# Statistical Machine Translation, Part II

## Can a Machine Translation Without Knowing What a Verb Is?

## Kevin Knight

USC/Information Sciences Institute
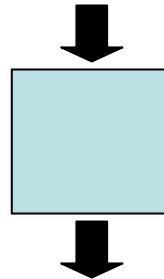USC/Computer Science Department

# Topics

- Quick review of statistical machine translation essentials
  - bilingual text
  - phrase substitution models
  - language models
  - decoding
- Syntactic Approaches
  - syntax-based translation models
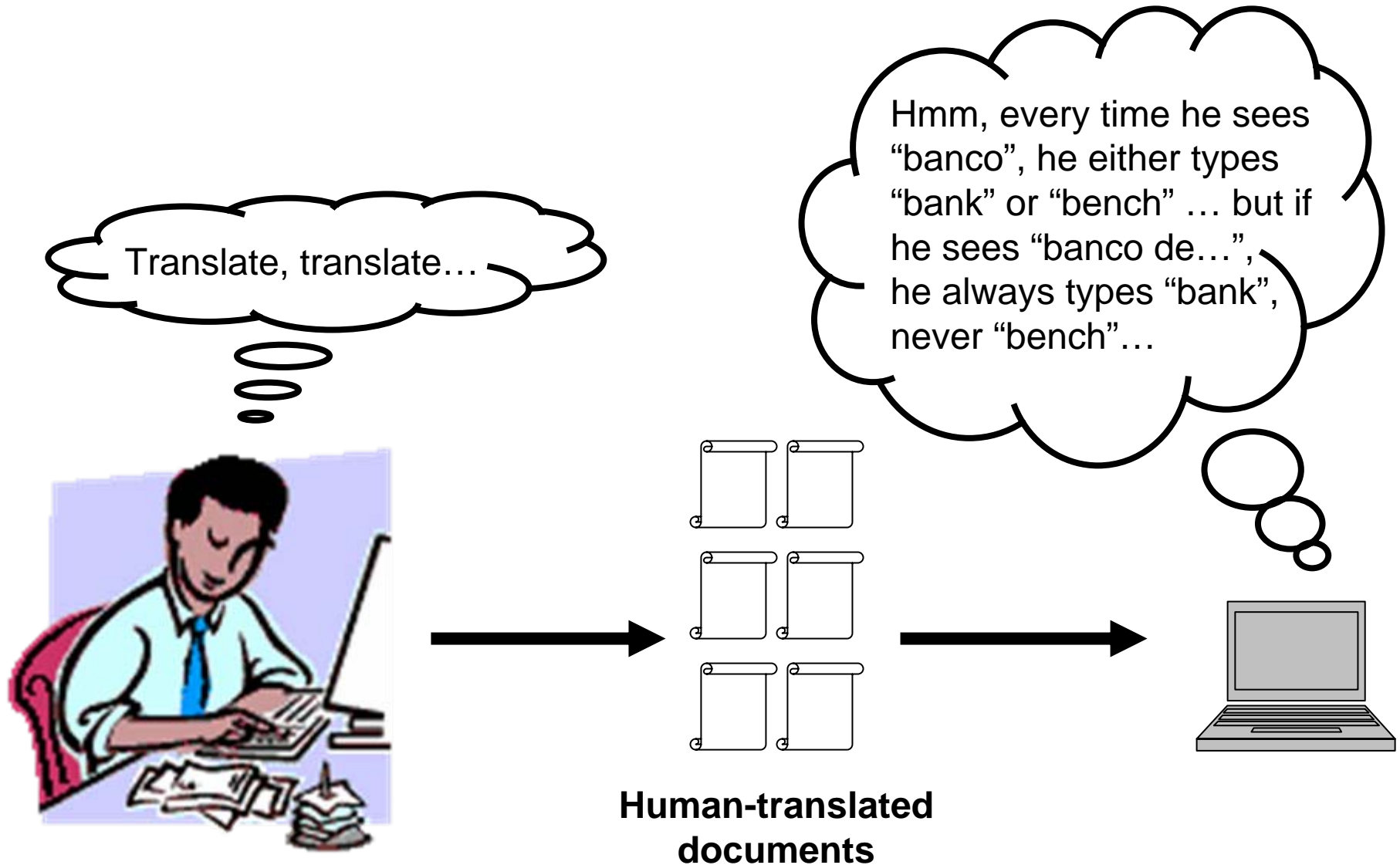  - learning syntactic rules from data
  - decoding

# Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。



The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

# Statistical Machine Translation



Translate, translate…

Hmm, every time he sees "banco", he either types "bank" or "bench" … but if he sees "banco de…", he always types "bank", never "bench"…

**Human-translated documents**

# Spanish/English corpus

| | |
|---|---|
| 1a. Garcia and associates .<br>1b. Garcia y asociados . | 7a. the clients and the associates are enemies .<br>7b. los clients y los asociados son enemigos . |
| 2a. Carlos Garcia has three associates .<br>2b. Carlos Garcia tiene tres asociados . | 8a. the company has three groups .<br>8b. la empresa tiene tres grupos . |
| 3a. his associates are not strong .<br>3b. sus asociados no son fuertes . | 9a. its groups are in Europe .<br>9b. sus grupos estan en Europa . |
| 4a. Garcia has a company also .<br>4b. Garcia tambien tiene una empresa . | 10a. the modern groups sell strong pharmaceuticals .<br>10b. los grupos modernos venden medicinas fuertes . |
| 5a. its clients are angry .<br>5b. sus clientes estan enfadados . | 11a. the groups do not sell zenzanine .<br>11b. los grupos no venden zanzanina . |
| 6a. the associates are also angry .<br>6b. los asociados tambien estan enfadados . | 12a. the small groups are not modern .<br>12b. los grupos pequenos no son modernos . |

# Spanish/English corpus

Translate:  Clients do not sell pharmaceuticals in Europe.

| | |
|---|---|
| 1a. Garcia and associates . <br> 1b. Garcia y asociados . | 7a. the clients and the associates are enemies . <br> 7b. los clients y los asociados son enemigos . |
| 2a. Carlos Garcia has three associates . <br> 2b. Carlos Garcia tiene tres asociados . | 8a. the company has three groups . <br> 8b. la empresa tiene tres grupos . |
| 3a. his associates are not strong . <br> 3b. sus asociados no son fuertes . | 9a. its groups are in Europe . <br> 9b. sus grupos estan en Europa . |
| 4a. Garcia has a company also . <br> 4b. Garcia tambien tiene una empresa . | 10a. the modern groups sell strong pharmaceuticals . <br> 10b. los grupos modernos venden medicinas fuertes . |
| 5a. its clients are angry . <br> 5b. sus clientes estan enfadados . | 11a. the groups do not sell zenzanine . <br> 11b. los grupos no venden zanzanina . |
| 6a. the associates are also angry . <br> 6b. los asociados tambien estan enfadados . | 12a. the small groups are not modern . <br> 12b. los grupos pequenos no son modernos . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:     farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:  farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok **farok** ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok **farok** izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:    farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok **farok** ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat **jjat** bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok **farok** izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat **jjat** quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:   **farok** crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok **farok** ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok **farok** izok stok . | 11a. lalok nok **crrrok** hihok yorok zanzanok . |
| | ??? |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:   **farok** crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok **hihok** ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok **hihok** yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok **hihok** mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:    farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok **hihok** ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok **yorok** ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok **hihok** **yorok** zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok **hihok** mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:   **farok** crrrok **hihok yorok** clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok **clok** . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:   **farok** crrrok **hihok yorok** clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok **clok** .  ??? |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:  **farok** crrrok **hihok yorok** clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:  **farok** crrrok **hihok yorok clok** kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

process of elimination

# Centauri/Arcturan [Knight 97]

Your assignment, translate this to Arcturan:     **farok** crrrok **hihok yorok clok** kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat .     cognate? |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight 97]

Your assignment, put these words in order:   { jjat, arrat, mat, bat, oloat, at-yurp }

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

zero
fertility

# Ready-to-Use Online Bilingual Data



Millions of words (English side)

Legend: Chinese/English, Arabic/English, French/English

(Data stripped of formatting, in sentence-pair format, available from the Linguistic Data Consortium at UPenn).

# Ready-to-Use Online Bilingual Data



Millions of words
(English side)

**Chinese/English**

**Arabic/English**

**French/English**

+ 1m-20m words for
<u>many</u> language pairs

(Data stripped of formatting, in sentence-pair format, available
from the Linguistic Data Consortium at UPenn).

