

# Statistical Machine Translation with Weighted Grammars

Matthias BÜchse

Chair of Foundations of Programming  
Institute of Theoretical Computer Science  
Technische Universität Dresden

September 12, 2011

# Outline

## Statistical Machine Translation

- Introduction

- Building an SMT System

- Training

- Evaluation

- Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

- Core: Grammar

- Feature Selection

- Decoding

- Further Grammar Classes

## My Dissertation

- Key Aspects

- Publications

# Outline

## Statistical Machine Translation

### Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

Further Grammar Classes

## My Dissertation

Key Aspects

Publications

# Ambiguity in Natural Language

Example (Chiang 2007)

- ▶ e = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶ f = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

# Ambiguity in Natural Language

Example (Chiang 2007)

- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Is  $e$  an adequate translation of  $f$ ?

# Ambiguity in Natural Language

Example (Chiang 2007)

- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

Is  $e$  an adequate translation of  $f$ ?



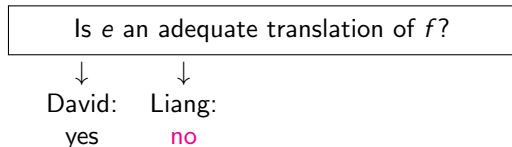
David:

yes

# Ambiguity in Natural Language

Example (Chiang 2007)

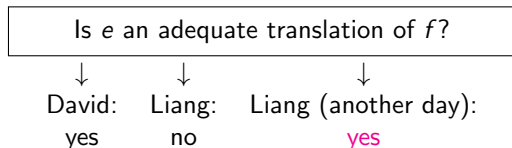
- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”



# Ambiguity in Natural Language

Example (Chiang 2007)

- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”





# Ambiguity in Natural Language

Example (Chiang 2007)

- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

What is  $P(e | f)$ ?

# Ambiguity in Natural Language

Example (Chiang 2007)

- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

What is  $P(e | f)$ ?

↓  
.42

# Search (aka Decoding) in SMT

From: German ▼  To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



 Listen

**German to English translation**

I would not buy this carpet.

 Listen

**New!** Click the words above to view alternate translations. [Dismiss](#)

given  $p \in \mathcal{M}(E | F)$

# Search (aka Decoding) in SMT

From: German ▼  To: English ▼

Ich möchte diesen Teppich nicht kaufen.



 Listen

**German to English translation**

I would not buy this carpet.

 Listen

**New!** Click the words above to view alternate translations. [Dismiss](#)

given  $p \in \mathcal{M}(E | F)$  ( $\sum_e p(e | f) = 1$  for every  $f$ )

# Search (aka Decoding) in SMT

From: German ▼  To: English ▼

Ich möchte diesen Teppich nicht kaufen.



 Listen

**German to English translation**

I would not buy this carpet.

 Listen

**New!** Click the words above to view alternate translations. [Dismiss](#)

given  $p \in \mathcal{M}(E | F)$  and  $f \in F$ ,

# Search (aka Decoding) in SMT

From: German ▼ To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



Listen

German to English translation

I would not buy this carpet.

Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given  $p \in \mathcal{M}(E | F)$  and  $f \in F$ ,  
determine  $\hat{e}$  such that for every  $e \in E$ :

$$p(\hat{e} | f) \geq p(e | f) . \quad (\text{not unique})$$

# Search (aka Decoding) in SMT

From: German ▼ To: English ▼ Translate

Ich möchte diesen Teppich nicht kaufen.



Listen

German to English translation

I would not buy this carpet.

Listen

New! Click the words above to view alternate translations. [Dismiss](#)

given  $p \in \mathcal{M}(E | F)$  and  $f \in F$ ,  
determine  $\hat{e}$  such that for every  $e \in E$ :

$$p(\hat{e} | f) \geq p(e | f) . \quad (\text{not unique})$$

for short

$$\hat{e} = \operatorname{argmax}_{e \in E} p(e | f) .$$

# Outline

## Statistical Machine Translation

Introduction

**Building an SMT System**

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

Further Grammar Classes

## My Dissertation

Key Aspects

Publications



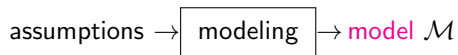
# Building an SMT System

assumptions

training data

test data

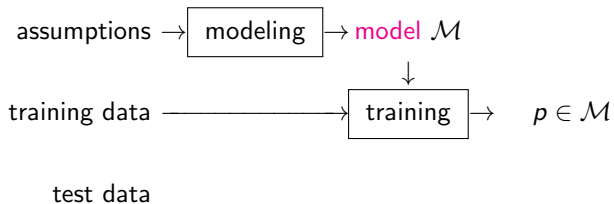
# Building an SMT System



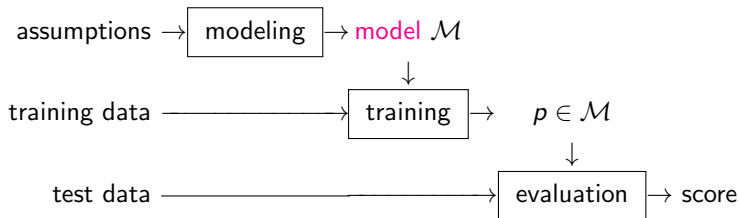
training data

test data

# Building an SMT System



# Building an SMT System



# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

**Training**

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

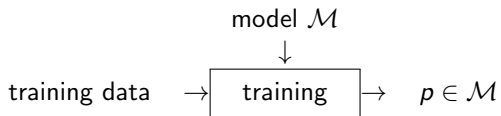
Further Grammar Classes

## My Dissertation

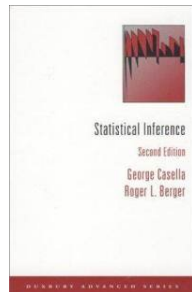
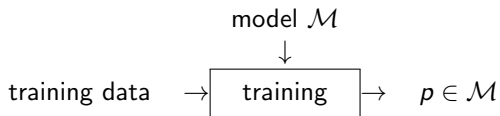
Key Aspects

