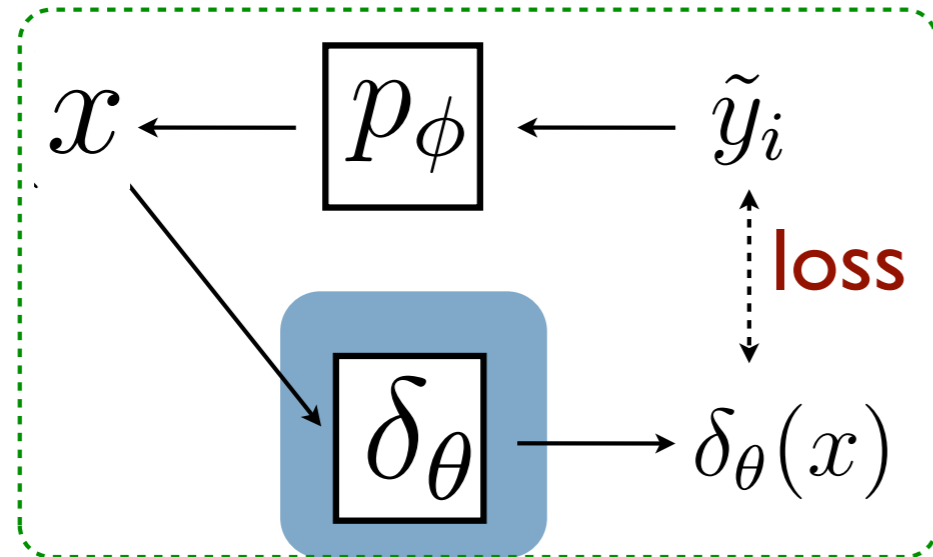
$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$p_\phi$: reverse model

$x$: imputed input

$\delta_\theta$: forward system

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

$$\text{loss}$$

$$\delta_\theta \longrightarrow \delta_\theta(x)$$

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$



- Minimum Imputed Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$p_\phi$: reverse model
$x$: imputed input
$\delta_\theta$: forward system



27

# Unsupervised: Minimum Imputed Risk

- Minimum Empirical Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$



x $\longrightarrow$ $\boxed{\delta_\theta}$ $\longrightarrow$ $\delta_\theta(x)$ $\longleftrightarrow$ loss $\longleftrightarrow$ y

- Minimum Imputed Risk Training

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$p_\phi$: reverse model

$x$: imputed input

$\delta_\theta$: forward system

Round trip translation



$x$ $\longleftarrow$ $\boxed{p_\phi}$ $\longleftarrow$ $\tilde{y}_i$

$\boxed{\delta_\theta}$ $\longrightarrow$ $\delta_\theta(x)$

loss

# Unsupervised: Minimum Imputed Risk

- **Minimum Empirical Risk Training**

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N}\sum_{i=1}^{N} \boxed{\mathrm{L}(\delta_\theta(x_i), \tilde{y}_i)}$$



- **Minimum Imputed Risk Training**

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N}\sum_{i=1}^{N} \boxed{\sum_x p_\phi(x \mid \tilde{y}_i)\mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

$p_\phi$: reverse model

$x$: imputed input

$\delta_\theta$: forward system

Round trip translation

Speech recognition?

# Training Reverse Model $p_\phi$

# Training Reverse Model $p_\phi$



Our goal is to train a good forward system $\delta_\theta$

# Training Reverse Model $p_\phi$



Our goal is to train a good forward system $\delta_\theta$

$p_\phi$ and $\delta_\theta$ are parameterized and trained separately

# Training Reverse Model $p_\phi$



Our goal is to train a good forward system $\delta_\theta$

$p_\phi$ and $\delta_\theta$ are parameterized and trained separately

$p_\phi$ is fixed when training $\delta_\theta$

# Approximating $p_\phi(x \mid \tilde{y}_i)$

# Approximating $p_\phi(x \mid \tilde{y}_i)$



SCFG

# Approximating $p_\phi(x \mid \tilde{y}_i)$



exponentially many x, stored in a hypergraph

SCFG

loss

# Approximating $p_\phi(x \mid \tilde{y}_i)$



exponentially
many x, stored in
a hypergraph

$x$   $p_\phi$   $\tilde{y}_i$

loss

$\delta_\theta$   $\longrightarrow$   $\delta_\theta(x)$

SCFG     SCFG

# Approximating $p_\phi(x \mid \tilde{y}_i)$



exponentially many x, stored in a hypergraph

SCFG

SCFG

CFG is not closed under composition!

# Approximating $p_\phi(x \mid \tilde{y}_i)$



exponentially many x, stored in a hypergraph

SCFG

SCFG

CFG is not closed under composition!

- Approximations
  - k-best
  - sampling
  - lattice

# Approximating $p_\phi(x \mid \tilde{y}_i)$



exponentially many x, stored in a hypergraph

SCFG

SCFG

CFG is not closed under composition!

- Approximations
    - k-best
    - sampling
    - lattice

variational approximation

+

lattice decoding (Dyer et al., 2008)

# The Forward System $\delta_\theta(x)$

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

loss

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

# The Forward System $\delta_\theta(x)$



$$\delta_\theta(x) = \underset{y}{\mathrm{argmax}}\, p_\theta(y \mid x)$$

- **Deterministic** Decoding
  - use **one-best** translation

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)$$

# The Forward System $\delta_\theta(x)$



$$\delta_\theta(x) = \underset{y}{\mathrm{argmax}}\, p_\theta(y \mid x)$$

- **Deterministic** Decoding
  - use **one-best** translation

the objective is not differentiable ☹

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)$$

# The Forward System $\delta_\theta(x)$



$$\delta_\theta(x) = \underset{y}{\operatorname{argmax}} \, p_\theta(y \mid x)$$

the objective is not differentiable ☹

- ## Deterministic Decoding
  - use **one-best** translation

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \mathrm{L}(\delta_\theta(x), \tilde{y}_i)$$

- ## Randomized Decoding
  - use a **distribution** of translations

$$\theta^* \;=\; \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \sum_{y} p_\theta(y \mid x) \mathrm{L}(y, \tilde{y}_i)$$

# The Forward System $\delta_\theta(x)$



$$\delta_\theta(x) = \underset{y}{\arg\max}\, p_\theta(y \mid x)$$

- ## Deterministic Decoding
  - use **one-best** translation

the objective is not differentiable

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \sum_x p_\phi(x \mid \tilde{y}_i) \boxed{\mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

- ## Randomized Decoding
  - use a **distribution** of translations

$$\theta^* \;=\; \underset{\theta}{\arg\min} \frac{1}{N} \sum_{i=1}^N \sum_x p_\phi(x \mid \tilde{y}_i) \boxed{\sum_y p_\theta(y \mid x) \mathrm{L}(y, \tilde{y}_i)}$$

# The Forward System $\delta_\theta(x)$



$$\delta_\theta(x) = \underset{y}{\operatorname{argmax}} \, p_\theta(y \mid x)$$

the objective is not differentiable 😦

- ## Deterministic Decoding
  - use **one-best** translation

$$\theta^* \;=\; \arg\min_\theta \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \boxed{\mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

- ## Randomized Decoding
  - use a **distribution** of translations

expected loss

$$\theta^* \;=\; \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^{N} \sum_{x} p_\phi(x \mid \tilde{y}_i) \boxed{\sum_{y} p_\theta(y \mid x) \mathrm{L}(y, \tilde{y}_i)}$$

30

# The Forward System $\delta_\theta(x)$



$$\delta_\theta(x) = \operatorname*{argmax}_y p_\theta(y \mid x)$$

- ## Deterministic Decoding
  - use **one-best** translation

  the objective is not differentiable

  $$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \sum_x p_\phi(x \mid \tilde{y}_i) \boxed{\mathrm{L}(\delta_\theta(x), \tilde{y}_i)}$$

- ## Randomized Decoding
  - use a **distribution** of translations

  differentiable

  expected loss

  $$\theta^* = \operatorname*{argmin}_\theta \frac{1}{N} \sum_{i=1}^N \sum_x p_\phi(x \mid \tilde{y}_i) \boxed{\sum_y p_\theta(y \mid x) \mathrm{L}(y, \tilde{y}_i)}$$

# Experiments

- Supervised Training
  - require bitext

- Unsupervised Training
  - require monolingual English

- Semi-supervised Training
  - interpolation of supervised and unsupervised

# Semi-supervised Training

# Semi-supervised Training

| Training scenario | Test BLEU |
|---|---|
| Sup, (200, 200*16) | 47.6 |

# Semi-supervised Training

| Training scenario | Test BLEU |
|---|---|
| Sup, (200, 200*16) | 47.6 |

# Semi-supervised Training

| Training scenario | Test BLEU |
|---|---|
| Sup, (200, 200*16) | 47.6 |



40K sent. pairs

# Semi-supervised Training

| Training scenario | Test BLEU |
|---|---|
| Sup, (200, 200*16) | 47.6 |

**40K sent. pairs**

# Semi-supervised Training

| Training scenario | Test BLEU |
|---|---|
| Sup, (200, 200*16) | 47.6 |



40K sent. pairs

Bilingual Data → Generative Training → Translation Models

Held-out Bilingual Data

Monolingual English → Generative Training → Language Models

Discriminative Training → Optimal Weights

551 features

Unseen Sentences → Decoding → Translation Outputs

# Semi-supervised Training

| Training scenario | Test BLEU |
|---|---|
| Sup, (200, 200*16) | 47.6 |
| +Unsup, 100*16 Eng sentences | 49.0 |
| +Unsup, 200*16 Eng sentences | 48.9 |

**40K sent. pairs**

**Adding unsupervised data helps!**



**551 features**

# Supervised vs. Unsupervised



Sup
Un-sup

Unsupervised training performs as well as (and often better than) the supervised one!

# Supervised vs. Unsupervised

**Sup**
**Un-sup**



Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

# Supervised vs. Unsupervised

**Sup** (purple)
**Un-sup** (blue)



Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

|  | Chinese | English |
|---|---|---|
| Sup | 100 | 16*100 |

# Supervised vs. Unsupervised

**Sup**
**Un-sup**



Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

|  | Chinese | English |
|---|---|---|
| Sup | 100 | 16*100 |
| Unsup | 16*100 | 16*16*100 |

# Supervised vs. Unsupervised

**Sup**
**Un-sup**



Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

|       | Chinese | English |
|-------|---------|---------|
| Sup   | 100     | 16*100  |
| Unsup | 16*100  | 16*16*100 |

But, fair comparison!

# Supervised vs. Unsupervised

**Sup**
**Un-sup**



**BLEU**

**Data Size**

Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

|  | Chinese | English |
|---|---|---|
| Sup | 100 | 16*100 |
| Unsup | 16*100 | 16*16*100 |

But, fair comparison!

- More experiments

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

loss

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

# Supervised vs. Unsupervised



**Sup**
**Un-sup**

BLEU / Data Size

Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

|  | Chinese | English |
|---|---|---|
| Sup | 100 | 16*100 |
| Unsup | 16*100 | 16*16*100 |

But, fair comparison!

- More experiments
  - different k-best size

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

$$\downarrow \text{loss}$$

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

# Supervised vs. Unsupervised

**Sup**
**Un-sup**



BLEU

Data Size

Unsupervised training performs as well as (and often better than) the supervised one!

Unsupervised uses **16** times of data as supervised. For example,

|       | Chinese | English      |
|-------|---------|--------------|
| Sup   | 100     | 16*100       |
| Unsup | 16*100  | 16*16*100    |

But, fair comparison!

- More experiments
  - different k-best size
  - different reverse model



$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

$$\updownarrow \text{loss}$$

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training

  ‣ minimum imputed risk

  ‣ **contrastive language model estimation**

- **Variational Decoding**

- **First- and Second-order Expectation Semirings**

# Language Modeling

# Language Modeling

- Language Model  $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

$$p(y) \;=\; \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

$$p(y) \;=\; \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

a set of n-grams occurred in y

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

$$p(y) \;=\; \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

Locally
normalized

a set of n-grams occurred in y

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

  $$p(y) = \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

  Locally
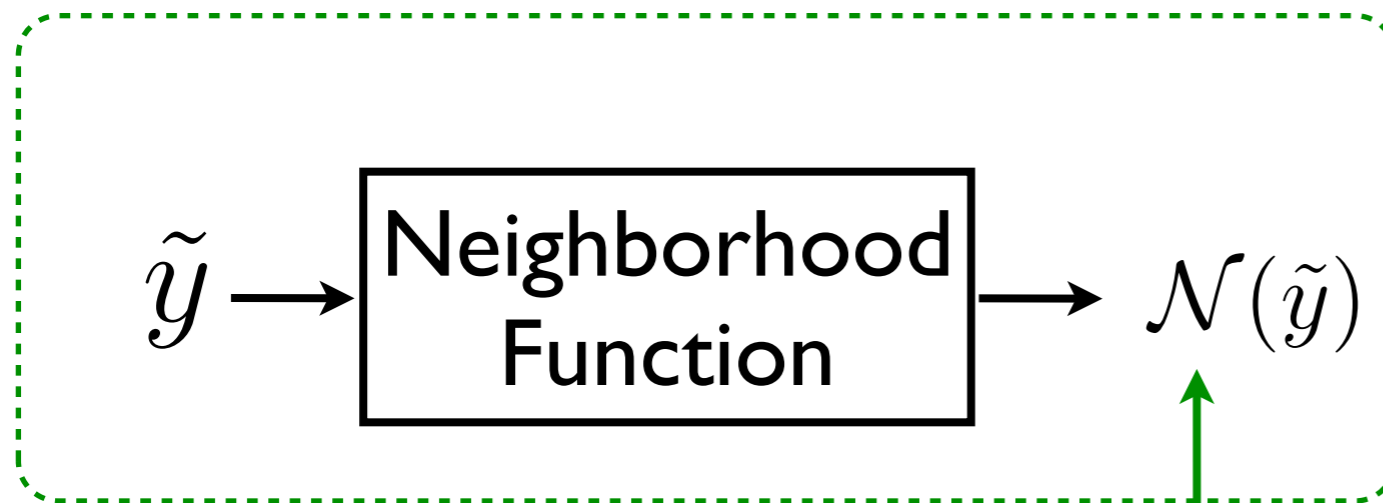  normalized

  $\longrightarrow$ a set of n-grams occurred in y

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

  $$p_\theta(y) = \frac{e^{f(y) \cdot \theta}}{Z(*)}$$

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

  $$p(y) \;=\; \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

  Locally normalized $\longrightarrow$ a set of n-grams occurred in y

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

  $$p_\theta(y) \;=\; \frac{e^{f(y) \cdot \theta}}{Z(*)}$$

  Globally normalized

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

  $$p(y) = \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

  a set of n-grams occurred in y
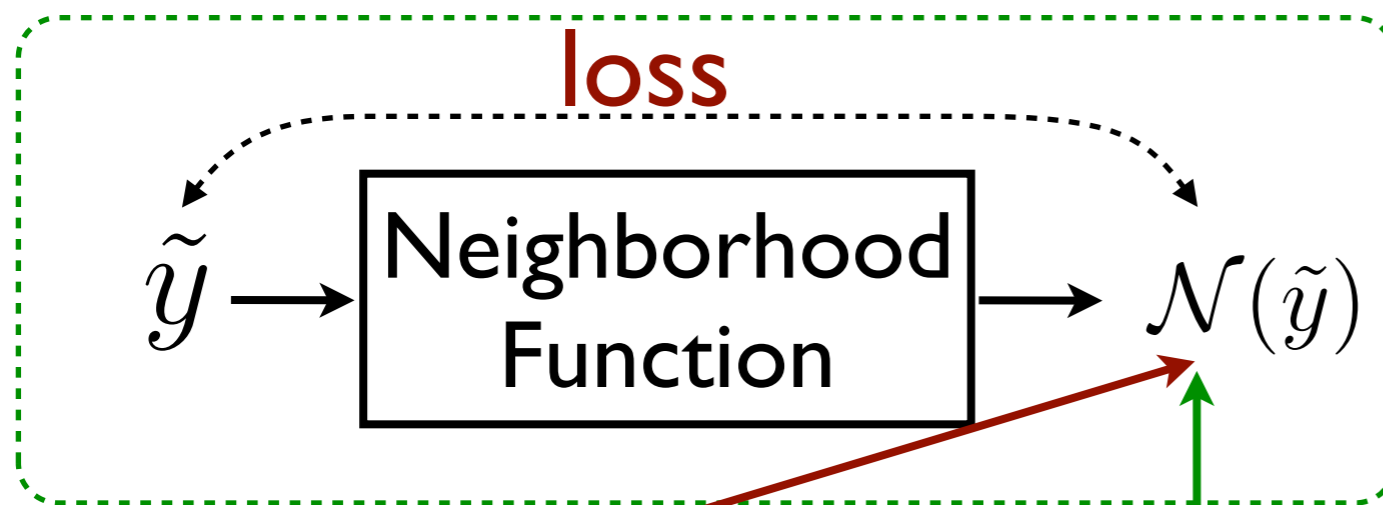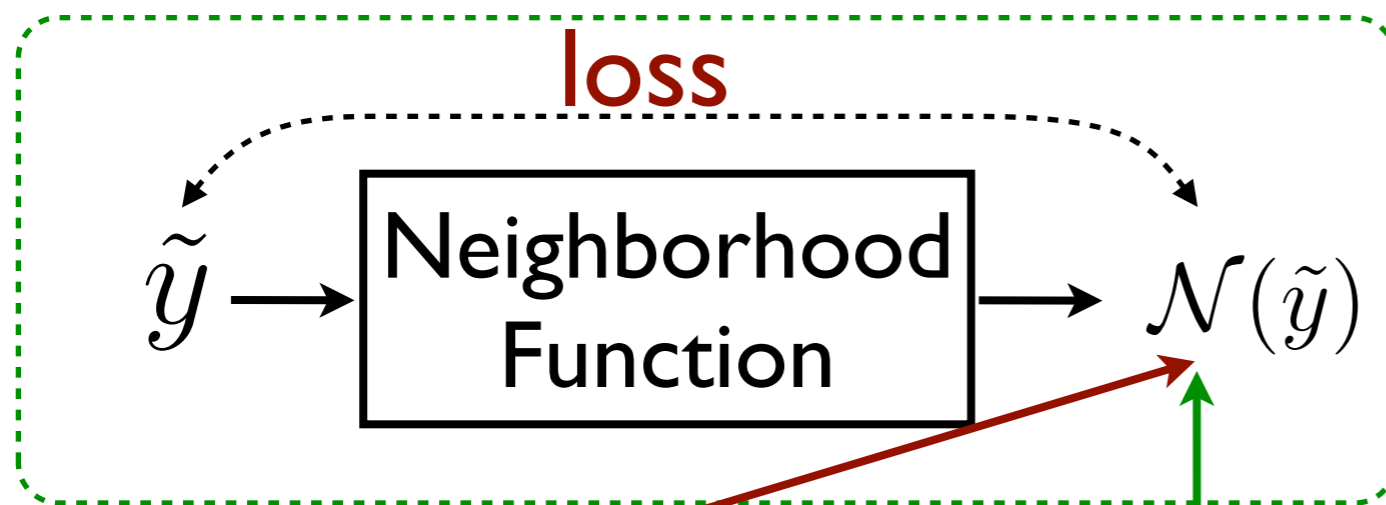
  Locally normalized

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

  $$p_\theta(y) = \frac{e^{f(y) \cdot \theta}}{Z(*)} \qquad Z(*) \overset{\text{def}}{=} \sum_{y' \in \Sigma^*} e^{f(y') \cdot \theta}$$

  Globally normalized

35

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence **y**

  - typically use an **n**-gram model

$$p(y) \;=\; \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

  a set of **n**-grams occurred in **y**

  Locally
  normalized

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \;=\; \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad\qquad Z(*) \;\overset{\text{def}}{=}\; \sum_{y' \in \Sigma^*} e^{f(y')\cdot\theta}$$

  Globally
  normalized

  All English sentences with any length!

