

Statistical Machine Translation

Kevin Knight

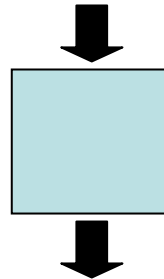
CS562

August 22, 2006



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.

Why People Get Into This Field

- Passion about understanding how human language works
 - What makes one sequence of words grammatical, and another not?
- Interest in foreign languages
 - What's the difference between English and Chinese?
- Desire to change the world
 - How will the world be different when the language barrier disappears?

Thousands of Languages Are Spoken

MANDARIN	885,000,000
SPANISH	332,000,000
ENGLISH	322,000,000
BENGALI	189,000,000

TURKISH	59,000,000
URDU	58,000,000
MIN NAN (China)	49,000,000
JINYU (China)	45,000,000

HINDI	182,000,000
PORTUGUESE	170,000,000
RUSSIAN	170,000,000
JAPANESE	125,000,000
GERMAN	98,000,000



GUJARATI	44,000,000
POLISH	44,000,000
ARABIC	42,500,000
UKRAINIAN	41,000,000

WU (China)	77,175,000
JAVANESE	75,500,800
KOREAN	75,000,000
FRENCH	72,000,000
VIETNAMESE	67,662,000

ITALIAN	37,000,000
XIANG (China)	36,015,000
MALAYALAM	34,022,000
HAKKA (China)	34,000,000

TELUGU	66,350,000
YUE (China)	66,000,000
MARATHI	64,783,000
TAMIL	63,075,000

KANNADA	33,663,000
ORIYA	31,000,000
PANJABI	30,000,000
SUNDA	27,000,000

Source: Ethnologue

Why It's Hard

- Each word has tons of meanings
 - I'll **get** a cup of coffee → ?
 - I didn't **get** that joke → ?
 - I **get** up at 8am → ?
 - I **get** nervous → ?
 - Yeah, I **get** around ... → ?
- Each word has zillions of contexts
- Word order is different

Why It's Hard

- Output must be a grammatical, sensible, never-before-uttered sentence!
- Computers **consume** lots of human language
 - Google, Yahoo, Altavista ...
 - Speech recognizers ...
- More challenging to also **produce** human language
 - What makes one sequence of words grammatical, and another not?

Rapid Progress in the Field of Statistical Machine Translation



insistent Wednesday may
recurred her trips to Libya
tomorrow for flying

Cairo 6-4 (AFP) - an official
announced today in the
Egyptian lines company for
flying Tuesday is a company "
insistent for flying " may
resumed a consideration of a
day Wednesday tomorrow her
trips to Libya of Security Council
decision trace international the
imposed ban comment .

Egyptair Has Tomorrow to
Resume Its Flights to Libya

Cairo 4-6 (AFP) - said an official
at the Egyptian Aviation
Company today that the
company egyptair may resume
as of tomorrow, Wednesday its
flights to Libya after the
International Security Council
resolution to the suspension of
the embargo imposed on Libya.

Statistical MT system results in annual
US government NIST evaluations.

2005



Documentation Logout

Welcome, System Administrator Admin Manager

news
broadcast

Administration

Search: bin laden

Showing: 1-5 of 9

NWI : International Newfirst

Thu Feb 26 13:14:20 EST 2004
- Michael Eisner
- Rio De Janeiro, Mexico City, Madrid
- United Technologies, Walt Disney, General Motors

Aljazeera : News Bulletin

Tue Feb 24 19:00:20 EST 2004
حصني مبارك أسامة بن لادن عبد الله بن عبد العزيز
الفرجين بن دناير الكاظمية
الجزيرة تنظيم القاعدة

Aljazeera : News Bulletin

Tue Feb 24 19:02:20 EST 2004

Aljazeera : News Bulletin

Wed Feb 25 12:14:20 EST 2004
أسامة بن لادن من بجي بالقطار شوة سلطان
الولايات المتحدة والجزائر القاهرة
تنظيم القاعدة وزير الدفاع الإسباني

Aljazeera : News Bulletin

Wed Feb 25 12:00:20 EST 2004
أسامة بن لادن خاسر من بجي بالقطار
الولايات المتحدة بن دناير العراق

news broadcast

Clip: evitap2-vl1_77728128_1077728308_1250kbps 19:09

Previous Clip Storyboard Next Clip

Show: Words

Stop Highlighting Translate

foreign language speech recognition

English translation

which touched several schools and villages, a number of victims of the first keep mwqtan [موقان] to While emerged with proceeds from the final, had been sent from within the island, one of the centres to shelter Al-Hussein City Pakistani military spokesman confirmed the two, General Sultan, said the Ouane northwest Pakistan carried out by the army against the Taliban and al-Qaeda had expired today and resulted in the arrest of a number of foreigners without specify the nationality of those on both, but he denied the presence of senior leaders of al-Qaeda in the area of aggression which returned to the front events of a new force has undergone intensive Pakistani army operations in search of al-Qaeda elements which they suspect of the Pakistani authorities their presence in this region bordering the border with Afghanistan, are repeated attacks on US forces found in the region of al-Qaeda island of al-Qaeda of al-Qaeda ruled leaders in this region and agnened sources on the involvement of military helicopters and heavy weapons during the operation which had been arrested during which some foreigners without specifying their nationality Comrade elaborated to Libya arrested during the operation some people including foreigners ascertained the presence of some foreigners in the region through documents and passports, which were found during the operation comes fire those operation after two addressed walsydan [والسبان] announced by US forces inside Afghanistan, which have paid some analysts to believe that these operations could

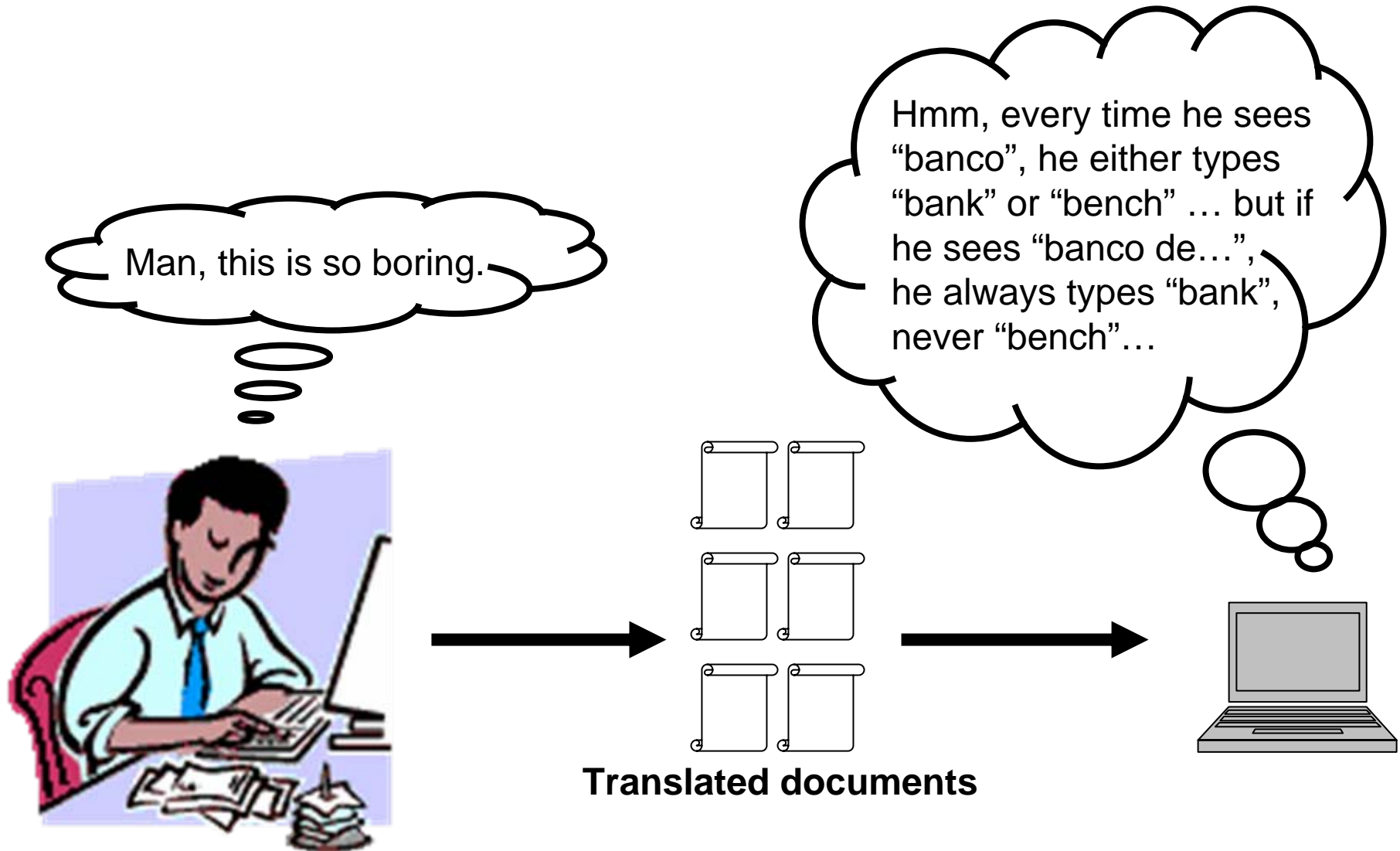
foreign language speech recognition

English translation

searchable archive

00:18:32:04 00:18:32:28 00:18:40:10 00:18:45:25 00:18:59:27 00:19:22:17 00:19:24:11 00:19:24:16 00:19:26:16 00:19:26:23 00:19:29:13 00:19:34:22 00:19:59:09 00:20:00:00

Data-Driven Machine Translation



Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok crrrok hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
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, 1997]

Your assignment, translate this to Arcturan: **farok** crrok hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrok 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, 1997]

Your assignment, translate this to Arcturan: **farok** crrok hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrok 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, 1997]

Your assignment, translate this to Arcturan: **farok** **crrok** hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok . /	11a. lalok nok crrok 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, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok** yorok klok 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
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, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok** **yorok** klok 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok . /	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok .
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, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok yorok** **clock** 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok . /	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok clock . /
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, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok yorok** **clock** 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 anak plok sprok .	8a. lalok brok anak plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok clock .
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, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok yorok** **clock** 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 eneak .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok . /
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok . / /	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilaat .	9b. totat nnat quat oloak at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok clock . / / /
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilaat .
5a. wiwok farok izok stok . /	11a. lalok nok crrrok hihok yorok zanzanak . / / /
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, 1997]

Your assignment, translate this to Arcturan: **farok** crrrok **hihok yorok** **clock** 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 eneak .
2a. ok-drubel ok-voon anak plok sprok .	8a. lalok brok anak plok nok . /
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok . / /	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hila .	9b. totat nnat quat oloak at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok clock . X /
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hila . / process of elimination
5a. wiwok farok izok stok . /	11a. lalok nok crrrok hihok yorok zanzanak . / /
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, 1997]

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

1a. ok-voon ororok sprok . 1b. at-voon bichat dat .	7a. lalok farok ororok lalok sprok izok enemok . / 7b. wat jjat bichat wat dat vat eneak .
2a. ok-drubel ok-voon anak plok sprok . 2b. at-drubel at-voon pippat rrat dat .	8a. lalok brok anak plok nok . / 8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok . / 3b. totat dat arrat vat hilat .	9a. wiwok nok izok kantok ok-yurp . 9b. totat nnat quat oloak at-yurp .
4a. ok-voon anak drok brok jok . 4b. at-voon krat pippat sat lat .	10a. lalok mok nok yorok ghrok klok . / 10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok . / 5b. totat jjat quat cat .	11a. lalok nok crrok hihok yorok zanzanak . / 11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok . 6b. wat dat krat quat cat .	12a. lalok rarok nok izok hihok mok . / 12b. wat nnat forat arrat vat gat .

