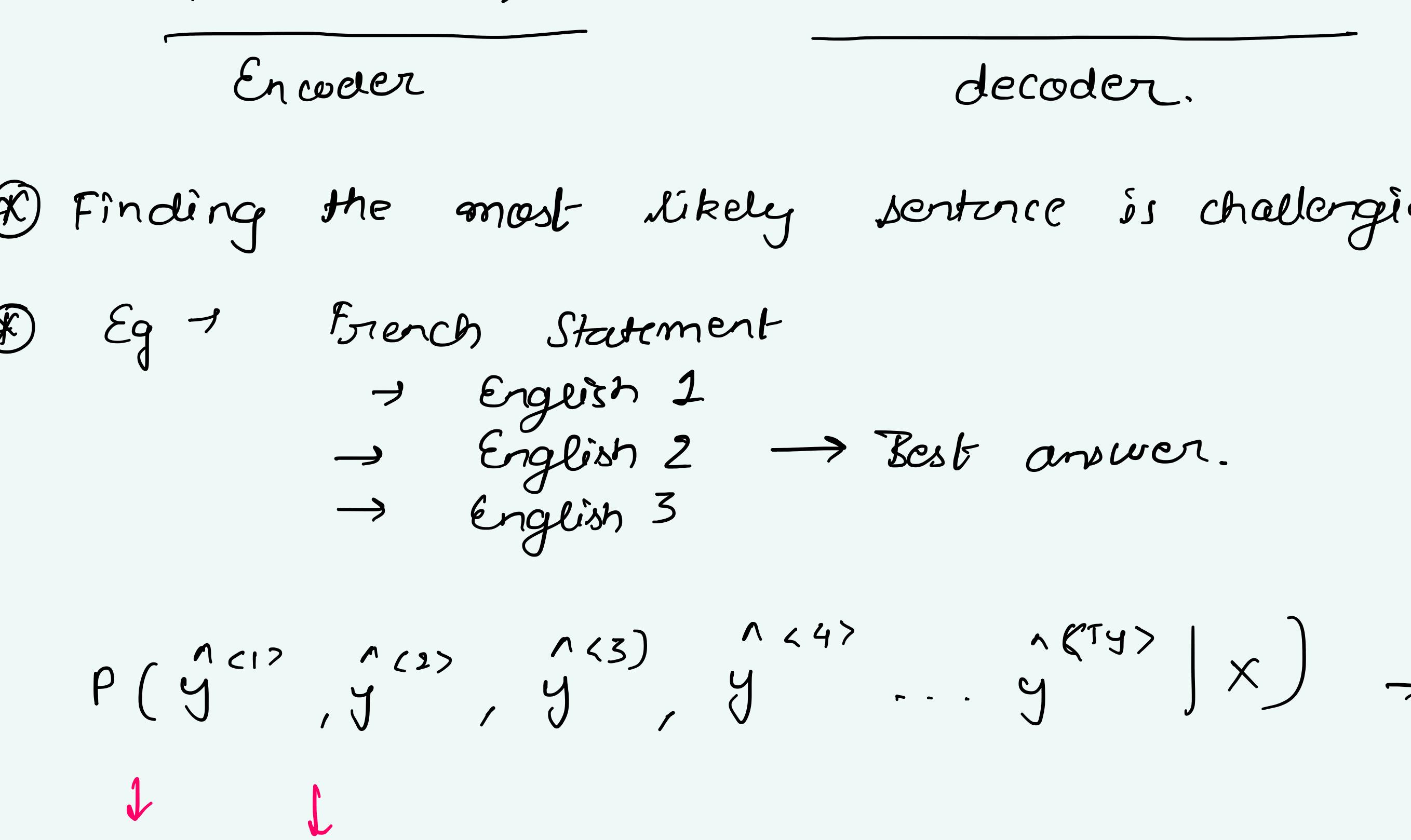


## Sequence to Sequence model :-

French sentence :-  
 ↳ English Sentence



④ Finding the most likely sentence is challenging.

⑤ Eg → French Statement  
 → English 1  
 → English 2 → Best answer.  
 → English 3

$$P(y^{<1>} | x) \cdot P(y^{<2>} | x) \cdot P(y^{<3>} | x) \cdot P(y^{<4>} | x) \dots P(y^{<Ty>} | x) \rightarrow \text{Maximize it}$$

Jane is visiting Africa in September

Jane is going to be visiting Africa in September.

$$P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$$

\* Selecting one word at a time is expensive and inaccurate.