Publications

# Training = Point Estimation

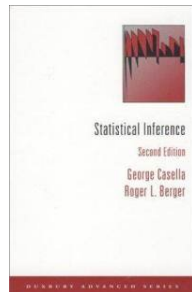
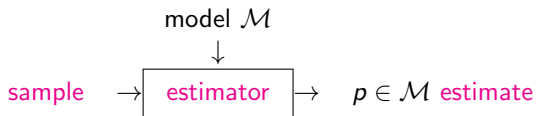


# Training = Point Estimation



Statistical Inference, Ch. 7:  
Point Estimation  
(Casella and Berger 2002)

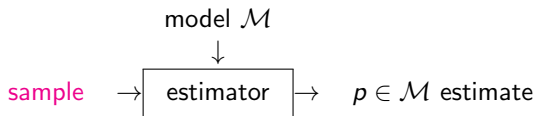
# Training = Point Estimation



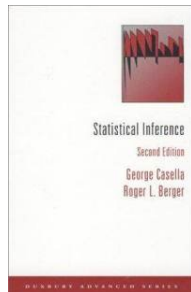
Statistical Inference, Ch. 7:  
Point Estimation  
(Casella and Berger 2002)



# Training = Point Estimation

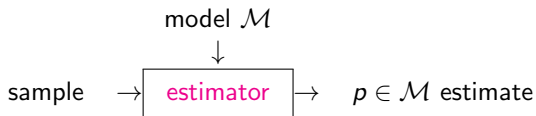


- sample
- ▶ representative of all translation situations
  - ▶ e. g., Hong Kong Hansards (parliament proceedings)



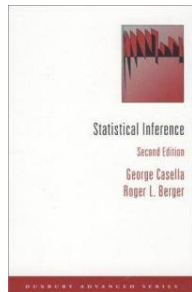
Statistical Inference, Ch. 7:  
Point Estimation  
(Casella and Berger 2002)

# Training = Point Estimation



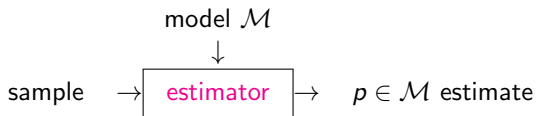
- sample
- ▶ representative of all translation situations
  - ▶ e. g., Hong Kong Hansards (parliament proceedings)

- estimators
- ▶ maximum likelihood
  - ▶ maximum a-posteriori
  - ▶ minimum risk
  - ▶ maximum entropy



Statistical Inference, Ch. 7:  
Point Estimation  
(Casella and Berger 2002)

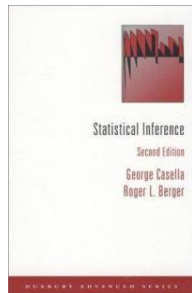
# Training = Point Estimation



- sample
- ▶ representative of all translation situations
  - ▶ e. g., Hong Kong Hansards (parliament proceedings)

- estimators
- ▶ maximum likelihood
  - ▶ maximum a-posteriori
  - ▶ minimum risk
  - ▶ maximum entropy

- properties
- ▶ loss function optimality (Statistical Decision Theory)
  - ▶ asymptotic optimality (consistency)



Statistical Inference, Ch. 7:  
Point Estimation  
(Casella and Berger 2002)

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

**Evaluation**

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

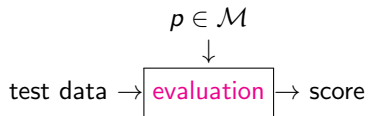
Further Grammar Classes

## My Dissertation

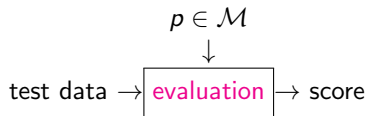
Key Aspects

Publications

# Evaluation

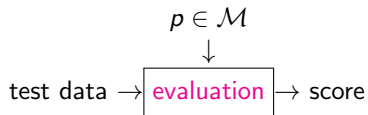


# Evaluation



$$\text{training data} \cap \text{test data} = \emptyset$$

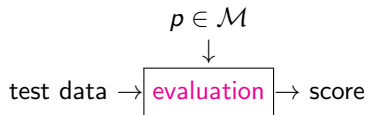
# Evaluation



$$\text{training data} \cap \text{test data} = \emptyset$$

method	how it works	goal
BLEU (Papineni et al. 2002)	$n$ -gram precision, brevity penalty	higher values

# Evaluation

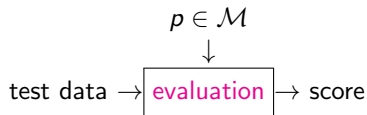


$$\text{training data} \cap \text{test data} = \emptyset$$

method	how it works	goal
BLEU (Papineni et al. 2002)	$n$ -gram precision, brevity penalty	higher values
TER (Snover et al. 2006)	edit distance	lower values



# Evaluation



$$\text{training data} \cap \text{test data} = \emptyset$$

method	how it works	goal
BLEU (Papineni et al. 2002)	$n$ -gram precision, brevity penalty	higher values
TER (Snover et al. 2006)	edit distance	lower values
Meteor (Banerjee and Lavie 2005)	alignments, unigram precision, unigram recall, penalty	higher values

