

## Statistical Machine Translation (SMT)

Core idea: learn a probabilistic model from data.

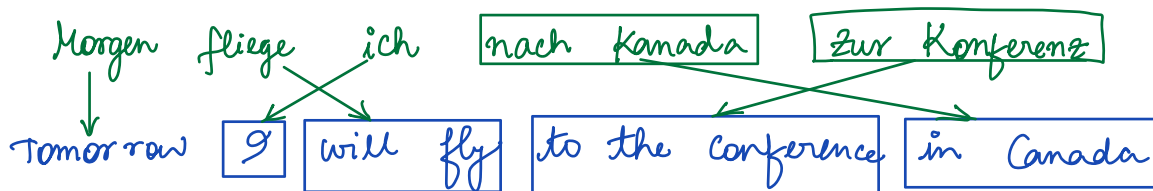
French  $\rightarrow$  English  
 $x$   $y$

$$\operatorname{argmax}_y P(y|x)$$

Use Bayes rule to break this down into:

$$\operatorname{argmax}_y \underbrace{P(x|y)}_{\text{Translation Model}} \underbrace{P(y)}_{\text{Language Model}}$$

## Alignment



## Spurious words

Japan	Le Japon
Shaken	Secoue
by	par
two	deux
new	nouveaux

## one to many

And	Le
the	programme
program	a
has	etc
been	mis

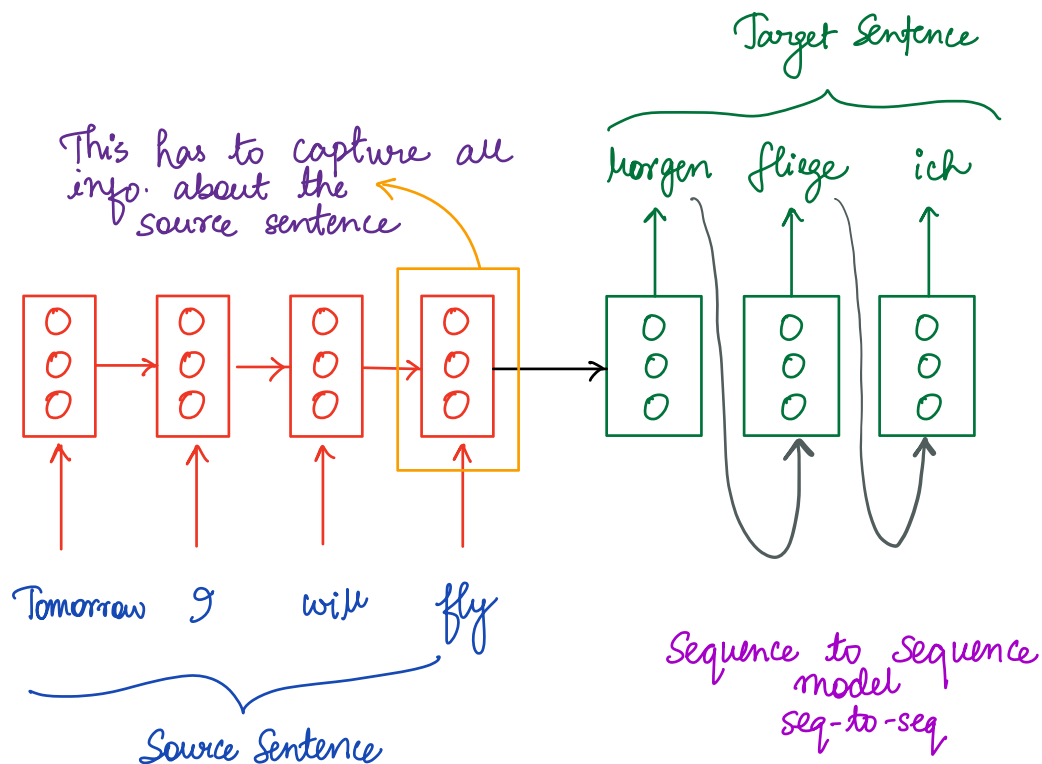
quakes

Seis'mes

implemented  $\left\{ \begin{array}{l} \text{en} \\ \text{application} \end{array} \right.$

