

Machine Translation: a historical perspective

David Kauchak
CS159 – Fall 2020

Philipp Koehn
School of Informatics
University of Edinburgh

Some slides adapted from
Kevin Knight
USC/Information Sciences Institute
USC/Computer Science Department

Dan Klein
Computer Science Department
UC Berkeley

1

Admin

Assignment 5 out

Quiz #2

2

Language translation



3

MT Systems

Where have you seen machine translation systems?

4

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裔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。



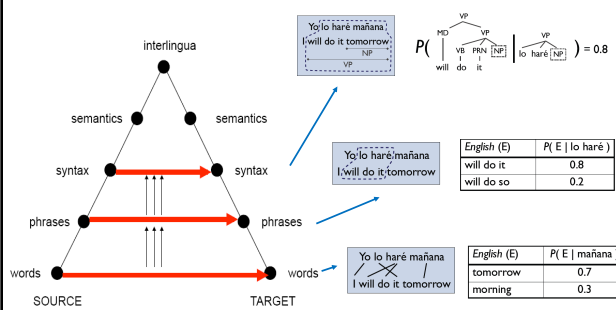
The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

A good test for natural language processing.

Requires capabilities in both interpretation and generation.

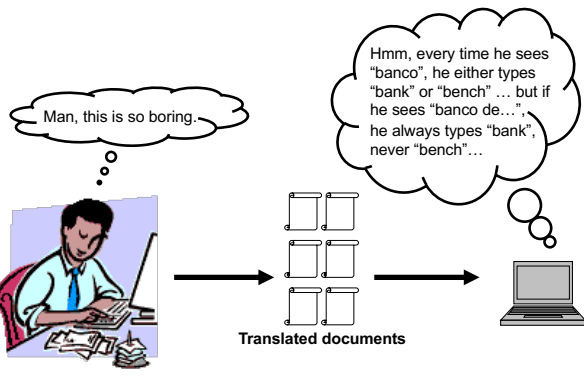
5

Levels of Transfer



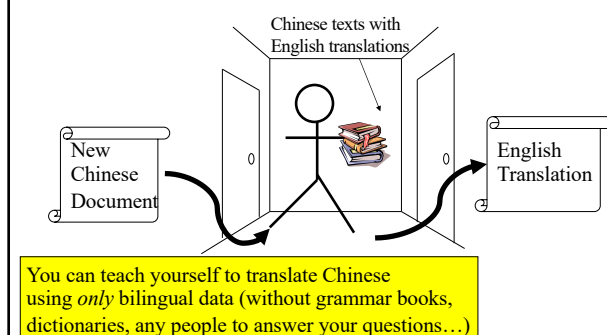
6

Data-Driven Machine Translation



9

Welcome to the Chinese Room



10

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok 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 errok 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 .

11

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok 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 errok 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 .

12

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok 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 errok 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 .

13

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok errok 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 errok 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 .

14

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **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 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 errrok 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 .

15

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **hihok** **yorok** **elok** **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 clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok 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 .

16

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **hihok** **yorok** **elok** **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 clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok 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 .

17

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **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 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 errrok 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 .

18

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **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 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 errrok 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 .

19

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **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 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 errrok 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 .

20

Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: **farok** **errrok** **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 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 errrok 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 .

21

Centauri/Arcturan [Knight, 1997]

Your assignment, put these words in order: { **jjat**, **arrat**, **mat**, **bat**, **olloat**, **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 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 errrok 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 .

22

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 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 .

