

# Machine Translation: IBM alignment models, practicalities, evaluation

Bill MacCartney

CS224N / Ling 284

[Based on slides by Chris Manning,  
Kevin Knight, Dan Klein, Dan Jurafsky]

# Quiz question

When working out expectations in mixture models using the EM algorithm, as in the last lecture:

Let  $y_{i,j}$  be the binary variable indicating that mixture component  $j$  generated data item  $i$ .

What is the value of  $\sum_i \sum_j E(y_{i,j} | x_i, \Theta)$  ?

- (a) 0
- (b) 1
- (c) -1
- (d)  $k$
- (e)  $n$

[Note:  $x$ ,  $\Theta$ ,  $k$ , and  $n$  follow the same notation as in lecture.]

# Last week's quiz question

How many possible (IBM-style) alignments?

	<i>i</i>		<i>j</i>		<i>a</i>
Mary	1		1	Maria	-
did	2		2	no	-
not	3		3	daba	-
slap	4		4	una	-
			5	botefada	-
			6	a	-
the	5		7	la	-
green	6		8	bruja	-
witch	7		9	verde	-

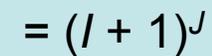
A. 362880

B. 4782969

C. 40353607

D. 43046721

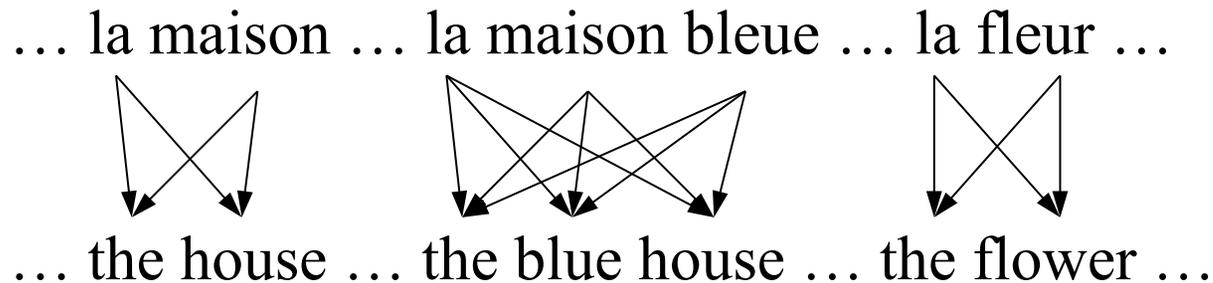
E. 134217728


$$= (I + 1)^J$$

# Unsupervised Word Alignment

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

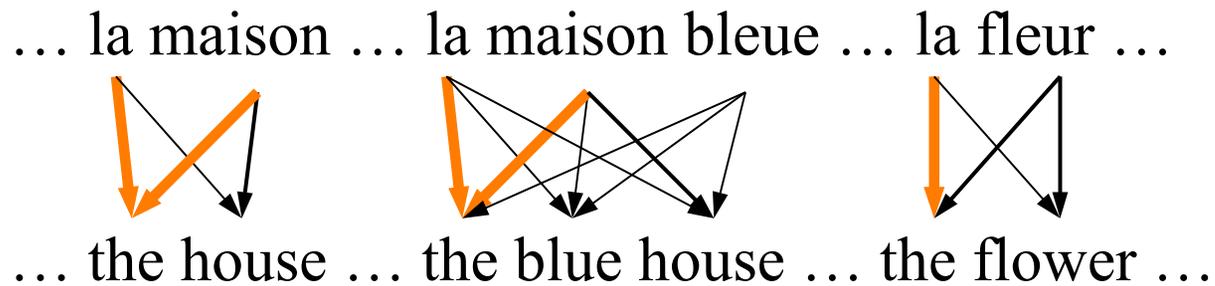
Output: *alignments*: pairs of translated *words*



All word alignments equally likely

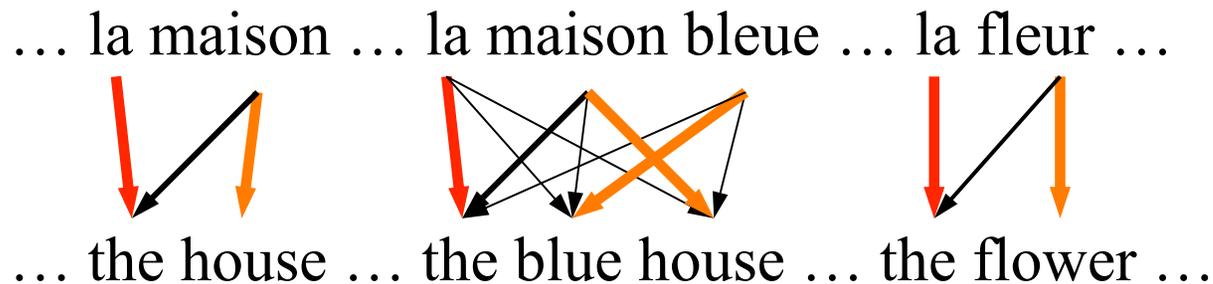
All  $P(\text{french-word} \mid \text{english-word})$  equally likely

