

# Applications of Tree Automata Theory

## Lecture IV: Machine Translation — Basics

Andreas Maletti

Institute of Computer Science  
Universität Leipzig, Germany

*on leave from:* Institute for Natural Language Processing  
Universität Stuttgart, Germany

`maletti@ims.uni-stuttgart.de`

Yekaterinburg — August 24, 2014

# Roadmap

- 1 Theory of Tree Automata
- 2 Parsing — Basics and Evaluation
- 3 Parsing — Advanced Topics
- 4 Machine Translation — Basics and Evaluation
- 5 Theory of Tree Transducers
- 6 Machine Translation — Advanced Topics

Always ask questions right away!

Foundations

# Statistical Machine Translation

## Definition

A **statistical machine translation system** is a (usually fully automatic) computer system that translates based on **statistical models** learnt from **parallel corpora**.

# Parallel Corpus

## Definition

A **parallel corpus** is a (linguistic) resource, in which text in one language is presented together with sentence-by-sentence translations into another language

## Notes

- parallel corpus = sentence-aligned bi-text
- parallel corpus  $\neq$  comparable corpus

## Definition (Sentence alignment)

A **sentence alignment** is an injective (partial) mapping  
 $f: \mathbb{N} \rightarrow \mathbb{N}$  (between sentence numbers)

# Parallel Corpus

## Definition (Sentence alignment)

A **sentence alignment** is an injective (partial) mapping  
 $f: \mathbb{N} \rightarrow \mathbb{N}$  (between sentence numbers)

## Definition (Parallel corpus [formal])

A **parallel corpus** consists of three (partial) mappings:

- $s: \mathbb{N} \rightarrow \Sigma^*$  (source sentences)
- sentence alignment  $f: \mathbb{N} \rightarrow \mathbb{N}$
- $t: \mathbb{N} \rightarrow \Delta^*$  (target sentences)

## Example (Aligned sentences)

### 1 nice example:

- “We can help countries catch up, but not by putting their neighbors on hold”
- “Wir können Ländern beim Aufholen helfen, aber nicht, indem wir ihre Nachbarn in den Wartesaal schicken”



## Example (Aligned sentences)

### 1 nice example:

- “We can help countries catch up, but not by putting their neighbors on hold”
- “Wir können Ländern beim Aufholen helfen, aber nicht, indem wir ihre Nachbarn in den Wartesaal schicken”

### 2 questionable example:

- “We must bear in mind the Community as a whole”
- “Wir müssen uns davor hüten, alles vergemeinschaften zu wollen”

## English-Russian example

- “Indeed, Republican lawyers identified only 300 cases of electoral fraud in the United States in a decade.”
- “К тому же юристы республиканцев насчитали только 300 случаев электоральных фальсификаций в Соединенных Штатах за десять лет.”

# Sentences

## Definition

A **sentence** is a sequence of **tokens**

(the result of tokenization)

# Sentences

## Definition

A **sentence** is a sequence of **tokens**

(the result of tokenization)

## Notes

- normally the atomic unit is the word (dictionary entry)
- **notable exceptions**: Japanese, Chinese, etc.

# Sentences

## Example

- Sentence:

We must bear in mind the Community as a whole

- Token sequence: length: 10

We	must	bear	in	mind	the	Community	as	a	whole
----	------	------	----	------	-----	-----------	----	---	-------

## Definition

Given a sentence  $S$  of length  $n$ , we write

- $\text{pos}(S) = \{1, \dots, n\}$  positions in  $S$
- $S[i]$  with  $i \in \text{pos}(S)$  for its  $i$ -th token
- $S^{-1}[w] = \{i \in \text{pos}(S) \mid S[i] = w\}$  its inverse
- $S[i, j]$  with  $i, j \in \text{pos}(S)$ ,  $i \leq j$  for its  $[i, j]$ -span  
$$S[i, j] = S[i]S[i + 1] \cdots S[j]$$

## Example

$S =$  We must bear in mind the Community as a whole

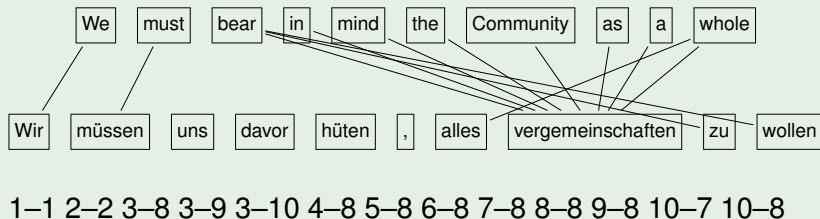
- $\text{pos}(S) = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$
- $S[4] = \text{in}$
- $S^{-1}[\text{must}] = \{2\}$
- $S[6, 8] = \text{the Community as}$

# Word Alignment

## Definition

A **word alignment** for two sentences  $S_1$  and  $S_2$  is a relation  $\rho \subseteq \text{pos}(S_1) \times \text{pos}(S_2)$  on its positions

## Example





# Word Alignment

## Nicer example



## English-Russian example



1-1 2-2 3-3 3-4 4-5 5-5 9-7 7-8 10-9 10-10

## EUROPARL German-English parallel corpus

- 1,920,209 parallel sentences
- 44,548,491 words in German
- 47,818,827 words in English
- sentence-aligned, but not word-aligned
- from parliament proceedings

# Parallel Corpus

## MULTIUN Chinese-English parallel corpus

- 9,564,315 parallel sentences
- 256,720,000 words in English
- sentence-aligned, but not word-aligned
- from official UN documents

## YANDEX parallel corpus

- $\approx$  1M parallel sentences
- Russian-English

Phrase-based Models

# Phrase-based Models

## Key points [OCH, NEY, 2004]

- **phrase** as basic translation unit
- dominant model for many language pairs  
incl. Russian-English
- can deal with (arbitrary) many-to-many alignments
- words internally move together

# Phrase-based Models

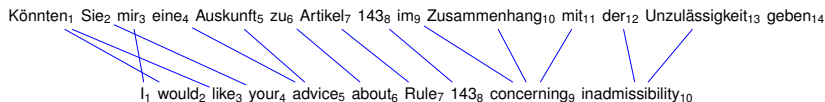
## Key points [OCH, NEY, 2004]

- **phrase** as basic translation unit
- dominant model for many language pairs  
incl. Russian-English
- can deal with (arbitrary) many-to-many alignments
- words internally move together