23

Data available

Many languages

- Europarl corpus has all European languages
 - <http://www.statmt.org/europarl/>
 - From a few hundred thousand sentences to a few million
- French/English from French parliamentary proceedings
- Lots of Chinese/English and Arabic/English from government projects/interests
 - Chinese-English: Hundreds of millions of sentence pairs)
 - Arabic-English: ~One hundred million sentence pairs
- Smaller corpora in many, many other languages

Lots of monolingual data available in many languages

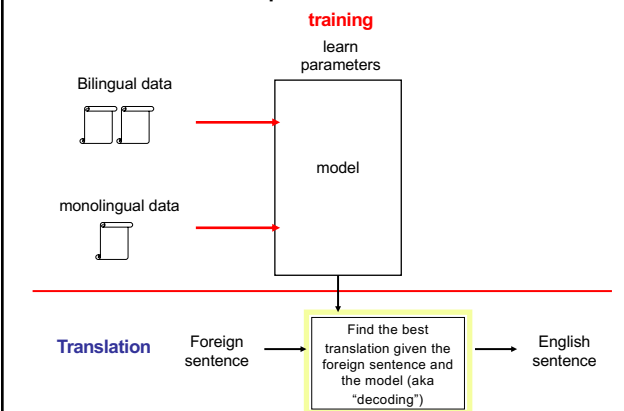
Even less data with multiple translations available

Available in limited domains

- most data is either news or government proceedings
- some other domains recently, like blogs

24

Historical Perspective: statistical MT



25

Why the historical perspective?

Allow us to contrast the difference with more recent approaches

- Historically, we had to model all phenomena explicitly
- Newer (network models) learn this automatically!

Still used in some low-resource situations

Can be useful models for related applications (e.g., word alignment)

26

Statistical MT

We will model the translation process probabilistically

Given a foreign sentence to translate, for any possible English sentence, we want to know the probability that the sentence is a translation of the foreign sentence

If we can find the most probable English sentence, we're done

$$p(\text{english sentence} \mid \text{foreign sentence})$$

27

Translation

Probabilistic model: $p(e \mid f)$ $p(\text{English} \mid \text{Foreign})$

What is the translation problem then?

$$\text{translation}(f) = \arg_e \max p(e \mid f)$$

28

Noisy channel model

$$p(e \mid f) = \frac{p(f \mid e)p(e)}{p(f)} \quad \text{Bayes' rule}$$

$p(f)$ probability of the foreign sentence

$p(e)$ language model: what are likely English word sequences?

$p(f \mid e)$ translation model: how does the translation process happen? probability of the translated English sentence given the foreign sentence

29

Noisy channel model

$$p(e \mid f) = p(f \mid e)p(e) \quad \text{Bayes' rule}$$

~~$p(f)$ probability of the foreign sentence~~ why?

$p(e)$ language model: what are likely English word sequences?

$p(f \mid e)$ translation model: how does the translation process happen? probability of the translated English sentence given the foreign sentence

30

Noisy channel model

$$p(e | f) = p(f | e)p(e) \quad \text{Bayes' rule}$$

~~$p(f)$~~ probability of the foreign sentence why?

$$\text{translation}(f) = \arg_e \max \frac{p(f | e)p(e)}{p(f)} = \arg_e \max p(f | e)p(e)$$

this is a constant
for any given f

31

Noisy channel model

$$\text{model } p(e | f) \propto p(f | e)p(e)$$

translation model

how do English
sentences get
translated to
foreign?

language model

what do English
sentences look
like?

32

Translation model

The models define probabilities over inputs

$$p(f | e)$$

Morgen fliege ich nach Kanada zur Konferenz

Tomorrow I will fly to the conference in Canada

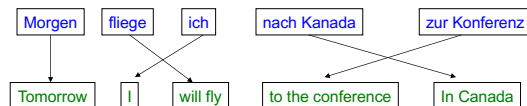
What is the probability that the English sentence is
a translation of the foreign sentence?

33

Translation model

The models define probabilities over inputs

$$p(f | e)$$



- What is the probability of a foreign word being translated as a particular English word?
- What is the probability of a foreign phrase being translated as a particular English phrase?
- What is the probability of a word/phrase changing ordering?
- What is the probability of a foreign word/phrase disappearing?
- What is the probability of an English word/phrase appearing?

