# Statistical NLP
## Spring 2007

University of California
C A L
N L P
Berkeley

Lecture 10: Word Alignment

Dan Klein – UC Berkeley

---

## Machine Translation: Examples

**Atlanta, preso il killer del palazzo di Giustizia**

**ATLANTA** - La grande paura che per 26 ore ha attanagliato Atlanta è finita: Brian Nichols, l'uomo che aveva ucciso tre persone a palazzo di Giustizia e che ~~~~~~~~~~~~~~~~~~~~~, s'è consegnato alla polizia, dopo avere cercato rifugio nell'alloggio di una donna in un complesso d'appartamenti alla periferia della città. Per tutto il giorno, il centro della città, sede della ~~~~~~~~~ dei Giochi 1996, cuore di una popolosa area metropolitana, era rimasto paralizzato.

**Atlanta, taken the killer of the palace of Justice**

**ATLANTA** - The great fear that for 26 hours has gripped Atlanta is ended: Brian Nichols, the man who had killed three persons to palace of Justice and that ~~~~~~~~~~~~~~~~~~~~~, s' is delivered to the police, after to have tried shelter in the lodging of one woman in a complex of apartments to the periphery of the city. For all the day, the center of the city, center of the ~~~~~~~~~ and of Giochi 1996, heart of one popolosa metropolitan area, was remained paralyzed.

---

## Machine Translation

Madame la présidente, votre présidence de cette institution a été marquante.
Mrs Fontaine, your presidency of this institution has been outstanding.
Madam President, president of this house has been discoveries.
Madam President, your presidency of this institution has been impressive.

Je vais maintenant m'exprimer brièvement en irlandais.
I shall now speak briefly in Irish .
I will now speak briefly in Ireland .
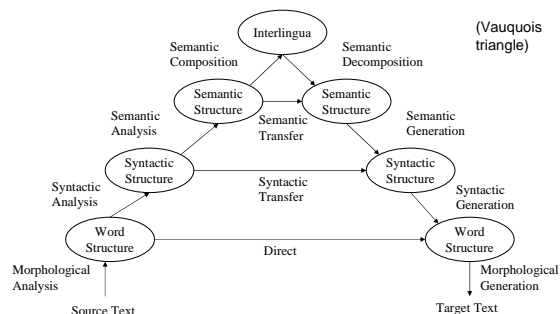I will now speak briefly in Irish .

Nous trouvons en vous un président tel que nous le souhaitions.
We think that you are the type of president that we want.
We are in you a president as the wanted.
We are in you a president as we the wanted.

---

## History

- 1950's: Intensive research activity in MT
- 1960's: Direct word-for-word replacement
- 1966 (ALPAC): NRC Report on MT
  - Conclusion: MT no longer worthy of serious scientific investigation.
- 1966-1975: `Recovery period'
- 1975-1985: Resurgence (Europe, Japan)
- 1985-present: Gradual Resurgence (US)

http://ourworld.compuserve.com/homepages/WJHutchins/MTS-93.htm

---

## Levels of Transfer



(Vauquois triangle)

Interlingua

Semantic Composition — Semantic Decomposition

Semantic Structure — Semantic Structure

Semantic Analysis — Semantic Transfer — Semantic Generation

Syntactic Structure — Syntactic Structure

Syntactic Analysis — Syntactic Transfer — Syntactic Generation

Word Structure — Direct — Word Structure

Morphological Analysis — Morphological Generation

Source Text — Target Text

---

## General Approaches

- Rule-based approaches
  - Expert system-like rewrite systems
  - Interlingua methods (analyze and generate)
  - Lexicons come from humans
  - Can be very fast, and can accumulate a lot of knowledge over time (e.g. Systran)

- Statistical approaches
  - Word-to-word translation
  - Phrase-based translation
  - Syntax-based translation (tree-to-tree, tree-to-string)
  - Trained on parallel corpora
  - Usually noisy-channel (at least in spirit)

## The Coding View

- "One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.' "
  - Warren Weaver (1955:18, quoting a letter he wrote in 1947)

---

## MT System Components

Language Model

Translation Model

source P(e) → e → channel P(f|e) → f

best e ← decoder ← observed f

$$\text{argmax } P(e|f) = \text{argmax } P(f|e)P(e)$$
$$e \qquad\qquad e$$

*Finds an English translation which is both fluent and semantically faithful to the French source*

---

## Today

- The components of a simple MT system
  - You already know about the LM
  - Word-alignment based TMs
    - IBM models 1 and 2, HMM model
  - A simple decoder

- Next few classes
  - More complex word-level and phrase-level TMs
  - Tree-to-tree and tree-to-string TMs
  - More sophisticated decoders

---

## Word Alignment

**x**

What is the anticipated cost of collecting fees under the new proposal?

