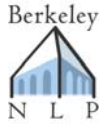


Natural Language Processing



Machine Translation III

Dan Klein – UC Berkeley

Phrase-Based MT

Phrase-Based Translation Overview

Input: lo haré rápidamente | *The decoder... tries different segmentations,*

Translations: I'll do it quickly | *translates phrase by phrase,*
 quickly I'll do it | *and considers reorderings.*

Objective: $\arg \max_e [P(f|e) \cdot P(e)]$

$$\arg \max_e \left[\prod_{(e,f)} P(\bar{f}|\bar{e}) \cdot \prod_{i=1}^{|\bar{e}|} P(e_i|e_{i-1}, e_{i-2}) \right]$$

Phrase-Based Decoding

这 7人 中包括 来自 法国 和 俄罗斯 的 宇航 员 .

| | | | | | | | | | | | |
|-------|---|---------|-----------|------|--------|--------|--------|-----|-----|---------------|---------------|
| the | 7 | people | including | by | some | and | the | the | the | astronauts | . |
| is | 7 | people | included | by | france | and | the | the | the | international | astronautical |
| who | 7 | of | including | the | france | and | the | the | the | frank | of |
| those | 7 | among | including | from | the | france | and | of | the | space | members |
| that | 7 | persons | including | from | the | of | france | and | of | the | astronauts |
| | | | | | | | | | | | |

Decoder design is important: [Koehn et al. 03]