34

Translation model

The models define probabilities over inputs

$$p(f | e)$$

$$p(\text{Morgen fliege ich nach Kanada zur Konferenz} | \text{Tomorrow I will fly to the conference in Canada}) = 0.1$$

$$p(\text{Morgen fliege ich nach Kanada zur Konferenz} | \text{I like peanut butter and jelly}) = 0.0001$$

35

Language model

The models define probabilities over inputs

$$p(e)$$

Tomorrow I will fly to the conference in Canada

36

What is a probability distribution?

A probability distribution defines the probability over a space of possible inputs

For the language model, what is the space of possible inputs?

- A language model describes the probability over **ALL** possible combinations of English words

For the translation model, what is the space of possible inputs?

- **ALL** possible combinations of foreign words with **ALL** possible combinations of English words

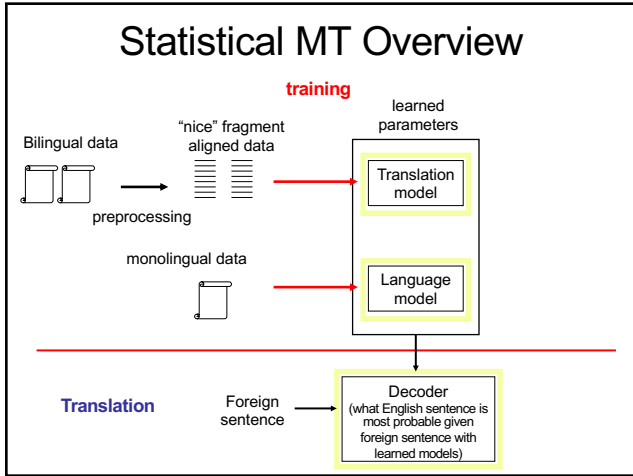
37

One way to think about it...

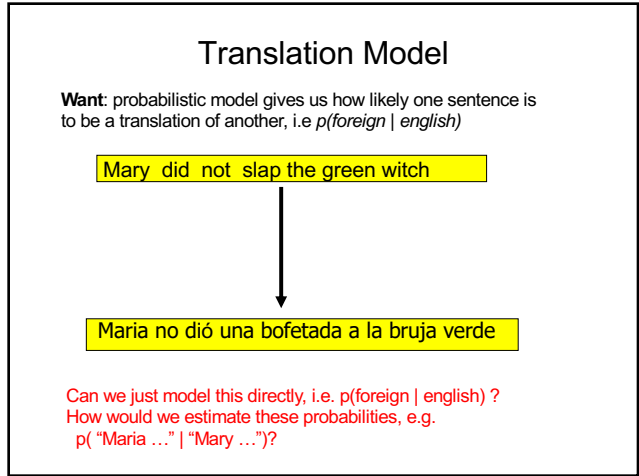


Que hambre tengo yo → What hunger have I,
Hungry I am so, → I am so hungry
I am so hungry,
Have I that hunger ...