En vertu des nouvelles propositions, quel est le coût prévu de perception des droits?

**z**

What
is
the
anticipated
cost
of
collecting
fees
under
the
new
proposal
?

En
vertu
de
les
nouvelles
propositions
,
quel
est
le
coût
prévu
de
perception
de
les
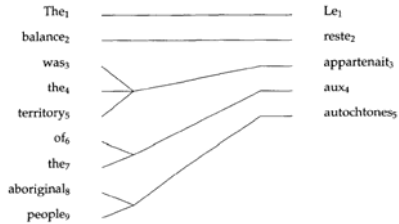droits
?

---

## Unsupervised Word Alignment

- Input: a **bitext**: pairs of translated sentences

  nous acceptons votre opinion .
  we accept your view .

- Output: **alignments**: pairs of translated words
  - When words have unique sources, can represent as a (forward) alignment function a from French to English positions

  nous
  acceptons
  votre
  opinion
  .

  we accept your view

---

## 1-to-Many Alignments

$\text{And}_1 \quad \text{the}_2 \quad \text{program}_3 \quad \text{has}_4 \quad \text{been}_5 \quad \text{implemented}_6$

$\text{Le}_1 \quad \text{programme}_2 \quad \text{a}_3 \quad \text{été}_4 \quad \text{mis}_5 \quad \text{en}_6 \quad \text{application}_7$

## Many-to-1 Alignments



## Many-to-Many Alignments



## A Word-Level TM?

- What might a model of P(f|e) look like?

$e = e_1 \dots e_I$    And₁    the₂    program₃    has₄    been₅    implemented₆

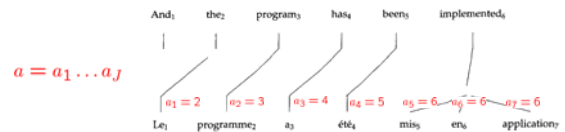$f = f_1 \dots f_J$    Le₁    programme₂    a₃    été₄    mis₅    en₆    application₇

$$P(f|e) = \prod_j P(f_j|e_1 \dots e_I)$$

*What can go wrong here?*

*How to estimate this?*

## IBM Model 1 (Brown 93)

- Alignments: a hidden vector called an *alignment* specifies which English source is responsible for each French target word.



$$P(f,a|e) = \prod_j P(a_j = i)P(f_j|e_i)$$

$$= \prod_j \frac{1}{I+1}P(f_j|e_i)$$

$$P(f|e) = \sum_a P(f,a|e)$$

## IBM Model 1

- Obvious first stab: greedy matchings
- Better approach: re-estimated generative models
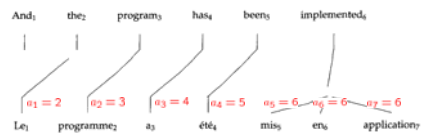
$$P(f|e) = \sum_a P(f,a|e)$$

$$P(f,a|e) = \prod_j P(a_j = i|e)P(f_j|e_i)$$

$$P(a_j = i|e,f) = \frac{P(f_j|e_i)}{\sum_{i'} P(f_j|e_{i'})}$$

- Basic idea: pick a source for each word, update co-occurrence statistics, repeat

## IBM Model 1 [Brown et al, 93]

- Alignments: a hidden vector called an *alignment* specifies which English source is responsible for each French target word. $a = a_1 \dots a_J$
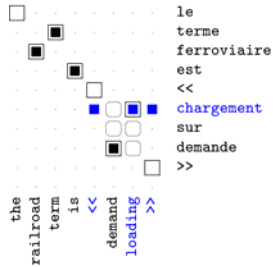


$$P(f,a|e) = \prod_j P(a_j = i)P(f_j|e_i)$$

$$= \prod_j \frac{1}{I+1}P(f_j|e_i)$$

## Problems with Model 1

- There's a reason they designed models 2-5!
- Problems: alignments jump around, align everything to rare words
- Experimental setup:
  - Training data: 1.1M sentences of French-English text, Canadian Hansards
  - Evaluation metric: alignment error Rate (AER)
  - Evaluation data: 447 hand-aligned sentences

```
          le
          terme
          ferroviaire
          est
          <<
          chargement
          sur
          demande
          >>
the railroad term is << demand loading >>
```