## Definition (Phrase)

**phrase** is a contiguous sequence of words  
(usually given by span  $[i, i']$ )

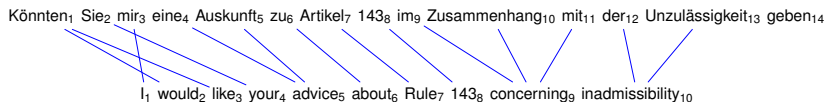
# Phrase-based Models



## Algorithm

- 1 phrase pair  $([j, j'], [i, i'])$  **consistently aligned** if
  - $\ell' \in [i, i']$  for all  $\ell \in [j, j']$  and  $(\ell, \ell') \in A$
  - $\ell \in [j, j']$  for all  $\ell' \in [i, i']$  and  $(\ell, \ell') \in A$

# Phrase-based Models

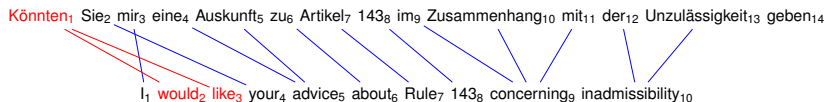


## Algorithm

- 1 phrase pair  $([j, j'], [i, i'])$  **consistently aligned** if
  - $\ell' \in [i, i']$  for all  $\ell \in [j, j']$  and  $(\ell, \ell') \in A$
  - $\ell \in [j, j']$  for all  $\ell' \in [i, i']$  and  $(\ell, \ell') \in A$
- 2 extract all consistently aligned phrase pairs
- 3 (restrict length of phrases based on corpus size)



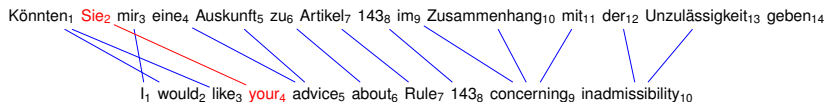
# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

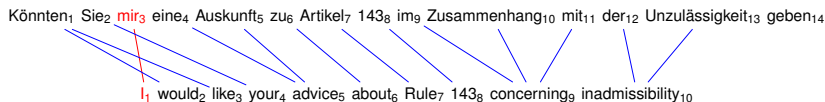
# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

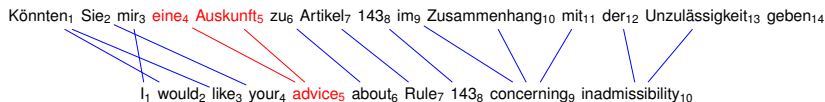
# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

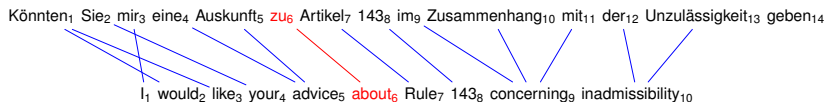
# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

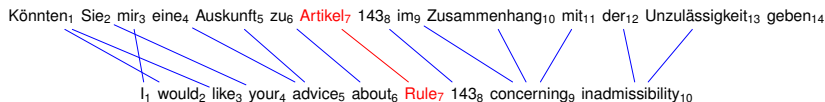
# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

# Phrase-based Models



Formally:

$([1,1], [2,3])$	$([2,2], [4,4])$	$([3,3], [1,1])$
$([4,5], [5,5])$	$([6,6], [6,6])$	$([7,7], [7,7])$
$([8,8], [8,8])$	$([9,11], [9,9])$	$([12,13], [10,10])$

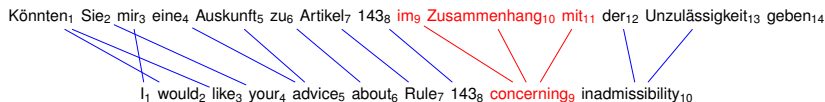
# Phrase-based Models



Formally:

$([1,1], [2,3])$	$([2,2], [4,4])$	$([3,3], [1,1])$
$([4,5], [5,5])$	$([6,6], [6,6])$	$([7,7], [7,7])$
$([8,8], [8,8])$	$([9,11], [9,9])$	$([12,13], [10,10])$

# Phrase-based Models

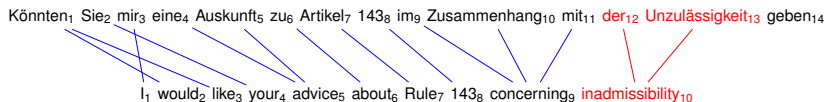


Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])



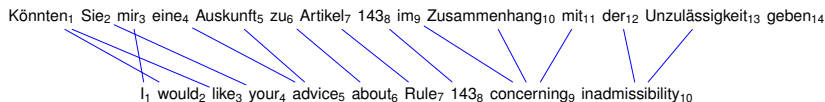
# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

# Phrase-based Models



Formally:

([1,1], [2,3])	([2,2], [4,4])	([3,3], [1,1])
([4,5], [5,5])	([6,6], [6,6])	([7,7], [7,7])
([8,8], [8,8])	([9,11], [9,9])	([12,13], [10,10])

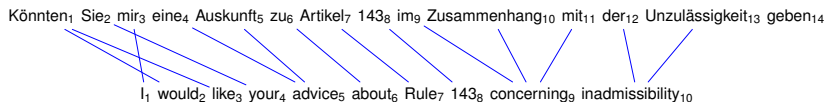
For better readability:

Könnten — would like  
eine Auskunft — advice  
143 — 143

Sie — your  
zu — about  
im Zusammenhang mit — concerning

mir — I  
Artikel — Rule  
der Unzulässigkeit — inadmissibility

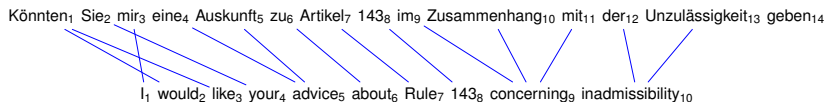
# Phrase-based Models



## Notes

- these were only **minimal** phrase pairs
- extract all (sensible) combinations of these
- e.g.,  $([1, 1], [2, 3])$  and  $([2, 2], [4, 4])$  yield  $([1, 2], [2, 4])$

