

CSCI 5832

Natural Language Processing

Lecture 25
Jim Martin

4/24/07

CSCI 5832 Spring 2007

1

Machine Translation

Slides stolen from Kevin Knight (USC/ISI)

4/24/07

CSCI 5832 Spring 2007

2

Today 4/24

- Machine translation framework
- State of the art
- Evaluation methods
- Word-based models

4/24/07

CSCI 5832 Spring 2007

3

Recent Progress

2002

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.



2003

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.

4/24/07

CSCI 5832 Spring 2007

4

Commercial Applications

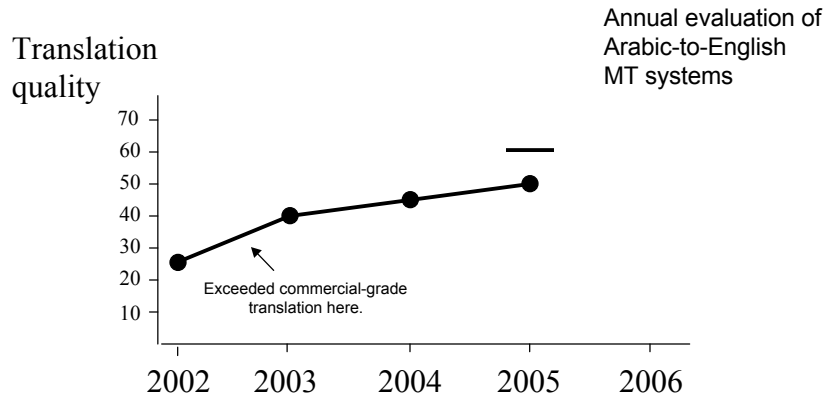
The screenshot shows a web-based interface for a news broadcast. The main content area displays a news article with a video player. Annotations with yellow boxes and pink arrows point to specific features:

- news broadcast**: Points to the top navigation bar.
- foreign language speech recognition**: Points to the 'Words' section on the right, which lists words in Arabic and English.
- English translation**: Points to the main text area, which shows an English translation of the Arabic text.
- searchable archive**: Points to the left sidebar, which contains a list of news bulletins.

Statistical Machine Translation

The diagram illustrates the concept of Statistical Machine Translation. It shows a man sitting at a desk with a laptop, looking bored. A thought bubble above him says, "Man, this is so boring." Another thought bubble above him says, "Hmm, every time he sees 'banco', he either types 'bank' or 'bench' ... but if he sees 'banco de...', he always types 'bank', never 'bench'...". Below the man, there are three boxes representing translated documents, and a laptop icon. The text "Translated documents" is written below the boxes. The date "4/24/07" is written below the man, and "CSCI 5832 Spring 2007" is written below the boxes. The number "6" is written below the laptop icon.

