# LSTM and variants
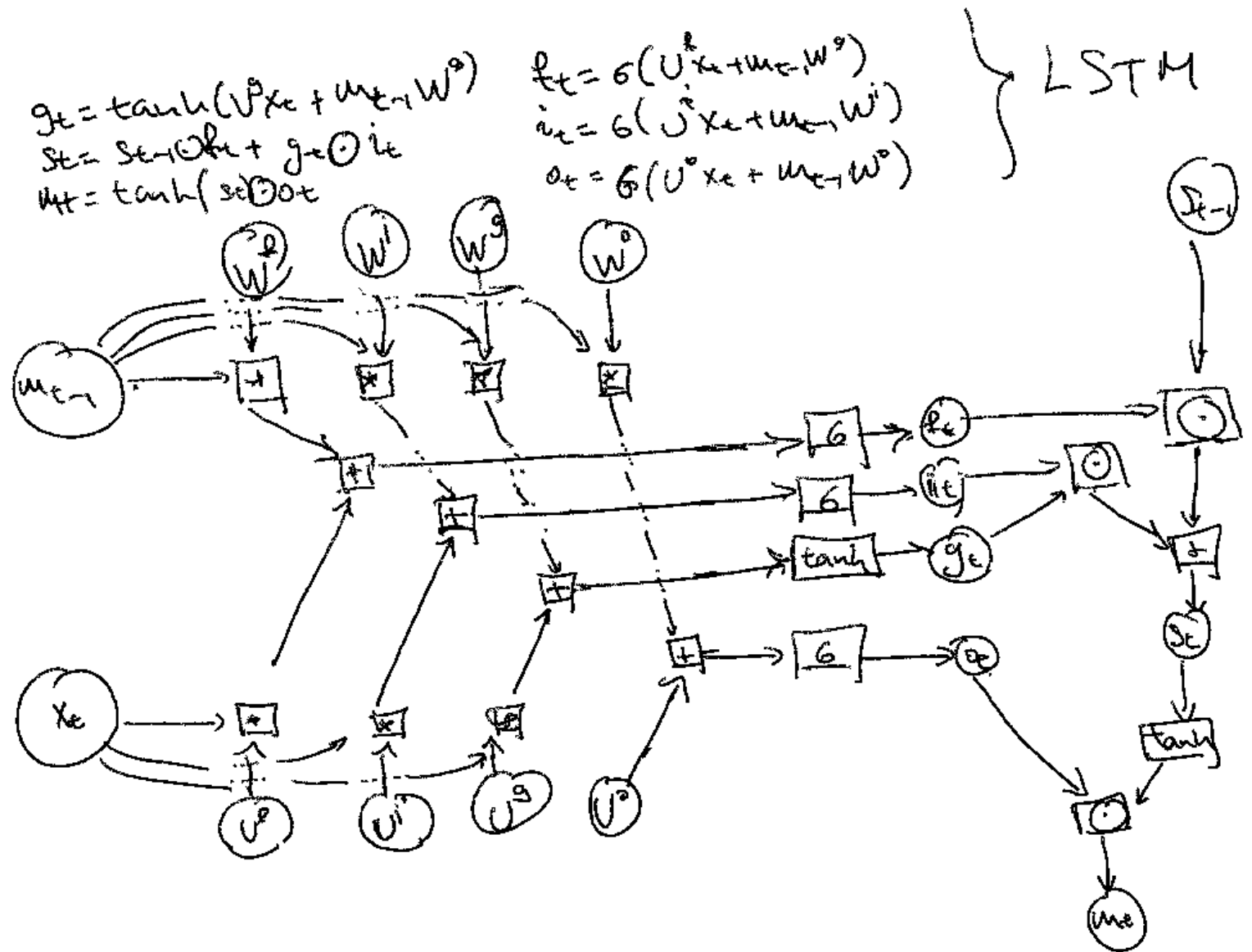
# Recap: vanishing gradients in RNNs

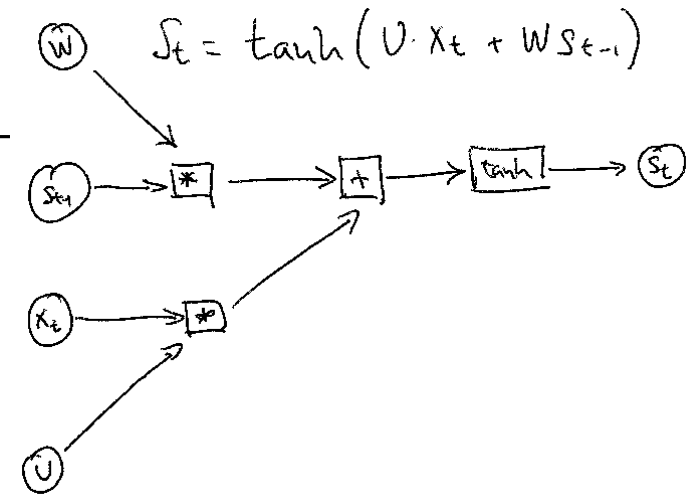o RNN equations

$$s_t = \tanh(U \cdot x_t + W \cdot s_{t-1})$$

o The RNN gradient of the memory parameters

$$\frac{\partial \mathcal{L}_t}{\partial W} = \sum_{i=0}^{t} \frac{\partial \mathcal{L}_t}{\partial y_t} \frac{\partial y_t}{\partial s_t} \left( \prod_{j=i+1}^{t} \frac{\partial s_j}{\partial s_{j-1}} \right) \frac{\partial s_i}{\partial W}$$

o Chain multiplications: $\frac{\partial s_j}{\partial s_{j-1}} < 1 \rightarrow$ vanishing gradients

# RNN → LSTM: Key idea

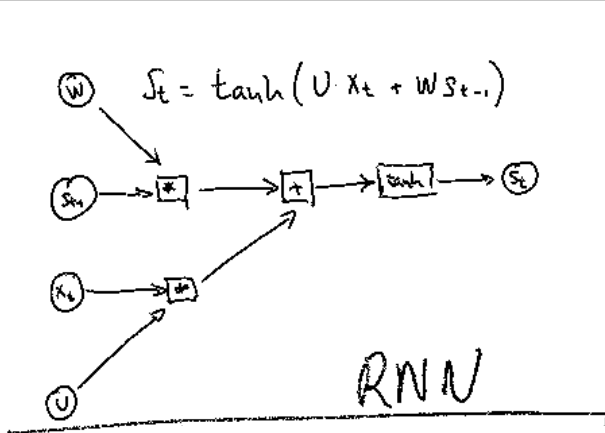$s_t = \tanh(V \cdot x_t + W s_{t-1})$