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

  $$p(y) \;=\; \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

  a set of n-grams occurred in y

  Locally normalized

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

  $$p_\theta(y) \;=\; \frac{e^{f(y)\cdot\theta}}{\mathrm{Z}(*)} \qquad\qquad \mathrm{Z}(*) \stackrel{\mathrm{def}}{=} \sum_{y' \in \Sigma^*} e^{f(y')\cdot\theta}$$

  Globally normalized

  All English sentences with any length!

  Sampling

35

# Language Modeling

- Language Model $p_\theta(y)$

  - assign a probability to an English sentence y

  - typically use an n-gram model

  $$p(y) = \prod_{w \in W_n} p(r(w) \mid h(w))^{c_w(y)}$$

  Locally normalized $\longrightarrow$ a set of n-grams occurred in y

- Global Log-linear Model
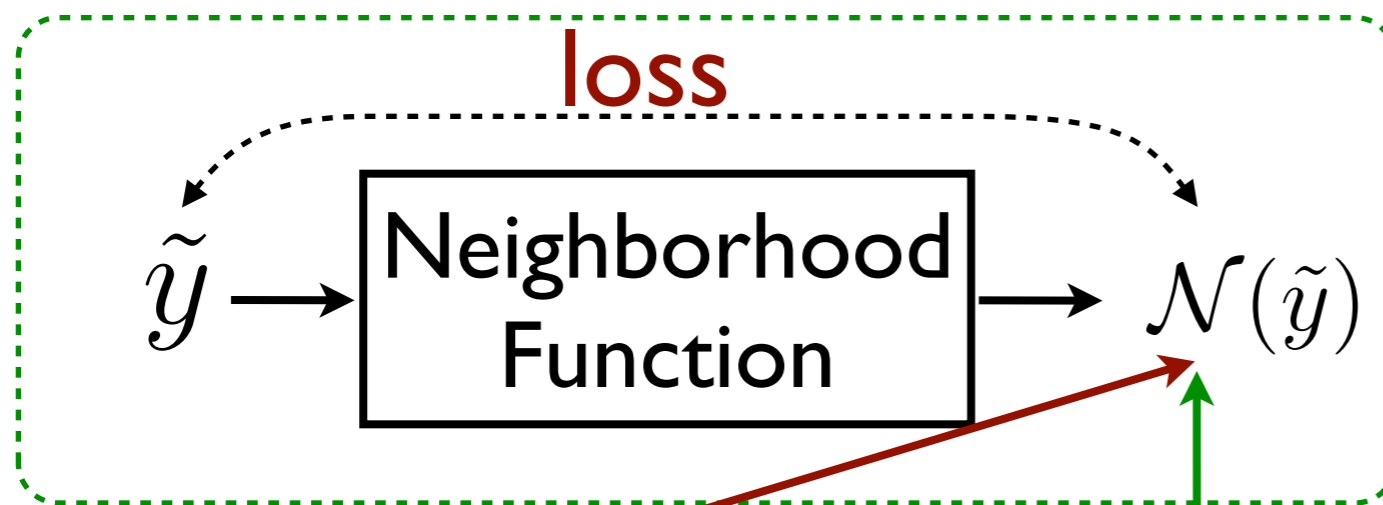  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

  $$p_\theta(y) = \frac{e^{f(y) \cdot \theta}}{Z(*)} \qquad Z(*) \overset{\text{def}}{=} \sum_{y' \in \Sigma^*} e^{f(y') \cdot \theta}$$

  Globally normalized

  All English sentences with any length!

  Sampling  slow ☹

35

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \ = \ \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad\qquad Z(*) \ \stackrel{\text{def}}{=} \ \sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}$$

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \quad = \quad \frac{e^{f(y) \cdot \theta}}{Z(*)} \qquad\qquad Z(*) \overset{\text{def}}{=} \sum_{y' \in \Sigma^*} e^{f(y') \cdot \theta}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \quad = \quad \frac{e^{f(\tilde{y}) \cdot \theta}}{Z(\tilde{y})} \quad = \quad \frac{e^{f(\tilde{y}) \cdot \theta}}{\sum_{y' \in \mathcal{N}(\tilde{y})} e^{f(y') \cdot \theta}}$$

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \;=\; \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad\qquad Z(*) \stackrel{\mathrm{def}}{=} \sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}$$
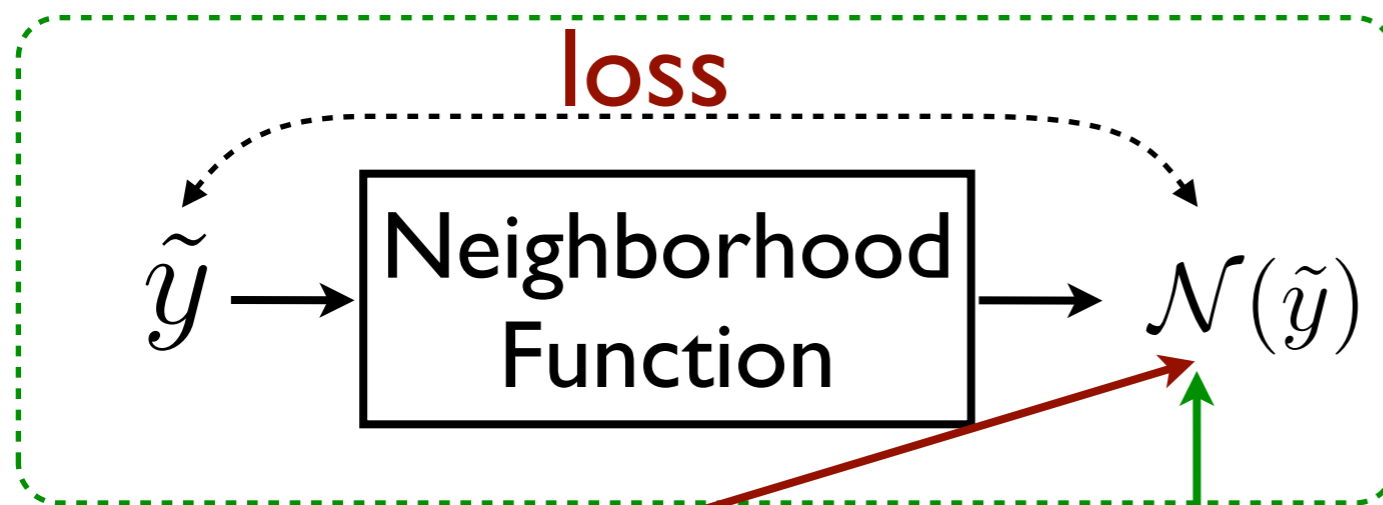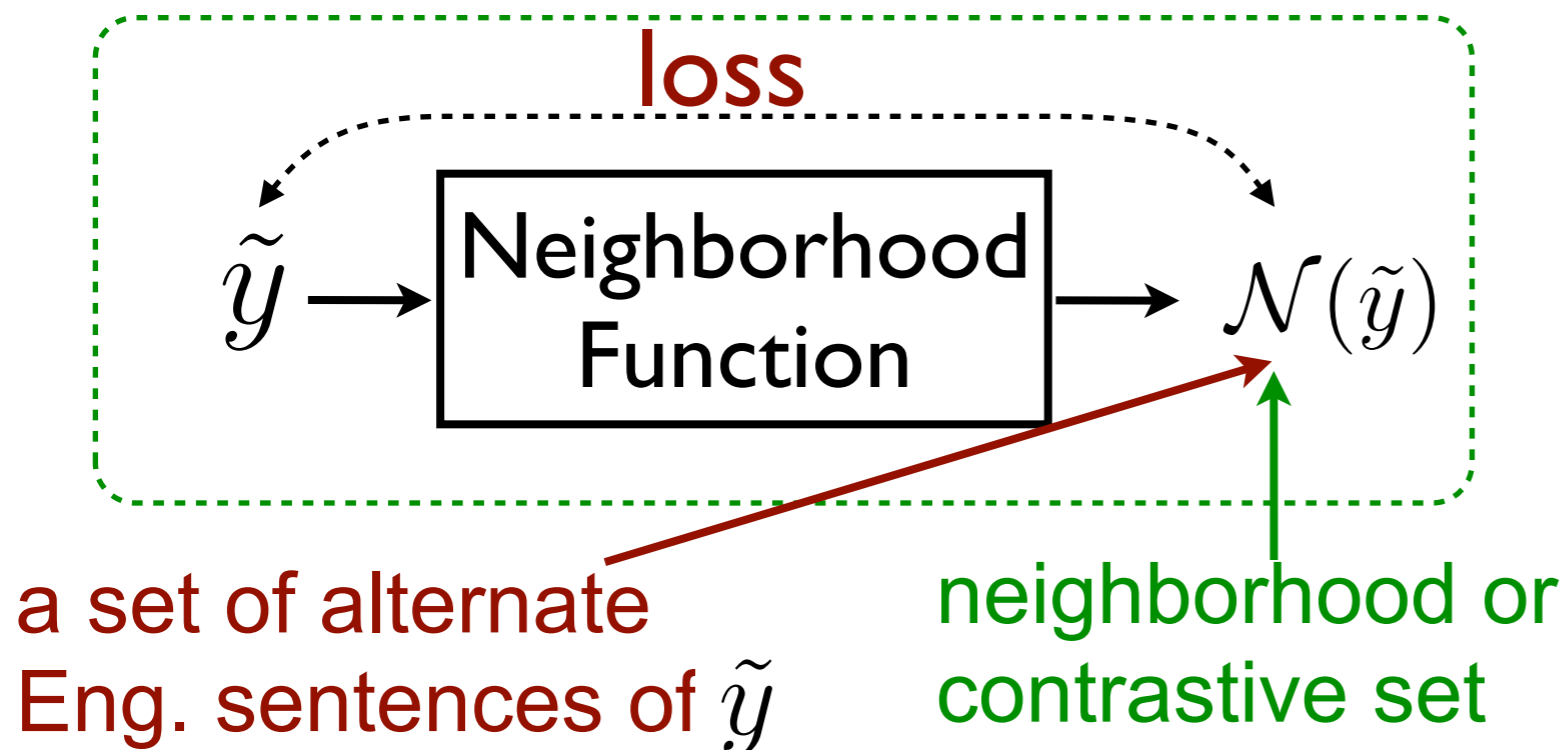
- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}$$

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \;=\; \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad\qquad Z(*) \;\overset{\text{def}}{=}\; \sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}$$

$$\tilde{y}$$

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \;=\; \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad\qquad Z(*) \stackrel{\text{def}}{=} \sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}$$

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

# Contrastive Estimation

- Global Log-linear Model
(whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \quad = \quad \frac{e^{f(y)\cdot\theta}}{\mathrm{Z}(*)} \qquad\qquad \mathrm{Z}(*) \stackrel{\mathrm{def}}{=} \sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \quad = \quad \frac{e^{f(\tilde{y})\cdot\theta}}{\mathrm{Z}(\tilde{y})} \quad = \quad \frac{e^{f(\tilde{y})\cdot\theta}}{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}$$

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

neighborhood or
contrastive set

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) = \frac{e^{f(y) \cdot \theta}}{Z(*)} \qquad Z(*) \stackrel{\text{def}}{=} \sum_{y' \in \Sigma^*} e^{f(y') \cdot \theta}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) = \frac{e^{f(\tilde{y}) \cdot \theta}}{Z(\tilde{y})} = \frac{e^{f(\tilde{y}) \cdot \theta}}{\sum_{y' \in \mathcal{N}(\tilde{y})} e^{f(y') \cdot \theta}}$$

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

a set of alternate
Eng. sentences of $\tilde{y}$

neighborhood or
contrastive set

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \;=\; \frac{e^{f(y) \cdot \theta}}{Z(*)} \qquad\qquad Z(*) \overset{\text{def}}{=} \sum_{y' \in \Sigma^*} e^{f(y') \cdot \theta}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \;=\; \frac{e^{f(\tilde{y}) \cdot \theta}}{Z(\tilde{y})} \;=\; \frac{e^{f(\tilde{y}) \cdot \theta}}{\sum_{y' \in \mathcal{N}(\tilde{y})} e^{f(y') \cdot \theta}}$$



loss

$\tilde{y} \longrightarrow$ | Neighborhood Function | $\longrightarrow \mathcal{N}(\tilde{y})$

a set of alternate Eng. sentences of $\tilde{y}$

neighborhood or contrastive set

36

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) = \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad Z(*) \stackrel{\text{def}}{=} \boxed{\sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) = \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} = \frac{e^{f(\tilde{y})\cdot\theta}}{\boxed{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}}$$

loss

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

a set of alternate
Eng. sentences of $\tilde{y}$

neighborhood or
contrastive set

# Contrastive Estimation

- Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \;=\; \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad Z(*) \stackrel{\text{def}}{=} \boxed{\sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}}$$

- Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} \;=\; \frac{e^{f(\tilde{y})\cdot\theta}}{\boxed{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}}$$

improve both speed and accuracy



loss

$\tilde{y} \longrightarrow$ Neighborhood Function $\longrightarrow \mathcal{N}(\tilde{y})$

a set of alternate Eng. sentences of $\tilde{y}$

neighborhood or contrastive set

# Contrastive Estimation

- ## Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

  $$p_\theta(y) = \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad Z(*) \stackrel{\text{def}}{=} \boxed{\sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}}$$

- ## Contrastive Estimation (CE) (Smith and Eisner, 2005)

  $$p_\theta(\tilde{y}) = \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} = \frac{e^{f(\tilde{y})\cdot\theta}}{\boxed{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}}$$

  improve both speed and accuracy

  not proposed for language modeling

  loss

  $$\tilde{y} \longrightarrow \boxed{\text{Neighborhood Function}} \longrightarrow \mathcal{N}(\tilde{y})$$

  a set of alternate Eng. sentences of $\tilde{y}$

  neighborhood or contrastive set

# Contrastive Estimation

- ## Global Log-linear Model
  (whole-sentence maximum-entropy LM) (Rosenfeld et al., 2001)

$$p_\theta(y) \ = \ \frac{e^{f(y)\cdot\theta}}{Z(*)} \qquad Z(*) \ \overset{\text{def}}{=} \ \boxed{\sum_{y'\in\Sigma^*} e^{f(y')\cdot\theta}}$$

- ## Contrastive Estimation (CE) (Smith and Eisner, 2005)

$$p_\theta(\tilde{y}) \ = \ \frac{e^{f(\tilde{y})\cdot\theta}}{Z(\tilde{y})} \ = \ \frac{e^{f(\tilde{y})\cdot\theta}}{\boxed{\sum_{y'\in\mathcal{N}(\tilde{y})} e^{f(y')\cdot\theta}}}$$

improve both speed and accuracy

loss

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

not proposed for language modeling

train to recover the original English as much as possible

a set of alternate Eng. sentences of $\tilde{y}$

neighborhood or contrastive set

36

# Contrastive Language Model Estimation

# Contrastive Language Model Estimation

- Step-1: extract a confusion grammar (CG)
  - an English-to-English SCFG

# Contrastive Language Model Estimation

- Step-1: extract a confusion grammar (CG)
  - an English-to-English SCFG

neighborhood function

# Contrastive Language Model Estimation

- Step-1: extract a confusion grammar (CG)
  - an English-to-English SCFG

  neighborhood function

- Step-2: for each English sentence, generate a contrastive set (or neighborhood) using the CG

# Contrastive Language Model Estimation

- Step-1: extract a confusion grammar (CG)
  - an English-to-English SCFG

  neighborhood function

- Step-2: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- Step-3: discriminative training

# Contrastive Language Model Estimation

- **Step-1**: extract a confusion grammar (CG)
  - an English-to-English SCFG

neighborhood function

- **Step-2**: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- **Step-3**: discriminative training

# Contrastive Language Model Estimation

- **Step-1**: extract a confusion grammar (CG)
  - an English-to-English SCFG

$$X \quad \rightarrow \quad \langle \, \text{lead to} \, , \; \text{result in} \, \rangle$$

neighborhood function

- **Step-2**: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- **Step-3**: discriminative training

# Contrastive Language Model Estimation

- **Step-1**: extract a confusion grammar (CG)
  - an English-to-English SCFG

  $$X \quad \rightarrow \quad \langle \text{lead to}, \ \text{result in} \rangle$$

  neighborhood function

  paraphrase

- **Step-2**: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- **Step-3**: discriminative training

# Contrastive Language Model Estimation

- **Step-1**: extract a confusion grammar (CG)
  - an English-to-English SCFG

$$X \;\rightarrow\; \langle \text{lead to}\,,\; \text{result in}\,\rangle$$

$$X \;\rightarrow\; \langle\, X_0 \text{ at beijing}\,,\; \text{beijing 's } X_0 \,\rangle$$

neighborhood function

paraphrase

- **Step-2**: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- **Step-3**: discriminative training

# Contrastive Language Model Estimation

- **Step-1**: extract a confusion grammar (CG)
  - an English-to-English SCFG

  *neighborhood function*

  **paraphrase**

  $$X \rightarrow \langle \text{lead to} , \text{ result in} \rangle$$

  $$X \rightarrow \langle X_0 \text{ at beijing} , \text{ beijing 's } X_0 \rangle$$

  $$X \rightarrow \langle X_0 \text{ of } X_1 , X_0 \text{ of the } X_1 \rangle$$

- **Step-2**: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- **Step-3**: discriminative training

