

# Statistical Machine Translation

Le Anh Cuong

# Reading

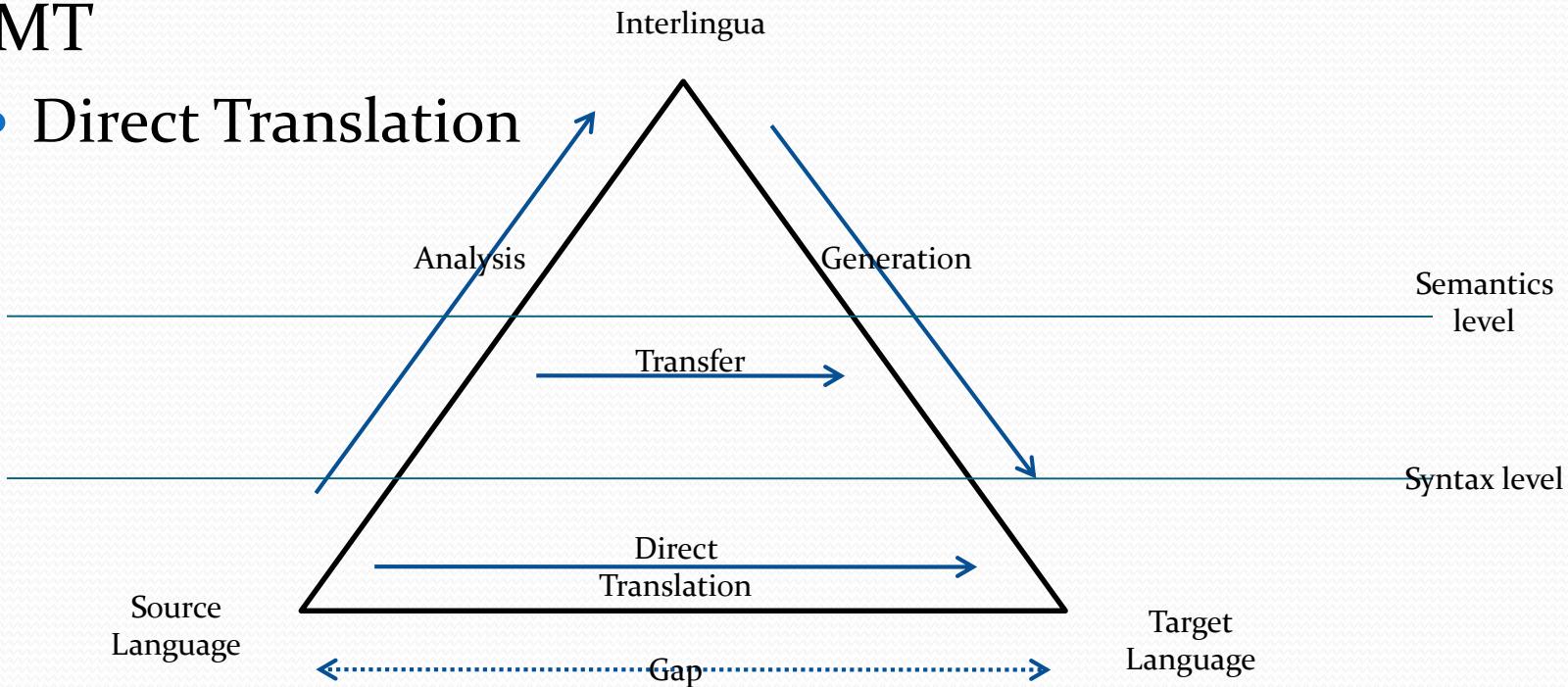
- Chapter 25 [1]
- Chapter 13 [2]

# Contents

- **What Samsung requires to survey**
  - Rule-base and Statistical approach
  - Alignment Model
  - Decoding Algorithms
  - Open Sources
  - Evaluation Methods
  - Using Syntax Info. in SMT

# Machine Translation Pyramid

- RBMT (Rule-Based machine Translation)
  - Analysis, Structure transfer, and Generation
- SMT
  - Direct Translation

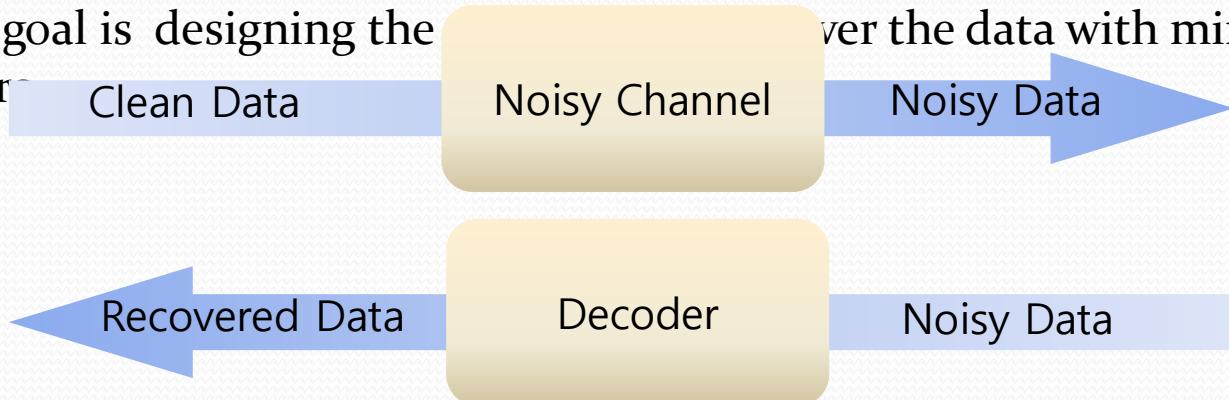


# Comparison

	RBMT	SMT
Approach	Analytic	Empirical
Based on	Transfer rules	Statistical Evidence
Analysis level	Various (morpheme ~ interlingua)	Generally, almost not
Translation Speed	Fast	(Relatively) Slow
Required Knowledge	Linguistic knowledge Dictionary (Ontology) (Conceptual and cultural differences)	Parallel Texts (morphology for spacing)
Adaptability	Low	High

# SMT: Noisy Channel Model

- Noisy Channel
  - Encoder
  - We Get the data through noisy channel
  - The clean (original) data can not be observed directly
  - Noisy channel adds some noise to the data
- Decoder
  - Estimates the original data from the noisy data
  - Recovered data may contain some errors
  - Our goal is designing the error correction algorithm to recover the data with minimum errors



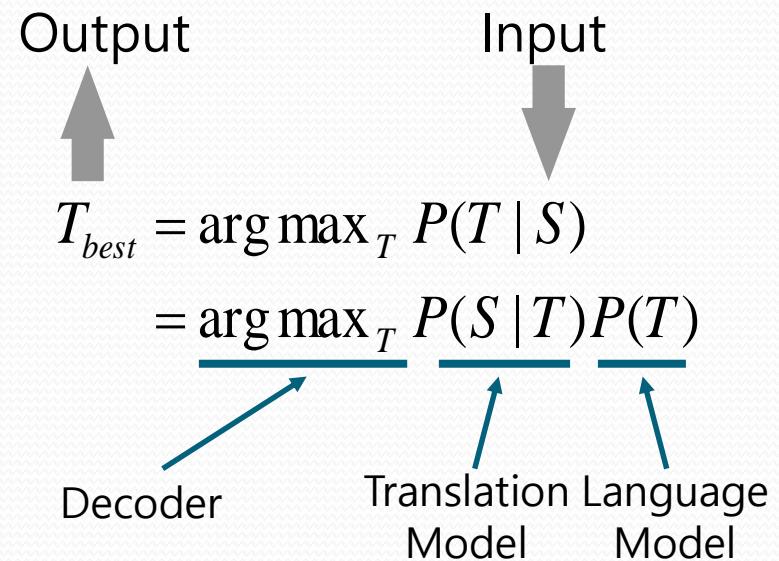
# Noisy Channel Model

- We want to recover original pattern for a given noisy pattern
  - Choose most probable A given B
  - A : Original data
  - B : Noisy data
  - A' : Estimation of original data

$$\begin{aligned}A' &= \arg \max_A \Pr(A | B) \\&= \arg \max_A \frac{\Pr(B | A) \Pr(A)}{\Pr(B)} \\&= \arg \max_A \Pr(B | A) \Pr(A)\end{aligned}$$

# Noisy Channel Model in SMT

- Given a source sentence  $S$  find  $T$  that maximizes probability of  $T$  given  $S$  [Brown et al 1988, 1990]
- Language model:
  - Role : Making fluent sentence
  - Model for target language
- Translation Model:
  - Role : Making correct translation
  - Model for both languages
- Decoder
  - Role : Find a sentence which gives best score
  - We use  $P(S|T)P(T)$  rather than  $P(T|S)$ .



# Noisy Channel model in SMT

- Estimating  $P(S|T)$  or  $P(T|S)$  at a sentence level is impossible
  - Data sparseness problem
    - → Estimate  $P(S|T)$  or  $P(T|S)$  at a smaller level. (Typically, words)
  - Assuming independence of translation units, our approximation is
$$P(T | S) \cong \prod_i P(T_i | S_i)$$
$$P(S | T)P(T) \cong \prod_i P(S_i | T_i)P(T_i)$$
- Actually the assumption is not true, we lose many information
- But,  $P(T)$  may model dependency on previous term
  - $P(T)$  may recover some portion of the lost information
$$P(T_i) \cong P(T_i | T_{i-1} T_{i-2} \dots)$$

# Log-linear Model

- Recently, Log-linear model is popular

- Maximize Sum of logarithms

$$T_{best} = \arg \max_T P(S | T)P(T)$$

$$= \arg \max_T [\log P(S | T) + \log P(T)]$$

- Introduce additional features and weights

$$T_{best} = \arg \max_T [w_1 \log P(S | T) + w_2 \log P(T) + w_3 \log F(S, T)]$$

- Generally, we write

$$T_{best} = \arg \max_T \sum_i w_i f_i(S, T)$$

# Parallel Corpus

- Two or more texts written in different languages have same meaning
- We need alignments at least sentence level
- Example

이곳에서 아침식사를 할 수 있습니까?  
아침 식사는 얼마 지요?  
취사가 가능 합니까?  
짐을 여기 두어도 될까요?  
금요일에 체크아웃 하려고 합니다.  
관광 안내소는 어디 있습니까?  
관광 안내를 바랍니다.  
이곳에서 관광 가이드 한 명을 고용 할 수 있나요?  
한국말을 하는 가이드를 구할 수 있나요?

Can I have breakfast here ?  
How much for breakfast ?  
Can I cook for myself ?  
Can I leave my baggage here ?  
I 'm going to leave on Friday .  
Where can I find tourist information ?  
Can I get some information , please ?  
Can I hire a tour guide here ?  
Is there a Korean-speaking guide available ?

# Word Alignment

- IBM-Model 1~5 [Brown et. al 1993]
  - Finding Best alignment
  - Estimating  $P(S|T)$

가장 가까운 버스 정류장은 어디에  
있습니까 ?

Where is the nearest bus stop ?