$$y_t = \text{softmax}(V \cdot s_t)$$
$$s_t = \tanh(U \cdot x_t + W \cdot s_{t-1})$$

o Setting $\frac{\partial s_j}{\partial s_{j-1}} = 1 \rightarrow$ no vanishing and exploding gradients

o Remove immediate nonlinear relation between $s_t$ and $s_{t-1}$
  ◦ Replace tanh between $s_t$ and $s_{t-1}$ with identity

o Also, avoid continuous overwriting of state
  ◦ Modulate the importance of new input by a gate
  ◦ Modulate the importance of new output by a gate
  ◦ Modulate the importance of past memories by a gate

# Designing an LSTM



$$s_t = \tanh(U \cdot x_t + W s_{t-1})$$
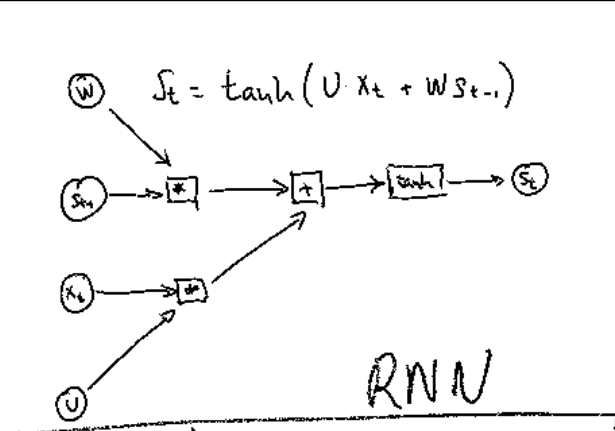
RNN

