

Statistical NLP

Spring 2011

Berkeley

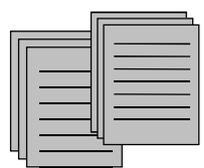
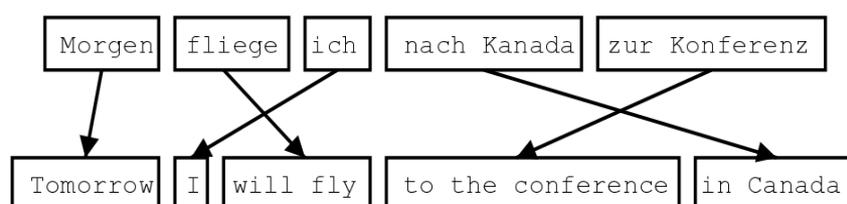


N L P

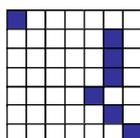
Lecture 8: Word Alignment

Dan Klein – UC Berkeley

Phrase-Based Systems



Sentence-aligned corpus



Word alignments



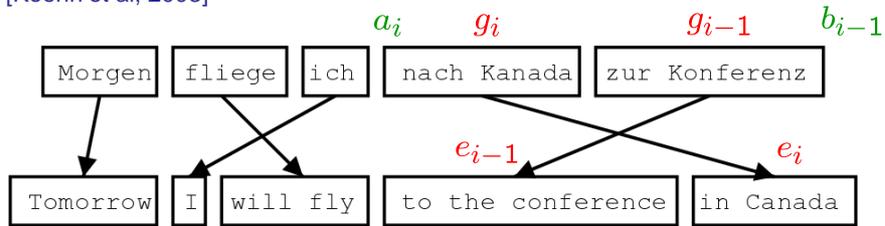
cat		chat		0.9	
the	cat		le chat		0.8
dog		chien		0.8	
house		maison		0.6	
my house		ma maison		0.9	
language		langue		0.9	
...					

Phrase table
(translation model)

Many slides and examples from Philipp Koehn or John DeNero

The Pharaoh "Model"

[Koehn et al, 2003]



$$P(e|g) = P(\{\bar{g}_i\}|g) \prod_i \phi(\bar{e}_i|\bar{g}_i) d(a_i - b_{i-1})$$

↓ Segmentation
 ↓ Translation
 ↓ Distortion

Phrase-Based Decoding

Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a slap		by		green witch	
	no		slap		to the			
	did not give				to			
					the			
				slap		the witch		

Phrase Translation

Maria	no	dio	una	bofetada	a	la	bruja	verde
-------	----	-----	-----	----------	---	----	-------	-------

<u>Mary</u>	<u>not</u>	<u>give</u>	<u>a</u>	<u>slap</u>	<u>to</u>	<u>the</u>	<u>witch</u>	<u>green</u>
	<u>did not</u>		<u>a slap</u>		<u>by</u>		<u>green witch</u>	
	<u>no</u>		<u>slap</u>		<u>to the</u>			
	<u>did not give</u>				<u>to</u>			
					<u>the</u>			
				<u>slap</u>		<u>the witch</u>		

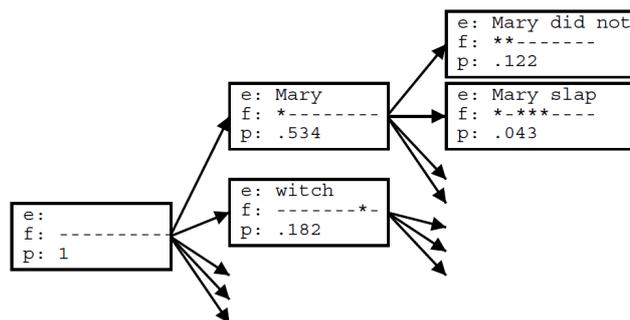
- If monotonic, almost an HMM; technically a semi-HMM

```

for (fPosition in 1...|f|)
  for (lastPosition < fPosition)
    for (eContext in eContexts)
      for (eOption in translations[fPosition])
        ... combine hypothesis for (lastPosition ending in eContext) with eOption
  