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

**Adding Sophistication**

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

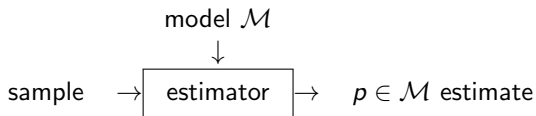
Further Grammar Classes

## My Dissertation

Key Aspects

Publications

# Training = Point Estimation



sample ▶ representative of all translation situations

▶ e. g., Hong Kong Hansards (parliament proceedings)

estimators ▶ maximum likelihood

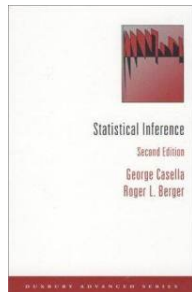
▶ maximum a-posteriori

▶ minimum risk

▶ maximum entropy

properties ▶ loss function optimality (Statistical Decision Theory)

▶ asymptotic optimality (consistency)



Statistical Inference, Ch. 7:  
Point Estimation  
(Casella and Berger 2002)

# Samples Are Flawed

- ▶ data sparsity

# Samples Are Flawed

- ▶ data sparsity
- ▶ noise

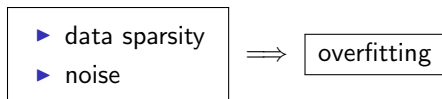
“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

# Samples Are Flawed

- ▶ data sparsity
- ▶ noise

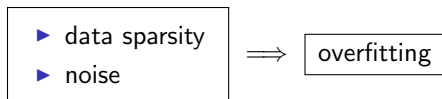
“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

# Samples Are Flawed



“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

# Samples Are Flawed

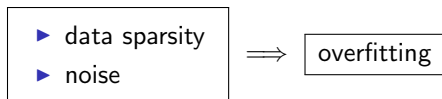


“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

Approaches: {



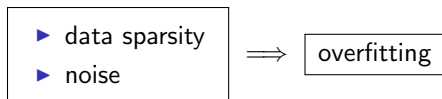
# Samples Are Flawed



“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

Approaches: {  
▶ use prior (smoothing)

# Samples Are Flawed

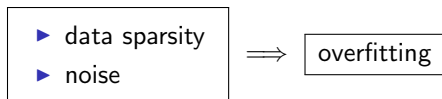


“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

Approaches: {

- ▶ use prior (smoothing)
- ▶ use ad-hoc estimator (heuristics)

# Samples Are Flawed



“There is no right answer in translation, but there are wrong answers. Translators make a lot of mistakes.” Kevin Knight

- Approaches: {
- ▶ use prior (smoothing)
  - ▶ use ad-hoc estimator (heuristics)
  - ▶ model appropriately

# Log-linear Model

(Berger, Della Pietra, and Della Pietra 1996; Och and Ney 2002)

- ▶  $h_1, \dots, h_n: E \times F \rightarrow \mathbb{R}$  features

# Log-linear Model

(Berger, Della Pietra, and Della Pietra 1996; Och and Ney 2002)

- ▶  $h_1, \dots, h_n: E \times F \rightarrow \mathbb{R}$  features
- ▶ maximum-entropy formulation:

$$\mathcal{M} = \{p \in \mathcal{M}(E | F) \mid \forall i: \mathbb{E}_p[h_i] = \mathbb{E}_{\tilde{p}}[h_i]\}$$

# Log-linear Model

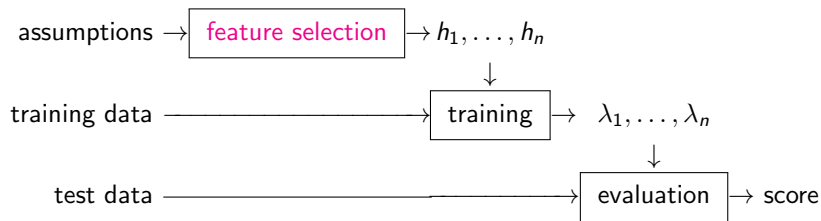
(Berger, Della Pietra, and Della Pietra 1996; Och and Ney 2002)

- ▶  $h_1, \dots, h_n: E \times F \rightarrow \mathbb{R}$  features
- ▶ maximum-likelihood formulation:

$$\mathcal{M} = \{p \in \mathcal{M}(E | F) \mid \exists \lambda_1, \dots, \lambda_n \in \mathbb{R}:$$

$$p(e | f) = \frac{\exp \sum_i \lambda_i \cdot h_i(e, f)}{\sum_{e' \in E} \exp \sum_i \lambda_i \cdot h_i(e', f)} \}$$

# Building an SMT System with Log-linear Models



# Feature Selection (Examples)

- ▶ translation model feature

$$h_1(e, f) = \log p_{\text{TM}}(e | f)$$

estimated from a model  $\mathcal{M}_{\text{TM}} \subseteq \mathcal{M}(E | F)$  via a sample of translations



# Feature Selection (Examples)

- ▶ translation model feature

$$h_1(e, f) = \log p_{\text{TM}}(e | f)$$

estimated from a model  $\mathcal{M}_{\text{TM}} \subseteq \mathcal{M}(E | F)$  via a sample of translations

- ▶ language model feature

$$h_2(e, f) = \log p_{\text{LM}}(e)$$

estimated from a model  $\mathcal{M}_{\text{LM}} \subseteq \mathcal{M}(E)$  via a sample of sentences

