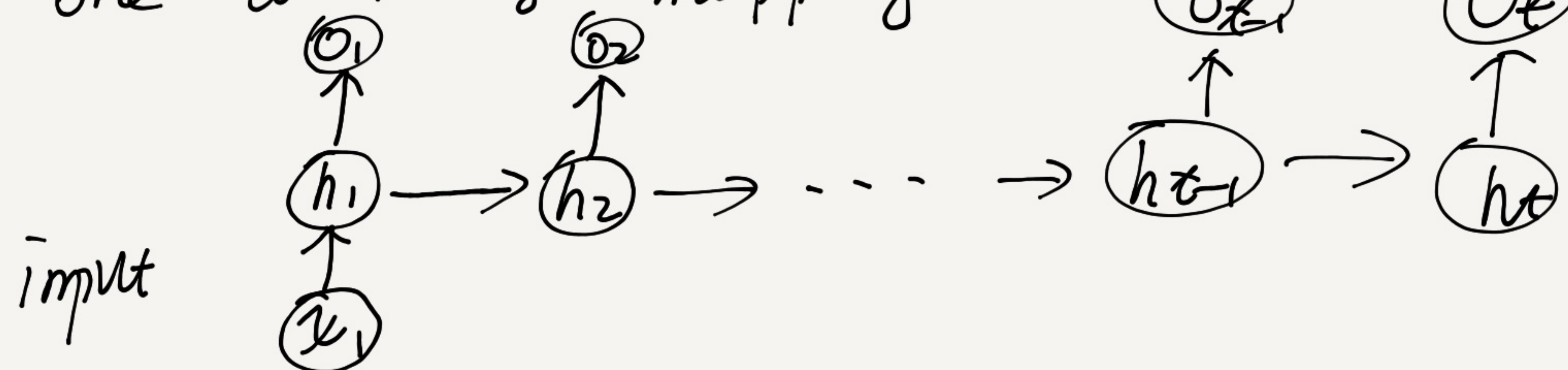


# RNN architectures.

## 1. Basic architectures.



### (1) one-to-many mapping

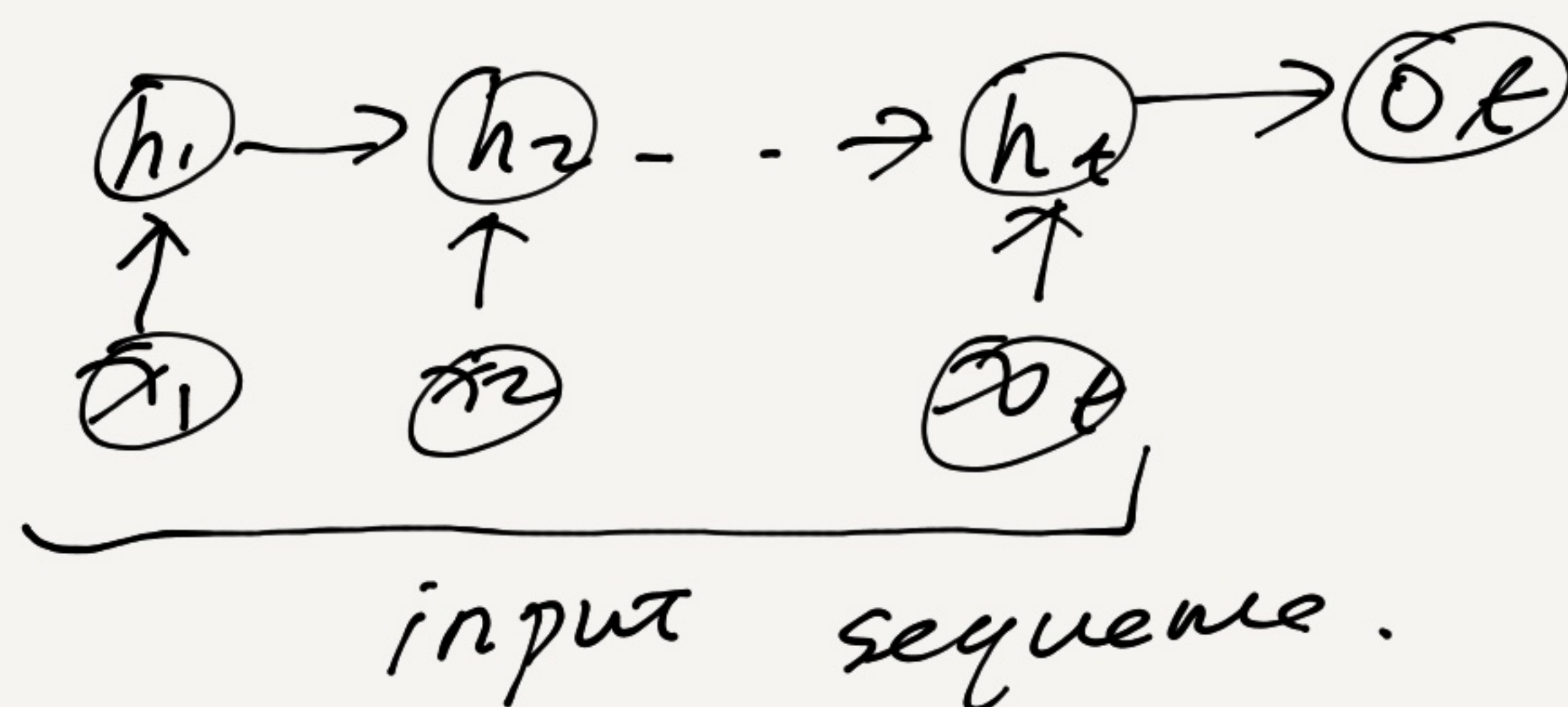


