

## Conditional language model

$$\begin{aligned} & p(Y), \text{ where } Y = (y_1, y_2, \dots, y_T) \\ & \dots \\ & p(Y|X) \text{ where } X = (x_1, x_2, \dots, x_{T'}) \end{aligned}$$

## Machine translation

- Learning

$$\max_{\theta} \frac{1}{N} \sum_{n=1}^N \log p(Y^n | X^n; \theta)$$

- Inference / generation

$$\arg \max_{\substack{Y \\ \dots}} \log p(Y | \underbrace{X}_{\text{given}})$$

## Parametrization

Autoregressive modeling

$$p(Y|X) = \prod_{t=1}^{T_Y} p(y_t | y_{<t}, X)$$

$$p(y_t | y_{<t}, X)$$

↓ context

• Input:  $(y_1, \dots, y_{t-1}) + (x_1, x_1, \dots, x_{T_X})$

• output:  $y_t \in V_Y = V$

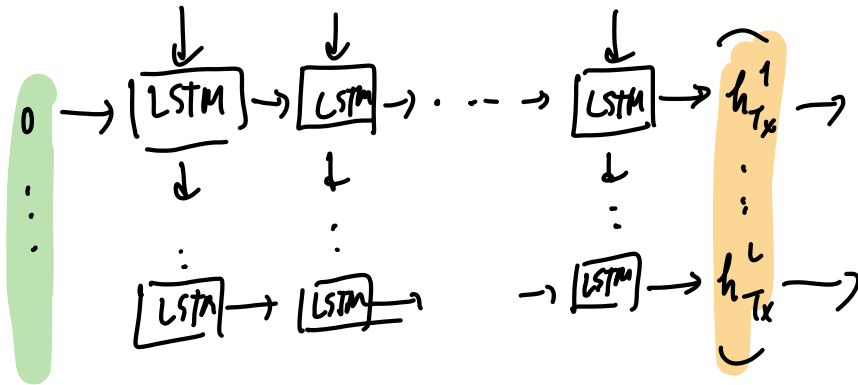
$$F(y_1, \dots, y_{t-1}, x_1, \dots, x_{T_X}, y_t)$$

$$= b_{y_t} + \underbrace{e(y_t)^T}_{\substack{\uparrow \\ \|e\|_2 = d}} G(y_1, \dots, y_{t-1}, x_1, \dots, x_{T_X})$$

Encoder : LSTM network

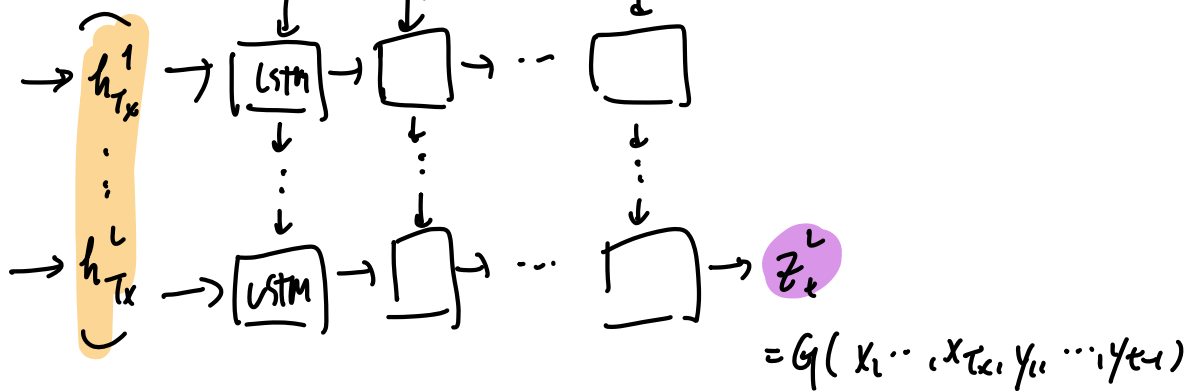
-  $e(x_t) \in \mathbb{R}^d$  for all  $t=1, \dots, T_x$