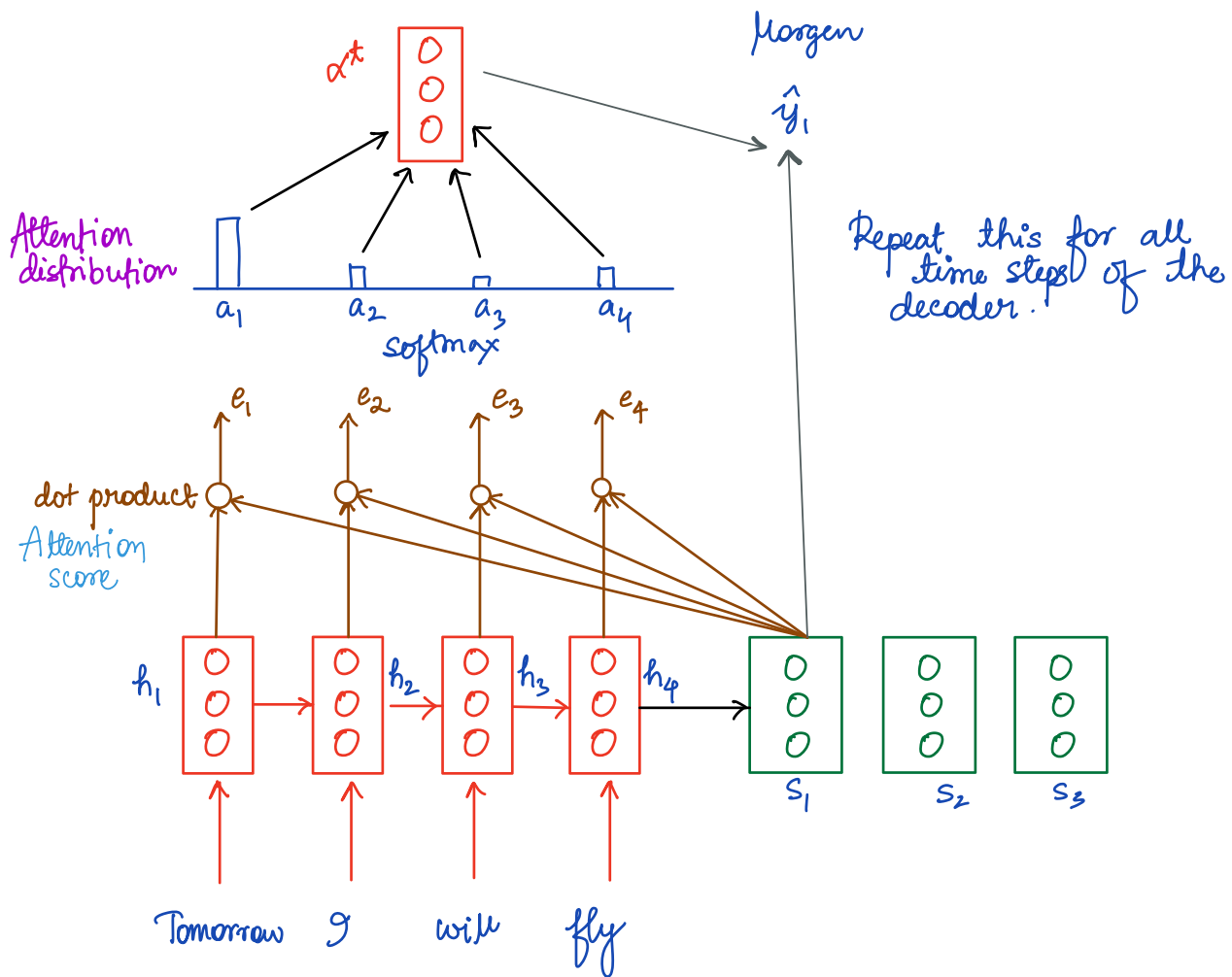
Alignment is complex.

## Neural Machine Translation (NMT)



Attention provides a solution to this problem.

Idea: use direct connections to the encoder to focus on a specific part of the input.



at time step (decoder)  $t$ .

$$e^t = \begin{bmatrix} s_t^T h_1 & s_t^T h_2 & \dots & s_t^T h_N \end{bmatrix}$$

$$e^t = [e_1^t \quad e_2^t \quad \dots \quad e_N^t]$$

$$a^t = \text{softmax}(e^t)$$

$$a^t = [a_1^t, a_2^t, \dots, a_N^t]$$

$$\alpha^t = \sum_{t=1}^N a_t h_t$$

$$[\alpha^t; s^t] \rightarrow \hat{y}^t$$