# Bilingual Text (200m words)

English strings

- - - - - - - - - - - - - - - -

…

## Bilingual text

Chinese strings

…

- - - - - - - - - - - - - - - -

# Bilingual Text (200m words)

English
strings

Word
alignments

Chinese
strings

Word-Aligned bilingual text

# Bilingual Text (200m words)

English strings

Word alignments

Chinese strings

Word-Aligned bilingual text

...

...

**Phrase Pair Extraction** [Och & Ney, 2004]

Vast Database of Phrase Pairs

California
ISI NLP

# Phrase-Based Translation

| 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 russian | the | the astronauts | | , |
|---|---|---|---|---|---|---|---|---|---|---|
| it | 7 people included | | by france | | and the | the russian | | international astronautical | of rapporteur . | |
| this | 7 out | including the | **from** | the french | and the russian | | the fifth | | . | |
| these | 7 among | including from | | the french and | | of the russian | of | space | members | . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace | members . | |
| | 7 include | | from the | of france and | | russian | | **astronauts** | | . the |
| | 7 numbers include | **from france** | | and russian | | of astronauts who | | | . '' | |
| | 7 populations include | those from france | | and russian | | astronauts . | | | | |
| | 7 deportees included | come from | **france** | **and russia** | | in | astronautical | personnel | ; | |
| | 7 philtrum | including those from | **france and** | | **russia** | a space | | **member** | | |
| | | including representatives from | france and the | | **russia** | astronaut | | | | |
| | | include | came from | **france and russia** | | by cosmonauts | | | | |
| | | include representatives from | french | **and russia** | | cosmonauts | | | | |
| | | include | came from france | and russia 's | | cosmonauts . | | | | |
| | **includes** | coming from | french and | | russia 's | cosmonaut | | | | |
| | | | french and russian | | 's | astronavigation | | member . | | |
| | | | french | **and russia** | **astronauts** | | | | | |
| | | | | and russia 's | | | | special rapporteur | | |
| | | | | , and | **russia** | | | rapporteur | | |
| | | | | , and russia | | | | rapporteur . | | |
| | | | | , and russia | | | | | | |
| | | | | or | russia 's | | | | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

# Phrase-Based Translation

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 russian | the | the astronauts | | , |
|-----|----------|-----------|---------|---|-----|-------------|-----|----------------|---|---|
| it | 7 people included | | by france | | and the | the russian | | international astronautical | of rapporteur . | |
| this | 7 out | including the | from | the french | and the russian | | the fifth | | . | |
| these | 7 among | including from | | the french and | | of the russian | of | space | members | . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace | members . | |
| | 7 include | | from the | of france and | | russian | | astronauts | | . the |
| | 7 numbers include | | from france | | and russian | | of astronauts who | | . " |
| | 7 populations include | | those from france | | and russian | | astronauts . | | |
| | 7 deportees included | | come from | france | and russia | | in | astronautical | personnel | ; |
| | 7 philtrum | including those from | | france and | russia | | a space | | member | |
| | | including representatives from | | france and the | russia | | astronaut | | | |
| | | include | came from | france and russia | | | by cosmonauts | | | |
| | | include representatives from | french | and russia | | cosmonauts | | | | |
| | | include | came from france | and russia 's | | cosmonauts . | | | | |
| | | includes | coming from | french and | russia 's | cosmonaut | | | | |
| | | | french and russian | | 's | astronavigation | member . | | | |
| | | | french | and russia | astronauts | | | | | |
| | | | and russia 's | | | special rapporteur | | | | |
| | | | , and | russia | | rapporteur | | | | |
| | | | , and russia | | | rapporteur . | | | | |
| | | | , and russia | | | | | | | |
| | | | or | russia 's | | | | | | |

Table 1: #11# the seven - member crew includes astronauts from france and
russia .