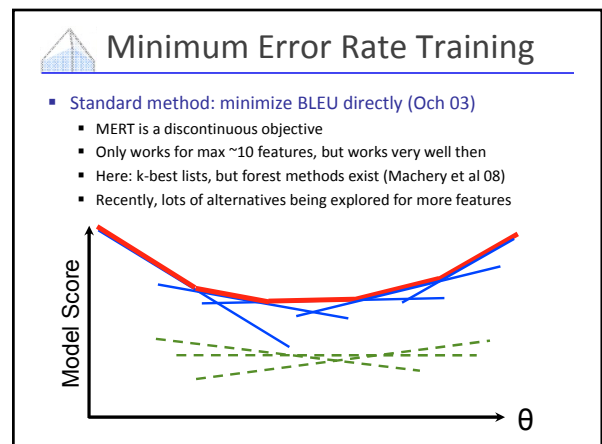
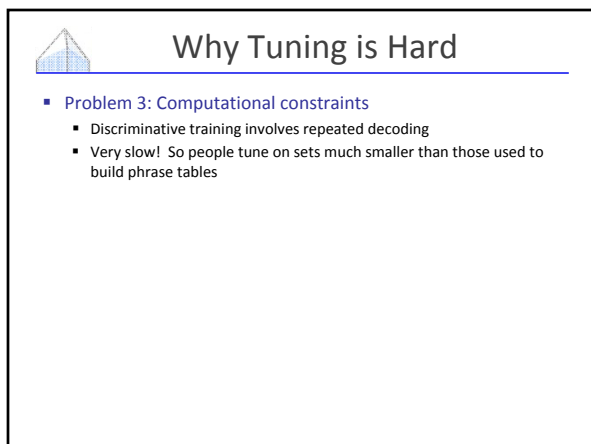
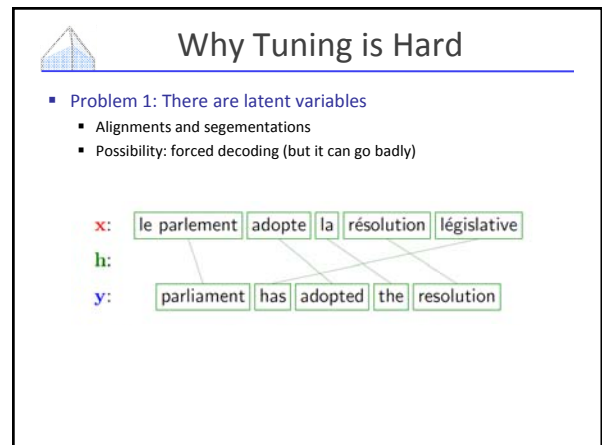
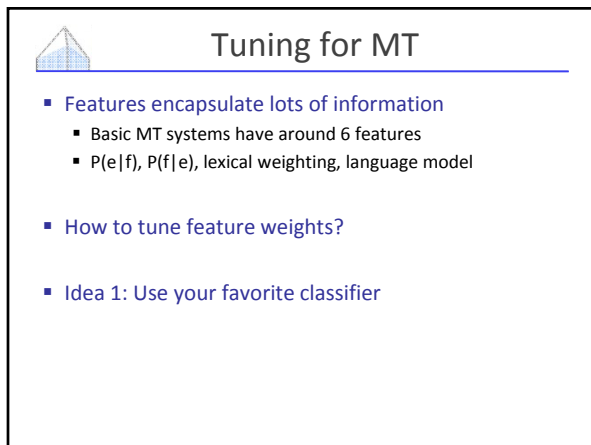
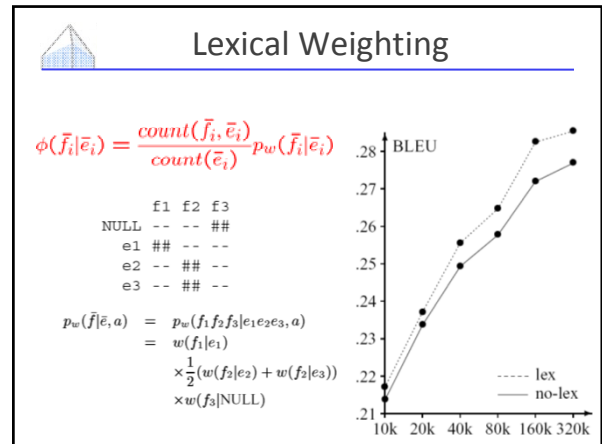
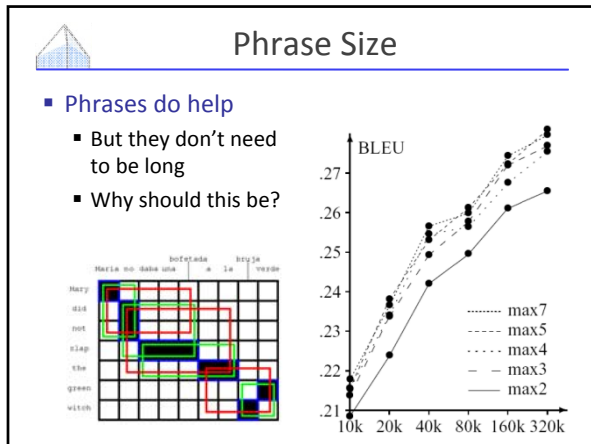
Phrase-Based Decoding

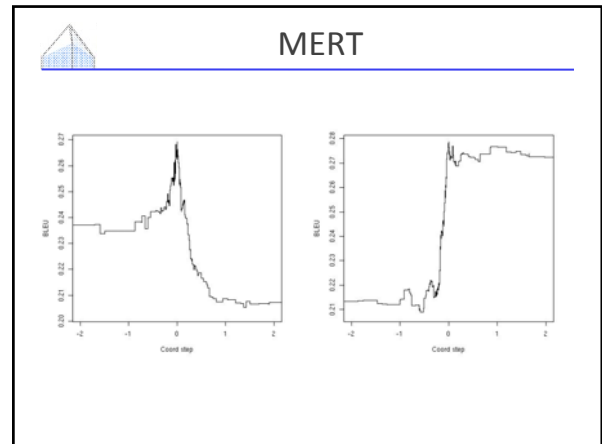
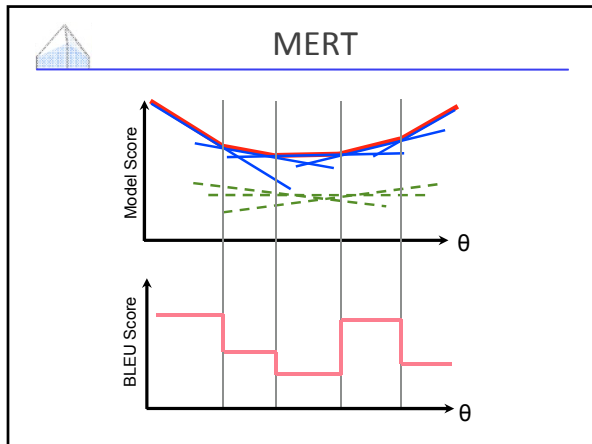
| | | | | | | | | |
|-------|--------------|------|------|----------|-----|-------|-------|-------|
| Maria | no | dio | una | bofetada | a | la | bruja | verde |
| Mary | not | give | a | slap | to | the | witch | green |
| | did not | slap | by | to | the | green | witch | |
| | no | slap | to | the | | | | |
| | did not give | | to | | | | | |
| | | | slap | | to | the | | |
| | | | | slap | | the | witch | |