# Unsupervised Word Alignment



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

# Unsupervised Word Alignment

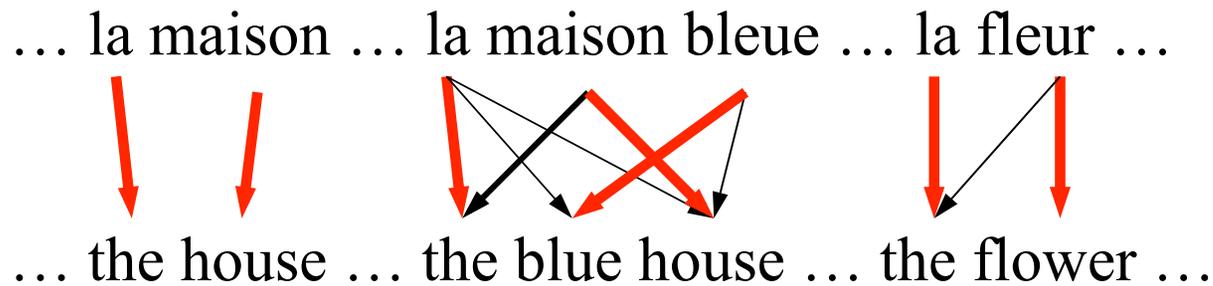


“maison” co-occurs with both “the” and “house”, but  $P(\text{maison} | \text{house})$  can be raised without limit, to 1.0, while  $P(\text{maison} | \text{the})$  is limited because of “la”

$$\sum_f P(f | \text{the}) = 1$$

(pigeonhole principle)

# Unsupervised Word Alignment



settling down after another iteration

That was the idea of IBM Model 1. For details, see following slides and:

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

# IBM Model 1 generative story

	Le	programme	a	été	mis	en	application
And							
the	■						
program		■					
has			■				
been				■			
implemented					■	■	■
$a_j$	2	3	4	5	6	6	6

Given English sentence  $e_1, e_2, \dots, e_l$

Choose length  $J$  for French sentence

For each  $j$  in 1 to  $J$ :

- Choose  $a_j$  uniformly from  $0, 1, \dots, l$
- Choose  $f_j$  by translating  $e_{a_j}$

# IBM Model 1 parameters

$$P(f, a|e) = P(J|I) \prod_j P(a_j) P(f_j|e_{a_j})$$

$$= \epsilon \prod_j P(a_j) P(f_j|e_{a_j})$$

$$= \epsilon \prod_j \frac{1}{I+1} P(f_j|e_{a_j})$$

$$= \frac{\epsilon}{(I+1)^J} \prod_j P(f_j|e_{a_j})$$

	Le	programme	a	été	mis	en	application
And							
the							
program							
has							
been							
implemented							
$a_j$	2	3	4	5	6	6	6

# Applying Model 1\*

$P(f, a | e)$  can be used as a *translation model* or an *alignment model*

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

As alignment model 
$$P(a|e, f) = \frac{P(f, a|e)}{P(f|e)}$$
$$= \frac{P(f, a|e)}{\sum_{a'} P(f, a'|e)}$$

\* Actually, any  $P(f, a | e)$ , e.g., any IBM model

# Applying Model 1 *efficiently*

(see Knight 99, section 31)

And the program has been implemented

	Le programme a été mis en application						
	0.01	0.02	0.01	0.02	0.01	0.03	0.01
	0.49	0.01	0.02	0.03	0.01	0.02	0.01
	0.01	0.53	0.01	0.02	0.01	0.03	0.01
	0.02	0.01	0.48	0.11	0.06	0.04	0.01
	0.02	0.01	0.17	0.39	0.01	0.02	0.01
	0.01	0.02	0.01	0.01	0.4	0.12	0.57

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

$$\propto \sum_a \prod_j P(f_j|e_{a_j})$$

exponential?

$$= \sum_{a_1=0}^I \dots \sum_{a_J=0}^I \prod_j P(f_j|e_{a_j})$$

$$\propto \prod_j \sum_i P(f_j|e_i)$$

quadratic!

# Model 1: Word alignment learning with Expectation-Maximization (EM)

- Start with  $P(f^p|e^q)$  uniform, including  $P(f^p|\text{NULL})$
- For each sentence
  - For each French position  $j$ 
    - Calculate posterior over English positions  $P(a_j | e, f)$

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

- Increment count of word  $f_j$  with word  $e_{a_j}$ 
      - $C(f_j|e_i) += P(a_j = i | f, e)$
- Renormalize counts to give probs  $P(f^p | e^q) = \frac{C(f^p | e^q)}{\sum_{f^x} C(f^x | e^q)}$
- Iterate until convergence

# IBM StatMT Translation Models

- IBM1 – lexical probabilities only
  - IBM2 – lexicon plus absolute position
  - HMM – lexicon plus relative position
  - IBM3 – plus fertilities
  - IBM4 – inverted relative position alignment
  - IBM5 – non-deficient version of model 4
- 
- All the models we discuss today handle 0:1, 1:0, 1:1, 1:n alignments *only*

