

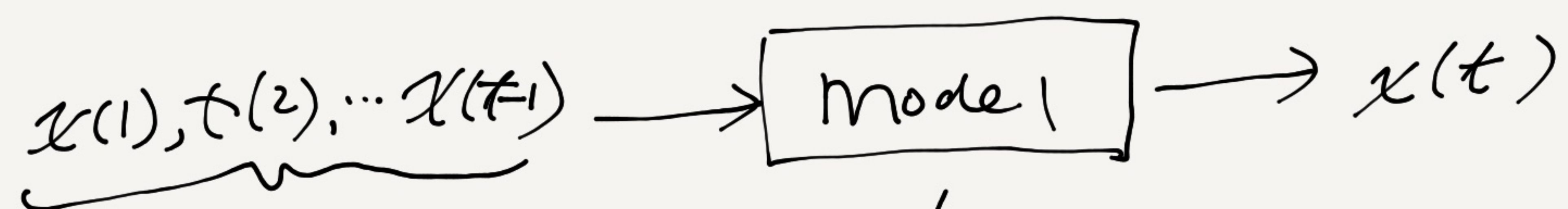
Recurrent Neural Networks (RNNs)

1. Models predict the future from the past.

Input: sequence: $x(1), x(2), x(3), \dots, x(t), \dots, x(T)$

t : timepoint

$x(t)$: t -th input value



↳ model is a dynamic model.

$f(x, t)$

2. Applications

- ① voice recognition / speech recognition
 - ② machine translation.
 - ③ weather prediction.
 - ④ stock market prediction
- } →

3. models

▷ Markov models (Markov chain, Hidden Markov model, (HMM)

Markov random field (MRF))

Sequence: $x(1), x(2), \dots, x(T)$

$$p(x(t) | x(1), x(2), \dots, x(t-1)) \approx \underbrace{p(x(t) | x(t-1))}_{\text{frequency of word in a sentence}} \quad x(t-1) \longrightarrow x(t)$$
$$= \frac{\text{count}(x(t-1), x(t))}{\text{count}(x(t-1))}$$

