$$h_t = u_t \odot \underbrace{\sigma(Uh_{t-1} + Wx_t + b)}_{= \hat{h}_t} + (1-u_t) \odot h_{t-1}$$

- $h_1 = u_1 \odot \hat{h}_1 + (1-u_1) \odot h_0$

- $h_2 = u_2 \odot \hat{h}_2 + (1-u_2) \odot u_1 \odot \hat{h}_1 + (1-u_2) \odot (1-u_1) \odot h_0$

- $h_3 = u_3 \odot \hat{h}_3 + (1-u_3) \odot u_2 \odot \hat{h}_2 + (1-u_3) \odot (1-u_2) \odot u_1 \odot \hat{h}_1$
$$+ (1-u_3) \odot (1-u_v) \odot (1-u_1) \odot h_0$$

$\vdots$

- $h_t = \sum_{t'=1}^{t} u_{t'} \odot \underbrace{\left(\prod_{t''=t'}^{t-1} (1-u_{t''})\right)}_{W} \odot \underbrace{\hat{h}_{t'}}_{h}$

$$= \sum_{t'=1}^{t} W_{t'}(x_{t'}, h_{t'-1}) \odot \hat{h}_{t'}(x_{t'}, h_{t'-1})$$

- Let's break the recurrence / temporal dependency.

① what if each $\hat{h}_{t'}$ was computed independently?
$$h_t = \sum_{t'=1}^{t} W_{t'}(x_{t'}, h_{t'-1}) \odot \hat{h}(x_{t'}, t')$$
positional embedding.

② what if each $W_{t'}$ was computed independently?
$$h_t = \sum_{t'=1}^{t} W(x_{t'}, x_t, t', t) \odot \hat{h}(x_{t'}, t').$$

Implementation.

- $W(x_{t'}, x_t, t', t) = \dfrac{\exp\{Q(x_t, t)^T K(x_{t'}, t')\}}{\sum_{t'=1}^{t} \exp\{Q(x_t, t)^T K(x_{t'}, t')\}}$

query        key

$$h(x_{t'}, t') = \underline{V}(x_{t'} + \underline{p(t')})$$

value.     positional embedding

• Multiple (attention) heads

$$h_{t,m} = \sum_{t'=1}^{t} W_m(x_{t'}, x_t) \odot \hat{h}_m(x_{t'})$$

$$h_t = [h_{t,1}; h_{t,2}; \cdots; h_{t,M}]^T$$

• Nonlinear fusion

$$h_t = f([h_{t,1}; \cdots; h_{t,M}]^T)$$

     nonlinear

$$f(\tilde{h}_t) = \max(0, U_f \max(0, W_f \tilde{h}_t + b_f) + c_f) + x_t.$$

• Self-attention : non-"causal" attention

$$h_{t,m} = \sum_{t'=1}^{T} W_m(x_{t'}, x_t, t'; t) \odot \hat{h}_m(x_{t'}, t')$$

Imputing

$$\underset{y}{\arg\max} \, \log \, p(x_{pre}, y, x_{post})$$

     seq   seq   seq

     prefix   ↑   suffix

        missing .

$$(\Leftrightarrow) \, \underset{y}{\arg\max} \, \log \, p(y | x_{pre}, x_{post}) + \log \, p(x_{pre}, x_{post})$$

$$(x_1, x_2, \cdots, x_t, y, x_{t+|y|}, \cdots, x_T)$$

$$(X_1, X_2, \cdots, X_t, \langle mask \rangle, \langle mask \rangle, \cdots, \langle mask \rangle, X_{t+|y|}, \cdots, X_T)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{|y|}$$

$$\Downarrow$$

$$Y$$

① $\begin{cases} \text{masked-out indices} : M_1, M_2, \cdots, M_{T_m} \\ \text{observed indices} : O_1, O_2, \cdots, O_{T_o} \\ \qquad T_m + T_o = T \end{cases}$

② $\}$ Corrupt $(X)$   $\langle mask \rangle$

③ $p(X_{m_1}, \cdots, X_{m_{T_m}} \mid \text{Corrupt}(X))$

$$= \prod_{i=1}^{T_m} p(X_{m_i} \mid \text{Corrupt}(X))$$

Objective funct:

$$J(\theta) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}_{M,O \sim \gamma(T_n)} \left[ \sum_{i=1}^{T_m^n} \log p(X_{m_i}^n \mid X_{O_1}^n \cdots X_{O_{T_o}}^n; \theta) \right]$$

$$\Longrightarrow BERT, \cdots$$

$$\overbrace{\qquad\qquad\qquad\qquad\qquad}$$

$[\text{masked}]$ language modeling

① $p(X) = \prod_{t=1}^{T} p(X_t \mid \underbrace{X_1, \cdots, X_{t-1}}_{\text{observed}}, \underbrace{X_{t+1} = \langle mask \rangle, \cdots, X_T = \langle mask \rangle}_{\text{masked}})$

$\underset{\text{masked}}{\cup}$

② $F(X) = \sum_{t=1}^{T} \log p(X_t \mid X_{1:t-1}, X_t = \langle mask \rangle, X_{t+1:T})$

$\Downarrow$                                    pseudo log likelihood

neg. energ.

$$p(x) = \frac{\exp\{\bar{F}(x)\}}{\sum\limits_{x' \in L} \exp\{\bar{F}(x')\}}$$

③ Implicit way.