Monotonic Word Translation

| | | | | | | | | |
|-------|--------------|------|------|----------|-----|-------|-------|-------|
| Maria | no | dio | una | bofetada | a | la | bruja | verde |
| Mary | not | give | a | slap | to | the | witch | green |
| | did not | slap | by | to | the | green | witch | |
| | no | slap | to | the | | | | |
| | did not give | | to | | | | | |
| | | | slap | | to | the | | |
| | | | | slap | | the | witch | |

- Cost is LM * TM
- It's an HMM?
 - $P(e_1, e_2)$
 - $P(f|e)$
- State includes
 - Exposed English
 - Position in foreign
- Dynamic program loop?
 - for (fPosition in 1...f)
 - for (eContext in allEContexts)
 - for (eContext in translations[fPosition])
 - score = scores[fPosition-1][eContext] * LM(eContext+eOption) * TM(eOption, fWord[fPosition])
 - scores[fPosition][eContext[2]+eOption] = max score





Translating with Tree Transducers

Input

lo haré de muy buen grado .

Output

Grammar

Translating with Tree Transducers

Input

lo haré de muy buen grado .

Output

Grammar

ADV \rightarrow < de muy buen grado ; gladly >

Syntactic Models

Translating with Tree Transducers

Input

lo haré de muy buen grado .

Output

ADV
|
gladly

Grammar

ADV \rightarrow < de muy buen grado ; gladly >

University of California Berkeley

Translating with Tree Transducers

| Input | Output |
|---|---|
| lo haré de muy buen grado . | I will do it gladly . |

Grammar

$s \rightarrow \langle \text{lo haré ADV . ; I will do it ADV .} \rangle$
 $ADV \rightarrow \langle \text{de muy buen grado ; gladly} \rangle$

University of California Berkeley

Translating with Tree Transducers

| Input | Output |
|---|--------|
| lo haré de muy buen grado . | |

Grammar

$s \rightarrow \langle \text{lo haré ADV . ; I will do it ADV .} \rangle$
 $ADV \rightarrow \langle \text{de muy buen grado ; gladly} \rangle$

University of California Berkeley

Translating with Tree Transducers

| Input | Output |
|---|--------|
| lo haré de muy buen grado . | |

Grammar

$s \rightarrow \langle \text{lo haré ADV . ; I will do it ADV .} \rangle$
 $ADV \rightarrow \langle \text{de muy buen grado ; gladly} \rangle$

University of California Berkeley

Translating with Tree Transducers

| Input | Output |
|---|---|
| lo haré de muy buen grado . | I will do it gladly . |

Grammar

$s \rightarrow \langle \text{lo haré ADV . ; I will do it ADV .} \rangle$
 $ADV \rightarrow \langle \text{de muy buen grado ; gladly} \rangle$

University of California Berkeley

Translating with Tree Transducers

| Input | Output |
|---|---|
| lo haré de muy buen grado . | I will do it gladly . |

Grammar

$VP \rightarrow \langle \text{lo haré ADV ; will do it ADV} \rangle$
 $s \rightarrow \langle \text{lo haré ADV . ; I will do it ADV .} \rangle$
 $ADV \rightarrow \langle \text{de muy buen grado ; gladly} \rangle$