Application: image captioning.

$x_1$ : an image

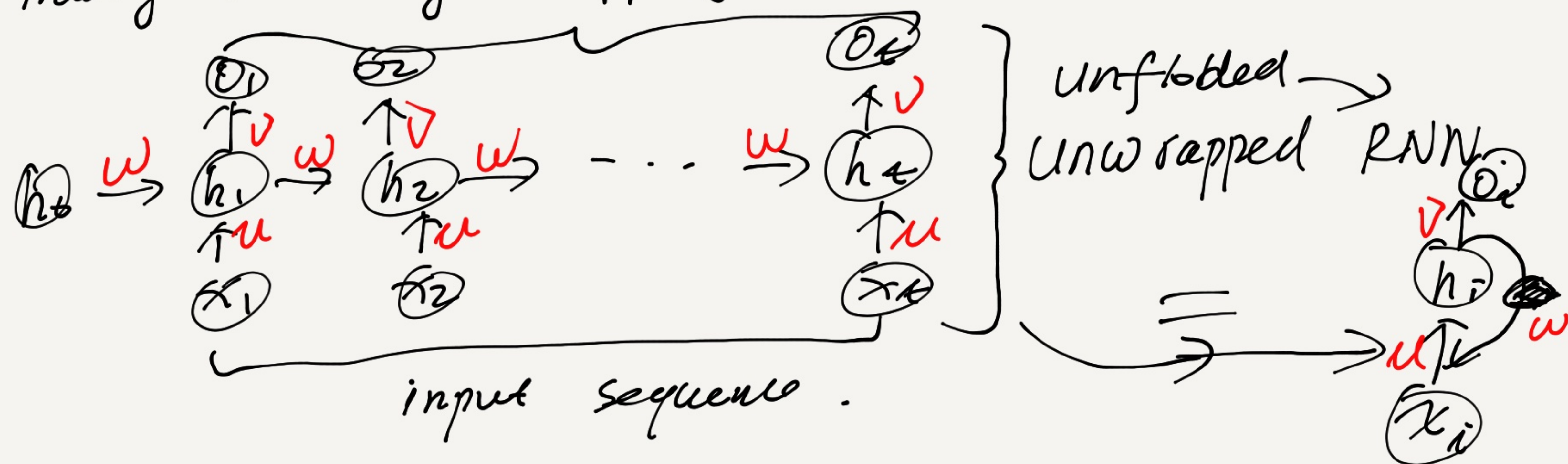
$O_1, \dots, O_t$ : sequence of words.