Where)

$P(\text{버스} | \text{bus})$

$P(\text{정류장} | \text{stop})$

# N-gram Language Model

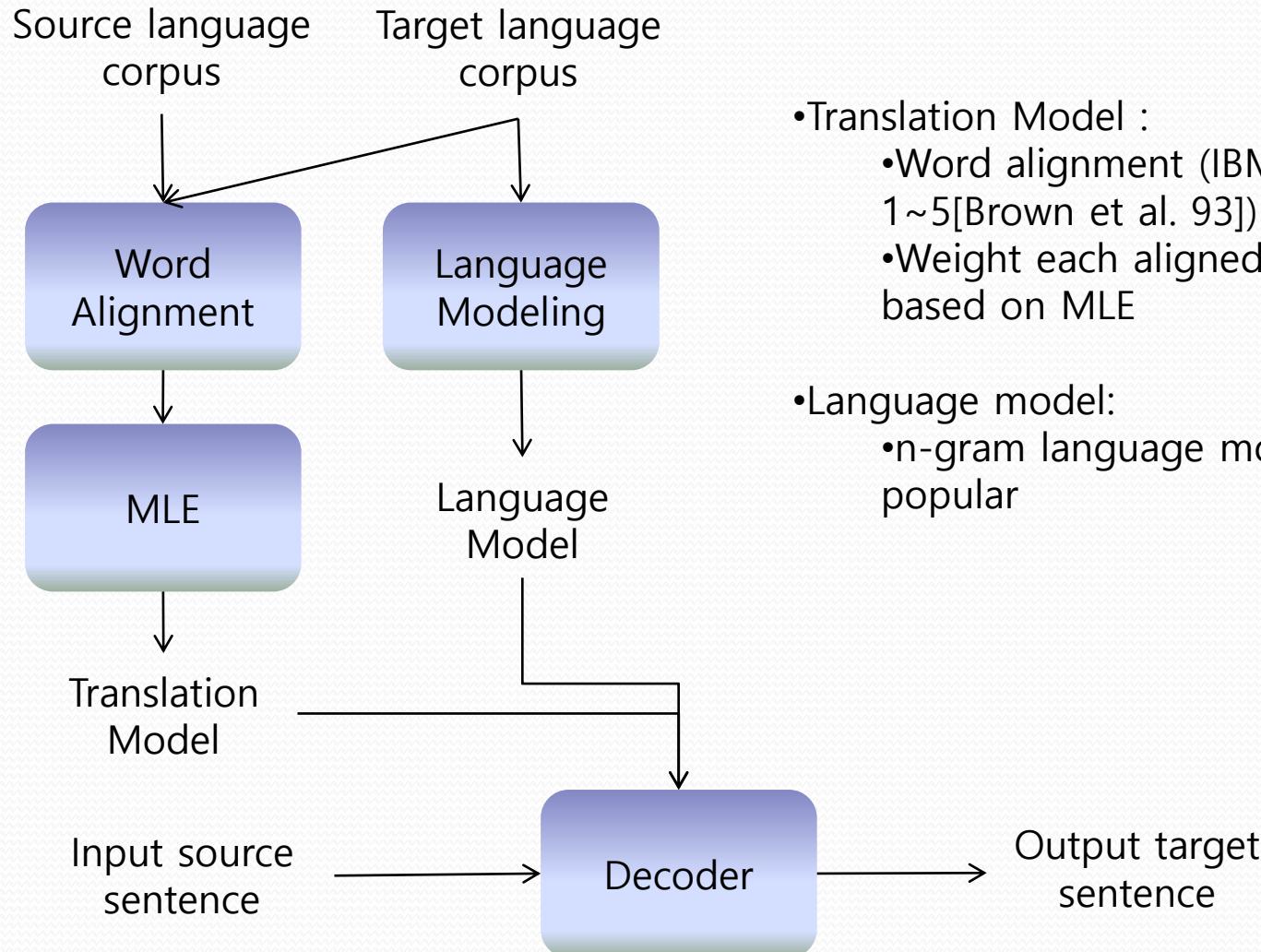
- Probability of next words given history
  - We can not store all the word sequences if the length is not limited
  - The words sequence is very scarce if the length is long
    - For the very scarce data, the probability or statistics is meaningless
- Approximation
  - Assume that the probability of a word is independent to too far history
  - Use limited history
    - 0 history – Unigram
    - 1 history – Bigram
    - 2 history – Trigram
    - ...
- N- gram is most popular method for scoring sentences

$$P(w_n | w_1^{n-1}) \approx P(w_n)$$

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1})$$

$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-2} w_{n-1})$$

# Data flow



- Translation Model :

- Word alignment (IBM Model 1~5[Brown et al. 93])
- Weight each aligned words based on MLE

- Language model:

- n-gram language model is popular

# Contents

- What Samsung requires to survey
  - Rule-base and Statistical approach
  - Alignment Model
  - Decoding Algorithms
  - Open Sources
  - Evaluation Methods
  - Using Syntax Info. in SMT

# Alignment Model

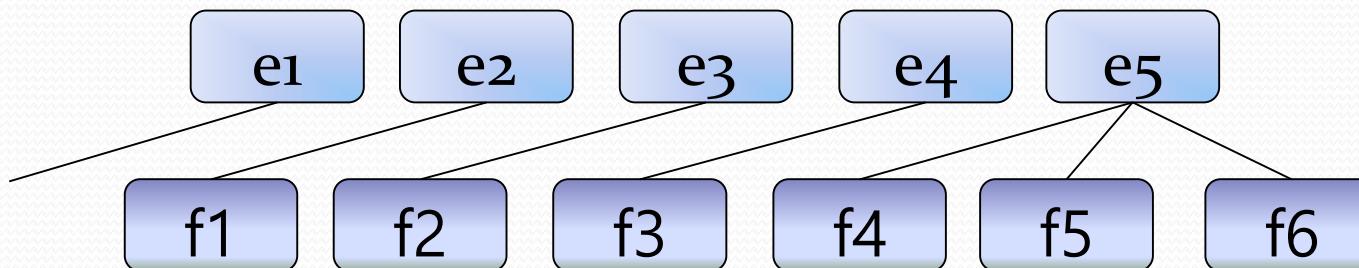
- GIZA++
  - IBM Translation Model 1~5
  - HMM Alignment Model
- Phrase Level Alignment
- Other Word Alignment Models

# IBM Translation Model Outline

- Goal  $\Pr(f \mid e)$ 
  - Modeling the conditional probability distribution
    - f : French sentence (or source sentence)
    - e: English sentence (or target sentence)
- Models [Brown et al. 1993]
  - A series of five translation models: Model1 ~ Model5
  - Train Model1
  - Train Model2 with the result of Model1 training
  - ...
  - Train Model5 with the result of Model4 training
- Algorithm
  - Apply EM algorithm to estimate parameters

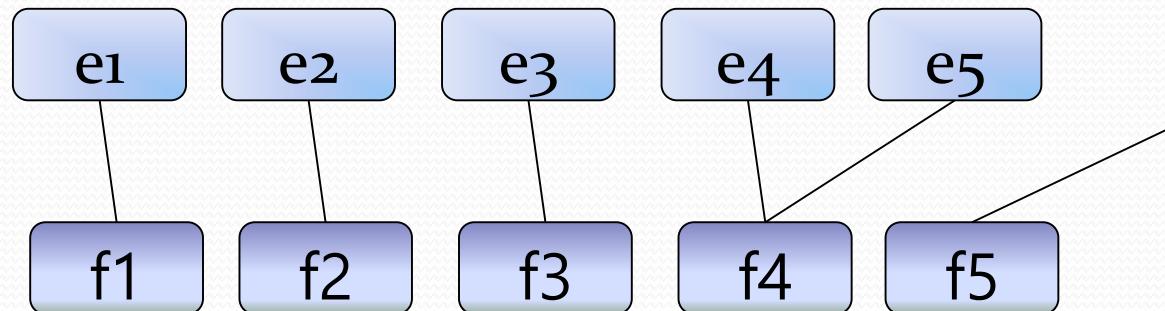
# Word Alignment

- Type 1
  - An alignment with  $t+1$  independent English words
  - $n : 1$  alignment
  - # of possible alignments
    - with  $l$  English words and  $m$  French words
    - for each French words,  $l+1$  alignments are possible including NULL
  - IBM Models use this restriction



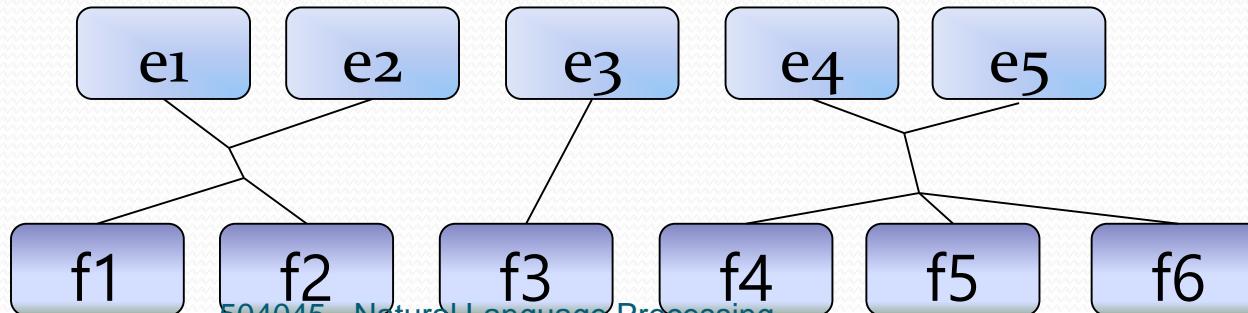
# Word Alignment

- Type 2
  - An alignment with independent French words
  - $1 : n$  alignment  $(m+1)^l$
  - # of possible alignments
    - with  $l$  English words and  $m$  French words
    - for each English words,  $m+1$  alignments are possible including NULL
  - Inverse direction



# Word Alignment

- Type 3
  - An alignment with independent English words
  - $n : n'$  alignment, general alignment
  - # of possible alignments : Very large
    - with  $m$  English words and  $l$  French words
    - For the first English word  $2^{l+1}$  alignments are possible
    - For the second English word  $2^{l+1-c}$  alignments are possible where  $c$  is number of French words aligned with the first English word
  - Alignment for Phrase based Machine Translation



# Word Alignment: variable

- Our goal is estimating the conditional probability

$$\Pr(\mathbf{f} \mid \mathbf{e})$$

- We can introduce a hidden variable  $\mathbf{a}$  (that is word alignment)

$$\Pr(\mathbf{f} \mid \mathbf{e}) = \sum_{\mathbf{a}} \Pr(\mathbf{f}, \mathbf{a} \mid \mathbf{e})$$

- Assume that each French word has exactly one connection
  - The word alignment  $\mathbf{a}$  can be represented by a series

$$a_1^m \equiv a_1 a_2 \dots a_m$$

- Values are between 0 and  $l$ , where  $l$  is the length of English sentence
- if  $a_2 = 4$ , it means
  - “The French word at position 2 aligned with the English word at position 4”
- Position 0 is reserved for the null word

# Model 1,2 Likelihood: Exact Equation

- A possible form of exact equation
  - “Exact” means that it is not an approximation

$$\Pr(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \underbrace{\Pr(m | \mathbf{e})}_{\text{1}} \prod_{j=1}^m \underbrace{\Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e})}_{\text{2}} \underbrace{\Pr(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e})}_{\text{3}}$$

- Part 1. Choose the length of the French string
  - Given English string
- Part 2. Choose where to connect
  - Given English string
  - Given length of French string
  - Given history
- Part 3. Choose identity of the word
  - Given English string
  - Given length of French string
  - Given history
  - Given current word alignment (Part 2)

# Alignment Process: Model 1,2

- An alignment process corresponding to exact equation on previous page
- Choose a length for French string  $f$
- **for**  $i = 1$  to  $m$
- **begin**
  - Decide which position in  $e$  is connected to  $f_i$
  - Decide which identity of  $f_i$  is
- **end**

$$\Pr(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \Pr(m | \mathbf{e}) \prod_{j=1}^m \Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e}) \Pr(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e})$$

# Model 1

- Exact equation
  - Too complex → We need some approximation
- Approximation  $\Pr(m | \mathbf{e})$ 
  - Part 1:
    - Assume that it is independent of  $m$  and  $\mathbf{e}$
    - $\rightarrow \epsilon$ , constant  
 $\Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e})$
  - Part 2:
    - Depends only on the length of the English string
    - $\rightarrow$ , uniform  
 $\Pr(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e})$
  - Part 3:
    - Depends only on the French word and corresponding English word:
    - $\rightarrow$ , translation probability