University of California Berkeley

Translating with Tree Transducers

| Input | Output |
|---|--------|
| lo haré de muy buen grado . | |

Grammar

$VP \rightarrow \langle \text{lo haré ADV ; will do it ADV} \rangle$
 $s \rightarrow \langle \text{lo haré ADV . ; I will do it ADV .} \rangle$
 $ADV \rightarrow \langle \text{de muy buen grado ; gladly} \rangle$

Learning Grammars for Translation

Grammar Rules

~~(haré ; will-do)~~

VP →
(lo haré de ... grado ;
will do it gladly)

VP →
(lo haré ADV ;
will do it ADV)

The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

Rules matching an example 40-word sentence

| Size of the source-side yield | Rule Count |
|-------------------------------|------------|
| 1 | ~10,000 |
| 2 | ~15,000 |
| 3 | ~40,000 |
| 4 | ~65,000 |
| 5 | ~85,000 |
| 6 | ~70,000 |
| 7 | ~55,000 |
| 8 | ~40,000 |
| 9 | ~25,000 |
| 10+ | ~10,000 |

The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

Rules matching an example 40-word sentence

| Size of the source-side yield | Rule Count |
|-------------------------------|------------|
| 1 | ~10,000 |
| 2 | ~15,000 |
| 3 | ~40,000 |
| 4 | ~65,000 |
| 5 | ~85,000 |
| 6 | ~70,000 |
| 7 | ~55,000 |
| 8 | ~40,000 |
| 9 | ~25,000 |
| 10+ | ~10,000 |

$s \rightarrow NP VP ; NP VP$

The Size of Tree Transducer Grammars

Extracted a transducer grammar from a 220 million word bitext

Relativized the grammar to each test sentence

Kept all rules with at most 6 non-terminals

Rules matching an example 40-word sentence

| Size of the source-side yield | Rule Count |
|-------------------------------|------------|
| 1 | ~10,000 |
| 2 | ~15,000 |
| 3 | ~40,000 |
| 4 | ~65,000 |
| 5 | ~85,000 |
| 6 | ~70,000 |
| 7 | ~55,000 |
| 8 | ~40,000 |
| 9 | ~25,000 |
| 10+ | ~10,000 |

$s \rightarrow NP VP ; NP VP$

$s \rightarrow NP no es ni ADJP ni ADJP . ; NP isn't ADJP or ADJP .$

Syntactic Decoding

University of California Berkeley

Tree Transducer Grammars

S
 No se olvide de subir un canto rodado en Colorado

Synchronous Grammar

NNP → Colorado ; Colorado
NN → canto rodado ; boulder
S → No se olvide de subir un **NN** en **NNP** ; Don't forget to climb a **NN** in **NNP**

Output

S
 Don't forget to climb a boulder in Colorado

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length:

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]:

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

Binary rule: $X \rightarrow Y Z$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

Binary rule: $X \rightarrow Y Z$
 Split points: $i < k < j$
 Operations: $O(j - i)$
 Time scales with: Grammar constant

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

$S \rightarrow$ No se **VB** de subir un **NN** en **NNP**

$_i$ No se olvide de subir un canto rodado en Colorado $_j$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

$S \rightarrow$ No se **VB** de subir un **NN** en **NNP**

$_i$ No se olvide de subir un canto rodado en Colorado $_j$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

$S \rightarrow$ No se **VB** de subir un **NN** en **NNP**

$_i$ No se olvide de subir un canto rodado en Colorado $_j$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

$S \rightarrow$ No se **VB** de subir un **NN** en **NNP**

$_i$ No se olvide de subir un canto rodado en Colorado $_j$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

$S \rightarrow$ No se **VB** de subir un **NN** en **NNP**

$_i$ No se olvide de subir un canto rodado en Colorado $_j$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span [i,j]: Apply all grammar rules to [i,j]

$S \rightarrow$ No se **VB** de subir un **NN** en **NNP**

$_i$ No se olvide de subir un canto rodado en Colorado $_j$

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span $[i,j]$: Apply all grammar rules to $[i,j]$