## Evaluating TMs

- How do we measure TM quality?
  - Method 1: use in an end-to-end translation system
    - Hard to measure translation quality
    - Option: human judges
    - Option: reference translations (NIST, BLEU scores)
  - Method 2: measure quality of the alignments produced
    - Easy to measure
    - Hard to know what the gold alignments should be
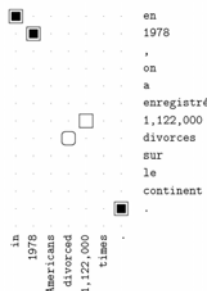    - May not correlate with translation quality (like perplexity in LMs)

## Alignment Error Rate

- Alignment Error Rate

☐ = Sure

◯ = Possible

■ = Predicted

```
          en
          1978
          ,
          on
          a
          enregistré
          1,122,000
          divorces
          sur
          le
          continent
          .
in 1978 Americans divorced 1,122,000 times
```

$$AER(A, S, P) = \left(1 - \frac{|A \cap S| + |A \cap P|}{|A| + |S|}\right)$$
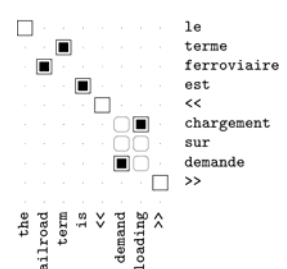
$$= \left(1 - \frac{3+3}{3+4}\right) = \frac{1}{7}$$

## Intersected Model 1

- Post-intersection: standard practice to train models in each direction then intersect their predictions [Och and Ney, 03]
- Second model is basically a filter on the first
  - Precision jumps, recall drops
  - End up not guessing hard alignments

| Model | P/R | AER |
|-------|-----|-----|
| Model 1 E→F | 82/58 | 30.6 |
| Model 1 F→E | 85/58 | 28.7 |
| Model 1 AND | 96/46 | 34.8 |

```
          le
          terme
          ferroviaire
          est
          <<
          chargement
          sur
          demande
          >>
the railroad term is << demand loading >>
```
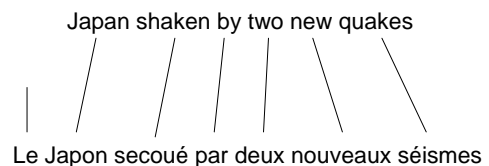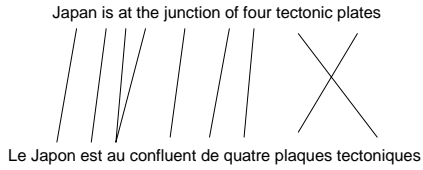
## Joint Training?

- Overall:
  - Similar high precision to post-intersection
  - But recall is much higher
  - More confident about positing non-null alignments

| Model | P/R | AER |
|-------|-----|-----|
| Model 1 E→F | 82/58 | 30.6 |
| Model 1 F→E | 85/58 | 28.7 |
| Model 1 AND | 96/46 | 34.8 |
| Model 1 INT | 93/69 | 19.5 |

## Monotonic Translation

Japan shaken by two new quakes

Le Japon secoué par deux nouveaux séismes

## Local Order Change

Japan is at the junction of four tectonic plates

Le Japon est au confluent de quatre plaques tectoniques

## IBM Model 2

- Alignments tend to the diagonal (broadly at least)

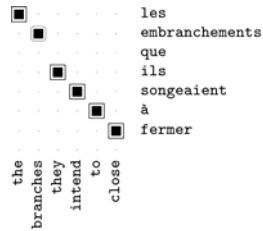$$P(f, a|e) = \prod_j P(a_j = i|j, I, J)P(f_j|e_i)$$

$$P(i - j\frac{I}{J})$$

$$\frac{1}{Z}e^{-\alpha(i-j\frac{I}{J})}$$

- Other schemes for biasing alignments towards the diagonal:
  - Relative alignment
  - Asymmetric distances
  - Learning a multinomial over distances