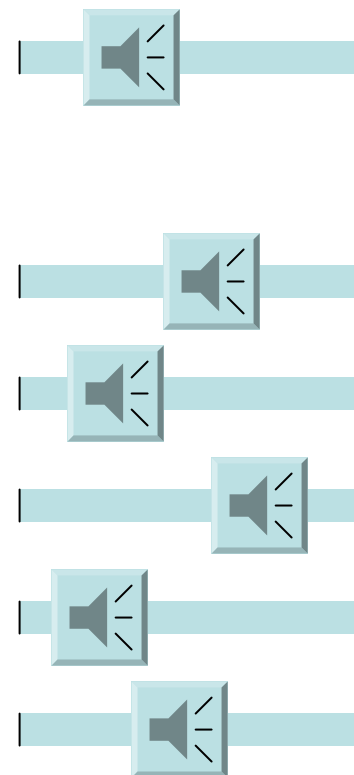Scoring: Try to use phrase pairs that have been frequently observed.
Try to output a sentence with frequent English word sequences.

# Phrase-Based Translation

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 russian | the | the astronauts | | , |
| it | 7 people included | | by france | | and the | the russian | | international astronautical | of rapporteur . | |
| this | 7 out | including the | from | the french | and the russian | | the fifth | | . | |
| these | 7 among | including from | | the french and | | of the russian | of | space | members | . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace | members | |
| | 7 include | | from the | of france and | | russian | | astronauts | | . the |
| | 7 numbers include | from france | | | and russian | | of astronauts who | | . " |
| | 7 populations include | those from france | | | and russian | | astronauts . | | |
| | 7 deportees included | come from | france | and russia | | in | astronautical | personnel | ; |
| | 7 philtrum | including those from | france and | russia | | a space | | member | |
| | | including representatives from | france and the | russia | | astronaut | | |
| | | include | came from | france and russia | | by cosmonauts | | |
| | | include representatives from | french | and russia | | cosmonauts | | |
| | | include | came from france | and russia 's | | cosmonauts . | | |
| | | includes | coming from | french and | russia 's | cosmonaut | | |
| | | | french and russian | | 's | astronavigation | member . | |
| | | | french | and russia | astronauts | | | |
| | | | | and russia 's | | | special rapporteur | |
| | | | | , and | russia | | rapporteur | |
| | | | | , and russia | | | rapporteur . | |
| | | | | , and russia | | | | |
| | | | | or | russia 's | | | |

Table 1: #11# the seven - member crew includes astronauts from france and
russia .

Scoring:  Try to use phrase pairs that have been frequently observed.
        Try to output a sentence with frequent English word sequences.

# Phrase-Based Translation

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 russian | the | the astronauts | | , |
| it | 7 people included | | by france | | and the | the russian | | international astronautical | of rapporteur . | |
| this | 7 out | including the | from | the french | and the russian | | the fifth | | . | |
| these | 7 among | including from | | the french and | | of the russian | of | space | members | . |
| that | 7 persons | including from the | | of france | and to | russian | of the | aerospace | members | |
| | 7 include | | from the | of france and | | russian | | astronauts | | . the |
| | 7 numbers include | from france | | and russian | | of astronauts who | | | . |
| | 7 populations include | those from france | | and russian | | astronauts . | | |
| | 7 deportees included | come from | france | and russia | | in | astronautical | personnel | ; |
| | 7 philtrum | including those from | france and | russia | a space | | astronaut | member | |
| | | including representatives from | france and the | | russia | | | | |
| | | include | came from | france and russia | | by cosmonauts | | | |
| | | include representatives from | french | and russia | | cosmonauts | | | |
| | | include | came from france | and russia 's | | cosmonauts . | | | |
| | | includes | coming from | french and | russia 's | | cosmonaut | | |
| | | | french and russian | | 's | astronavigation | member . | | |
| | | | french | and russia | astronauts | | | | |
| | | | | and russia 's | | | special rapporteur | | |
| | | | | , and | russia | | rapporteur | | |
| | | | | , and russia | | | rapporteur . | | |
| | | | | , and russia | | | | | |
| | | | | or | russia 's | | | | |

Table 1: #11# the seven - member crew includes astronauts from france and russia .

Scoring:  Try to use phrase pairs that have been frequently observed.
　　　　　Try to output a sentence with frequent English word sequences.

