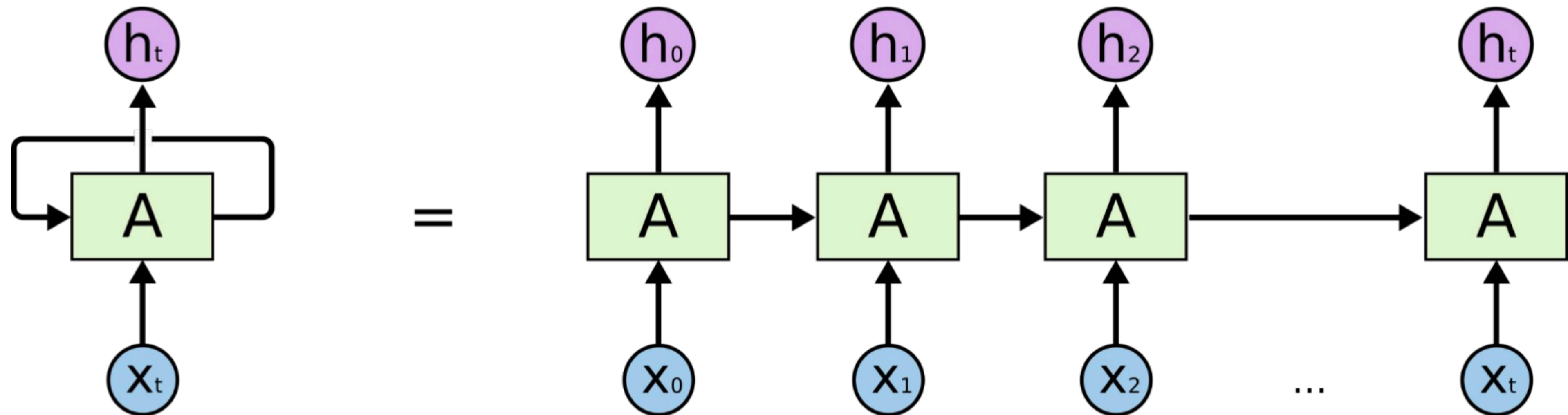


# **Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation**

Authors: Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Yoshua Bengio

# Background - RNN

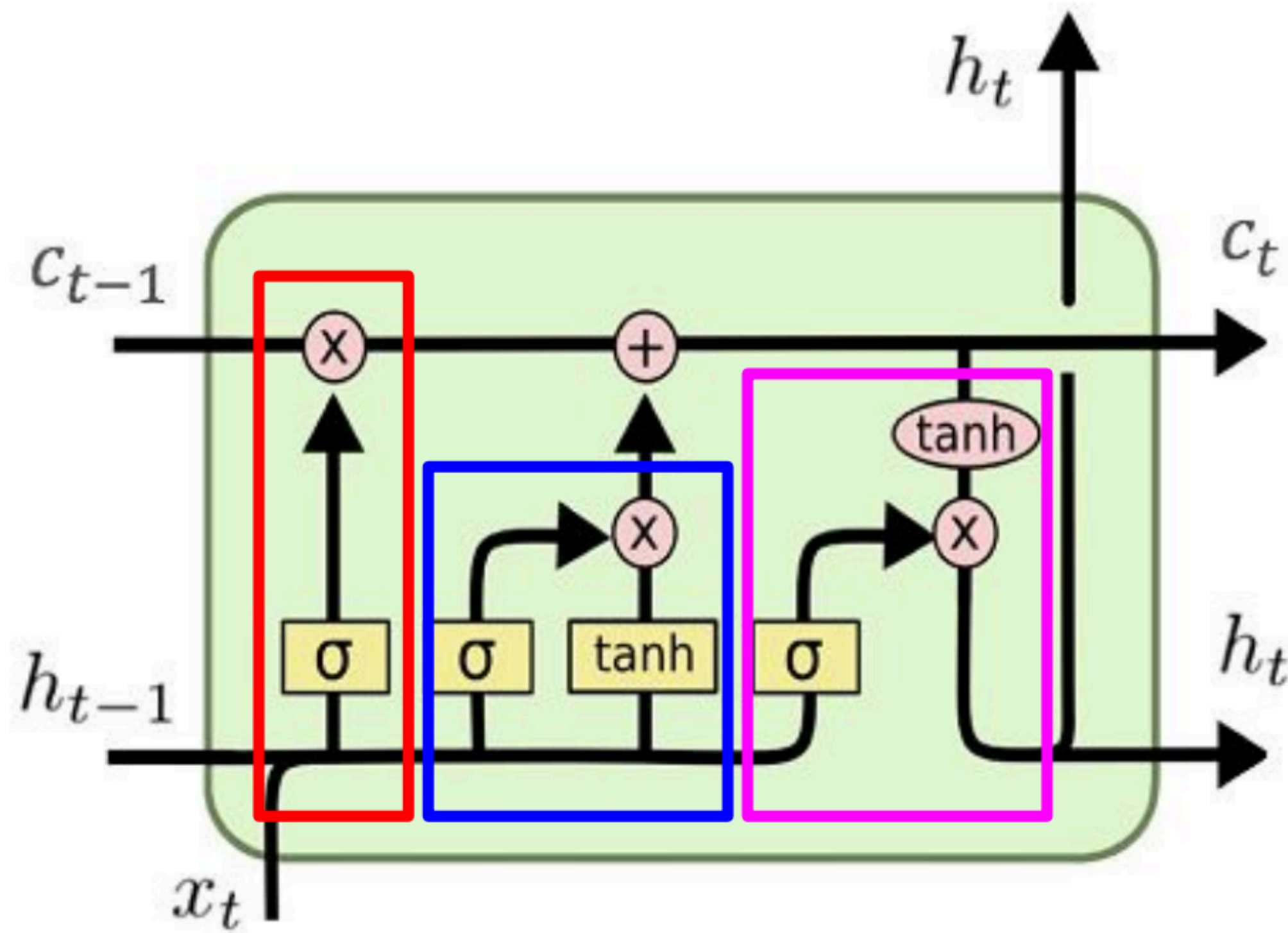


designed to process sequential data by maintaining an internal memory state

feedback connection: hidden state at each time step used as an another input

vanishing/exploding gradient problem

# Background - LSTM



- **Forget Gate**
  - how much **old** information to **remember**
- **Input Gate**
  - how much **new** information to **add**
- **Output Gate**
  - how much information to **expose** to next **output** or **hidden state**