$$g_t = U \cdot x_t + W s_{t-1}$$

$$s_t = \tanh(g_t)$$

1. Drop nonlinearity between $s_t$ and $s_{t-1}$

$$g_t = U \cdot x_t + W s_{t-1}$$

$$s_t = g_t$$

# Designing an LSTM



$$s_t = \tanh(U \cdot x_t + W s_{t-1})$$

RNN

$$g_t = U \cdot x_t + W s_{t-1}$$

$$s_t = \tanh(g_t)$$

1. Drop nonlinearity between $s_t$ and $s_{t-1}$
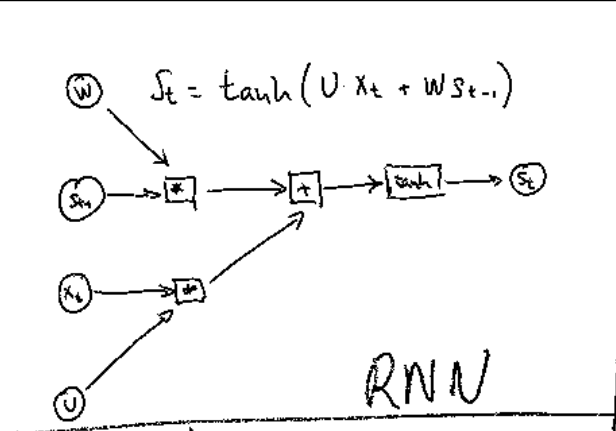
$$g_t = U \cdot x_t + W s_{t-1}$$

$$s_t = g_t$$

2. Now there is no nonlinearity in the neural network. Not good.
Add a nonlinearity on the term that does not depend on $s_{t-1}$

$$g_t = \tanh(U x_t)$$

$$s_t = s_{t-1} + g_t$$

# Designing an LSTM

$$s_t = \tanh(U \cdot x_t + W s_{t-1})$$



RNN

$$g_t = U \cdot x_t + W s_{t-1}$$

$$s_t = \tanh(g_t)$$

1. Drop nonlinearity between $s_t$ and $s_{t-1}$

$$g_t = U \cdot x_t + W s_{t-1}$$

$$s_t = g_t$$

2. Now there is no nonlinearity in the neural network. Not good.
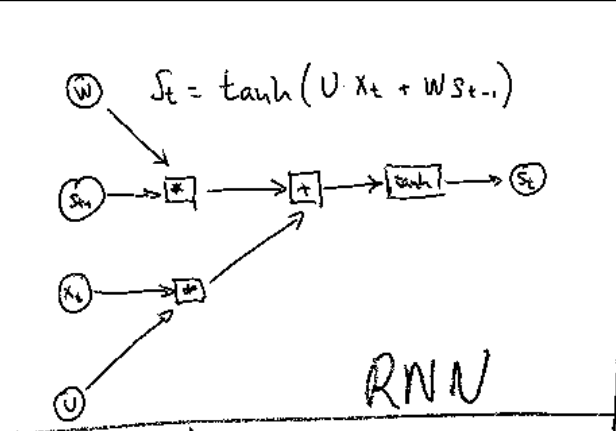Add a nonlinearity on the term that does not depend on $s_{t-1}$

$$g_t = \tanh(U x_t)$$

$$s_t = s_{t-1} + g_t$$

3. No vanishing gradients, but also no way to moderate input/memory
Add different gates to moderate/attend with $(0,1)$ values the key variables

$$g_t = \tanh(U^g x_t)$$

$$s_t = s_{t-1} \odot f_t + g_t \odot i_t$$

$$f_t = \sigma(U^f x_t)$$

$$i_t = \sigma(U^i x_t)$$

# Designing an LSTM

$S_t = \tanh(U \cdot X_t + W S_{t-1})$

RNN

$g_t = U \cdot X_t + W S_{t-1}$

$S_t = \tanh(g_t)$

1. Drop nonlinearity between $S_t$ and $g_t$

$g_t = U \cdot X_t + W S_{t-1}$

$S_t = g_t$

2. Now there is no nonlinearity in the neural network. Not good. Add a nonlinearity on the term that does not depend on $S_{t-1}$