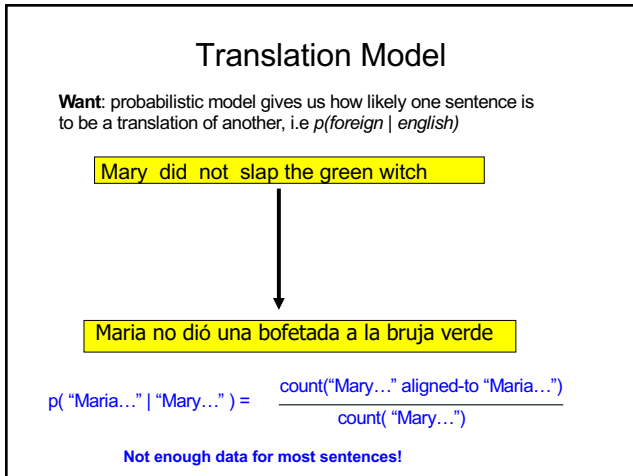
38



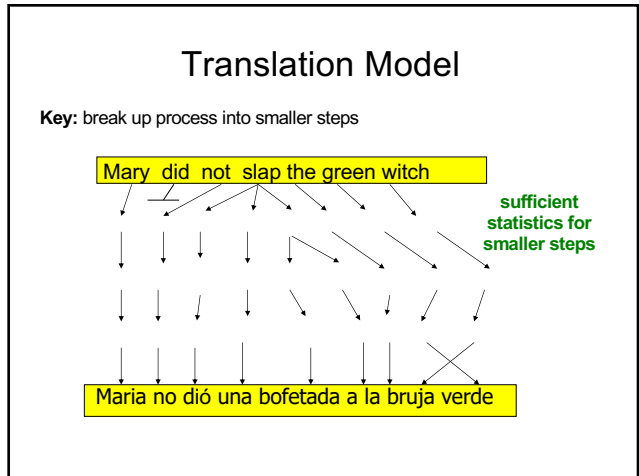
39



45

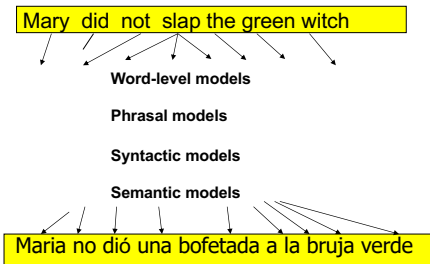


46



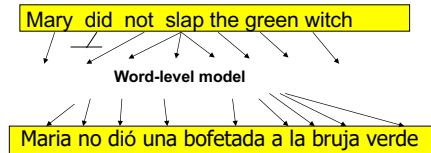
47

What kind of Translation Model?



48

IBM Word-level models

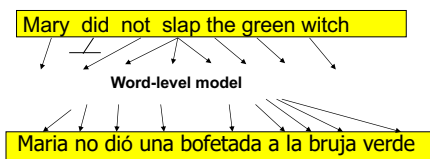


Generative story: description of how the translation happens

1. Each English word gets translated as 0 or more Foreign words
2. Some additional foreign words get inserted
3. Foreign words then get shuffled

49

IBM Word-level models



Each foreign word is *aligned* to exactly one English word.

Key idea: decompose $p(\text{foreign} | \text{english})$ into word translation probabilities of the form $p(\text{foreign_word} | \text{english_word})$

IBM described 5 different levels of models with increasing complexity (and decreasing independence assumptions)

50

Some notation

$E = e_1 e_2 \dots e_{|E|}$ English sentence with length $|E|$

$F = f_1 f_2 \dots f_{|F|}$ Foreign sentence with length $|F|$

Mary did not slap the green witch
 $e_1 \quad e_2 \quad e_3 \quad e_4 \quad e_5 \quad e_6 \quad e_7$

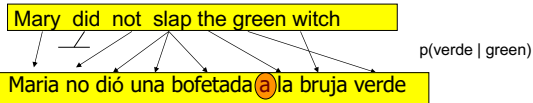
$f_1 \quad f_2 \quad f_3 \quad f_4 \quad f_5 \quad f_6 \quad f_7 \quad f_8 \quad f_9$

Maria no dió una bofetada a la bruja verde

Translation model: $p(F | E) = p(f_1 f_2 \dots f_{|F|} | e_1 e_2 \dots e_{|E|})$

51

Word models: IBM Model 1



$p(\text{verde} \mid \text{green})$

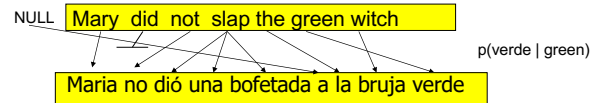
Each foreign word is aligned to exactly one English word

This is the **ONLY** thing we model!

Does the model handle foreign words that are not aligned, e.g. "a"?