# Problem Description - SMT(statistical machine translation)

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

maximizes the conditional probability given a source sequence

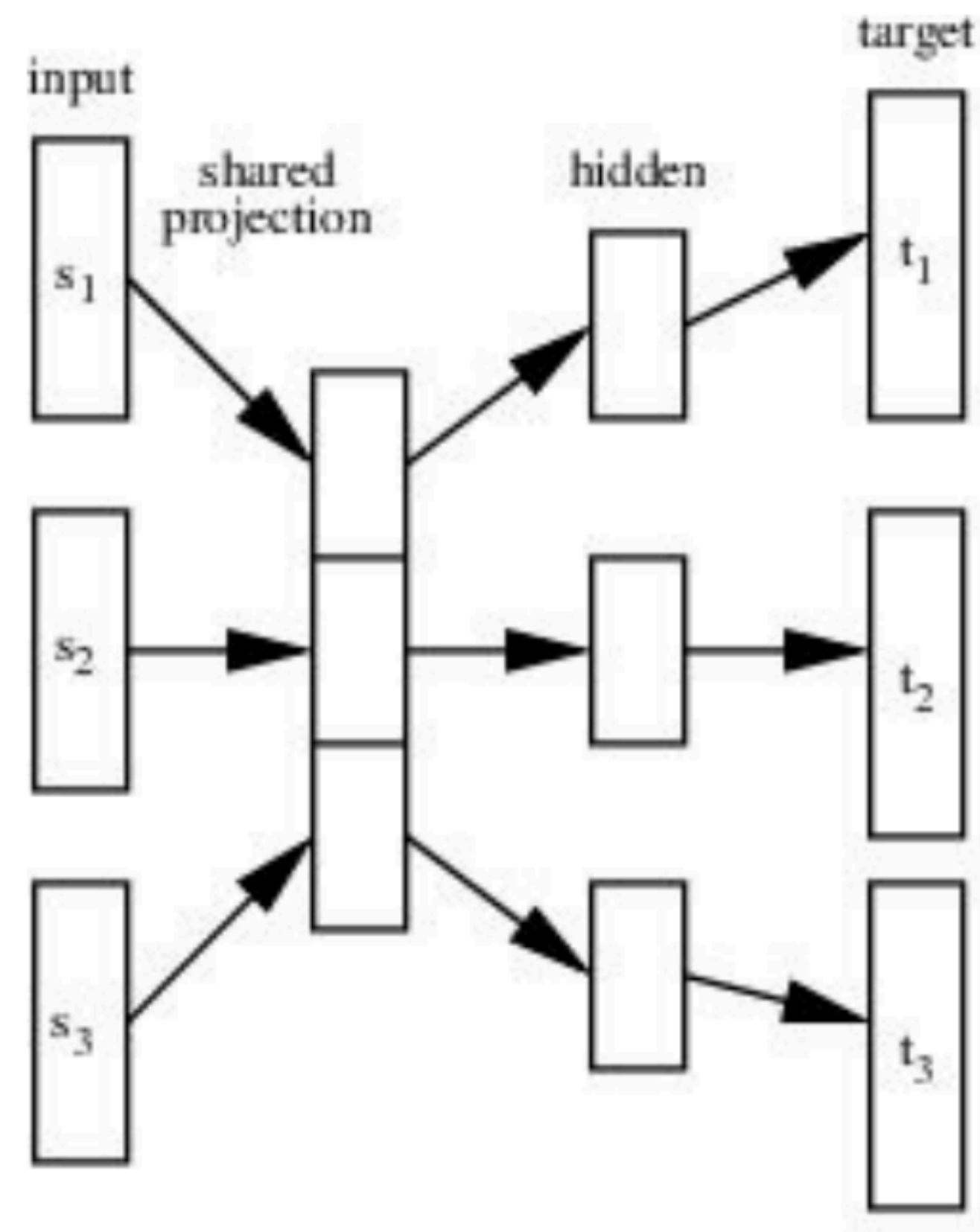
maximize

$$\log p(\mathbf{f} \mid \mathbf{e}) = \sum_{n=1}^N w_n f_n(\mathbf{f}, \mathbf{e})$$

log-linear model with features and weights is optimized to maximize the BLEU score

$$\text{BLEU} = \min \left( 1, \frac{\text{output-length}}{\text{reference-length}} \right) \left( \prod_{i=1}^4 \text{precision}_i \right)^{\frac{1}{4}}$$