## Example

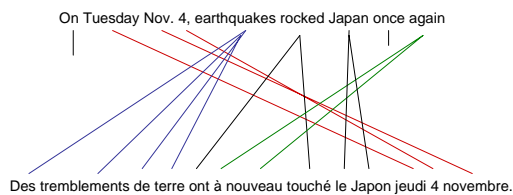| | |
|---|---|
| ■ | les |
| ■ | embranchements |
| | que |
| ■ | ils |
| ■ | songeaient |
| ■ | à |
| ■ | fermer |

the branches they intend to close

## EM for Models 1/2

- Model 1 Parameters:
  Translation probabilities (1+2)   $P(f_j|e_i)$
  Distortion parameters (2 only)    $P(a_j = i|j, I, J)$
- Start with $P(f_j|e_i)$ uniform, including $P(f_j|null)$
- For each sentence:
  - For each French position j
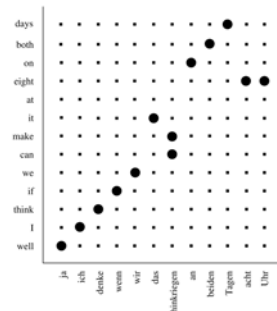    - Calculate posterior over English positions

$$P(a_j = i|f, e) = \frac{P(a_j = i|j, I, J)P(f_j|e_i)}{\sum_{i'} P(a_j = i'|j, I, J)P(f_j|e'_i)}$$

  - (or just use best single alignment)
  - Increment count of word $f_j$ with word $e_i$ by these amounts
  - Also re-estimate distortion probabilities for model 2
- Iterate until convergence

## Phrase Movement

On Tuesday Nov. 4, earthquakes rocked Japan once again

Des tremblements de terre ont à nouveau touché le Japon jeudi 4 novembre.

## Phrase Movement

days
both
on
eight
at
it
make
can
we
if
think
I
well
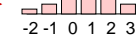
ja ich denke wenn wir das hinkriegen an beiden Tagen acht Uhr

## The HMM Model

- Model 2 preferred global monotonicity
- We want local monotonicity:
  - Most jumps are small
- HMM model (Vogel 96)

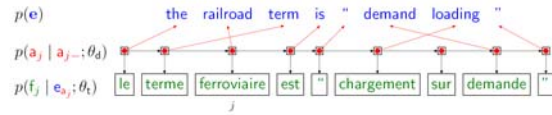| f | t(f \| e) |
|---|---|
| nationale | 0.469 |
| national | 0.418 |
| nationaux | 0.054 |
| nationales | 0.029 |

$$P(f, a|e) = \prod_j P(a_j|a_{j-1})P(f_j|e_i)$$

$$P(a_j - a_{j-1})$$

-2 -1 0 1 2 3

- Re-estimate using the forward-backward algorithm
- Handling nulls requires some care
- What are we still missing?

## The HMM Model



$p(e)$   the  railroad  term  is  " demand  loading  "

$p(a_j \mid a_{j-}; \theta_d)$

$p(f_j \mid e_{a_j}; \theta_t)$   le  terme  ferroviaire  est  "  chargement  sur  demande  "

Distortion $\theta_d$

$p(\; ) = 0.6$
$p(\; ) = 0.2$
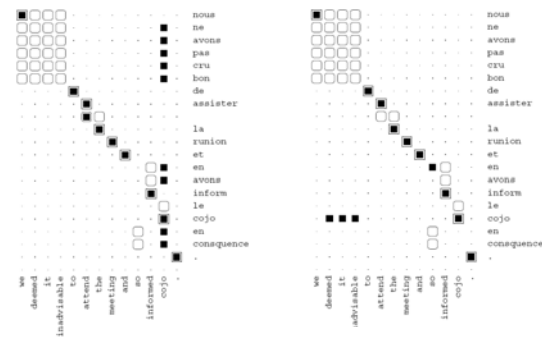$p(\; ) = 0.1$
...

Translation $\theta_t$

$p(\text{ the} \rightarrow \text{le} \;) = 0.53$
$p(\text{ the} \rightarrow \text{la} \;) = 0.24$
$p(\text{ railroad} \rightarrow \text{ferroviaire} \;) = 0.19$
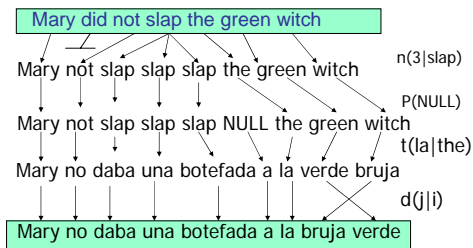$p(\text{ NULL} \rightarrow \text{le} \;) = 0.12$
...