# Components

- Training algorithms
  - Word alignment, phrase pair extraction…
    - P(chinese | english) = product of conditional phrase pair probabilities
  - English n-gram models…
    - P(english) = product of trigram probabilities
    - P(w3 | w1 w2)
- Decoding algorithm
  - argmax e  P(chinese | english)   =
    argmax e  P(english) * P(chinese | english)

# Features and Tuning

- English trigram language model

- Phrase pairs
  - Corpus probability of phrase pair
  - Bad-phrase spotter
  - Word-drop spotter
  - "Move Me" preference

- English output length

We compute a total score for each possible translation -- a linear weighted combination of these six values. This generalizes the formula from the previous slide, if we switch to log probs.

# Features and Tuning

- English trigram language model

- Phrase pairs
  - Corpus probability of phrase pair
  - Bad-phrase spotter
  - Word-drop spotter
  - "Move Me" preference

- English output length

We compute a total score for each possible translation -- a linear weighted combination of these six values. This generalizes the formula from the previous slide, if we switch to log probs.

# Features and Tuning

- English trigram language model

- Phrase pairs
  - Corpus probability of phrase pair
  - Bad-phrase spotter
  - Word-drop spotter
  - "Move Me" preference

- English output length

We compute a total score for each possible translation -- a linear weighted combination of these six values. This generalizes the formula from the previous slide, if we switch to log probs.

# Features and Tuning

- English trigram language model

- Phrase pairs
  - Corpus probability of phrase pair
  - Bad-phrase spotter
  - Word-drop spotter
  - "Move Me" preference

- English output length

These six weights (plus about six more) are set by [Och 03]'s **Maximum BLEU training**, which optimizes similarity of MT outputs to human reference translations.

Hill-climbing with MaxBleu [Och 2003]

$W_{TM}$ fixed at 1.0

Translation accuracy

$W_{length}$

$W_{NGLM}$

plot by Emil Ettelaie

# These Ideas Work!



Translation Quality (BLEU) vs. NIST Common Evaluations (Arabic/English) — Phrase-based MT Progress

# Some Lessons

- The simpler, the better
- It takes a long time just to get the bugs out!
- Every change has to be carefully checked
- **Good ideas** often don't help
- Have to try lots of things
- It's highly experimental

# This Kind of MT Research is Highly Experimental

# Two Ways to Improve Translation Systems

Quality of resulting
translation system

more data

Amount of bilingual training data

# Two Ways to Improve Translation Systems



Quality of resulting translation system

better algorithms

more data

Amount of bilingual training data

# Can a machine translate between Chinese and English without knowing what a verb is?

- Of course
- But the output is often bad

  "Frequent high-tech exports are bright spots for foreign trade growth of Guangdong has made important contributions."

- This phrase-based story is a little bit crazy

# Syntax

Maybe we need some grammar?

# MT Research Landscape



Working on syntax-based approach to translation (nouns, verbs, prepositional phrases…)

# MT Research Landscape

Syntax will never work!
We're better off without syntax!
Syntax has been *shown* to make things worse!
It has never worked in speech recognition!
You are crazy!

Language Engineers

Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases…)

# MT Research Landscape

Syntax will never work!
We're better off without syntax!
Syntax has been *shown* to make things worse!
It has never worked in speech recognition!
You are crazy!

Syntax will never work!
You need *semantics*!
Language is about the world!
You are crazy!

AI Fellows

Language Engineers

Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases…)

# MT Research Landscape

**Syntax will never work!**
We're better off without syntax!
Syntax has been *shown* to make things worse!
It has never *ever* worked in speech recognition!
**You are crazy!**

**Syntax will never work!**
You need *semantics*!
Language is about the world!
**You are crazy!**

AI Fellows

Language Engineers

Working on syntax-based approach
to translation (nouns, verbs,
prepositional phrases…)

# MT Progress

# Syntax Started to Be Helpful in 2006

**Translation Accuracy**



Chinese/English

Phrase-based

sentences < 16 words
(NIST-03/04)

Phrase-based

all sentences
(NIST-2003)

45

40

35

30

apr
2005
may jun jul aug sept oct nov dec jan
2006
feb mar apr may jun july jan
2007
feb

# How to Add Syntax?

- Automatically parse training data
  - Add syntactic features to phrase-based system
  - Syntactically re-order source sentences into target-language word order (for training and decoding)
  - Build tree-to-tree translation systems
    - many approaches!
  - Build tree-to-string translation systems
    - many approaches!
  - Build string-to-tree translation systems
    - many approaches!