# Related works - neural networks in machine translation

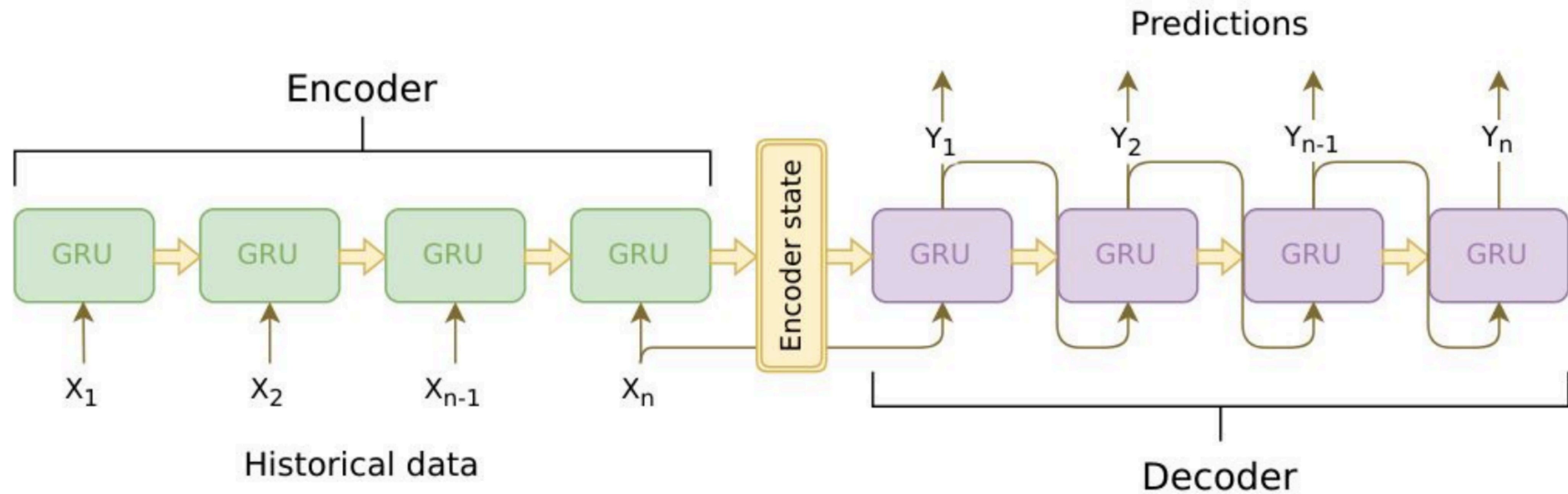


(Schwenk, 2012)