# Feature Selection (Examples)

- ▶ translation model feature

$$h_1(e, f) = \log p_{\text{TM}}(e | f)$$

estimated from a model  $\mathcal{M}_{\text{TM}} \subseteq \mathcal{M}(E | F)$  via a sample of translations

- ▶ language model feature

$$h_2(e, f) = \log p_{\text{LM}}(e)$$

estimated from a model  $\mathcal{M}_{\text{LM}} \subseteq \mathcal{M}(E)$  via a sample of sentences

- ▶ length feature

$$h_3(e, f) = |e|$$

## More Recent Estimators

- ▶ minimum-error-rate training (Och 2003)
- ▶ minimum-risk annealing (Smith and Eisner 2006)
- ▶ large-margin methods (Chiang)

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

**Core: Grammar**

Feature Selection

Decoding

Further Grammar Classes

## My Dissertation

Key Aspects

Publications

## Example (Chiang 2007)

### Recall

- ▶  $e$  = “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f$  = “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

## Example (Chiang 2007)

### Recall

- ▶  $e =$  “Australia is one of the few countries that have diplomatic relations with North Korea”
- ▶  $f =$  “Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi”

### Synchronous Context-Free Grammar $G$ :

$\pi_1: S \rightarrow \langle S X, S X \rangle$

$\pi_2: S \rightarrow \langle X, X \rangle$

$\pi_3: X \rightarrow \langle \text{yu } X_{[1]} \text{ you } X_{[2]}, \text{ have } X_{[2]} \text{ with } X_{[1]} \rangle$

$\pi_4: X \rightarrow \langle X_{[1]} \text{ de } X_{[2]}, \text{ the } X_{[2]} \text{ that } X_{[1]} \rangle$

$\pi_5: X \rightarrow \langle X \text{ zhiyi, one of } X \rangle$

$\pi_6: X \rightarrow \langle \text{Aozhou, Australia} \rangle$

$\pi_7: X \rightarrow \langle \text{Beihan, North Korea} \rangle$

$\pi_8: X \rightarrow \langle \text{shi, is} \rangle$

$\pi_9: X \rightarrow \langle \text{bangjiao, diplomatic relations} \rangle$

$\pi_{10}: X \rightarrow \langle \text{shaoshu guojia, few countries} \rangle$

# Derivation

$\langle S, S \rangle$

# Derivation

$$\begin{array}{l} \langle S, S \rangle \\ \xRightarrow{\pi_1} \langle S X, S X \rangle \end{array}$$



# Derivation

$$\begin{array}{l} \langle S, S \rangle \\ \xRightarrow{\pi_1} \langle S X, S X \rangle \end{array}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \end{aligned}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{\boxed{1}} X_{\boxed{2}}, S X_{\boxed{1}} X_{\boxed{2}} \rangle \end{aligned}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \end{aligned}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \end{aligned}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \end{aligned}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xrightarrow{\pi_1} & \langle S X, S X \rangle \\ \xrightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xrightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \end{aligned}$$

# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xRightarrow{\pi_1} & \langle S X, S X \rangle \\ \xRightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_8} & \langle \text{Aozhou } \text{shi } X, \text{Australia } \text{is } X \rangle \end{aligned}$$



# Derivation

$$\begin{aligned} & \langle S, S \rangle \\ \xRightarrow{\pi_1} & \langle S X, S X \rangle \\ \xRightarrow{\pi_1} & \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_2} & \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_6} & \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle \\ \xRightarrow{\pi_8} & \langle \text{Aozhou shi } X, \text{Australia is } X \rangle \end{aligned}$$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi}, \text{Australia is one of } X \rangle$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
 $\text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
 $\quad \text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the  $X_{[2]}$  that  $X_{[1]}$   $\rangle$   
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi } \text{yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the  $X_{[2]}$  that have  $X_{[0]}$  with  $X_{[1]}$   $\rangle$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the  $X_{[2]}$  that  $X_{[1]}$   $\rangle$   
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the  $X_{[2]}$  that have  $X_{[0]}$  with  $X_{[1]}$   $\rangle$

# Derivation

$\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that } X\_{[1]} \rangle  
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } X\_{[0]} \text{ with } X\_{[1]} \rangle  
 $\xRightarrow{\pi_7} \langle \text{Aozhou shi yu } \text{Beihan} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } X\_{[0]} \text{ with } \text{North Korea} \rangle



# Derivation

$\langle S, S \rangle$   
 $\xrightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xrightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xrightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xrightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xrightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xrightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xrightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that } X\_{[1]} \rangle  
 $\xrightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } X\_{[0]} \text{ with } X\_{[1]} \rangle  
 $\xrightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } X\_{[0]} \text{ with North Korea} \rangle

# Derivation

- $\langle S, S \rangle$   
 $\xrightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xrightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xrightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xrightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xrightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xrightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xrightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that } X\_{[1]} \rangle  
 $\xrightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } X\_{[0]} \text{ with } X\_{[1]} \rangle  
 $\xrightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } X\_{[0]} \text{ with North Korea} \rangle  
 $\xrightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you } \text{bangjiao} \text{ de } X_{[2]} \text{ zhiyi,}$   
    Australia is one of the } X\_{[2]} \text{ that have } \text{diplomatic relations} \text{ with North Korea} \rangle

# Derivation

- $\langle S, S \rangle$
- $\xRightarrow{\pi_1} \langle S X, S X \rangle$
- $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$
- $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$
- $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$
- $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$
- $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$
- $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
 $\text{Australia is one of the } X_{[2]} \text{ that } X_{[1]} \rangle$
- $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
 $\text{Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with } X_{[1]} \rangle$
- $\xRightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
 $\text{Australia is one of the } X_{[2]} \text{ that have } X_{[0]} \text{ with North Korea} \rangle$
- $\xRightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you bangjiao de } X_{[2]} \text{ zhiyi,}$   
 $\text{Australia is one of the } X_{[2]} \text{ that have diplomatic relations with North Korea} \rangle$

