

①

## Statistical Machine Translation (SMT)

In SMT we aim to translate the source language sentence into a target language sentence. We attempt to solve this problem as a conditional probability of  $T$  given  $S$ .  $P(T|S)$ .

Now the MT system has to figure out that  $\hat{T}$  that maximizes this conditional probability.

$$\hat{T} = \underset{T}{\operatorname{argmax}} (T|S)$$

$$\text{Using bayes rule } P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$\hat{T} = \underset{T}{\operatorname{argmax}} \frac{P(S|T)P(T)}{P(S)}$$

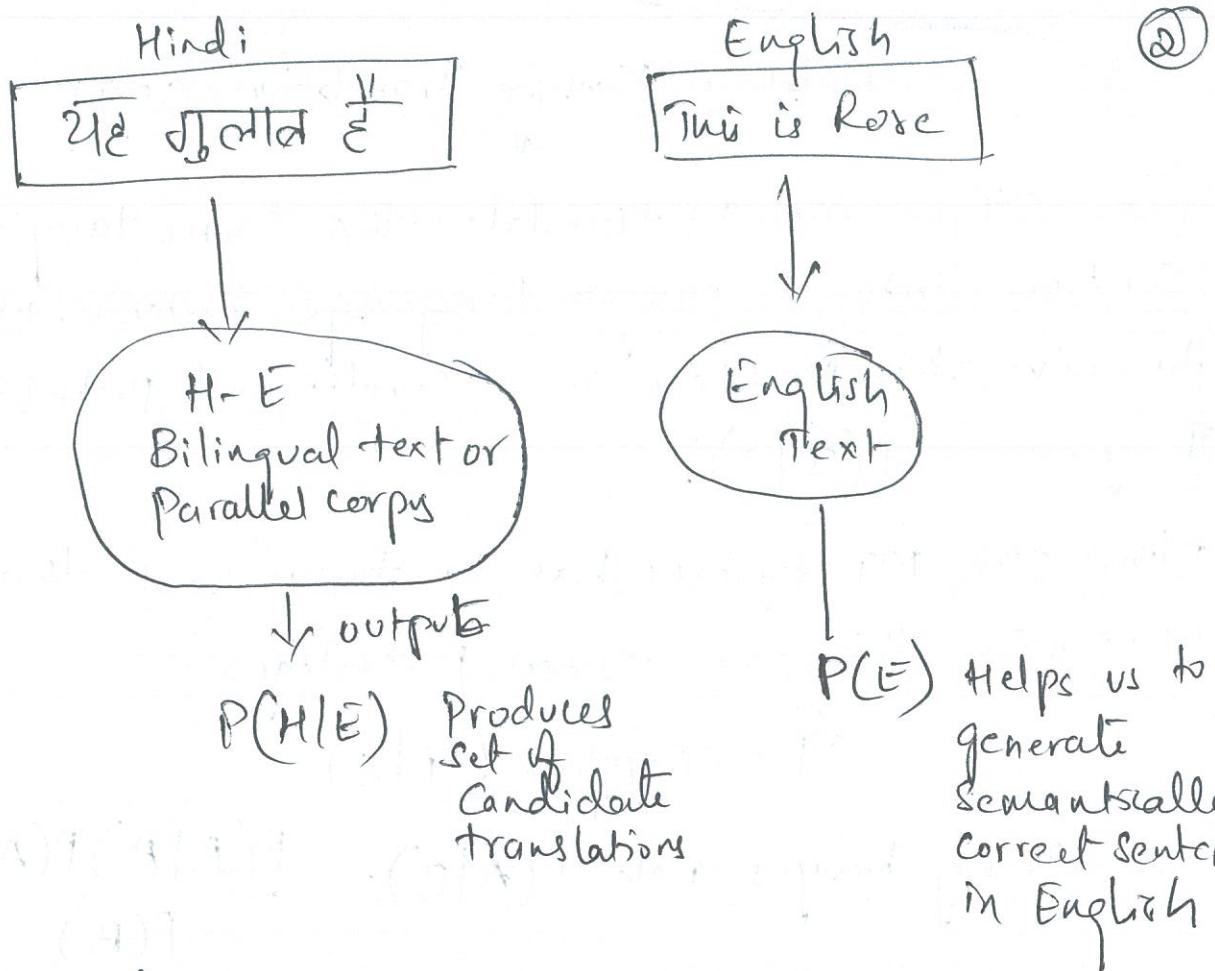
We decompose this problem into 2 subproblems and since  $P(S)$  in the denominator is not a function of  $T$  and is a constant which we multiply to all  $T$ 's we ignore it.

$$\hat{T} = \underset{T}{\operatorname{argmax}} P(S|T)P(T)$$

$P(S|T)$  is called a translation model & used to find the probability that  $S$  comes from  $T$ .