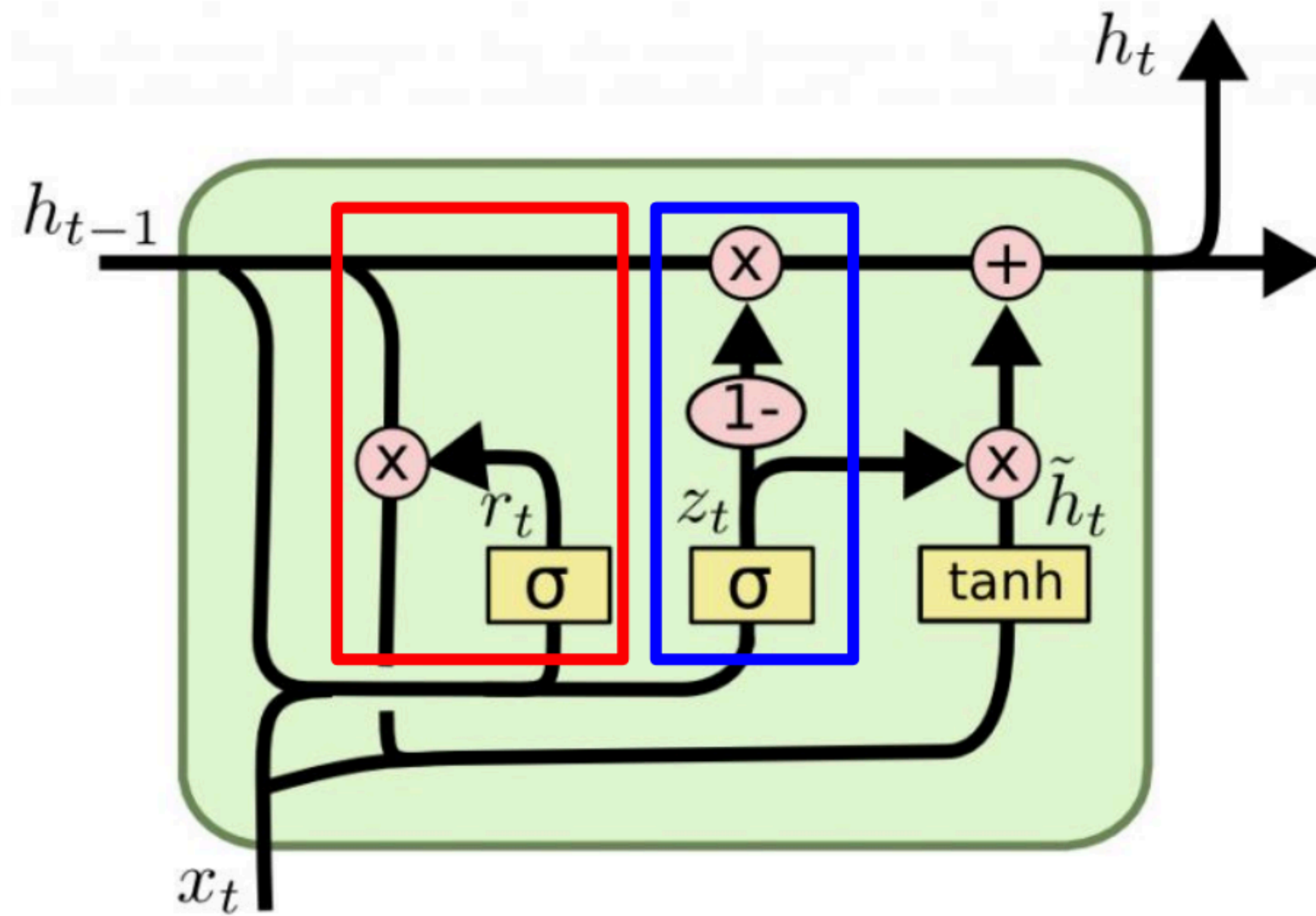
- Fixed input/output length
- Prediction of single word at a time
- Bag-of-words: absence of order information