# Model 1

- Likelihood function

$$\Pr(\mathbf{f}, \mathbf{a} | \mathbf{e}) = \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m t(f_j | e_{a_j})$$

- The Simplest model
- The order of the words in  $\mathbf{e}$  and  $\mathbf{f}$  does not affect the likelihood
- Model 1 likelihood function has only one maximum
  - Model 1 always finds global maximum

# Model 2

- Approximation  $\Pr(m | \mathbf{e})$ 
  - Part 1:
    - Assume that it is independent of  $m$  and  $e$
    - $\rightarrow \epsilon$ , constant
    - same to Model 1
  - Part 2:  $\Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, \mathbf{e})$ 
    - Depends on
      - Positions of French word and corresponding English word ( $j, a_j$ )
      - Length of French string and English string ( $m, l$ )
    - We introduce alignment probabilities
$$a(a_j | j, m, l) = \Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, l)$$
- Part 3:  $\Pr(f_j | a_1^j, f_1^{j-1}, m, \mathbf{e})$ 
  - Depends only on the French word and corresponding English word  $f_j, e_{a_j}$
  - $\rightarrow t(f_j | e_{a_j})$ , translation probability
  - same to Model 1

# Model 2

- Likelihood function

$$\Pr(\mathbf{f} \mid \mathbf{e}) = \varepsilon \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j \mid e_{a_j}) a(a_j \mid j, m, l)$$

- Alignment probability is introduced compared to Model 1
- Model 1 is a special case of Model 2

# Fertility

- Definition : Fertility of  $e$ 
  - A random variable  $\Phi_e$  that corresponds to the number of French words to which  $e$  is connected in a randomly selected alignment
- Modeling the Fertility
  - Model 1 and 2 : not clear
  - Model 3, 4 and 5 : Parameterize fertilities directly
- Tablet
  - A list of French words to connect to each English word
- Tableau
  - The collection of tablets, A random variable
  - $T_i$  : the tablet for  $i^{\text{th}}$  English word
  - $T_{ik}$  :  $k^{\text{th}}$  French word in the  $i^{\text{th}}$  tablet

# Model 3, 4 and 5 Likelihood: Exact Equation

- The Joint likelihood for a tableau,  $\tau$ , and a permutation,  $\pi$

$$\Pr(\tau, \pi | \mathbf{e}) = \prod_{i=1}^l \Pr(\phi_i | \phi_1^{i-1}, \mathbf{e}) \Pr(\phi_0 | \phi_1^l, \mathbf{e}) \times \prod_{i=0}^l \prod_{k=1}^{\phi_i} \Pr(\tau_{ik} | \tau_{i1}^{k-1}, \tau_{i1}^{i-1}, \phi_0^l, \mathbf{e}) \times \\ \prod_{i=1}^l \prod_{k=1}^{\phi_i} \Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_{i1}^{i-1}, \tau_0^l, \phi_0^l, \mathbf{e}) \times \prod_{k=1}^{\phi_0} \Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, \mathbf{e})$$

- Knowing  $\tau$  and  $\pi$  determine a French string and an alignment
  - Different  $\tau$  and  $\pi$  may lead to same pair  $f, a$
  - $\langle f, a \rangle$ : pairs of  $\tau$  and  $\pi$  that lead pair  $f$  and  $a$
- From the above, we have  $\Pr(f, a | \mathbf{e}) = \sum_{(\tau, \pi) \in \langle f, a \rangle} \Pr(\tau, \pi | \mathbf{e})$

# Alignment Process: Model 3, 4, 5

- for each English word
- begin
  - Decide the fertility of the word
  - Get a list of French words to connect to the word
- end
- Permute words in tableau to generate  $f$

$$\Pr(\tau, \pi | \mathbf{e}) = \prod_{i=1}^l \Pr(\phi_i | \phi_1^{i-1}, \mathbf{e}) \Pr(\phi_0 | \phi_1^l, \mathbf{e}) \times \prod_{i=0}^l \prod_{k=1}^{\phi_i} \Pr(\tau_{ik} | \tau_{i1}^{k-1}, \tau_0^{i-1}, \phi_0^l, \mathbf{e}) \times \prod_{i=1}^l \prod_{k=1}^{\phi_i} \Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_1^{i-1}, \tau_0^l, \phi_0^l, \mathbf{e}) \times \prod_{k=1}^{\phi_0} \Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, \mathbf{e})$$

# Model-3

- Approximation  $\Pr(\phi_i | \phi_1^{i-1}, \mathbf{e})$ 
  - Part 1 :
    - Depends only on  $\phi_i$  and  $e_i$
    - Fertility probability
    - $\rightarrow n(\phi_i | e_i)$
  - Part 2 :  $\Pr(\tau_{ik} | \tau_{i1}^{k-1}, \tau_0^{i-1}, \phi_0^l, \mathbf{e})$ 
    - Depends only on  $\tau_{ik}$  and  $e_i$
    - Translation probability
    - $\rightarrow t(f | e_i)$
  - Part 3 :  $\Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_1^{i-1}, \tau_0^l, \phi_0^l, \mathbf{e})$ 
    - Depends only on  $\pi_i$ ,  $i$ ,  $m$ ,  $l$
    - Distortion probability
    - $\rightarrow d(j | i, m, l)$

# Model-3

- Approximation (cont')  $\Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, \mathbf{e})$

- Part 4 :

$$\begin{aligned} & \frac{1}{\phi_0 - k} \\ \bullet \rightarrow & \prod_{k=1}^{\phi_0} \Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, \mathbf{e}) = \frac{1}{\phi_0!} \end{aligned}$$

$$\Pr(\phi_0 | \phi_1^l, \mathbf{e})$$

- Part 5 :

$$\binom{\phi_1 + \dots + \phi_l}{\phi_0} p_0^{\phi_1 + \dots + \phi_l - \phi_0} p_1^{\phi_0} = \binom{m - \phi_0}{\phi_0} p_0^{m - 2\phi_0} p_1^{\phi_0}$$

From the above approximations ...

$$\Pr(\mathbf{f} | \mathbf{e}) = \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \binom{m - \phi_0}{\phi_0} p_0^{m - 2\phi_0} p_1^{\phi_0} \prod_{i=1}^l \phi_i! n(\phi_i | e_i) \times \prod_{j=1}^m t(f_j | e_{a_j}) d(j | a_j, m, l)$$

$$\Pr(f | e) = \sum_a \Pr(f, a | e) = \sum_a \sum_{(\tau, \pi) \in \langle f, a \rangle} \Pr(\tau, \pi | e)$$

$$\Pr(\tau, \pi | e) = \prod_{i=1}^l \Pr(\phi_i | \phi_1^{i-1}, e) \Pr(\phi_0 | \phi_1^l, e) \times \prod_{i=0}^l \prod_{k=1}^{\phi_i} \Pr(\tau_{ik} | \tau_{i1}^{k-1}, \tau_0^{i-1}, \phi_0^l, e) \times$$

$$\prod_{i=1}^l \prod_{k=1}^{\phi_i} \Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_1^{i-1}, \tau_0^l, \phi_0^l, e) \times \prod_{k=1}^{\phi_0} \Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, e)$$

$$\rightarrow \Pr(f | e) = \sum_a \sum_{(\tau, \pi) \in \langle f, a \rangle} \prod_{i=1}^l \Pr(\phi_i | \phi_1^{i-1}, e) \Pr(\phi_0 | \phi_1^l, e) \times \prod_{i=0}^l \prod_{k=1}^{\phi_i} \Pr(\tau_{ik} | \tau_{i1}^{k-1}, \tau_0^{i-1}, \phi_0^l, e) \times$$

$$\prod_{i=1}^l \prod_{k=1}^{\phi_i} \Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_1^{i-1}, \tau_0^l, \phi_0^l, e) \times \prod_{k=1}^{\phi_0} \Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, e)$$

$$\Pr(f | e) = \sum_{a_1} \dots \sum_{a_m} \sum_{(\tau, \pi) \in \langle f, a \rangle} \prod_{i=1}^l \Pr(\phi_i | \phi_1^{i-1}, e) \Pr(\phi_0 | \phi_1^l, e) \times \prod_{i=0}^l \prod_{k=1}^{\phi_i} \Pr(\tau_{ik} | \tau_{i1}^{k-1}, \tau_0^{i-1}, \phi_0^l, e) \times$$

$$\prod_{i=1}^l \prod_{k=1}^{\phi_i} \Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_1^{i-1}, \tau_0^l, \phi_0^l, e) \times \prod_{k=1}^{\phi_0} \Pr(\pi_{0k} | \pi_{01}^{k-1}, \pi_1^l, \tau_0^l, \phi_0^l, e)$$

$a \equiv a_1, a_2, \dots$   
의 조합으로 봄

Model 3  
approximation

$$\Pr(f | e) = \sum_{a_1} \dots \sum_{a_m} \sum_{(\tau, \pi) \in \langle f, a \rangle} \prod_{i=1}^l n(\phi_i | e_i) \binom{m - \phi_0}{\phi_0} p_0^{\phi_0} p_1^{\phi_0} \times \prod_{j=1}^m t(f_j, e_{a_j}) \times$$

$$\phi_0! \prod_{j=1}^m d(j | a_j, m, l) \times \frac{1}{\phi_0!}$$

Distortion에서  
도출 된 성분과  
Part 4가 서로  
상쇄됨(약분)

$$\Pr(f | e) = \sum_{a_1} \dots \sum_{a_m} \sum_{(\tau, \pi) \in \langle f, a \rangle} \prod_{i=1}^l n(\phi_i | e_i) \binom{m - \phi_0}{\phi_0} p_0^{m-2\phi_0} p_1^{\phi_0} \times \prod_{j=1}^m t(f_j, e_{a_j}) \times \\ \prod_{j=1}^m d(j | a_j, m, l)$$

In Model 3 approximation,

If  $(\tau_1, \pi_1)$  and  $(\tau_2, \pi_2)$  represents same  $\langle f, a \rangle$ ,

then  $\Pr(\tau_1, \pi_1 | e) = \Pr(\tau_2, \pi_2 | e)$

The number of elements in  $\langle f, a \rangle$  is  $\prod_{i=0}^l \phi_i!$

$$\Pr(f | e) = \sum_{a_1} \dots \sum_{a_m} \prod_{i=1}^l \phi_i! \prod_{i=1}^l n(\phi_i | e_i) \binom{m - \phi_0}{\phi_0} p_0^{m-2\phi_0} p_1^{\phi_0} \times \prod_{j=1}^m t(f_j, e_{a_j}) \times \\ \prod_{j=1}^m d(j | a_j, m, l)$$

식을 정리하면

$$\Pr(f | e) = \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \binom{m - \phi_0}{\phi_0} p_0^{m-2\phi_0} p_1^{\phi_0} \prod_{i=1}^l \phi_i! n(\phi_i | e_i) \times \prod_{j=1}^m t(f_j | e_{a_j}) d(j | a_j, m, l)$$