```

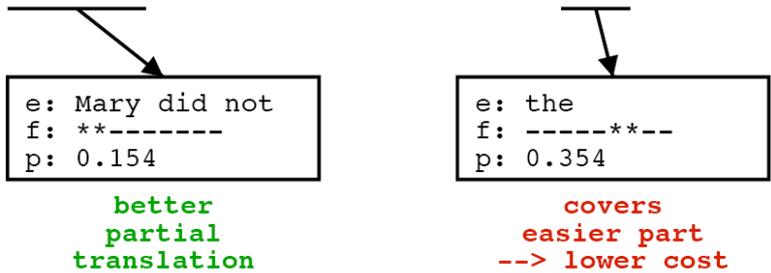
- If distortion... now what?

Non-Monotonic Phrasal MT



Pruning: Beams + Forward Costs

Maria no dio una bofetada a la bruja verde



- Problem: easy partial analyses are cheaper
 - Solution 1: use beams per foreign subset
 - Solution 2: estimate forward costs (A*-like)

The Pharaoh Decoder

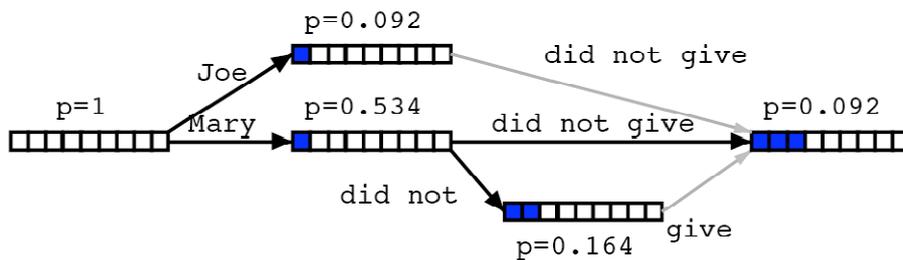
Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a slap		by		green witch	
	no		slap		to the			
	did not give				to			
				slap		the		
						the witch		

Maria	no	dio una bofetada	a la	bruja	verde
-------	----	------------------	------	-------	-------

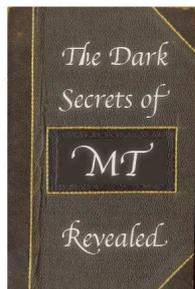
Mary	did not	slap	the	green	witch
------	---------	------	-----	-------	-------

Hypothesis Lattices

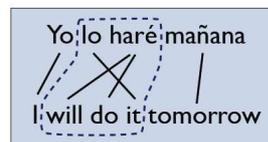
Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a slap		by		green	witch
	no		slap		to the			
	did not give				to			
					the			
			slap			the	witch	