# Method -RNN Encoder & Decoder



# Method -GRU(Gate Recurrent Unit)



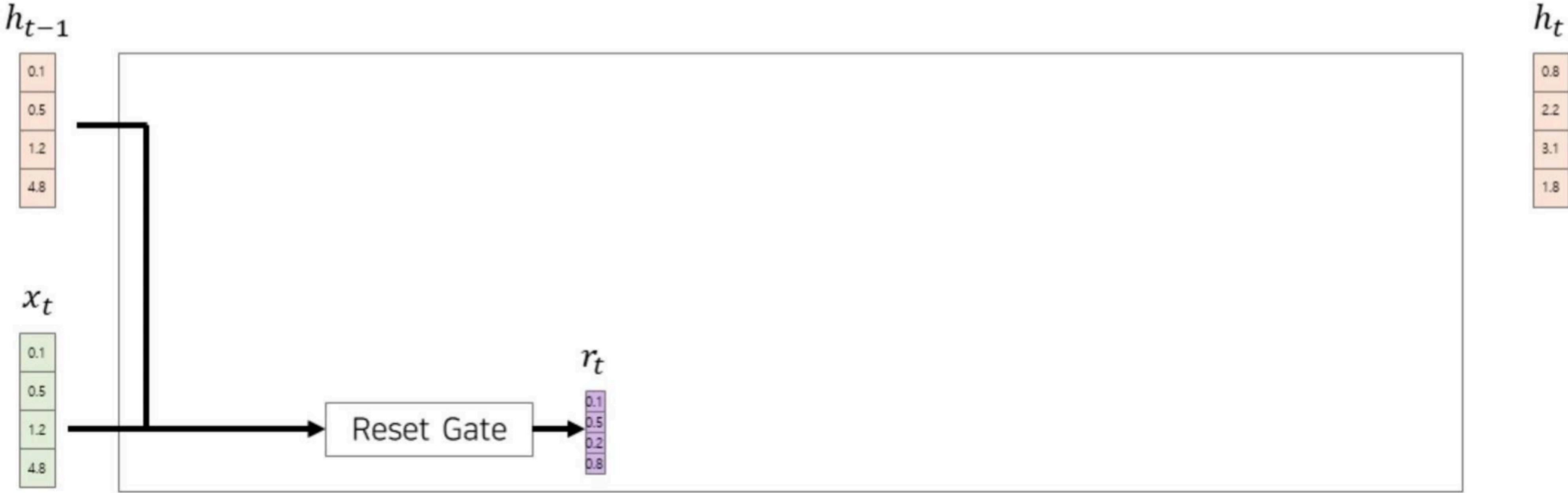
- **Reset Gate**
  - how much **previous** information to **forget**
- **Update Gate**
  - how much previous information from **reset gate** to **carry**