$S \rightarrow$ No se VB de subir un NN en NNP

i No se olvide de subir un canto rodado en Colorado j

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span $[i,j]$: Apply all grammar rules to $[i,j]$

$S \rightarrow$ No se VB de subir un NN en NNP

i No se olvide de subir un canto rodado en Colorado j

Many untransformed lexical rules can be applied in linear time

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span $[i,j]$: Apply all grammar rules to $[i,j]$

$S \rightarrow$ No se VP NP PP

i No se olvide de subir un canto rodado en Colorado j

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span $[i,j]$: Apply all grammar rules to $[i,j]$

$S \rightarrow$ No se VP NP PP

i No se olvide de subir un canto rodado en Colorado j

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span $[i,j]$: Apply all grammar rules to $[i,j]$

$S \rightarrow$ No se VP NP PP

i No se olvide de subir un canto rodado en Colorado j

University of California Berkeley

CKY-style Bottom-up Parsing

For each span length: For each span $[i,j]$: Apply all grammar rules to $[i,j]$

$S \rightarrow$ No se VP NP PP

i No se olvide de subir un canto rodado en Colorado j

Problem: Applying adjacent non-terminals is slow



Eliminating Non-terminal Sequences

Lexical Normal Form (LNF)

- (a) lexical rules have at most one adjacent non-terminal
- (b) all unlexicalized rules are binary.

Original rule: $S \rightarrow \text{No se } VB \text{ } VB \text{ un } NN \text{ } PP$

Transformed rules: $S \rightarrow \text{No se } VB\sim VB \text{ un } NN\sim PP$
 $VB\sim VB \rightarrow VB \text{ } VB$
 $NN\sim PP \rightarrow NN \text{ } PP$

- Parsing stages:
- Lexical rules are applied by matching
 - Unlexicalized rules are applied by iterating over split points



Speeding up Lexical Rule Application

Problem: Lexical rules can apply to many spans

$S \rightarrow \text{No se olvide de subir } NP$



Speeding up Lexical Rule Application

Problem: Lexical rules can apply to many spans

$S \rightarrow \text{No se olvide de subir } NP$



Speeding up Lexical Rule Application

Problem: Lexical rules can apply to many spans

$S \rightarrow \text{No se olvide de subir } NP$



Speeding up Lexical Rule Application

Problem: Lexical rules can apply to many spans

$S \rightarrow \text{No se olvide de subir } NP$



Flexible Syntax

Soft Syntactic MT: From Chiang 2010



reference: An official from Japan's science and technology ministry said, "We are highly encouraged by Abraham's comment."

Hiero: Officials of the Japanese ministry of education and science, "said Abraham speeches, we are deeply encouraged by."

string-to-tree: Japan's ministry of education, culture, sports, science and technology, "Abraham's statement, which is most encouraging," the official said.

Previous work

| | | | |
|--|------------------|-----------------------------------|---|
| | string-to-string | ITG (Wu 1997) | Hiero (Chiang 2005) |
| | string-to-tree | Yamada & Knight 2001 | Galley et al 2004/2006 |
| | tree-to-string | | Huang et al 2006 Y Liu et al 2006 |
| | tree-to-tree | DOT (Poutsma 2000) Eisner 2003 | Stat-XFER (Lavie et al 2008) M Zhang et al. 2008 Y Liu et al., 2009 |



Hiero Rules

- $S \rightarrow \langle S_{\square} X_{\square}, S_{\square} X_{\square} \rangle$
- $S \rightarrow \langle X_{\square}, X_{\square} \rangle$
- $X \rightarrow \langle \text{yu } X_{\square} \text{ you } X_{\square}, \text{have } X_{\square} \text{ with } X_{\square} \rangle$
- $X \rightarrow \langle X_{\square} \text{ de } X_{\square}, \text{the } X_{\square} \text{ that } X_{\square} \rangle$
- $X \rightarrow \langle X_{\square} \text{ zhiyi, one of } X_{\square} \rangle$
- $X \rightarrow \langle \text{Aozhou, Australia} \rangle$
- $X \rightarrow \langle \text{shi, is} \rangle$
- $X \rightarrow \langle \text{shaoshu guojia, few countries} \rangle$
- $X \rightarrow \langle \text{bangjiao, diplomatic relations} \rangle$
- $X \rightarrow \langle \text{Bei Han, North Korea} \rangle$