### (2) many-to-one mapping





### ③ many-to-many mapping



input:  $x_1 \rightarrow h_1$  (hidden state)

$x_2, h_1 \rightarrow h_2$

$x_3, h_2 \rightarrow h_3$

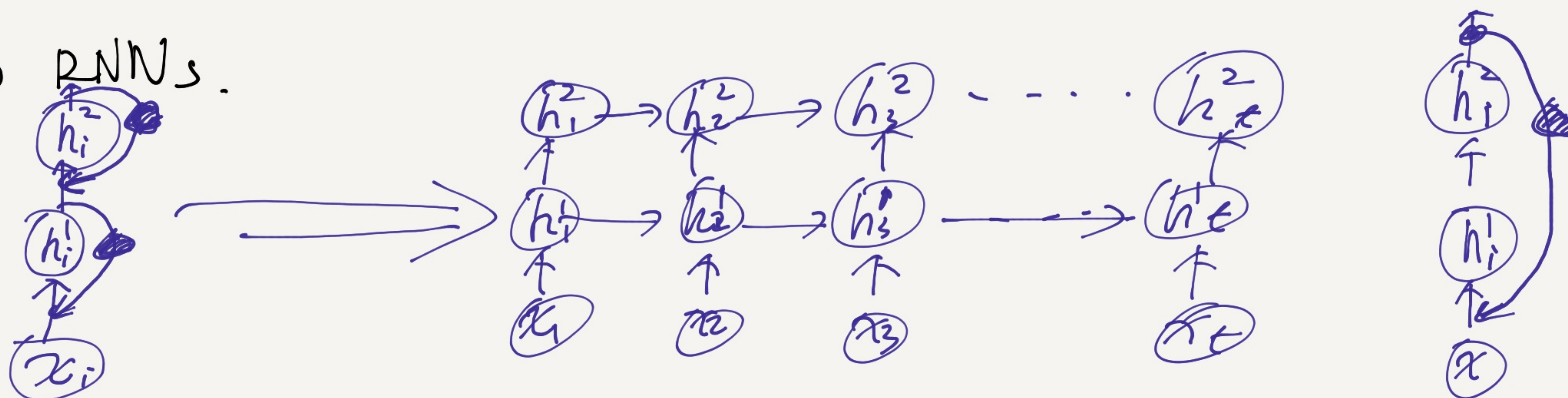
▷ RNNs share weight across time.

▷ carry over past information to  $(h)$

future by using the feedback connection.

Machine translation:

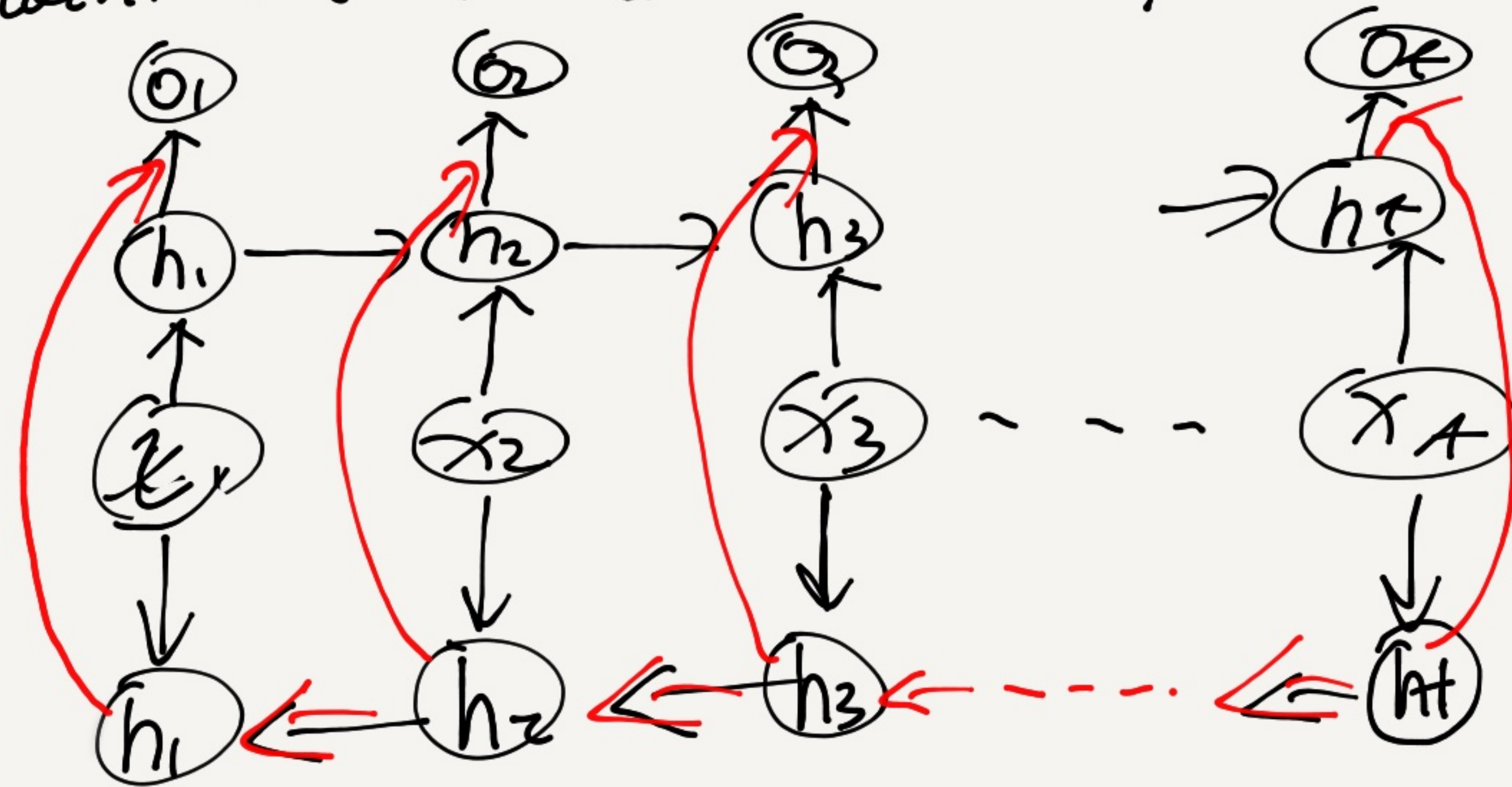
### 2 Deep RNNs.





3 Bidirectional RNNs use both past and future values in a sequence.

machine translation and speech recognition.