Centauri/Arcturan [Knight, 1997]

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 anak plok sprok .	8a. lalok brok anak plok nok . /
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghrok . / /	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anak drok brok jok .	10a. lalok mok nok yorok ghrok klok . / / /
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 . / zero
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat . fertility / / /
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 .

It's Really Spanish/English

Clients do not sell pharmaceuticals in Europe => Clientes no venden medicinas en Europa

1a. Garcia and associates .
1b. Garcia y asociados .

7a. the clients and the associates are enemies .
7b. los clients y los asociados son enemigos .

2a. Carlos Garcia has three associates .
2b. Carlos Garcia tiene tres asociados .

8a. the company has three groups .
8b. la empresa tiene tres grupos .

3a. his associates are not strong .
3b. sus asociados no son fuertes .

9a. its groups are in Europe .
9b. sus grupos estan en Europa .

4a. Garcia has a company also .
4b. Garcia tambien tiene una empresa .

10a. the modern groups sell strong pharmaceuticals .
10b. los grupos modernos venden medicinas fuertes .

5a. its clients are angry .
5b. sus clientes estan enfadados .

11a. the groups do not sell zenzanine .
11b. los grupos no venden zanzanina .

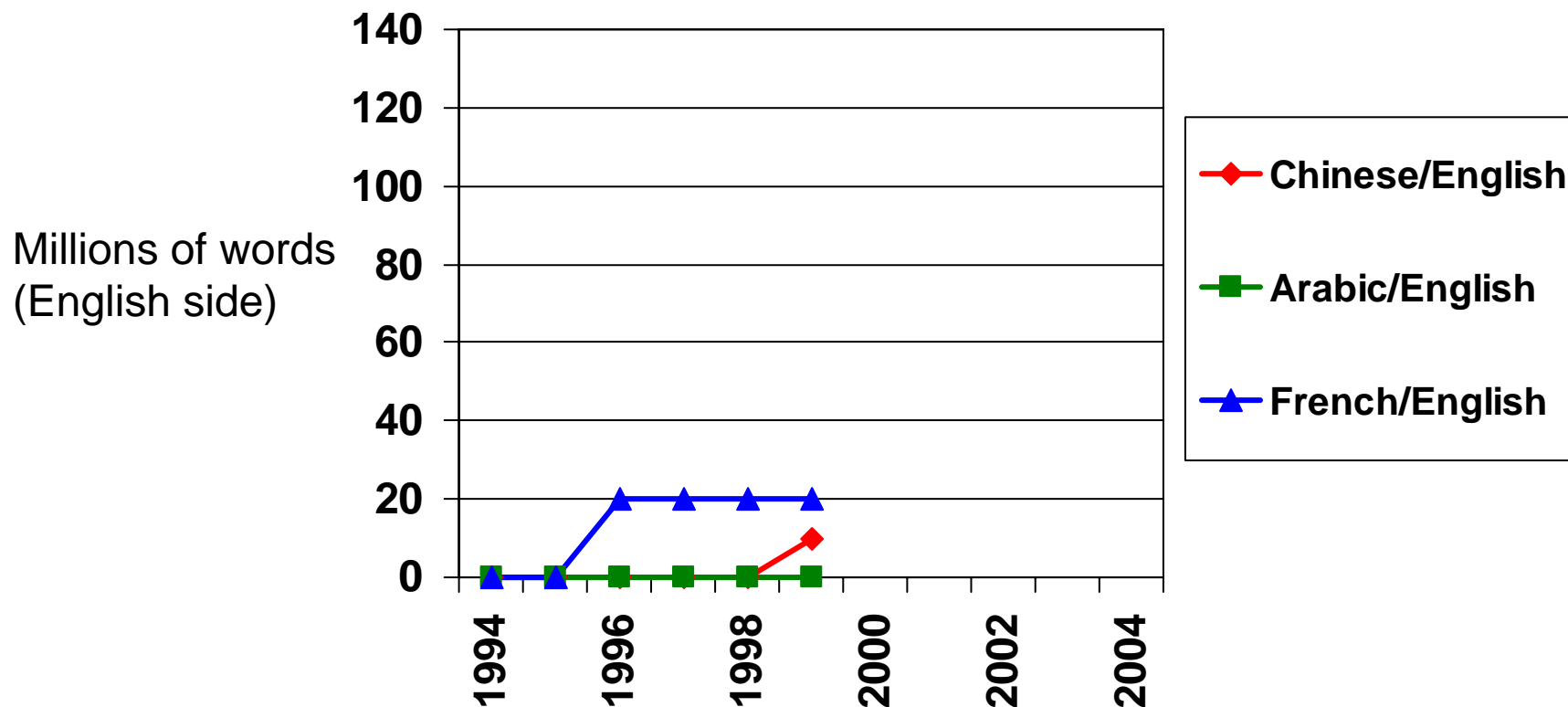
6a. the associates are also angry .
6b. los asociados tambien estan enfadados .

12a. the small groups are not modern .
12b. los grupos pequenos no son modernos .

Data for Statistical Machine Translation

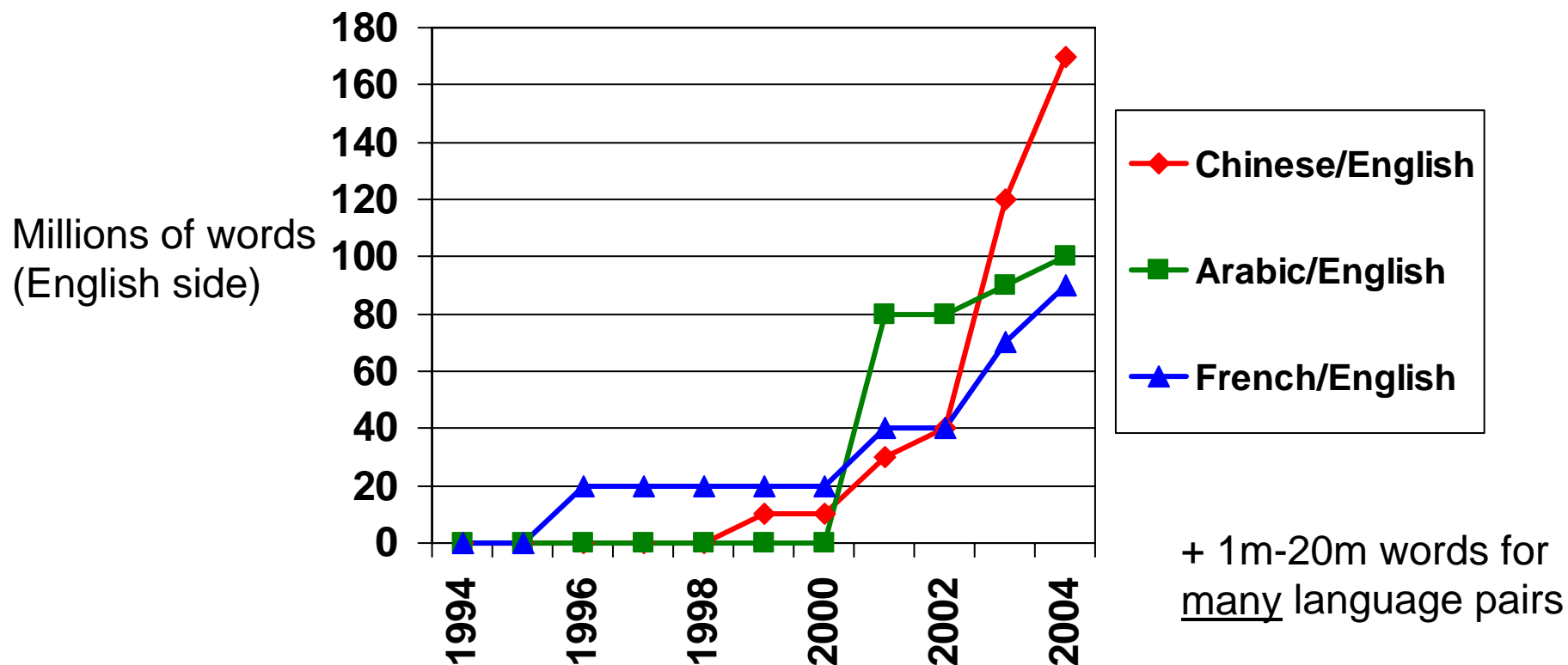
and data preparation

Ready-to-Use Online Bilingual Data



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

Ready-to-Use Online Bilingual Data



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

Sentence Alignment

The old man is
happy. He has
fished many times.
His wife talks to
him. The fish are
jumping. The
sharks await.

El viejo está feliz
porque ha pescado
muchos veces. Su
mujer habla con él.
Los tiburones
esperan.

Sentence Alignment

1. The old man is happy.
2. He has fished many times.
3. His wife talks to him.
4. The fish are jumping.
5. The sharks await.

1. El viejo está feliz porque ha pescado muchos veces.
2. Su mujer habla con él.
3. Los tiburones esperan.

Sentence Alignment

-
1. The old man is happy.
2. He has fished many times.
3. His wife talks to him.
4. The fish are jumping.
5. The sharks await.
1. El viejo está feliz porque ha pescado muchos veces.
2. Su mujer habla con él.
3. Los tiburones esperan.
- The diagram illustrates the alignment of five English sentences with three Spanish sentences. Arrows indicate the following connections: English sentence 1 aligns with Spanish sentence 1; English sentence 2 aligns with Spanish sentence 1; English sentence 3 aligns with Spanish sentence 2; English sentence 4 aligns with Spanish sentence 3; and English sentence 5 aligns with Spanish sentence 3.

Sentence Alignment

- | | | |
|---|----|---|
| 1. The old man is
happy. He has
fished many
times. | —— | 1. El viejo está feliz
porque ha
pescado muchos
veces. |
| 2. His wife talks to
him. | —— | 2. Su mujer habla
con él. |
| 3. The sharks await. | —— | 3. Los tiburones
esperan. |

Note that unaligned sentences are thrown out, and sentences are merged in n-to-m alignments ($n, m > 0$).

Tokenization (or Segmentation)

- English

- Input (some byte stream):

`"There," said Bob.`

- Output (7 “tokens” or “words”):

`" There , " said Bob .`

- Chinese

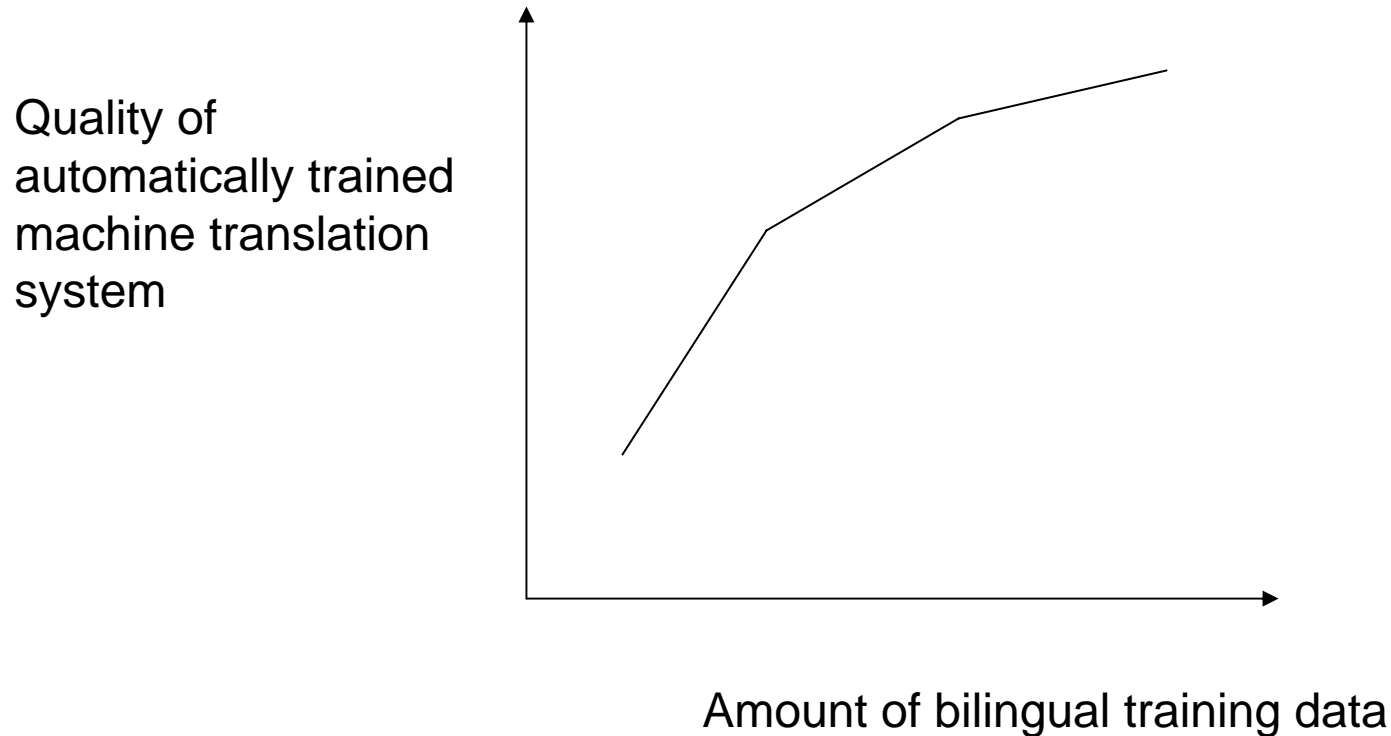
- Input (byte stream):

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。

- Output:

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。

It Is Possible to Draw Learning Curves: How Much Data Do We Need?