# Model 4

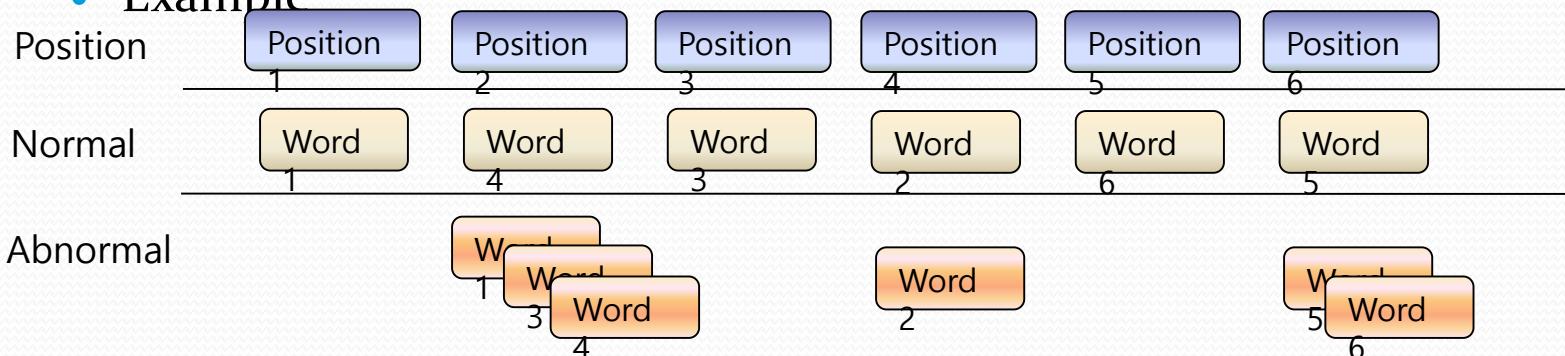
- A Problem of Model 3
  - Phrases should be considered as a unit
  - Model 3 does not account well for phrases
    - Every words are moved independently
  - Model 4 modifies Model 3 to consider phrases
    - Modeling the phrase property
    - Modifying distortion model

# Model 4

- Solution
  - Replace the distortion model by two sets of parameters
    - A parameter for head of each cept
    -
  - A parameter for remaining part of each cept
$$\Pr(\Pi_{[i]1} = j | \pi_1^{[i]-1}, \tau_0^l, \phi_0^l, \mathbf{e}) = d_1(j - \otimes_{i-1} | A(e_{[i-1]}), B(f_j))$$
  - Terms
    - [i] : The position in the English string of the  $i^{\text{th}}$  one-word cept.
    - : Ceiling of the average value of the positions in the French string of the words from  $i^{\text{th}}$  tablet
    - A(.),B(.) : function that changes words into some class of vocabulary
  - cept :
    - Roughly, a unit of concepts , it can be a translation unit
    - A set of English words connected to a French word in a particular alignment

# Deficiency

- A problem with distortion probability (for Model 3, 4)
  - We assumed  $\Pr(\pi_{ik} | \pi_{i1}^{k-1}, \pi_1^{i-1}, \tau_0^l, \phi_0^l, e)$  depends only on  $j, i, m$ , and  $l$
  - The distortion probabilities for assigning positions to later words do not depend on the positions assigned to earlier words
    - Multiple words can have same position
    - Empty position also possible
    - Example



- The model assigns some probability mass to the abnormal string
- For this problem we say that it is “deficient”

# Model-5

- Model 3 and Model 4 are deficient
- Model 5 remove the deficiency
- Restriction to assigning position
  - We need to avoid unavailable positions when we assign the positions
    - Model 3 and 4 do not consider this point
  - First, we define  $v_j$ 
    - the number of available position up to j including
$$v_j = v(j, \tau_1^{[i]-1}, \tau_{[i]1}^{k-1})$$

0, if  $j$  is not available  
1, if  $j$  is available
- Rewrite the distortion probabilities of model 4
  - For head  $d_1(v_j | B(f_j), v_{\oplus i=1}, v_m - \phi_{[i]} + 1) (1 - \delta(v_j, v_{j-1}))$
  - For remaining parts
$$d_{>1}(v_j - v_{\pi[i]k-1} | B(f_j), v_m - v_{\pi[i]k-1} - \phi_{[i]} + k) (1 - \delta(v_j, v_{j-1}))$$

# Summary of IBM Model 1~5

	Probability models			Other features
	Translation	Alignment / Distortion	Length/ Fertility	
Model 1	$t(f_j   e_{a_j})$	$(l+1)^{-1}$	Constant	Unique maxima
Model 2	"	$a(a_j   j, m, l)$	Constant	-
Model 3	"	$d(j   i, m, l)$ $(\phi_0 - k)^{-1}$	$n(\phi   e_i)$ $_{m-\phi_0} C_{\phi_0} p_0^{m-2\phi_0} p_1^{\phi_0}$	Introduced fertility
Model 4	"	$d_{>1}(j - \pi_{[i]k-1}   B(f_j))$ $d_1(j - \odot_{i-1}   A(e_{[i-1]}), B(f_j))$ $(\phi_0 - k)^{-1}$	"	Phrase property
Model 5	"	$d_1(v_j   B(f_j), v_{\oplus_{i=1}}, v_m - \phi_{[i]} + 1)(1 - \delta(v_j, v_{j-1}))$ $d_{>1}(v_j - v_{\pi[i]k-1}   B(f_j), v_m - v_{\pi[i]k-1} - \phi_{[i]} + k)(1 - \delta(v_j, v_{j-1}))$ $(\phi_0 - k)^{-1}$	"	Removed deficiency

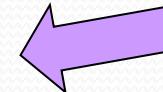
# Phrase Extraction

- Phrase level alignment
- Pharaoh's process
  - Get word alignments in both directions
    - From the GIZA++, IBM Model 4
    - Bi-directional word alignment (source to target, target to source)
  - Intersect the word alignments
  - Expand the intersection to the union
    - Use some heuristic to resolve conflict
    - Pharaoh presents 6 heuristics
  - Extract all possible phrase pairs which consistent with word alignment
  - Assign probabilities to the phrase pairs
    - Count the phrase co-occurrences
    - Divide it by count of occurrence of phrase e

# Bidirectional alignment

- Intersection and Union

Inter-sect	주	한	잔	주	세 요	.
A						
Draft	█					
Beer						
,				█		
Please						█
.						█



E-k	생 맥 주	한	잔	주	세 요	.
A						
Draft	█					
Beer		█	█			
,					█	
Please					█	
.						█

Intersection



Inter-sect	생 맥 주	한	잔	주	세 요	.
A	█					
Draft	█					
Beer		█	█			
,				█		
Please					█	
.						█

K-E	생 맥 주	한	잔	주	세 요	.
A	█					
Draft	█					
Beer		█	█			
,					█	
Please					█	
.						█

GIZA++  
results

# Phrase Extraction

- Learning all phrase pairs that are consistent with the word alignment
  - One should limit the maximum length of phrases

Inter-sect	생 맥 주	한	잔	주	세요	.
A						
Draft						
Beer						
,						
Please						
.						

- (A Draft | 생맥주) ( Beer | 한 잔 ) (, | 주) (Please | 세요) (. | .)

# Phrase Extraction

- Learning all phrase pairs that are consistent with the word alignment

Inter- sect	생 맥 주	한	잔	주	세요	.
A	■					
Draft	■					
Beer		■	■	■		
,		■		■	■	
Please				■	■	■
.					■	■

- (A Draft | 생맥주) (Beer | 한 잔) (, | 주) (Please | 세요) (. | .)
- (A Draft Beer | 생맥주 한 잔) (Beer , | 한 잔 주) (, Please | 주 세요) (Please . | 세요 .)

# Phrase Extraction

- Learning all phrase pairs that are consistent with the word alignment

Inter-sect	생 맥 주	한	잔	주	세 요	.
A	■					
Draft	■					
Beer		■	■	■		
,				■		
Please					■	■
.						■

- (A Draft | 생맥주) (Beer | 한 잔) (, | 주) (Please | 세요) (, | .)
- (A Draft Beer | 생맥주 한 잔) (Beer , | 한 잔 주) (, Please | 주 세요) (Please . | 세요 .)
- (A Draft Beer , | 생맥주 한 잔 주) (Beer , Please | 한 잔 주 세요) (, Please | 주 세요 .)

# Phrase Extraction

- Learning all phrase pairs that are consistent with the word alignment

Inter-sect	생 맥 주	한	잔	주	세요	.
A	blue					
Draft	blue					
Beer	blue	blue	blue	red		
,	blue	blue	blue	blue		
Please	green			blue	red	red
.				blue	blue	blue

- (A Draft | 생맥주) (Beer | 한 잔) (, | 주) (Please | 세요) (, | .)
- (A Draft Beer | 생맥주 한 잔) (Beer, | 한 잔 주) (, Please | 주 세요) (Please . | 세요 .)
- (A Draft Beer, | 생맥주 한 잔 주) (Beer, Please | 한 잔 주 세요) (, Please | 주 세요 .)
- (A Draft Beer, Please | 생맥주 한 잔 주 세요) (Beer, Please . | 한 잔 주 세요 .)

# Phrase Extraction

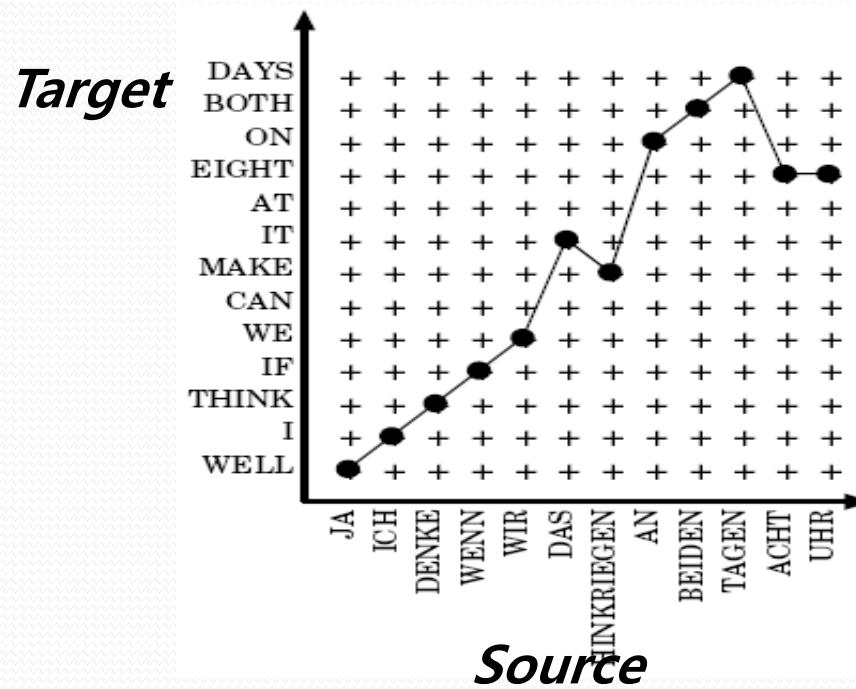
- Learning all phrase pairs that are consistent with the word alignment

Inter-sect	생 맥 주	한	잔	주	세 요	.
A	█					
Draft	█					
Beer		█	█	█		
,	█	█	█	█		
Please				█	█	█
.		█			█	█

- (A Draft | 생맥주) (Beer | 한 잔) ( | 주) (Please | 세요) ( | .)
- (A Draft Beer | 생맥주 한 잔) (Beer , | 한 잔 주) (, Please | 주 세요) (Please . | 세요 .)
- (A Draft Beer , | 생맥주 한 잔 주) (Beer , Please | 한 잔 주 세요) (, Please | 주 세요 .)
- (A Draft Beer , Please | 생맥주 한 잔 주 세요) (Beer , Please . | 한 잔 주 세요 .)
- (A Draft Beer , Please . | 생맥주 한 잔 주 세요 .)