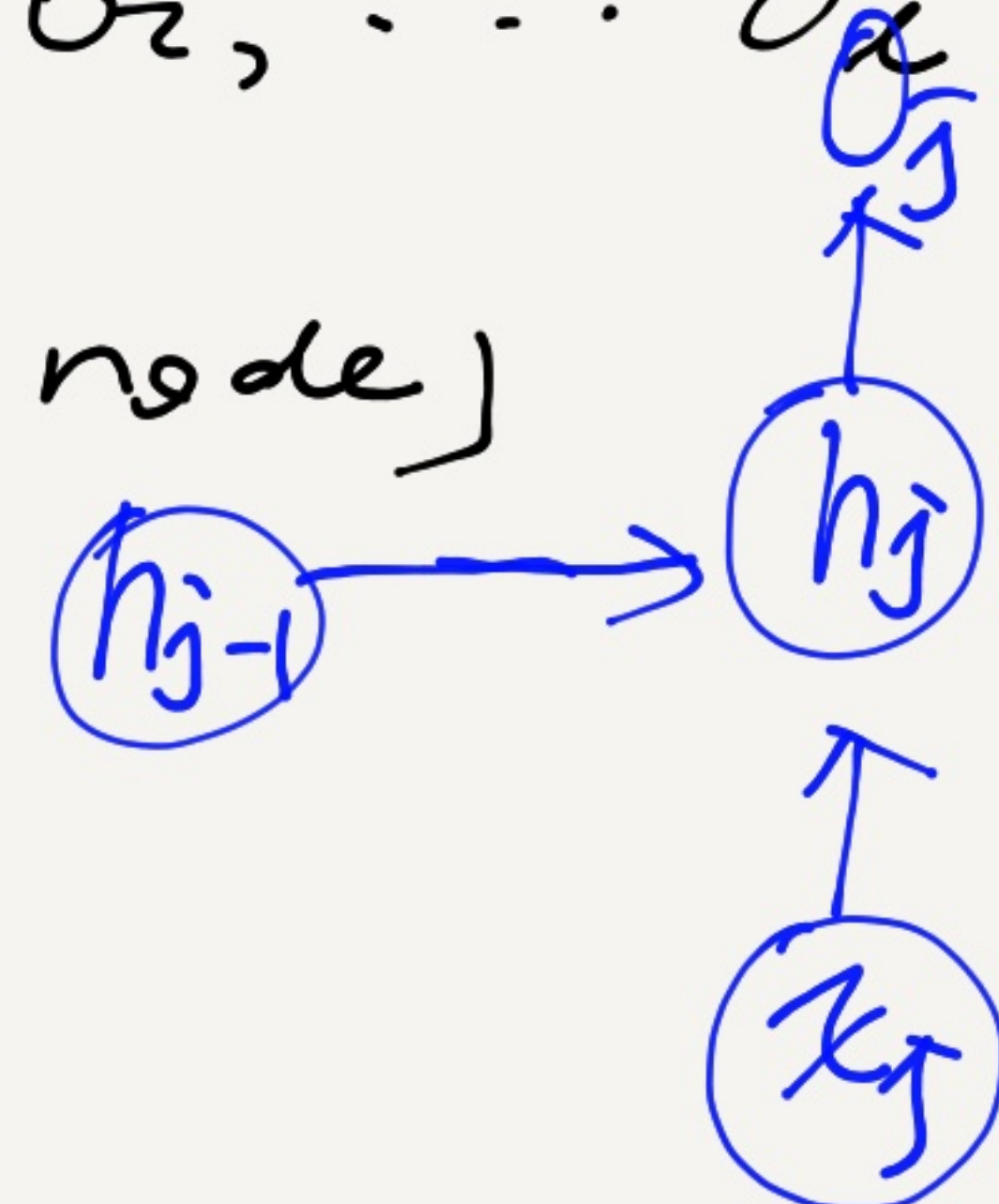
BERT (LLM)

4. Optimize RNNs: (BP through time)

$$L_{seq} = \sum_{j=1}^n -y_j \cdot \log o_j \quad (\text{cross-entropy between two sequences})$$

$$o_j = g_o(u h_j + b_o) \quad (\text{output node}) \quad \begin{cases} y_1, y_2, \dots, y_n \quad (\text{target}) \\ o_1, o_2, \dots, o_n \quad (\text{pred.}) \end{cases}$$