$(e(x_1), e(x_2), \dots, e(x_{T_x}))$



Decoder : LSTM network

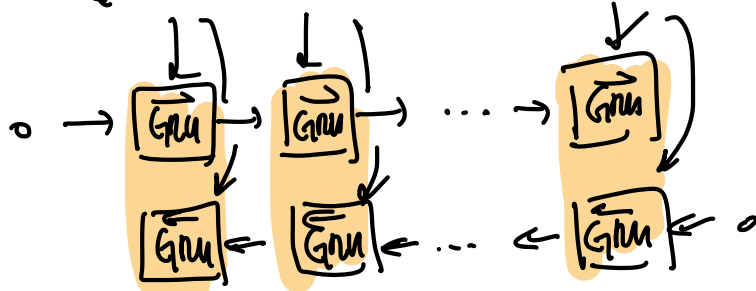
$(e(y_1), e(y_2), \dots, e(y_{t-1}))$



Bahdanau et al. [2014].

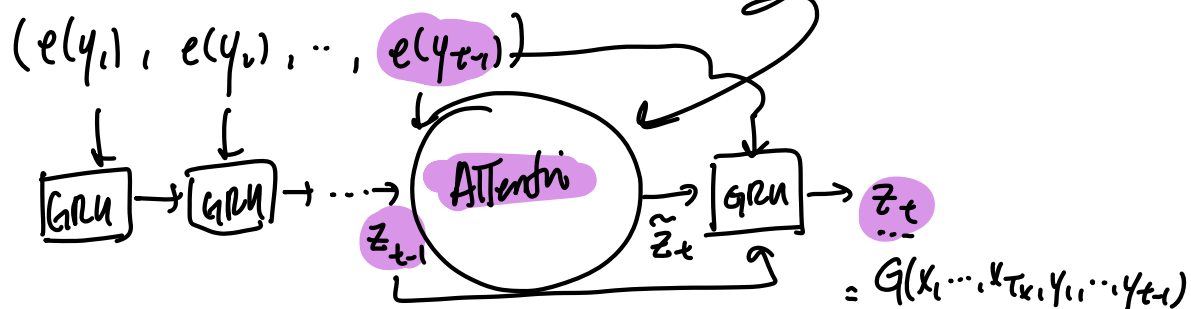
Encoder : bidirectional GRU/LSTM network.

$(e(x_1), e(x_2), \dots, e(x_{T_x}))$



$$(h_1, h_2, \dots, h_{T_x})$$

Decoder: unidirectional GRU (LSTM) network



① mix  $e(y_{t+1})$  &  $z_{t+1}$  into  $Q = \text{mix}(e(y_{t+1}), z_{t+1})$

②  $V_{t'} = K_{t'} = h_{t'}, t' = 1, \dots, T_x$

③  $d_{t'} \propto \underbrace{\exp\{Q^T K_{t'}\}}_{\text{softmax}}$

④  $\tilde{z}_t = \sum_{t'=1}^{T_x} d_{t'} h_{t'}$

Generative

$$\underset{y \in \mathcal{L}}{\text{argmax}} \log p(y|x)$$

$$p(y|x) = \prod_{t=1}^{T_y} p(y_t | y_{<t}, x)$$



greedy decoding

$$\hat{y}_t = \underset{y_t \in \mathcal{V}}{\text{argmax}} \log p(y_t | \hat{y}_{<t}, x)$$

learn search.

$$\{y_{\leq t}^1 \dots y_{\leq t}^k\} = \arg \max_{\substack{k=1, \dots, K \\ i=1, \dots, |V|}} \text{top-}k$$

$$\log p(\hat{y}_{\leq t}^k | x) + \log p(y_t = i | \hat{y}_{\leq t}^k, x)$$