# HMM Alignment Model

- Goal : Improve IBM Model 1-2
- Idea : relative position model



# HMM Alignment Model

- The alignment model  $\Pr(f_1^J, a_1^J | e_1^I)$  can be structured without loss of generality as follows: (IBM Model 1,2)

$$\begin{aligned}\Pr(f_1^J, a_1^J | e_1^I) &= \Pr(J | e_1^I) \prod_{j=1}^J \Pr(f_j, a_j | f_1^{j-1}, a_1^{j-1}, e_1^I) \\ &= \Pr(J | e_1^I) \prod_{j=1}^J \Pr(a_j | f_1^{j-1}, a_1^{j-1}, e_1^I) \Pr(f_j | f_1^{j-1}, a_1^{j-1}, e_1^I)\end{aligned}$$

- In HMM alignment model
  - A first-order dependence for the alignments  $a_j$
  - The lexicon probability depends only on the word at position  $a_j$

$$\Pr(f_j | f_1^{j-1}, a_1^{j-1}, e_1^I) = \Pr(f_j | e_{a_j})$$

# HMM Alignment Model

- Alignment probability:

- A simple length model  $p(J | e_1^I) = p(J | I)$

$$\Pr(f_1^J | e_1^I) = p(J | I) \sum_{a_1^J} \prod_{j=1}^J p(a_j | a_{j-1}, I) p(f_j | e_{a_j})$$

- Alignment depends on relative position

$$p(i | i', I) = \frac{c(i - i')}{\sum_{i''=1}^I c(i'' - i')}$$

- Maximum approximation:

$$\Pr(f_1^J | e_1^I) \cong p(J | I) \max_{a_1^J} \prod_{j=1}^J p(a_j | a_{j-1}, I) p(f_j | e_{a_j})$$

# Other Alignment Method

- Heuristic Method
  - Dictionary Look up
  - Transliteration and string similarity
  - Nearest aligned neighbor (alignment locality)
  - POS affinities
  - ...
- Hybrid Method
  - Combing two or more methods
  - Intersection, Union, Voting, ...
- Variants of IBM Models and HMM Model

# Contents

- What Samsung requires to survey
  - Rule-base and Statistical approach
  - Alignment Model
  - Decoding Algorithms
  - Open Sources
  - Evaluation Methods
  - Using Syntax Info. in SMT

# Decoding Algorithms

- Beam Search Style
  - Phrase-based Systems
  - Pharaoh, Moses and its variants
- CFG Parsing style
  - Syntax-based Systems, SMT by parsing
  - Hiero, GenPar, ...

# Pharaoh Decoding

- Translation options
  - In a sentence of length n, there are  $\sum_{i=1}^n i$  phrases
  - A translation option is a possible translation of a phrase

1      2      3      4      5

**나      는      소년      입니다 .**

I      boy      is      .

me      a boy am ?

I      am a boy

am a

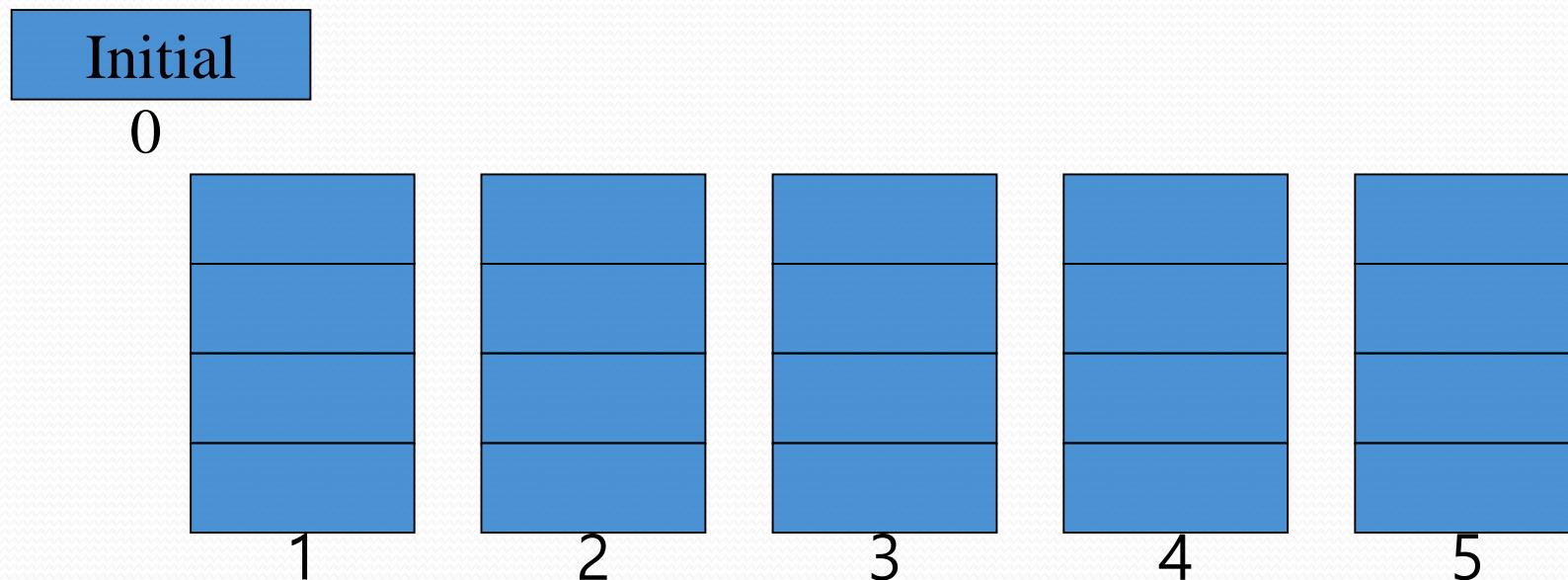
am a boy .

# Pharaoh Decoding

- Hypothesis
  - Hypothesis is a Partial translation taken by applying some translation options
  - Contents (data structure)
    - A back link to the previous state
    - The foreign words covered so far
    - The last two native words generated
    - The end of the last foreign phrase covered
    - The last added native phrase
    - The cost so far
    - An estimate of the future cost

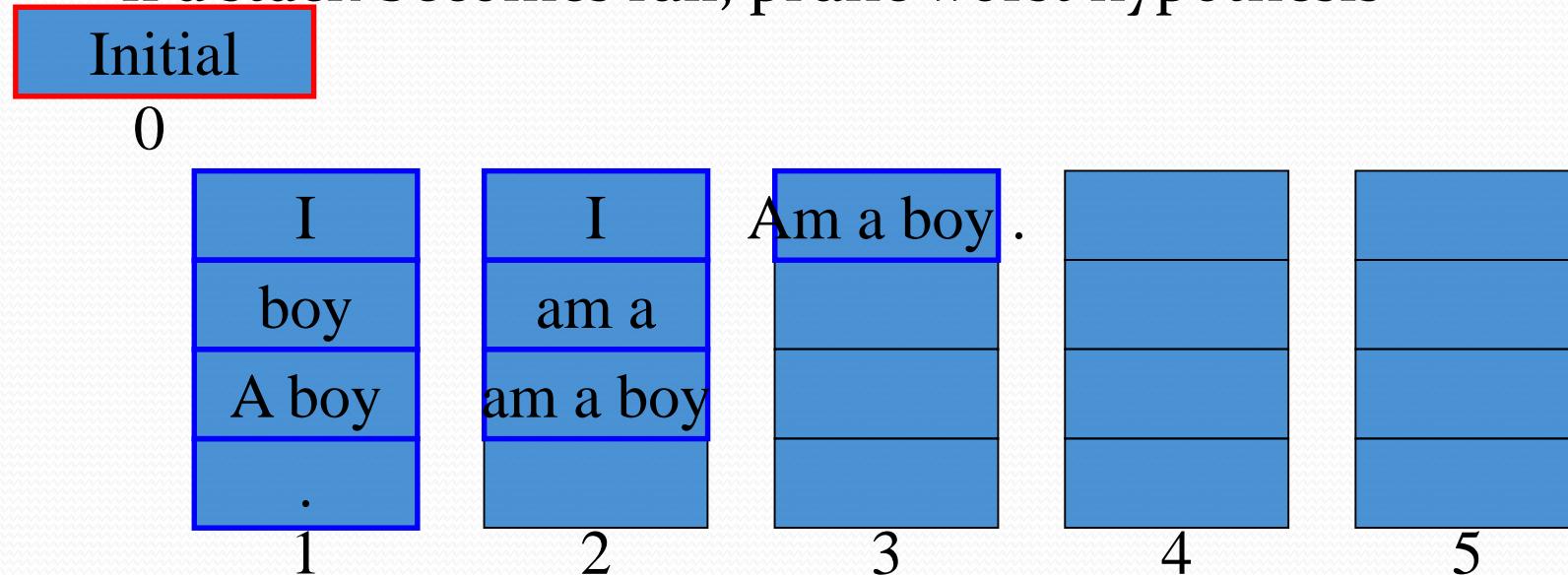
# Pharaoh Decoding

- Decoding
  - Initialize hypothesis stack
  - Create initial hypothesis



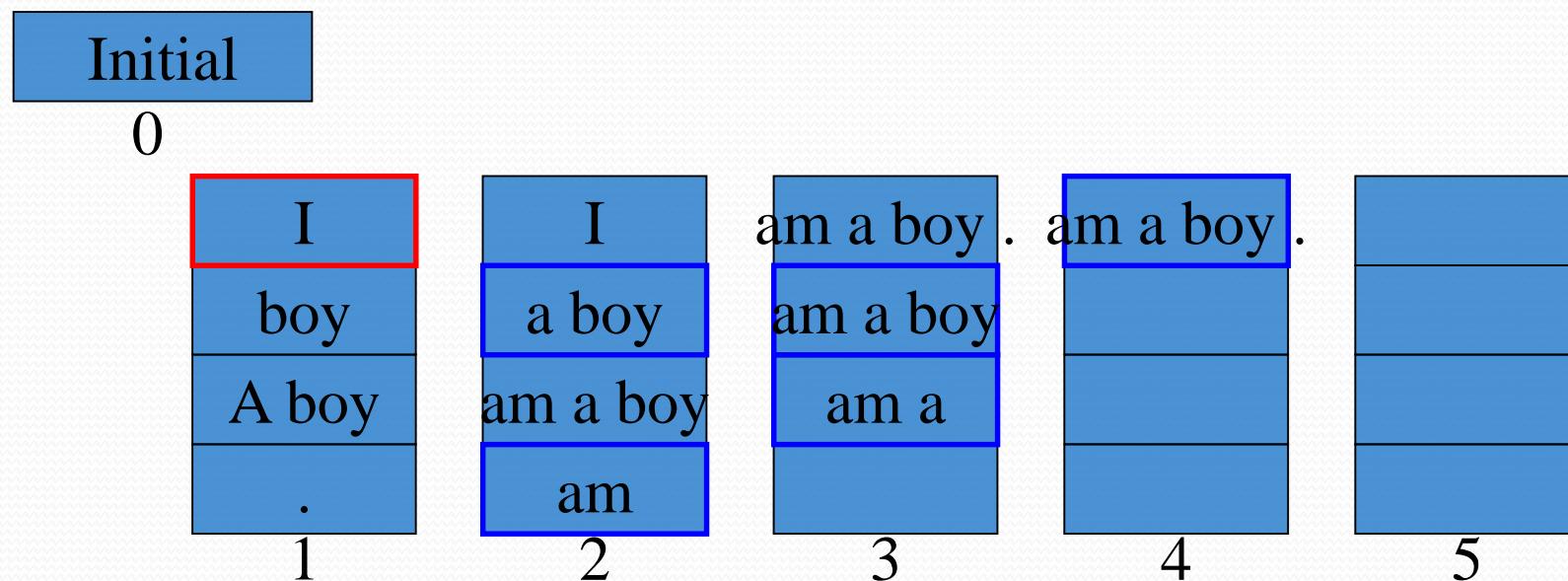
# Pharaoh Decoding

- Decoding
  - Derive new hypothesis from previous hypothesis by applying possible translation options
  - If a stack becomes full, prune worst hypothesis



