

# **Statistical Machine Translation: IBM Models 1 and 2**

**Cognition Volume 116, Issue 2**

Michael Collins

**<sup>a</sup>Department of Psychology, Rice University, Houston, TX, USA**

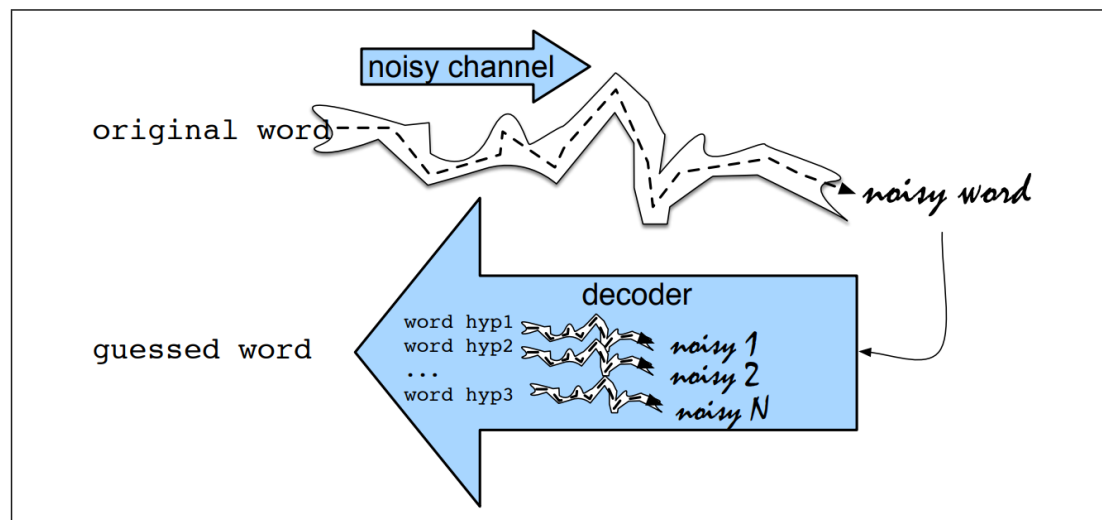
**<sup>b</sup>Department of Chinese, Translation and Linguistics,  
City University of Hong Kong, Kowloon, Hong Kong**

## • IBM Model

- In the noisy-channel approach, the output of the translation model on a new sentence

$$\arg \max_{e \in E} p(e|f) = \arg \max_{e \in E} \frac{p(e)p(f|e)}{\sum_f p(e)p(f|e)} = \arg \max_{e \in E} p(e)p(f|e)$$

- A language model: a probability  $p(e)$  for any sentence  $e = e_1 \dots e_l$  in English
- A translation model: a conditional probability  $p(f|e)$  to any French/English pair of sentences



### Noisy Channel Model

- Imagine that the surface form we see is actually a "distorted" form of an original word passed through a noisy channel
- Decoder passes each hypothesis through a model of this channel and picks the word that best matches the surface noisy word

# Appendix

## • Alignment

- A case where  $l = 6, m = 7$
- Source Sentence  $e = \textit{And the programme has been implemented}$
- Target Sentence  $f = \textit{Le programme a ete mis en application}$
- Alignment variables  $a_1, a_2, \dots, a_7 = \langle 2, 3, 4, 5, 6, 6, 6 \rangle$

### Alignment

<i>Le</i>	$\Rightarrow$	the
<i>Programme</i>	$\Rightarrow$	program
<i>a</i>	$\Rightarrow$	has
<i>ete</i>	$\Rightarrow$	been
<i>mis</i>	$\Rightarrow$	implemented
<i>en</i>	$\Rightarrow$	implemented
<i>application</i>	$\Rightarrow$	implemented

# Appendix

## • Alignment Models: IBM Model 2

- $t(f|e)$  : Conditional probability of generating target word  $f$  from source word  $e$
- $q(j|i, l, m)$  : Probability of alignment variable  $a_i$  (The value  $j$ , source & target  $l$  &  $m$ )

$$p(f_1 \dots f_m, a_1 \dots a_m | e_1 \dots e_l, m) = \prod_{i=1}^m q(a_i | i, l, m) t(f_i | e_{a_i})$$