$g_t = \tanh(U X_t)$

$S_t = S_{t-1} + g_t$

3. No vanishing gradients, but also no way to moderate input/memory. Add different gates to modulate/attend with $(0,1)$ values the key variables

$g_t = \tanh(U^g X_t)$

$S_t = S_{t-1} \odot f_t + g_t \odot i_t$

$f_t = \sigma(U^f X_t)$

$i_t = \sigma(U^i X_t)$

4. All good, but the output plays no role in the modulation. Add an output variable
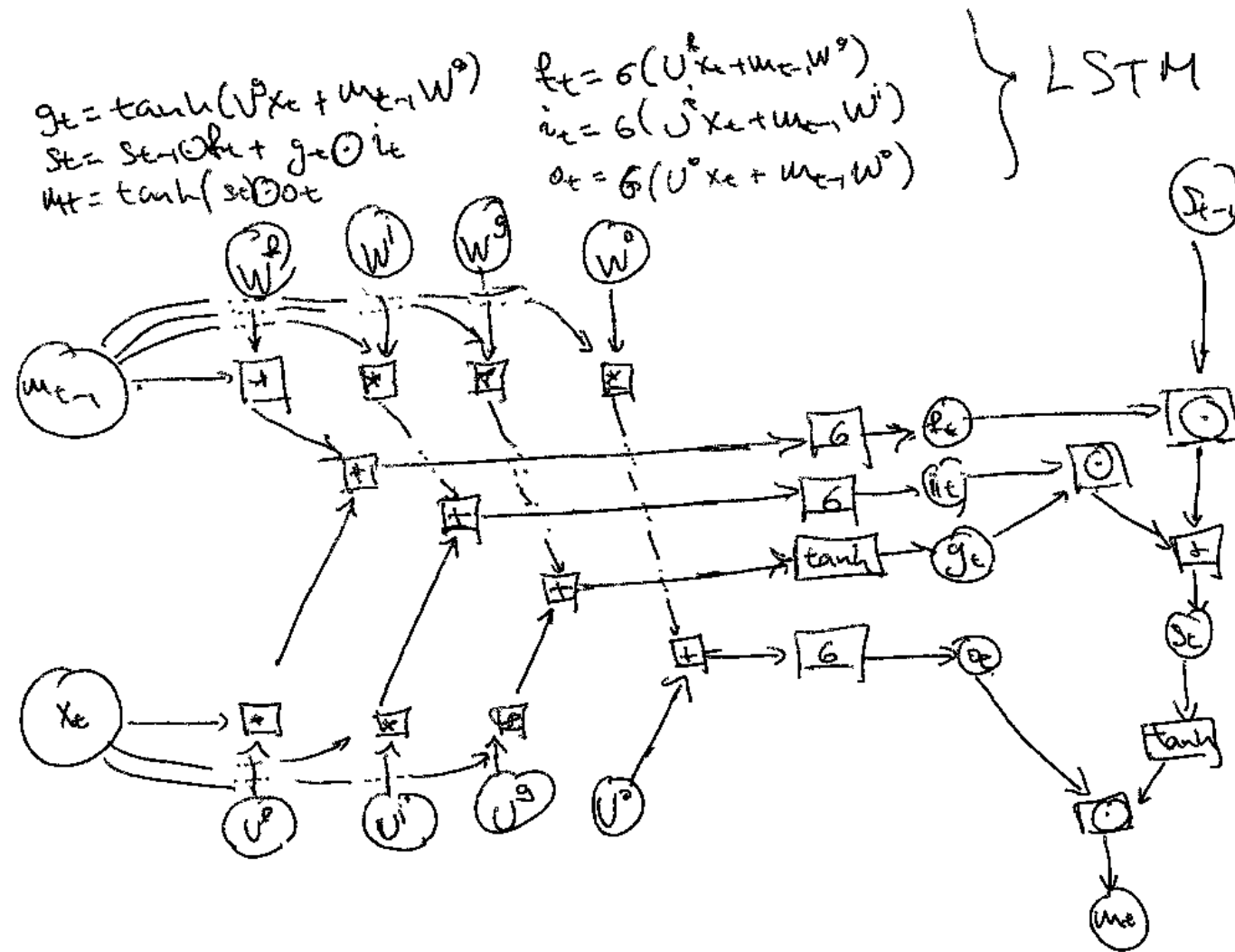
$g_t = \tanh(U^g X_t + m_{t-1} W^g)$