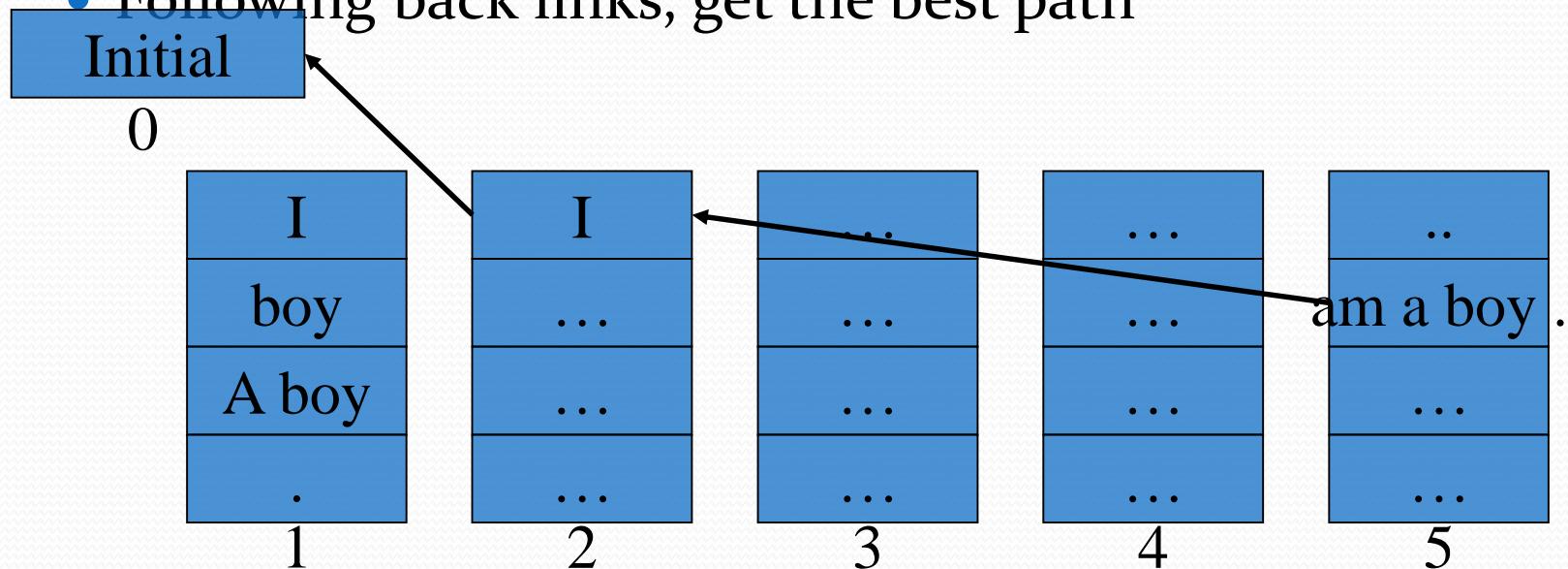
# Pharaoh Decoding

- Decoding
  - Apply translation options for each hypothesis in stack1~4.



# Pharaoh Decoding

- Decoding
  - After processing last element of stack4
  - Fine the best hypothesis in the stack5
  - Following back links, get the best path

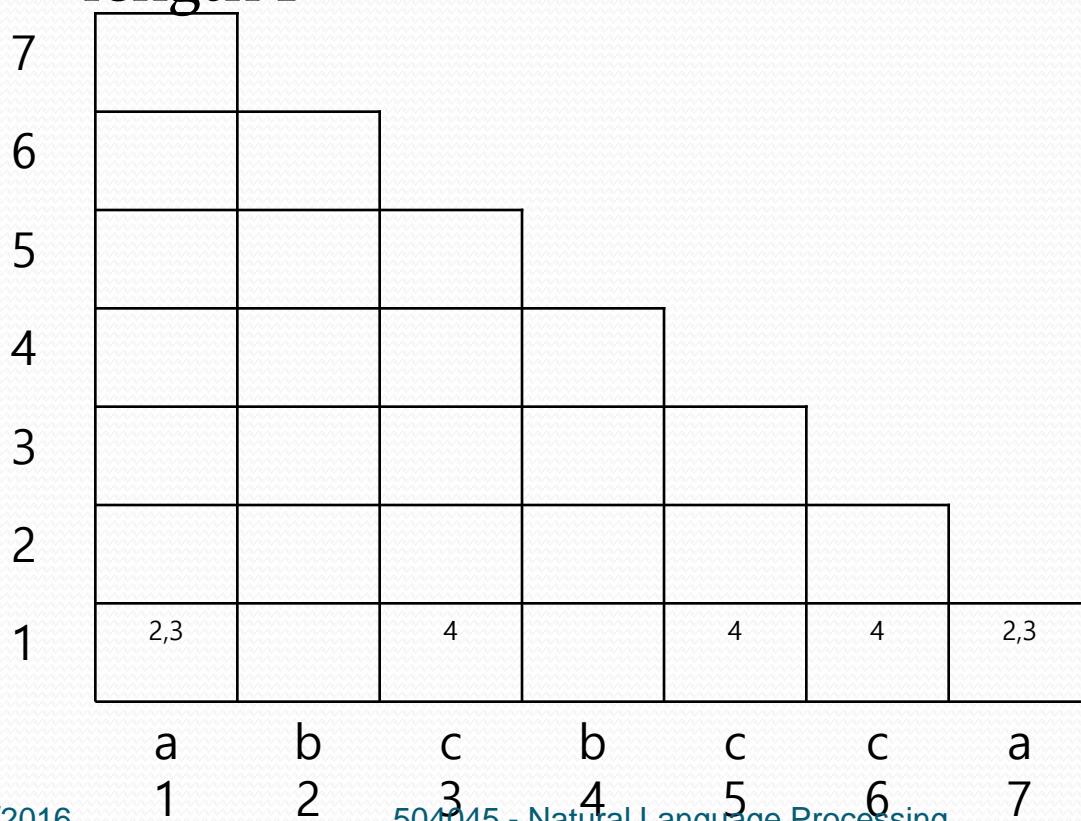


# Hiero Decoding

- Goal
  - CKY Algorithm
    - Reject or Accept a string for given grammar rules
  - Decoding for Translation
    - Ultimate goal: Get Most probable string
    - Practical goal: Get Most probable derivation

# Hiero Decoding

- Example parsing
  - length 1

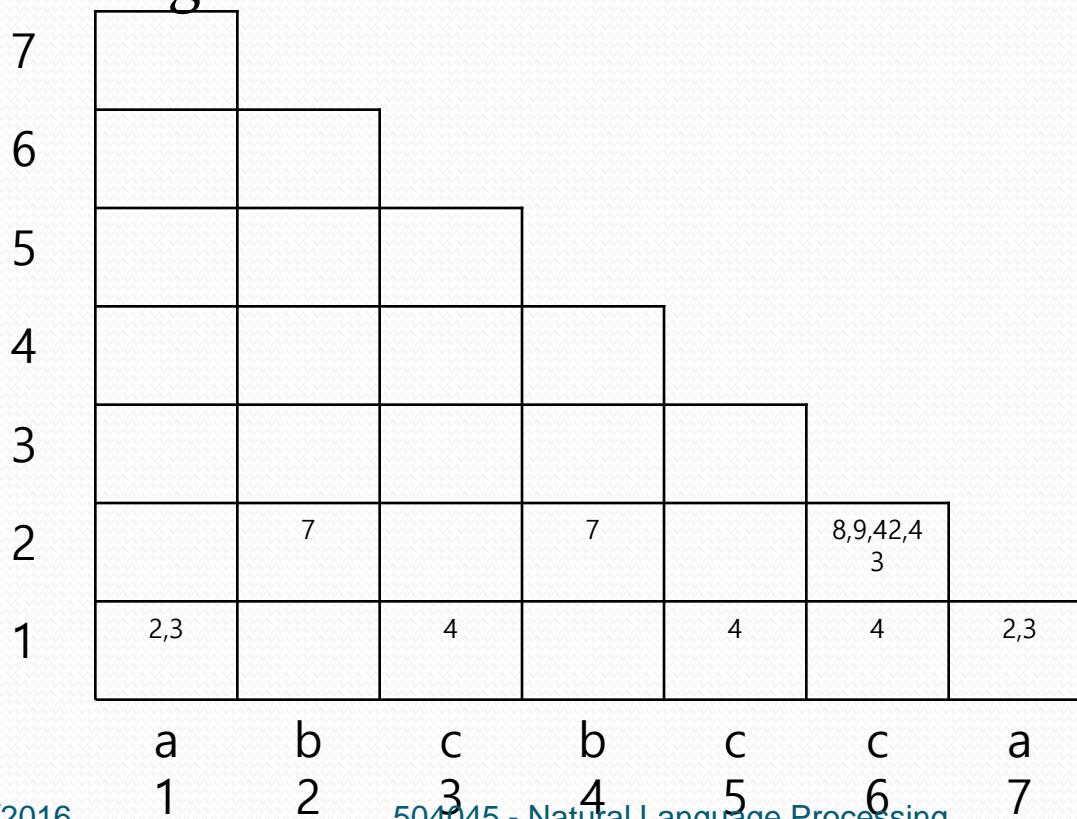


Example grammar

R0:S $\rightarrow$ (X,X) :1.0  
R1:S $\rightarrow$ (SX,SX) :0.7  
R2:X $\rightarrow$ (a,  $\neg$ ) : 0.4  
R3:X $\rightarrow$ (a,  $\sqsubset$ ) : 0.3  
R4:X $\rightarrow$ (c,  $\sqcup$ ) :0.7  
R5:X $\rightarrow$ (abXb,X  $\neg$   
 $\sqsubset$ ):0.1  
R6:X $\rightarrow$ (bcX,  $\neg$ X  $\sqsubset$ )  
:0.2  
R7:X $\rightarrow$ (bc  $\neg$ ) :0.3

# Hiero Decoding

- Example parsing
  - length 2

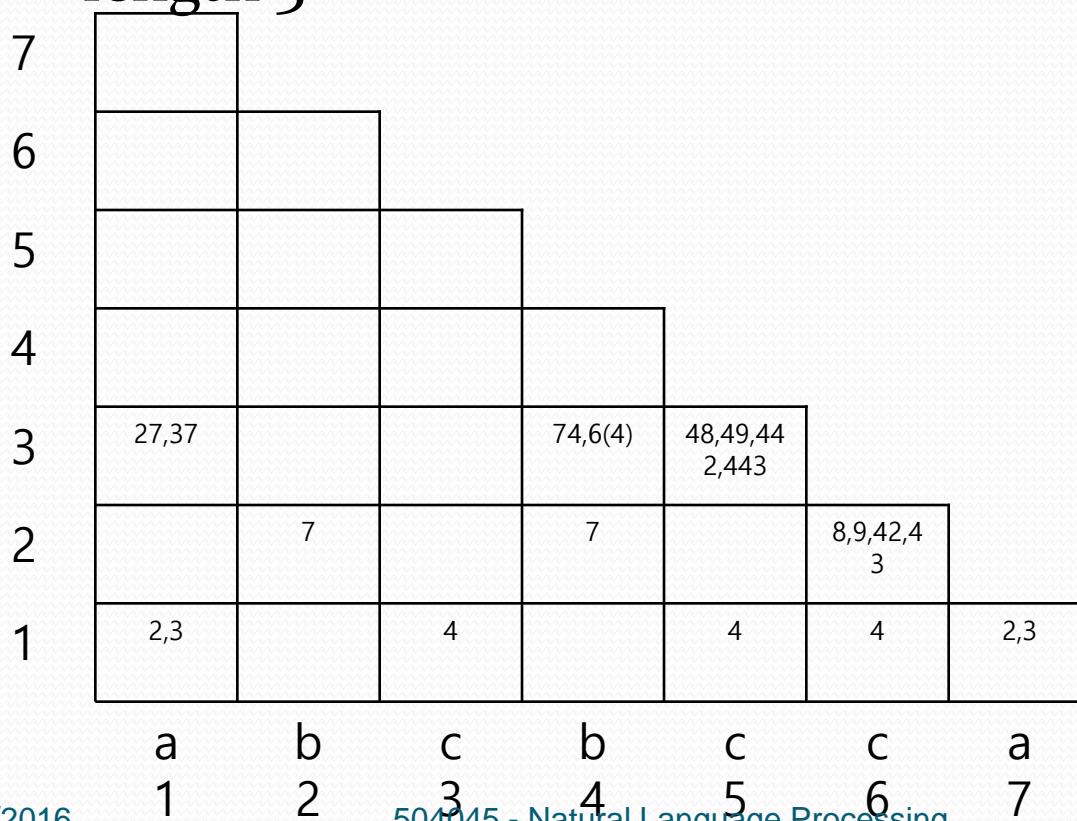


Example grammar

R0:S → (X,X) :1.0  
R1:S → (SX,SX) :0.7  
R2:X → (a, ⊁) : 0.4  
R3:X → (a, ⊚) : 0.3  
R4:X → (c, ⊂) :0.7  
R5:X → (abXb,X ⊁  
⊚ ):0.1  
R6:X → (bcX, ⊁ X ⊚ )  
:0.2  
R7:X → (bc ⊁ )·0.3