# Contrastive Language Model Estimation

- **Step-1**: extract a confusion grammar (CG)
  - an English-to-English SCFG

$$X \rightarrow \langle \text{lead to} , \text{ result in} \rangle \qquad \text{paraphrase}$$

$$X \rightarrow \langle X_0 \text{ at beijing} , \text{ beijing 's } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ of } X_1 , X_0 \text{ of the } X_1 \rangle \qquad \text{insertion}$$

neighborhood function

- **Step-2**: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- **Step-3**: discriminative training

# Contrastive Language Model Estimation

- Step-1: extract a confusion grammar (CG)
  - an English-to-English SCFG

  neighborhood function

  $$X \rightarrow \langle \text{lead to} , \text{ result in} \rangle \qquad \text{paraphrase}$$

  $$X \rightarrow \langle X_0 \text{ at beijing} , \text{beijing 's } X_0 \rangle$$

  $$X \rightarrow \langle X_0 \text{ of } X_1 , X_0 \text{ of the } X_1 \rangle \qquad \text{insertion}$$

  $$X \rightarrow \langle X_0 \text{ 's } X_1 , X_1 \text{ of } X_0 \rangle$$

- Step-2: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- Step-3: discriminative training

# Contrastive Language Model Estimation

- Step-1: extract a confusion grammar (CG)
  - an English-to-English SCFG

  neighborhood function

$$X \rightarrow \langle \text{lead to} , \text{ result in} \rangle \qquad \text{paraphrase}$$

$$X \rightarrow \langle X_0 \text{ at beijing} , \text{ beijing 's } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ of } X_1 , X_0 \text{ of the } X_1 \rangle \qquad \text{insertion}$$

$$X \rightarrow \langle X_0 \text{ 's } X_1 , X_1 \text{ of } X_0 \rangle \qquad \text{re-ordering}$$

- Step-2: for each English sentence, generate a contrastive set (or neighborhood) using the CG

- Step-3: discriminative training

# Step-1: Extracting a Confusion Grammar (CG)

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|:---:|:---:|

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|---|---|

$$X \rightarrow \langle \text{ mao}, \text{ a cat} \rangle$$
$$X \rightarrow \langle \text{ mao}, \text{ the cat} \rangle$$

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|---|---|
| $X \rightarrow \langle \text{mao},\ \text{a cat} \rangle$ | $X \rightarrow \langle \text{a cat},\ \text{the cat} \rangle$ |
| $X \rightarrow \langle \text{mao},\ \text{the cat} \rangle$ | $X \rightarrow \langle \text{the cat},\ \text{a cat} \rangle$ |

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|---|---|

$$X \rightarrow \langle \text{ mao, a cat} \rangle$$

$$X \rightarrow \langle \text{ mao, the cat} \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ on } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle \text{ a cat, the cat} \rangle$$

$$X \rightarrow \langle \text{ the cat, a cat} \rangle$$

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|---|---|

$$X \rightarrow \langle \text{ mao, a cat} \rangle$$

$$X \rightarrow \langle \text{ mao, the cat} \rangle$$

$$X \rightarrow \langle \text{ a cat, the cat} \rangle$$

$$X \rightarrow \langle \text{ the cat, a cat} \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ on } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ of } X_1, X_1 \text{ on } X_0 \rangle$$

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|---|---|

$$X \rightarrow \langle \text{ mao}, \text{ a cat} \rangle \qquad X \rightarrow \langle \text{ a cat}, \text{ the cat} \rangle$$

$$X \rightarrow \langle \text{ mao}, \text{ the cat} \rangle \qquad X \rightarrow \langle \text{ the cat}, \text{ a cat} \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ on } X_1 \rangle \qquad X \rightarrow \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle \qquad X \rightarrow \langle X_0 \text{ of } X_1, X_1 \text{ on } X_0 \rangle$$

CG captures the confusion an MT system will
have when translating an input.

# Step-1: Extracting a Confusion Grammar (CG)

- Deriving a CG from a bilingual grammar
  - use Chinese side as pivots

| Bilingual Rule | Confusion Rule |
|---|---|

$$X \rightarrow \langle \text{ mao, a cat} \rangle$$

$$X \rightarrow \langle \text{ mao, the cat} \rangle$$

$$X \rightarrow \langle \text{ a cat, the cat} \rangle$$

$$X \rightarrow \langle \text{ the cat, a cat} \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ on } X_1 \rangle$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$$

$$X \rightarrow \langle X_0 \text{ of } X_1, X_1 \text{ on } X_0 \rangle$$

CG captures the confusion an MT system will have when translating an input.

Our neighborhood function is **learned** and **MT-specific**.

# Step-2: Generating Contrastive Sets

a cat on the mat

# Step-2: Generating Contrastive Sets

a cat on the mat

**CG**

# Step-2: Generating Contrastive Sets

a cat on the mat

CG

S | 0,5

$S \to \langle \mathbf{X_0}, \mathbf{X_0} \rangle$

X | 0,5

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ on } X_0 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_0 \text{ 's } X_1 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_0 X_1 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$

X | 0,2

X | 3,5

$X \to \langle \mathbf{a\ cat}, \mathbf{the\ cat} \rangle$

$X \to \langle \mathbf{the\ mat}, \mathbf{the\ mat} \rangle$

$a_0$  $cat_1$         $on_2$         $the_3$  $mat_4$

39

# Step-2: Generating Contrastive Sets

a cat on the mat

CG

**Contrastive set:**

the cat the mat

the cat 's the mat

the mat on the cat

the mat of the cat

| S | 0,5 |

$S \to \langle \mathbf{X_0}, \mathbf{X_0} \rangle$

| X | 0,5 |

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ on } X_0 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_0 \text{ 's } X_1 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_0 X_1 \rangle$

| X | 0,2 |

| X | 3,5 |

$X \to \langle \mathbf{a \ cat}, \mathbf{the \ cat} \rangle$

$X \to \langle \mathbf{the \ mat}, \mathbf{the \ mat} \rangle$

$a_0 \quad cat_1$

$on_2$

$the_3 \ mat_4$

# Step-2: Generating Contrastive Sets

a cat on the mat

CG

Contrastive set:

the cat the mat
the cat 's the mat
the mat on the cat
the mat of the cat

| S | 0,5 |

$S \to \langle X_0, X_0 \rangle$

| X | 0,5 |

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ on } X_0 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_0 \text{ 's } X_1 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_0 \ X_1 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$

| X | 0,2 |

| X | 3,5 |

$X \to \langle \mathbf{a\ cat}, \mathbf{the\ cat} \rangle$

$X \to \langle \mathbf{the\ mat}, \mathbf{the\ mat} \rangle$

$a_0$ cat$_1$      on$_2$      the$_3$ mat$_4$

$S \to \langle X_0, X_0 \rangle$

$X \to \langle X_0 \text{ on } X_1, X_1 \text{ of } X_0 \rangle$

$X \to \langle \mathbf{a\ cat}, \mathbf{the\ cat} \rangle$    $X \to \langle \mathbf{the\ mat}, \mathbf{the\ mat} \rangle$

**a cat**      **on**      **the mat**

Translating "dianzi shang de mao"?

# Step-3: Discriminative Training

- Training Objective

$$\theta^* \;=\; \arg\min_{\theta} \sum_i \sum_{y \in \mathcal{N}(\tilde{y}_i)} \mathrm{L}(y, \tilde{y}_i) p_\theta(y \mid \tilde{y}_i)$$

# <span style="color:magenta">Step-3</span>: <span style="color:blue">Discriminative Training</span>

- Training Objective

$$\theta^* = \arg\min_\theta \sum_i \sum_{y \in \mathcal{N}(\tilde{y}_i)} \mathrm{L}(y, \tilde{y}_i) p_\theta(y \mid \tilde{y}_i)$$

<span style="color:darkred">contrastive set</span>

# Step-3: Discriminative Training

- ## Training Objective

expected loss

$$\theta^* = \arg\min_{\theta} \sum_i \boxed{\sum_{y \in \mathcal{N}(\tilde{y}_i)} \mathrm{L}(y, \tilde{y}_i) p_{\theta}(y \mid \tilde{y}_i)}$$

contrastive set

# Step-3: Discriminative Training

- Training Objective

expected loss

$$\theta^* \;=\; \arg\min_{\theta} \sum_i \boxed{\sum_{y \in \mathcal{N}(\tilde{y}_i)} \mathrm{L}(y, \tilde{y}_i) p_{\theta}(y \mid \tilde{y}_i)}$$

contrastive set

CE maximizes the
conditional likelihood

# Step-3: Discriminative Training

- Training Objective

expected loss

$$\theta^* = \arg\min_\theta \sum_i \boxed{\sum_{y \in \mathcal{N}(\tilde{y}_i)} \mathrm{L}(y, \tilde{y}_i) p_\theta(y \mid \tilde{y}_i)}$$

contrastive set

CE maximizes the
conditional likelihood

- Iterative Training

  - Step-2: for each English sentence, generate a contrastive set (or neighborhood) using the CG

  - Step-3: discriminative training

# Applying the Contrastive Model

# Applying the Contrastive Model

- We can use the contrastive model as a regular language model

# Applying the Contrastive Model

- We can use the contrastive model as a regular language model

- We can incorporate the contrastive model into an end-to-end MT system as a feature

# Applying the Contrastive Model

- We can use the contrastive model as a regular language model

- We can incorporate the contrastive model into an end-to-end MT system as a feature

- We may also use the contrastive model to generate paraphrase sentences
  (if the loss function measures semantic similarity)

  - the rules in CG are symmetric

# Test on Synthesized Hypergraphs of English Data

# Test on Synthesized Hypergraphs of English Data

**Monolingual**

**English**

**Confusion grammar**

# Test on Synthesized Hypergraphs of English Data

**Monolingual English**

**Confusion grammar**

**Training**

# Test on Synthesized Hypergraphs of English Data

**Monolingual English**

**Confusion grammar**

**Training**

**Contrastive LM**

# Test on Synthesized Hypergraphs of English Data

**Monolingual English**

**Confusion grammar**

**Training**

**Contrastive LM**

**English Sentence**

**Parsing**

**Hypergraph (Neighborhood)**

# Test on Synthesized Hypergraphs of English Data

**Monolingual English**

**Confusion grammar**

**Training** → **Contrastive LM**

**English Sentence** → **Parsing** → **Hypergraph (Neighborhood)** → **Rank** → **One-best English**

# Test on Synthesized Hypergraphs of English Data

**Monolingual English**

**Confusion grammar**

**Training** → **Contrastive LM**

**English Sentence** → **Parsing** → **Hypergraph (Neighborhood)** → **Rank** → **One-best English**

BLEU Score?

# Test on Synthesized Hypergraphs of English Data

# Results on Synthesized Hypergraphs

# Results on Synthesized Hypergraphs

# Results on Synthesized Hypergraphs

# Results on Synthesized Hypergraphs



BLEU

26
24
22
20
18
16
14
12
10

BLM    +WP    +RuleBigram

**Features**

baseline LM (5-gram)

word penalty

# Results on Synthesized Hypergraphs

# Results on Synthesized Hypergraphs



BLEU

26
24
22
20
18
16
14
12
10

BLM    +WP    +RuleBigram

**Features**

baseline LM (5-gram)

word penalty

# Results on Synthesized Hypergraphs



baseline LM (5-gram)

word penalty

- **Target side** of a confusion rule

  "*on the $X_1$ issue of $X_2$*"

# Results on Synthesized Hypergraphs

**BLEU**

26
24
22
20
18
16
14
12
10

BLM    +WP    +RuleBigram

**Features**

baseline LM (5-gram)

word penalty

- Target side of a confusion rule

  *"on the $X_1$ issue of $X_2$"*

- Rule bigram features

  *"on the"*    *"the X"*    *"X issue"*

  *"issue of"*    *"of X"*

# Results on Synthesized Hypergraphs



BLEU (y-axis: 10, 12, 14, 16, 18, 20, 22, 24, 26)

Features: BLM, +WP, +RuleBigram

The contrastive LM better **recovers** the original English than a regular n-gram LM.

baseline LM (5-gram)

word penalty

- Target side of a confusion rule

  *"on the $X_1$ issue of $X_2$"*

- Rule bigram features

  *"on the"*    *"the X"*    *"X issue"*

  *"issue of"*    *"of X"*

# Results on Synthesized Hypergraphs



The contrastive LM better **recovers** the original English than a regular n-gram LM.

All the features look at only the target sides of confusion rules

- Target side of a confusion rule

  *"on the $X_1$ issue of $X_2$"*

- Rule bigram features

  *"on the"*      *"the X"*      *"X issue"*

  *"issue of"*      *"of X"*

# Results on MT Test Set

# Results on MT Test Set

BLEU

50
49
48
47
46
45

Baseline      +CLM

**Features**

# Results on MT Test Set



BLEU

50
49
48
47
46
45

Baseline      +CLM

**Features**

Bilingual Data → Generative Training → Translation Models

Monolingual English → Generative Training → Language Models

Held-out Bilingual Data

Discriminative Training → Optimal Weights

Unseen Sentences → Decoding → Translation Outputs

# Results on MT Test Set



BLEU

50
49
48
47
46
45

Baseline        +CLM

**Features**

Bilingual Data → Generative Training → Translation Models

Monolingual English → Generative Training → Language Models

Held-out Bilingual Data

→ Discriminative Training → Optimal Weights

Unseen Sentences → Decoding → Translation Outputs

Add CLM as a feature

# Results on MT Test Set



The contrastive LM helps to improve MT performance.

Add CLM as a feature

# Adding Features on the CG itself

- On English Set

- On MT Set

# Adding Features on the CG itself

- ## On English Set



- ## On MT Set

# Adding Features on the CG itself

- ## On English Set



- ## On MT Set

# Adding Features on the CG itself

- ## On English Set



- ## On MT Set

# Adding Features on the CG itself

- ## On English Set



- ## On MT Set

# Adding Features on the CG itself

- ## On English Set



- ## On MT Set



glue rules or regular confusion rules?

$$S \rightarrow \langle S_0\ X_1, S_0\ X_1 \rangle$$
$$S \rightarrow \langle X_0, X_0 \rangle$$

# Adding Features on the CG itself

- ## On English Set



one big feature

glue rules or regular confusion rules?

$$S \;\to\; \langle\, S_0\, X_1\,,\, S_0\, X_1\,\rangle$$
$$S \;\to\; \langle\, X_0\,,\, X_0\,\rangle$$

- ## On MT Set

# Adding Features on the CG itself

- **On English Set**



one big feature

glue rules or regular confusion rules?

$$S \rightarrow \langle S_0 \, X_1 \,, S_0 \, X_1 \rangle$$

$$S \rightarrow \langle X_0 \,, X_0 \rangle$$

- **On MT Set**

# Adding Features on the CG itself

- ## On English Set



Paraphrasing model

one big feature

glue rules or regular confusion rules?

$$S \rightarrow \langle S_0\, X_1\,, S_0\, X_1 \rangle$$
$$S \rightarrow \langle X_0\,, X_0 \rangle$$

- ## On MT Set

# Adding Features on the CG itself

- ## On English Set



**Paraphrasing model**

**one big feature**

- ## On MT Set



glue rules or regular confusion rules?

$$S \rightarrow \langle S_0\, X_1\,, S_0\, X_1 \rangle$$

$$S \rightarrow \langle X_0\,, X_0 \rangle$$

# Adding Features on the CG itself

- ## On English Set



Paraphrasing model

one big feature

- ## On MT Set



glue rules or regular confusion rules?

$$S \rightarrow \langle S_0\ X_1, S_0\ X_1 \rangle$$
$$S \rightarrow \langle X_0, X_0 \rangle$$

# Adding Features on the CG itself

- ## On English Set



Paraphrasing model

one big feature

glue rules or regular confusion rules?

$$S \rightarrow \langle S_0\, X_1 , S_0\, X_1 \rangle$$
$$S \rightarrow \langle X_0 , X_0 \rangle$$

- ## On MT Set

# Summary for Discriminative Training

- **Supervised**: Minimum Empirical Risk

require bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\quad \text{loss} \quad} y$$

- **Unsupervised**: Minimum Imputed Risk

require monolingual English

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

loss

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

- **Unsupervised**: Contrastive LM Estimation

require monolingual English

loss

$$\tilde{y} \longrightarrow \boxed{\text{Neighborhood Function}} \longrightarrow \mathcal{N}(\tilde{y})$$

# Summary for Discriminative Training

- **Supervised**: Minimum Empirical Risk

require bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

- **Unsupervised**: Minimum Imputed Risk

require monolingual English

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$
$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \quad \text{loss}$$

- **Unsupervised**: Contrastive LM Estimation

require monolingual English

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood} \\ \text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y}) \quad \text{loss}$$

# Summary for Discriminative Training

- ## Supervised Training

require
bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

- ## Unsupervised: Minimum Imputed Risk

require
monolingual
English

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$
$$\downarrow \qquad\qquad \updownarrow \text{loss}$$
$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

- ## Unsupervised: Contrastive LM Estimation

require
monolingual
English

loss

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

# Summary for Discriminative Training

- ## Supervised Training

require bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}}\dashleftarrow\dashrightarrow y$$

- ## Unsupervised: Minimum Imputed Risk

require monolingual English

require a reverse model

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

$$\text{loss}$$

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

- ## Unsupervised: Contrastive LM Estimation

require monolingual English

loss

$$\tilde{y} \longrightarrow \boxed{\text{Neighborhood Function}} \longrightarrow \mathcal{N}(\tilde{y})$$

# Summary for Discriminative Training

- ## Supervised Training

require
bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

- ## Unsupervised: Minimum Imputed Risk

require
monolingual
English

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

loss

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

require a reverse model

can have both TM and
LM features

- ## Unsupervised: Contrastive LM Estimation

require
monolingual
English

loss

$$\tilde{y} \longrightarrow \boxed{\text{Neighborhood Function}} \longrightarrow \mathcal{N}(\tilde{y})$$

# Summary for Discriminative Training

- ## Supervised Training

require
bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

- ## Unsupervised: Minimum Imputed Risk

require
monolingual
English

require a reverse model

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

$$\updownarrow \text{loss}$$

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

can have both TM and
LM features

- ## Unsupervised: Contrastive LM Estimation

require
monolingual
English

loss

$$\tilde{y} \longrightarrow \boxed{\begin{array}{c}\text{Neighborhood}\\\text{Function}\end{array}} \longrightarrow \mathcal{N}(\tilde{y})$$

can have LM features
only

# Summary for Discriminative Training

- ## Supervised Training

require
bitext

$$x \longrightarrow \boxed{\delta_\theta} \longrightarrow \delta_\theta(x) \xleftarrow{\text{loss}} y$$