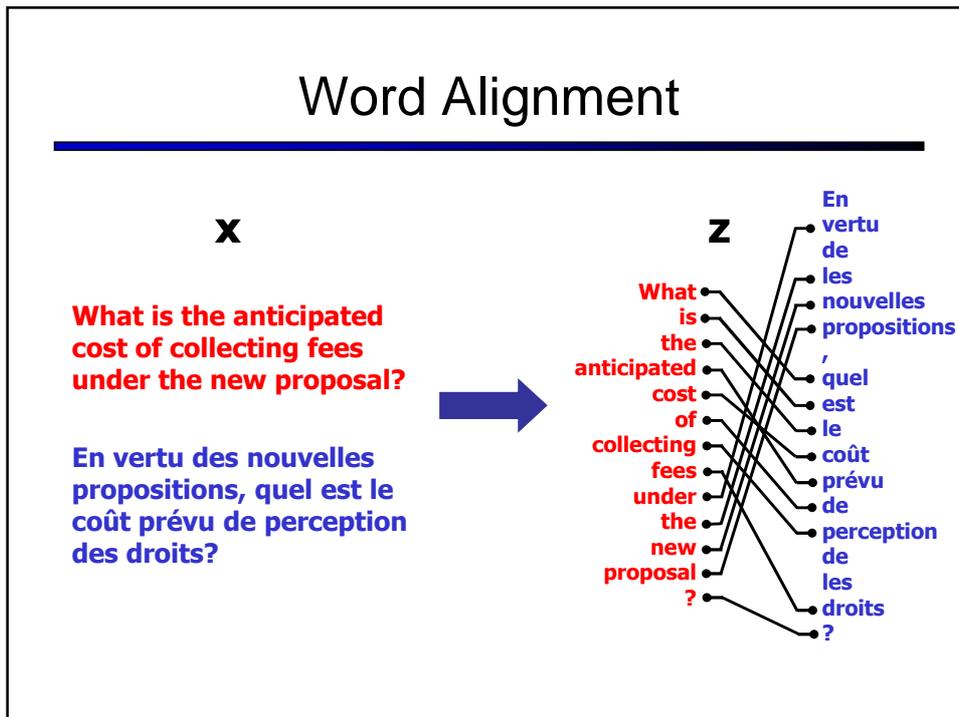
Word Alignment



- 1) Align words with a probabilistic model
- 2) Infer presence of larger structures from this alignment
- 3) Translate with the larger structures



Word Alignment

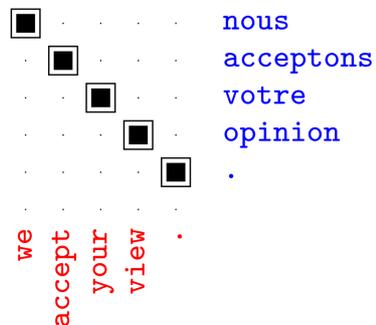


Unsupervised Word Alignment

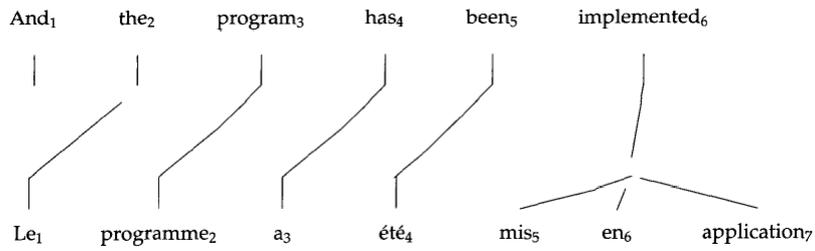
- Input: a *bitext*: pairs of translated sentences

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

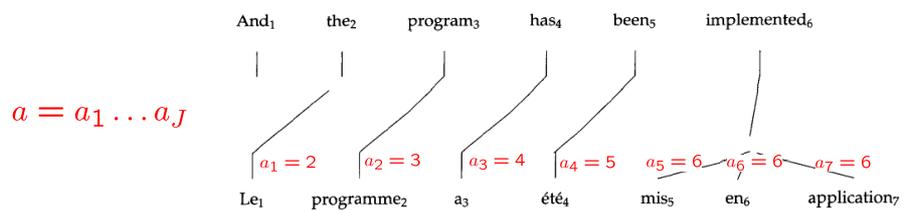


1-to-Many Alignments



IBM Model 1 (Brown 93)