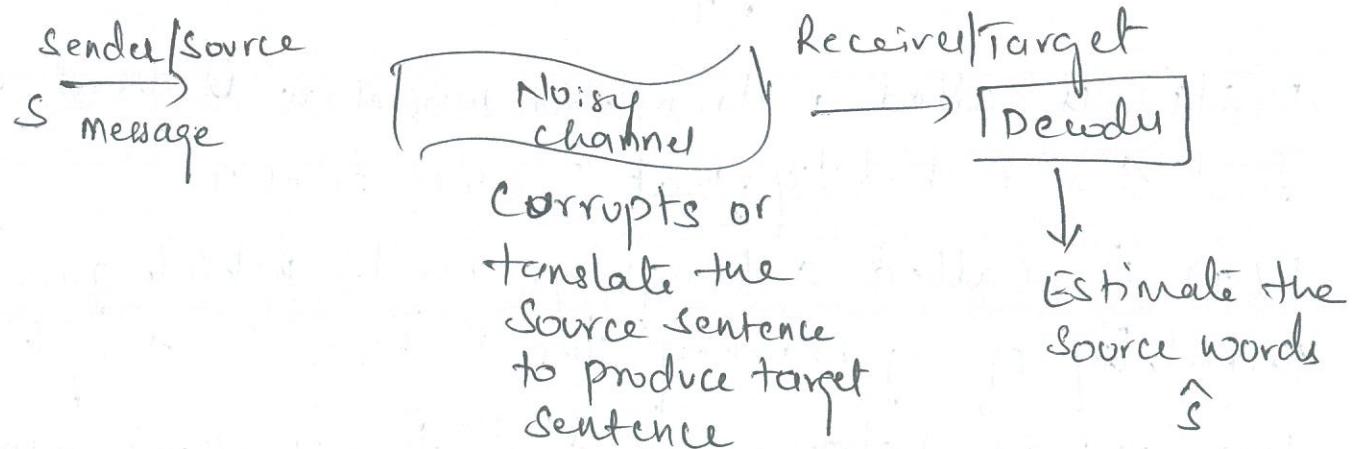
$P(T)$  is called a language model which gives us a probability of generating the sentence in target language.

Decoder: It helps us in find the  $\hat{T}$  that maximizes the product of both the above subproblems.



$$\hat{E} = \underset{E}{\operatorname{argmax}} P(H|E) P(E)$$

But we have an interpretation problem in this model. Our main goal was to translate  $P(T|S)$  and according to bayes rule we are estimating  $P(S|T)$ . The noisy channel model comes handy to say that instead of  $P(T|S)$  it is similar to  $P(S|T)$



At the receiver end we need to guess the original message which was in the source language.

(3)

Now our decoding algorithm will iterate over a set of possible source sentences which could be the original message (argmax). The sentence that we picked up in source language will have a probability of being in the target language  $P(T)$ . Hence when we multiply TM & LM we get the probability of this target to be the original sentence of this source sentence's.

Consider that we are translating from English to Hindi originally. Now due to noisy channel model our TM will be  $P(E|H)$ .

O/P sentences	$\sum_E P(E)$	$\sum_H P(H E)$	$\sum_{E,H} P(H E)P(E)$
This nose is low	low	high	low
is this nose	low	high	low
is nose this	low	high	low
This pink is	low	high	low
pink is good	low	low	low
This is nose	high	high	high

The argmax is a product of TM & LM and then rank all candidate target sentences based on this product and pick the one that has highest probability.

(4)

Language Model (LM): Predict the probability of generating the sentence in a given language. In our SMT we estimate  $P(T)$ .

