

CS 316: Introduction to Deep Learning

Advanced Language Models – LSTM & GRU

Dr Abdul Samad

Gradients

- Long chain of dependencies for backprop
 - Need to keep a lot of intermediate values in memory
 - Butterfly effect style dependencies
 - Gradients can vanish or diverge (more on this later)
- Clipping to prevent divergence

$$\mathbf{g} \leftarrow \min \left(1, \frac{\theta}{\|\mathbf{g}\|} \right) \mathbf{g}$$

rescales to gradient of size at most θ

Recurrent Neural Networks (with hidden state)

- Hidden State update

$$h_t = f(h_{t-1}, x_{t-1}, w)$$

- Observation update

$$o_t = g(h_t, w)$$

Analysis of gradient in RNN

Consider a simple RNN

$$\begin{aligned}h_t &= f(x_t, h_{t-1}, w_h) \\ o_t &= g(h_t, w_o)\end{aligned}$$

Loss is calculated as

$$L(x_1, \dots, x_T, y_1, \dots, y_T, w_h, w_o) = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$

Analysis of gradient in RNN

$$\begin{aligned}\frac{\partial L}{\partial w_h} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial w_h} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\partial l(y_t, o_t)}{\partial o_t} \frac{\partial g(h_t, w_0)}{\partial h_t} \frac{\partial h_t}{\partial w_h}\end{aligned}$$

Analysis of gradient in RNN

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{w}_h} = \frac{\partial f(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{w}_h)}{\partial \mathbf{w}_h} + \frac{\partial f(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{w}_h)}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{w}_h}$$

Assume we have three sequences $\{a_t\}, \{b_t\}, \{c_t\}$ satisfying $a_0 = 0$ and $a_t = b_t + c_t a_{t-1}$ for $t = 1, 2, \dots$

$$a_t = b_t + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t c_j \right) b_i$$

$$a_t = \frac{\partial \mathbf{h}_t}{\partial \mathbf{w}_h}$$

$$b_t = \frac{\partial f(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{w}_h)}{\partial \mathbf{w}_h}$$

$$c_t = \frac{\partial f(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{w}_h)}{\partial \mathbf{h}_{t-1}}$$