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

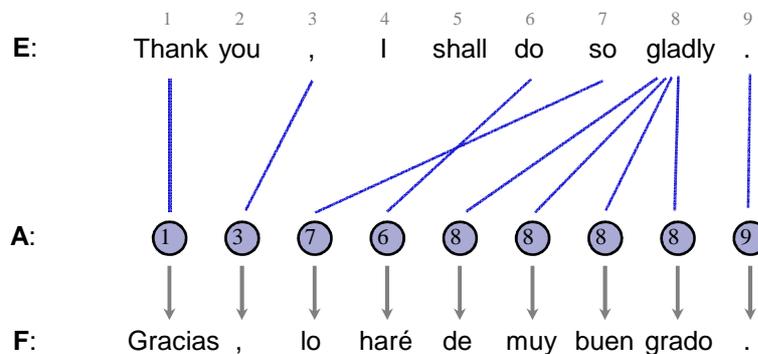


$$P(f, a|e) = \prod_j P(a_j = i) P(f_j|e_i)$$

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

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

IBM Models 1/2



Model Parameters

Emissions: $P(F_1 = \text{Gracias} \mid E_{A_1} = \text{Thank})$ Transitions: $P(A_2 = 3)$

Evaluating TMs

- How do we measure quality of a word-to-word model?
 - Method 1: use in an end-to-end translation system
 - Hard to measure translation quality
 - Option: human judges
 - Option: reference translations (NIST, BLEU)
 - Option: combinations (HTER)
 - Actually, no one uses word-to-word models alone as TMs
 - Method 2: measure quality of the alignments produced
 - Easy to measure
 - Hard to know what the gold alignments should be
 - Often does not correlate well with translation quality (like perplexity in LMs)

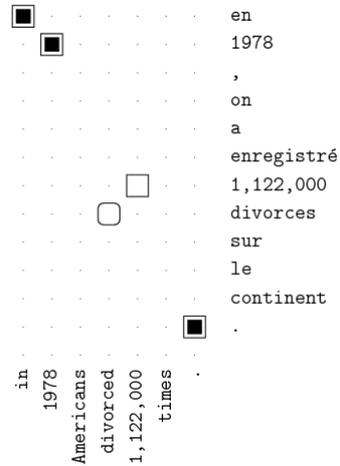
Alignment Error Rate

Alignment Error Rate

- = Sure
- = Possible
- = Predicted

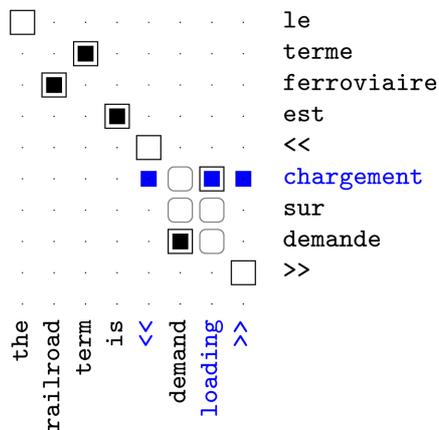
$$AER(A, S, P) = \left(1 - \frac{|A \cap S| + |A \cap P|}{|A| + |S|}\right)$$

$$= \left(1 - \frac{3 + 3}{3 + 4}\right) = \frac{1}{7}$$



Problems with Model 1

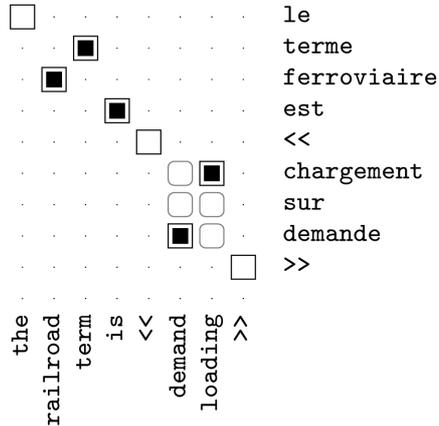
- There's a reason they designed models 2-5!
- Problems: alignments jump around, align everything to rare words
- Experimental setup:
 - Training data: 1.1M sentences of French-English text, Canadian Hansards
 - Evaluation metric: alignment error Rate (AER)
 - Evaluation data: 447 hand-aligned sentences



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



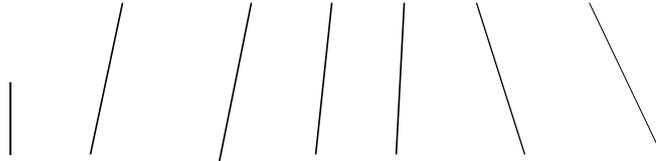
Joint Training?