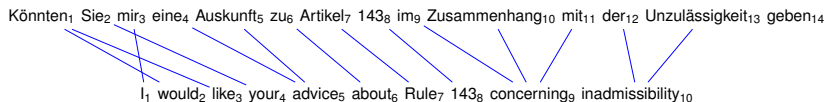
# Phrase-based Models



## Notes

- these were only **minimal** phrase pairs
- extract all (sensible) combinations of these
- e.g.,  $([1, 1], [2, 3])$  and  $([2, 2], [4, 4])$  yield  $([1, 2], [2, 4])$
- unaligned words can be added to neighboring phrases
- e.g.,  $([12, 13], [10, 10])$  extends to  $([12, 14], [10, 10])$

# Phrase-based Models

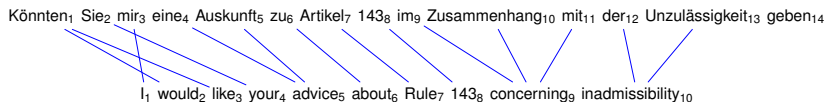


## Notes

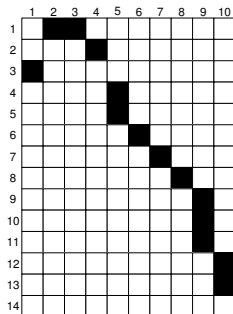
- these were only **minimal** phrase pairs
- extract all (sensible) combinations of these
- e.g.,  $([1, 1], [2, 3])$  and  $([2, 2], [4, 4])$  yield  $([1, 2], [2, 4])$
- unaligned words can be added to neighboring phrases
- e.g.,  $([12, 13], [10, 10])$  extends to  $([12, 14], [10, 10])$

Könnten Sie — would like your      der Unzulässigkeit geben — inadmissibility

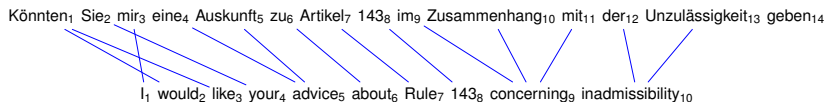
# Phrase-based Models



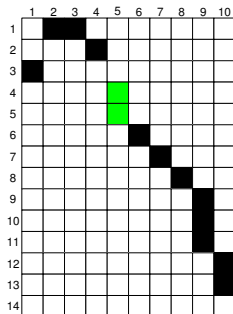
Alternative representation (rectangles):



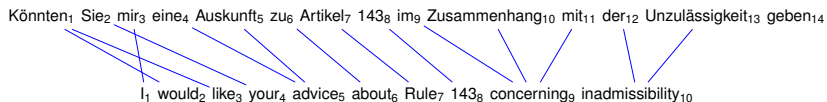
# Phrase-based Models



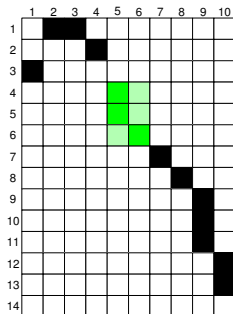
Alternative representation (rectangles):



# Phrase-based Models

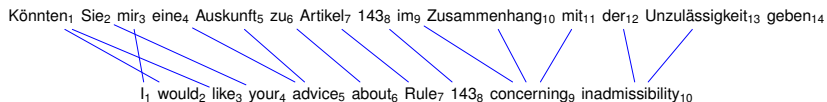


Alternative representation (rectangles):

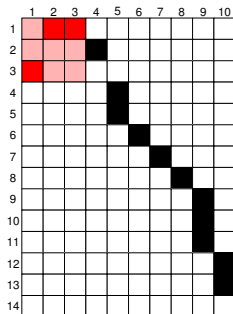




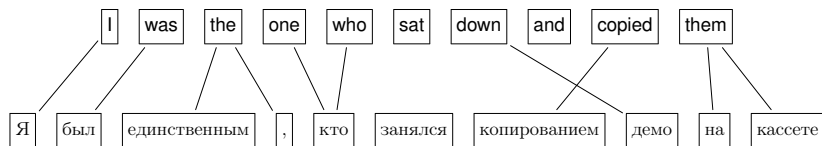
# Phrase-based Models



Alternative representation (rectangles):



# Phrase-based Models



## Extractable phrase pairs

I was — Я был

down and copied them — копированием демо на кассете

the one who sat — единственным, кто занялся

## Notes

- **phrases** are not linguistic phrases!  
(noun phrases, verb phrases, etc.)
- non-linguistic phrase pair: **Sinne der** — **keeping with the**

## Notes

- **phrases** are not linguistic phrases!  
(noun phrases, verb phrases, etc.)
- non-linguistic phrase pair: **Sinne der — keeping with the**
- having only linguistic phrases lowers translation quality

# Phrase-based Models

## Notes

- **phrases** are not linguistic phrases!  
(noun phrases, verb phrases, etc.)
- non-linguistic phrase pair: **Sinne der** — **keeping with the**
- having only linguistic phrases lowers translation quality
- phrase translation table typically huge  
(much larger than parallel corpus)

Hierarchical Phrase-based Models

# Phrase-based Models

Input

Er<sub>1</sub> hat<sub>2</sub> ein<sub>3</sub> neues<sub>4</sub> ,<sub>5</sub> sparsames<sub>6</sub> Auto<sub>7</sub> gekauft<sub>8</sub>

He<sub>1</sub> bought<sub>2</sub> a<sub>3</sub> new<sub>4</sub> fuel-efficient<sub>5</sub> car<sub>6</sub>



# Phrase-based Models

Input

Er<sub>1</sub> hat<sub>2</sub> ein<sub>3</sub> neues<sub>4</sub> ,<sub>5</sub> sparsames<sub>6</sub> Auto<sub>7</sub> gekauft<sub>8</sub>

He<sub>1</sub> bought<sub>2</sub> a<sub>3</sub> new<sub>4</sub> fuel-efficient<sub>5</sub> car<sub>6</sub>



# Phrase-based Models

Input

Er<sub>1</sub> hat<sub>2</sub> ein<sub>3</sub> neues<sub>4</sub> ,<sub>5</sub> sparsames<sub>6</sub> Auto<sub>7</sub> gekauft<sub>8</sub>

He<sub>1</sub> bought<sub>2</sub> a<sub>3</sub> new<sub>4</sub> fuel-efficient<sub>5</sub> car<sub>6</sub>