52

Word models: IBM Model 1



$p(\text{verde} \mid \text{green})$

Each foreign word is aligned to exactly one English word

This is the **ONLY** thing we model!

Include a "NULL" English word and align to this to account for deletion

53

Word models: IBM Model 1

generative story -> probabilistic model

- Key idea: introduce "hidden variables" to model the word alignment

$$p(f_1 f_2 \dots f_{|F|} \mid e_1 e_2 \dots e_{|E|})$$

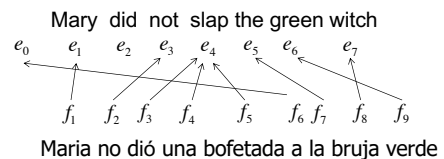


$$p(f_1 f_2 \dots f_{|F|}, a_1 a_2 \dots a_{|F|} \mid e_1 e_2 \dots e_{|E|})$$

- one variable for each foreign word
- a_i corresponds to the i th foreign word
- each a_i can take a value $0 \dots |E|$

54

Alignment variables



a_1	1
a_2	3
a_3	4
a_4	4
a_5	4
a_6	0
a_7	5
a_8	7
a_9	6

55

Alignment variables

And the program has been implemented

e_0 e_1 e_2 e_3 e_4 e_5 e_6

Alignment?

f_1 f_2 f_3 f_4 f_5 f_6 f_7

Le programme a ete mis en application

56

Alignment variables

And the program has been implemented

e_0 e_1 e_2 e_3 e_4 e_5 e_6
 f_1 f_2 f_3 f_4 f_5 f_6 f_7

Le programme a ete mis en application

a_1	?
a_2	?
a_3	?
a_4	?
a_5	?
a_6	?
a_7	?

57

Alignment variables

And the program has been implemented

e_0 e_1 e_2 e_3 e_4 e_5 e_6

f_1 f_2 f_3 f_4 f_5 f_6 f_7

Le programme a ete mis en application

a_1	2
a_2	3
a_3	4
a_4	5
a_5	6
a_6	6
a_7	6

58

Probabilistic model

$$p(f_1 f_2 \dots f_{|F|} | e_1 e_2 \dots e_{|E|}) \stackrel{?}{=} p(f_1 f_2 \dots f_{|F|}, a_1 a_2 \dots a_{|F|} | e_1 e_2 \dots e_{|E|})$$

NO!

$$p(f_1 f_2 \dots f_{|F|}, a_1 a_2 \dots a_{|F|} | e_1 e_2 \dots e_{|E|}) \longrightarrow p(f_1 f_2 \dots f_{|F|} | e_1 e_2 \dots e_{|E|})$$

How do we get rid of variables?

59

Joint distribution

NLPPass, EngPass	P(NLPPass, EngPass)
true, true	.88
true, false	.01
false, true	.04
false, false	.07

What is P(ENGPass)?

60

Joint distribution

NLPPass, EngPass	P(NLPPass, EngPass)
true, true	.88
true, false	.01
false, true	.04
false, false	.07

0.92

How did you figure that out?

61

Joint distribution

$$P(x) = \sum_{y \in Y} p(x, y)$$

Called "marginalization", aka summing over a variable

NLPPass, EngPass	P(NLPPass, EngPass)
true, true	.88
true, false	.01
false, true	.04
false, false	.07

62

Probabilistic model

$$p(f_1, f_2, \dots, f_{|F|} | e_1, e_2, \dots, e_{|E|}) = \sum_{a_1} \sum_{a_2} \dots \sum_{a_{|F|}} p(f_1, f_2, \dots, f_{|F|}, a_1, a_2, \dots, a_{|F|} | e_1, e_2, \dots, e_{|E|})$$

Sum over all possible values, i.e. marginalize out the alignment variables

63

Independence assumptions

IBM Model 1:

$$p(f_1 f_2 \dots f_{|f|}, a_1 a_2 \dots a_{|f|} | e_1 e_2 \dots e_{|e|}) = \prod_{i=1}^{|f|} p(f_i | e_{a_i})$$

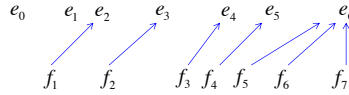
What independence assumptions are we making?