# Method -GRU(Gate Recurrent Unit)

## Reset Gate

$$\sigma \left( \begin{bmatrix} & & & \\ & W_r & & \\ & & & \\ & & & \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} + \begin{bmatrix} & & & \\ & U_r & & \\ & & & \\ & & & \end{bmatrix} \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} \right) = \begin{bmatrix} 0.1 \\ 0.5 \\ 0.2 \\ 0.8 \end{bmatrix}$$

$$r_j = \sigma \left( [\mathbf{W}_r \mathbf{x}]_j + [\mathbf{U}_r \mathbf{h}_{\langle t-1 \rangle}]_j \right)$$



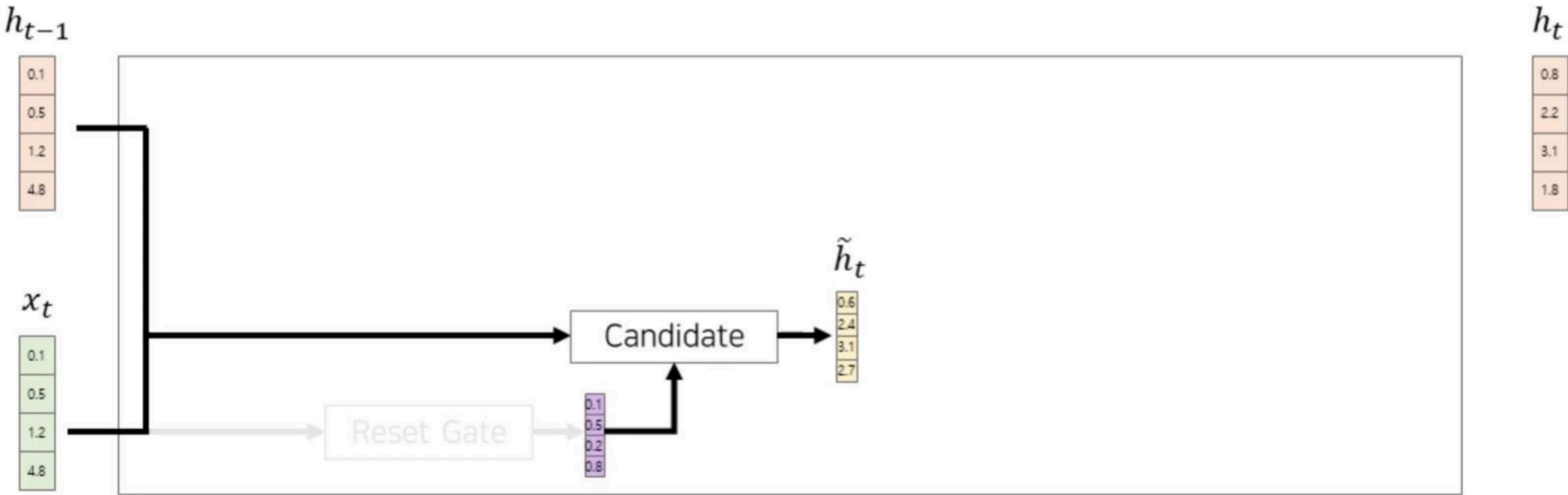


# Method -GRU(Gate Recurrent Unit)

## Candidate Gate

$$\tau \left( \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} W \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} * \begin{bmatrix} 0.1 \\ 0.5 \\ 0.2 \\ 0.8 \end{bmatrix} + \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} U \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} \right) = \begin{bmatrix} 0.6 \\ 2.4 \\ 3.1 \\ 2.7 \end{bmatrix}$$

$$\tilde{h}_j^{(t)} = \phi \left( [\mathbf{W}\mathbf{x}]_j + [\mathbf{U}(\mathbf{r} \odot \mathbf{h}_{\langle t-1 \rangle})]_j \right)$$