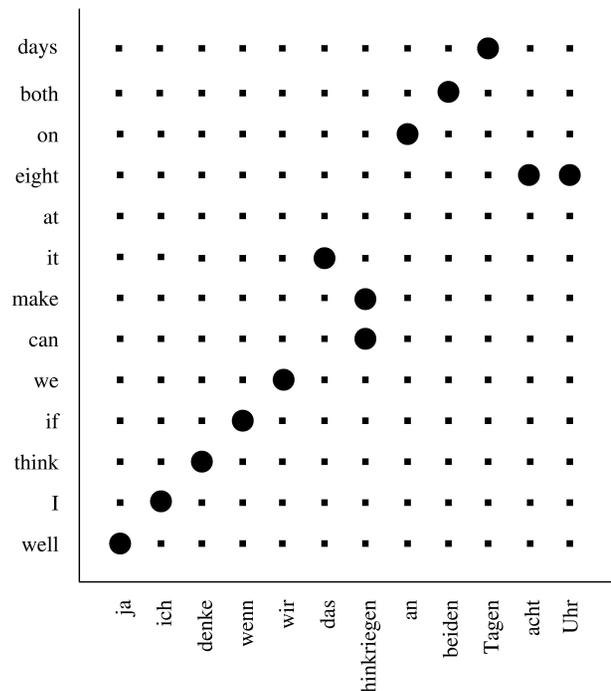
[Brown et al. 93, Vogel et al. 96]

# Comparative results

[Och & Ney 03]		Size of training corpus			
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

# IBM models 1,2,3,4,5

- In Model 2, the placement of a word in the French depends on where it was in the English

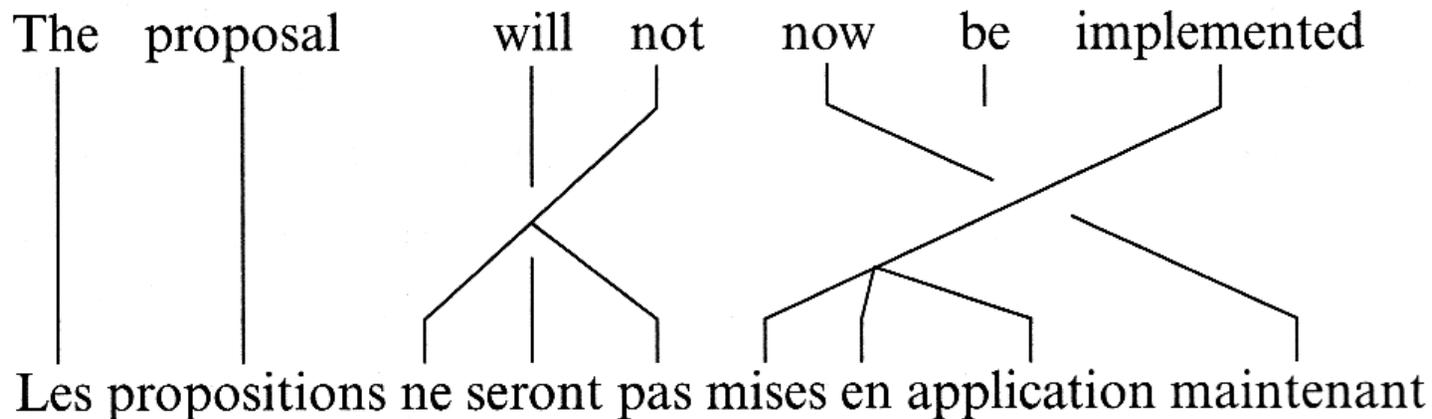


- Unlike Model 1, Model 2 captures the intuition that translations should usually “lie along the diagonal”.

- The main focus of PA #2.

# IBM models 1,2,3,4,5

- In model 3 we model how many French words an English word can produce, using a concept called fertility



**Figure 32.3**

Alignment example.

# Examples: translation & fertility

*the*

f	$t(f   e)$	$\phi$	$n(\phi   e)$
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*