- ## Unsupervised: Minimum Imputed Risk

require
monolingual
English

$$x \longleftarrow \boxed{p_\phi} \longleftarrow \tilde{y}_i$$

loss

$$\boxed{\delta_\theta} \longrightarrow \delta_\theta(x)$$

require a reverse model

can have both TM and
LM features

loss

- ## Unsupervised: Contrastive LM Estimation

require
monolingual
English

loss

$$\tilde{y} \longrightarrow \boxed{\text{Neighborhood Function}} \longrightarrow \mathcal{N}(\tilde{y})$$

can have LM features
only

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training

  ‣ minimum imputed risk

  ‣ contrastive language model estimation

- **Variational Decoding**

- **First- and Second-order Expectation Semirings**

| **decoding**<br>(e.g., mbr) | **training**<br>(e.g., mert) |
|---|---|
| **atomic inference operations**<br>(e.g., finding one-best, k-best or expectation,<br>inference can be *exact* or *approximate*) ||

# Variational Decoding

# Variational Decoding

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

# Variational Decoding

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

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

# Variational Decoding

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

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

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

# Variational Decoding

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

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

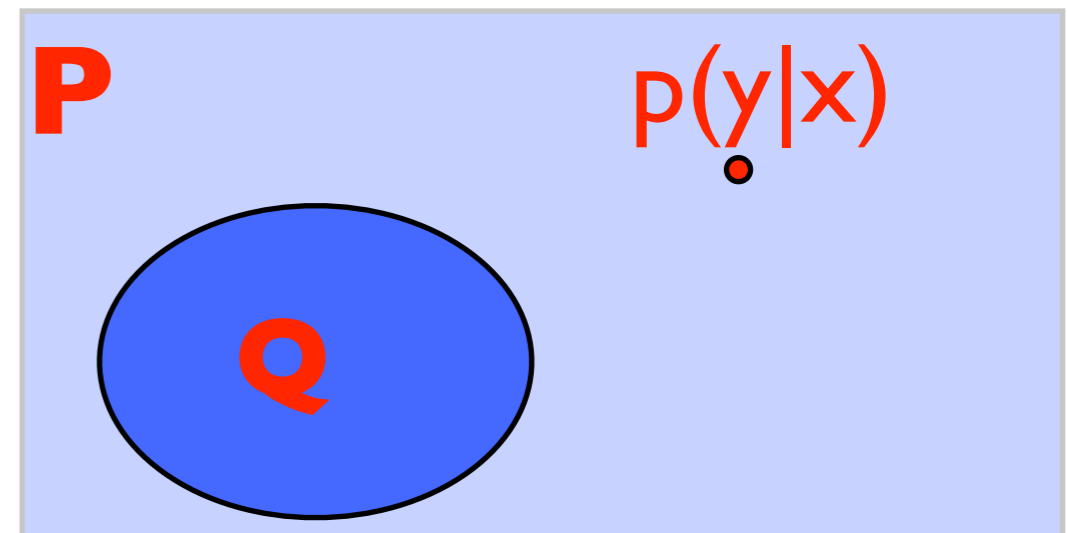- Then, we will use $q^*$ as a surrogate for $p$ in inference

# Variational Decoding

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

**intractable MAP decoding**

$$y^* \;=\; \arg\max_{y} \; p(y \mid x)$$

- Instead, we derive a simpler distribution q*

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

# Variational Decoding

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

**intractable MAP decoding**

$$y^* \;=\; \arg\max_y \; p(y \mid x) = \; \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

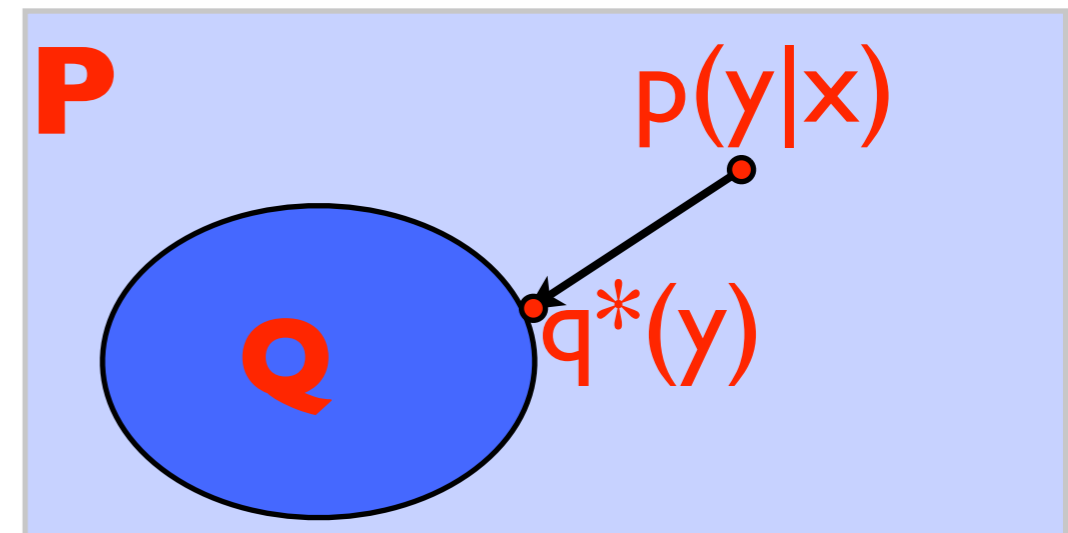- Then, we will use q* as a surrogate for p in inference

# Variational Decoding

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

**intractable MAP decoding** (Sima'an 1996)

$$y^* \;=\; \arg\max_y \; p(y \mid x) = \; \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

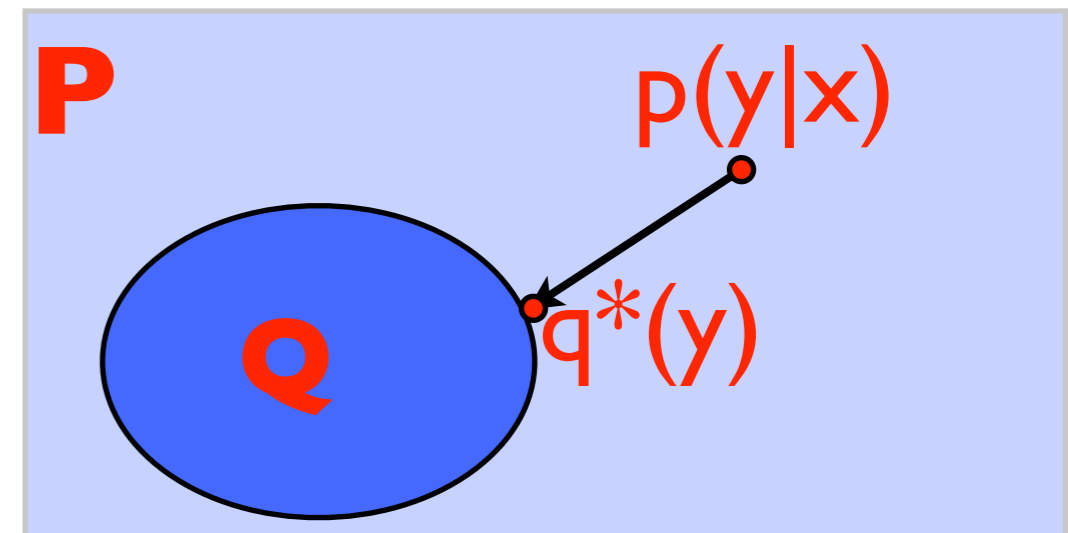- Then, we will use q* as a surrogate for p in inference

# Variational Decoding

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

**intractable MAP decoding**  (Sima'an 1996)

$$y^* \;=\; \arg\max_y \; p(y \mid x) = \; \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

**tractable estimation**

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

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

# Variational Decoding

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

**intractable MAP decoding**   (Sima'an 1996)

$$y^* \;=\; \arg\max_y \; p(y \mid x) = \; \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

**tractable estimation**

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

**P**

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

# Variational Decoding

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

**intractable MAP decoding**    (Sima'an 1996)

$$y^* \quad = \quad \arg\max_y \; p(y \mid x) = \quad \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

**tractable estimation**

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

P    p(y|x)

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

# Variational Decoding

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

**intractable MAP decoding** (Sima'an 1996)

$$y^* \;=\; \arg\max_y \; p(y \mid x) = \; \arg\max_y \; \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

**tractable estimation**

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

P   p(y|x)

Q

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

# Variational Decoding

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

**intractable MAP decoding**   (Sima'an 1996)

$$y^* \;=\; \arg\max_y\; p(y \mid x) = \; \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

**tractable estimation**

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

P

p(y|x)

Q

q*(y)

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

# Variational Decoding

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

**intractable MAP decoding**    (Sima'an 1996)

$$y^* \;=\; \arg\max_y \; p(y \mid x) = \; \arg\max_y \sum_{d \in \mathrm{D}(x,y)} p(d \mid x)$$

- Instead, we derive a simpler distribution q*

**tractable estimation**

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



P

p(y|x)

Q

q*(y)

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

**tractable decoding**

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

# 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 for the foreign sentence

# Variational Decoding for MT: an Overview

Sentence-specific decoding

Three steps:

**1** Generate a hypergraph for the foreign sentence

Foreign
sentence x

# Variational Decoding for MT: an Overview

Sentence-specific decoding

MAP decoding under P is intractable

Three steps:

**1** Generate a hypergraph for the foreign sentence



$$p(d \mid x)$$

$$p(y \mid x) = \sum_{d \in D(x,y)} p(d|x)$$

Foreign sentence x $\Rightarrow$ SMT $\Rightarrow$

$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4 \mid$ the $\cdots$ cat

$X \mid 0,4 \mid$ a $\cdots$ mat

$X \rightarrow \langle X_0$ de $X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_1$ of

$X \mid 0,2$

$X \mid 3,4 \mid$ a $\cdots$ cat

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi$_0$   shang$_1$     de$_2$    mao$_3$

# 1

Generate a hypergraph

$$p(d \mid x)$$

S 0,4

S→⟨X₀, X₀⟩

X 0,4 the ⋯ cat

X 0,4 a ⋯ mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨...₁, X₁ of X₀⟩

X 0,2 the ⋯ mat

X 3,4 a ⋯ cat

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀ shang₁ de₂ mao₃

54

**1**

S | 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X | 0,4 | the ··· cat

X | 0,4 | a ··· mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨... X₁, X₁ of X₀⟩

p(d | x)

X | 0,2 | the ··· mat

X | 3,4 | a ··· cat

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁          de₂     mao₃

Generate a hypergraph

54

**1**



Generate a hypergraph

**2**

**1**

S | 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X | 0,4 | the ⋯ cat          X | 0,4 | a ⋯ mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→                          ...₁, X₁ of X₀⟩

p(d | x)

X | 0,2 | the ⋯ mat          X | 3,4 | a ⋯ cat

X→⟨dianzi shang, the mat⟩     X→⟨mao, a cat⟩

dianzi₀  shang₁      de₂    mao₃

**Generate a hypergraph**

---

**2**

S | 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X | 0,4 | the ⋯ cat          X | 0,4 | a ⋯ mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→                          ...₁, X₁ of X₀⟩

p(d | x)

X | 0,2 | the ⋯ mat          X | 3,4 | a ⋯ cat

X→⟨dianzi shang, the mat⟩     X→⟨mao, a cat⟩

dianzi₀  shang₁      de₂    mao₃

54

**1** Generate a hypergraph

**2** Estimate a model from the hypergraph by minimizing KL

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

$q^*(y \mid x)$
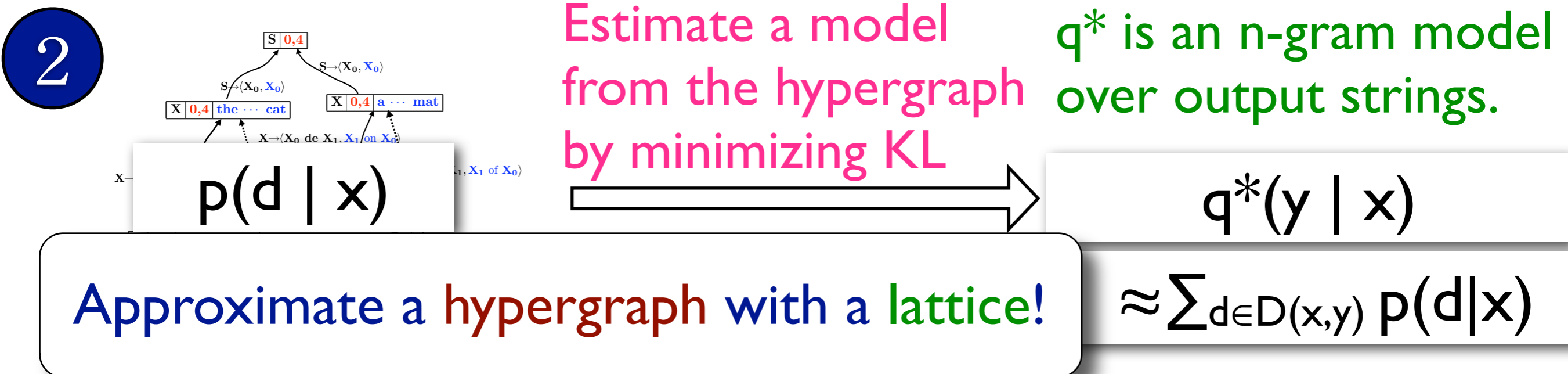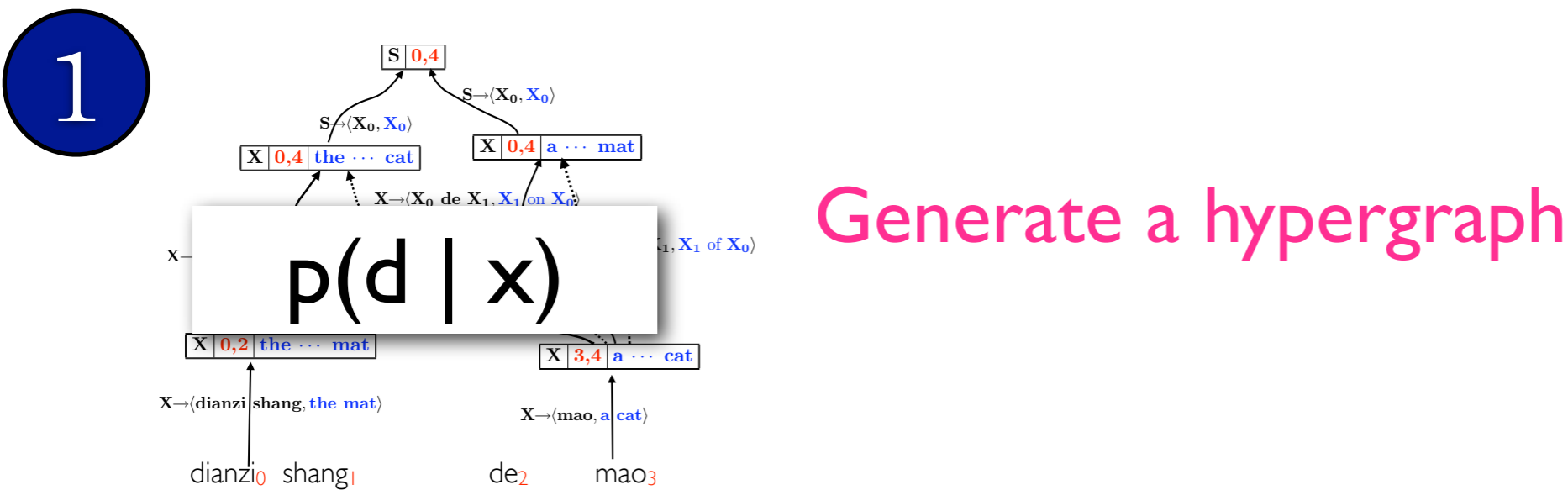
**1** Generate a hypergraph

**2** Estimate a model from the hypergraph by minimizing KL

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

$$q^*(y \mid x) \approx \sum_{d \in D(x,y)} p(d \mid x)$$

**①**

Generate a hypergraph

**②**

Estimate a model from the hypergraph by minimizing KL

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

$$q^*(y \mid x)$$

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

**③**

**1** Generate a hypergraph

**2** Estimate a model from the hypergraph by minimizing KL

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

$$q^*(y \mid x)$$

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

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

54

**1.** Generate a hypergraph

**2.** Estimate a model from the hypergraph by minimizing KL

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

$$q^*(y \mid x)$$

$$\approx \sum_{d \in D(x,y)} p(d \mid x)$$

Approximate a hypergraph with a lattice!

**3.** Decode using q* on the hypergraph

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training

  ‣ minimum imputed risk

  ‣ contrastive language model estimation

- Variational Decoding

- **First- and Second-order Expectation Semirings**

| **decoding** (e.g., mbr) | **training** (e.g., mert) |
|---|---|
| **atomic inference operations** (e.g., finding one-best, k-best or expectation, inference can be *exact* or *approximate*) ||

A semiring framework to compute all of these

**Probabilistic Hypergraph**

- "Decoding" quantities:
  - Viterbi
  - K-best
  - Counting
  - ......

- First-order expectations:
  - expectation
    - entropy
    - expected loss
    - cross-entropy
    - KL divergence
    - feature expectations
  - first-order gradient of $Z$

- Second-order expectations:
  - expectation over product
    - interaction between features
  - Hessian matrix of $Z$
    - second-order gradient descent
  - gradient of expectation
    - gradient of expected loss or entropy

53

**A semiring framework to compute all of these**

Probabilistic Hypergraph

- "Decoding" quantities:
  - Viterbi
  - K-best
  - Counting
  - ......

- First-order expectations:
  - expectation
    - entropy
    - expected loss
  - cross-entropy
  - KL divergence
  - feature expectations
  - first-order gradient of Z

- Second-order expectations:
  - expectation over product
    - interaction between features
  - Hessian matrix of Z
    - second-order gradient descent
  - gradient of expectation
    - gradient of expected loss or entropy

# A semiring framework to compute all of these

**S** | **0,4**

S→⟨**X₀**, **X₀**⟩

S→⟨**X₀**, **X₀**⟩

**X** | **0,4** | the ⋯ cat

**X** | **0,4** | a ⋯ mat

X→⟨**X** ...

