# SEQ2SEQ Modeling

SEQ 2 LABEL : $\quad P(y^{(T)} \mid x^{(1)}, \ldots, x^{(T+1)})$

SEQ 2 SEQ : $\quad P(y^{(1)}, \ldots, y^{(T_y)} \mid x^{(1)}, \ldots, x^{(T_x)})$

↳ APPLICATIONS

① $\underline{NMT}$     INPUT    FRENCH SENT
              OUTPUT    ENGLISH SENT

② TEXT      INPUT   PREFIX
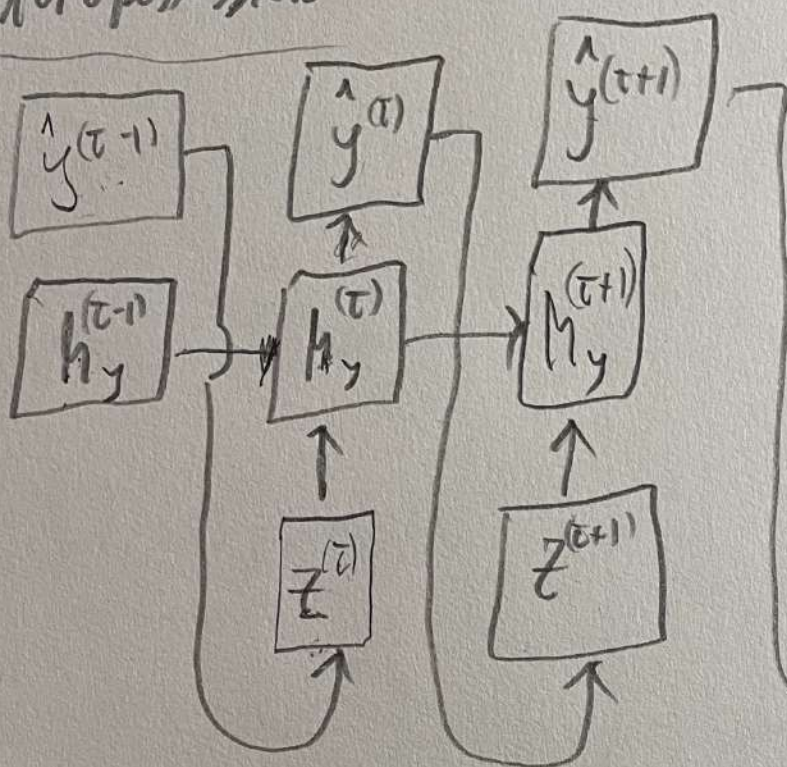   GEN       OUTPUT  SUFFIX

↳ FACTOR THE JOINT DIST $P(y^{(1)} \ldots y^{(T_y)} \mid x^{(1)} \ldots x^{(T_x)})$

INTO $\displaystyle\prod_{t=1}^{T_y} P(y^{(t)} \mid x^{(1)}, \ldots, x^{T_x}, y^{(1)}, \ldots, y^{(t-1)})$

↳

# Autoregression



$$\hat{y}^{(0)} = \langle START \rangle$$
$$\hat{y}^{(T_y)} = \langle STOP \rangle$$

## SEQ 2 SEQ Algorithm

① $h_y^{(0)} = h_x^{(T_x)} = f_h (x^{(1)}, \ldots x^{(T_x)})$

② Set $T=1$, $\hat{y}^{(0)} = \langle START \rangle$

③ while TRUE:

   $h^{(\tau)} = f_h (h_y^{(\tau-1)}, y^{(\tau-1)}) \rightarrow$ (evolving Hidden State

   $\hat{y}^{(\tau)} = f_y (h_y^{(\tau)}) \qquad \rightarrow$ PREDICT $\tau^{th}$ token

   if $\hat{y}^{(\tau)} = \langle STOP \rangle$ ;

   $\qquad$ Break

   $\tau = \tau + 1$

| SEQ 2 LABEL | SEQ2SEQ |
|---|---|
| $h^{(\tau)} = f_h(h^{(\tau-1)}, \underline{x^{(\tau)}})$ | $h^{(\tau)} = f_h(h^{(\tau-1)}, \underline{\hat{y}^{(\tau-1)}})$ |
| COMES FROM THE DATA! | COMES FROM THE MODEL! |