Things are Consistently Improving

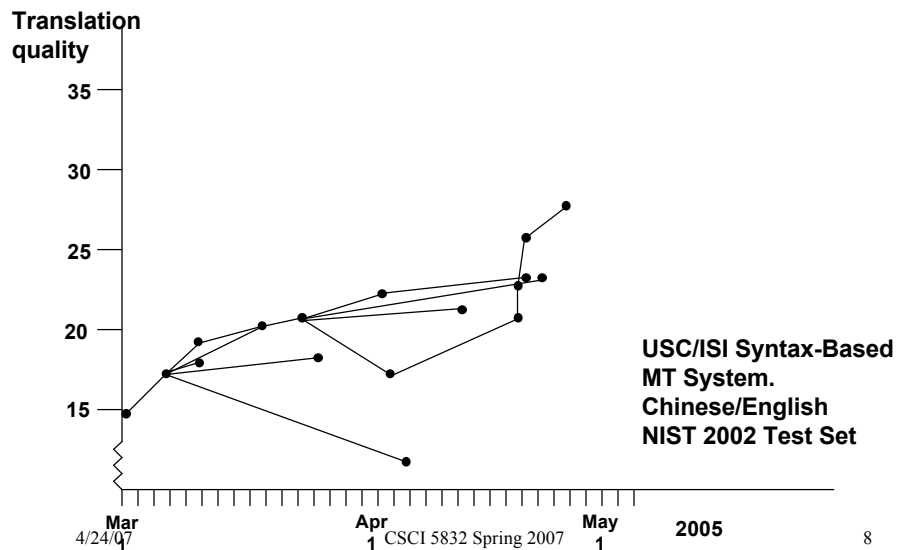


4/24/07

CSCI 5832 Spring 2007

7

Progress Driven by Experiments with Empirical Measures of Success



Mar 4/24/07

Apr 1 CSCI 5832 Spring 2007

May 1 2005

8

Current Approaches

- Same old noisy channel model...
- If we're translating French to English the French we're seeing is just a weird garbled version of English
- There must have been some process that generated the French from the original English
- The key is to **decode** the garbles back into the original English by...
- $\text{Argmax } P(E | F)$ by Bayes
- A very old idea

4/24/07

CSCI 5832 Spring 2007

9

Warren Weaver (1947)

When I look at an article in Russian, I say to myself: This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.



4/24/07

CSCI 5832 Spring 2007

10

MT History

- **Earliest applications of NLP/AI**
- **Didn't work for a lot of reasons**
- **Led to an NLP/AI winter where funding was hard to come by for a long time.**
- **MT was a topic no one would touch in respected circles...**

4/24/07

CSCI 5832 Spring 2007

11

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 ghirok .	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 ghirok 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 .

4/24/07

CSCI 5832 Spring 2007

12

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 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 .

4/24/07

CSCI 5832 Spring 2007

13

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 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 .

4/24/07

CSCI 5832 Spring 2007

14

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 ghirok .	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 ghirok 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 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 .

4/24/07

CSCI 5832 Spring 2007

15

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 ghirok .	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 ghirok 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 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 .

4/24/07

CSCI 5832 Spring 2007

16

Centauri/Arcturan [Knight, 1997]

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 eneac .
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 ghirok .	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 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 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** 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 eneac .
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 ghirok .	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 ghirok 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 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** 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 ghirok .	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 ghirok 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 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** 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 ghirok .	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 ghirok 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 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 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 eneāt .
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 ghirok .	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 ghirok 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 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 .

process of elimination

CSCI 5832 Spring 2007

21

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 eneāt .
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 ghirok .	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 ghirok 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 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 .

cognate?

CSCI 5832 Spring 2007

22

Centauri/Arcturan [Knight, 1997]

Your assignment, put these words in order: { jjat, arrat, mat, bat, oloat, at-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 eneat .
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 oloat 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 crrrok hihok yorok zanzanok . 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 .

4/24/07

CSCI 5832 Spring 2007

23

Spanish/English text

Translate: Clients do not sell pharmaceuticals in Europe.

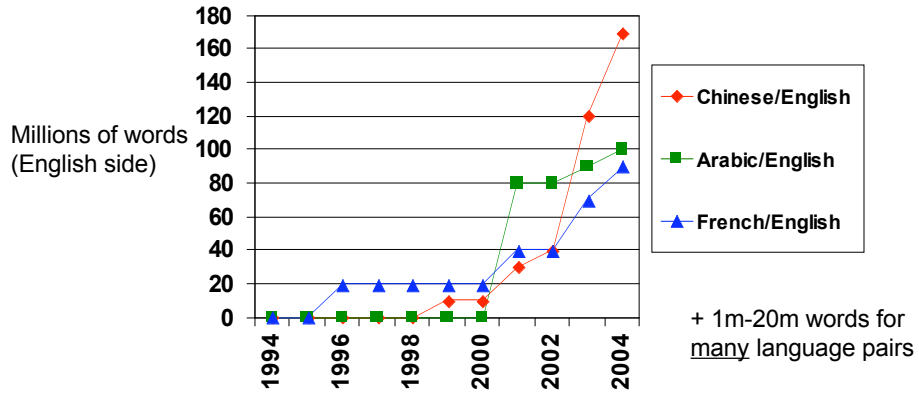
1a. Garcia and associates . 1b. Garcia y asociados .	7a. the clients and the associates are enemies . 7b. los clientes y los asociados son enemigos .
2a. Carlos Garcia has three associates . 2b. Carlos Garcia tiene tres asociados .	8a. the company has three groups . 8b. la empresa tiene tres grupos .
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 .