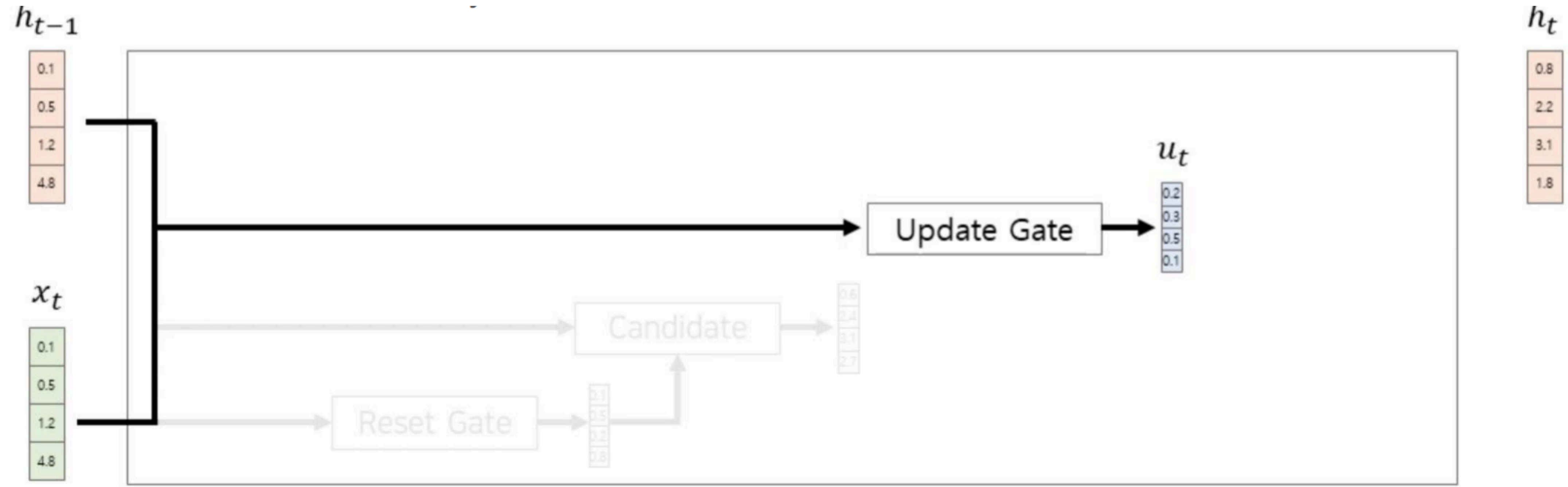


# Method -GRU(Gate Recurrent Unit)

## Update Gate

$$\sigma \left( \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} W_u \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} + \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} U_u \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} \right) = \begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \\ 0.1 \end{bmatrix}$$

$$z_j = \sigma \left( [\mathbf{W}_z \mathbf{x}]_j + [\mathbf{U}_z \mathbf{h}_{\langle t-1 \rangle}]_j \right)$$



# Method -GRU(Gate Recurrent Unit)

## Hidden State

$$\begin{bmatrix} 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \end{bmatrix} - \begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \\ 0.1 \end{bmatrix} * \begin{bmatrix} 0.1 \\ 0.5 \\ 1.2 \\ 4.8 \end{bmatrix} + \begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \\ 0.1 \end{bmatrix} * \begin{bmatrix} 0.6 \\ 2.4 \\ 3.1 \\ 2.7 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 2.2 \\ 3.1 \\ 1.8 \end{bmatrix}$$

$$h_j^{(t)} = z_j h_j^{(t-1)} + (1 - z_j) \tilde{h}_j^{(t)}$$