# Hiero Decoding

- Example parsing
  - length 3



Example grammar

R0:S $\rightarrow$ (X,X) :1.0  
R1:S $\rightarrow$ (SX,SX) :0.7  
R2:X $\rightarrow$ (a,  $\sqcap$ ) : 0.4  
R3:X $\rightarrow$ (a,  $\sqsubset$ ) : 0.3  
R4:X $\rightarrow$ (c,  $\sqcup$ ) :0.7  
R5:X $\rightarrow$ (abXb,X $\sqcap$   
 $\sqsubset$ ):0.1  
R6:X $\rightarrow$ (bcX,  $\sqcap$ X $\sqsubset$ )  
:0.2  
R7:X $\rightarrow$ (bc  $\sqcap$ ):0.3

# Hiero Decoding

- Example parsing

- length 4

7						
6						
5						
5(4)	6(7),77		6(8,9,42, 43),...			
27,37			74,6(4)	48,49,44 2,443		
	7		7		8,9,42,4 3	
2,3		4		4	4	2,3
a	b	c	b	c	c	a
1	2	3	4	5	6	7

6(8,9,42,43),78,  
79,742,743,742,  
6(4)2,743,6(4)3

Example grammar

R0:S → (X,X) : 1.0  
R1:S → (SX,SX) : 0.7  
R2:X → (a, ⊁) : 0.4  
R3:X → (a, ⊂) : 0.3  
R4:X → (c, ⊂) : 0.7  
R5:X → (abXb,X ⊁  
⊂):0.1  
R6:X → (bcX, ⊁X ⊂)  
:0.2  
R7:X → (bc ⊁):0.3

# Hiero Decoding

- Example parsing

- End

7	276(8,9, 42,43)...						
6	26(6(8,9, 42,43)...) 76(8,9,4 2,43)...						
5	26(7),....	6(6(8,9,4 2,43) ....	46(8,9,4 2,43)....				
4	5(4)	6(7),77		6(8,9,42, 43),...			
3	27,37			74,6(4)	48,49,44 2,443		
2		7		7		8,9,42,4 3	
1	2,3		4		4	4	2,3
	a	b	c	b	c	c	a
	1	2	3	4	5	6	7

Example grammar

R0:S → (X,X) : 1.0  
R1:S → (SX,SX) : 0.7  
R2:X → (a, ⊁) : 0.4  
R3:X → (a, ⊂) : 0.3  
R4:X → (c, ⊂) : 0.7  
R5:X → (abXb,X ⊁  
⊂):0.1  
R6:X → (bcX, ⊁X ⊂)  
:0.2  
R7:X → (bc ⊁):0.3

# Hiero Decoding

- Read Out the first cell
  - 276(8,9,42,43), ....
  - 276(8)
    - string: ㄱ ㄱ ㄱ ㄷ ㄱ ㄷ
    - score:  $0.4^*0.3^*0.2^*0.4^*$ LM Score
  - 276(9)
    - string: ㄱ ㄱ ㄱ ㄷ ㄱ ㄴ ㄷ
  - 276(42)
    - string: ㄱ ㄱ ㄱ ㄴ ㄱ ㄷ
  - 276(43)
    - string: ㄱ ㄱ ㄱ ㄴ ㄷ ㄷ
- Other derivations

Example grammar

R0:S → (X,X) :1.0  
R1:S → (SX,SX) :0.7  
R2:X → (a, ㄱ) : 0.4  
R3:X → (a, ㄷ) : 0.3  
R4:X → (c, ㄴ) :0.7  
R5:X → (abXb,X ㄷ  
ㄷ):0.1  
R6:X → (bcX, ㄱ X ㄷ)  
:0.2  
R7:X → (bc ㄱ ) :0.3

# Contents

- What Samsung requires to survey
  - Rule-base and Statistical approach
  - Alignment Model
  - Decoding Algorithms
  - Open Sources
  - Evaluation Methods
  - Using Syntax Info. in SMT

# Open Sources

- GIZA++
  - Franz Josef Och, 2000
  - Most SMT researchers use GIZA++
  - Much research on alignment start from IBM Model and HMM alignment model
  - A C++ Implementation of
    - IBM model 1~5
    - HMM alignment model
    - Smoothing for fertility, distortion/alignment parameters
    - Some improvements of IBM and HMM models
  - License : **GPL**
  - <http://www.fjoch.com/GIZA++.html>

# Open Sources

- Sri-LM
  - A. Stolcke, 2002
  - Implements State-of-the-art LM techniques
    - The latest update 1.5.3 (2007.07)
  - A C++ Implementation of
    - N-gram language modeling
    - Kneser-Ney discounting
    - Witten-Bell discounting
    - ....
  - License: **SRILM Research Community License**
  - <http://www.speech.sri.com/projects/srilm/>

# Open Sources

- Moses
  - Philipp Koehn et. al. 2007
  - State-of-the art SMT system
  - C++ & Perl implementation of
    - Phrase-based SMT ( Pharaoh )
    - Factor phrase-based decoder
    - Minimum error rate training
    - Translation Model training
  - License : **LGPL**
  - <http://www.statmt.org/moses/>

# Open Sources

- GenPar Toolkit
  - C++ implementation
  - Translation Performance is not so good
  - SMT by parsing
    - Training by parallel parsing on target and source language
    - Decoding by CKY style parsing algorithm
  - License : **GPL 2.0 or later**
  - <http://nlp.cs.nyu.edu/GenPar/>

# Open Sources

- Phramer
  - Marian Olteanu, 2006
  - Java implementation of
    - Phrase-based machine translation
    - MERT training of MT
  - License : **Free** (Copyright (c) 2006-2007, Marian Olteanu All rights reserved. )
  - <http://www.utdallas.edu/~mgo031000/phramer/>

# Contents

- Part1: What Samsung requires to survey
  - Rule-base and Statistical approach
  - Alignment Model
  - Decoding Algorithms
  - Open Sources
  - Evaluation Methods
  - Using Syntax Info. in SMT

# Automatic Evaluation

- Advantages of automatic evaluation
  - Fast, Low Cost
  - Objective
- Evaluation methods
  - BLEU Score: Bi-Lingual Evaluation Understudy Score
    - Geometric mean of modified n-gram precision
  - NIST Score:
    - Arithmetic mean of modified n-gram precision
  - METEOR Score: Metric for Evaluation of Translation With Explicit Ordering
  - WER :Word Error Rate
  - PER :Position independent word Error Rate
  - TER :Translation Error Rate
  - Others ..

# Automatic Evaluation: examples

- BLEU score
  - Most famous metric
  - Range 0~1.
  - Higher score means better translation
  - Typically, consider up to 4-gram
    - denote BLEU-4 score

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \leq r \end{cases}$$

c : length of candidate translation

$$\text{BLEU} = BP \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right)$$

BP: factor related to the length of candidate translation

$p_n$ : n-gram precision, ignoring duplicate count

# Automatic Evaluation: examples

- NIST score
  - Similar to BLEU metric
  - Higher score means better translation
  - Arithmetic mean of n-gram precisions
  - $\beta$  is chosen to make the BP =0.5 when  $L_{sys}/L_{ref} = 2/3$

$$BP = \exp(\beta \log_2 \min(\frac{L_{sys}}{L_{ref}}, 1))$$

$L_{sys}$  : length of candidate translation

$$\text{Info}(w_1 \dots w_n) = \log_2 \left( \frac{\text{the \# of occurrences of } w_1 \dots w_n \text{ in reference sentences}}{\text{the \# of occurrences of } w_1 \dots w_n} \right)$$

$$\text{NIST} = BP \cdot \sum_{n=1}^N \frac{\sum_{\substack{\text{co-occurred words} \\ \text{words in output}}} \text{Info}(w_1 \dots w_n)}{\sum_{\substack{1 \\ \text{words in output}}}}$$

# Automatic Evaluation: examples

- METEOR Score
  - Not very popular
  - Based on uni-gram precision and recall
  - Chunks : A sequence of uni-grams those are adjacent in both reference and system output

$$F\text{-mean} = \frac{10 \text{ Precision}}{\text{Recall} + 9 \text{ Precision}}$$

$$\text{Penalty} = 0.5 \left( \frac{\# \text{ of chunks}}{\# \text{ of unigrams matched}} \right)^3$$

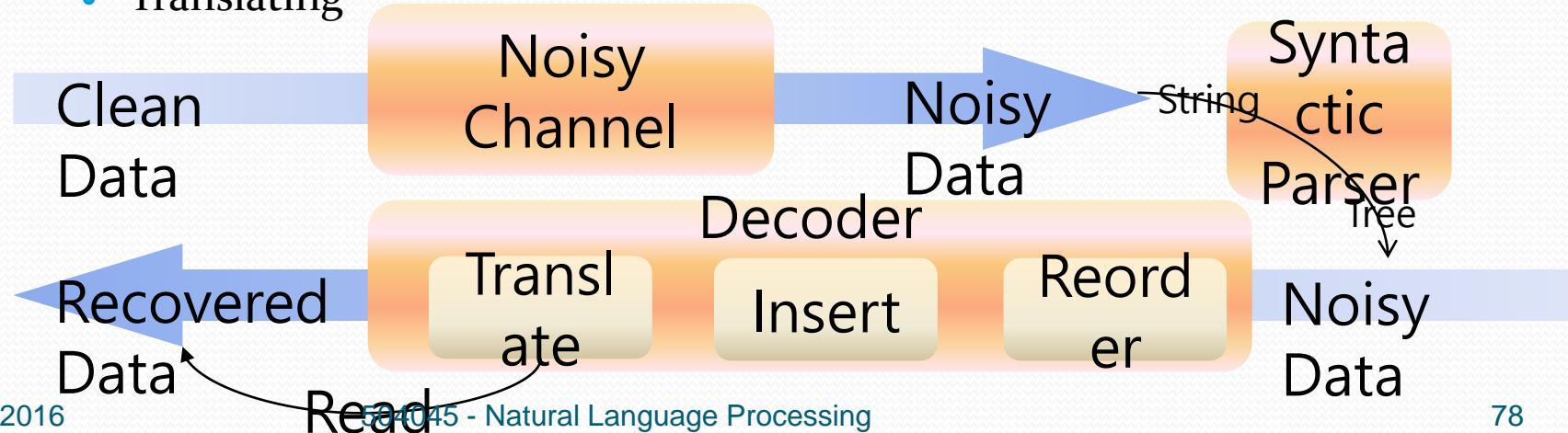
$$\text{Meteor} = F\text{-mean} (1 - \text{Penalty})$$