A Sentence generation can be modelled as a sequence of words using chain rule.

$$P(s) = P(w_1, \dots, w_n) \\ = P(w_1) P(w_2 | w_1) P(w_3 | w_2, w_1) \dots P(w_n | w_{n-1}, w_{n-2})$$

↑  
This is to  
lengthy

To get rid of this lengthy subsequence we apply Markov assumption

$$P(w_i | w_{i-1}, w_{i-2}, \dots, w_1) = P(w_i | w_{i-1}, \dots, w_{i-n+1})$$

In Unigram model we assume each word is independent of each other. Hence if we are to generate a sentence of 4 words it is modelled as

$$P(w_1, w_2, w_3, w_4) = P(w_1) \times P(w_2) \times P(w_3) \times P(w_4)$$

The drawback of this model is that the word order does not matter and hence the sentences may not be semantically correct.

(3)  
Bigram model

$$P(w_1, w_2, w_3, w_4) = P(w_1 | \text{begin of sentence}) P(w_2 | w_1) P(w_3 | w_2) P(w_4 | w_3)$$

The probability of each word depends only on the previous word.

Trigram model

$$P(w_1, w_2, w_3, w_4) = P(w_1 | \text{begin of sentence}) P(w_2 | w_1) P(w_3 | w_2, w_1) \\ P(w_4 | w_3, w_2)$$

The longer the conditional prob window the better syntactically correct sentences will be generated. But too longer subsequences can be generated with nearly zero probability.

How to compute these parameters  $P(w_3 | w_2, w_1)$ ?

$$P(w_3 | w_2, w_1) = \frac{\text{Count}(w_3, w_2, w_1)}{\text{Count}(w_2, w_1)}$$

All the parameters can be computed in this form for the language model from the target language.

Now let's see how to model Translation Model

How to assign values for  $P(S|T)$ ?

$$P(S|T) = \frac{\text{Count}(S, T)}{\text{Count}(T)}$$

These estimates needs to be done at sentence level, but this is not possible since <sup>the</sup> number of such exact sentences may be very low. Hence we model this using a hidden variable  $\alpha$  (alignment) that represent alignment between the individual words in a sentence as marginalization over all alignments.

Now let's see how to model Translation?

How do we estimate values of  $P(S|T)$ ?

$$P(S|T) = \frac{\text{count}(S, T)}{\text{count}(T)}$$

What we are looking for is word to word mapping to estimate these counts. If we have word aligned text we could easily estimate  $P(S|T)$ . Unfortunately we have the sentences aligned in the parallel corpus.

Using these sentences we model a hidden variable  $a$  (alignment) that represent how individual words in a sentence can be mapped.

$P(S|T)$  involves two choices

1. choice of length  $l_s$  for the source sentence (constraint)
2. choice of words  $s_1, s_2, s_3 \dots s_{l_s}$

From now on we will take the length  $l_s$  to be fixed and we try to model

$$P(s_1, s_2, s_3 \dots s_{l_s} | t_1, t_2, t_3 \dots t_{l_t}, l_s)$$

It is difficult to estimate this directly. Hence we introduce alignments variable. An alignment variable  $a$  identifies which target word each source word originated from.

Formally an alignment variable can takes values from  $\{0, 1, \dots, l_t\}$  and its length is  $l_s$ .

For Example if the translation model is  $P(H|E)$

1 2 3

(T) E = This is rose

(S) H = ये गुलाब हैं

In this case the length of Hindi sentence is 3

$$\therefore a = \langle 1 3 2 \rangle$$

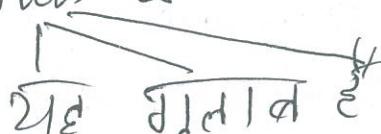
(T) E = Ram hit the boy with stick

(S) H = राम ने बच्चे को छला

In this example length of source sentence (H) is 7  
 $\therefore a = \langle 1 1 6 5 4 3 2 \rangle$

Note that each source word is aligned to exactly one target word. This alignment is many to one i.e more than one source word can be aligned to a single target word.

Another possible alignment could be

$a = \langle 1, 1, 1 \rangle$       This is rose  


This alignment is bad

(8)

Given the length of source and target sentences there would be  $(l_t + 1)^{l_s}$  possible alignments.