What information is lost?

64

$$p(f_1 f_2 \dots f_{|f|}, a_1 a_2 \dots a_{|f|} | e_1 e_2 \dots e_{|e|}) = \prod_{i=1}^{|f|} p(f_i | e_{a_i})$$

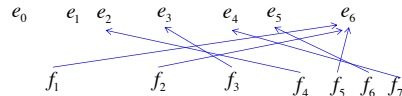
And the program has been implemented



Le programme a ete mis en application

Are the probabilities any different under model 1?

And the program has been implemented

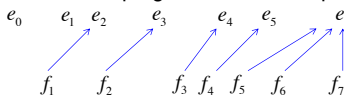


application en programme Le mis ete a

65

$$p(f_1 f_2 \dots f_{|f|}, a_1 a_2 \dots a_{|f|} | e_1 e_2 \dots e_{|e|}) = \prod_{i=1}^{|f|} p(f_i | e_{a_i})$$

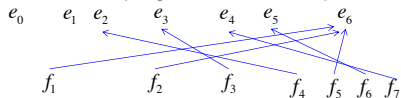
And the program has been implemented



Le programme a ete mis en application

No. Model 1 ignores word order!

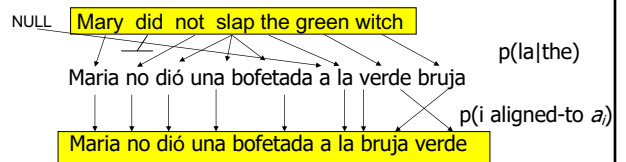
And the program has been implemented



application en programme Le mis ete a

66

IBM Model 2



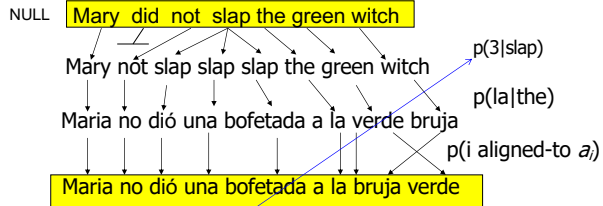
$$p(f_1 f_2 \dots f_{|f|}, a_1 a_2 \dots a_{|f|} | e_1 e_2 \dots e_{|e|}) = \prod_{i=1}^{|f|} p(i \text{ aligned-to } a_i) p(f_i | e_{a_i})$$

Models word movement by position, e.g.

- Words don't tend to move too much
- Words at the beginning move less than words at the end

67

IBM Model 3



Incorporates "fertility": how likely a particular English word is to produce multiple foreign words

68

Word-level models

Problems/concerns?

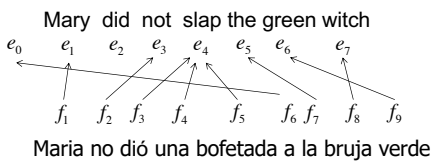
- Multiple English words for one French word
 - IBM models can do one-to-many (fertility) but not many-to-one
- Phrasal Translation
 - "real estate", "note that", "interest in"
- Syntactic Transformations
 - Verb at the beginning in Arabic
 - Translation model penalizes any proposed re-ordering
 - Language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

69

Benefits of word-level model

Rarely used in practice for modern MT systems

Why talk about them?



Two key side effects of training a word-level model:

- Word-level alignment
- $p(f | e)$: translation dictionary

70

Training a word-level model

$$p(f_1, f_2, \dots, f_{|F|}, a_1, a_2, \dots, a_{|F|} | e_1, e_2, \dots, e_{|E|}) = \prod_{i=1}^{|F|} p(f_i | e_{a_i})$$

Where do these come from?

Have to learn them!

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

71