MT Evaluation

MT Evaluation

- No single right answer to compare against!
- Manual Evaluation:
 - SSER (subjective sentence error rate)
 - Correct/Incorrect
 - Error categorization
- Testing in an application that uses MT as one component
 - Question answering from foreign language documents
- Automatic:
 - WER (word error rate)
 - **BLEU (Bilingual Evaluation Understudy)**

BLEU Evaluation Metric

(Papineni et al, ACL-2002)

Reference (human) translation:

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 .

Machine translation:

The American [?] international airport and its the office all receives one calls self the sand Arab rich business [?] and so on electronic mail , which sends out ; The threat will be able after public place and so on the airport to start the biochemistry attack , [?] highly alerts after the maintenance.

- N-gram precision (score is between 0 & 1)
 - What percentage of machine n-grams can be found in the reference translation?
 - An n-gram is an sequence of n words
 - Not allowed to use same portion of reference translation twice (can't cheat by typing out "the the the the the")
 - Brevity penalty
 - Can't just type out single word "the" (precision 1.0!)
- *** Amazingly hard to "game" the system (i.e., find a way to change machine output so that BLEU goes up, but quality doesn't)

BLEU Evaluation Metric

(Papineni et al, ACL-2002)

Reference (human) translation:

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 .

Machine translation:

The American [?] international airport and its the office all receives one calls self the sand Arab rich business [?] and so on electronic mail , which sends out ; The threat will be able after public place and so on the airport to start the biochemistry attack , [?] highly alerts after the maintenance.

- BLEU4 formula
(counts n-grams up to length 4)

$$\exp (1.0 * \log p_1 + 0.5 * \log p_2 + 0.25 * \log p_3 + 0.125 * \log p_4 - \max(\text{words-in-reference} / \text{words-in-machine} - 1, 0))$$

p1 = 1-gram precision

P2 = 2-gram precision

P3 = 3-gram precision

P4 = 4-gram precision

Multiple Reference Translations

Reference translation 1:

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.

Reference translation 2:

Guam International Airport and its offices are maintaining a high state of alert after receiving an e-mail that was from a person claiming to be the wealthy Saudi Arabian businessman Bin Laden and that threatened to launch a biological and chemical attack on the airport and other public places.

Machine translation:

The American [?] international airport and its the office all receives one calls self the sand Arab rich business [?] and so on electronic mail, which sends out; The threat will be able after public place and so on the airport to start the biochemistry attack, [?] highly alerts after the maintenance.

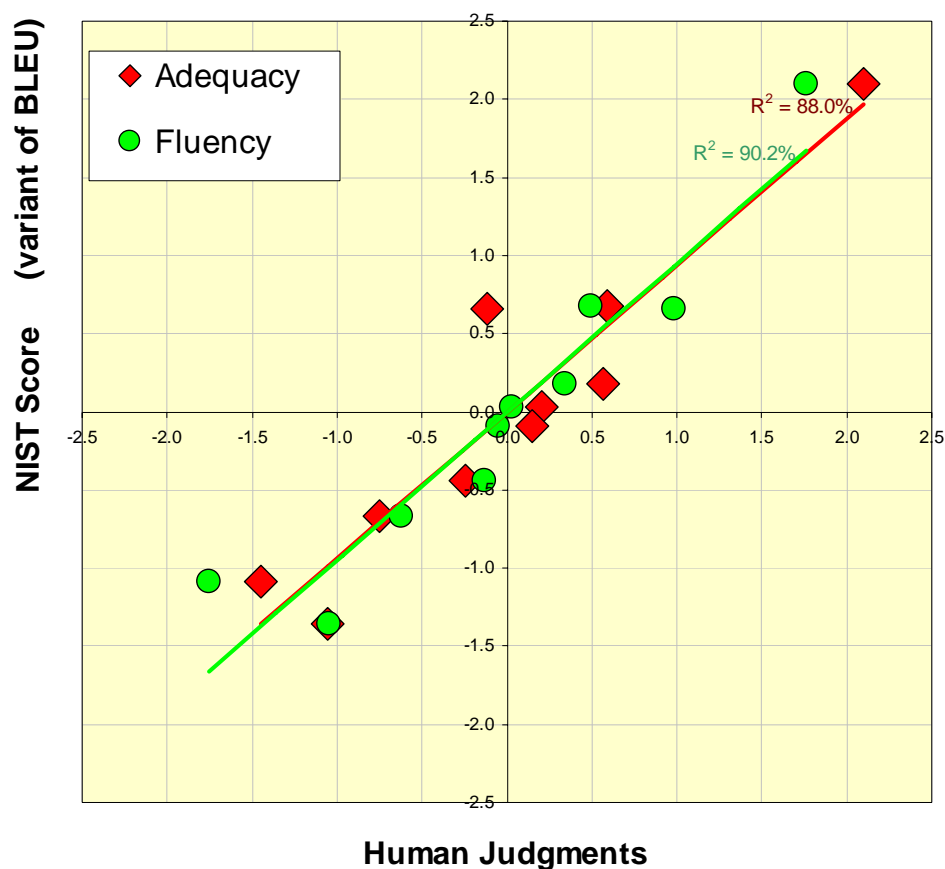
Reference translation 3:

The US International Airport of Guam and its office has received an email from a self-claimed Arabian millionaire named Laden, which threatens to launch a biochemical attack on such public places as airport. Guam authority has been on alert.

Reference translation 4:

US Guam International Airport and its office received an email from Mr. Bin Laden and other rich businessman from Saudi Arabia. They said there would be biochemistry air raid to Guam Airport and other public places. Guam needs to be in high precaution about this matter.

BLEU Tends to Predict Human Judgments



slide from G. Doddington (NIST)

BLEU in Action

枪手被警方击毙。

(Foreign Original)

the gunman was shot to death by the police .

(Reference Translation)

the gunman was police kill .

#1

wounded police jaya of

#2

the gunman was shot dead by the police .

#3

the gunman arrested by police kill .

#4

the gunmen were killed .

#5

the gunman was shot to death by the police .

#6

gunmen were killed by police ?SUB>0 ?SUB>0

#7

al by the police .

#8

the ringer is killed by the police .

#9

police killed the gunman .

#10

BLEU in Action

枪手被警方击毙。

(Foreign Original)

the gunman was shot to death by the police .

(Reference Translation)

the gunman was police kill .

#1

wounded police jaya of

#2

the gunman was shot dead by the police .

#3

the gunman arrested by police kill .

#4

the gunmen were killed .

#5

the gunman was shot to death by the police .

#6

gunmen were killed by police ?SUB>0 ?SUB>0

#7

al by the police .

#8

the ringer is killed by the police .

#9

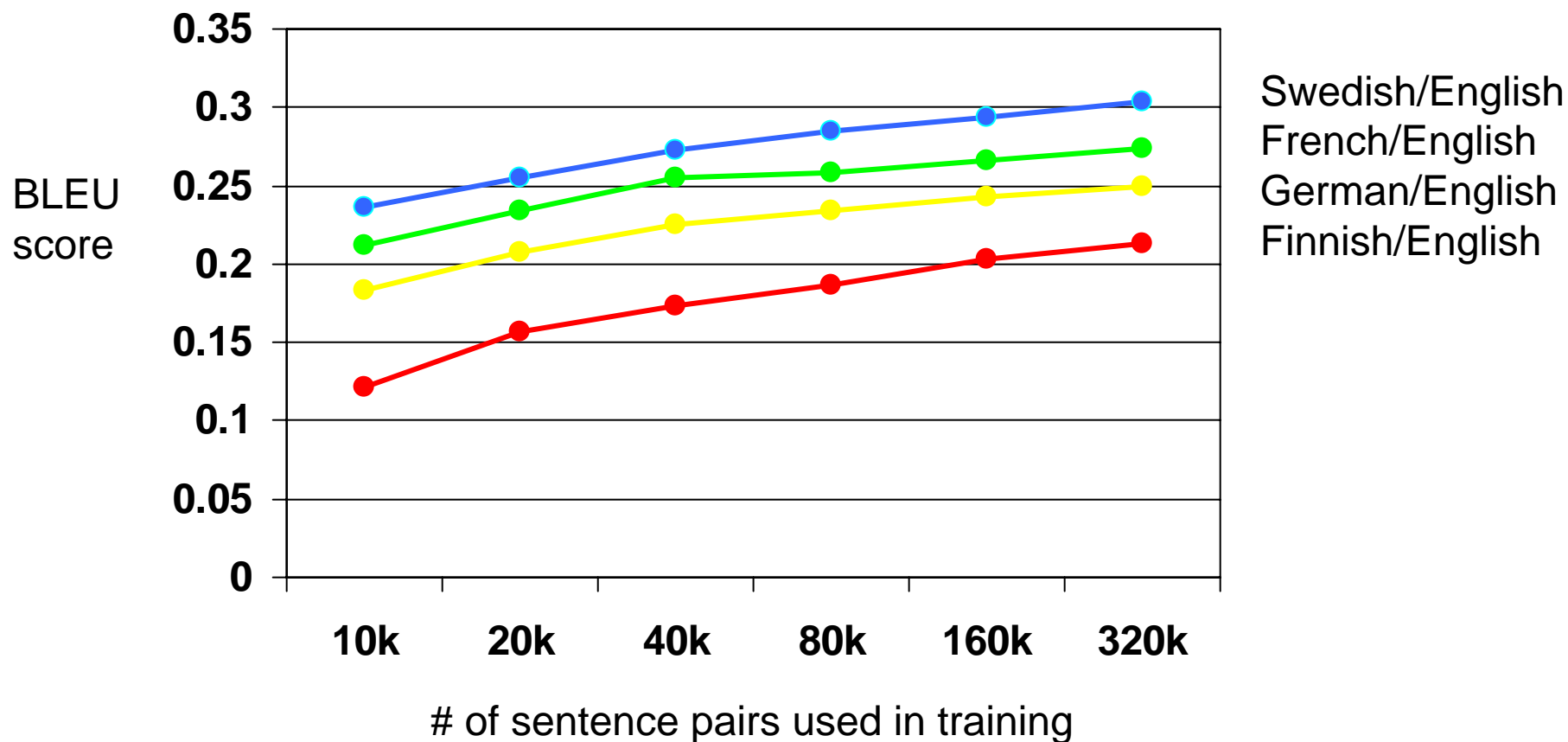
police killed the gunman .