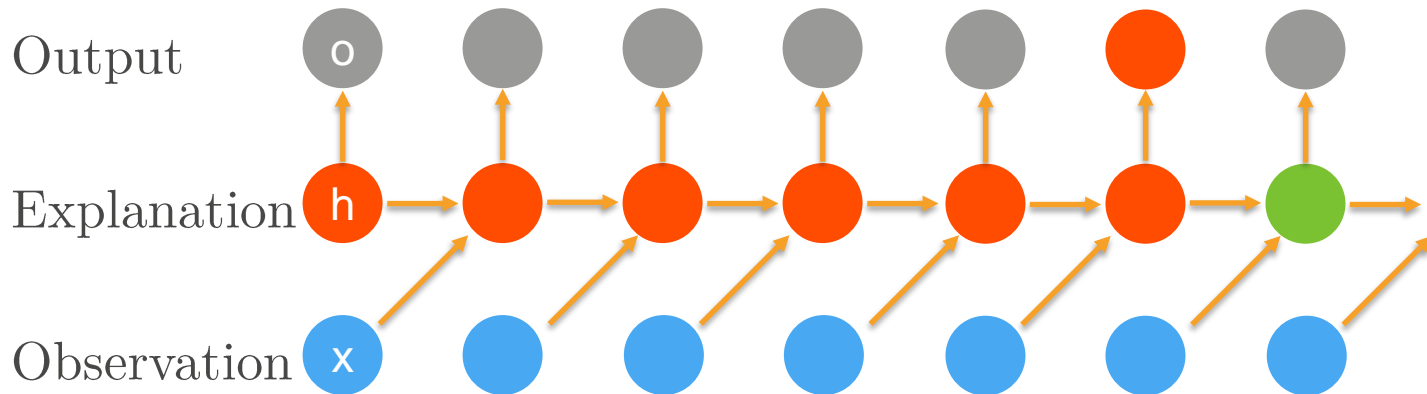
Analysis of gradient of RNN

$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$

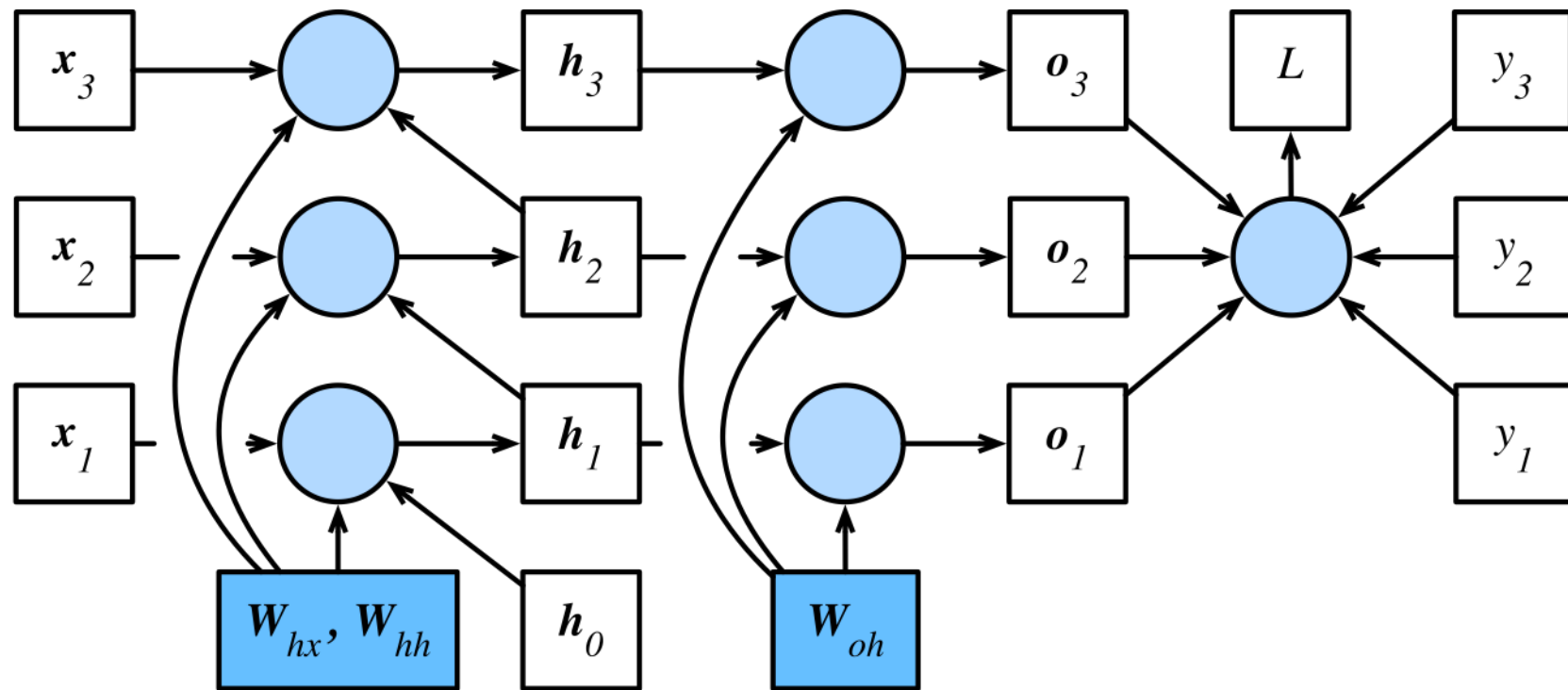
Latent State Gradient $\partial_w h_t$

- Gradient Recursion

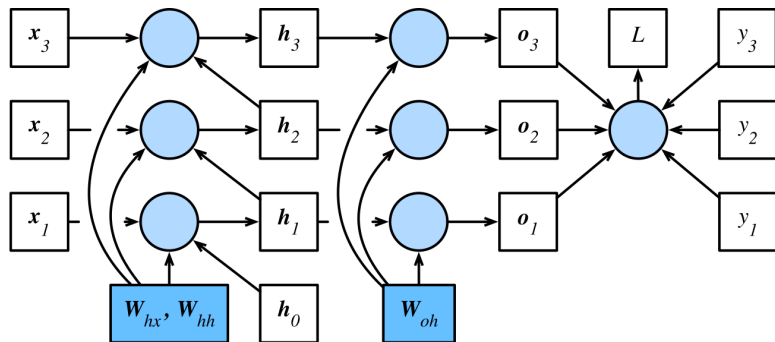
$$\frac{\partial h_t}{\partial w_h} = \frac{\partial f(x_t, h_{t-1}, w_h)}{\partial w_h} + \sum_{i=1}^{t-1} \left(\prod_{j=i+1}^t \frac{\partial f(x_j, h_{j-1}, w_h)}{\partial h_{j-1}} \right) \frac{\partial f(x_i, h_{i-1}, w_h)}{\partial w_h}$$



Computational Graph



Class Activity



$$\frac{\partial L}{\partial o_t} =$$