- Let's take just one approach & investigate

# Phrase-Based Output

枪手　被　警方　击毙　.

Gunman of police killed .

*Decoder*
*Hypothesis #1*

# Phrase-Based Output

枪手 被 警方 击毙 .

Gunman of police attack .

*Decoder*
*Hypothesis #7*

# Phrase-Based Output

枪手　被　警方　击毙 .

Gunman by police killed .

*Decoder
Hypothesis #12*

# Phrase-Based Output

枪手　被　警方　击毙　.

Killed gunman by police .

*Decoder*
*Hypothesis #134*

# Phrase-Based Output

枪手　被　警方　击毙 .

Gunman killed the police .

*Decoder*
*Hypothesis #9,329*

# Phrase-Based Output

枪手 被 警方 击毙 .

Gunman killed by police .

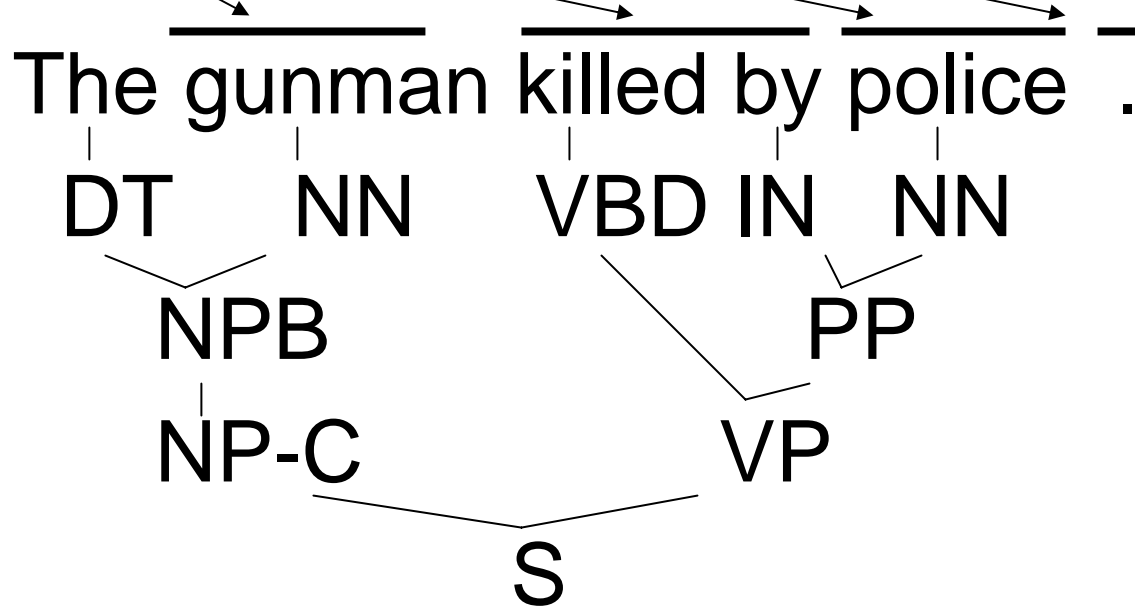highest scoring output, phrase-based model

*Decoder Hypothesis #50,654*

**Problematic:**
- VBD "killed" needs a direct object
- VBN "killed" needs an auxiliary verb ("was")
- countable "gunman" needs an article ("a", "the")
- "passive marker" in Chinese controls re-ordering

**Can't enforce/encourage any of this!**

# Syntax-Based Output

枪手　被　警方　击毙 ．

The gunman killed by police ．

DT　NN　　VBD IN　NN

NPB　　　　　　PP

NP-C　　　　　VP

S

*Decoder Hypothesis #1*

# Syntax-Based Output

枪手　被　警方　击毙 .

Gunman  by police shot  .

NN    IN  NN  VBD

NPB         PP

NP-C              VP

S

*Decoder Hypothesis #16*

# Syntax-Based Output

枪手　被　警方　击毙　.

The gunman was killed by police .    *Decoder Hypothesis #1923*

DT　NN　AUX VBN IN　NN

NPB

NP-C

VP

PP

S

highest scoring output, syntax-based model

# Syntax-Based Output

- Better modeling of target language structure
  - Always a verb
  - Verb is always in the right place

- Better handling of function words
  - They often don't translate
  - They control translation

- Better generalization in translation patterns