#10

green = 4-gram match (good!)

red = word not matched (bad!)

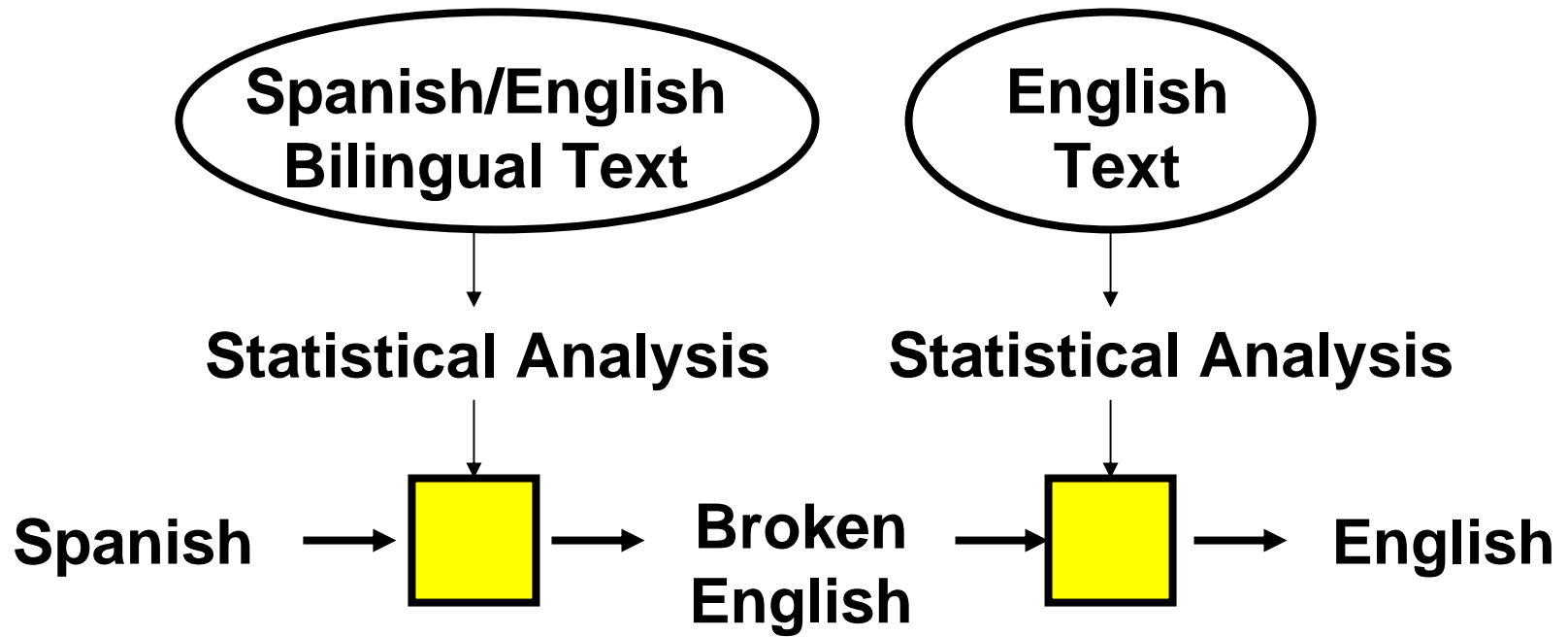
Sample Learning Curves



Experiments by
Philipp Koehn

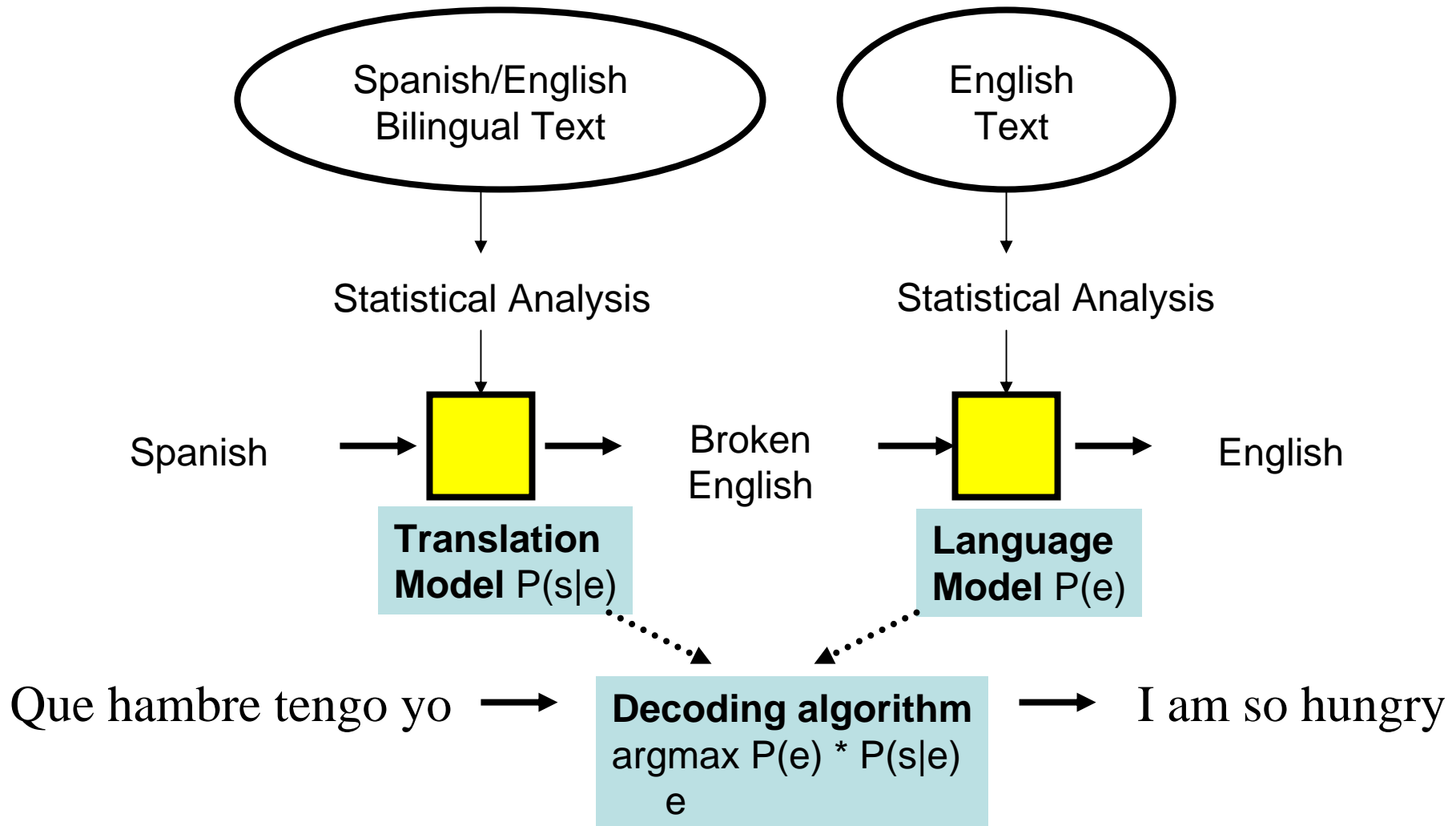
Word-Based Statistical MT

Statistical MT Systems



Que hambre tengo yo → What hunger have I,
Hungry I am so,
I am so hungry, → I am so hungry
Have I that hunger ...

Statistical MT Systems



Three Problems for Statistical MT

- Language model
 - Given an English string e , assigns $P(e)$ by formula
 - good English string \rightarrow high $P(e)$
 - random word sequence \rightarrow low $P(e)$
- Translation model
 - Given a pair of strings $\langle f, e \rangle$, assigns $P(f | e)$ by formula
 - $\langle f, e \rangle$ look like translations \rightarrow high $P(f | e)$
 - $\langle f, e \rangle$ don't look like translations \rightarrow low $P(f | e)$
- Decoding algorithm
 - Given a language model, a translation model, and a new sentence f ... find translation e maximizing $P(e) * P(f | e)$

Language Modeling

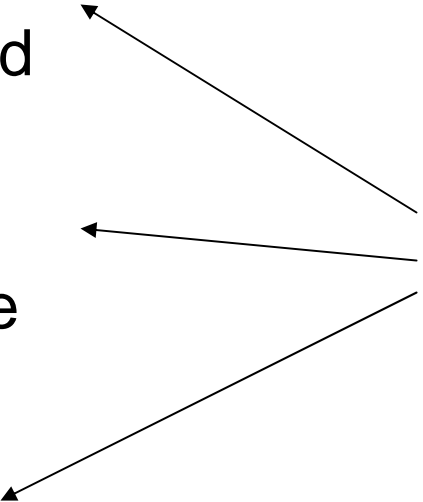
Goal of a language model for MT:

He is on the soccer field
He is in the soccer field

Is table the on cup the
The cup is on the table

American shrine
American company

Need to make these decisions, because translation model may not have a lot of context information!



The diagram consists of three arrows pointing from the explanatory text on the right to the ambiguous sentence pairs on the left. The top arrow points from the explanatory text to the 'He is on the soccer field' / 'He is in the soccer field' pair. The middle arrow points from the explanatory text to the 'Is table the on cup the' / 'The cup is on the table' pair. The bottom arrow points from the explanatory text to the 'American shrine' / 'American company' pair.

The Classic Language Model

Word Bigrams

Process model of English:

Generate each word based only on the previous word.

$P(\text{I saw water on the table}) =$

$P(\text{I} \mid \text{START}) \cdot$

$P(\text{saw} \mid \text{I}) \cdot$

$P(\text{water} \mid \text{saw}) \cdot$

$P(\text{on} \mid \text{water}) \cdot$

$P(\text{the} \mid \text{on}) \cdot$

$P(\text{table} \mid \text{the}) \cdot$

$P(\text{END} \mid \text{table})$

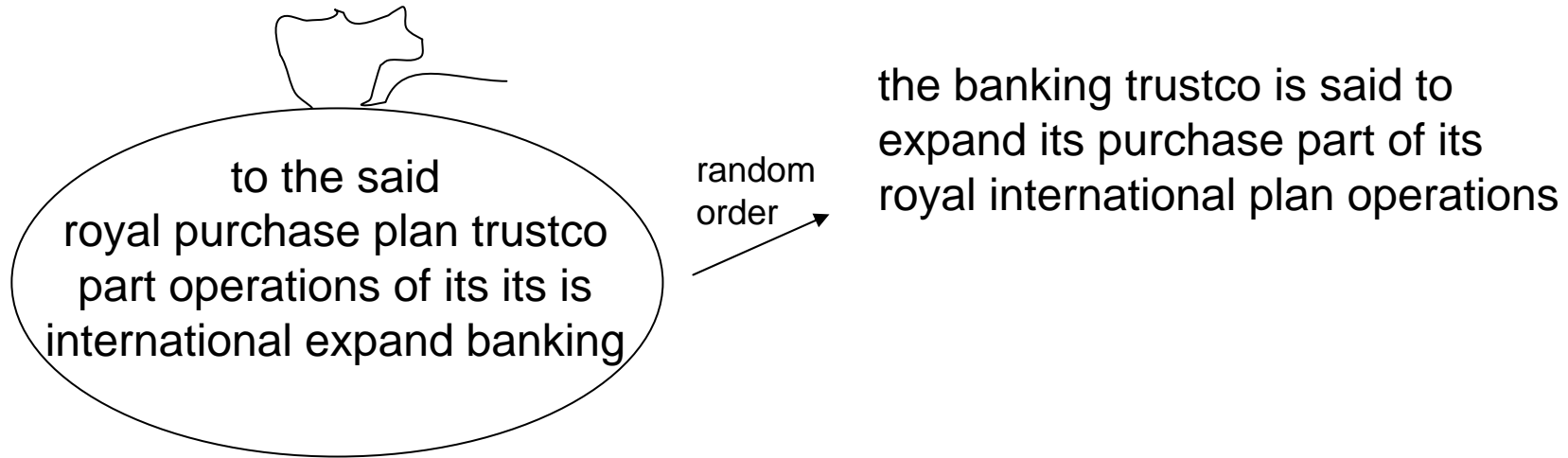
Probabilities can be tabulated
from an online English corpus ...
just like Weaver's Turkish case.