$S_t = S_{t-1} \odot f_t + g_t \odot i_t$

$m_t = \tanh(S_t) \odot o_t$

$f_t = \sigma(U^f X_t + m_{t-1} W^f)$

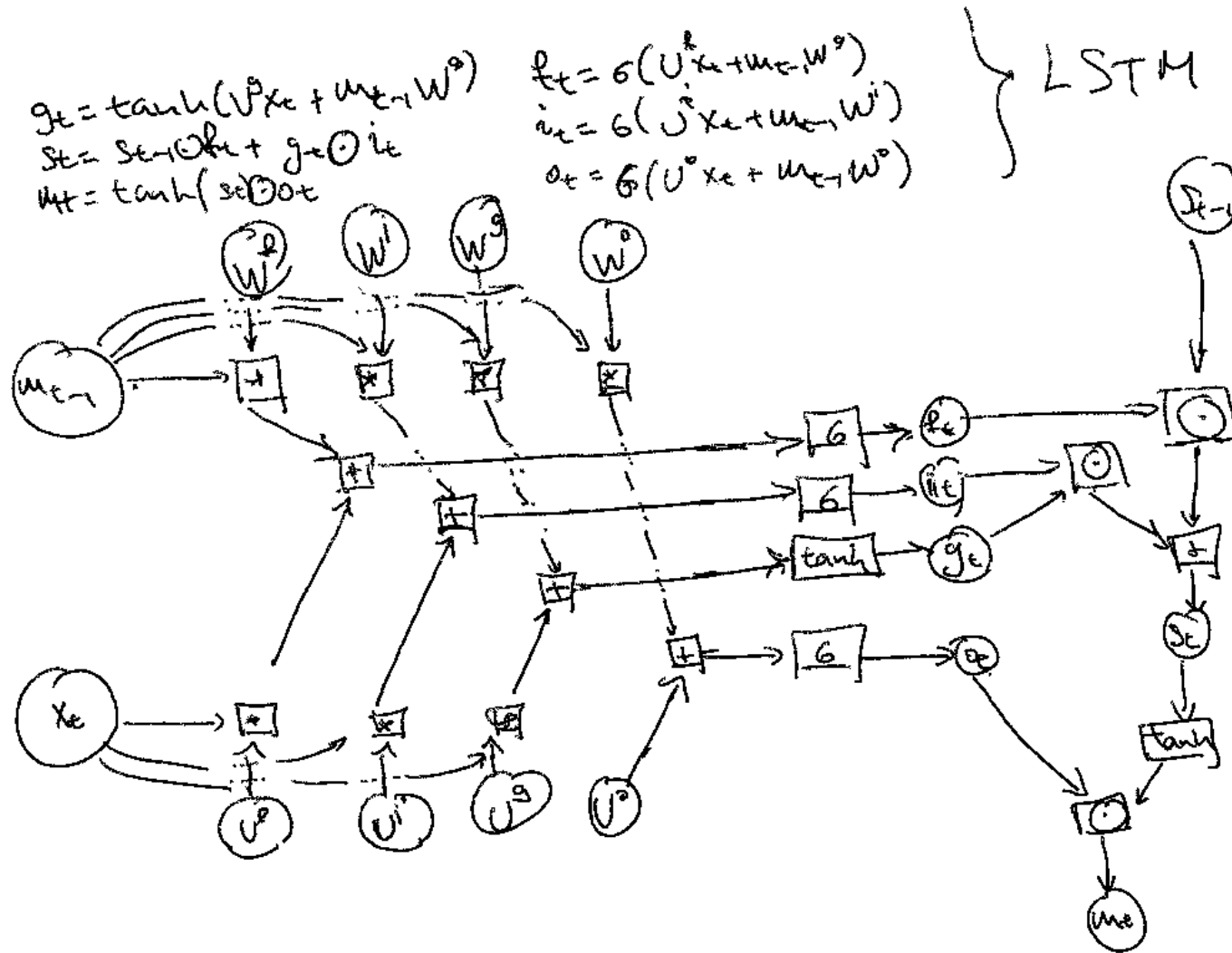$i_t = \sigma(U^i X_t + m_{t-1} W^i)$