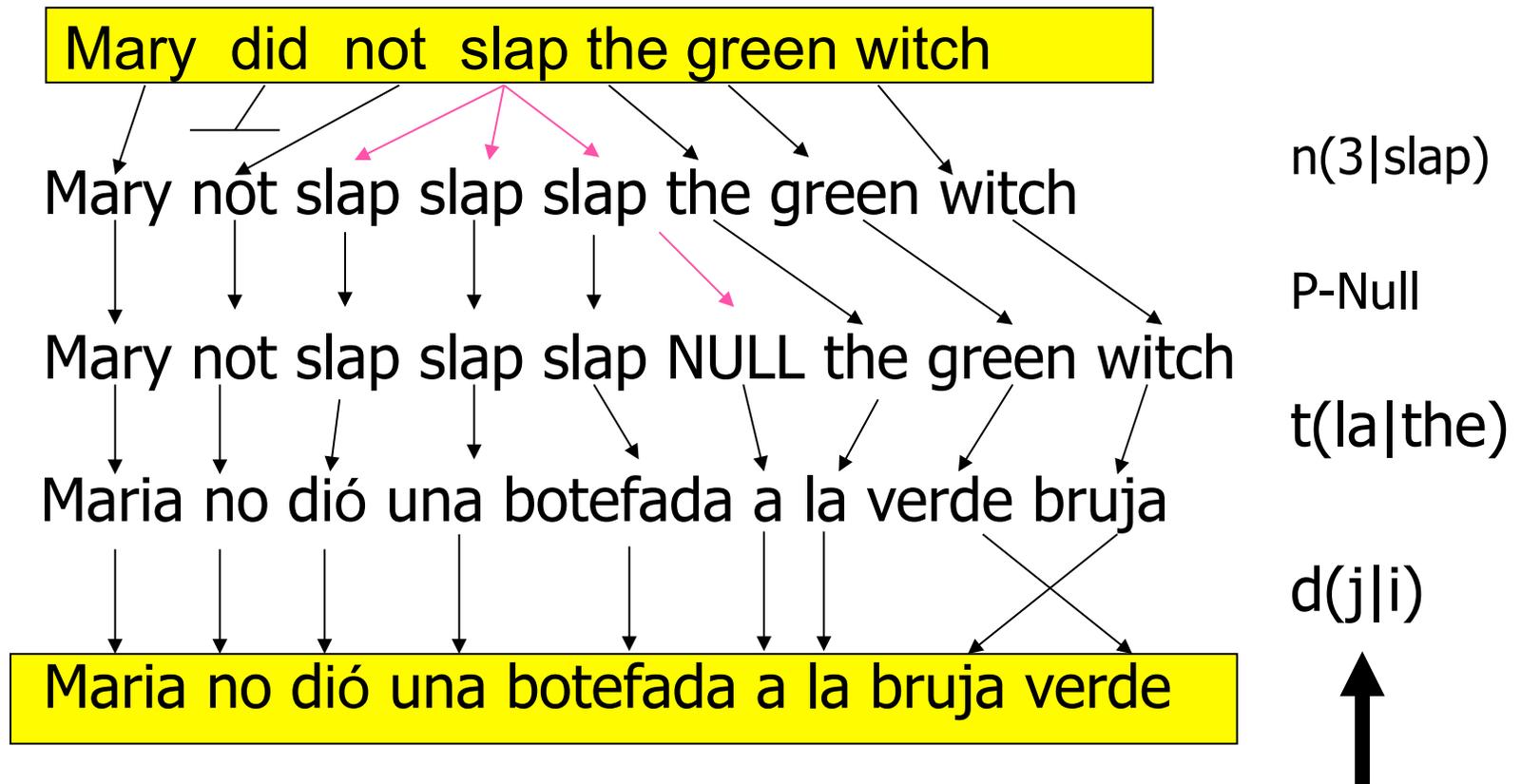
f	$t(f   e)$	$\phi$	$n(\phi   e)$
ne	0.497	2	0.735
pas	0.442	0	0.154
non	0.029	1	0.107
rien	0.011		

*farmers*

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

# IBM Model 3, Brown et al., 1993

Generative approach:



Probabilities can be learned from raw bilingual text.

# IBM Model 3 (from Knight 1999)

- For each word  $e_i$  in English sentence, choose a **fertility**  $\Phi_i$ . The choice of  $\Phi_i$  depends only on  $e_i$ , not other words or  $\Phi$ 's.
- For each word  $e_i$ , generate  $\Phi_i$  French words. Choice of French word depends only on English word  $e_i$ , not English context or any French words.
- Permute all the French words. Each French word gets assigned absolute target position slot (1,2,3, etc). Choice of French word position dependent only on absolute position of English word generating it.

# Model 3 parameters

- What are the parameters for this model?
- **Words**:  $P(\text{casa} \mid \text{house})$
- **Spurious words**:  $P(a \mid \text{NULL})$
- **Fertilities**:  $n(1 \mid \text{house})$ : prob that “house” will produce 1 foreign word whenever ‘house’ appears.
- **Distortions**:  $d(5 \mid 2)$  prob. that word in position 2 of English sentence generates word in position 5 of foreign translation
  - Actually, distortions are  $d(5 \mid 2, 4, 6)$  where 4 is length of English sentence, 6 is foreign length

# Spurious words

- We could have  $n(3 \mid \text{NULL})$  (probability of being exactly 3 spurious words in a foreign translation)
- But instead, of  $n(0 \mid \text{NULL})$ ,  $n(1 \mid \text{NULL})$ , ...  $n(25 \mid \text{NULL})$ , have a single parameter  $p_1$
- After assign fertilities to non-NULL English words we want to generate (say)  $z$  foreign words.
- As we generate each of  $z$  words, we optionally toss in spurious foreign word with probability  $p_1$
- Probability of not adding spurious word  $p_0 = 1 - p_1$

# Distortion probabilities for spurious words

- Can't just have  $d(5 | 0, 4, 6)$ , i.e. chance that NULL word will end up in position 5.
- Why? These are spurious words! Could occur anywhere!! Too hard to predict
- Instead,
  - Use normal-word distortion parameters to choose positions for normally-generated foreign words
  - Put NULL-generated words into empty slots left over
  - If three NULL-generated words, and three empty slots, then there are  $3!$ , or six, ways for slotting them all in
  - We'll assign a probability of  $1/6$  for each way

# Real Model 3

- For each word  $e_i$  in English sentence, choose fertility  $\Phi_i$  with prob  $n(\Phi_i | e_i)$
- Choose number  $\Phi_0$  of spurious foreign words to be generated from  $e_0 = \text{NULL}$  using  $p_1$  and sum of fertilities from step 1
- Let  $m$  be sum of fertilities for all words including NULL
- For each  $i=0,1,2,\dots,L$ ,  $k=1,2,\dots,\Phi_i$ :
  - choose foreign word  $\tau_{ik}$  with probability  $t(\tau_{ik} | e_i)$
- For each  $i=1,2,\dots,L$ ,  $k=1,2,\dots,\Phi_i$ :
  - choose target foreign position  $\pi_{ik}$  with prob  $d(\pi_{ik} | i, L, m)$
- For each  $k=1,2,\dots,\Phi_0$  choose position  $\pi_{0k}$  from  $\Phi_0 - k + 1$  remaining vacant positions in  $1,2,\dots,m$  for total prob of  $1 / \Phi_0!$
- Output foreign sentence with words  $\tau_{ik}$  in positions  $\pi_{ik}$  ( $0 \leq i \leq L, 1 \leq k \leq \Phi_i$ )

# Learning Model 3 parameters

- $n, t, p_1, d$
- Again, if we had complete data of English strings and step-by-step rewritings into Spanish, we could:
  - Compute  $n(0 | \text{did})$  by locating every instance of “did”, and seeing how many words it translates to
  - $t(\text{maison} | \text{house})$  how many of all French words generated by “house” were “maison”
  - $d(5 | 2, 4, 6)$  out of all times some word2 was translated, how many times did it become word5?