# Derivation

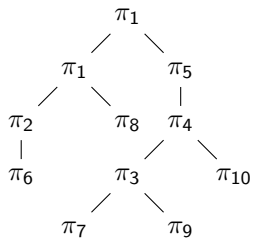
- $\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that  $X_{[1]}$   $\rangle$   
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that have  $X_{[0]}$  with  $X_{[1]}$   $\rangle$   
 $\xRightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that have  $X_{[0]}$  with North Korea  $\rangle$   
 $\xRightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you bangjiao de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that have diplomatic relations with North Korea  $\rangle$   
 $\xRightarrow{\pi_{10}} \langle \text{Aozhou shi yu Beihan you bangjiao de } \text{shaoshu guojia} \text{ zhiyi,}$   
Australia is one of the few countries that have diplomatic relations with N. K.  $\rangle$

# Derivation

- $\langle S, S \rangle$   
 $\xRightarrow{\pi_1} \langle S X, S X \rangle$   
 $\xRightarrow{\pi_1} \langle S X_{[1]} X_{[2]}, S X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_2} \langle X_{[0]} X_{[1]} X_{[2]}, X_{[0]} X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_6} \langle \text{Aozhou } X_{[1]} X_{[2]}, \text{Australia } X_{[1]} X_{[2]} \rangle$   
 $\xRightarrow{\pi_8} \langle \text{Aozhou shi } X, \text{Australia is } X \rangle$   
 $\xRightarrow{\pi_5} \langle \text{Aozhou shi } X \text{ zhiyi, Australia is one of } X \rangle$   
 $\xRightarrow{\pi_4} \langle \text{Aozhou shi } X_{[1]} \text{ de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that  $X_{[1]}$   $\rangle$   
 $\xRightarrow{\pi_3} \langle \text{Aozhou shi yu } X_{[1]} \text{ you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that have  $X_{[0]}$  with  $X_{[1]}$   $\rangle$   
 $\xRightarrow{\pi_7} \langle \text{Aozhou shi yu Beihan you } X_{[0]} \text{ de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that have  $X_{[0]}$  with North Korea  $\rangle$   
 $\xRightarrow{\pi_9} \langle \text{Aozhou shi yu Beihan you bangjiao de } X_{[2]} \text{ zhiyi,}$   
Australia is one of the  $X_{[2]}$  that have diplomatic relations with North Korea  $\rangle$   
 $\xRightarrow{\pi_{10}} \langle \text{Aozhou shi yu Beihan you bangjiao de shaoshu guojia zhiyi,}$   
Australia is one of the few countries that have diplomatic relations with N. K.  $\rangle$

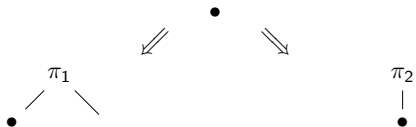
# Inducing a Model

## Derivation Tree



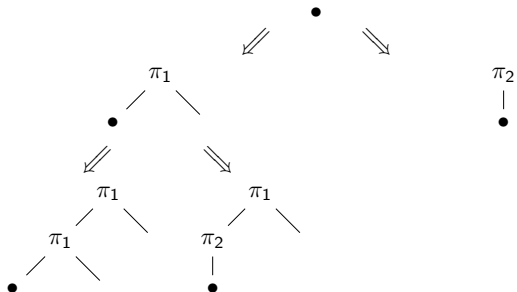
# Inducing a Model

## Decision Tree



# Inducing a Model

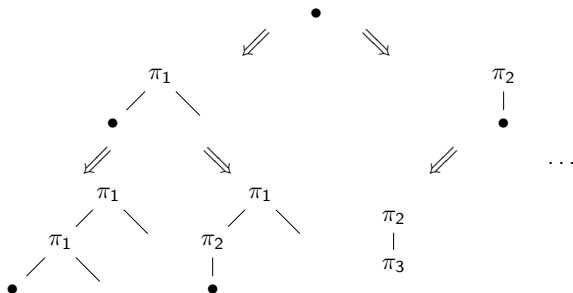
## Decision Tree





# Inducing a Model

## Decision Tree



# Inducing a Model

## Assignments

grammar  $G$  } “distribution”  $p'(e, d, f)$   
assignment  $p'$  }

model  $\mathcal{M}_G = \{p \in \mathcal{M}(E, F) \mid \exists p' : p(e, f) = \sum_d p'(e, d, f)\}$

# Training: Rule Extraction

## Word-aligned Training Pair

⟨30 duonianlai de youhao hezuo, friendly cooperation over the last 30 years⟩

	friendly	cooperation	over	the	last	30	years
30						■	
duonianlai			■	■	■		■
de							
youhao	■						
hezuo		■					

Figure 2 of (Chiang 2007)

# Training: Rule Extraction

## Initial Rules

X → ⟨30 duonianlai de youhao hezuo, friendly cooperation over the last 30 years⟩

X → ⟨duonianlai, over the last 30 years⟩

X → ⟨youhao hezuo, friendly cooperation⟩

X → ⟨30, 30⟩

X → ⟨youhao, friendly⟩

X → ⟨hezuo, cooperation⟩

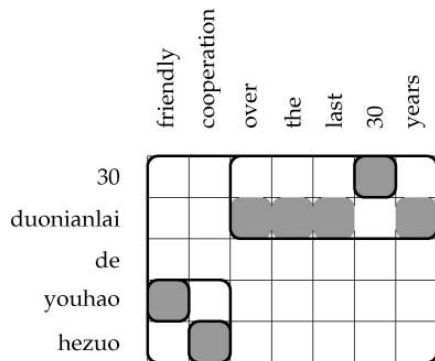


Figure 2 of (Chiang 2007)