4/24/07

CSCI 5832 Spring 2007

24

Bilingual Training Data



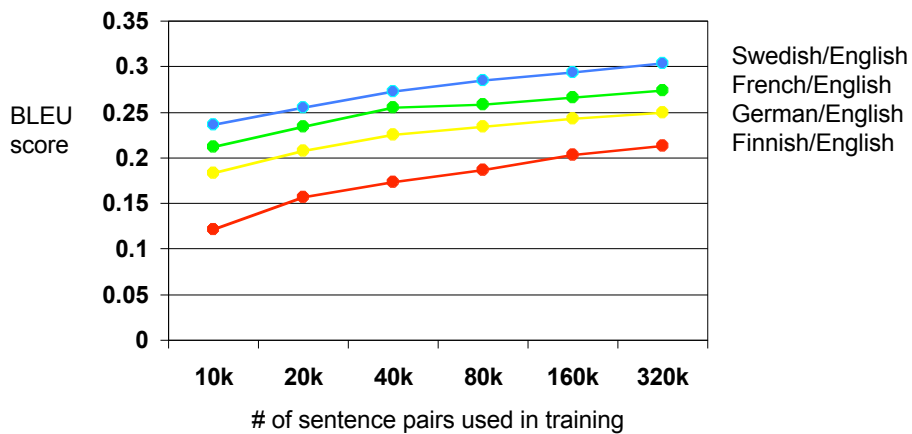
(Data stripped of formatting, in sentence-pair format, available from the Linguistic Data Consortium at UPenn).

4/24/07

CSCI 5832 Spring 2007

25

Sample Learning Curves



4/24/07

CSCI 5832 Spring 2007

Experiments by
Philipp Koehn

26

MT Evaluation

Traditionally difficult because there is no "right answer".

20 human translators will translate the same sentence 20 different ways.

4/24/07

CSCI 5832 Spring 2007

27

Evaluation Metric (BLEU)

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:

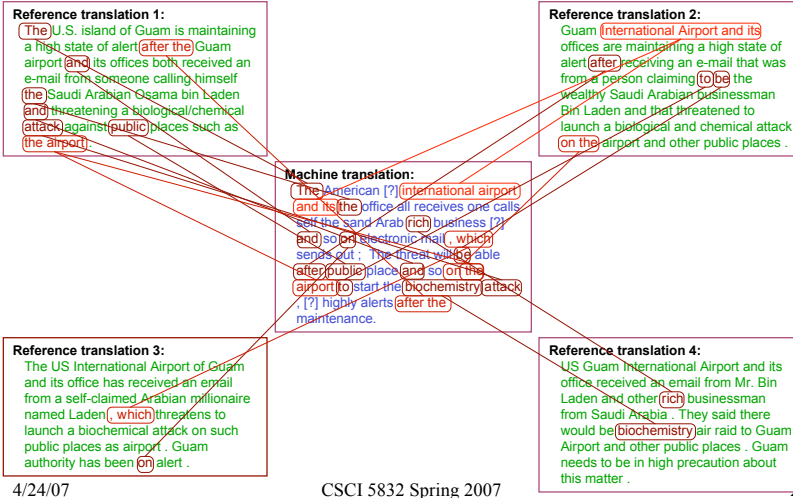
The American [?] international airport and its the office at receives one calls self the sand Arab rich business [?] and so on electronic mail , which sends out ; The threat will be able after public place and so on the airport to start the biochemistry attack , [?] highly alerts after the maintenance.

- 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)
 - Contra doesn't hold. Can find perfectly good translations that hurt, or don't help, BLEU

CSCI 5832 Spring 2007

28

Multiple Reference Translations



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 | #7 |
| al by the police . | #8 |
| the ringer is killed by the police . | #9 |
| police killed the gunman . | #10 |

4/24/07

CSCI 5832 Spring 2007

30

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 #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!)