$o_t = \sigma(U^o X_t + m_{t-1} W^o)$

LSTM

# An LSTM, graphically

# LSTM: Forward propagation step by step
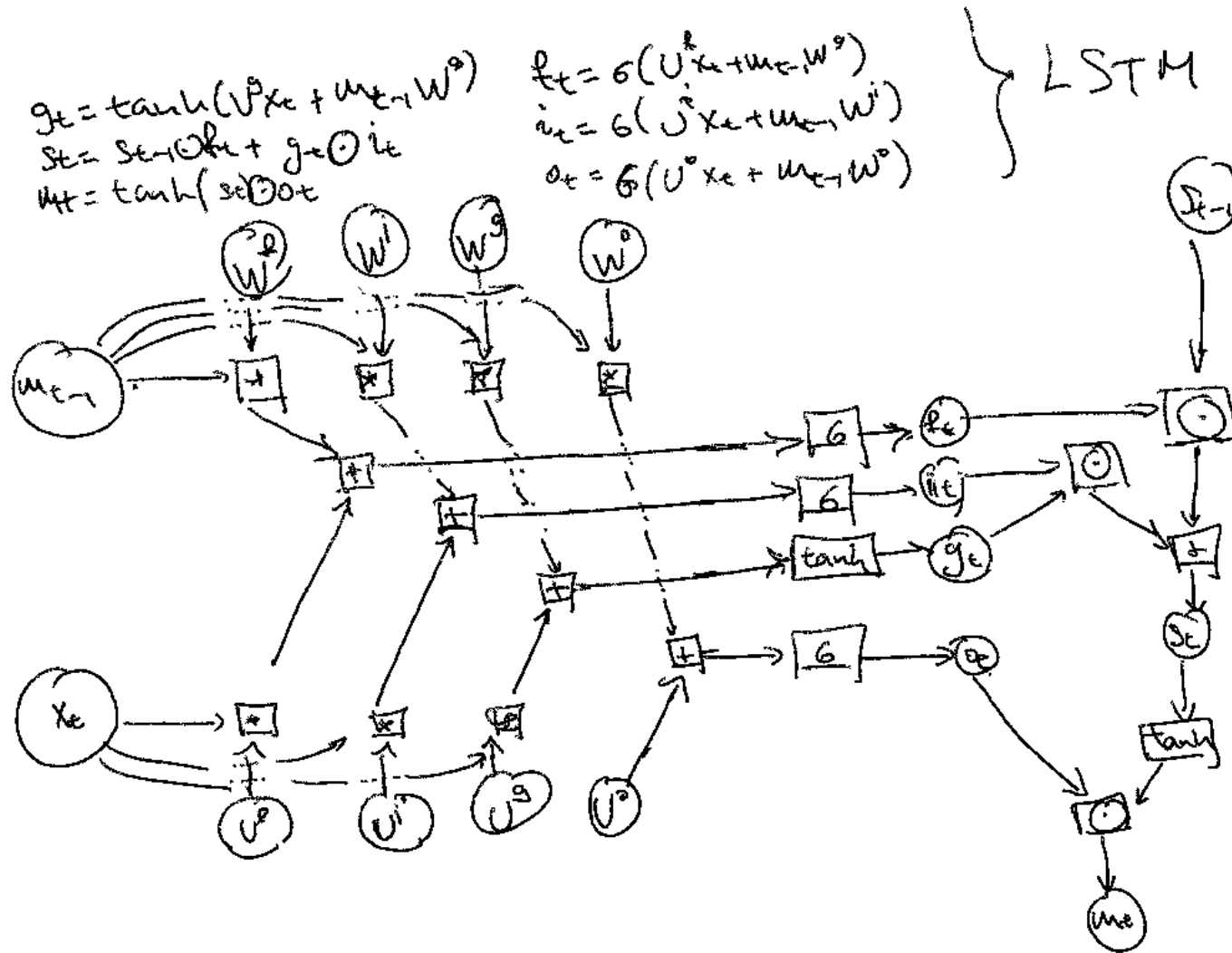


LSTM speaking