# Since we don't have word-aligned data...

- We bootstrap alignments from incomplete data
- From a sentence-aligned bilingual corpus
  - 1) Assume some startup values for  $n$ ,  $d$ ,  $\Phi$ , etc.
  - 2) Use values for  $n$ ,  $d$ ,  $\Phi$ , etc. to use Model 3 to work out chances of different possible alignments. Use these alignments to retrain  $n$ ,  $d$ ,  $\Phi$ , etc.
  - 3) Go to 2
- This is a more complicated case of the EM algorithm

# IBM models 1,2,3,4,5

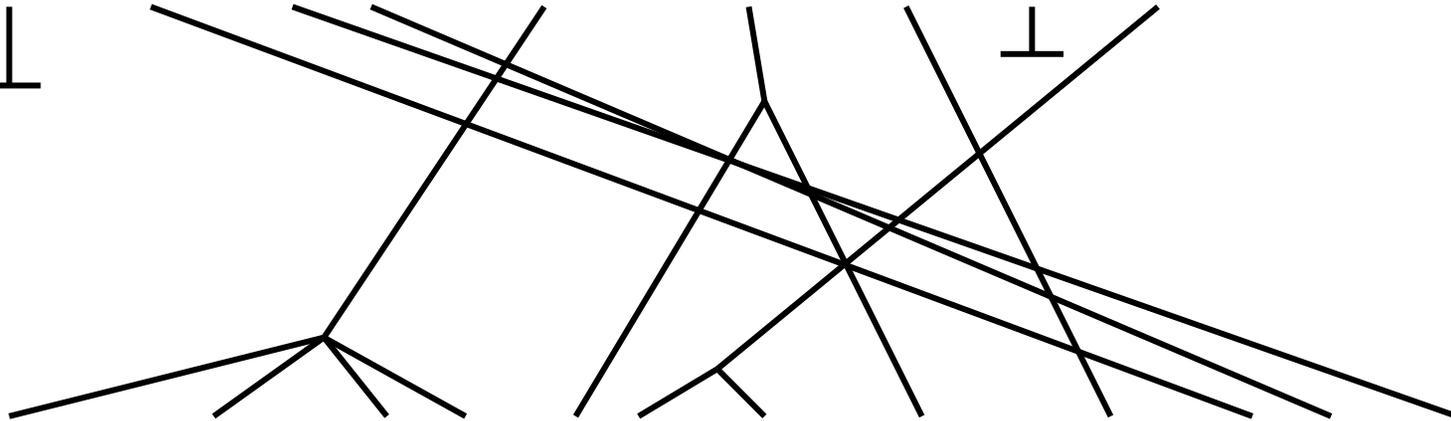
- In model 4 the placement of later French words produced by an English word depends on what happened to earlier French words generated by that same English word

# Alignments: linguistics

On Tuesday Nov. 4, earthquakes rocked Japan once again



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



# IBM models 1,2,3,4,5

- In model 5 they do non-deficient alignment. That is, you can't put probability mass on impossible things.

# Sample Translation Probabilities

## Translation Model

e	f	P(f   e)
national	nationale	0.47
	national	0.42
	nationaux	0.05
	nationales	0.03
the	le	0.50
	la	0.21
	les	0.16
	l'	0.09
	ce	0.02
	cette	0.01
farmers	agriculteurs	0.44
	les	0.42
	cultivateurs	0.05
	producteurs	0.02

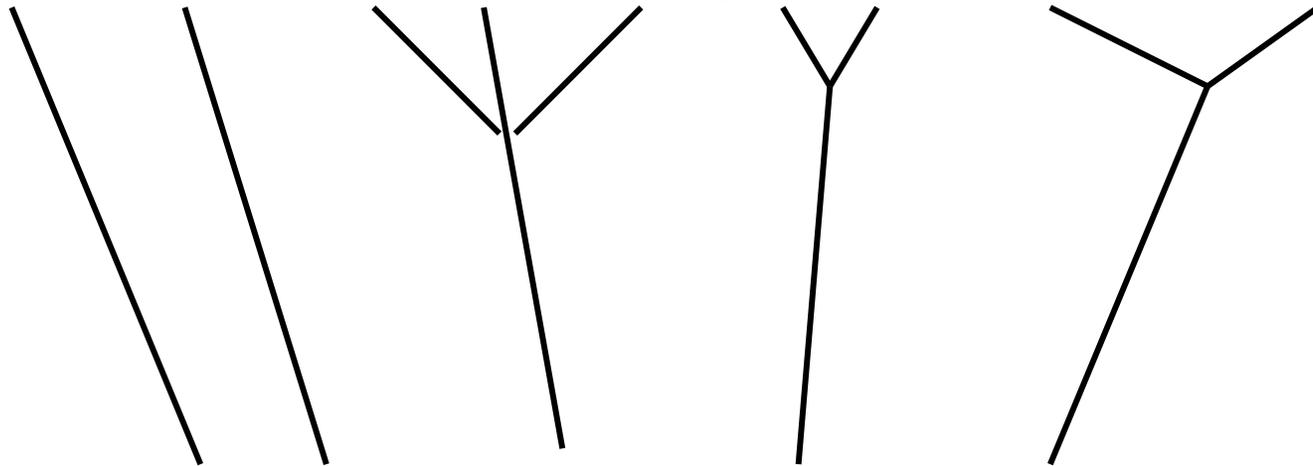
[Brown et al 93]