### Alignment

<i>Le</i>	$\Rightarrow$	the
<i>Programme</i>	$\Rightarrow$	program
<i>a</i>	$\Rightarrow$	has
<i>ete</i>	$\Rightarrow$	been
<i>mis</i>	$\Rightarrow$	implemented
<i>en</i>	$\Rightarrow$	implemented
<i>application</i>	$\Rightarrow$	implemented

$$\begin{aligned}
 & p(f_1 \dots f_m, a_1 \dots a_m | e_1 \dots e_l, m) \\
 = & q(2|1, 6, 7) \times t(Le|the) \\
 & \times q(3|2, 6, 7) \times t(Programme|program) \\
 & \times q(4|3, 6, 7) \times t(a|has) \\
 & \times q(5|4, 6, 7) \times t(ete|been) \\
 & \times q(6|5, 6, 7) \times t(mis|implemented) \\
 & \times q(6|6, 6, 7) \times t(en|implemented) \\
 & \times q(6|7, 6, 7) \times t(application|implemented)
 \end{aligned}$$

## • Independence Assumptions in IBM Model 2

- The chain rule of probabilities to decompose this into two terms

$$\begin{aligned}
 & P(F_1 = f_1 \dots F_m = f_m, A_1 = a_1 \dots A_m = a_m | E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &= P(A_1 = a_1 \dots A_m = a_m | E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &\quad \times P(F_1 = f_1 \dots F_m = f_m | A_1 = a_1 \dots A_m = a_m, E_1 = e_1 \dots E_l = e_l, L = l, M = m)
 \end{aligned}$$

- Independence assumptions  $q(a_i|i, l, m)$ 
  - It is independent of the English words  $E_1 \dots E_l$ , and of the other alignment variables

$$\begin{aligned}
 & P(A_1 = a_1 \dots A_m = a_m | E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &= \prod_{i=1}^m P(A_i = a_i | A_1 = a_1 \dots A_{i-1} = a_{i-1}, E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &= \prod_{i=1}^m P(A_i = a_i | L = l, M = m)
 \end{aligned}$$

$$P(A_i = a_i | L = l, M = m) = q(a_i|i, l, m)$$

## • Independence Assumptions in IBM Model 2

- The chain rule of probabilities to decompose this into two terms

$$\begin{aligned}
 & P(F_1 = f_1 \dots F_m = f_m, A_1 = a_1 \dots A_m = a_m | E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &= P(A_1 = a_1 \dots A_m = a_m | E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &\quad \times P(F_1 = f_1 \dots F_m = f_m | A_1 = a_1 \dots A_m = a_m, E_1 = e_1 \dots E_l = e_l, L = l, M = m)
 \end{aligned}$$

- Independence assumptions  $t(f_i | e_{a_i})$

- It is independent of the English words  $E_1 \dots E_l$ , and of the other alignment variables

$$\begin{aligned}
 & P(F_1 = f_1 \dots F_m = f_m | A_1 = a_1 \dots A_m = a_m, E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &= \prod_{i=1}^m P(F_i = f_i | F_1 = f_1 \dots F_{i-1} = f_{i-1}, A_1 = a_1 \dots A_m = a_m, E_1 = e_1 \dots E_l = e_l, L = l, M = m) \\
 &= \prod_{i=1}^m P(F_i = f_i | E_{a_i} = e_{a_i})
 \end{aligned}$$

$$P(F_i = f_i | E_{a_i} = e_{a_i}) = t(f_i | e_{a_i})$$

# Appendix

## • Alignment

- A case where  $l = 6, m = 7$
- Source Sentence  $e = \textit{And the programme has been implemented}$
- Target Sentence  $f = \textit{Le programme a ete mis en application}$
- Alignment variables  $a_1, a_2, \dots, a_7 = \langle 2, 3, 4, 5, 6, 6, 6 \rangle$

### Alignment

<i>Le</i>	$\Rightarrow$	the
<i>Programme</i>	$\Rightarrow$	program
<i>a</i>	$\Rightarrow$	has
<i>ete</i>	$\Rightarrow$	been
<i>mis</i>	$\Rightarrow$	implemented
<i>en</i>	$\Rightarrow$	implemented
<i>application</i>	$\Rightarrow$	implemented

## • Alignment Models: IBM Model 3

- Mapping relationship between source and target sentence
- Fertility
  - The number of target words generated by a single source word
- Distortion
  - How the position of words changes from source to target

### Phrase alignment (many-to-many alignment)

Elle					
lui					
a					
donné					
un					
livre					
	she	gave	him	a	book
<b>Fertility</b>	1	2	1	1	1

- Fertility is modeled using a probability distribution  $n(\phi|f)$
- $\phi$  = the number of target words produced by the source word  $f$
- Distortion is represented by a probability distribution  $d(j|i, l, m)$
- A word at position  $i$  in the source sentence
- A word at position  $j$  in the target sentence
- The lengths  $l$  and  $m$  of the source and target sentences