# Training: Rule Extraction

## Example Rule

$X \rightarrow \langle X_{[1]} \text{ duonianlai de } X_{[2]}, X_{[2]} \text{ over the last } X_{[1]} \text{ years} \rangle$

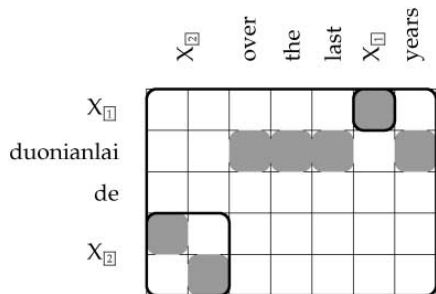


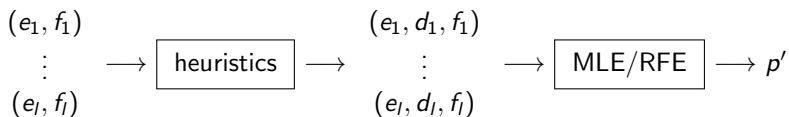
Figure 2 of (Chiang 2007)

# Training: Maximum-Likelihood Estimation

- ▶ training data: sequence  $\mathcal{D} = (e_1, f_1), \dots, (e_l, f_l)$
- ▶ model:  $\mathcal{M}_G = \{p \in \mathcal{M}(E, F) \mid \exists p' : p(e, f) = \sum_d p'(e, d, f)\}$
- ▶ likelihood:  $P(\mathcal{D} \mid p') = \prod_j \sum_d p'(e_j, d, f_j)$
- ▶ use EM algorithm

# Training: Reality

Chiang (2007), DeNeeffe and Knight (2009):



# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

**Feature Selection**

Decoding

Further Grammar Classes

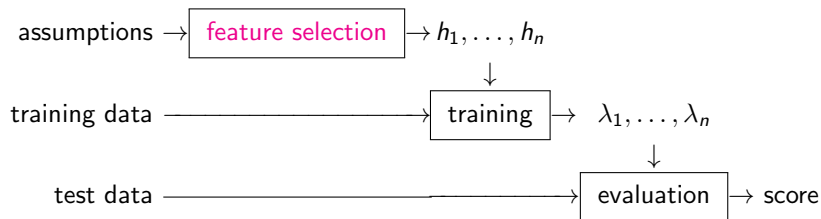
## My Dissertation

Key Aspects

Publications



# Building an SMT System with Log-linear Models



# Feature Selection

Hidden variable: derivation of  $(e, f)$  in  $G$

# Feature Selection

Hidden variable: derivation of  $(e, f)$  in  $G$



- ▶  $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶  $h_1(e, f) = \log \sum_d p_G(e, d, f)$
- ▶  $h_2(e, f) = \log p(e)$
- ▶  $h_3(e, f) = |e|$

# Feature Selection

Hidden variable: derivation of  $(e, f)$  in  $G$



- ▶  $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶  $h_1(e, f) = \log \sum_d p_G(e, d, f)$
- ▶  $h_2(e, f) = \log p(e)$
- ▶  $h_3(e, f) = |e|$
- ▶  $h_4(e, f) = \log \sum_a p(f, a | e)$   
alignment-based model

# Feature Selection

Hidden variable: derivation of  $(e, f)$  in  $G$



Chiang (2007):

- ▶  $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶  $h_1(e, f) = \log \sum_d p_G(e, d, f)$
- ▶  $h_2(e, f) = \log p(e)$
- ▶  $h_3(e, f) = |e|$
- ▶  $h_4(e, f) = \log \sum_a p(f, a | e)$   
alignment-based model

- ▶  $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶  $h_1(e, d, f) = \log p_G(e, d | f)$
- ▶  $h_2(e, d, f) = \log p_G(f, d | e)$
- ▶  $h_3(e, d, f) = \log p(e)$
- ▶  $h_4(e, d, f) = |e|$

# Feature Selection

Hidden variable: derivation of  $(e, f)$  in  $G$



Chiang (2007):

- ▶  $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶  $h_1(e, f) = \log \sum_d p_G(e, d, f)$
- ▶  $h_2(e, f) = \log p(e)$
- ▶  $h_3(e, f) = |e|$
- ▶  $h_4(e, f) = \log \sum_a p(f, a | e)$   
alignment-based model

- ▶  $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶  $h_1(e, d, f) = \log p_G(e, d | f)$
- ▶  $h_2(e, d, f) = \log p_G(f, d | e)$
- ▶  $h_3(e, d, f) = \log p(e)$
- ▶  $h_4(e, d, f) = |e|$
- ▶  $h_5(e, d, f) = \text{lexical weights}$

# Feature Selection

Hidden variable: derivation of  $(e, f)$  in  $G$



Chiang (2007):

- ▶  $\mathcal{M} \subseteq \mathcal{M}(E | F)$
- ▶  $h_1(e, f) = \log \sum_d p_G(e, d, f)$
- ▶  $h_2(e, f) = \log p(e)$
- ▶  $h_3(e, f) = |e|$
- ▶  $h_4(e, f) = \log \sum_a p(f, a | e)$   
alignment-based model