From [Chiang et al. 2005]

STSG extraction

1. Phrases

- * respect word alignments
- * are syntactic constituents on both sides

2. Phrase pairs form rules

3. Subtract phrases to form rules



STSG extraction

1. Phrases

- * respect word alignments
- * are syntactic constituents on both sides

2. Phrase pairs form rules

3. Subtract phrases to form rules



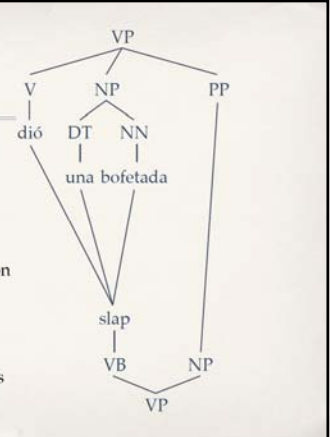
STSG extraction

1. Phrases

- * respect word alignments
- * are syntactic constituents on both sides

2. Phrase pairs form rules

3. Subtract phrases to form rules

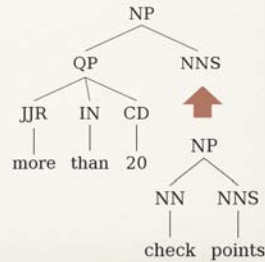


Why is tree-to-tree hard?

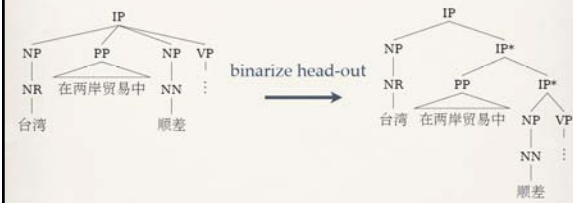
too few rules



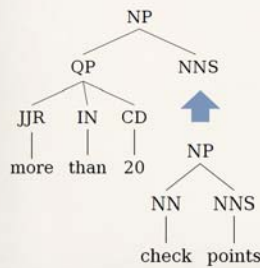
too few derivations



Extracting more rules

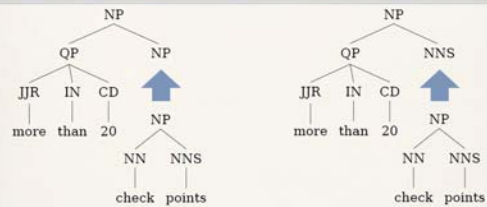


Allow more derivations



- ♦ STSG: allow only matching substitutions
- ♦ Hiero-like: allow any substitutions
- ♦ Let the model learn to choose:
 - ♦ matching substitutions
 - ♦ mismatching substitutions
 - ♦ monotone phrase-based

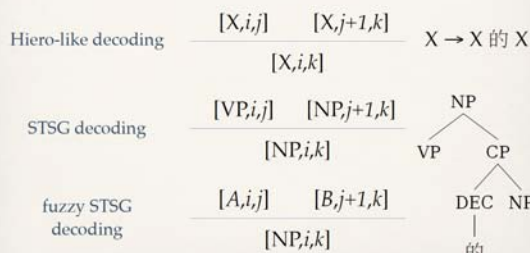
Allow more derivations



fire subst:NP→NP
fire subst:match

fire subst:NNS→NP
fire subst:unmatch

Allow more derivations



Results

| extraction | Chinese-English | | | Arabic-English | | |
|-------------------------|-----------------|-------|------|----------------|-------|------|
| | rules | feats | BLEU | rules | feats | BLEU |
| Hiero | 440M | 1k | 23.7 | 790M | 1k | 48.9 |
| fuzzy STSG | 50M | 5k | 23.9 | 38M | 5k | 47.5 |
| fuzzy STSG +binarize | 64M | 5k | 24.3 | 40M | 6k | 48.1 |
| fuzzy STSG +SAMT | 440M | 160k | 24.3 | 790M | 130k | 49.7 |