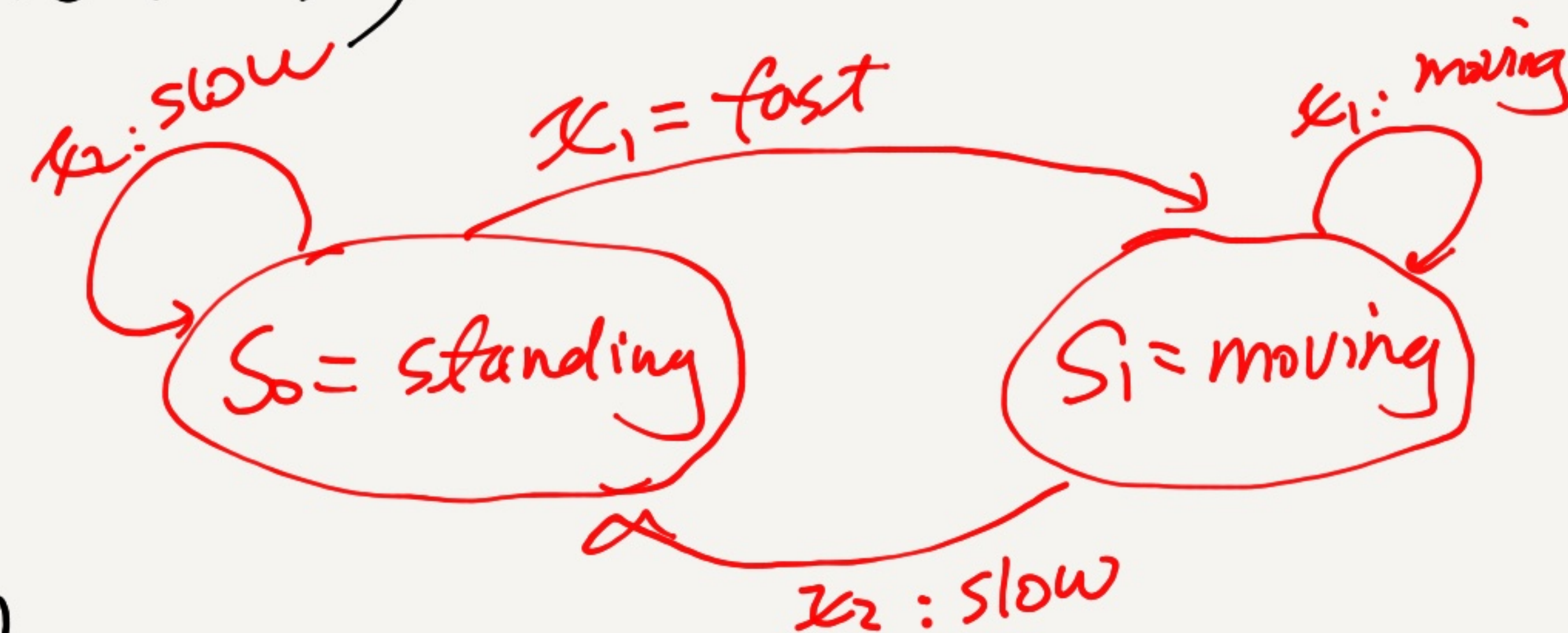
2) State machines

Set of states: $S = \{s_0, s_1, \dots, s_m\}$

Set of input: $X = \{x_1, x_2, \dots, x_n\}$

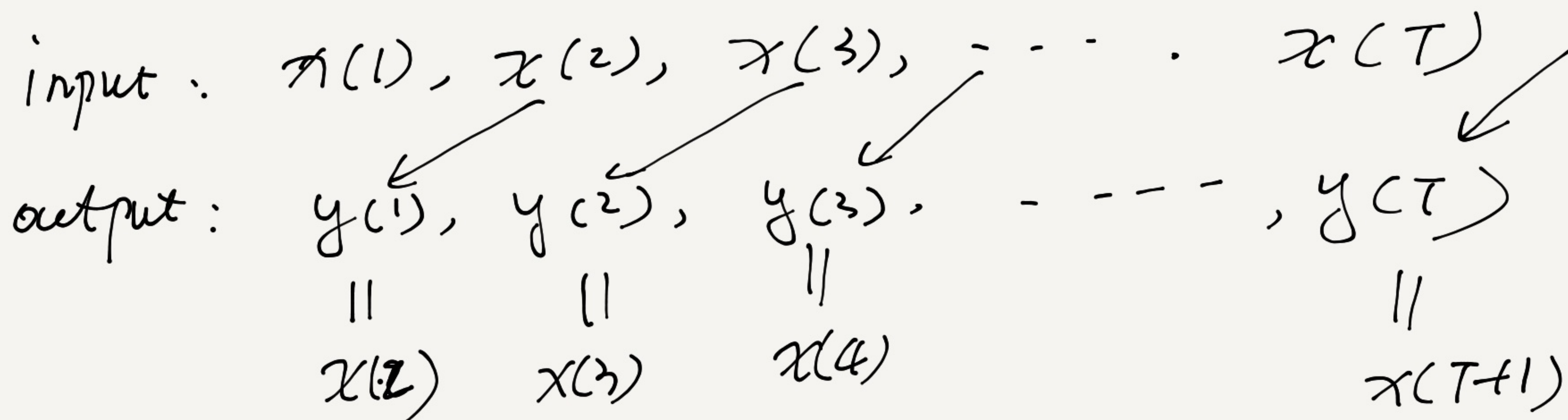
State transition: $f: S \times X \rightarrow S: f(s_{t-1}, x_t) = s_t$

Output: $g: S \rightarrow Y: y_t = g(s_t)$



RNNs are state machines:

4. How can we design a NN to predict the future?



① use one timepoint to predict next timepoint

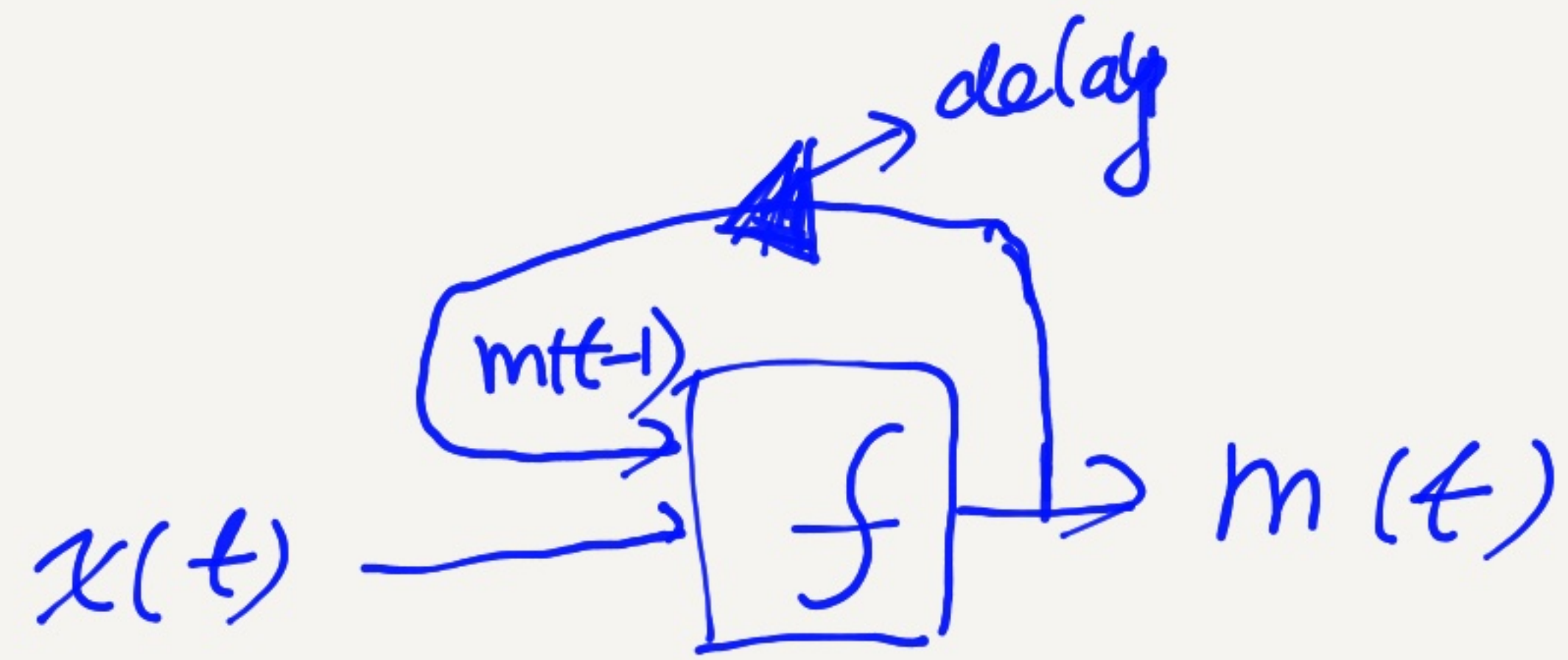
$$x(t-1) \rightarrow \boxed{f} \rightarrow y(t-1) = x(t)$$