**X** | **0,2** | the ⋯ mat

X→⟨**dianzi shang**, **the mat**⟩

dianzi₀ shang₁

## Probabilistic Hyp...

- "Decoding" quantities:

## Recipe to compute a quantity:

- First-orde...
  - expectat...
    - entropy
    - **expected loss**
    - cross-entropy
    - KL divergence
    - feature expectations
  - first-order gradient of Z
    - interaction between features
  - Hessian matrix of Z
    - second-order gradient descent
  - gradient of expectation
    - gradient of expected loss or entropy

53

A semiring framework to compute all of these

S 0,4

$S \to \langle X_0, X_0 \rangle$

$S \to \langle X_0, X_0 \rangle$

X 0,4 the ⋯ cat

X 0,4 a ⋯ mat

$X \to \langle X$

X 0,2 the ⋯ mat

$X \to \langle \text{dianzi shang}, \text{the mat} \rangle$

dianzi$_0$ shang$_1$

Probabilistic Hyp

• "Decoding" quantities:

Recipe to compute a quantity:

• Choose a semiring

- interaction between features

- Hessian matrix of Z

- second-order gradient descent

- gradient of expectation

- gradient of expected loss or entropy

• First-orde

- expectat

- entropy

- expected loss

- cross-entropy

- KL divergence

- feature expectations

- first-order gradient of Z

A semiring framework to compute all of these



**S** | 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

**X** | 0,4 | the ⋯ cat

**X** | 0,4 | a ⋯ mat

X→⟨X

Probabilistic
Hype

**X** | 0,2 | the ⋯ mat

X→⟨dianzi shang, the mat⟩

dianzi₀  shang₁

- "Decoding" quantities:

Recipe to compute a quantity:

- Choose a semiring

- Specific a semiring weight for each hyperedge

- First-orde
  - **expectat**
    - **entropy**
    - **expected loss**
    - **cross-entropy**
    - **KL divergence**
    - **feature expectations**
  - **first-order gradient of** Z

- interaction between features
  - **Hessian matrix of** Z
  - second-order gradient descent
  - **gradient of expectation**
    - **gradient of expected loss or entropy**

53

A semiring framework to compute all of these

S 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X 0,4 the ⋯ cat

X 0,4 a ⋯ mat

X→⟨X...

Probabilistic
Hype...

X 0,2 the ⋯ mat

X→⟨dianzi shang, the mat⟩

dianzi₀ shang₁

- "Decoding" quantities:

## Recipe to compute a quantity:

- Choose a semiring

- Specific a semiring weight for each hyperedge

- Run the inside algorithm

- First-orde...
  - expectat...
    - entropy
    - expected loss
    - cross-entropy
    - KL divergence
    - feature expectations
  - first-order gradient of Z

- interaction between features
  - Hessian matrix of Z
  - second-order gradient descent
  - gradient of expectation
    - gradient of expected loss or entropy

53

# Applications of Expectation Semirings: a Summary

| Quantity | $k_e$ | $k_{\mathbf{root}}$ | Final |
|---|---|---|---|
| **Expectation** | $\langle p_e, p_e r_e \rangle$ | $\langle Z, \bar{r} \rangle$ | $\bar{r}/Z$ |
| Entropy | $r_e \stackrel{\text{def}}{=} \log p_e$, so $k_e = \langle p_e, p_e \log p_e \rangle$ | $\langle Z, \bar{r} \rangle$ | $\log Z - \bar{r}/Z$ |
| Cross-entropy | $\langle q_e \rangle$ <br> $r_e \stackrel{\text{def}}{=} \log q_e$, so $k_e = \langle p_e, p_e \log q_e \rangle$ | $\langle Z_q \rangle$ <br> $\langle Z_p, \bar{r} \rangle$ | $\log Z_q - \bar{r}/Z_p$ |
| Bayes risk | $r_e \stackrel{\text{def}}{=} \mathrm{L}_e$, so $k_e = \langle p_e, p_e \mathrm{L}_e \rangle$ | $\langle Z, \bar{r} \rangle$ | $\bar{r}/Z$ |
| **First-order gradient** | $\langle p_e, \nabla p_e \rangle$ | $\langle Z, \nabla Z \rangle$ | $\nabla Z$ |
| **Covariance matrix** | $\langle p_e, p_e r_e, p_e s_e, p_e r_e s_e \rangle$ | $\langle Z, \bar{r}, \bar{s}, \bar{t} \rangle$ | $\dfrac{\bar{t}}{Z} - \dfrac{\bar{r}\,\bar{s}^{\mathbf{T}}}{Z^2}$ |
| **Hessian matrix** | $\langle p_e, \nabla p_e, \nabla p_e, \nabla^2 p_e \rangle$ | $\langle Z, \nabla Z, \nabla Z, \nabla^2 Z \rangle$ | $\nabla^2 Z$ |
| **Gradient of expectation** | $\langle p_e, p_e r_e, \nabla p_e, (\nabla p_e) r_e + p_e(\nabla r_e) \rangle$ | $\langle Z, \bar{r}, \nabla Z, \nabla \bar{r} \rangle$ | $\dfrac{Z \nabla \bar{r} - \bar{r} \nabla Z}{Z^2}$ |
| Gradient of entropy | $\langle p_e, p_e \log p_e, \nabla p_e, (1 + \log p_e)\nabla p_e \rangle$ | $\langle Z, \bar{r}, \nabla Z, \nabla \bar{r} \rangle$ | $\dfrac{\nabla Z}{Z} - \dfrac{Z \nabla \bar{r} - \bar{r} \nabla Z}{Z^2}$ |
| Gradient of risk | $\langle p_e, p_e \mathrm{L}_e, \nabla p_e, \mathrm{L}_e \nabla p_e \rangle$ | $\langle Z, \bar{r}, \nabla Z, \nabla \bar{r} \rangle$ | $\dfrac{Z \nabla \bar{r} - \bar{r} \nabla Z}{Z^2}$ |

# Inference, Training and Decoding on Hypergraphs

- Unsupervised Discriminative Training

  ▸ minimum imputed risk (In Preparation)

  ▸ contrastive language model estimation (In Preparation)

- Variational Decoding
  (Li et al., ACL 2009)

- First- and Second-order Expectation Semirings
  (Li and Eisner, EMNLP 2009)

# My Other MT Research

- **Training methods (supervised)**

  - Discriminative forest reranking with Perceptron
    (Li and Khudanpur, GALE book chapter 2009)

  - Discriminative n-gram language models
    (Li and Khudanpur, AMTA 2008)

- **Algorithms**

  - Oracle extraction from hypergraphs
    (Li and Khudanpur, NAACL 2009)

  - Efficient intersection between n-gram LM and CFG
    (Li and Khudanpur, ACL SSST 2008)

- **Others**
  - System combination (Smith et al., GALE book chapter 2009)
  - Unsupervised translation induction for Chinese abbreviations (Li and Yarowsky, ACL 2008)

# Research other than MT

- **Information extraction**

  - Relation extraction between formal and informal phrases (Li and Yarowsky, EMNLP 2008)

- **Spoken dialog management**

  - Optimal dialog in consumer-rating systems using a POMDP (Li et al., SIGDial 2008)

# Joshua project

- An open-source parsing-based MT toolkit (Li et al. 2009)
  - support Hiero (Chiang, 2007) and SAMT (Venugopal et al., 2007)
- Team members
  - **Zhifei Li**, Chris Callison-Burch, Chris Dyer, Sanjeev Khudanpur, Wren Thornton, Jonathan Weese, Juri Ganitkevitch, Lane Schwartz, and Omar Zaidan



Only rely on word-aligner and SRI LM!

All the methods presented have been implemented in Joshua!

# Thank you!
## XieXie!
## 谢谢!

# Decoding over a hypergraph

# Decoding over a hypergraph



## Given a hypergraph of possible translations
(generated for a given foreign sentence by already-trained model)

# Decoding over a hypergraph



Given a hypergraph of possible translations
(generated for a given foreign sentence by already-trained model)

Pick a single translation to output
(why not just pick the tree with the highest weight?)

# Spurious Ambiguity

- Statistical models in MT exhibit **spurious ambiguity**

  - Many **different derivations** (e.g., trees or segmentations) generate the **same translation string**

- Tree-based MT systems

  - derivation tree ambiguity

- Regular phrase-based MT systems

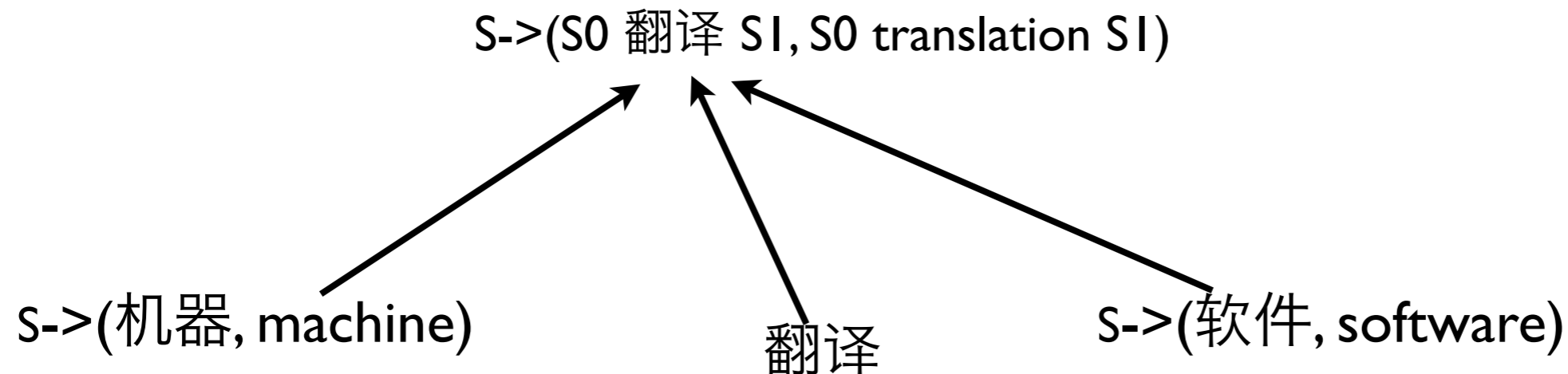  - phrase segmentation ambiguity

# Spurious Ambiguity in Derivation Trees

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件  machine translation software

jiqi fanyi yuanjian

# Spurious Ambiguity in Derivation Trees

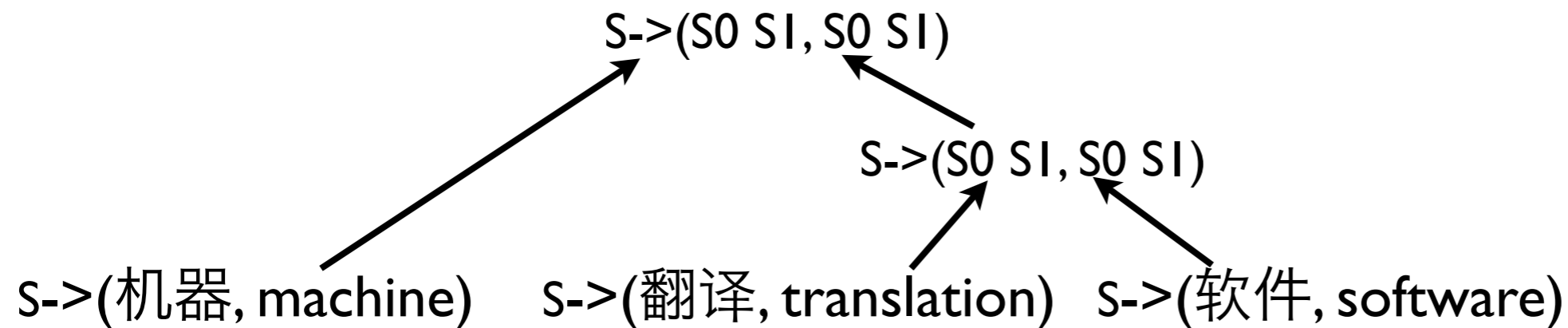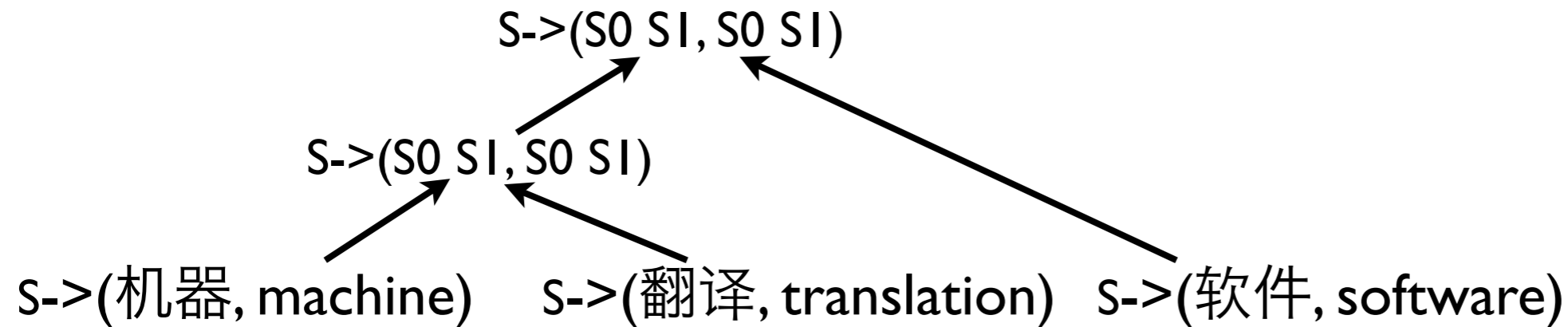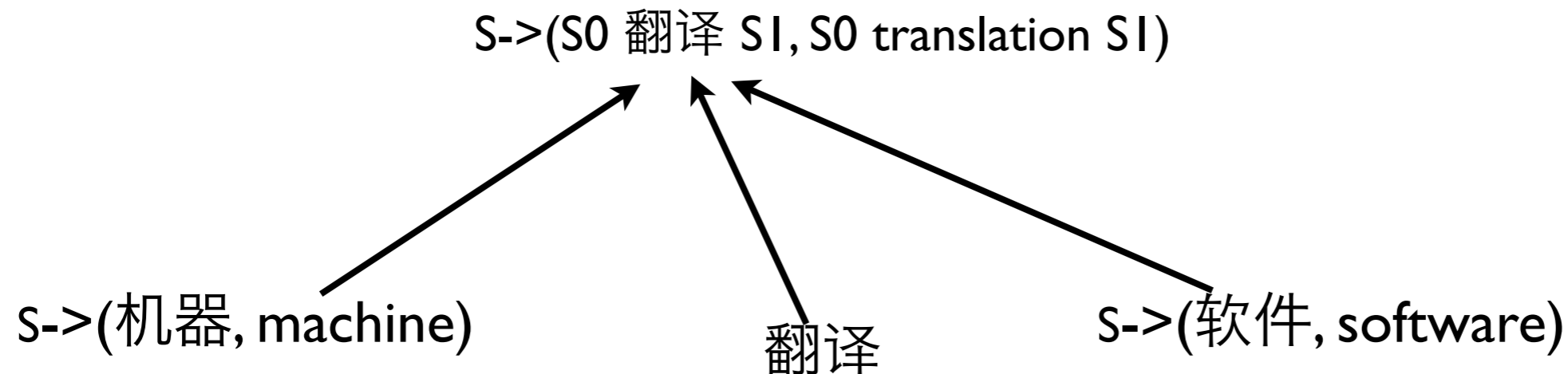机器 翻译 软件   machine translation software
jiqi fanyi yuanjian

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)    S->(翻译, translation)   S->(软件, software)

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件　machine translation software

jiqi fanyi yuanjian

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)　　S->(翻译, translation)　S->(软件, software)

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)　　S->(翻译, translation)　S->(软件, software)

# Spurious Ambiguity in Derivation Trees

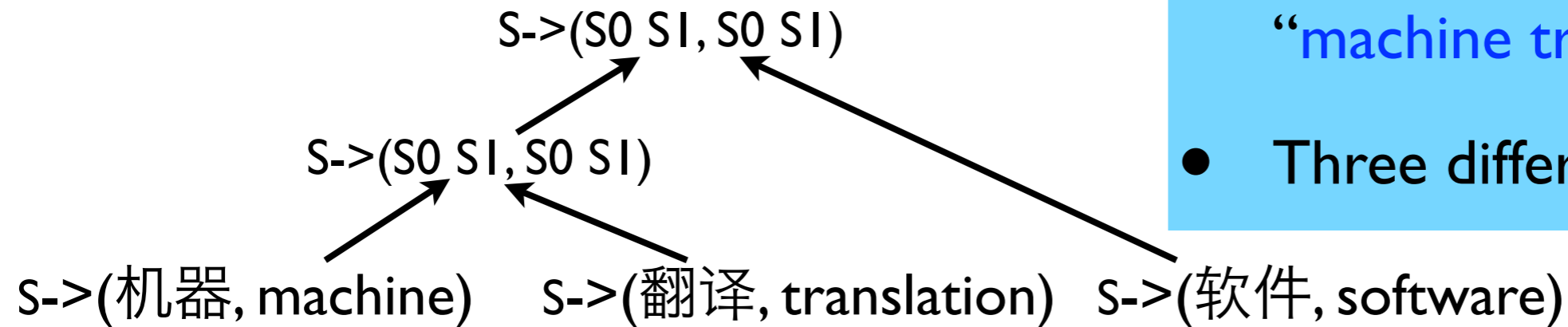机器 翻译 软件   machine translation software

jiqi fanyi yuanjian

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)   S->(翻译, translation)   S->(软件, software)

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)   S->(翻译, translation)   S->(软件, software)

S->(S0 翻译 S1, S0 translation S1)

S->(机器, machine)   翻译   S->(软件, software)
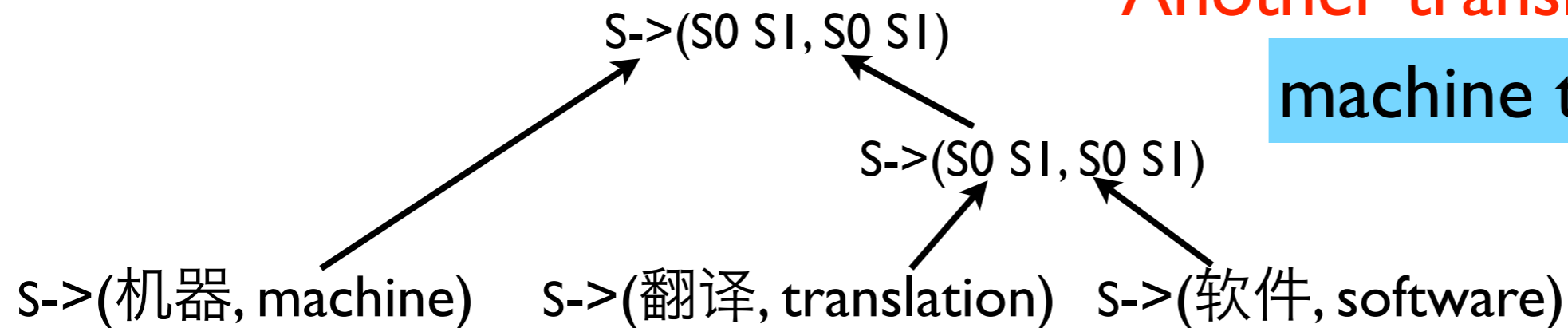
# Spurious Ambiguity in Derivation Trees