- ▶  $\mathcal{M} \subseteq \mathcal{M}(E, D | F)$
- ▶  $h_1(e, d, f) = \log p_G(e, d | f)$
- ▶  $h_2(e, d, f) = \log p_G(f, d | e)$
- ▶  $h_3(e, d, f) = \log p(e)$
- ▶  $h_4(e, d, f) = |e|$
- ▶  $h_5(e, d, f) = \text{lexical weights}$
- ▶  $h_6(e, d, f) = \#_{\text{ex}} d$
- ▶  $h_7(e, d, f) = \#_{\text{gl}} d$

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

**Decoding**

Further Grammar Classes

## My Dissertation

Key Aspects

Publications



# Decoding

$$\hat{e} = \operatorname{argmax}_e p(e | f)$$

# Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f)\end{aligned}$$

# Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)}\end{aligned}$$

# Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

# Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

**NP-hard!** (Casacuberta and Higuera 2000;  
Li, Eisner, and Khudanpur 2009)

# Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

NP-hard! (Casacuberta and Higuera 2000;  
Li, Eisner, and Khudanpur 2009)

$$\approx \operatorname{argmax}_e \sum_{d \in D_{\text{fin}}} \exp \sum_i \lambda_i \cdot h_i(e, d, f)$$

# Decoding

$$\begin{aligned}\hat{e} &= \operatorname{argmax}_e p(e | f) \\ &= \operatorname{argmax}_e \sum_d p(e, d | f) \\ &= \operatorname{argmax}_e \sum_d \frac{\exp \sum_i \lambda_i \cdot h_i(e, d, f)}{\sum_{e', d'} \exp \sum_i \lambda_i \cdot h_i(e', d', f)} \\ &= \operatorname{argmax}_e \sum_d \exp \sum_i \lambda_i \cdot h_i(e, d, f)\end{aligned}$$

NP-hard! (Casacuberta and Higuera 2000;  
Li, Eisner, and Khudanpur 2009)

$$\approx \operatorname{argmax}_e \sum_{d \in D_{\text{fin}}} \exp \sum_i \lambda_i \cdot h_i(e, d, f)$$

where  $D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$

## Computing $D_{\text{fin}}$

$$D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$$



## Computing $D_{\text{fin}}$

$$D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$$

$\max\{ 0.9, 0.99, 0.999, \dots \} ?$

## Computing $D_{\text{fin}}$

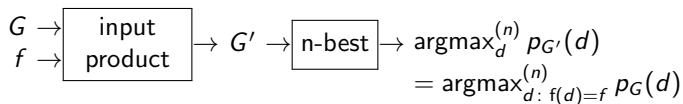
$$D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$$

$$G' \rightarrow \boxed{\text{n-best}} \rightarrow \operatorname{argmax}_d^{(n)} p_{G'}(d)$$

n-best  $\approx O(|R'| + |N'| \cdot n \cdot \log n)$  (Huang and Chiang 2005)

## Computing $D_{\text{fin}}$

$$D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$$

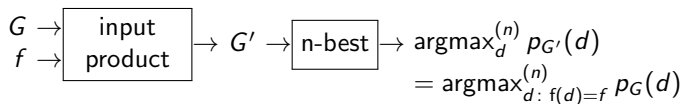


n-best  $\approx O(|R'| + |N'| \cdot n \cdot \log n)$  (Huang and Chiang 2005)

input product  $\approx O(|R| \cdot |f|^c)$

## Computing $D_{\text{fin}}$

$$D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$$



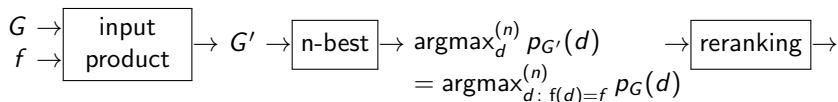
n-best  $\approx O(|R'| + |N'| \cdot n \cdot \log n)$  (Huang and Chiang 2005)

input product  $\approx O(|R| \cdot |f|^c)$

$\rightsquigarrow$  Cube Pruning (Chiang 2007)

## Computing $D_{\text{fin}}$

$$D_{\text{fin}} = \operatorname{argmax}_d^{(n)} \exp \sum_i \lambda_i \cdot h_i(e(d), d, f)$$



n-best  $\approx O(|R'| + |N'| \cdot n \cdot \log n)$  (Huang and Chiang 2005)

input product  $\approx O(|R| \cdot |f|^c)$

$\rightsquigarrow$  Cube Pruning (Chiang 2007)