4/24/07

CSCI 5832 Spring 2007

31

NIST 2006 Results

Arabic-to-English Results

Large Data Track

NIST Subset

Overall BLEU Scores

Site ID	BLEU-4
google	0.4281
ibm	0.3954
isi	0.3908
rwth	0.3906
aptek*#	0.3874
lw	0.3741
bbn	0.3690
ntt	0.3680
itcirst	0.3466
cmu-uka	0.3369
umd-jhu	0.3333
edinburgh*#	0.3303
sakhr	0.3296

Chinese-to-English Results

Large Data Track

NIST Subset

Overall BLEU Scores

Site ID	BLEU-4
isi	0.3393
google	0.3316
lw	0.3278
rwth	0.3022
ict	0.2913
edinburgh*#	0.2830
bbn	0.2781
nrc	0.2762
itcirst	0.2749
umd-jhu	0.2704

4/24/07

CSCI 5832 Spring 2007

32

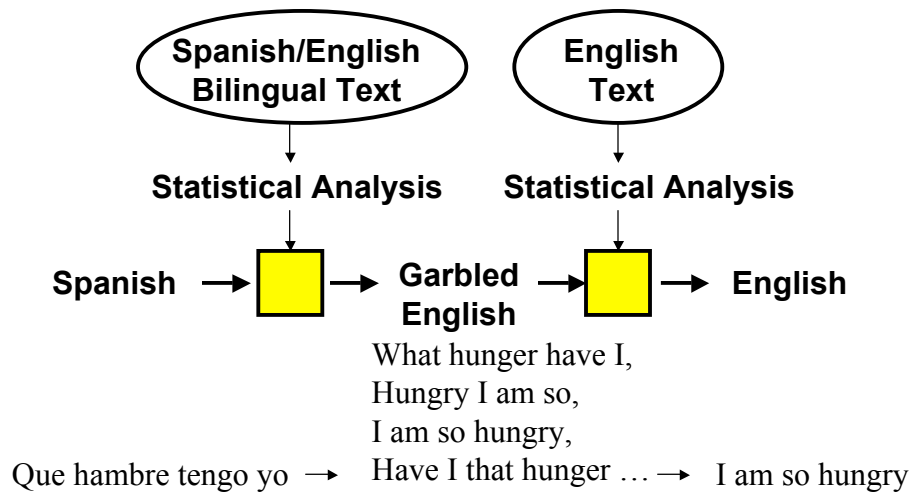
Word-Based Statistical MT

4/24/07

CSCI 5832 Spring 2007

33

Statistical MT Systems

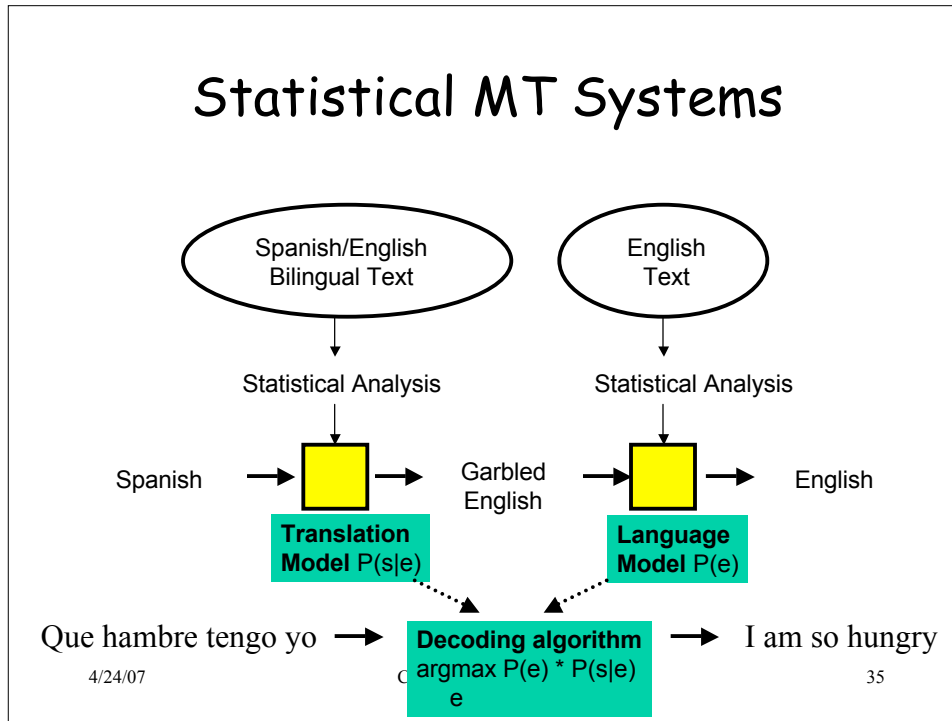


4/24/07

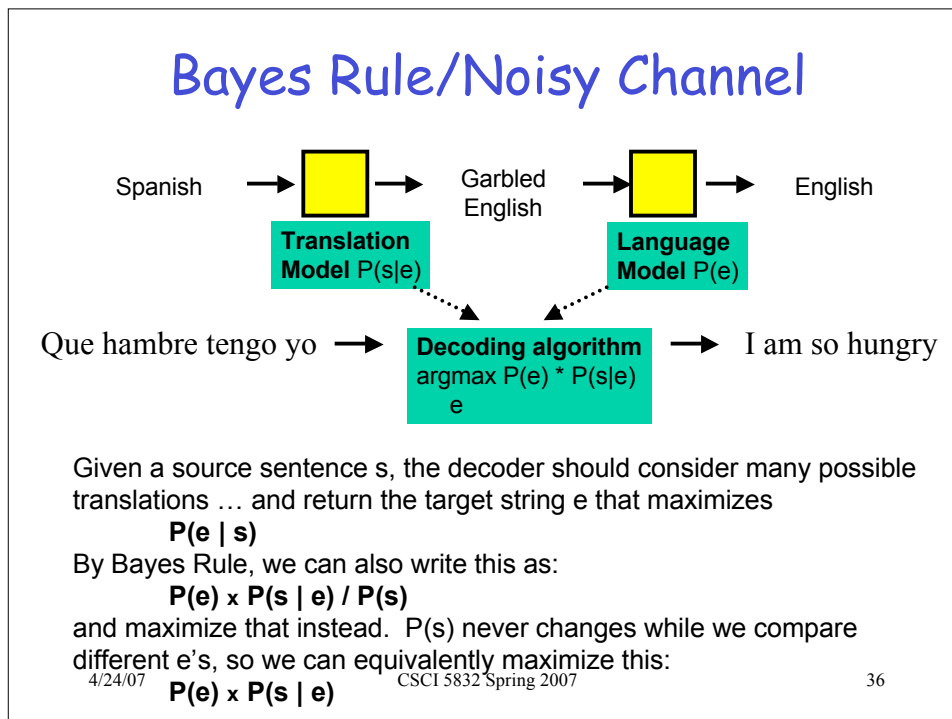
CSCI 5832 Spring 2007

34

Statistical MT Systems



Bayes Rule/Noisy Channel



Three Sub-Problems of Statistical MT

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

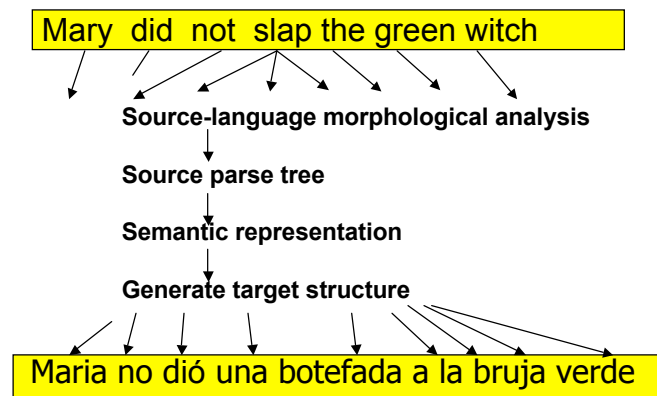
4/24/07

CSCI 5832 Spring 2007

37

Translation Model

Generative story:



4/24/07

CSCI 5832 Spring 2007

38

Translation Model?