# Why all the models?

- We don't start with aligned text, so we have to get initial alignments from somewhere.
- Model 1 is words only, and is relatively easy and fast to train.
- We are working in a space with many local maxima, so output of model 1 can be a good place to start model 2. Etc.
- The sequence of models allows a better model to be found faster [the intuition is like deterministic annealing].

# Alignments: impossible (in IBM)

The balance was the territory of the aboriginal people

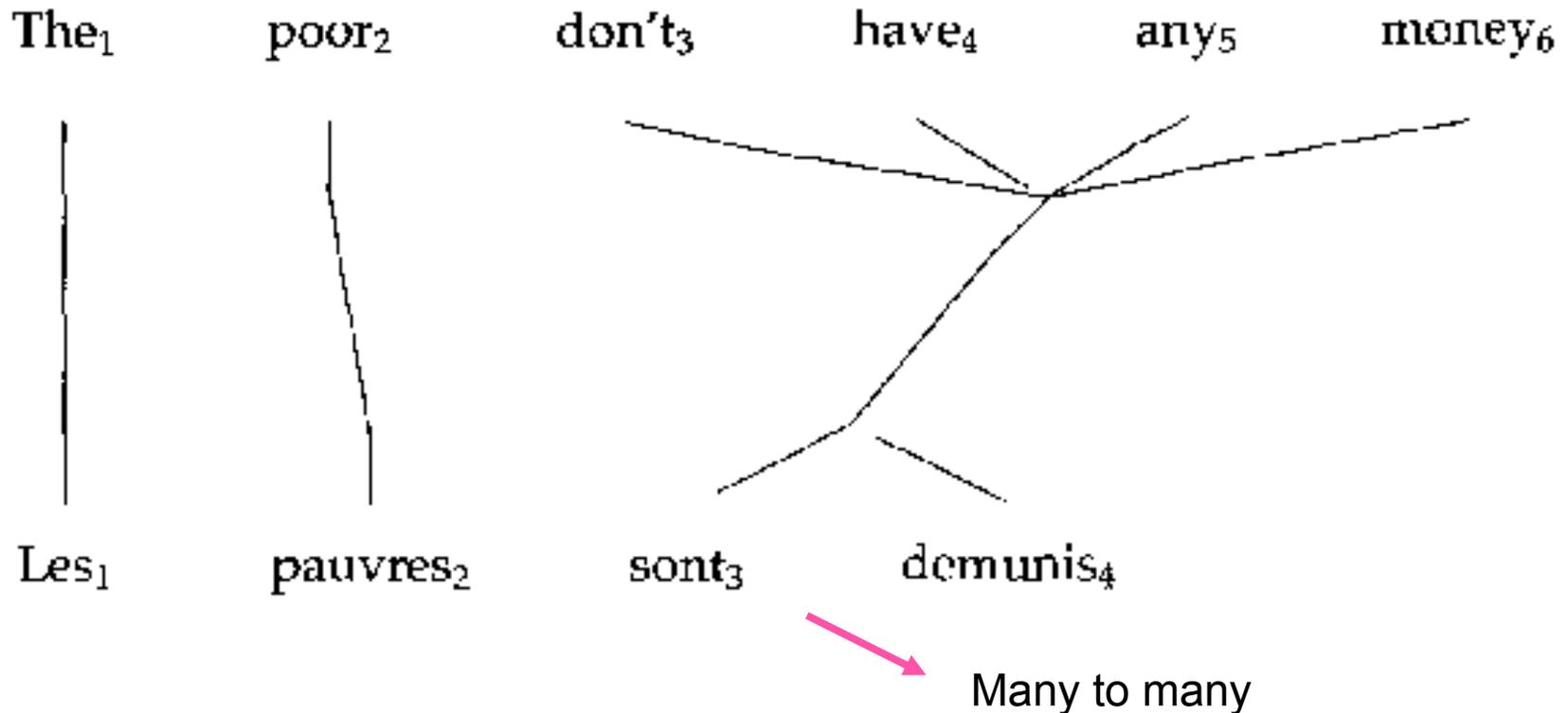


Le reste appartenait aux autochtones



Several words translated as one

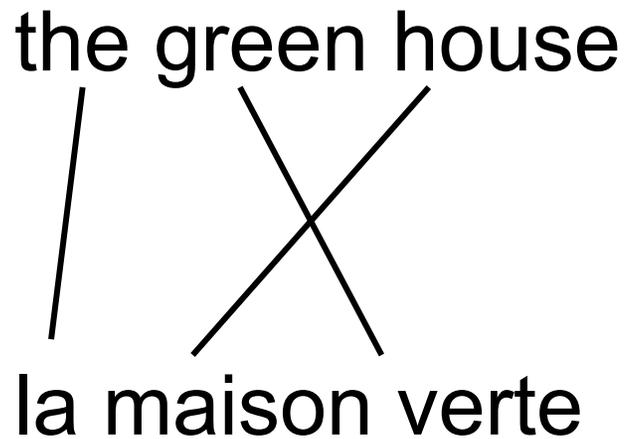
# Alignments: impossible (in IBM)



- A minimal aligned subset of words is called a 'cept' in the IBM work; often a 'bead' or '(aligned) statistical phrase' elsewhere.