⚡: This only model local context; sensitive to noise

⚡: This is static model. It has no long-range / long-term dependence.

② $x(1), x(2), \dots, x(t) \rightarrow \boxed{f(x)} \rightarrow y(t-1)$

③ Design RNNs based on the idea of moving average.



$m(t-1)$: encodes the past sequence.

input sequence. moving average

x_1

$$m_1 = x_1$$

x_2

$$m_2 = \frac{x_1 + x_2}{2} = \frac{m_1}{2} + \frac{x_2}{2}$$

x_3

$$m_3 = \frac{2m_2 + x_3}{3} = \frac{2}{3}m_2 + \frac{x_3}{3}$$

\vdots

x_t

current state

state

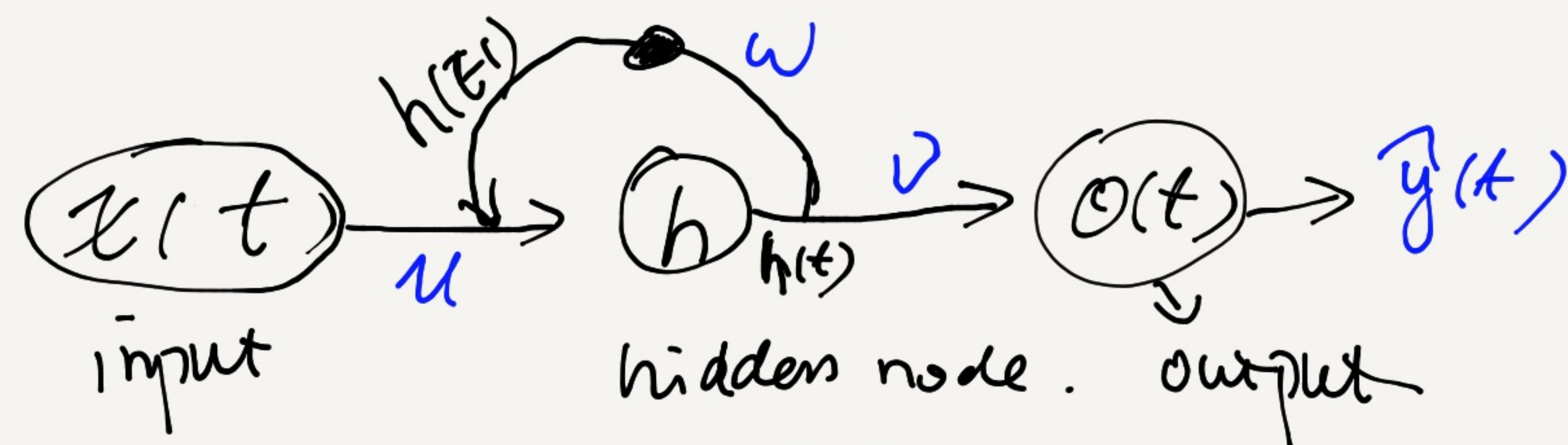
$$m_t = \frac{(t-1) \cdot m_{t-1}}{t} + \frac{x_t}{t}$$

5. RNN concepts.

1) RNN is a family of NNs for processing sequential data.

Both the input and output are sequences.

2) standard RNN.



$$h(t) = g_H(\text{activation})$$

$$\text{net}_H = \underbrace{u^T x(t)}_{\text{input}} + \underbrace{w^T \cdot h(t-1)}_{\text{bias}} + b_h$$

$$\hat{y}(t) = g_o(\text{activation})$$

The feed back connection gives the network the ability to model long-term dependence.