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

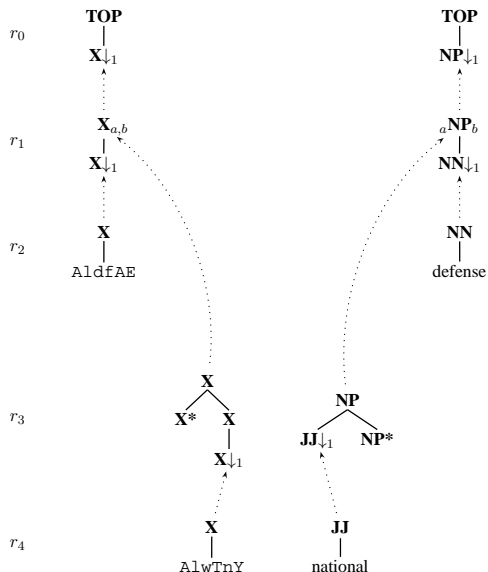
Further Grammar Classes

## My Dissertation

Key Aspects

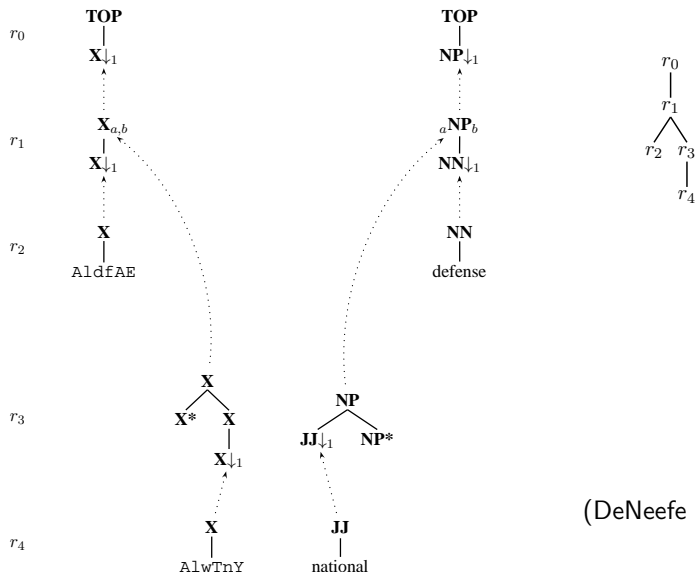
Publications

# Example: Synchronous Tree-Insertion Grammar



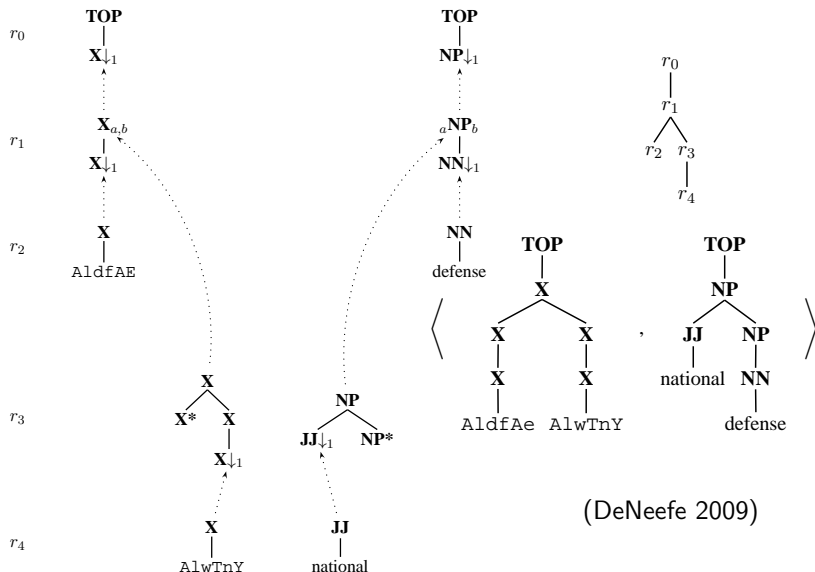
(DeNeefe 2009)

# Example: Synchronous Tree-Insertion Grammar





# Example: Synchronous Tree-Insertion Grammar



# Overview

- ▶ Extended Top-down Tree-to-String Transducers (Galley et al. 2004; Huang, Knight, and Joshi 2006; Graehl, Knight, and May 2008)
- ▶ Extended Multi-Bottom-up Tree Transducers (Engelfriet, Lilin, and Maletti 2008)
- ▶ Synchronous Context-Free Grammars (Chiang 2007)
- ▶ Synchronous Tree-Substitution Grammars (Eisner 2003)
- ▶ Synchronous Tree-Sequence-Substitution Grammars (Zhang et al. 2008)
- ▶ Synchronous Tree-Insertion Grammars (Nesson, Shieber, and Rush 2006; DeNeefe and Knight 2009)

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

Further Grammar Classes

## My Dissertation

Key Aspects

Publications

# Key Aspects

- ▶ Translation Model: Synchronous Tree-Adjoining Grammars
  - ▶ formulation with states and explicit sites
  - ▶ formal properties: closure under input/output product
  - ▶ algorithms for training, decoding

# Key Aspects

- ▶ Translation Model: Synchronous Tree-Adjoining Grammars
  - ▶ formulation with states and explicit sites
  - ▶ formal properties: closure under input/output product
  - ▶ algorithms for training, decoding
- ▶ Auxiliary Algorithms for Decoding
  - ▶ n-best algorithm
  - ▶ determinization of weighted tree automata

# Outline

## Statistical Machine Translation

Introduction

Building an SMT System

Training

Evaluation

Adding Sophistication

## Weighted Grammars for Statistical Machine Translation

Core: Grammar

Feature Selection

Decoding

Further Grammar Classes

## My Dissertation

Key Aspects

Publications

## Publications so far

- ▶ MB, Jonathan May, and Heiko Vogler (2010).  
“Determinization of Weighted Tree Automata using Factorizations”.  
In: Journal of Automata, Languages and Combinatorics. (accepted for publication)
- ▶ MB, Daniel Geisler, Torsten Stüber, and Heiko Vogler (2010).  
“n-Best Parsing Revisited”.  
In: Proceedings of the 2010 Workshop on Applications of Tree Automata in Natural Language Processing, ACL 2010. Uppsala, Sweden, 16 July 2010, pp. 46–54.
- ▶ MB, Mark-Jan Nederhof, and Heiko Vogler (2011).  
“Tree Parsing with Synchronous Tree-Adjoining Grammars”.  
In: Proceedings of IWPT. (accepted for publication)