机器 翻译 软件    machine translation software

jiqi fanyi yuanjian

- Same output:
  "machine translation software"

- Three different derivation trees

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)    S->(翻译, translation)    S->(软件, software)

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)    S->(翻译, translation)    S->(软件, software)

S->(S0 翻译 S1, S0 translation S1)

S->(机器, machine)    翻译    S->(软件, software)

# Spurious Ambiguity in Derivation Trees

机器 翻译 软件 machine translation software

jiqi fanyi yuanjian

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)    S->(翻译, translation)   S->(软件, software)

- Same output:
  "machine translation software"

- Three different derivation trees

S->(S0 S1, S0 S1)

S->(S0 S1, S0 S1)

S->(机器, machine)    S->(翻译, translation)   S->(软件, software)

Another translation:

machine transfer software

S->(S0 翻译 S1, S0 translation S1)

S->(机器, machine)         翻译         S->(软件, software)

65

# MAP, Viterbi and N-best Approximations

# MAP, Viterbi and N-best Approximations

- Exact MAP decoding

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

# MAP, Viterbi and N-best Approximations

- **Exact MAP decoding**

$$y^* = \arg\max_{y \in \text{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

# MAP, Viterbi and N-best Approximations

- **Exact MAP decoding**

$$y^* = \arg \max_{y \in \mathrm{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

- **Viterbi approximation**

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

# MAP, Viterbi and N-best Approximations

- Exact MAP decoding

$$y^* = \arg \max_{y \in \mathrm{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

- Viterbi approximation

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

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

# MAP, Viterbi and N-best Approximations

- Exact MAP decoding

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

NP-hard (Sima'an 1996)

$$= \quad \arg \max_{y \in \mathrm{Trans}(x)} \boxed{\sum}_{d \in \mathrm{D}(x,y)} p(y,d|x)$$

- Viterbi approximation

$$y^* \quad = \quad \arg \max_{y \in \mathrm{Trans}(x)} \boxed{\max}_{d \in \mathrm{D}(x,y)} p(y,d|x)$$

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

# MAP, Viterbi and N-best Approximations

- Exact MAP decoding

$$y^* = \arg\max_{y \in \mathrm{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

- Viterbi approximation

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

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

# MAP, Viterbi and N-best Approximations

- **Exact MAP decoding**

$$y^* = \arg \max_{y \in \mathrm{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

- **Viterbi approximation**

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

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

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

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

66

# MAP, Viterbi and N-best Approximations

- **Exact MAP decoding**

$$y^* = \arg\max_{y \in \mathrm{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

- **Viterbi approximation**

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

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

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

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

66

# MAP, Viterbi and N-best Approximations

- **Exact MAP decoding**

$$y^* = \arg\max_{y \in \mathrm{Trans}(x)} p(y|x) \qquad \text{NP-hard (Sima'an 1996)}$$

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

- **Viterbi approximation**

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

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

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

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

# MAP vs. Approximations

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

# MAP vs. Approximations

| translation string | MAP | Viterbi | 4-best crunching | derivation | probability |
|---|---|---|---|---|---|
| | | | | | 0.16 |
| red translation | 0.28 | 0.16 | 0.16 | | 0.14 |
| 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 |

- Exact MAP decoding under spurious ambiguity is intractable on HG

67

# MAP vs. Approximations

| translation string | MAP | Viterbi | 4-best crunching | derivation | probability |
|---|---|---|---|---|---|
| | | | | ──── (red) | 0.16 |
| red translation | 0.28 | **0.16** | 0.16 | ──── (blue) | 0.14 |
| blue translation | 0.28 | 0.14 | **0.28** | ──── (blue) | 0.14 |
| | | | | ──── (green) | 0.13 |
| green translation | **0.44** | 0.13 | 0.13 | ──── (red) | 0.12 |
| | | | | ──── (green) | 0.11 |
| | | | | ──── (green) | 0.10 |
| | | | | ──── (green) | 0.10 |

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

# MAP vs. Approximations

| translation string | MAP | Viterbi | 4-best crunching | derivation | probability |
|---|---|---|---|---|---|
| | | | | ───── (red) | 0.16 |
| red translation | 0.28 | **0.16** | 0.16 | ───── (blue) | 0.14 |
| blue translation | 0.28 | 0.14 | **0.28** | ───── (blue) | 0.14 |
| green translation | **0.44** | 0.13 | 0.13 | ───── (green) | 0.13 |
| | | | | ───── (red) | 0.12 |
| | | | | ───── (green) | 0.11 |
| | | | | ───── (green) | 0.10 |
| | | | | ───── (green) | 0.10 |

- Exact MAP decoding under spurious ambiguity is intractable on HG

- 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

**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 for the foreign sentence

# Variational Decoding for MT: an Overview

Sentence-specific decoding

Three steps:

**1** Generate a hypergraph for the foreign sentence

Foreign
sentence x

# Variational Decoding for MT: an Overview

Sentence-specific decoding

Three steps:

**1** Generate a hypergraph for the foreign sentence

Foreign sentence x $\Rightarrow$ SMT $\Rightarrow$

# Variational Decoding for MT: an Overview

Sentence-specific decoding

MAP decoding under P is intractable

Three steps:

**1** Generate a hypergraph for the foreign sentence

Foreign sentence x $\Rightarrow$ SMT $\Rightarrow$

$\mathbf{S} \mid \mathbf{0,4}$

$\mathbf{S} \rightarrow \langle \mathbf{X_0}, \mathbf{X_0} \rangle$

$\mathbf{S} \rightarrow \langle \mathbf{X_0}, \mathbf{X_0} \rangle$

$\mathbf{X} \mid \mathbf{0,4} \mid \mathbf{the} \cdots \mathbf{cat}$

$\mathbf{X} \mid \mathbf{0,4} \mid \mathbf{a} \cdots \mathbf{mat}$

$p(y, d \mid x)$

$p(y \mid x)$

$\mathbf{X} \rightarrow \langle \mathbf{X_0} \text{ de } \mathbf{X_1}, \mathbf{X_0} \mathbf{X_1} \rangle$

$\mathbf{X} \rightarrow \langle \mathbf{X_0} \text{ de } \mathbf{X_1}, \mathbf{X_1} \text{ of}$

$\mathbf{X} \mid \mathbf{0,2} \mid \mathbf{the} \cdots \mathbf{mat}$

$\mathbf{X} \mid \mathbf{3,4} \mid \mathbf{a} \cdots \mathbf{cat}$

$\mathbf{X} \rightarrow \langle \mathbf{dianzi\ shang}, \mathbf{the\ mat} \rangle$

$\mathbf{X} \rightarrow \langle \mathbf{mao}, \mathbf{a\ cat} \rangle$

$\text{dianzi}_0 \quad \text{shang}_1 \qquad \text{de}_2 \qquad \text{mao}_3$

**1**

Generate a hypergraph

$$p(d \mid x)$$

S 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X 0,4 the ⋯ cat

X 0,4 a ⋯ mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→⟨..., X₁ of X₀⟩

X 0,2 the ⋯ mat

X 3,4 a ⋯ cat

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁          de₂     mao₃

69

**1**

Generate a hypergraph

p(d | x)

**1**

$$S \mid 0,4$$

$$S \rightarrow \langle X_0, X_0 \rangle$$

$$S \rightarrow \langle X_0, X_0 \rangle$$

$$X \mid 0,4 \mid \text{the} \cdots \text{cat}$$

$$X \mid 0,4 \mid a \cdots \text{mat}$$

$$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$$

$$X \rightarrow \langle \ldots, X_1, X_1 \text{ of } X_0 \rangle$$

p(d | x)

$$X \mid 0,2 \mid \text{the} \cdots \text{mat}$$

$$X \mid 3,4 \mid a \cdots \text{cat}$$

$$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$$

$$X \rightarrow \langle \text{mao}, a \text{ cat} \rangle$$

dianzi₀  shang₁          de₂      mao₃

## Generate a hypergraph

**2**

**1**

S | 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X | 0,4 | the ⋯ cat

X | 0,4 | a ⋯ mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→ ⟨..., X₁, X₁ of X₀⟩

$$p(d \mid x)$$

X | 0,2 | the ⋯ mat

X | 3,4 | a ⋯ cat

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁          de₂     mao₃

Generate a hypergraph

---

**2**

S | 0,4

S→⟨X₀, X₀⟩

S→⟨X₀, X₀⟩

X | 0,4 | the ⋯ cat

X | 0,4 | a ⋯ mat

X→⟨X₀ de X₁, X₁ on X₀⟩

X→ ⟨..., X₁, X₁ of X₀⟩

$$p(d \mid x)$$

X | 0,2 | the ⋯ mat

X | 3,4 | a ⋯ cat

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁          de₂     mao₃

69

**1** Generate a hypergraph

**2** Estimate a model from the hypergraph

$p(d \mid x) \Rightarrow q^*(y \mid x)$

69

**1**

Generate a hypergraph

**2**

Estimate a model from the hypergraph

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

**1** Generate a hypergraph

**2** Estimate a model from the hypergraph

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

$$q^*(y \mid x)$$

$$\approx \sum_{d \in D(x,y)} p(d \mid x)$$

**1** Generate a hypergraph

$p(d \mid x)$

**2** Estimate a model from the hypergraph

$p(d \mid x)$

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

$$q^*(y \mid x)$$

$$\approx \sum_{d \in D(x,y)} p(d \mid x)$$

**3**

**1** Generate a hypergraph

**2** Estimate a model from the hypergraph

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

$$q^*(y \mid x)$$

$$\approx \sum_{d \in D(x,y)} p(d \mid x)$$

**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 <span style="color:red">p</span>, but it is intractable

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

# Variational Inference

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

$$y^* \quad = \quad \arg\max_{y} p(y|x)$$

- Instead, we derive a simpler distribution q*

# Variational Inference

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

$$y^* \quad = \quad \arg \max_{y} p(y|x)$$

- Instead, we derive a simpler distribution q*

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

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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

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

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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

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

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

# 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*

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

**P**

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

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

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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

**P**        P

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

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

# 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*

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



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

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

# 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*

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



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

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

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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



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

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

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

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

→ a family of distributions

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

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

$$= \quad \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p \log \frac{p}{q}$$

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

$$
\begin{aligned}
q^* \quad = \quad & \arg\min_{q \in Q} \mathrm{KL}(p||q) \longrightarrow \text{a family of distributions} \\[2ex]
= \quad & \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p\log\frac{p}{q} \\[2ex]
= \quad & \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (p\log p - p\log q)
\end{aligned}
$$

# Variational Approximation

- q*: an approximation having minimum distance to p

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

$$= \quad \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p\log\frac{p}{q}$$

$$= \quad \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (p\log p - p\log q)$$

constant

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

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

constant

# Variational Approximation

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

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

constant

- Three questions

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

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

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p \log \frac{p}{q}$$

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (\underbrace{p \log p}_{\text{constant}} - p \log q)$$

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

- Three questions

  - how to parameterize **q**?

# Variational Approximation

- q*: an approximation having minimum distance to p

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

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p\log\frac{p}{q}$$

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (p\log p - p\log q)$$

constant

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

- Three questions

  - how to parameterize q?

  - how to estimate q*?

# Variational Approximation

- q*: an approximation having minimum distance to p

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

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p\log\frac{p}{q}$$

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (\underbrace{p\log p} - p\log q) \longrightarrow \text{constant}$$

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

- Three questions

  - how to parameterize q?

  - how to estimate q*?

  - how to use q* for decoding?

# Variational Approximation

- q*: an approximation having minimum distance to p

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

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p \log \frac{p}{q}$$

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (p \log p - p \log q)$$

constant

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

- Three questions

  - how to parameterize q?

    an n-gram model

  - how to estimate q*?

  - how to use q* for decoding?

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

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

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} p\log\frac{p}{q}$$

$$= \arg\min_{q \in Q} \sum_{y \in \mathrm{Trans}(x)} (p\log p - p\log q) \quad \text{constant}$$

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

- Three questions

  - how to parameterize **q**?

    an n-gram model

  - how to estimate **q\***?

    compute expected n-gram counts and normalize

  - how to use **q\*** for decoding?

# Variational Approximation

- **q\***: an approximation having minimum distance to **p**

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

$$= \arg\min_{q \in Q} \sum_{y \in \text{Trans}(x)} p\log\frac{p}{q}$$

$$= \arg\min_{q \in Q} \sum_{y \in \text{Trans}(x)} (p\log p - p\log q) \quad \longrightarrow \text{constant}$$

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

- Three questions

  - how to parameterize **q**?

    an n-gram model

  - how to estimate **q\***?

    compute expected n-gram counts and normalize

  - how to use **q\*** for decoding?

    score the hypergraph with the n-gram model

72

# KL divergences under different variational models

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

# KL divergences under different variational models

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

| **Measure** bits/word | $\overline{\overline{\mathrm{H}}}(p)$ | $\overline{\mathrm{KL}}(p\|\cdot)$ | | | |
|---|---|---|---|---|---|
| | | $q_1^*$ | $q_2^*$ | $q_3^*$ | $q_4^*$ |
| 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

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

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

- Larger n ==> better approximation q_n ==> smaller KL divergence from p

- The reduction of KL divergence happens mostly when switching from unigram to bigram

# BLEU Results on Chinese-English NIST MT 2004 Tasks

| Decoding scheme | BLEU |
| --- | --- |
| Viterbi | 35.4 |
| MBR ($K{=}1000$) | 35.8 |
| Crunching ($N{=}10000$) | 35.7 |
| Crunching+MBR ($N{=}10000$) | 35.8 |
| Variational (1to4gram+wp+vt) | **36.6** |

(Kumar and Byrne, 2004)

(May and Knight, 2006)

New!

- variational decoding improves over Viterbi, MBR, and crunching

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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



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

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

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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



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

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

# Variational Inference

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

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

- Instead, we derive a simpler distribution q*

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



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

**tractable**
$$y^* \quad = \quad \arg\max_y q^*(y \mid x)$$

# Outline

- Hypergraph as Hypothesis Space

- Unsupervised Discriminative Training

  ▸ minimum imputed risk

  ▸ contrastive language model estimation

- Variational Decoding

- **First- and Second-order Expectation Semirings**

| **decoding** (e.g., mbr) | **training** (e.g., mert) |
|---|---|
| **atomic inference operations** (e.g., finding one-best, k-best or expectation, inference can be *exact* or *approximate*) ||

$S \mid 0,4$

$S \to \langle X_0, X_0 \rangle$

$S \to \langle X_0, X_0 \rangle$

$X \mid 0,4 \mid$ the $\cdots$ cat

$X \mid 0,4 \mid$ a $\cdots$ mat

**Probabilistic Hypergraph**

$X_1$ of $X_0 \rangle$

$X \mid 0,2 \mid$ the $\cdots$ mat

$X \mid 3,4 \mid$ a $\cdots$ cat

$X \to \langle$ dianzi shang, the mat $\rangle$

$X \to \langle$ mao, a cat $\rangle$

dianzi$_0$  shang$_1$        de$_2$        mao$_3$

- "decoding" quantities:
  - Viterbi
  - K-best
  - Counting
  - ......

- First-order quantities:
  - expectation
    - entropy
    - Bayes risk
    - cross-entropy
    - KL divergence
    - feature expectations
  - first-order gradient of Z

- Second-order quantities:
  - expectation over product
    - interaction between features
  - Hessian matrix of Z
    - second-order gradient descent
  - gradient of expectation
    - gradient of entropy or Bayes risk

77

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▸ **choose a semiring**

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▸ choose a semiring



    - ▸ specify a weight for each hyperedge

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▸ choose a semiring



    - ▸ specify a weight for each hyperedge

    - ▸ run the inside algorithm

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▸ **choose a semiring**

$$\langle K, \oplus, \otimes \rangle$$

a **set** with **plus** and **times** operations



dianzi$_0$  shang$_1$        de$_2$    mao$_3$

- ▸ **specify a weight for each hyperedge**

- ▸ **run the inside algorithm**

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▶ choose a semiring

$$\langle K, \oplus, \otimes \rangle$$

a **set** with **plus** and **times** operations

e.g., integer numbers with regular **+** and **×**

    - ▶ specify a weight for each hyperedge

    - ▶ run the inside algorithm

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▶ choose a semiring

$$\langle K, \oplus, \otimes \rangle$$

a **set** with **plus** and **times** operations

e.g., integer numbers with regular **+** and **×**

    - ▶ specify a weight for each hyperedge

each weight is a semiring member

    - ▶ run the inside algorithm

# Compute Quantities on a Hypergraph: a Recipe

- Semiring-weighted inside algorithm
  - three steps:
    - ▶ choose a semiring

$$\langle K, \oplus, \otimes \rangle$$

a **set** with **plus** and **times** operations

e.g., integer numbers with regular **+** and **×**

Hypergraph labels: $v_5$, $v_3$, $v_4$, $v_1$, $v_2$ with hyperedges $e_7$, $e_8$, $e_3$, $e_4$, $e_5$, $e_6$, $e_1$, $e_2$

dianzi$_0$  shang$_1$  de$_2$  mao$_3$

- ▶ specify a weight for each hyperedge

each weight is a semiring member

- ▶ run the inside algorithm

complexity is O(hypergraph size)

# Semirings

- "Decoding" time semirings (Goodman, 1999)

  - counting, Viterbi, K-best, etc.

- "Training" time semirings

  - first-order expectation semirings (Eisner, 2002)

  - second-order expectation semirings (new)

- Applications of the Semirings (new)

  - entropy, risk, gradient of them, and many more

$v_5$

$e_7$   $e_8$

$v_3$   $v_4$

$e_3$   $e_4$   $e_5$   $e_6$

$v_1$   $v_2$

$e_1$   $e_2$

dianzi$_0$   shang$_1$   de$_2$   mao$_3$

How many trees?

four ☺

compute it
use a semiring?