# Phrase-based Models

Input

Er<sub>1</sub> hat<sub>2</sub> ein<sub>3</sub> neues<sub>4</sub> ,<sub>5</sub> sparsames<sub>6</sub> Auto<sub>7</sub> gekauft<sub>8</sub>  
He<sub>1</sub> bought<sub>2</sub> a<sub>3</sub> new<sub>4</sub> fuel-efficient<sub>5</sub> car<sub>6</sub>

Problem

Only rule translating **hat** or **gekauft** has very limited use:

*hat ein neues , sparsames Auto gekauft*  
— *bought a new fuel-efficient car*

# Hierarchical Phrase-based Models

## Analysis

- the restriction to phrases (i.e., contiguous segments) yields undesirably long phrase pairs

*hat ein neues , sparsames Auto gekauft*

— *bought a new fuel-efficient car*

# Hierarchical Phrase-based Models

## Analysis

- the restriction to phrases (i.e., contiguous segments) yields undesirably long phrase pairs
  - it would be better to allow phrases inside phrases
- **hierarchical phrases** [CHIANG, 2007]

hat *ein neues , sparsames Auto* gekauft

— bought *a new fuel-efficient car*

# Hierarchical Phrase-based Models

## Hierarchical phrases

- phrases with gaps (written as  $X$ )

*hat*  $X$  *gekauft* — *bought*  $X$

# Hierarchical Phrase-based Models

## Hierarchical phrases

- phrases with gaps (written as  $X$ )

*hat*  $X$  *gekauft* — *bought*  $X$

- gaps filled by other hierarchical phrases (maybe gaps)

*ein neues , sparsames Auto* — *a new fuel-efficient car*

*ein*  $X$  *Auto* — *a*  $X$  *car*

# Hierarchical Phrase-based Models

## Hierarchical phrases

- phrases with gaps (written as  $X$ )

*hat*  $X$  *gekauft* — *bought*  $X$

- gaps filled by other hierarchical phrases (maybe gaps)

*ein neues , sparsames Auto* — *a new fuel-efficient car*

*ein*  $X$  *Auto* — *a*  $X$  *car*

- put together we have

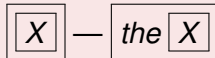
*hat* *ein*  $X$  *Auto* *gekauft* — *bought* *a*  $X$  *car*

# Hierarchical Phrase-based Models

## Restrictions

- must contain lexical item in source side

illegal:





# Hierarchical Phrase-based Models

## Restrictions

- must contain lexical item in source side

illegal:  $\boxed{X}$  — *the*  $\boxed{X}$

- at most 5 gaps (for efficiency: at most 2 gaps)

illegal:  $\boxed{X_1 \text{ hat } X_2 \text{ des } X_3}$  —  $\boxed{X_2 \text{ of } X_3 \text{ was } X_1}$

# Hierarchical Phrase-based Models

## Restrictions

- must contain lexical item in source side

illegal:  $\boxed{X}$  — *the*  $\boxed{X}$

- at most 5 gaps (for efficiency: at most 2 gaps)

illegal:  $\boxed{X_1 \text{ hat } X_2 \text{ des } X_3}$  —  $\boxed{X_2 \text{ of } X_3 \text{ was } X_1}$

- no adjacent gaps in source side

illegal:  $\boxed{ihm \ X_1 \ X_2}$  — *the*  $\boxed{X_2}$  *of*  $\boxed{X_1}$  *to him*

# Hierarchical Phrase-based Models

## Algorithm

- 1 Extract phrase-pairs as usual

*hat ein neues , sparsames Auto gekauft* — *bought a new fuel-efficient car*

*ein neues , sparsames Auto* — *a new fuel-efficient car*

... — ...

# Hierarchical Phrase-based Models

## Algorithm

- 1 Extract phrase-pairs as usual

hat *ein neues , sparsames Auto* gekauft — bought *a new fuel-efficient car*

*ein neues , sparsames Auto* — a new fuel-efficient car

... — ...

- 2 Find phrase-pair that contains another phrase-pair

# Hierarchical Phrase-based Models

## Algorithm

- 1 Extract phrase-pairs as usual

*hat ein neues , sparsames Auto gekauft* — *bought a new fuel-efficient car*

*ein neues , sparsames Auto* — *a new fuel-efficient car*

... — ...

- 2 Find phrase-pair that contains another phrase-pair
- 3 Remove inner phrase-pair and leave new  $X_i$  (if possible)

*hat*  $X_1$  *gekauft* — *bought*  $X_1$

# Hierarchical Phrase-based Models

## Algorithm

- 1 Extract phrase-pairs as usual

*hat ein neues , sparsames Auto gekauft* — *bought a new fuel-efficient car*

*ein neues , sparsames Auto* — *a new fuel-efficient car*

... — ...

- 2 Find phrase-pair that contains another phrase-pair
- 3 Remove inner phrase-pair and leave new  $X_i$  (if possible)

*hat*  $X_1$  *gekauft* — *bought*  $X_1$

- 4 if there were changes, then go to 2

Evaluation

## Translation

- **Input:**

Республиканская стратегия сопротивления повторному избранию обамы



## Translation

- **Input:**

Республиканская стратегия сопротивления повторному избранию обамы

- **Hierarchical phrase-based (YANDEX corpus):**

Republican strategy resistance renewal elect obama

## Translation

- **Input:**

Республиканская стратегия сопротивления повторному избранию обамы

- **Hierarchical phrase-based (YANDEX corpus):**

Republican strategy resistance renewal elect obama

- **GOOGLE Translate:**

The Republican strategy of resistance to the re-election of Obama

# Statistical Machine translation

## Translation

- **Input:**

Республиканская стратегия сопротивления повторному избранию обамы

- **Hierarchical phrase-based (YANDEX corpus):**

Republican strategy resistance renewal elect obama

- **GOOGLE Translate:**

The Republican strategy of resistance to the re-election of Obama

- **YANDEX Translate:**

The Republican strategy of resistance for the re-election of Obama

## Translation

### ■ Input:

Кроме того, эти законы сокращают период досрочного голосования, упраздняют пра во регистрации избирателя в день волеизъявления и отнимают право голоса у граждан, имеющих судимость.

## Translation

### ■ Input:

Кроме того, эти законы сокращают период досрочного голосования, упраздняют пра во регистрации избирателя в день волеизъявления и отнимают право голоса у граждан, имеющих судимость.