In IBM Model-1 the assumption is that all alignments ~~are~~ of  $T$  are equally likely given the length of source.

$$\forall a, P(a | T, l_s) = \frac{1}{(l_t + 1)^{l_s}}$$

It means that

$$P\left(\begin{array}{c} \text{This is rose} \\ \text{The dog eat } \end{array}\right) = P\left(\begin{array}{c} \text{This is rose} \\ \cancel{\text{The }} \cancel{\text{dog }} \text{ eat } \end{array}\right)$$

Given an alignment  $a$  and the target sentence  $T$ , what is the probability of a source sentence  $S$ ?

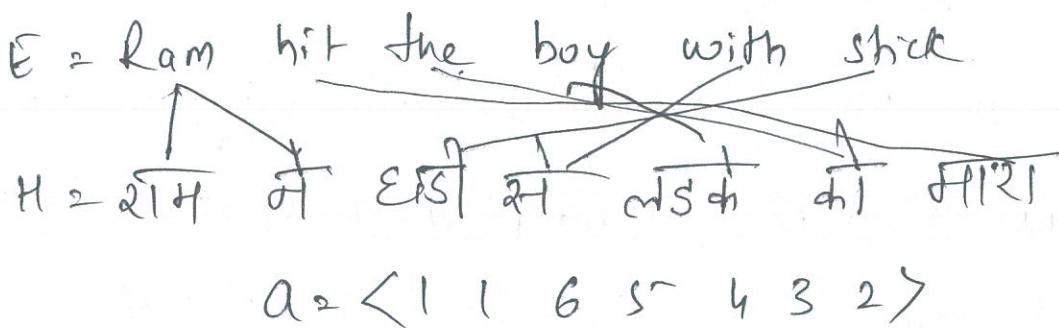
$$P(S|a, T)$$

In IBM Model-1

$$P(S|a, T) = \prod_{j=1}^{l_s} P(S_j | T_{a_j})$$

$P(S_j | T_{a_j})$  denotes the probability of  $j^{\text{th}}$  source word given that it was generated from  $a_j^{\text{th}}$  target word.

(9)



$$P(H|a, E) = P(\overline{\text{राम}} | \text{Ram}) P(\overline{\text{ने}} | \text{hit}) P(\overline{\text{ईस्ट}} | \text{stick}) \\ P(\overline{\text{लड़का}} | \text{boy}) P(\overline{\text{को}} | \text{with}) P(\overline{\text{मारा}} | \text{stick})$$