Generative story:



4/24/07

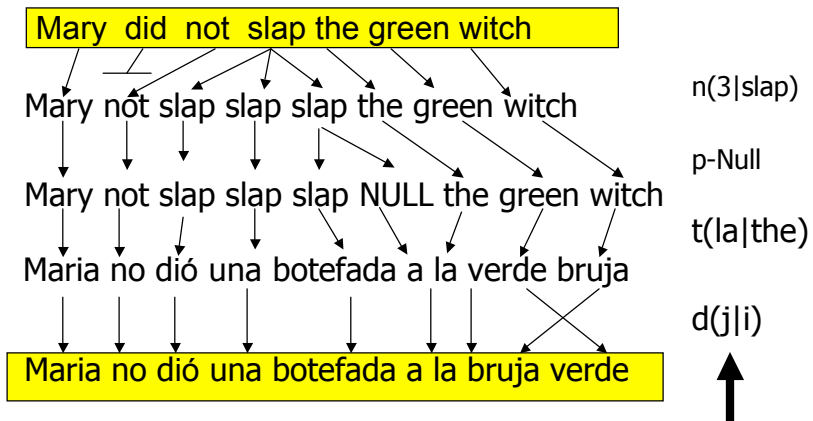
CSCI 5832 Spring 2007

39

The Classic Translation Model

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

Generative story:



4/24/07

CSCI 5832 Spring 2007

40

Parts List

- **We need probabilities for**
 - $n(x|y)$ The probability that word y will yield x outputs in the translation... (fertility)
 - p The probability of a null insertion
 - t The actual word translation probability table
 - $d(j|i)$ the probability that a word at position i will make an appearance at position j in the translation

4/24/07

CSCI 5832 Spring 2007

41

Parts List


- **Every one of these can be learned from a sentence aligned corpus...**
 - Ie. A corpus where sentences are paired but nothing else is specified
- **And the EM algorithm**

4/24/07

CSCI 5832 Spring 2007

42

Word Alignment

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

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


- Assume that All word alignments equally likely.
- I.e. that all $P(\text{french-word} \mid \text{english-word})$ are equal
- Recall that we want $P(f|e)$

4/24/07

CSCI 5832 Spring 2007

43

Word Alignment

... 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.


4/24/07

CSCI 5832 Spring 2007

44

Word Alignment

... la maison ... la maison bleue ... la fleur ...
... the house ... the blue house ... the flower ...




4/24/07

CSCI 5832 Spring 2007

45

Word Alignment

... la maison ... la maison bleue ... la fleur ...
... the house ... the blue house ... the flower ...



settling down after another iteration

4/24/07

CSCI 5832 Spring 2007

46

Word Alignment

... la maison ... la maison bleue ... la fleur ...
... the house ... the blue house ... the flower ...

Inherent hidden structure revealed by EM training

4/24/07

CSCI 5832 Spring 2007

47

Parts List

- Given a sentence alignment we can **induce a word alignment**
- Given that word alignment we can get the p , t , d and n parameters we need for the model.
- I.e. We can $\text{argmax} P(e|f)$ by maxing over $P(f|e)*P(e)...$ and we can do that by iterating over some large space of f possibilities.

4/24/07

CSCI 5832 Spring 2007

48

Decoding

- **Remember Viterbi? Just a fancier Viterbi**
 - Given foreign sentence f , find English sentence e that maximizes $P(e) \times P(f | e)$