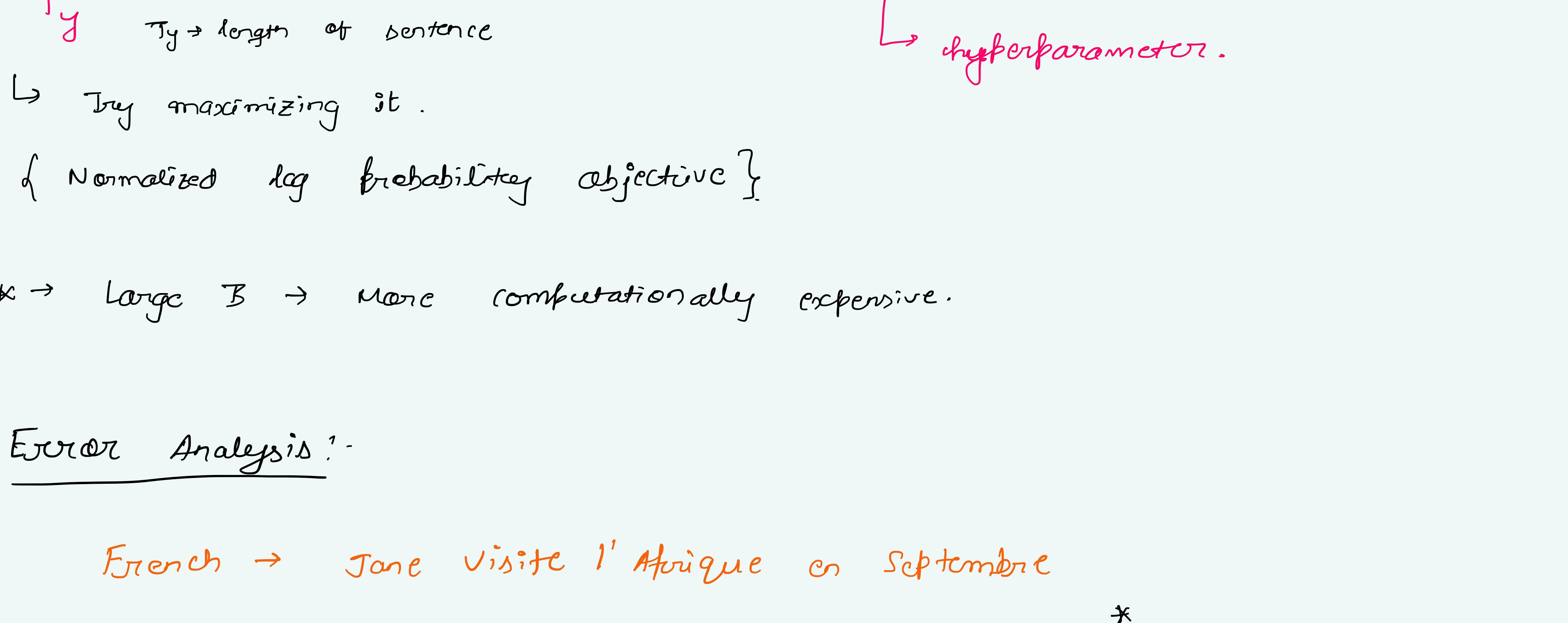
\* The final sentence might not be correct.

## Beam Search :- {most likely sentence}

Beam width = 3. (width = 1 is greedy approach)

It tries to pick the first word of English translation that it is going to output.

↳ Considering multiple outputs.



① Predict the first word using beam width (top 3)

② Consider their next word. ↳ Maximize this in second step.

$$P(1) \text{ in } \frac{P(y^{<1>} | x)}{P(y^{<1>} | x)}$$

$$P(2) \text{ Jane } \frac{P(y^{<1>} | x)}{P(y^{<1>} | x)}$$

$$P(3) \text{ Jane } \frac{P(y^{<1>} | x)}{P(y^{<1>} | x)}$$

Top-3 is never updated

↳ September → Rejected → \*\* Important

then the search moves forward, top 3 always updated accordingly

## Length Normalization:-

$$\log \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>})$$



$$\frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>})$$

$T_y \rightarrow$  length of sentence

$$\alpha = 0.7 \quad \alpha = 1 \quad \alpha = 0$$

↳ hyperparameters.

↳ Try maximizing it.

{ Normalized log probability objective }

\* → Large  $B$  → more computationally expensive.

## Error Analysis :-

French → Jane visite l'Afrique en Septembre

Human → Jane visits Africa in September

Algorithm → Jane visited Africa last september.

Case 1:-  $P(\hat{y}^* | x) > P(\hat{y} | x)$  Beam search is at fault

↳ increase beam width

Case 2:-  $P(\hat{y}^* | x) < P(\hat{y} | x)$  RNN is at fault.

↳ deeper layers → overfitting

↳ underfitting

\* Do a statistical analysis on the results.

## Attention Models:-



$$\alpha^{<t>} = (\alpha^{<t>}, \alpha^{<t>})$$

$$\sum_t \alpha^{<t>} = 1.$$

IMPORTANT

$$c^{<1>} = \sum_t \alpha^{<1,t>} * a^{<t>}$$

$$c^{<2>} = \sum_t \alpha^{<2,t>} * a^{<t>}$$

Computing attention  $\alpha^{<t,t>}$

$$\alpha^{<t,t>} = \frac{\exp(e^{<t,t>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

$$\alpha^{<t,t>} = \text{softmax}$$

$$\alpha^{<t,t>} \text{ and } e^{<t,t>} \text{ depend on } S^{<t-1>} \text{ and } a^{<t>}$$

④ Quadratic time to train

$$d \rightarrow 10 \rightarrow O(10^10) \rightarrow \text{very slow.}$$