### DBM-Model 1

To generate a source sentence  $S$  from target  $T$

1. Pick a length of  $S$
2. Pick an alignment with uniform probability

$$\frac{1}{(L_t+1)^{L_s}}$$

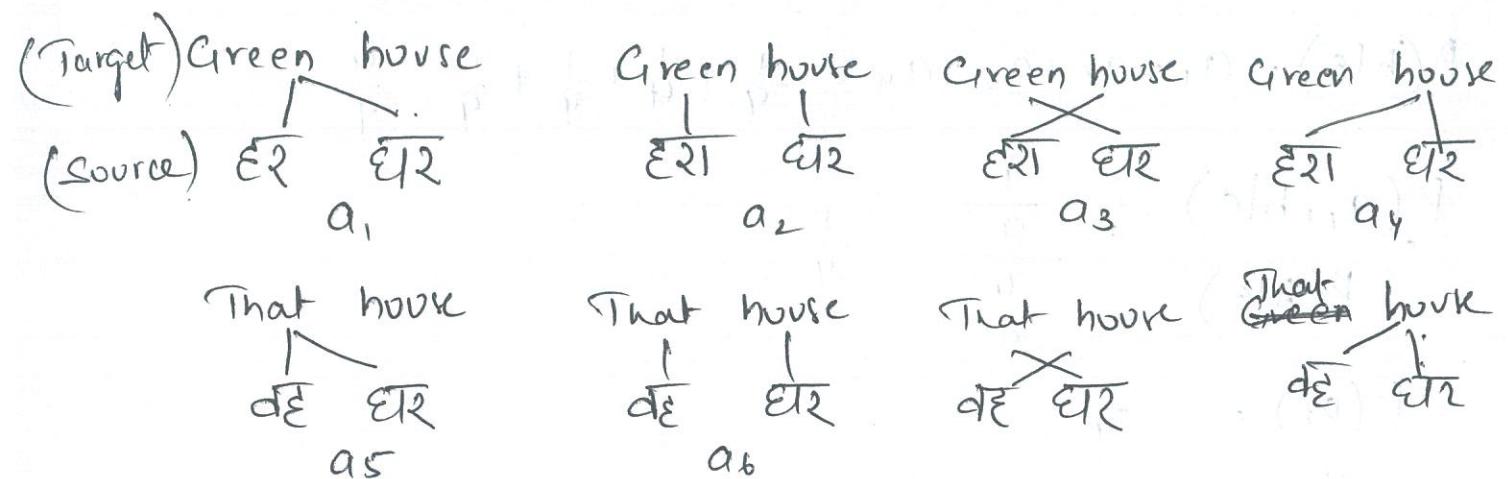
3. Sample source word with probability

$$P(S|a, T) = \prod_{j=1}^{L_s} P(s_j | T_{a_j})$$

$$\therefore P(s_j | a | T) = P(a | T) P(s_j | a, T)$$

$$= \frac{1}{(L_t+1)^{L_s}} \prod_{j=1}^{L_s} P(s_j | T_{a_j})$$

In order to estimate the translation model parameters we use the Expectation maximization algorithm. Here is a worked out example with 2 sample sentences.



Initialize all parameters  $t(h|e)$  uniformly

	$\widehat{E21}$	$\widehat{E12}$	$\widehat{d3}$
Green	$1/3$	$1/3$	$1/3$
house	$1/3$	$1/3$	$1/3$
That	$1/3$	$1/3$	$1/3$

Compute  $P(a_i, h|e)$  for each alignment  $a_1 - a_8$

$$P(a_1, h|e) = P(\widehat{E21} | \text{Green}) P(\widehat{E12} | \text{Green}) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$P(a_2, h|e) = P(\widehat{E21} | \text{Green}) P(\widehat{E12} | \text{house}) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$P(a_8, h|e) = P(\widehat{d3} | \text{house}) P(\widehat{E12} | \text{house}) = \frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

(11)

Normalize the alignment values

$$\text{i.e. } P(a|h,e) = \frac{P(a,h|e)}{P(h|e)}$$

$$P(h|e) = a_1 + a_2 + a_3 + a_4 = \frac{1}{9} + \frac{1}{9} + \frac{1}{9} + \frac{1}{9} = \frac{4}{9}$$

$$\frac{P(a_1, h|e)}{P(h|e)} = \frac{\frac{1}{9}}{\frac{4}{9}} = \frac{1}{4}$$

$$P(a_2) = \frac{1}{4}$$

;

;

$$P(a_8) = \frac{1}{4}$$

Compute fractional counts

$$c(\widehat{e1}, \text{Green}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (\overset{\text{using}}{a_1, a_2})$$

$$c(\widehat{e2}, \text{Green}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (a_1, a_3)$$

$$c(\widehat{a1}, \text{Green}) = 0 \quad (\text{These words never occurred many alignment})$$

$$c(\widehat{e2}, \text{house}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (a_3, a_4)$$

$$c(\widehat{e1}, \text{house}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (a_2, a_4)$$

$$c(\widehat{a1}, \text{house}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (a_7, a_8)$$

$$c(\widehat{e1}, \text{that}) = 0$$

$$c(\widehat{e2}, \text{that}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (a_5, a_7)$$

$$c(\widehat{a1}, \text{that}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad (a_5, a_6)$$

(12)

$$\text{total count (Green)} = \frac{1}{2} + \frac{1}{2} + 0 = 1$$

$$(\text{house}) = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2}$$

$$(\text{that}) = 0 + \frac{1}{2} + \frac{1}{2} = 1$$

Maximization step

$$t(\overline{e_21} | \text{Green}) = \frac{\frac{1}{2}}{\frac{1}{2}} = \frac{1}{2}$$

$$t(\overline{e_12} | \text{Green}) = \frac{\frac{1}{2}}{\frac{1}{2}} = \frac{1}{2}$$

$$t(\overline{d_2} | \text{Green}) = 0$$

$$t(\overline{e_21} | \text{house}) = \frac{\frac{1}{2}}{\frac{3}{2}} = \frac{1}{3}$$

$$t(\overline{e_12} | \text{house}) = \frac{\frac{1}{2}}{\frac{3}{2}} = \frac{1}{3}$$

$$t(\overline{d_2} | \text{house}) = \frac{\frac{1}{2}}{\frac{3}{2}} = \frac{1}{3}$$

$$t(\overline{e_21} | \text{that}) = 0$$

$$t(\overline{e_12} | \text{that}) = \frac{\frac{1}{2}}{\frac{1}{2}} = \frac{1}{2}$$

$$t(\overline{d_2} | \text{that}) = \frac{\frac{1}{2}}{\frac{1}{2}} = \frac{1}{2}$$

keep iterating over the Expectation & Maximization steps till all the translation parameters converge.