### ■ Hierarchical phrase-based (YANDEX corpus):

In addition, these laws reduce period early voting, упраздняют right registering voters in the day of voting and take the citizens have the right to vote, have a criminal record.

## Translation

### ■ Input:

Кроме того, эти законы сокращают период досрочного голосования, упраздняют пра во регистрации избирателя в день волеизъявления и отнимают право голоса у граждан, имеющих судимость.

### ■ Hierarchical phrase-based (YANDEX corpus):

In addition, these laws reduce period early voting, упраздняют right registering voters in the day of voting and take the citizens have the right to vote, have a criminal record.

### ■ GOOGLE Translate:

In addition, these laws reduce the early voting period, abolish the right-in voter registration on the day of expression and the right to take away votes from citizens who have a criminal record.

## Translation

### ■ Input:

Кроме того, эти законы сокращают период досрочного голосования, упраздняют пра во регистрации избирателя в день волеизъявления и отнимают право голоса у граждан, имеющих судимость.

### ■ Hierarchical phrase-based (YANDEX corpus):

In addition, these laws reduce period early voting, упраздняют right registering voters in the day of voting and take the citizens have the right to vote, have a criminal record.

### ■ GOOGLE Translate:

In addition, these laws reduce the early voting period, abolish the right-in voter registration on the day of expression and the right to take away votes from citizens who have a criminal record.

### ■ YANDEX Translate:

In addition, these laws reduce the early voting period will void the law of registration of the voter on the day of will and take away the right to vote of citizens with criminal records.

# BLEU — Bilingual Evaluation Understudy

## Approach [PAPINENI et al., 2002]

- compare translation output  $T$  to reference translation  $R$
- count  $n$ -gram matches (recall)
- normalize by count of all possible  $n$ -gram matches

$$\text{prec}_i = \frac{\text{match}_i(T, R)}{|T| - i + 1}$$



# BLEU — Bilingual Evaluation Understudy

## Approach [PAPINENI et al., 2002]

- compare translation output  $T$  to reference translation  $R$
- count  $n$ -gram matches (recall)
- normalize by count of all possible  $n$ -gram matches

$$\text{prec}_i = \frac{\text{match}_i(T, R)}{|T| - i + 1}$$

- BLEU- $n = \text{brev-pen} \cdot \sqrt[n]{\prod_{i=1}^n \text{prec}_i}$  (geometric mean)

# BLEU — Bilingual Evaluation Understudy

## Approach [PAPINENI et al., 2002]

- compare translation output  $T$  to reference translation  $R$
- count  $n$ -gram matches (recall)
- normalize by count of all possible  $n$ -gram matches

$$\text{prec}_i = \frac{\text{match}_i(T, R)}{|T| - i + 1}$$

- $\text{BLEU-}n = \text{brev-pen} \cdot \sqrt[n]{\prod_{i=1}^n \text{prec}_i}$  (geometric mean)
- adjustment factor **brevity penalty** (punishes short output)

$$\text{brev-pen} = \min\left(1, \frac{|T|}{|R|}\right)$$

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

Translation: (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

Reference: The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
1 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The **rotated** with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
1 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated **with** miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken **with** miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
2 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with **miniature** cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with **miniature** cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
3 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature **cameras** mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature **cameras** attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
4 of 29	of 28	of 27	of 26



# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras **mounted** on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
4 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
4 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on **helmets** recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop **helmets** , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
5 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets **recordings** are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
5 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings **are** checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , **are** monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
6 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are **checked** at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
6 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked **at** the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
6 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at **the** airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto **the** Internet .

1-grams	2-grams	3-grams	4-grams
7 of 29	of 28	of 27	of 26



# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the **airbase** in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
7 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
8 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in **Kandahar** and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in **Kandahar** and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
9 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar **and** then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar **and** then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
10 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and **then** sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and **then** transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
11 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then **sent** to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
11 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
12 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to **London** , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to **London** from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
13 of 29	of 28	of 27	of 26



# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
14 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , **where** they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from **where** they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
15 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where **they** are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where **they** are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
16 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they **are** on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they **are** uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
17 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
17 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on **the** Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto **the** Internet .

1-grams	2-grams	3-grams	4-grams
18 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the **Internet** .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the **Internet** .

1-grams	2-grams	3-grams	4-grams
19 of 29	of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	of 28	of 27	of 26



# Evaluation — BLEU Scoring

Translation: (Google Translate)

The rotated **with miniature** cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

Reference: The images , taken **with miniature** cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	1 of 28	of 27	of 26

# Evaluation — BLEU Scoring

Translation: (Google Translate)

The rotated with **miniature cameras** mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

Reference: The images , taken with **miniature cameras** attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	2 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase **in Kandahar** and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre **in Kandahar** and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	3 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in **Kandahar and** then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in **Kandahar and** then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	4 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar **and then** sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar **and then** transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	5 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	6 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , **where they** are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from **where they** are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	7 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where **they are** on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where **they are** uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	8 of 28	of 27	of 26



# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on **the Internet** .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto **the Internet** .

1-grams	2-grams	3-grams	4-grams
20 of 29	9 of 28	of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the **Internet** .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the **Internet** .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	of 27	of 26

# Evaluation — BLEU Scoring

Translation: (Google Translate)

The rotated **with miniature cameras** mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

Reference: The images , taken **with miniature cameras** attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	1 of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase **in Kandahar and** then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre **in Kandahar and** then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	2 of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in **Kandahar and then** sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in **Kandahar and then** transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	3 of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , **where they are** on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from **where they are** uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	4 of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on **the Internet** .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto **the Internet** .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	5 of 27	of 26

# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase **in Kandahar and then** sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre **in Kandahar and then** transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	5 of 27	1 of 26



# Evaluation — BLEU Scoring

**Translation:** (Google Translate)

The rotated with miniature cameras mounted on helmets recordings are checked at the airbase in Kandahar and then sent to London , where they are on the Internet .

**Reference:** The images , taken with miniature cameras attached to troop helmets , are monitored by a command centre in Kandahar and then transferred to London from where they are uploaded onto the Internet .