Trigram Language Model

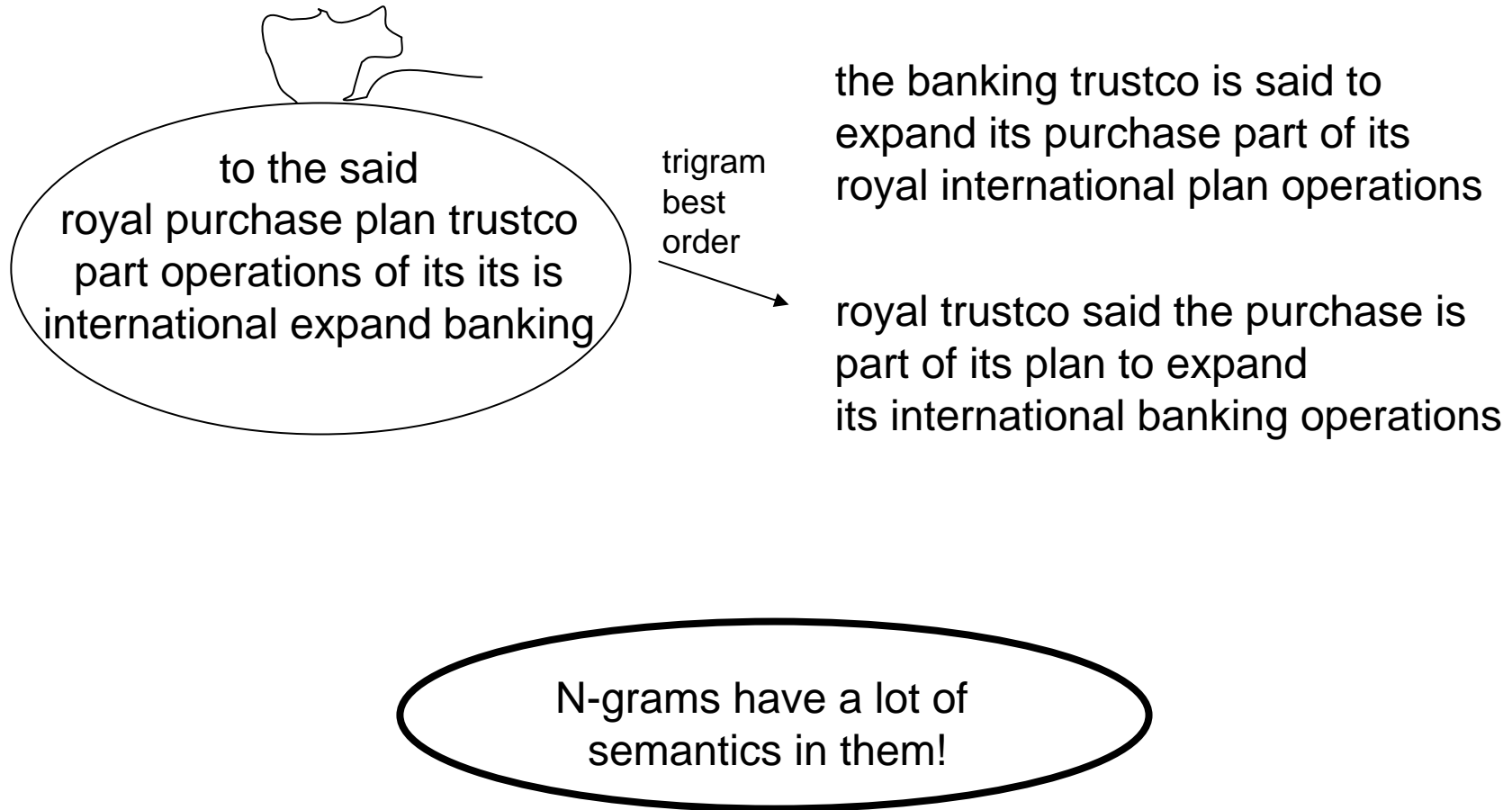


to the said
royal purchase plan trustco
part operations of its its is
international expand banking

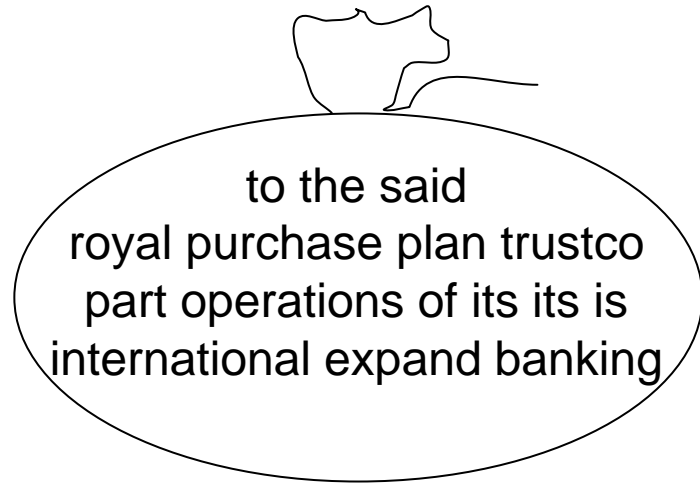
Trigram Language Model



Trigram Language Model

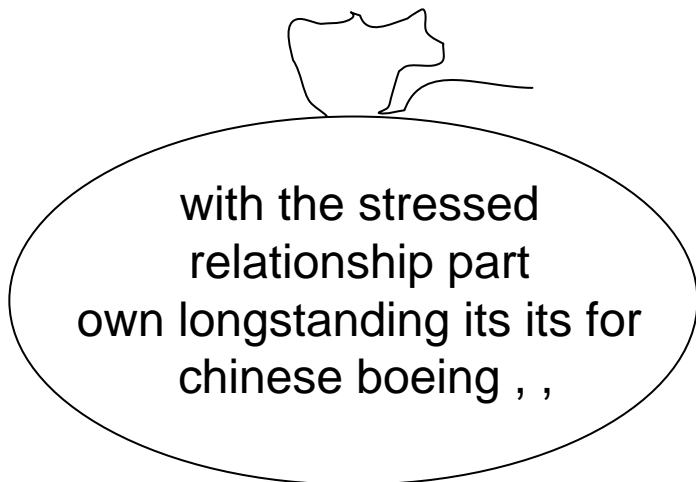


Trigram Language Model



the banking trustco is said to
expand its purchase part of its
royal international plan operations

royal trustco said the purchase is
part of its plan to expand
its international banking operations



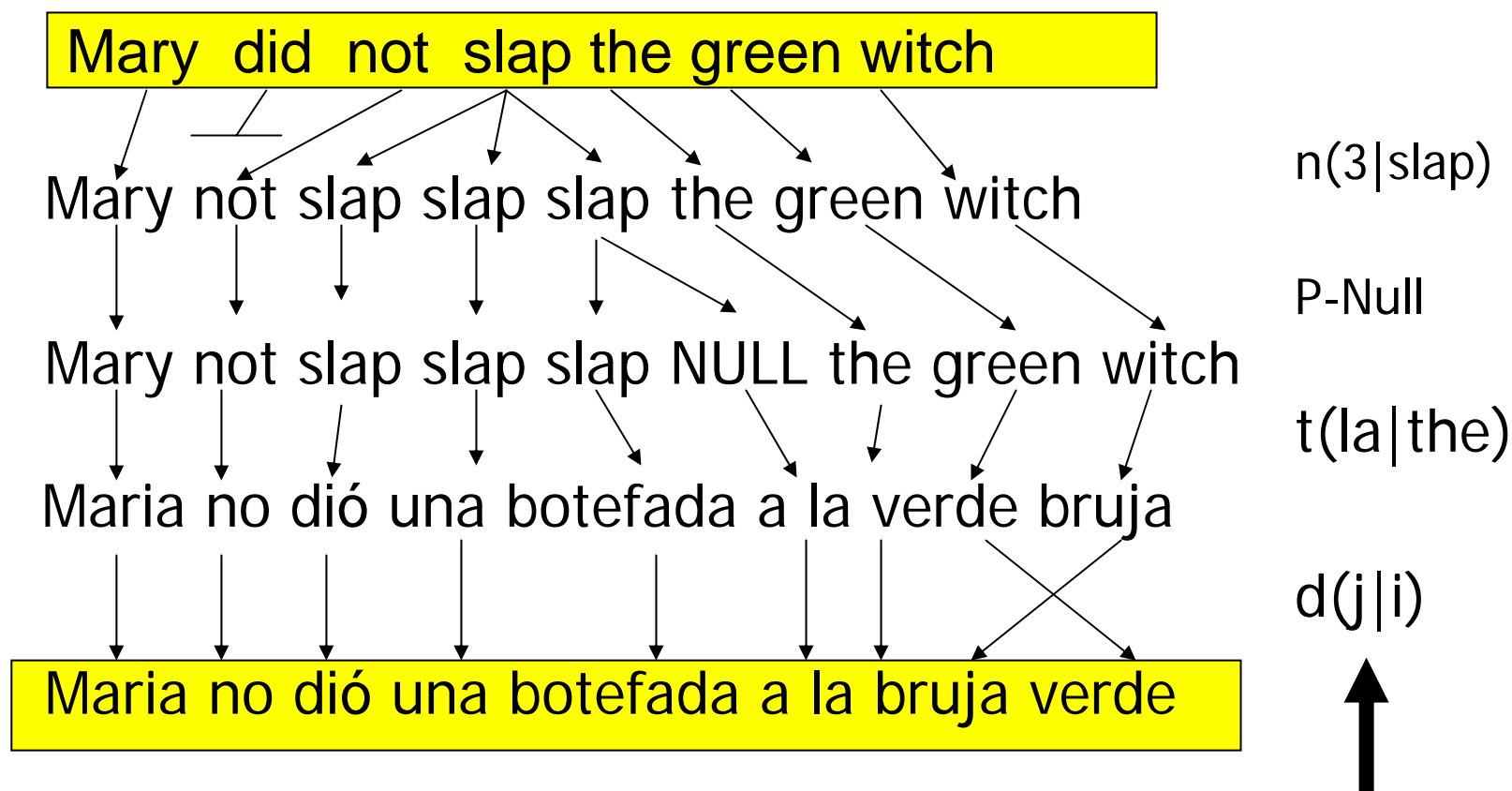
for its part, stressed the longstanding
relationship with its own, chinese boeing

*boeing, for its part, stressed its own
longstanding relationship with the chinese*

The Classic Translation Model

Word Substitution/Permutation [IBM Model 3, Brown et al., 1993]

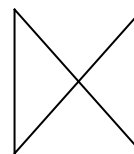
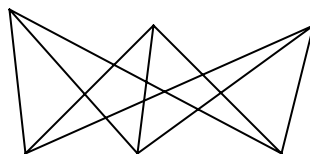
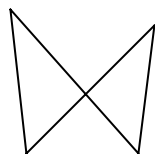
Generative story:



Probabilities can be learned from raw bilingual text.

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...



... the house ... the blue house ... the flower ...

All word alignments equally likely

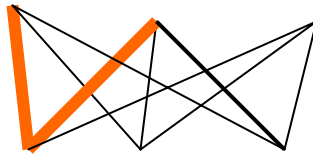
All $P(\text{french-word} \mid \text{english-word})$ equally likely

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...



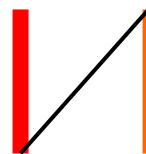
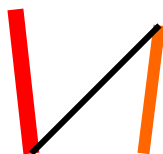
... the house ... the blue house ... the flower ...