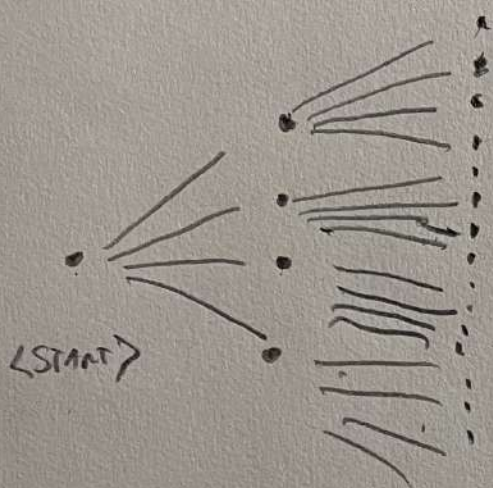## The INFERENCE GAP

INFERENCE TASK: $\hat{y}^{(1)} \dots \hat{y}^{(T_y)} = \underset{y^{(1)} \dots y^{T_y}}{ARGMAX} P(y^{(1)}, y^{(T_y)} | X^{(1)} \dots X^{(T_x)})$

TRAINING: $\hat{\theta} = \underset{\theta}{ARGMAX} \ P(y^{(t)} | f_h(h^{(t-1)}, y^{(t-1)}))$
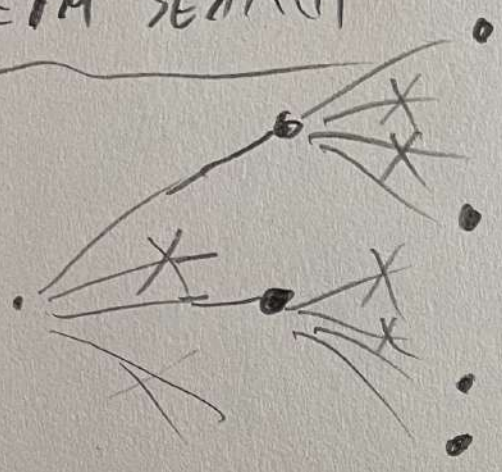
## SEARCH PROBLEM



&lt;START&gt;

$N = 10^5$

$T_y = 5$

$P(y^{(1)}, \dots, y^{(T_y)}) = \prod_{t=1}^{T_y} P(y^{(t)} | y^{(1)} \dots y^{(t-1)}) \Rightarrow MLE$

\# TRAJECTORIES: $N^{T_y}$

$\hookrightarrow$ E.g.: $(10^5)^5 = 10^{25}$ evaluations

# BEAM SEARCH



Reduce Branching factor from $N \Rightarrow b$, where $b \ll N$.

$$N = 10^5 \quad, \quad T_y = 5 \quad \Rightarrow 10^{25}$$

$$N = 10 \quad, \quad T_y = 5 \quad \Rightarrow 10^5$$

## INFERENCE GAP

During training: $L_{CE}\left( \hat{y}^{(t)}, P(y^{(t)} | \underline{y^{(1)} \cdots y^{(t-1)}}) \right)$

                                        ↑          ↑            ↑

                                  LABEL     PROB     LABELS

                                    "TEACHER Forcing"

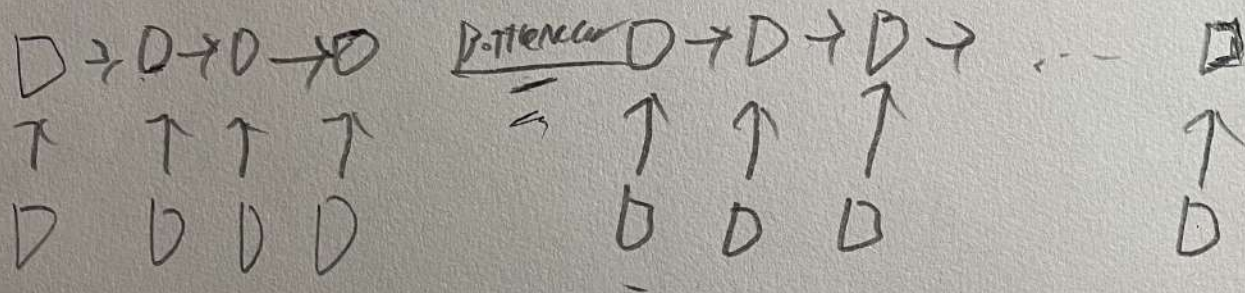During test: $P\left( y^{(t)} | \hat{y}^{(1)} \cdots \hat{y}^{(t-1)} \right)$