$$h_j = g_h(u x_j + \boxed{w h_{j-1}} + b_h) \quad (\text{hidden node})$$



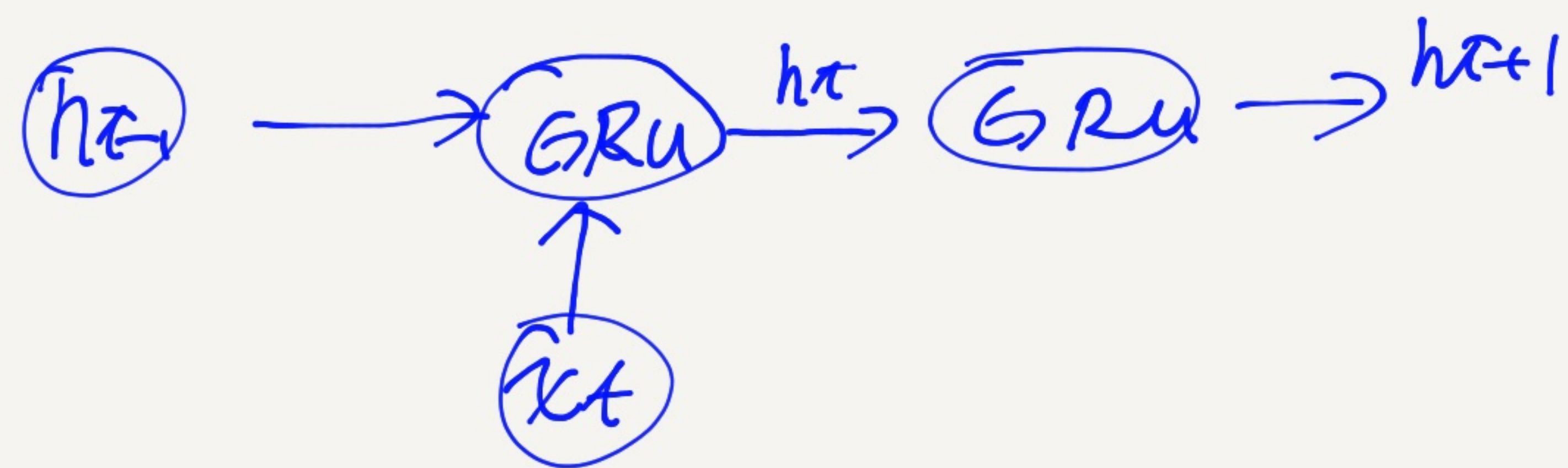


## 5. Two popular RNNs

1) LSTM (long short-term memory) 1997.

2) Gated Recurrent Unit (GRU)

① overall architecture.



② 3-component in GRU.

update gate:  $z_t$  : controls how much information from the past

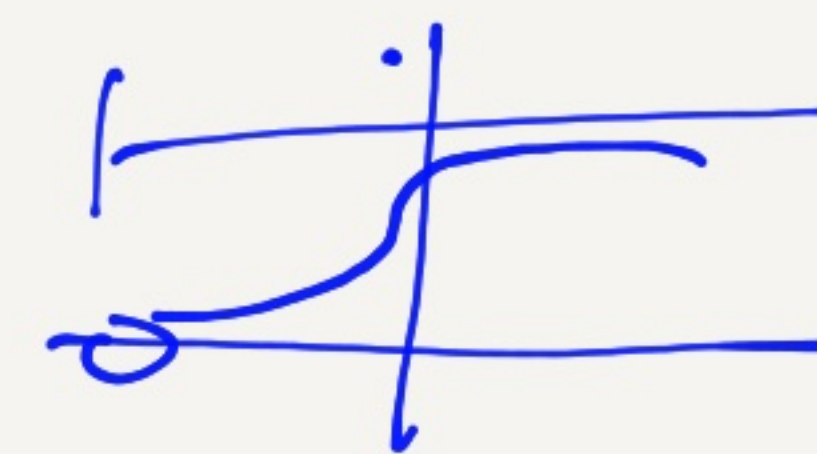
reset gate:  $r_t$  controls how much information to ignore from the past.