# Alignments: linguistics

the green house  
la maison verte



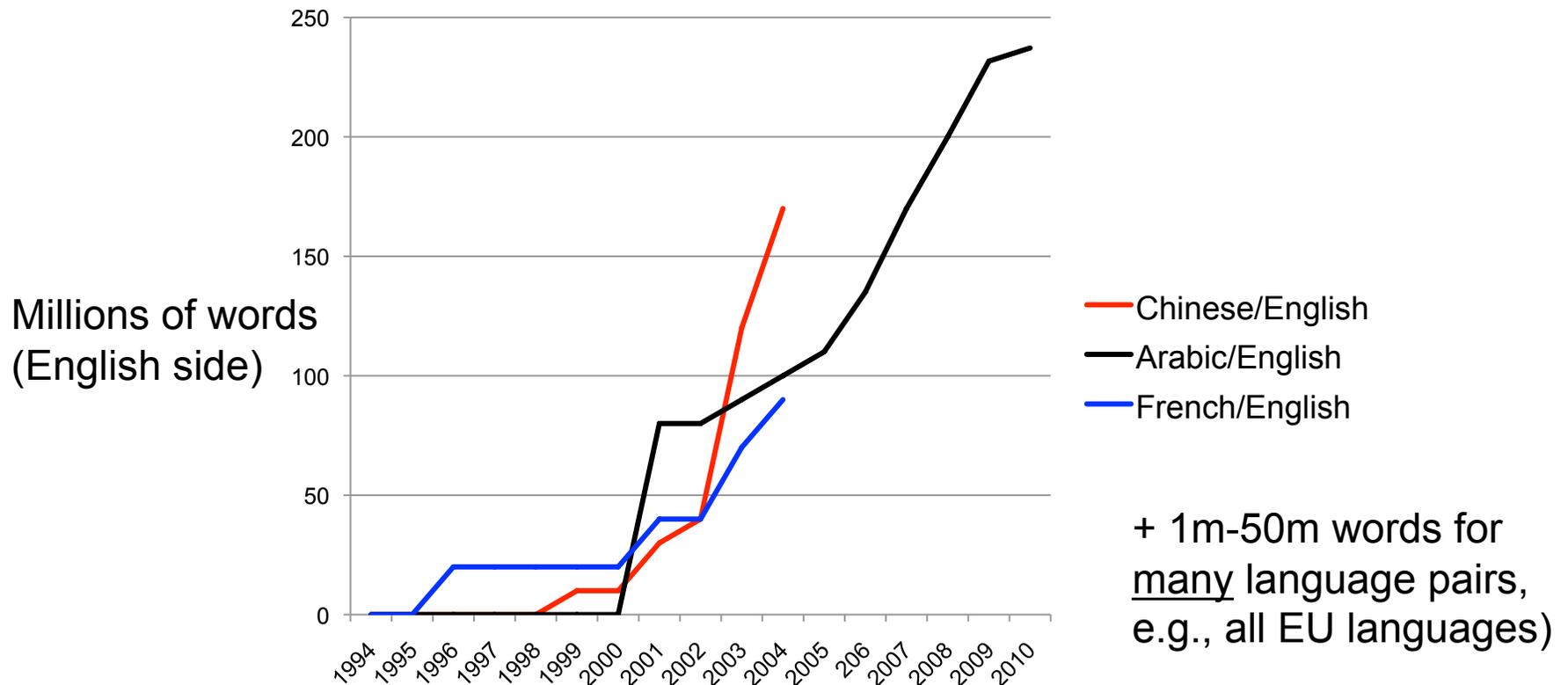
- There isn't enough linguistics to explain this in the translation model ... have to depend on the language model ... that may be unrealistic ... and may be harming our translation model

**Some practicalities**

# Getting Sentence Pair Data

- Really hard way: pay \$\$\$
  - Suppose one billion words of parallel data were sufficient
  - At 5 cents/word, that's \$50 million
- Pretty hard way: Find it, and then earn it!
  - De-formatting
  - Remove strange characters
  - Character code conversion
  - Document alignment
  - **Sentence alignment**
  - **Tokenization (also called Segmentation)**
- Easy way: Linguistic Data Consortium (LDC)

# Ready-to-Use Online Bilingual Data



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

# Tokenization (or Segmentation)

- English

- Input (some character stream):

"There," said Bob.

- Output (7 “tokens” or “words”):

" There , " said Bob .