How important is my input?

LSTM speaking

How important is my past state?

LSTM speaking

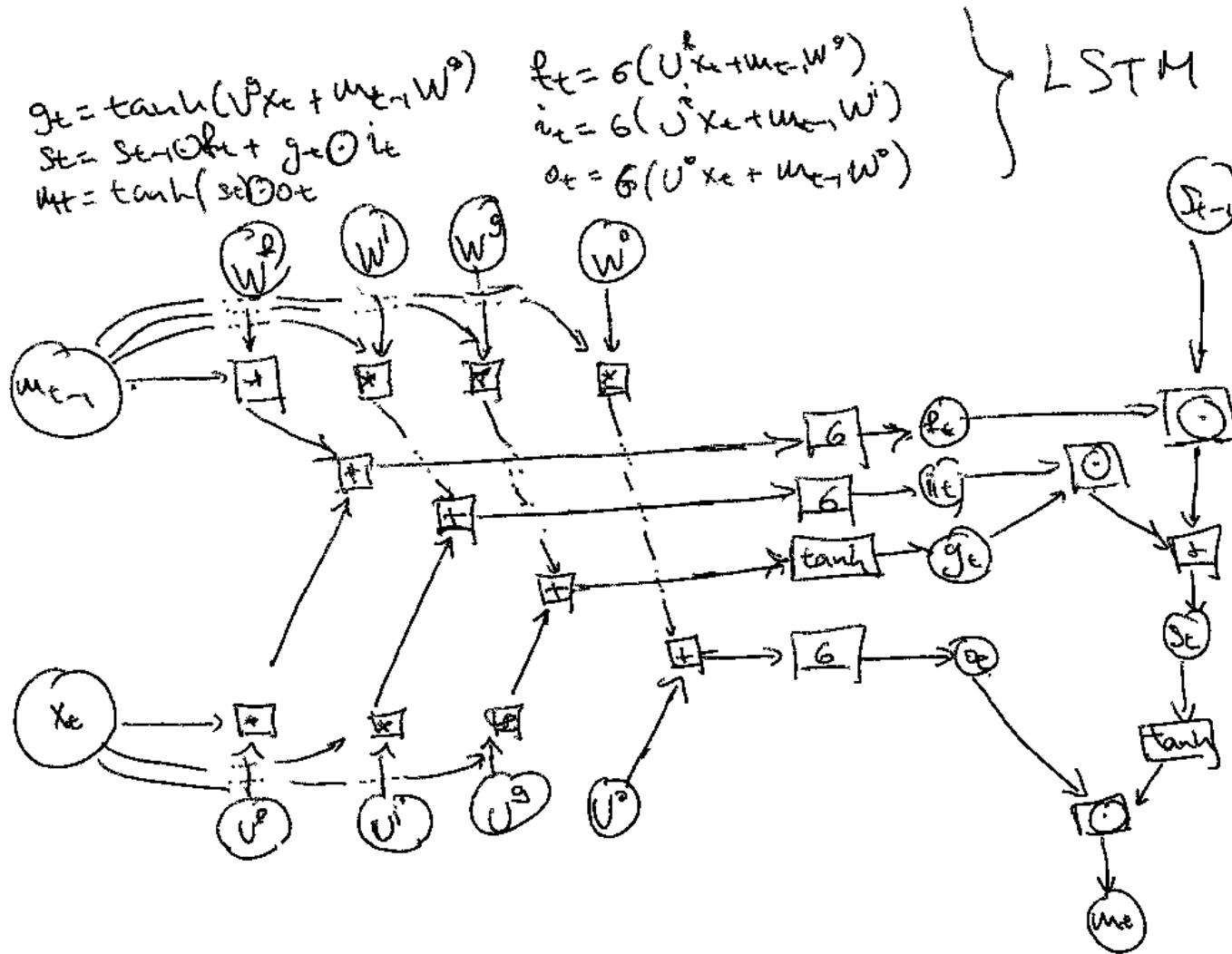What could be a relevant new memory?