1-grams	2-grams	3-grams	4-grams
20 of 29	10 of 28	5 of 27	1 of 26

$$\text{BLEU-4} = \frac{29}{34} \cdot \sqrt[4]{\frac{20 \cdot 10 \cdot 5 \cdot 1}{29 \cdot 28 \cdot 27 \cdot 26}} = 17.46\%$$

# Evaluation — BLEU Scoring

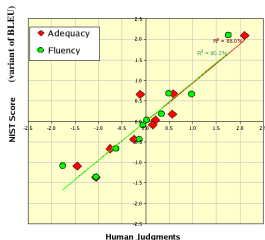
## Notes

- should be used with multiple references
- should be used for documents (not sentences)
- primary evaluation for MT systems
- correlates reasonably well with human judgements

# Evaluation — BLEU Scoring

## Notes

- should be used with multiple references
- should be used for documents (not sentences)
- primary evaluation for MT systems
- correlates reasonably well with human judgements



(Figure from [Koehn, 2010])

## Russian-to-English

System	BLEU-4
WMT '13 winner [PINO et al., 2013] (hierarchical phrase-based)	25.9
hierarchical phrase-based (vanilla)	21.9
MBOT string-to-tree (vanilla)	20.7
string-to-tree (vanilla)	19.8

WMT '13 winner: Hierarchical phrase-based translation with wFSTs.

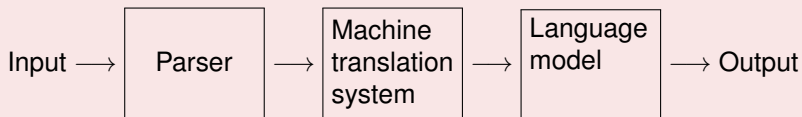
Pre-processing: STANFORD CoreNLP, Morfessor, Stem+POSTag (TreeTagger)

Post-processing: lattices 5-gram-LM rescored and union of lattices rescored via Lattice Minimum Bayes Risk.

Syntax-based Models

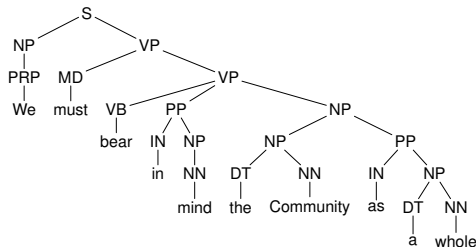
# Syntax-based Machine Translation

## Syntax-based systems



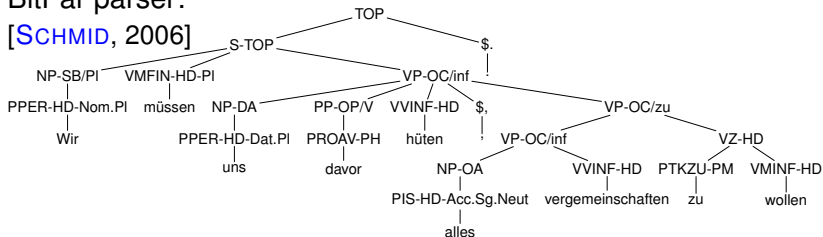
# Syntax-based Machine Translation

CHARNIAK parser: [CHARNIAK, JOHNSON, 2005]



BitPar parser:

[SCHMID, 2006]



# Syntax-based Machine Translation

## Arabic-English

*Yugoslav President Voislav signed for Serbia.*

و تولى التوقيع عن صربيا الرئيس اليوغوسلافي فويسلاف

Translit.: w twlY AltwqyE En SrbyA Alrjys AlywgwslAfy fwyslAf.

*And then the matter was decided, and everything was put in place.*

ف كان ان تم الحسم و وضعت الأمور في نصابها

Translit.: f kAn An tm AlHsm w wDEt Al>mwr fy nSAb hA.

*Below are the male and female winners in the different categories.*

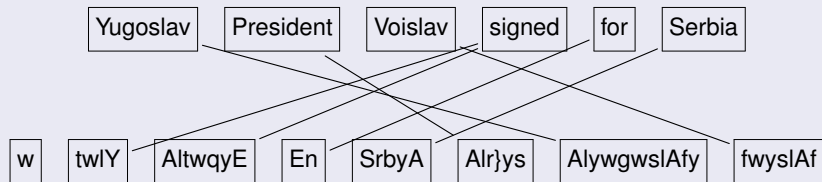
وهنا الأوائل و الأوليات في مختلف الفئات

Translit.: w hnA Al>wAjl w Al>wlyAt fy mxltf AlfjAt.



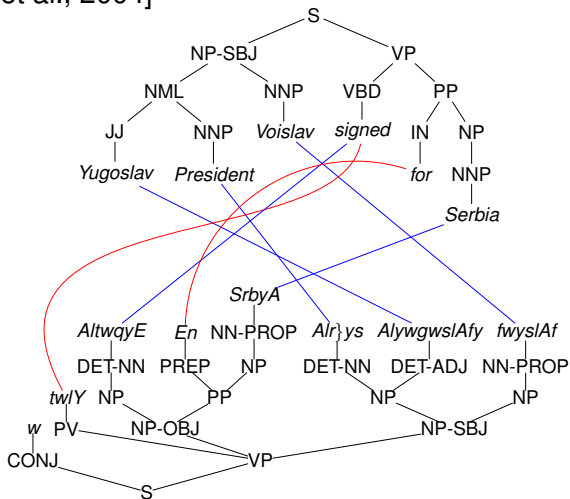
# Syntax-based Machine Translation

## Alignment



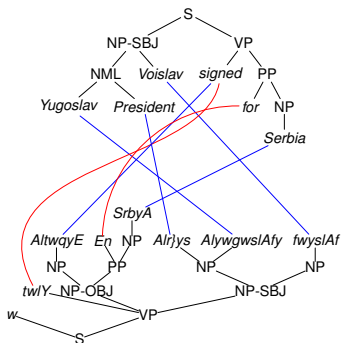
# Syntax-based Machine Translation

[GALLEY et al., 2004]