$v_5$

$e_7$   $e_8$

$v_3$   $v_4$

$e_3$   $e_4$

$e_5$   $e_6$

$v_1$

$v_2$

$e_1$

$e_2$

dianzi$_0$   shang$_1$          de$_2$      mao$_3$

**Three steps:**

# Compute the Number of Derivation Trees

**Three steps:**

▸ choose a semiring

**Three steps:**

▸ choose a semiring

counting semiring:
ordinary integers with
regular + and x

# Compute the Number of Derivation Trees

**Three steps:**

▸ choose a semiring

counting semiring:
ordinary integers with
regular + and x

▸ specify a weight for each hyperedge

# Compute the Number of Derivation Trees

**Three steps:**

▸ choose a semiring

    counting semiring:
      ordinary integers with
      regular + and x

▸ specify a weight for each hyperedge

# Compute the Number of Derivation Trees

## Three steps:

▸ choose a semiring

counting semiring:
ordinary integers with
regular + and x

▸ specify a weight for each hyperedge

▸ run the inside algorithm

**Bottom-up** process in computing the number of trees

$k(v_1) = k(e_1)$



**Bottom-up** process in computing the number of trees

$k(v_1) = 1$

$k(v_1) = k(e_1)$     $k(v_2) = k(e_2)$



**Bottom-up**
process in
computing the
number of trees

$v_5$

$1$ $e_7$     $e_8$ $1$

$v_3$     $v_4$

$e_3$     $e_4$     $1$     $1$

$1$     $1$     $e_5$     $e_6$

$v_1$   $k(v_1) = 1$     $v_2$   $k(v_2) = 1$

$e_1$ $1$     $1$ $e_2$

dianzi$_0$   shang$_1$     de$_2$     mao$_3$

$k(v_1) = k(e_1)$          $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$



**Bottom-up** process in computing the number of trees

$k(v_3)=2$   $v_3$

$k(v_1)=1$     $k(v_2)=1$

$e_1$ 1     $1\ e_2$

dianzi$_0$  shang$_1$          de$_2$   mao$_3$

$k(v_1) = k(e_1)$          $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$

$k(v_4) = k(e_5) \otimes k(v_1) \otimes k(v_2) \oplus k(e_6) \otimes k(v_1) \otimes k(v_2)$



**Bottom-up** process in computing the number of trees

$k(v_3) = 2$   $k(v_4) = 2$   $k(v_1) = 1$   $k(v_2) = 1$

dianzi$_0$  shang$_1$          de$_2$          mao$_3$

$k(v_1) = k(e_1)$        $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$

$k(v_4) = k(e_5) \otimes k(v_1) \otimes k(v_2) \oplus k(e_6) \otimes k(v_1) \otimes k(v_2)$

$k(v_5) = k(e_7) \otimes k(v_3) \oplus k(e_8) \otimes k(v_4)$

**Bottom-up** process in computing the number of trees

$k(v_1) = k(e_1)$         $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$

$k(v_4) = k(e_5) \otimes k(v_1) \otimes k(v_2) \oplus k(e_6) \otimes k(v_1) \otimes k(v_2)$

$k(v_5) = k(e_7) \otimes k(v_3) \oplus k(e_8) \otimes k(v_4)$



**Bottom-up** process in computing the number of trees

$v_5$    $k(v_5) = 4$

$e_7$    $e_8$

$k(v_3) = 2$  $v_3$    $v_4$    $k(v_4) = 2$

$e_3$  $e_4$    $e_5$  $e_6$

$v_1$  $k(v_1) = 1$    $v_2$  $k(v_2) = 1$

$e_1$    $e_2$

$dianzi_0$  $shang_1$    $de_2$  $mao_3$

82

$S \rightarrow \langle X_0, X_0 \rangle$

$S \mid 0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X \mid 0,4 \mid$ the $\cdots$ cat

$X \mid 0,4 \mid$ a $\cdots$ mat

$X \rightarrow \langle X_0$ de $X_1, X_1$ on $X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_0$ 's $X_1 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_0 \; X_1 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_1$ of $X_0 \rangle$

$X \mid 0,2 \mid$ the $\cdots$ mat

$X \mid 3,4 \mid$ a $\cdots$ cat

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi$_0$   shang$_1$            de$_2$   mao$_3$

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_0 \; X_1 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi$_0$ shang$_1$   de$_2$   mao$_3$

the mat a cat

p=3/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_0$ 's $X_1 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi$_0$ shang$_1$   de$_2$   mao$_3$

the mat 's a cat

p=1/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_1$ on $X_0 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi$_0$ shang$_1$   de$_2$   mao$_3$

a cat on the mat

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0$ de $X_1, X_1$ of $X_0 \rangle$

$X \rightarrow \langle$ dianzi shang, the mat $\rangle$

$X \rightarrow \langle$ mao, a cat $\rangle$

dianzi$_0$ shang$_1$   de$_2$   mao$_3$

a cat of the mat

83

expected translation length?

S→⟨X₀, X₀⟩

$S \to \langle X_0, X_0 \rangle$

$X \to \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \to \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \to \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$

$X \to \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \to \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \to \langle \text{mao}, \text{a cat} \rangle$

dianzi₀ shang₁    de₂    mao₃

p=2/8

the mat a cat

p=3/8

the mat 's a cat

p=1/8

a cat on the mat

p=2/8

a cat of the mat

83

expected translation length?

$S \rightarrow \langle X_0, X_0 \rangle$

$S|0,4$

$S \rightarrow \langle X_0, X_0 \rangle$

$X|0,4|\text{the} \cdots \text{cat}$

$X|0,4|\text{a} \cdots \text{mat}$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X|0,2|\text{the} \cdots \text{mat}$

$X|3,4|\text{a} \cdots \text{cat}$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀  shang₁

de₂  mao₃

p=3/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀  shang₁    de₂    mao₃

the mat 's a cat  5

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀  shang₁    de₂    mao₃

the mat a cat  4

p=1/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀  shang₁    de₂    mao₃

a cat on the mat  5

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀  shang₁    de₂    mao₃

a cat of the mat  5

83

expected translation length?

$$2/8 \times 4 + 6/8 \times 5 = 4.75$$

S|0,4

S→⟨X₀,X₀⟩

S→⟨X₀,X₀⟩

X|0,4|the ⋯ cat

X|0,4|a ⋯ mat

X→⟨X₀ de X₁,X₁ on X₀⟩

X→⟨X₀ de X₁,X₀ 's X₁⟩

X→⟨X₀ de X₁,X₀ X₁⟩

X→⟨X₀ de X₁,X₁ of X₀⟩

X|0,2|the ⋯ mat

X|3,4|a ⋯ cat

X→⟨dianzi shang,the mat⟩

X→⟨mao,a cat⟩

dianzi₀ shang₁

de₂ mao₃

p=2/8

S→⟨X₀,X₀⟩

X→⟨X₀ de X₁,X₀ X₁⟩

X→⟨dianzi shang,the mat⟩

X→⟨mao,a cat⟩

dianzi₀ shang₁ de₂ mao₃

the mat a cat   4

p=3/8

S→⟨X₀,X₀⟩

X→⟨X₀ de X₁,X₀ 's X₁⟩

X→⟨dianzi shang,the mat⟩

X→⟨mao,a cat⟩

dianzi₀ shang₁ de₂ mao₃

the mat 's a cat   5

p=1/8

S→⟨X₀,X₀⟩

X→⟨X₀ de X₁,X₁ on X₀⟩

X→⟨dianzi shang,the mat⟩

X→⟨mao,a cat⟩

dianzi₀ shang₁ de₂ mao₃

a cat on the mat   5

p=2/8

S→⟨X₀,X₀⟩

X→⟨X₀ de X₁,X₁ of X₀⟩

X→⟨dianzi shang,the mat⟩

X→⟨mao,a cat⟩

dianzi₀ shang₁ de₂ mao₃

a cat of the mat   5

83

expected translation length?

$$2/8 \times 4 + 6/8 \times 5 = 4.75$$

variance?

S→⟨X₀,X₀⟩

$S \rightarrow \langle X_0, X_0 \rangle$

p=2/8

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \ X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀ shang₁ de₂ mao₃

the mat a cat    4

p=3/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_0 \text{ 's } X_1 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀ shang₁ de₂ mao₃

the mat 's a cat    5

p=1/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ on } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀ shang₁ de₂ mao₃

a cat on the mat    5

p=2/8

$S \rightarrow \langle X_0, X_0 \rangle$

$X \rightarrow \langle X_0 \text{ de } X_1, X_1 \text{ of } X_0 \rangle$

$X \rightarrow \langle \text{dianzi shang}, \text{the mat} \rangle$

$X \rightarrow \langle \text{mao}, \text{a cat} \rangle$

dianzi₀ shang₁ de₂ mao₃

a cat of the mat    5

expected translation length?

$$2/8 \times 4 + 6/8 \times 5 = 4.75$$

variance?

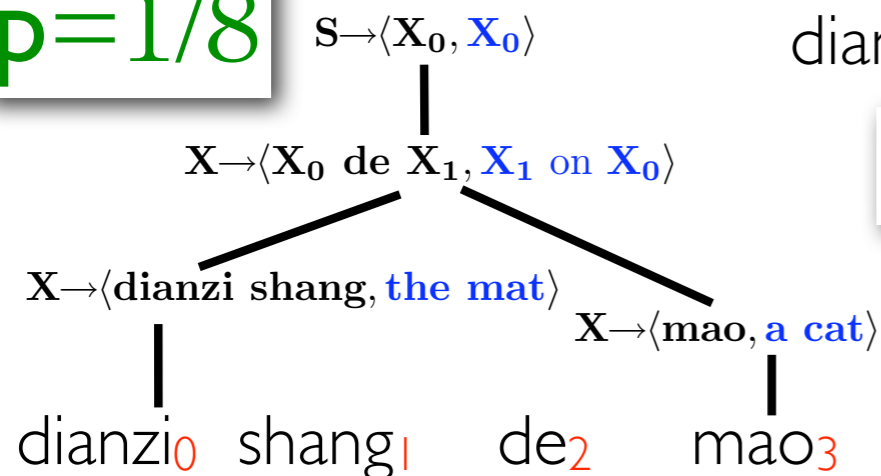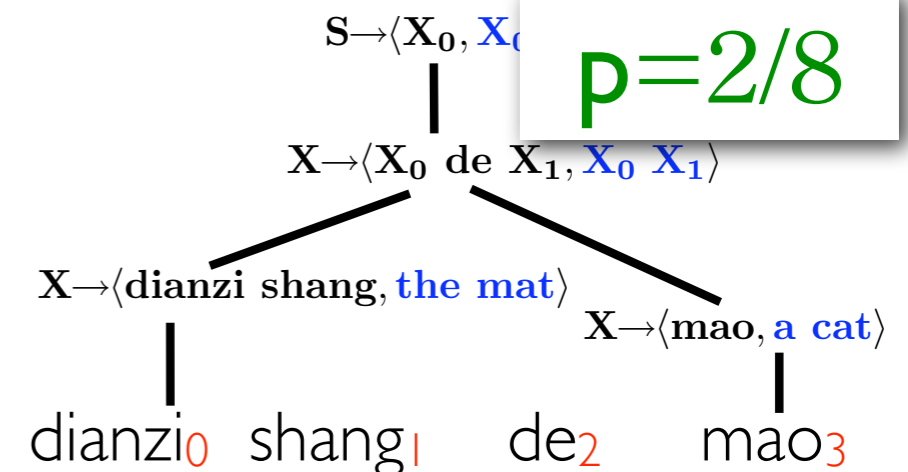$$2/8 \times (4\text{-}4.75)^2 + 6/8 \times (5\text{-}4.75)^2 \approx 0.19$$

p=2/8

the mat a cat    4
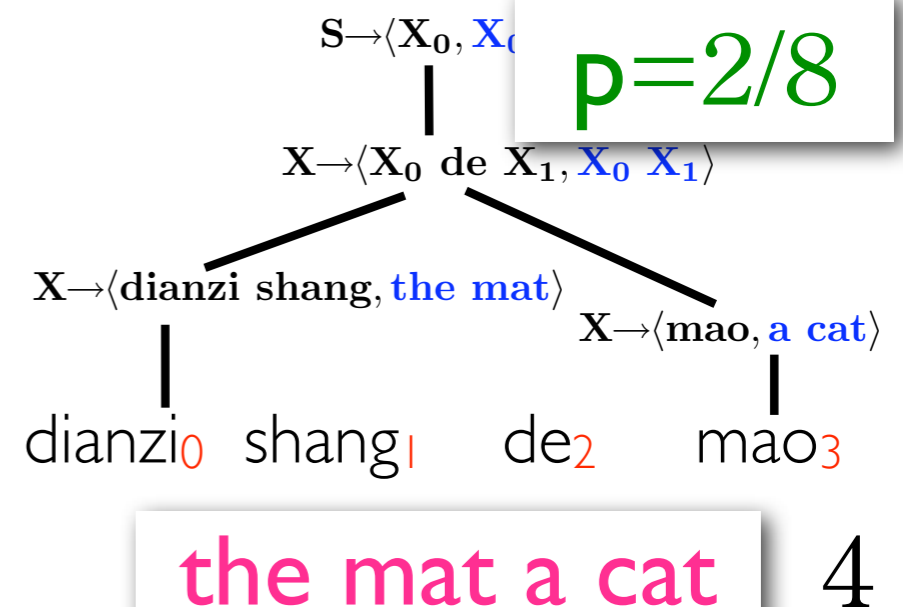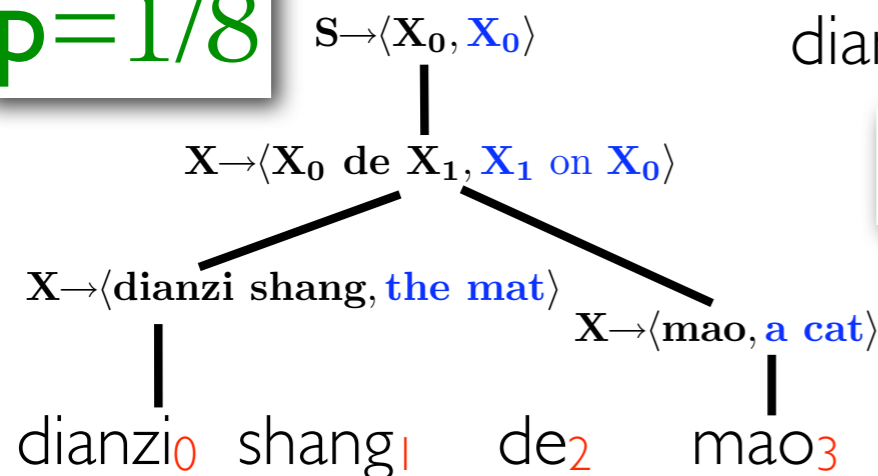
p=3/8

the mat 's a cat    5

p=1/8

a cat on the mat    5

p=2/8

a cat of the mat    5

# First- and Second-order Expectation Semirings

- each member is a **2**-tuple: $\langle p, r \rangle$

| | |
|---|---|
| $\langle p_1, r_1 \rangle \otimes \langle p_2, r_2 \rangle$ | $\langle p_1 p_2,\ p_1 r_2 + p_2 r_1 \rangle$ |
| $\langle p_1, r_1 \rangle \oplus \langle p_2, r_2 \rangle$ | $\langle p_1 + p_2,\ r_1 + r_2 \rangle$ |

Second-order:

- each member is a **4**-tuple: $\langle p, r, s, t \rangle$

| | |
|---|---|
| $\langle p_1, r_1, s_1, t_1 \rangle \otimes \langle p_2, r_2, s_2, t_2 \rangle$ | $\langle p_1 p_2,\ p_1 r_2 + p_2 r_1,\ p_1 s_2 + p_2 s_1,$ $p_1 t_2 + p_2 t_1 + r_1 s_2 + r_2 s_1 \rangle$ |
| $\langle p_1, r_1, s_1, t_1 \rangle \oplus \langle p_2, r_2, s_2, t_2 \rangle$ | $\langle p_1 + p_2,\ r_1 + r_2,\ s_1 + s_2,\ t_1 + t_2 \rangle$ |

$k(v_1) = k(e_1)$  $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$

$k(v_4) = k(e_5) \otimes k(v_1) \otimes k(v_2) \oplus k(e_6) \otimes k(v_1) \otimes k(v_2)$

$k(v_5) = k(e_7) \otimes k(v_3) \oplus k(e_8) \otimes k(v_4)$

$v_5$  $k(v_5) = \langle 8, 4.75 \rangle$

$k(v_3) = \langle 0, 1 \rangle$  $e_7 \langle 1, 0 \rangle$  $e_8 \langle 1, 0 \rangle$  $k(v_4) = \langle 1, 2 \rangle$

$v_3$  $v_4$

**First-order:**
each semiring member is a **2-tuple**

$e_3$  $e_4$ $\langle 3, 3 \rangle$  $\langle 1, 1 \rangle$ $\langle 2, 2 \rangle$

$\langle 2, 0 \rangle$  $e_5$  $e_6$

$k(v_2) = \langle 1, 2 \rangle$

$v_1$  $k(v_1) = \langle 1, 2 \rangle$  $v_2$

$e_1 \langle 1, 2 \rangle$  $\langle 1, 2 \rangle e_2$

dianzi$_0$  shang$_1$  de$_2$  mao$_3$

85

$k(v_1) = k(e_1)$     $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$

$k(v_4) = k(e_5) \otimes k(v_1) \otimes k(v_2) \oplus k(e_6) \otimes k(v_1) \otimes k(v_2)$

$k(v_5) = k(e_7) \otimes k(v_3) \oplus k(e_8) \otimes k(v_4)$

$v_5$  $k(v_5) = \langle 8, 4.5, 4.5, 5 \rangle$

$e_7$  $e_8 \langle 1,0,0,0 \rangle$

$\langle 1, 0,0,0 \rangle$

$k(v_3) = \langle 1,1,1,1 \rangle$  $v_3$     $v_4$  $k(v_4) = \langle 1,2,1,3 \rangle$

$\langle 1,1,1,1 \rangle$

**Second-order:**
each semiring member
is a **4-tuple**

$e_3$     $e_4$

$\langle 3,3,3,3 \rangle$     $\langle 2,2,2,2 \rangle$

$\langle 2,0,0,0 \rangle$     $e_5$   $e_6$

$v_1$  $k(v_1) = \langle 1,2,2,4 \rangle$     $v_2$  $k(v_2) = \langle 1,2,2,4 \rangle$