“la” and “the” observed to co-occur frequently,
so $P(\text{la} \mid \text{the})$ is increased.

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...



... the house ... the blue house ... the flower ...

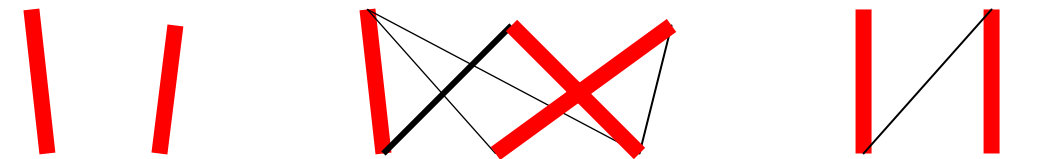
“house” co-occurs with both “la” and “maison”, but $P(\text{maison} \mid \text{house})$ can be raised without limit, to 1.0, while $P(\text{la} \mid \text{house})$ is limited because of “the”

(pigeonhole principle)

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...

... the house ... the blue house ... the flower ...




The diagram illustrates word alignment between the French sentence "... la maison ... la maison bleue ... la fleur ..." and the English sentence "... the house ... the blue house ... the flower ...". Red lines represent alignments that are settled or correct. For the first pair, "la maison" aligns with "the house". For the second pair, "la maison bleue" aligns with "the blue house", with "bleue" aligning to "blue". For the third pair, "la fleur" aligns with "the flower". Black lines and red X marks indicate rejected or incorrect alignments, such as "la maison" aligning with "the blue house" or "la fleur" aligning with "the house".

settling down after another iteration

Statistical Machine Translation

... la maison ... la maison bleue ... la fleur ...

... the house ... the blue house ... the flower ...



Inherent hidden structure revealed by EM training!

For details, see:

- “A Statistical MT Tutorial Workbook” (Knight, 1999).
- “The Mathematics of Statistical Machine Translation” (Brown et al, 1993)
- Software: GIZA++

Statistical Machine Translation

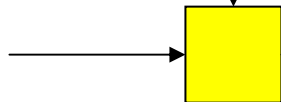
... la maison ... la maison bleue ... la fleur ...

... the house ... the blue house ... the flower ...

$P(\text{juste} \mid \text{fair}) = 0.411$
 $P(\text{juste} \mid \text{correct}) = 0.027$
 $P(\text{juste} \mid \text{right}) = 0.020$

...

new French
sentence

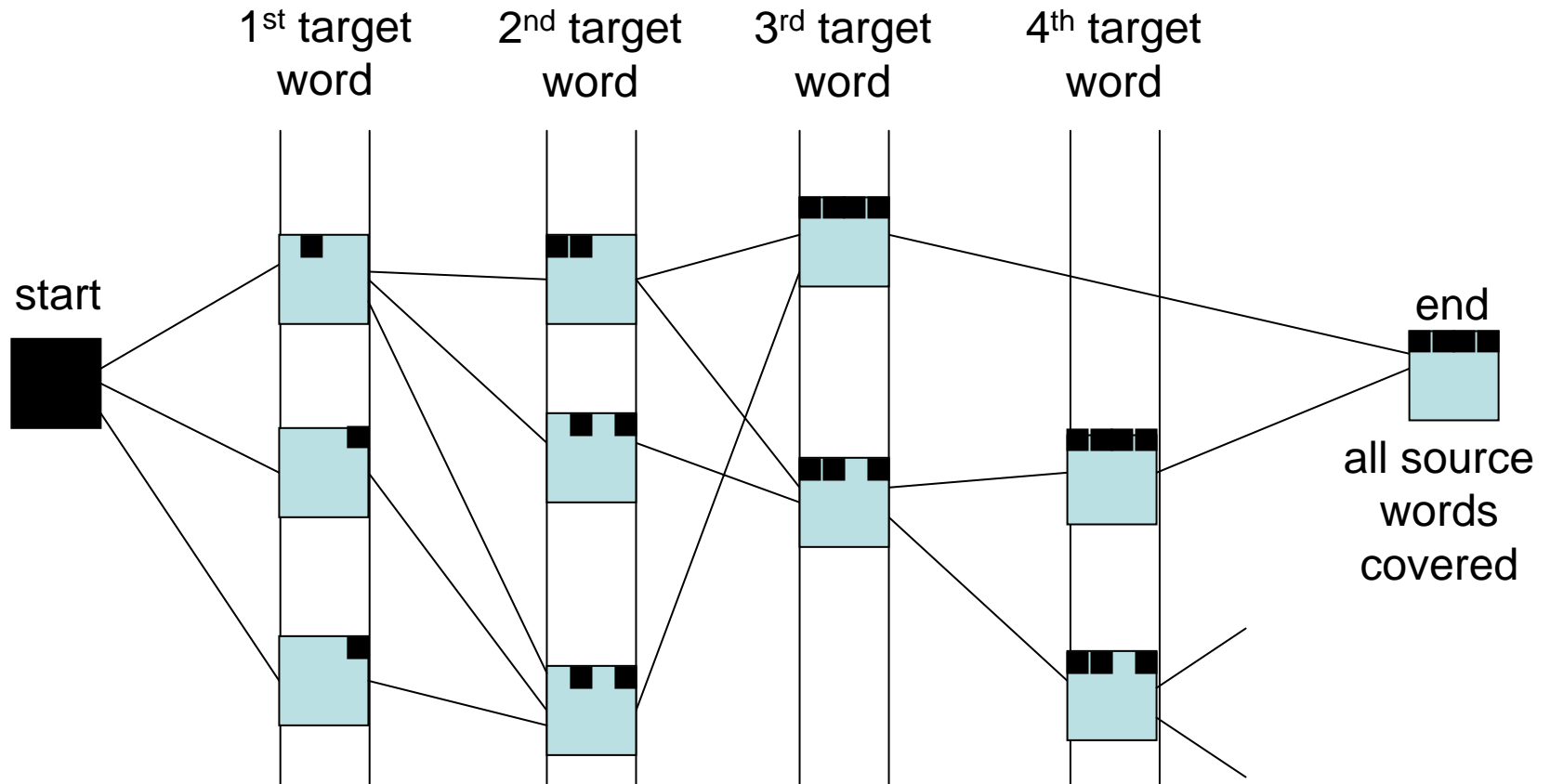


Possible English translations,
to be rescored by language model

Decoding for “Classic” Models

- Of all conceivable English word strings, find the one maximizing $P(e) \times P(f | e)$
- Decoding is an NP-complete challenge
 - (Knight, 1999)
- Several search strategies are available
- Each potential English output is called a *hypothesis*.

Dynamic Programming Beam Search

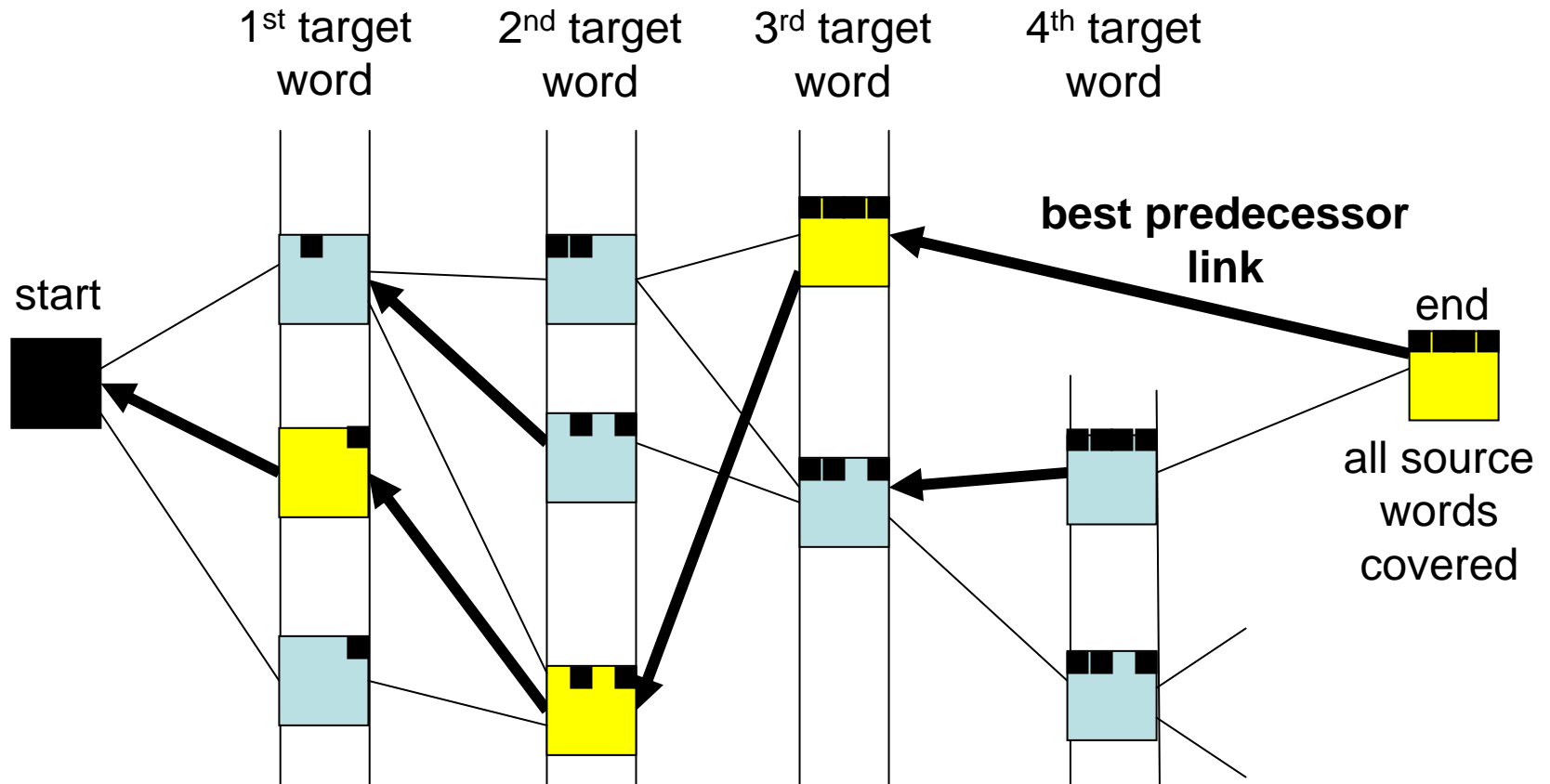


Each partial translation hypothesis contains:

- Last English word chosen + source words covered by it
- Next-to-last English word chosen
- Entire coverage vector (so far) of source sentence ■ ■ ■
- Language model and translation model scores (so far)

[Jelinek, 1969;
Brown et al, 1996 US Patent;
(Och, Ueffing, and Ney, 2001)]

Dynamic Programming Beam Search



Each partial translation hypothesis contains:

- Last English word chosen + source words covered by it
- Next-to-last English word chosen
- Entire coverage vector (so far) of source sentence ■ ■ ■
- Language model and translation model scores (so far)

[Jelinek, 1969;
Brown et al, 1996 US Patent;
(Och, Ueffing, and Ney, 2001)]

The Classic Results

- *la politique de la haine .* (Foreign Original)
- politics of hate . (Reference Translation)
- the policy of the hatred . (IBM4+N-grams+Stack)

- *nous avons signé le protocole .* (Foreign Original)
- we did sign the memorandum of agreement . (Reference Translation)
- we have signed the protocol . (IBM4+N-grams+Stack)

- *où était le plan solide ?* (Foreign Original)
- but where was the solid plan ? (Reference Translation)
- where was the economic base ? (IBM4+N-grams+Stack)