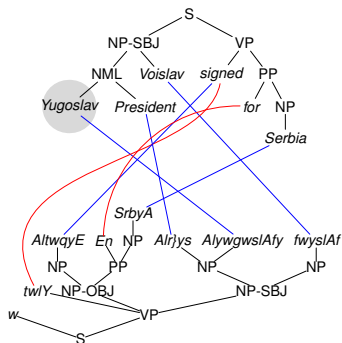


# Syntax-based Machine Translation

- Select next node bottom-up

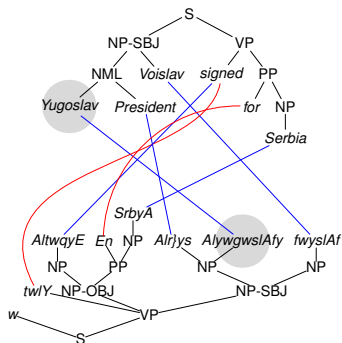


# Syntax-based Machine Translation



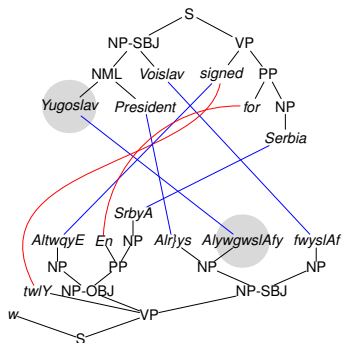
- Select next node bottom-up
- Identify maximal subtree of aligned nodes

# Syntax-based Machine Translation



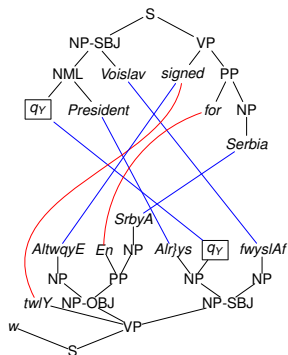
- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.

# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state

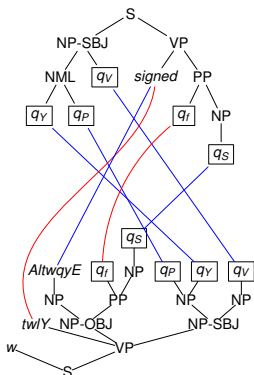
# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

*Yugoslav*  $\frac{q_Y}{AlywgwslAfy}$

# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

*Yugoslav*  $\xrightarrow{q_Y}$  *Alyw gwsl Afy*

*President*  $\xrightarrow{q_P}$  *Alr}ys*

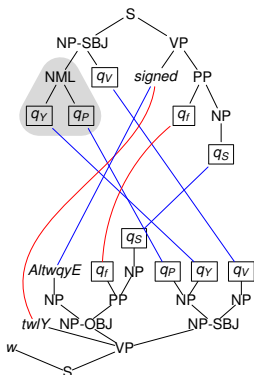
*Voislav*  $\xrightarrow{q_V}$  *fwysl Af*

*for*  $\xrightarrow{q_I}$  *En*

*Serbia*  $\xrightarrow{q_S}$  *SrbyA*



# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leaf state
- Repeat

*Yugoslav*  $\xrightarrow{q_Y}$  *Alyw gwsl Afy*

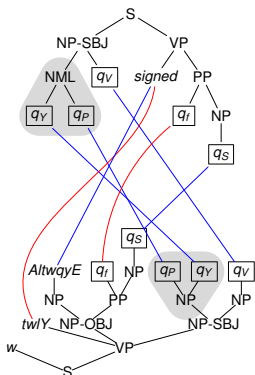
*President*  $\xrightarrow{q_P}$  *Alr}ys*

*Voislav*  $\xrightarrow{q_V}$  *fwysl Af*

*for*  $\xrightarrow{q_I}$  *En*

*Serbia*  $\xrightarrow{q_S}$  *SrbyA*

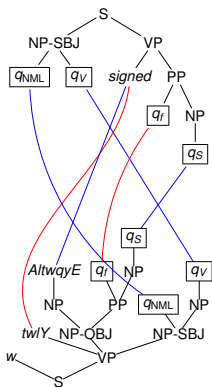
# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

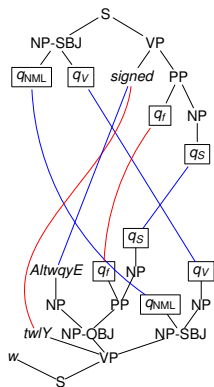
# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

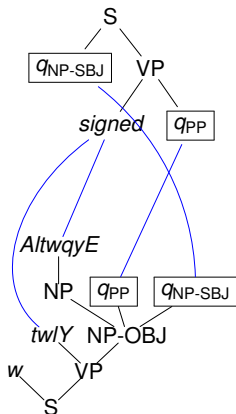
$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

$$\text{NP}(q_S) \xrightarrow{q_{\text{NP}}} \text{NP}(q_S)$$

$$\text{PP}(q_f, q_{\text{NP}}) \xrightarrow{q_{\text{PP}}} \text{PP}(q_f, q_{\text{NP}})$$

$$\text{NP-SBJ}(q_{\text{NML}}, q_V) \xrightarrow{q_{\text{NP-SBJ}}} \text{NP-SBJ}(q_{\text{NML}}, \text{NP}(q_V))$$

# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

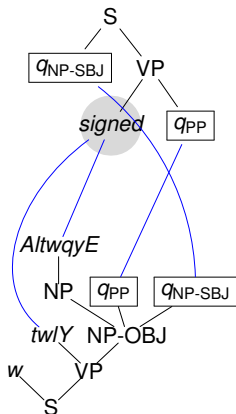
$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

$$\text{NP}(q_S) \xrightarrow{q_{\text{NP}}} \text{NP}(q_S)$$

$$\text{PP}(q_f, q_{NP}) \xrightarrow{q_{\text{PP}}} \text{PP}(q_f, q_{NP})$$

$$\text{NP-SBJ}(q_{\text{NML}}, q_V) \xrightarrow{q_{\text{NP-SBJ}}} \text{NP-SBJ}(q_{\text{NML}}, \text{NP}(q_V))$$

# Syntax-based Machine Translation



$$\text{NP-SBJ}(q_{\text{NML}}, q_V) \xrightarrow{q_{\text{NP-SBJ}}} \text{NP-SBJ}(q_{\text{NML}}, \text{NP}(q_V))$$

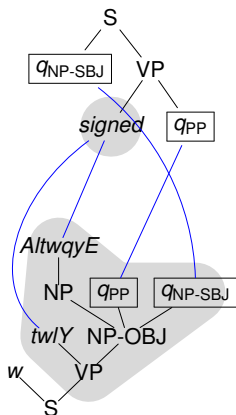
- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

$$\text{NP}(q_S) \xrightarrow{q_{\text{NP}}} \text{NP}(q_S)$$

$$\text{PP}(q_f, q_{\text{NP}}) \xrightarrow{q_{\text{PP}}} \text{PP}(q_f, q_{\text{NP}})$$