4/24/07

CSCI 5832 Spring 2007

49

Decoding

Que	hambre	tengo	yo
what	hunger	have	I
that	hungry	am	me
so		make	
where			

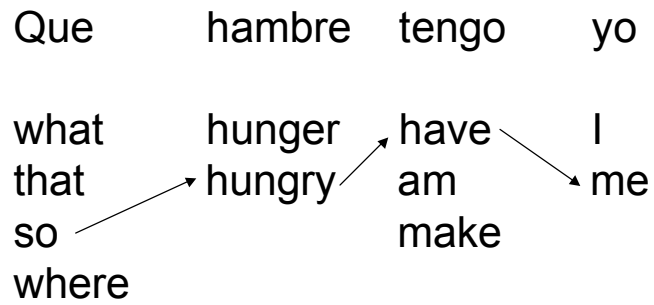
4/24/07

CSCI 5832 Spring 2007

50

Decoding

Que	hambre	tengo	yo
what	hunger	have	I
that	hungry	am	me
so		make	
where			



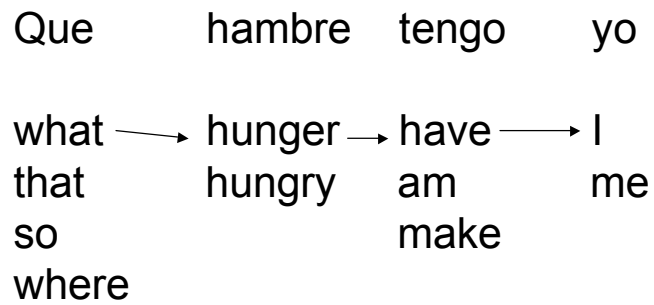
4/24/07

CSCI 5832 Spring 2007

51

Decoding

Que	hambre	tengo	yo
what	hunger	have	I
that	hungry	am	me
so		make	
where			



4/24/07

CSCI 5832 Spring 2007

52

Decoding

Que hambre tengo yo

what hunger have I

that hungry am me

so make

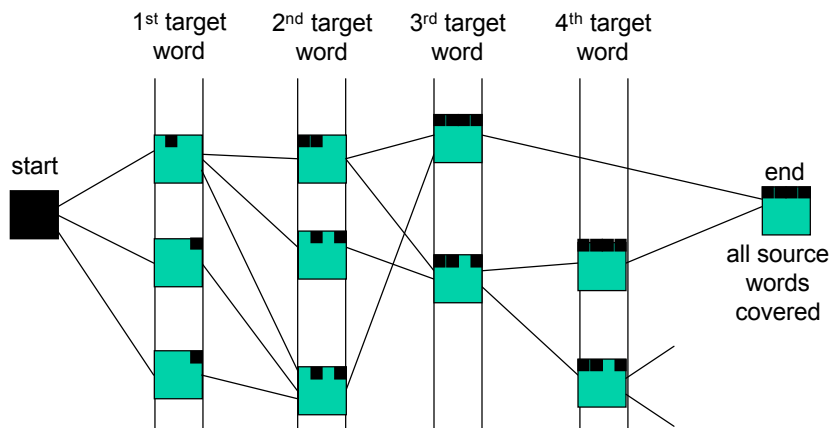
where

4/24/07

CSCI 5832 Spring 2007

53

Decoder: Actually Translates New Sentences

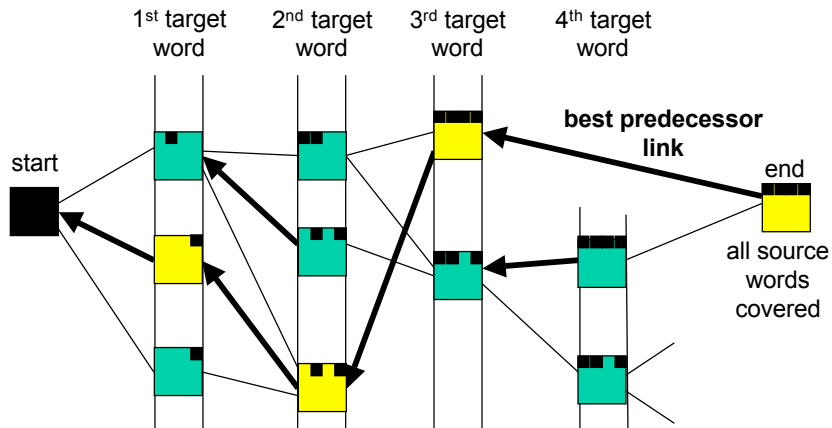


Each partial translation hypothesis contains:

- 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)]

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 foreign 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

4/24/07

CSCI 5832 Spring 2007

57

Next Time

- **EM alignment example**
- **Phrase-based translation**

4/24/07

CSCI 5832 Spring 2007

58