LSTM speaking

Ok, let's compute my new state

## LSTM speaking

Is my new state useful for output? Check first how important is to give an output.
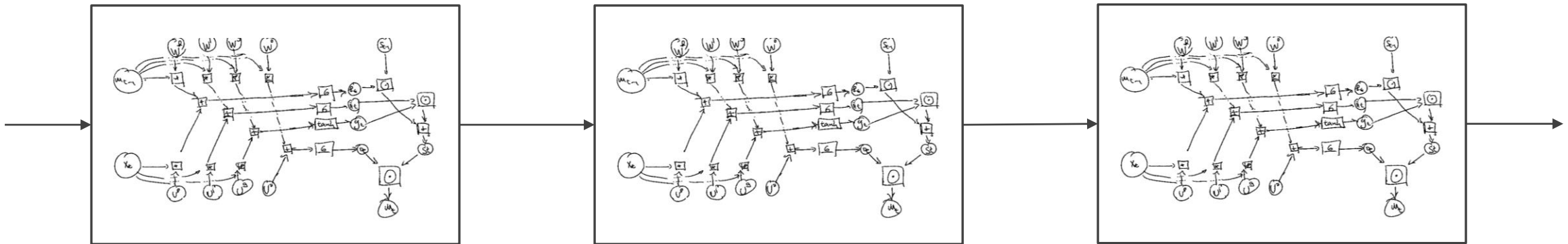
LSTM speaking

What is my new output?

# LSTM insights

o Comparing the state equations between RNN and LSTM
  ◦ RNN: $\boldsymbol{s}_t = \tanh(U \cdot \boldsymbol{x}_t + W \cdot \boldsymbol{s}_{t-1})$
  ◦ LSTM: $\boldsymbol{s}_t = \boldsymbol{s}_{t-1} \odot \boldsymbol{f}_t + \boldsymbol{g}_t \odot \boldsymbol{i}_t$, $\boldsymbol{m}_t = \tanh(\boldsymbol{s}_t) \odot \boldsymbol{o}_t$

o The LSTM also has indirect nonlinear relation between $\boldsymbol{s}_t$ and $\boldsymbol{s}_{t-1}$ via $\boldsymbol{m}_t$
  ◦ There is also direct linear relation → Strong gradients encouraged

o Use sigmoids for gating/squashing → (0,1) values
  ◦ Use tanh as module's recurrence nonlinearity, instead

Nice tutorial: http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Unrolling the LSTMs

o Just the same like for RNNs

o The internal equations are more complicated but the idea is the same

o Because of linear relation between $s_t$ and $s_{t-1}$ vanishing gradients mitigated
  ◦ LSTM captures short-term, as well as long-term correlations

# LSTM variants

o LSTM with peephole connections

o Gates have access also to the previous cell states $c_{t-1}$ (not only memories)

o Bi-directional recurrent networks

o Gated Recurrent Units (GRU)

o Phased LSTMs

o Skip LSTMs

o And many more …