# Syntax-based Machine Translation



- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

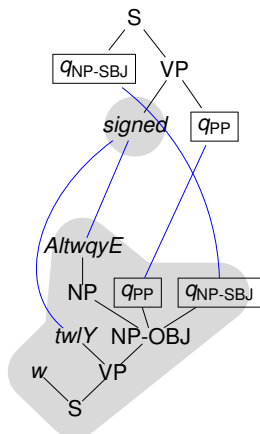
$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

$$\text{NP}(q_S) \xrightarrow{q_{\text{NP}}} \text{NP}(q_S)$$

$$\text{PP}(q_f, q_{\text{NP}}) \xrightarrow{q_{\text{PP}}} \text{PP}(q_f, q_{\text{NP}})$$

$$\text{NP-SBJ}(q_{\text{NML}}, q_V) \xrightarrow{q_{\text{NP-SBJ}}} \text{NP-SBJ}(q_{\text{NML}}, \text{NP}(q_V))$$

# Syntax-based Machine Translation



$$\text{NP-SBJ}(q_{\text{NML}}, q_V) \xrightarrow{q_{\text{NP-SBJ}}} \text{NP-SBJ}(q_{\text{NML}}, \text{NP}(q_V))$$

- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

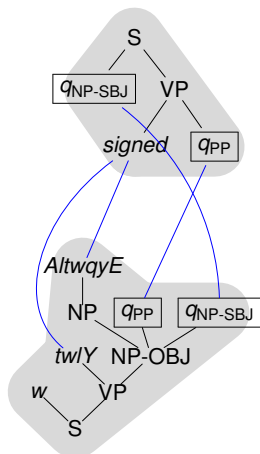
$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

$$\text{NP}(q_S) \xrightarrow{q_{\text{NP}}} \text{NP}(q_S)$$

$$\text{PP}(q_f, q_{\text{NP}}) \xrightarrow{q_{\text{PP}}} \text{PP}(q_f, q_{\text{NP}})$$



# Syntax-based Machine Translation



$$\text{NP-SBJ}(q_{\text{NML}}, q_V) \xrightarrow{q_{\text{NP-SBJ}}} \text{NP-SBJ}(q_{\text{NML}}, \text{NP}(q_V))$$

- Select next node bottom-up
- Identify maximal subtree of aligned nodes
- Identify subtree of nodes aligned to aligned nodes, etc.
- Extract rule and leave state
- Repeat

$$\text{NML}(q_Y, q_P) \xrightarrow{q_{\text{NML}}} \text{NP}(q_P, q_Y)$$

$$\text{NP}(q_S) \xrightarrow{q_{\text{NP}}} \text{NP}(q_S)$$

$$\text{PP}(q_f, q_{\text{NP}}) \xrightarrow{q_{\text{PP}}} \text{PP}(q_f, q_{\text{NP}})$$

# Syntax-based Machine Translation

## Rules

*Yugoslav*  $\frac{q_Y}{}$  *AlywgwslAfy*

*Voislav*  $\frac{q_V}{}$  *fwyslAf*

*Serbia*  $\frac{q_S}{}$  *SrbyA*

$NML(q_Y, q_P) \frac{q_{NML}}{}$   $NP(q_P, q_Y)$

$PP(q_f, q_{NP}) \frac{q_{PP}}{}$   $PP(q_f, q_{NP})$

$NP\text{-}SBJ(q_{NML}, q_V) \frac{q_{NP\text{-}SBJ}}{}$   $NP\text{-}SBJ(q_{NML}, NP(q_V))$

*President*  $\frac{q_P}{}$  *Alr}ys*

*for*  $\frac{q_f}{}$  *En*

$NP(q_S) \frac{q_{NP}}{}$   $NP(q_S)$

# Syntax-based Machine Translation

## Rules

*Yugoslav*  $\xrightarrow{q_Y}$  *AlywgwslAfy*

*Voislav*  $\xrightarrow{q_V}$  *fwyslAf*

*Serbia*  $\xrightarrow{q_S}$  *SrbyA*

$NML(q_Y, q_P) \xrightarrow{q_{NML}} NP(q_P, q_Y)$

$PP(q_f, q_{NP}) \xrightarrow{q_{PP}} PP(q_f, q_{NP})$

$NP-SBJ(q_{NML}, q_V) \xrightarrow{q_{NP-SBJ}} NP-SBJ(q_{NML}, NP(q_V))$

*President*  $\xrightarrow{q_P}$  *Alr}ys*

*for*  $\xrightarrow{q_f}$  *En*

$NP(q_S) \xrightarrow{q_{NP}} NP(q_S)$

→ Rules of an Extended Top-down Tree Transducer

# Extended Top-down Tree Transducer

## Advantages

- ✓ simple and natural model
- ✓ easy to train (from linguistic resources)  
[GRAEHL et al., 2008]
- ✓ symmetric

# Extended Top-down Tree Transducer

## Advantages

- ✓ simple and natural model
- ✓ easy to train (from linguistic resources)  
[GRAEHL et al., 2008]
- ✓ symmetric

## Generic implementation

- TIBURON [MAY, KNIGHT, 2006]

# Extended Top-down Tree Transducer

## Disadvantages (also of STSG)

- ✗ no discontinuities
- ✗ not binarizable  
[AHO, ULLMAN, 1972; ZHANG et al., 2006]
- ✗ inefficient input/output restriction  
[~, SATTA, 2010]
- ✗ not composable  
[ARNOLD, DAUCHET, 1982]

## Selected references



**CHIANG**: *Hierarchical Phrase-Based Translation*  
Comput. Linguist. 33, 2007



**GALLEY, HOPKINS, KNIGHT, MARCU**  
*What's in a Translation Rule?* Proc. NAACL, 2004



**KOEHN**: *Statistical Machine Translation*  
Cambridge University Press, 2010



**OCH, NEY**: *The Alignment Template Approach to Statistical Machine Translation.* Comput. Linguist. 30, 2004



**PAPINENI, ROUKOS, WARD, ZHU**: *BLEU: A Method for Automatic Evaluation of Machine Translation.* Proc. 40th ACL, 2002