# Contents

- What Samsung requires to survey
  - Rule-base and Statistical approach
  - Alignment Model
  - Decoding Algorithms
  - Open Sources
  - Evaluation Methods
  - **Using Syntax Info. in SMT**

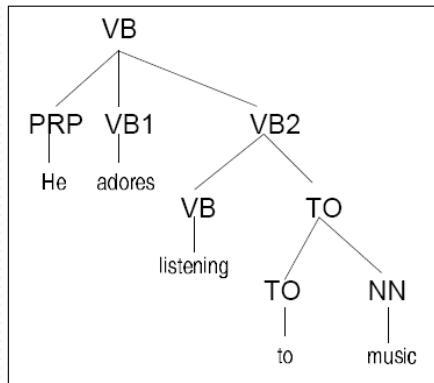
# Syntax-based Statistical Translation

- K. Yamada and K. Knight [2001] proposed a method
- Modified source-channel model
  - Input
    - Sentences → Parse trees
    - Input sentences are preprocessed by a syntactic parser
  - Channel operation
    - Reordering
    - Inserting
    - Translating



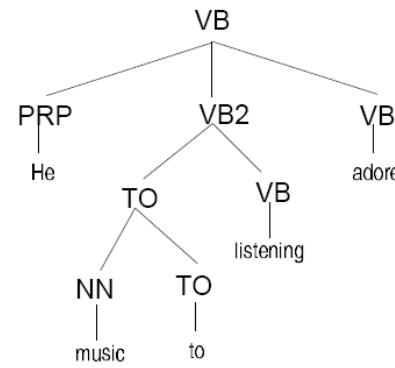
# Syntax-based MT: Process

- Original sentence is processed by a syntactic parser
- Start with a parse tree
- Step1: Reordering
  - Assume that only the sequence of child node labels influences the reordering
  - The probability of reordering is given by r-table (reordering model)
    - e.g.  $0.723(\text{PRP VB1 VB2} \rightarrow \text{PRP VB2 VB1}) * 0.749(\text{VB TO} \rightarrow \text{TO VB}) * 0.893(\text{TO NN} \rightarrow \text{NN TO}) = 0.484$



1. Channel Input

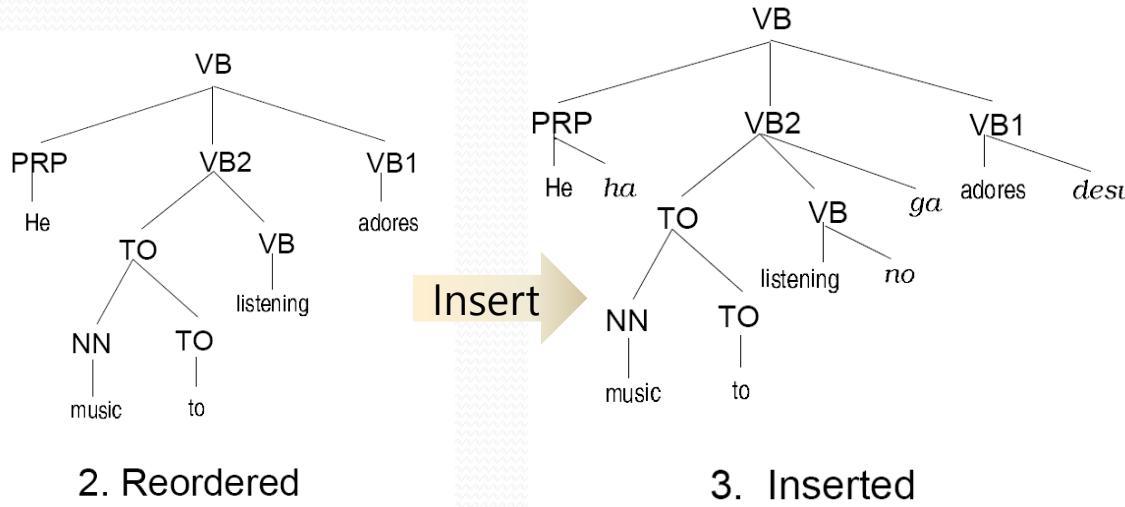
Reorder  
→



2. Reordered

Original order	Reordering	P(reorder)
PRP VB1 VB2	PRP VB1 VB2 <b>PRP VB2 VB1</b>	0.074 0.723
	VB1 PRP VB2	0.061
	VB1 VB2 PRP	0.037
	VB2 PRP VB1	0.083
	VB2 VB1 PRP	0.021
VB TO	VB TO <b>TO VB</b>	0.251 0.749
TO NN	TO NN <b>NN TO</b>	0.107 0.893
...	...	...

# Syntax-based MT: Process



Calculation  
of this  
example

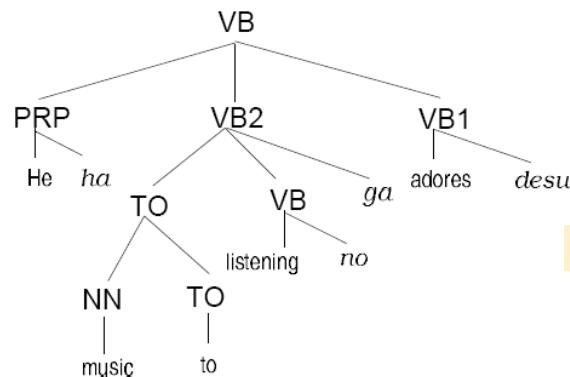
$$(0.652 * 0.219) * (0.252 * 0.094) \\ (0.252 * 0.62) * (0.252 * 0.0007) * (0.735 * 0.709) * (0.900 * 0.800) = 3.498e-9$$

- Insertion probability is defined by n-table
- n-table is divided into two
  - Table for position
  - Table for word identity

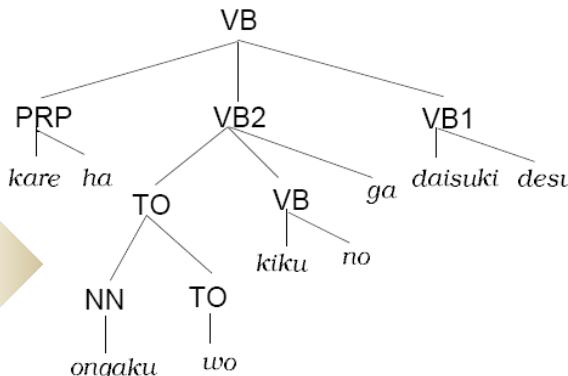
Parent	TOP	VB	VB	VB	TO	TO	...
Node	VB	VB	PRP	TO	TO	NN	...
P(None)	0.735	0.687	0.344	0.709	0.900	0.800	...
P(Left)	0.004	0.061	0.004	0.030	0.003	0.096	...
P(Right)	0.260	0.252	0.652	0.261	0.007	0.104	...

W	P(ins-w)
ha	0.219
ta	0.131
wo	0.099
no	0.094
ni	0.080
te	0.078
ga	0.062
...	...
desu	0.0007
...	...

# Syntax-based MT: Process



3. Inserted



4. Translated

Calculation  
of this  
example  
 $0.952 * 0.900 * 0.0038 * 0.333 * 1.000 = 0.0108$

- The translate operation is applied to each leaf
- This operation is dependent only on the word itself

E	adores	he	I	listening	Music	to	...
J	daisuki	1.000	kare NULL nani da shi ...	0.95 0.016 0.005 0.003 0.003 ...	NUL L watas i kare shi nani ....	0.471 0.111 0.055 0.021 0.020 0.333 0.333 0.333 0.900 0.100 ni NULL to no wo ....	0.216 0.204 0.133 0.046 0.038 ...

These figures came from the original paper [Yamada and Knight, "A Syntax Based Translation Model", 2001]

# Hierarchical Modeling

- Hierarchical organization of Natural Language
  - A sentence is derived by recursive application of some production rules
  - S yields NP and VP
  - VP may yield another NP
  - ...
- Traditional Statistical Systems
  - A sentence is generated by sequentially concatenating some phrases
  - We need to model the Hierarchical property of language

# Hiero

- Hiero
  - A Hierarchical Phrase based Statistical Machine Translation System
  - Automatically extracts production rules from un-annotated parallel texts
  - Finds the best derivation for a given sentence using Modified CKY beam search decoder
  - Grammar
    - Form of synchronous CFG
    - Can be Automatically extracted from parallel texts
  - Model
    - Use log-linear model
    - Assign a weight for each rule
    - Goal is finding total weight
$$w(X \rightarrow \langle \gamma, \alpha \rangle) = \prod_i \phi_i(X \rightarrow \langle \gamma, \alpha \rangle)^{\lambda_i}$$

# Synchronous Grammar

- A synchronous CFG
  - Consists of a pair of CFG rules with aligned non-terminal symbols
  - Derivation starts with a pair of start symbols
- A Partial Derivation

$\langle S_{\text{I}}, S_{\text{I}} \rangle \xrightarrow{1} \langle S_{\text{2}} X_{\text{3}}, S_{\text{2}} X_{\text{3}} \rangle$

$\xrightarrow{1} \langle S_{\text{4}} X_{\text{5}} X_{\text{3}}, S_{\text{4}} X_{\text{5}} X_{\text{3}} \rangle$

$\xrightarrow{2} \langle X_{\text{6}} X_{\text{5}} X_{\text{3}}, X_{\text{6}} X_{\text{5}} X_{\text{3}} \rangle$

$\xrightarrow{6} \langle \text{Aozhou } X_{\text{5}} X_{\text{3}}, \text{Australia } X_{\text{5}} X_{\text{3}} \rangle$

$\xrightarrow{7} \langle \text{Aozhou shi } X_{\text{3}}, \text{Australia is } X_{\text{3}} \rangle$

$\xrightarrow{5} \langle \text{Aozhou shi } X_{\text{7}} \text{ zhiyi, Australia is one of } X_{\text{7}} \rangle$

$\xrightarrow{4} \langle \text{Aozhou shi } X_{\text{8}} \text{ de } X_{\text{9}} \text{ zhiyi, Australia is one of the } X_{\text{9}} \text{ that } X_{\text{8}} \rangle$

$\xrightarrow{3} \langle \text{Aozhou shi } yu X_{\text{1}} \text{ you } X_{\text{2}} \text{ de } X_{\text{9}} \text{ zhiyi, Australia is one of the } X_{\text{9}} \text{ that have } X_{\text{2}} \text{ with } X_{\text{1}} \rangle$

## Grammar example

- |    |   |
|----|---|
| 1  | $S \rightarrow \langle S_{\text{I}} X_{\text{2}}, S_{\text{I}} X_{\text{2}} \rangle$  |
| 2  | $S \rightarrow \langle X_{\text{I}}, X_{\text{II}} \rangle$   |
| 3  | $X \rightarrow \langle \text{yu } X_{\text{II}} \text{ you } X_{\text{2}}, \text{have } X_{\text{2}} \text{ with } X_{\text{II}} \rangle$ |
| 4  | $X \rightarrow \langle X_{\text{II}} \text{ de } X_{\text{9}}, \text{the } X_{\text{2}} \text{ that } X_{\text{II}} \rangle$              |
| 5  | $X \rightarrow \langle X_{\text{II}} \text{ zhiyi, one of } X_{\text{II}} \rangle$  |
| 6  | $X \rightarrow \langle \text{Aozhou, Australia} \rangle$  |
| 7  | $X \rightarrow \langle \text{shi, is} \rangle$  |
| 8  | $X \rightarrow \langle \text{shaoshu guojia, few countries} \rangle$  |
| 9  | $X \rightarrow \langle \text{bangjiao, diplomatic relations} \rangle$   |
| 10 | $X \rightarrow \langle \text{Bei Han, North Korea} \rangle$   |