对外经济贸易合作部今天提供的数据表明，今年至十一月中国实际利用外资四百六十九点五九亿美元，其中包括外商直接投资四百点零七亿美元。

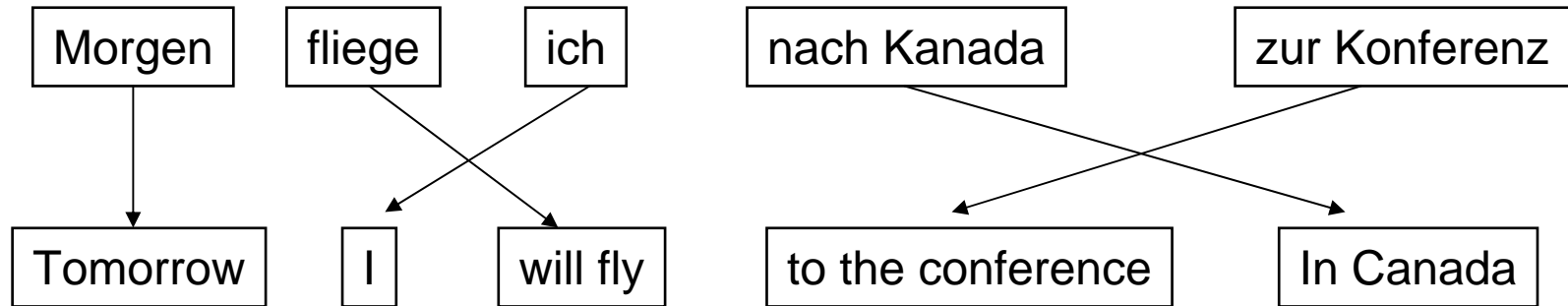
the Ministry of Foreign Trade and Economic Cooperation, including foreign direct investment 40.007 billion US dollars today provide data include that year to November china actually using foreign 46.959 billion US dollars and

Flaws of Word-Based MT

- Multiple English words for one French word
 - IBM models can do one-to-many (fertility) but not many-to-one
- Phrasal Translation
 - “real estate”, “note that”, “interest in”
- Syntactic Transformations
 - Verb at the beginning in Arabic
 - Translation model penalizes any proposed re-ordering
 - Language model not strong enough to force the verb to move to the right place

Phrase-Based Statistical MT

Phrase-Based Statistical MT



- Foreign input segmented in to phrases
 - “phrase” is any sequence of words
- Each phrase is probabilistically translated into English
 - $P(\text{to the conference} \mid \text{zur Konferenz})$
 - $P(\text{into the meeting} \mid \text{zur Konferenz})$
- Phrases are probabilistically re-ordered

See [Koehn et al, 2003] for an intro.

This is state-of-the-art!

Advantages of Phrase-Based

- Many-to-many mappings can handle non-compositional phrases
- Local context is very useful for disambiguating
 - “Interest rate” → ...
 - “Interest in” → ...
- The more data, the longer the learned phrases
 - Sometimes whole sentences

How to Learn the Phrase Translation Table?

- Start with word alignment, build phrases from that.

	Maria	no	dió	una	bofetada	a	la	bruja	verde
Mary									
did									
not									
slap									
the									
green									
witch									

Word Alignment Induced Phrases

	Maria	no	dió	una	bofetada	a	la	bruja	verde
Mary									
did									
not									
slap									
the									
green									
witch									

(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)

Word Alignment Induced Phrases

	Maria	no	dió	una	bofetada	a	la	bruja	verde
Mary									
did									
not									
slap									
the									
green									
witch									

(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)
(a la, the) (dió una bofetada a, slap the)

Word Alignment Induced Phrases

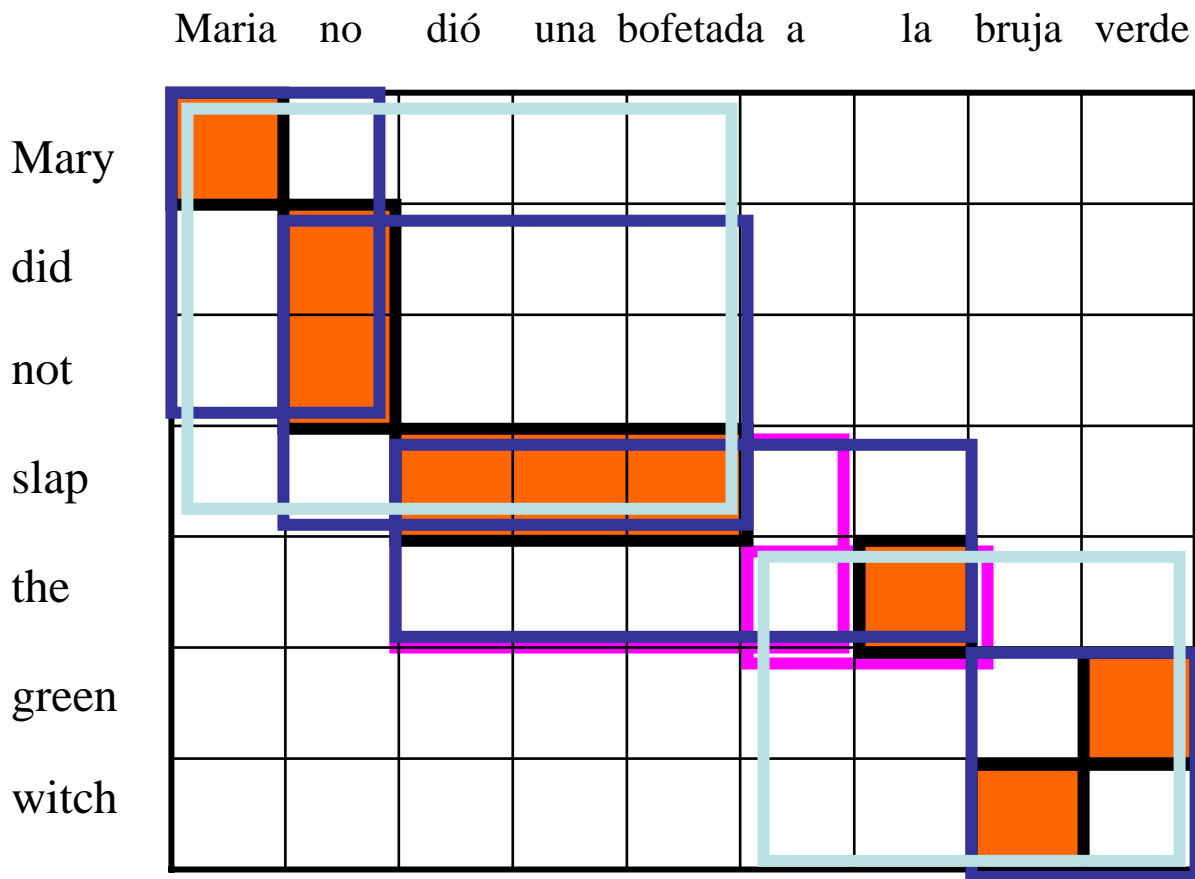
	Maria	no	dió	una	bofetada	a	la	bruja	verde
Mary									
did									
not									
slap									
the									
green									
witch									

(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)

(a la, the) (dió una bofetada a, slap the)

(Maria no, Mary did not) (no dió una bofetada, did not slap), (dió una bofetada a la, slap the)
(bruja verde, green witch)

Word Alignment Induced Phrases



(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)

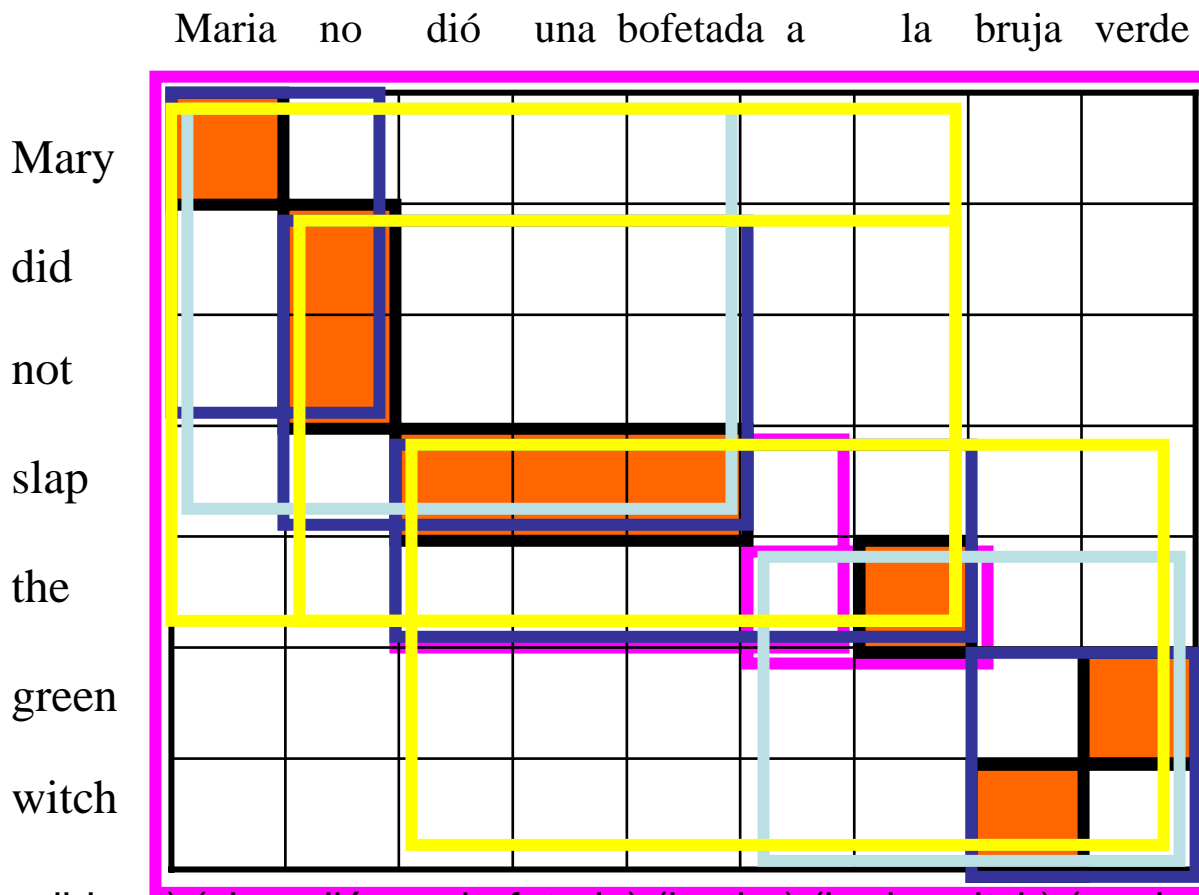
(a la, the) (dió una bofetada a, slap the)

(Maria no, Mary did not) (no dió una bofetada, did not slap), (dió una bofetada a la, slap the)

(bruja verde, green witch) (Maria no dió una bofetada, Mary did not slap)

(a la bruja verde, the green witch) ...

Word Alignment Induced Phrases



(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)

(a la, the) (dió una bofetada a, slap the)

(Maria no, Mary did not) (no dió una bofetada, did not slap), (dió una bofetada a la, slap the)

(bruja verde, green witch) (Maria no dió una bofetada, Mary did not slap)

(a la bruja verde, the green witch) ...