# Scheduled Sampling

$$L_{CE}\left(y^{(t)}, P(y^{(t)} \mid \tilde{y}^{(1)}, \dots, \tilde{y}^{(t-1)})\right)$$
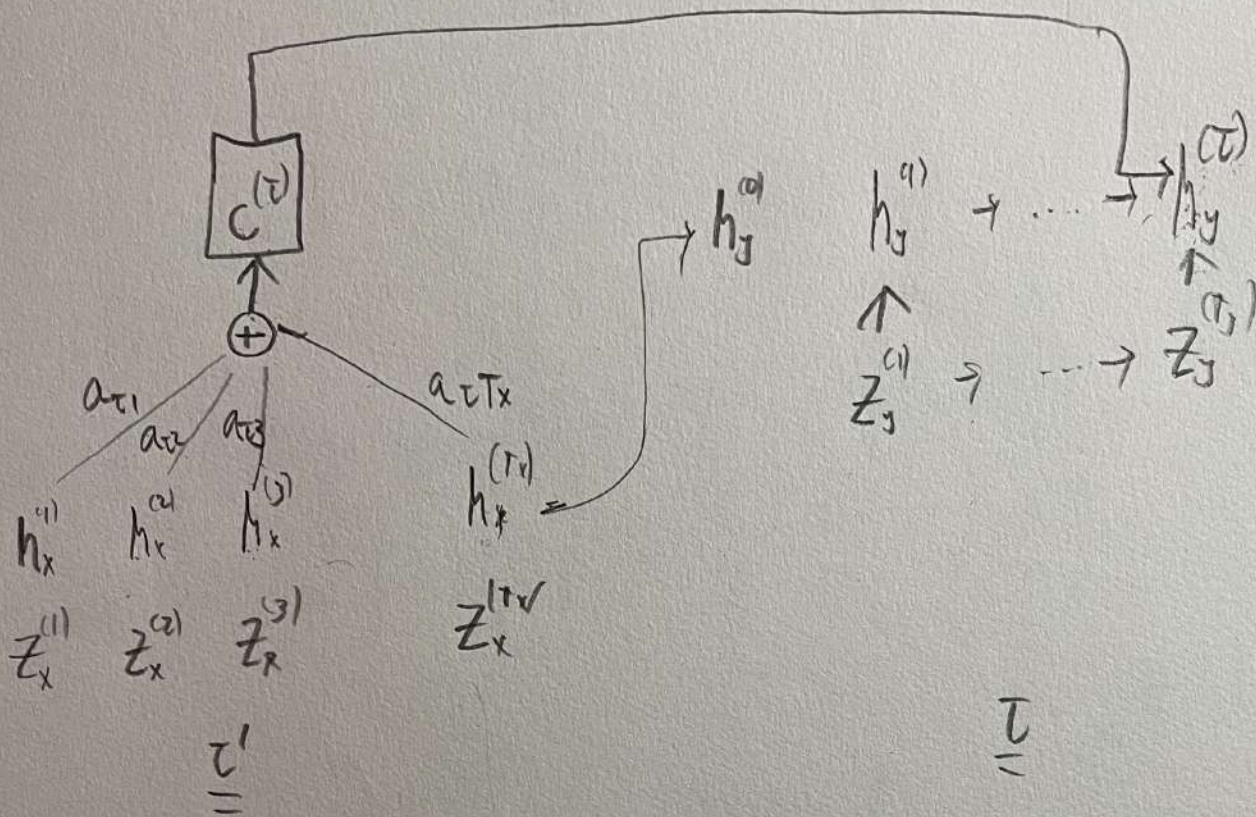
$$\tilde{y}^{(t)} = \begin{cases} \hat{y}^{(t)} & \text{with probability } P \\ y^{(t)} & \text{with probability } 1-P \end{cases}$$

---

## Information Bottleneck

Top diagram: $c^{(\tau)}$ box connects via $\bigoplus$ with weights $a_{\tau 1}, a_{\tau 2}, a_{\tau 3} \ldots a_{\tau T_x}$

$h_x^{(1)} \quad h_x^{(2)} \quad h_x^{(3)} \qquad h_x^{(T_x)} =$

$z_x^{(1)} \quad z_x^{(2)} \quad z_x^{(3)} \qquad z_x^{(T_x)}$

$\tau' =$

Right diagram:

$\to h_y^{(0)} \quad h_y^{(1)} \to \cdots \to h_y^{(\tau)}$

$\uparrow \qquad\qquad\qquad \uparrow$

$z_y^{(1)} \to \cdots \to z_y^{(\tau)}$

$\tau =$

$$\sum_{\tau'=1}^{T_x} a_{\tau, \tau'} = 1$$

$$\to A \in [0,1]^{T_y \times T_x}$$

$$c^{(\tau)} = \sum_{\tau'=1}^{T_x} a_{\tau, \tau'} h_x^{(\tau')}$$

Box:

$a_{\tau \tau'}$ is the degree of attention that $h_y^{(\tau)}$ has on $h_x^{\tau'}$.

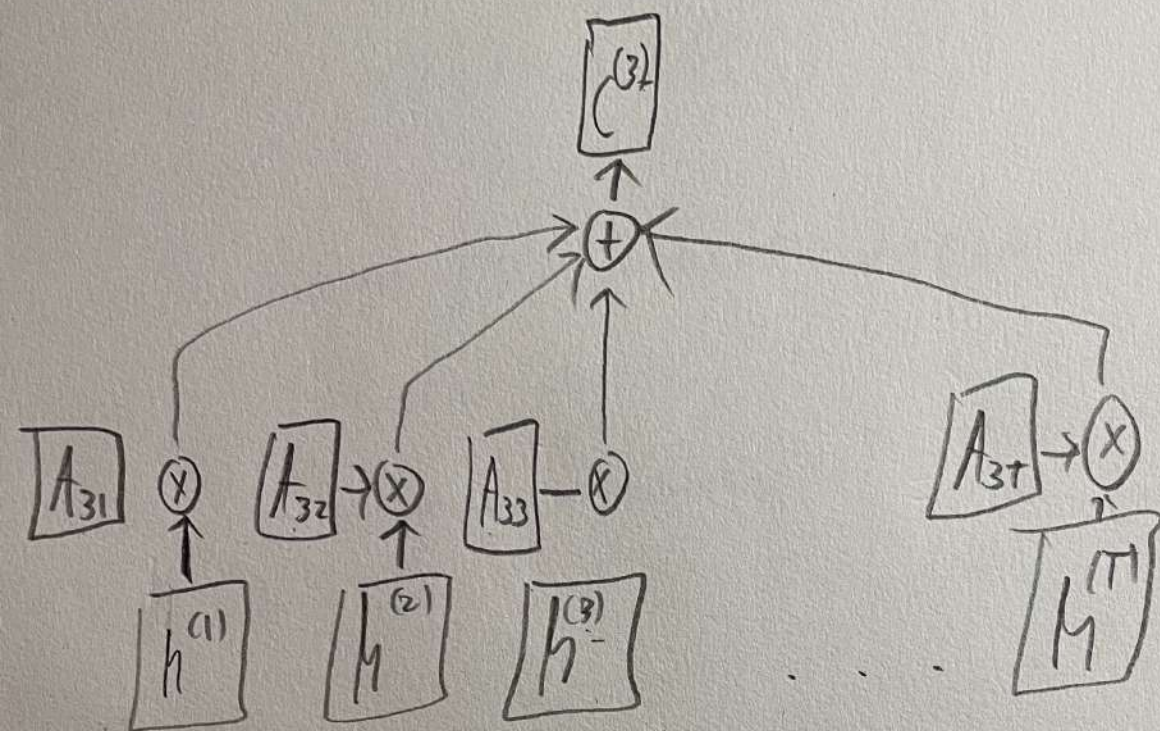$$a_{\tau \tau'} = \frac{e^{\alpha_{\tau, \tau'}}}{\sum_{\tau''}^{T_x} e^{\alpha_{\tau, \tau''}}}$$

$\alpha? \Rightarrow$ ① $\alpha_{\tau \tau'} = h_y^{(\tau)} \cdot h_x^{(\tau')}$

② $\alpha_{\tau \tau'} = f_{h_y}(h_y^{(\tau)}) \cdot f_{h_x}(h_x^{(\tau')})$

$\to f = NN$

# SELF ATTENTION



$$C^{(\tau)} = \sum_{\tau'} A_{\tau,\tau'} h^{(\tau')} \quad \text{where} \quad A_{\tau\tau'} = \frac{e^{h^{(\tau)} \cdot h^{(\tau')}}}{\sum_{\tau''} e^{h^{(\tau)} \cdot h^{(\tau'')}}}$$