## IBM Model-2

13

In Model-1 all alignment are equally likely which is wrong. Since languages having SOV (Subject Object Verb) Order like most Indian languages the verbs predominantly move to the End of the Sentence. Once the position of the noun is fixed adjectives qualify these nouns. Hence it is realistic to assume a probability distribution over position of words as

$$q_r(i|s, l_t, l_s)$$

$q_r$  can be interpreted as the probability that  $j^{th}$  Source word is connected to  $i^{th}$  Target word, given the sentence lengths of source & target sentence as  $l_t$  and  $l_s$  respectively.

$$P(s_1, s_2, s_3, \dots, s_{l_s}, a_1, a_2, \dots, a_{l_s} | t_1, t_2, \dots, t_{l_t}, l_s) =$$

$$\prod_{j=1}^{l_s} q_r(a_j | i, l_t, l_s) T(s_j | t_{a_j})$$

Target  $i$ : That house is beautiful  
 Source  $j$ :  $\cancel{d}\cancel{e}$   $\cancel{h}\cancel{e}\cancel{l}$   $\cancel{b}\cancel{e}\cancel{a}\cancel{u}\cancel{l}\cancel{i}\cancel{f}\cancel{l}\cancel{h}\cancel{l}$   $a = \langle 1\ 2\ 4\ 3 \rangle$

$$\begin{aligned}
 P(s|T) = & t(d|that) q_r(1|1, 4, 4) \times \\
 & t(e|house) q_r(2|2, 4, 4) \times \\
 & t(h|is) q_r(3|3, 4, 4) \times \\
 & t(l|beautiful) q_r(4|4, 4, 4)
 \end{aligned}$$

## IBM Model-3

(14)

The assumption in IBM model 1 & 2 is that each word in the Target generates one or many source words but their distribution is uniform whether it generates one or many words. This factor is handled as a fertility factor in Model 3. For each target word  $t$  in target sentence  $T$ , we choose a fertility  $\phi$ . The choice of  $\phi$  depends on  $t$  which is modeled as a probability distribution parameter  $n(\phi|t)$ . These parameters for each target word is computed from the training corpus.

We are given  $T = \{t_1, t_2, t_3, \dots, t_{l_T}\}$  want to model  $P(S|T)$ .

Since we have  $l_T$  no of target words we have to compute  $l_T + 1$  fertilities with probability

$$P(\{\phi_0, \phi_1, \phi_2, \dots, \phi_{l_T}\} | T)$$

where  $\phi_i$  stands for the no of source words  $t_i$  generates and  $n(\phi_i | t_i)$  is the probability that  $t_i$  generates  $\phi_i$  words.

For Example given the following Target Sentence

That house is beautiful.

we compute a probability distribution

$$n(0| \text{that}) = 0.1 \quad n(0| \text{house}) = 0.001 \quad n(0| \text{is}) = 0.1 \quad \dots$$

$$n(1| \text{that}) = 0.9 \quad n(1| \text{house}) = 0.9 \quad n(1| \text{is}) = 0.9$$

$$n(2| \text{that}) = 0 \quad n(2| \text{house}) = 0.09 \quad n(2| \text{is}) = 0$$

:

:

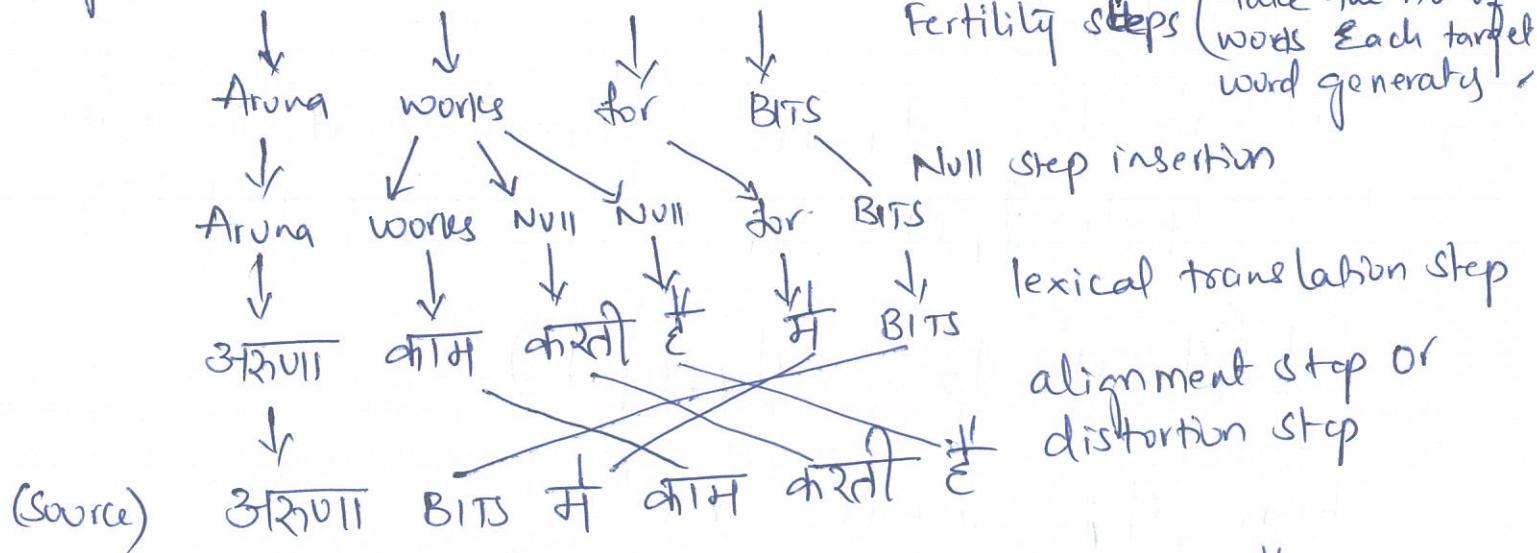
:

choose  $L_t+1$  fertilities  $\{\phi_0, \dots, \phi_{L_t}\}$  with probability

$$P(\{\phi_0, \phi_1, \phi_2, \dots, \phi_{L_t}\} | T) = P(\phi_0 | \phi_1, \phi_2, \dots, \phi_t) \underbrace{\prod_{i=1}^{L_t} n(\phi_i | n_i)}_{\text{Null word insertion}}$$

Let's see an example of the whole process of translation

(Target) Aruna works for BITS



If we want to add the source words ~~करती है~~, it can be generated by the special Null token.

We could model the fertility of the Null token in the same way as all other words by the conditional distribution  $n(\phi | \text{Null})$ . However the no of inserted words clearly depends on the sentence length, hence we model the Null word insertion as a special step.

$\therefore$  we represent a parameter  $P_0$  which is the probability of having no Null token inserted after every target word  $t$ .  $P_1$  is the probability of generating one Null token inserted after each target word  $t$ .

The number of source words generated from a Null token is represented as  $\phi_0$ .

$$\text{No of o/p words} = \sum_{i=1}^{L_t} \phi_i = L_s - \phi_0$$

The estimation of  $\phi_0$  is treated as a binomial distribution since each target word  $t$  may generate a null ( $P_1$ ) or not ( $P_0$ ).  $P(\phi_0)$  is the probability of generating  $\phi_0$  source words from a Null token is estimated as

$$P(\phi_0) = \binom{L_s - \phi_0}{\phi_0} P_1^{\phi_0} P_0^{L_s - 2\phi_0} \quad (\text{Using binomial distribution})$$

(17)

$$P(S|T, l_s) = \binom{l_s - \phi_0}{\phi_0} p_1^{\phi_0} p_0^{l_s - \phi_0} \prod_{i=1}^{l_t} \phi_i! n(\phi_i | t_i) \times \\ \prod_{j=1}^{l_s} q(a_j | j, l_t, l_s) t(s_j | t_{a_j})$$