$z_t$ :

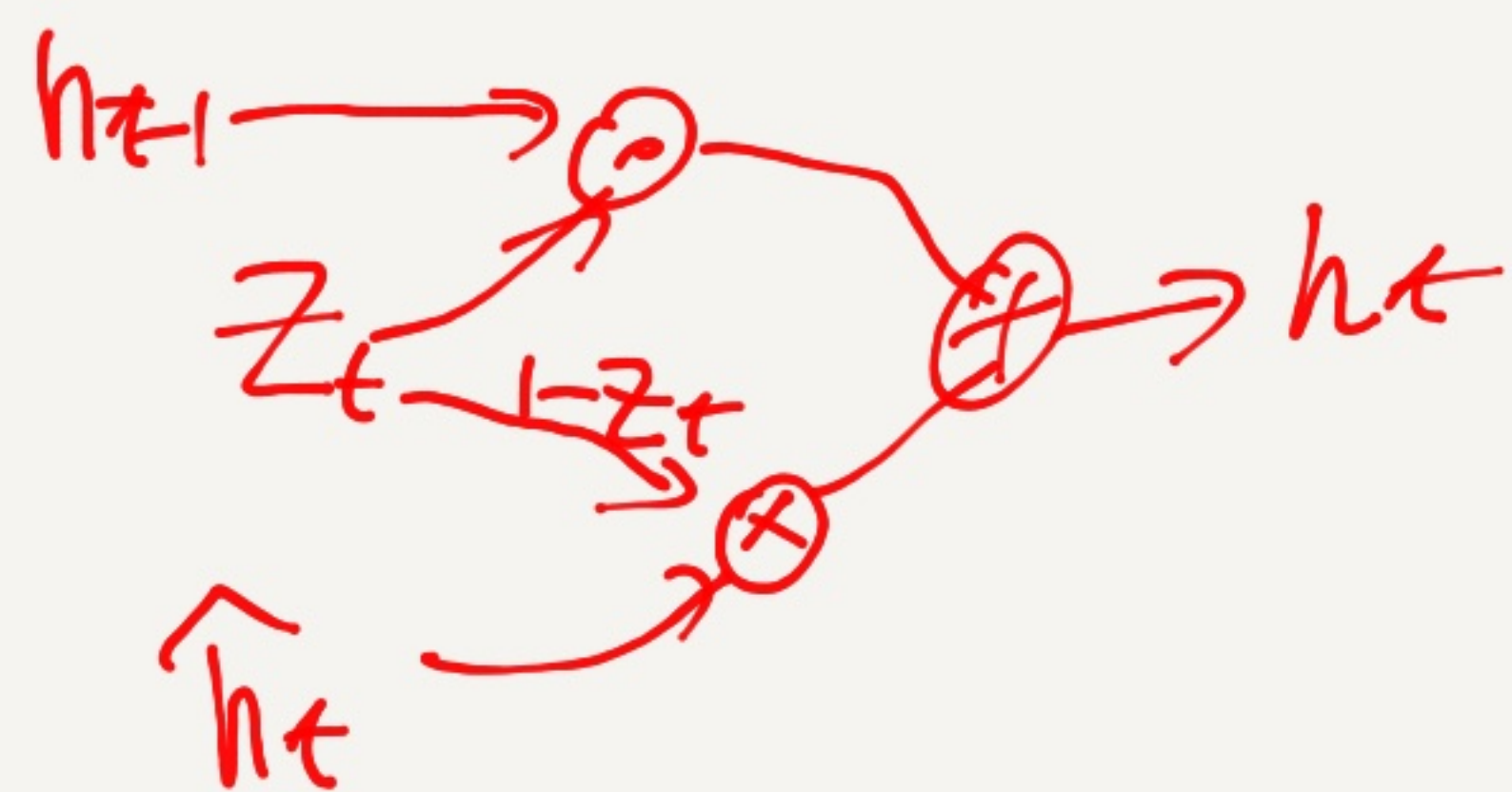
$$z_t = g_z (w_z \cdot h_{t-1} + u_z \cdot x_t)$$

$g_z$ : sigmoid  $\rightarrow [0, 1]$



update gate.

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \hat{h}_t$$



$\odot$ : element-wise multiplication.

Candidate:  $\hat{h}_t = g_h (w_h (\underline{x_t \odot h_{t-1}}) + u_h \cdot x_t)$

$r_t$ :

$$r_t = g_r (w_r \cdot h_{t-1} + u_r \cdot x_t)$$

$g_r$ : Sigmoid  $\rightarrow [0, 1]$

Reset gate