- 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

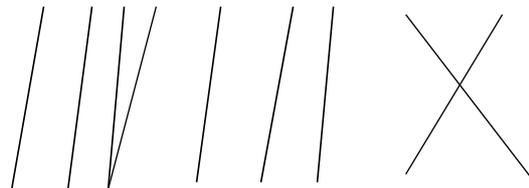
Japan shaken by two new quakes



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)

$$P(f, a|e) = \prod_j P(a_j = i|j, I, J) P(f_j|e_i) \\ P(\text{dist} = i - j \frac{I}{J}) \\ \frac{1}{Z} e^{-\alpha(i - j \frac{I}{J})}$$

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

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

Example

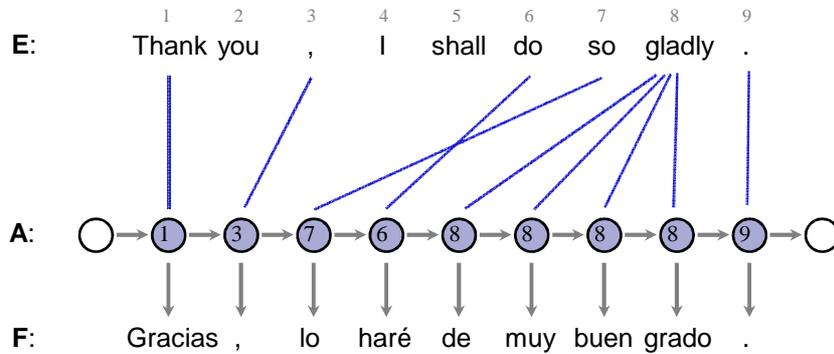
the	branches	they	intend	to	close	les	embranchements	que	ils	songeaient	à	fermer
						■						
							■					
								■				
									■			
										■		
											■	

Phrase Movement

On Tuesday Nov. 4, earthquakes rocked Japan once again

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

The HMM Model



Model Parameters

Emissions: $P(F_1 = \text{Gracias} \mid E_{A_1} = \text{Thank})$ Transitions: $P(A_2 = 3 \mid A_1 = 1)$

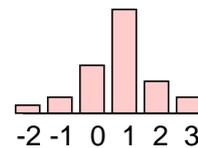
The HMM Model

- 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

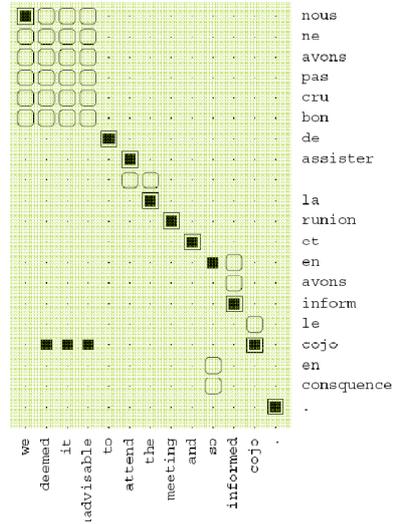
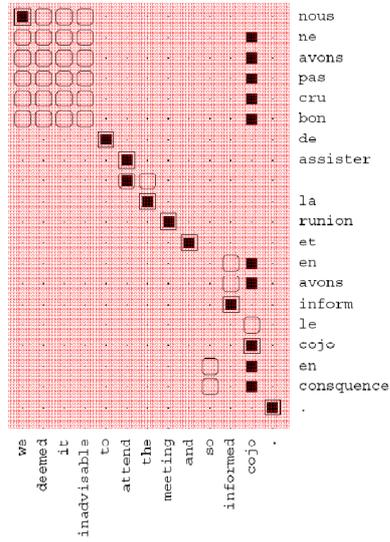
$$P(f, a|e) = \prod_j P(a_j|a_{j-1})P(f_j|e_i)$$

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



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

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