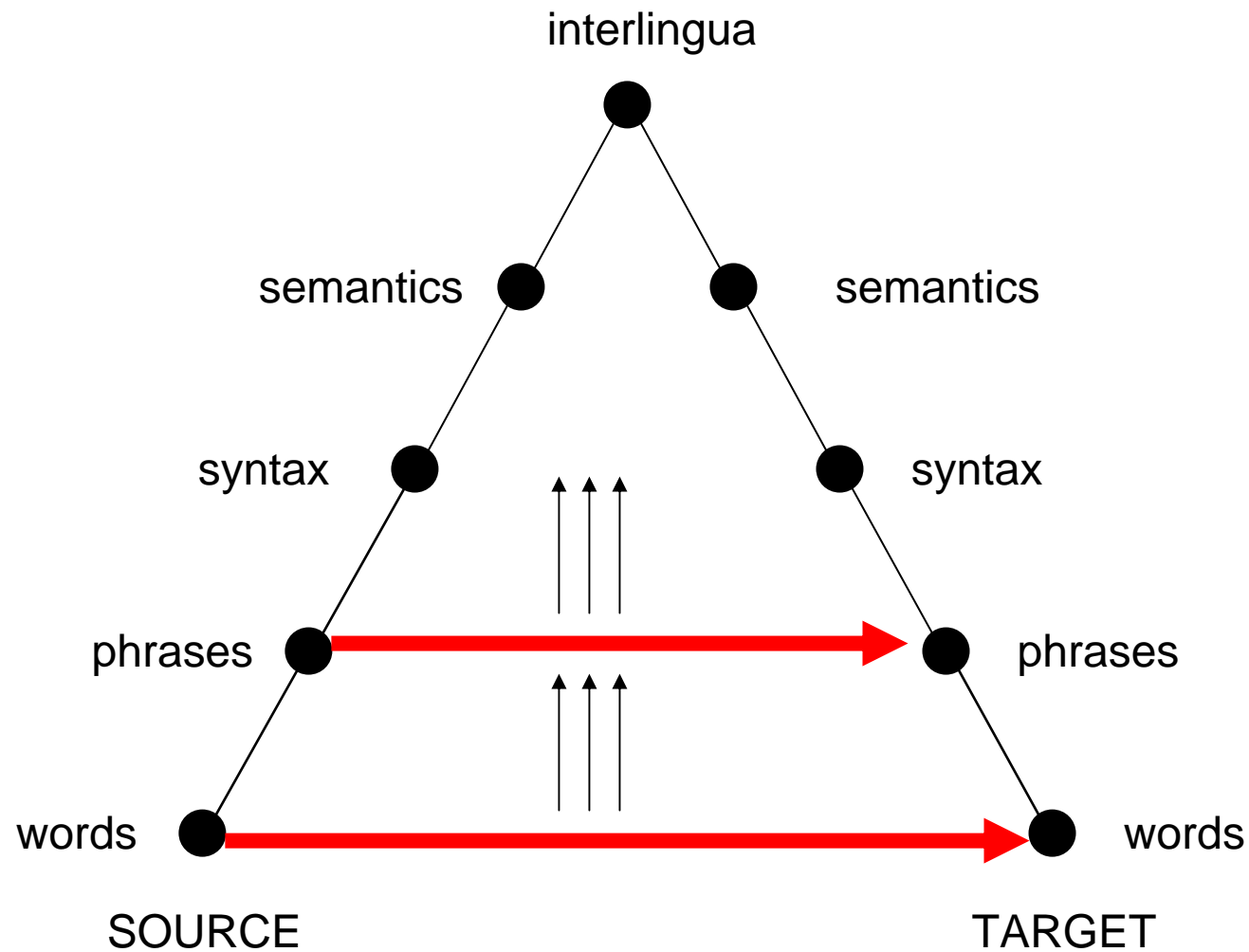
(Maria no dió una bofetada a la bruja verde, Mary did not slap the green witch)

Phrase-Based Models

- So, now we have a vast list of phrase pairs with frequencies
 - $P(F | E) = \text{count}(F, E) / \text{count}(E)$
- Billions of phrase pairs ready in our database!
- Translation accuracy is much better than with word-based methods

Syntax and Semantics in Statistical MT

MT Pyramid

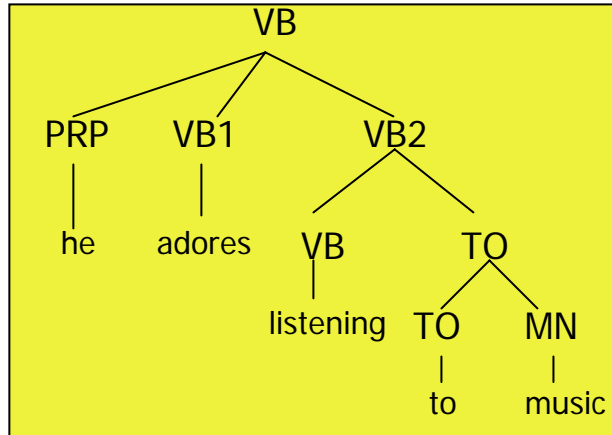


Why Syntax?

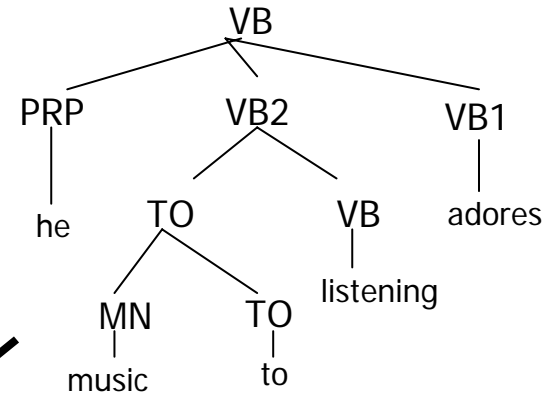
- Need much more grammatical output
- Need accurate control over re-ordering
- Need accurate insertion of function words
- Word translations need to depend on grammatically-related words

Yamada/Knight 01: Modeling and Training

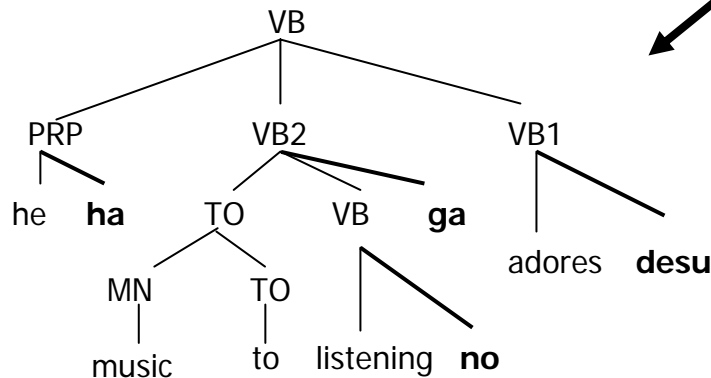
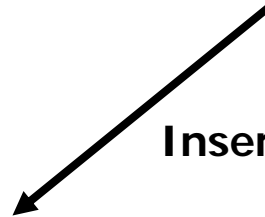
Parse Tree(E)



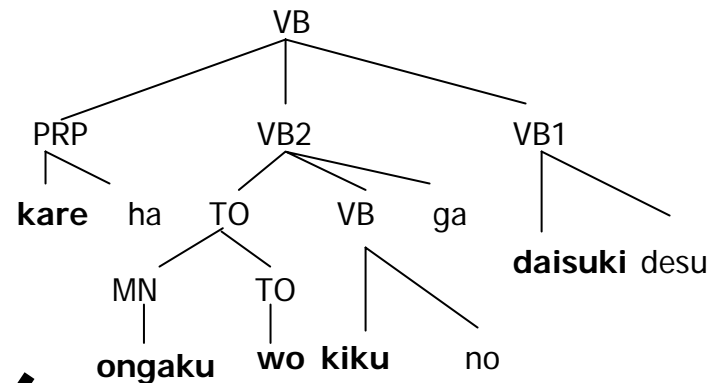
Reorder



Insert



Translate



Take Leaves



Sentence(J)

Kare ha ongaku wo kiku no ga daisuki desu

Japanese/English Reorder Table

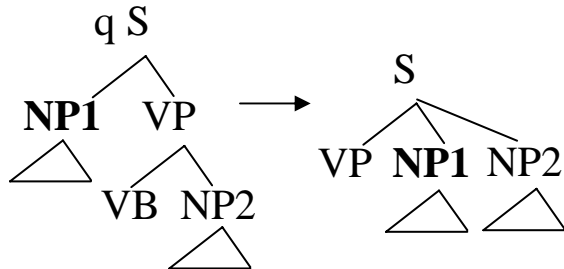
Original Order	Reordering	$P(\text{reorder} \text{original})$
PRP VB1 VB2	PRP VB1 VB2	0.074
	PRP VB2 VB1	0.723
	VB1 PRP VB2	0.061
	VB1 VB2 PRP	0.037
	VB2 PRP VB1	0.083
	VB2 VB1 PRP	0.021
VB TO	VB TO	0.107
	TO VB	0.893
TO NN	TO NN	0.251
	NN TO	0.749

For French/English, useful parameters like $P(N \text{ ADJ} | \text{ADJ } N)$.

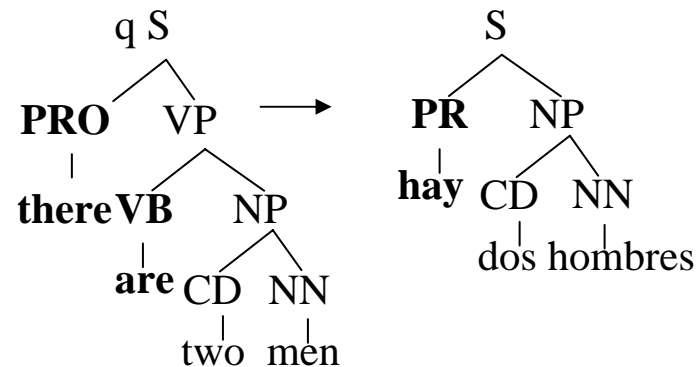
Syntax MT Models are Tree Automata

[Graehl & Knight 04]

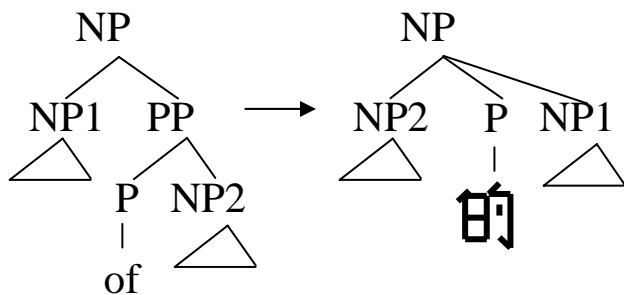
Non-local Re-Ordering (*English/Arabic*)



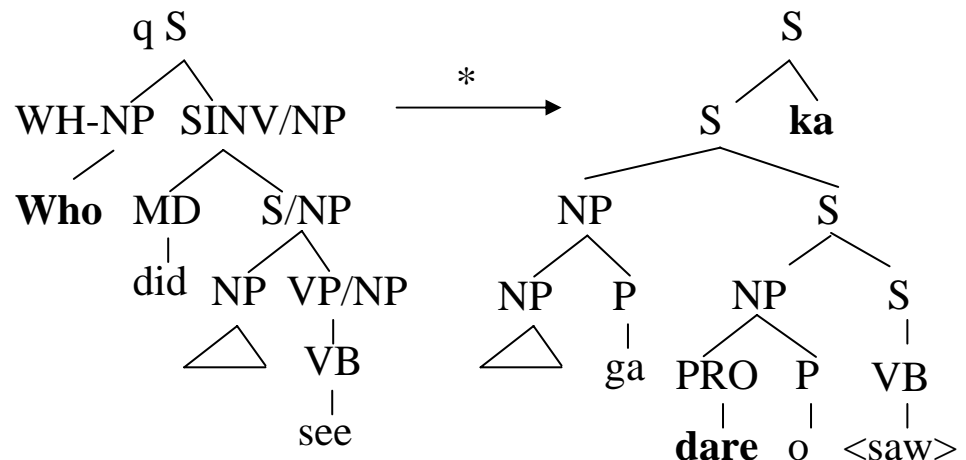
Non-constituent Phrasal Translation (*English/Spanish*)



Lexicalized Re-Ordering (*English/Chinese*)



Long-distance Re-Ordering (*English/Japanese*)



Present and Future

- Phrase-based models have been state-of-the-art for six years
 - Word alignments
 - Phrase pair extraction & probabilities
 - N-gram language models
 - Beam search decoding
 - Feature functions & learning weights
- In 2006, statistical models using syntax outperformed phrase-based models for the first time! Only just beginning...

the end