$e_1 \langle 1,2,2,4 \rangle$     $\langle 1,2,2,4 \rangle e_2$

dianzi$_0$  shang$_1$     de$_2$   mao$_3$

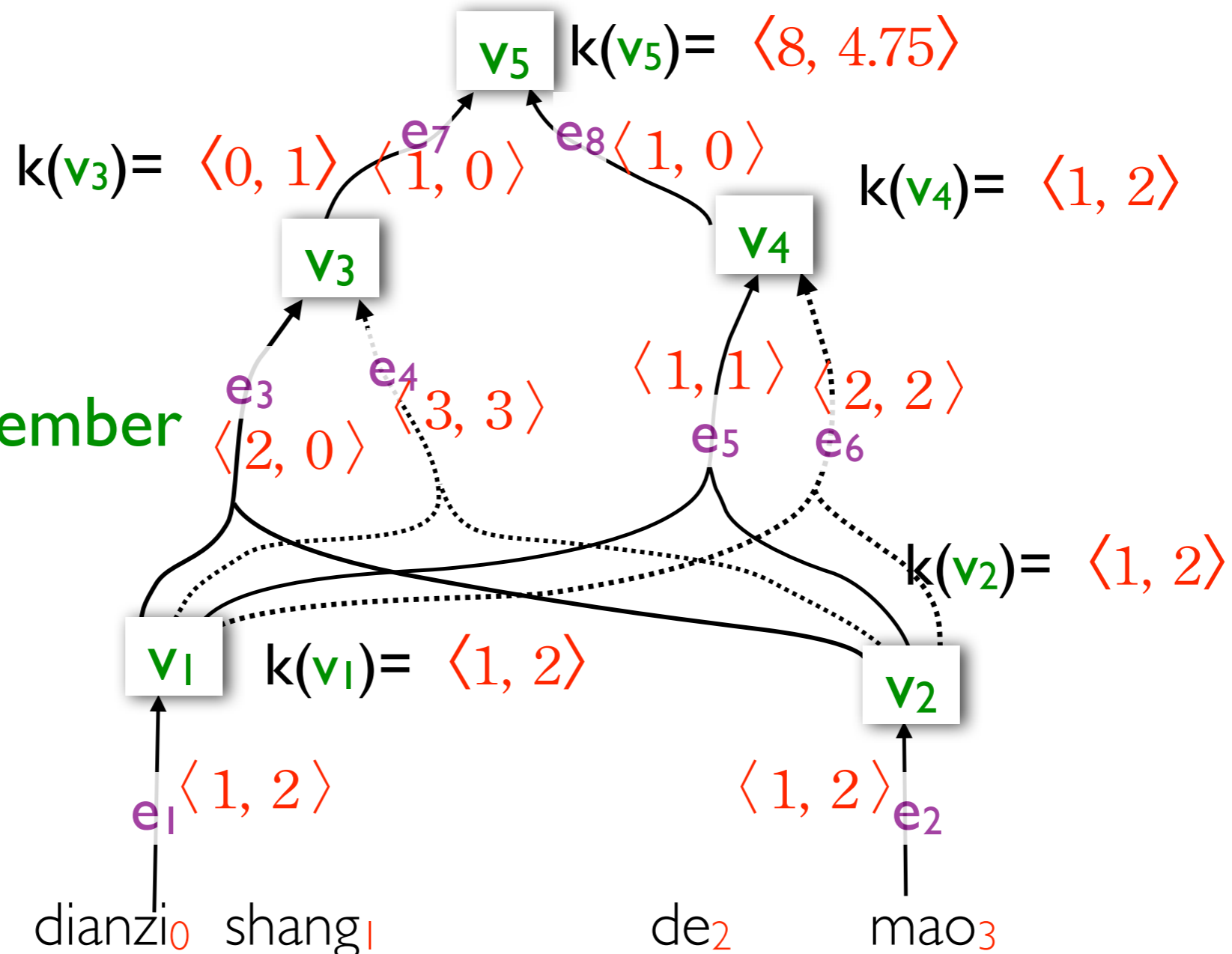$k(v_1) = k(e_1)$          $k(v_2) = k(e_2)$

$k(v_3) = k(e_3) \otimes k(v_1) \otimes k(v_2) \oplus k(e_4) \otimes k(v_1) \otimes k(v_2)$

$k(v_4) = k(e_5) \otimes k(v_1) \otimes k(v_2) \oplus k(e_6) \otimes k(v_1) \otimes k(v_2)$

$k(v_5) = k(e_7) \otimes k(v_3) \oplus k(e_8) \otimes k(v_4)$

$v_5$    $k(v_5) = \langle 8, 4.5, 4.5, 5 \rangle$

$e_7$    $e_8 \langle 1, 0, 0, 0 \rangle$

$\langle 1, 0, 0, 0 \rangle$

$k(v_3) = \langle 1, 1, 1, 1 \rangle$    $v_3$    $v_4$    $k(v_4) = \langle 1, 2, 1, 3 \rangle$

**Second-order:**
each semiring member
is a **4-tuple**

$e_3$    $e_4$

$\langle 1, 1, 1, 1 \rangle$

$\langle 3, 3, 3, 3 \rangle$    $e_5$    $\langle 2, 2, 2, 2 \rangle$    $e_6$

$\langle 2, 0, 0, 0 \rangle$

$v_1$    $k(v_1) = \langle 1, 2, 2, 4 \rangle$    $v_2$    $k(v_2) = \langle 1, 2, 2, 4 \rangle$

⚠️ WARNING
Fake

$e_1 \langle 1, 2, 2, 4 \rangle$    $\langle 1, 2, 2, 4 \rangle e_2$

dianzi$_0$  shang$_1$          de$_2$    mao$_3$

# Expectations on Hypergraphs

# Expectations on Hypergraphs

- Expectation over a hypergraph

# Expectations on Hypergraphs

- Expectation over a hypergraph

$$\overline{r} \overset{\text{def}}{=} \mathbb{E}_p[r] = \sum_{d \in \text{HG}} p(d)r(d)$$

- *r(d)* is a function over a derivation *d*

  e.g., the length of the translation yielded by *d*

# Expectations on Hypergraphs

- Expectation over a hypergraph

$$\overline{r} \stackrel{\text{def}}{=} \mathbb{E}_p[r] = \sum_{d \in \text{HG}} p(d)r(d)$$

exponential size

- *r(d)* is a function over a derivation *d*

  e.g., the length of the translation yielded by *d*

# Expectations on Hypergraphs

- Expectation over a hypergraph

$$\overline{r} \stackrel{\mathrm{def}}{=} \mathbb{E}_p[r] = \sum_{d \in \mathrm{HG}} p(d)r(d)$$

exponential size

- *r(d)* is a function over a derivation *d*

  e.g., the length of the translation yielded by *d*

- *r(d)* is additively decomposed

$$r(d) \stackrel{\mathrm{def}}{=} \sum_{e \in d} r_e$$

e.g., translation length is additively decomposed!

# Second-order Expectations on Hypergraphs

- Expectation of **products** over a hypergraph

$$\bar{t} \stackrel{\mathrm{def}}{=} \mathbb{E}_p[r \cdot s] = \sum_{d \in \mathrm{HG}} p(d) r(d) s(d)$$

exponential size

- r and s are additively decomposed

$$r(d) \stackrel{\mathrm{def}}{=} \sum_{e \in d} r_e$$

$$s(d) \stackrel{\mathrm{def}}{=} \sum_{e \in d} s_e$$

r and s can be identical or different functions.

# Compute expectation using expectation semiring:

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$\mathsf{p_e}$ : transition probability or log-linear score at edge $e$

$\mathsf{r_e}$ ?

### Entropy:

$$r_e \stackrel{\mathrm{def}}{=} \log p_e$$

## Compute expectation using expectation semiring:

$$k_e \overset{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$\mathsf{p_e}$ : transition probability or log-linear score at edge $e$

$\mathsf{r_e}$ ?

Entropy:

$$r_e \overset{\mathrm{def}}{=} \log p_e$$

Why?

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$\mathsf{p_e}$ : transition probability or log-linear score at edge $e$

$\mathsf{r_e}$ ?

**Entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log p_e$$

**Why?**     entropy is an **expectation**

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log p_e$$

**Why?** entropy is an **expectation**

$$\mathrm{H}(p) = \mathbb{E}_p[-\log p] = - \sum_{d \in \mathrm{HG}} p(d) \log p(d)$$

## Compute expectation using expectation semiring:

$$k_e \overset{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$\mathsf{p_e}$ : transition probability or log-linear score at edge $e$

$\mathsf{r_e}$ ?

**Entropy:**

$$r_e \overset{\mathrm{def}}{=} \log p_e$$

**Why?**    entropy is an **expectation**

$$\mathrm{H}(p) = \mathbb{E}_p[-\log p] = - \sum_{d \in \mathrm{HG}} p(d) \log p(d)$$

log *p(d)* is additively decomposed!

## Compute expectation using expectation semiring:

$$k_e \overset{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \overset{\mathrm{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \overset{\mathrm{def}}{=} \log q_e$$

**Why?**

**Compute expectation using expectation semiring:**

$$k_e \stackrel{\text{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \stackrel{\text{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \stackrel{\text{def}}{=} \log q_e$$

**Why?**    cross-entropy is an **expectation**

## Compute expectation using expectation semiring:

$$k_e \overset{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \overset{\mathrm{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \overset{\mathrm{def}}{=} \log q_e$$

**Why?**     cross-entropy is an **expectation**

$$H(p, q) = \mathbb{E}_p(-\log q) = -\sum_{d \in \mathrm{HG}} p(d) \log q(d)$$

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log q_e$$

**Why?**     cross-entropy is an **expectation**

$$\mathrm{H}(p, q) = \mathbb{E}_p(-\log q) = -\sum_{d \in \mathrm{HG}} p(d) \log q(d)$$

log $q(d)$ is additively decomposed!

## Compute expectation using expectation semiring:

$$k_e \overset{\text{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

Entropy:

$$r_e \overset{\text{def}}{=} \log p_e$$

Cross-entropy:

$$r_e \overset{\text{def}}{=} \log q_e$$

Bayes risk:

$$r_e \overset{\text{def}}{=} \text{loss at edge } e$$

Why?

## Compute expectation using expectation semiring:

$$k_e \overset{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \overset{\mathrm{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \overset{\mathrm{def}}{=} \log q_e$$

**Bayes risk:**

$$r_e \overset{\mathrm{def}}{=} \text{loss at edge } e$$

**Why?**   Bayes risk is an **expectation**

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log q_e$$

**Bayes risk:**

$$r_e \stackrel{\mathrm{def}}{=} \text{loss at edge } e$$

**Why?**    Bayes risk is an **expectation**

$$\mathrm{Risk} = \mathbb{E}_p(L) = -\sum_{d \in \mathrm{HG}} p(d) \cdot L(Y(d))$$

## Compute expectation using expectation semiring:

$$k_e \stackrel{\mathrm{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \stackrel{\mathrm{def}}{=} \log q_e$$

**Bayes risk:**

$$r_e \stackrel{\mathrm{def}}{=} \text{loss at edge } e$$

**Why?**     Bayes risk is an **expectation**

$$\text{Risk} = \mathbb{E}_p(L) = - \sum_{d \in \mathrm{HG}} p(d) \cdot L(Y(d))$$

$L(Y(d))$ is additively decomposed!

## Compute expectation using expectation semiring:

$$k_e \overset{\text{def}}{=} \langle p_e, p_e r_e \rangle$$

$p_e$ : transition probability or log-linear score at edge $e$

$r_e$ ?

**Entropy:**

$$r_e \overset{\text{def}}{=} \log p_e$$

**Cross-entropy:**

$$r_e \overset{\text{def}}{=} \log q_e$$

**Bayes risk:**

$$r_e \overset{\text{def}}{=} \text{loss at edge } e$$

**Why?**  Bayes risk is an **expectation**

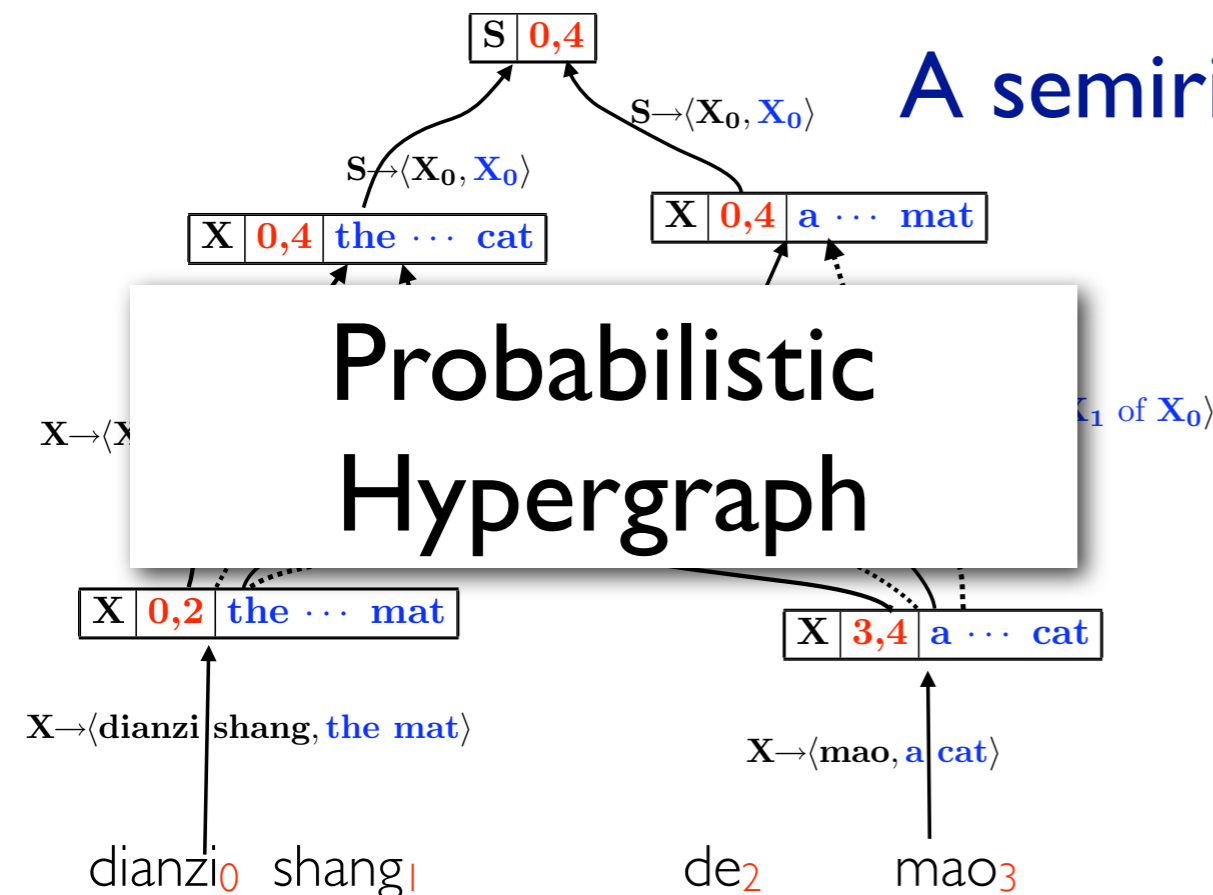$$\text{Risk} = \mathbb{E}_p(L) = - \sum_{d \in \text{HG}} p(d) \cdot L(Y(d))$$

$L(Y(d))$ is additively decomposed!  (Tromble et al. 2008)

91

# Applications of Expectation Semirings: a Summary

| Quantity | $k_e$ | $k_{\mathbf{root}}$ | Final |
|---|---|---|---|
| **Expectation** | $\langle p_e, p_e r_e \rangle$ | $\langle Z, \bar{r} \rangle$ | $\bar{r}/Z$ |
| Entropy | $r_e \overset{\text{def}}{=} \log p_e$, so $k_e = \langle p_e, p_e \log p_e \rangle$ | $\langle Z, \bar{r} \rangle$ | $\log Z - \bar{r}/Z$ |
| Cross-entropy | $\langle q_e \rangle$ <br> $r_e \overset{\text{def}}{=} \log q_e$, so $k_e = \langle p_e, p_e \log q_e \rangle$ | $\langle Z_q \rangle$ <br> $\langle Z_p, \bar{r} \rangle$ | $\log Z_q - \bar{r}/Z_p$ |
| Bayes risk | $r_e \overset{\text{def}}{=} \mathrm{L}_e$, so $k_e = \langle p_e, p_e \mathrm{L}_e \rangle$ | $\langle Z, \bar{r} \rangle$ | $\bar{r}/Z$ |
| **First-order gradient** | $\langle p_e, \nabla p_e \rangle$ | $\langle Z, \nabla Z \rangle$ | $\nabla Z$ |
| **Covariance matrix** | $\langle p_e, p_e r_e, p_e s_e, p_e r_e s_e \rangle$ | $\langle Z, \bar{r}, \bar{s}, \bar{t} \rangle$ | $\frac{\bar{t}}{Z} - \frac{\bar{r}\,\bar{s}^{\mathbf{T}}}{Z^2}$ |
| **Hessian matrix** | $\langle p_e, \nabla p_e, \nabla p_e, \nabla^2 p_e \rangle$ | $\langle Z, \nabla Z, \nabla Z, \nabla^2 Z \rangle$ | $\nabla^2 Z$ |
| **Gradient of expectation** | $\langle p_e, p_e r_e, \nabla p_e, (\nabla p_e) r_e + p_e (\nabla r_e) \rangle$ | $\langle Z, \bar{r}, \nabla Z, \nabla \bar{r} \rangle$ | $\frac{Z \nabla \bar{r} - \bar{r} \nabla Z}{Z^2}$ |
| Gradient of entropy | $\langle p_e, p_e \log p_e, \nabla p_e, (1 + \log p_e)\nabla p_e \rangle$ | $\langle Z, \bar{r}, \nabla Z, \nabla \bar{r} \rangle$ | $\frac{\nabla Z}{Z} - \frac{Z \nabla \bar{r} - \bar{r} \nabla Z}{Z^2}$ |
| Gradient of risk | $\langle p_e, p_e \mathrm{L}_e, \nabla p_e, \mathrm{L}_e \nabla p_e \rangle$ | $\langle Z, \bar{r}, \nabla Z, \nabla \bar{r} \rangle$ | $\frac{Z \nabla \bar{r} - \bar{r} \nabla Z}{Z^2}$ |

S  0,4

S→⟨X₀, X₀⟩

S↛⟨X₀, X₀⟩

X  0,4  the ⋯ cat

X  0,4  a ⋯ mat

Probabilistic Hypergraph

X→⟨X ...

X₁ of X₀⟩

X  0,2  the ⋯ mat

X  3,4  a ⋯ cat

X→⟨dianzi shang, the mat⟩

X→⟨mao, a cat⟩

dianzi₀  shang₁          de₂     mao₃

- "decoding" quantities:
  - Viterbi
  - K-best
  - Counting
  - ......

- First-order quantities:
  - expectation
    - entropy
    - Bayes risk
    - cross-entropy
    - KL divergence
    - feature expectations
  - first-order gradient of Z

- Second-order quantities:
  - Expectation over product
    - interaction between features
  - Hessian matrix of Z
    - second-order gradient descent
  - gradient of expectation
    - gradient of entropy or Bayes risk

93