$\hookrightarrow$ EACH ROW $\sum$ $A_{Row}$ $\sum / = 1$

# ATTENTION IS ALL YOU NEED — 2017

## TRICKS INTRODUCED :

① Scaled Dot Product Attention

RAW VALUE

SOFTMAX

$$h \in \mathbb{R}^D \implies \|h\| \text{ vs } D$$

$$D=2 \quad [1,1] = \sqrt{2}$$
$$D=3 \quad [1,1,1] = \sqrt{3} \qquad \sqrt{D}$$
$$D=4 \quad [1,1,1,1] = \sqrt{4}$$

$$\rightarrow D \text{ Scaled Dot Product} \rightarrow \boxed{\frac{1}{\sqrt{D}}} \text{ !}$$

② ATTENTION IS LIKE A "SOFT" LOOKUP TABLE

$$\text{for } c^{(i)} \Rightarrow \qquad \boxed{c^{(i)}}$$

$$\boxed{h^{(i)}}$$

③ Attention as a "Soft" Lookup

$h^{(4)} = $ Query
$h^{(1)} = $ key
$c$

$\boxed{c}$

$h^{(4)} = $ key, Query

$\boxed{h}$ $\boxed{h}$ $\boxed{h}$ $\boxed{h}$
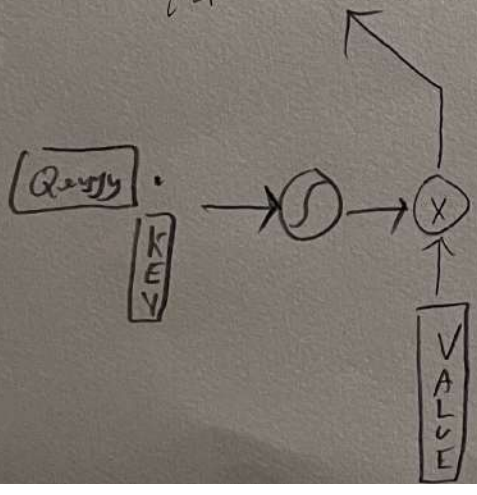
key   key   key   key

**hard**

$$Value = f[key]$$

**Soft**

$$Value = f[Query, key]$$

For Every $c^{(t)}$ There is:

① A Query $h^{(t)}$

② T Separate keys: $h^{(1)} \ldots h^{(T)}$

③ T Separate Values: $h^{(1)}, \ldots h^{(T)}$

$$c^{(t)} = DICT[Query^{(t)}] \quad \leftarrow \text{ "fuzzy" or "Soft" lookup}$$

$$= \sum_{t'=1}^{t} Value^{(t')} \cdot \sigma_{softmax}(Query^{(t)} \cdot key^{(t')})$$

$$= \sum_{t'=1}^{t} h^{(t')} \cdot \sigma_{softmax}(h^{(t)} \cdot h^{(t')})$$

$\boxed{Query} \cdot$
$\boxed{\begin{matrix}K\\E\\Y\end{matrix}}$ $\rightarrow \bigcirc S \rightarrow \bigotimes X$
$\boxed{\begin{matrix}V\\A\\L\\U\\E\end{matrix}}$

EACH INPUT h SERVES
AS THE key, Value, Query
DEPENDING ON WHICH POSITION
WE ARE AT.

④ LEARNABLE SET OF KEYS, VALUES, QUERIES

$$K_\tau = K h^{(\tau)} + b_r$$
$$q_\tau = Q h^{(\tau)} + b_q$$
$$V_\tau = V h^{(\tau)} + b_v$$

$$\Rightarrow \Theta = \{K, Q, V, b_r, b_q, b_v\}$$

⑤ Multihead SA:

$$K = \{K^{(1)} \ldots K^{(n)}\}$$
$$Q = \{Q^{(1)} \ldots Q^{(n)}\}$$
$$V = \{V^{(1)} \ldots V^{(n)}\}$$

$n$ = NUMBER OF ATTENTION HEADS

↳ $n$ CONTEXT VECTORS AT EACH POSITION IN THE SEQUENCE ⇒ $c^{(t)} = [c_1^{(t)}, \ldots, c_n^{(t)}]$

↳ $K, Q, V \in \mathbb{R}^{D/n}$ SUCH THAT WE CAN CONCAT TO GET $c^{(t)} \in \mathbb{R}^{D}$