## HMM Examples



## AER for HMMs

| Model | AER |
|---|---|
| Model 1 INT | 19.5 |
| HMM E→F | 11.4 |
| HMM F→E | 10.8 |
| HMM AND | 7.1 |
| HMM INT | 4.7 |
| GIZA M4 AND | 6.9 |

## IBM Models 3/4/5

Mary did not slap the green witch

Mary not slap slap slap the green witch      n(3|slap)

Mary not slap slap slap NULL the green witch      P(NULL)

Mary no daba una botefada a la verde bruja      t(la|the)

Mary no daba una botefada a la bruja verde      d(j|i)

[Al-Onaizan and Knight, 1998]

## Examples: Translation and Fertility

*the*

| f | t(f \| e) | φ | n(φ \| e) |
|---|---|---|---|
| le | 0.497 | 1 | 0.746 |
| la | 0.207 | 0 | 0.254 |
| les | 0.155 | | |
| l' | 0.086 | | |
| ce | 0.018 | | |
| cette | 0.011 | | |

*not*

| f | t(f \| e) | φ | n(φ \| e) |
|---|---|---|---|
| ne | 0.497 | 2 | 0.735 |
| pas | 0.442 | 0 | 0.154 |
| non | 0.029 | 1 | 0.107 |
| rien | 0.011 | | |

*farmers*

| f | t(f \| e) | φ | n(φ \| e) |
|---|---|---|---|
| agriculteurs | 0.442 | 2 | 0.731 |
| les | 0.418 | 1 | 0.228 |
| cultivateurs | 0.046 | 0 | 0.039 |
| producteurs | 0.021 | | |

## Example: Idioms

*nodding*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|---|---|---|---|
| signe | 0.164 | 4 | 0.342 |
| la | 0.123 | 3 | 0.293 |
| tête | 0.097 | 2 | 0.167 |
| oui | 0.086 | 1 | 0.163 |
| fait | 0.073 | 0 | 0.023 |
| que | 0.073 | | |
| hoche | 0.054 | | |
| hocher | 0.048 | | |
| faire | 0.030 | | |
| me | 0.024 | | |
| approuve | 0.019 | | |
| qui | 0.019 | | |
| un | 0.012 | | |
| faites | 0.011 | | |

## Example: Morphology

*should*

| f | $t(f \mid e)$ | $\phi$ | $n(\phi \mid e)$ |
|---|---|---|---|
| devrait | 0.330 | 1 | 0.649 |
| devraient | 0.123 | 0 | 0.336 |
| devrions | 0.109 | 2 | 0.014 |
| faudrait | 0.073 | | |
| faut | 0.058 | | |
| doit | 0.058 | | |
| aurait | 0.041 | | |
| doivent | 0.024 | | |
| devons | 0.017 | | |
| devrais | 0.013 | | |

## Some Results

- [Och and Ney 03]

| Model | Training scheme | 0.5K | 8K | 128K | 1.47M |
|---|---|---|---|---|---|
| Dice | | 50.9 | 43.4 | 39.6 | 38.9 |
| Dice+C | | 46.3 | 37.6 | 35.0 | 34.0 |
| Model 1 | $1^5$ | 40.6 | 33.6 | 28.6 | 25.9 |
| Model 2 | $1^5 2^5$ | 46.7 | 29.3 | 22.0 | 19.5 |
| HMM | $1^5 H^5$ | 26.3 | 23.3 | 15.0 | 10.8 |
| Model 3 | $1^5 2^5 3^3$ | 43.6 | 27.5 | 20.5 | 18.0 |
| | $1^5 H^5 3^3$ | 27.5 | 22.5 | 16.6 | 13.2 |
| Model 4 | $1^5 2^5 3^3 4^3$ | 41.7 | 25.1 | 17.3 | 14.1 |
| | $1^5 H^5 3^3 4^3$ | 26.1 | 20.2 | 13.1 | 9.4 |
| | $1^5 H^5 4^3$ | 26.3 | 21.8 | 13.3 | 9.3 |
| Model 5 | $1^5 H^5 4^3 5^3$ | 26.5 | 21.5 | 13.7 | 9.6 |
| | $1^5 H^5 3^3 4^3 5^3$ | 26.5 | 20.4 | 13.4 | 9.4 |
| Model 6 | $1^5 H^5 4^3 6^3$ | 26.0 | 21.6 | 12.8 | 8.8 |
| | $1^5 H^5 3^3 4^3 6^3$ | 25.9 | 20.3 | 12.5 | 8.7 |