- Chinese

- Input (char 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 。

# 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.

Done by similar Dynamic Programming or EM: see FSNLP ch. 13 for details

# MT Evaluation

(left over to 2011/01/24)

# Illustrative translation 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

# MT Evaluation

- Manual (the best!?):
  - SSER (subjective sentence error rate)
  - Correct/Incorrect
  - **Adequacy and Fluency** (5 or 7 point scales)
  - Error categorization
  - **Comparative ranking of translations**
- Testing in an application that uses MT as one sub-component
  - Question answering from foreign language documents
- Automatic metric:
  - WER (word error rate) – why problematic?
  - **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 match same portion of reference translation twice at a certain n-gram level (two MT words *airport* are only correct if two reference words *airport*; can't cheat by typing out "the the the the the")
  - Do count unigrams also in a bigram for unigram precision, etc.
- Brevity Penalty
  - Can't just type out single word "the" (precision 1.0!)
- It was thought quite hard to "game" the system (i.e., to 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.

- BLEU is a weighted geometric mean, with a brevity penalty factor added.
  - Note that it's precision-oriented
- BLEU4 formula  
(counts n-grams up to length 4)

$$\exp (1.0 * \log p1 + 0.5 * \log p2 + 0.25 * \log p3 + 0.125 * \log p4 - \max(\text{words-in-reference} / \text{words-in-machine} - 1, 0))$$

p1 = 1-gram precision

P2 = 2-gram precision

P3 = 3-gram precision

P4 = 4-gram precision

Note: only works at corpus level (zeroes kill it); there's a smoothed variant for sentence-level

# 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!)

# 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.

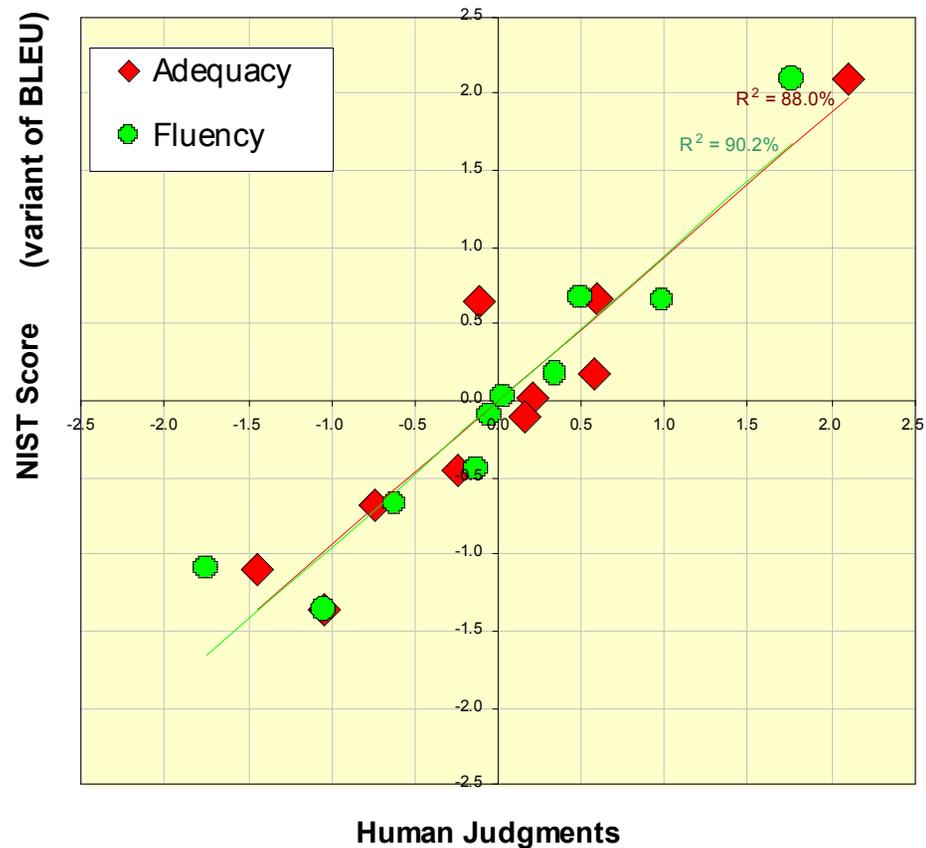
## Reference translation 3:

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 .

# Initial results showed that BLEU predicts human judgments well



slide from G. Doddington (NIST)

# Automatic evaluation of MT

- People started optimizing their systems to maximize BLEU score
  - BLEU scores improved rapidly
  - The correlation between BLEU and human judgments of quality went way, way down
  - StatMT BLEU scores now approach those of human translations but their true quality remains far below human translations
- Coming up with automatic MT evaluations has become its own research field
  - There are many proposals: TER, METEOR, MaxSim, SEPIA, our own RTE-MT
  - TERpA is a representative good one that handles some word choice variation.
- MT research really requires *some* automatic metric to allow a rapid development and evaluation cycle.

# Quiz question!

FOR MONDAY JANUARY 24<sup>TH</sup>

Hyp: The gunman was shot dead by police .

Ref1: The gunman was shot to death by the police .

Ref2: The cops shot the gunman dead .

Compute the unigram precision  $P_1$  and the trigram precision  $P_3$ .

(Note: punctuation tokens are counted, but not sentence boundary tokens.)

- (a)  $P_1 = 1.0$        $P_3 = 0.5$
- (b)  $P_1 = 1.0$        $P_3 = 0.333$
- (c)  $P_1 = 0.875$      $P_3 = 0.333$
- (d)  $P_1 = 0.875$      $P_3 = 0.167$
- (e)  $P_1 = 0.8$        $P_3 = 0.167$