$$\frac{\partial L}{\partial w_{qh}} =$$

$$\frac{\partial L}{\partial h_T} =$$

$$\frac{\partial L}{\partial w_{hx}} =$$

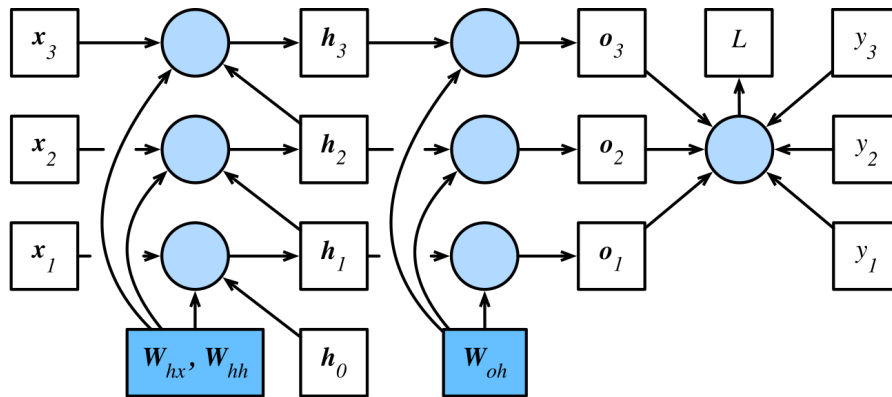
$$\frac{\partial L}{\partial w_{hh}} =$$

Toy Model

$$h_t = W_{hx}x_t + W_{hh}h_{t-1}$$

$$o_t = W_{qh}h_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(y_t, o_t)$$



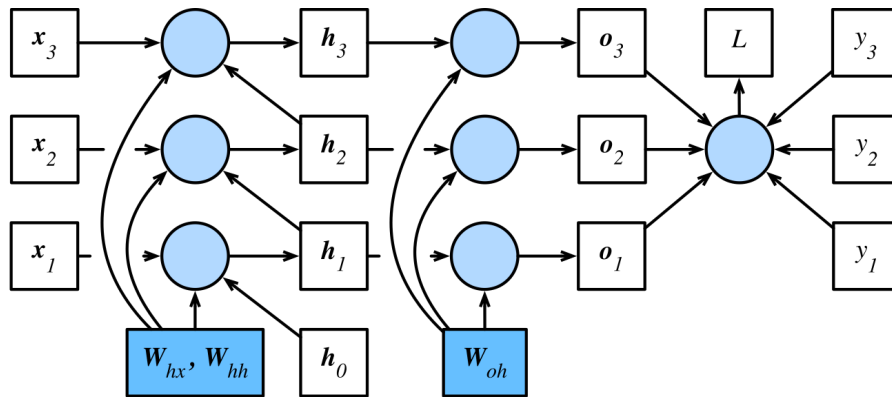
$$\frac{\partial L}{\partial o_t} = \frac{\partial l(o_t, y_t)}{T \cdot \partial o_t} \in R^q$$

Toy Model

$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{y}_t, \mathbf{o}_t)$$



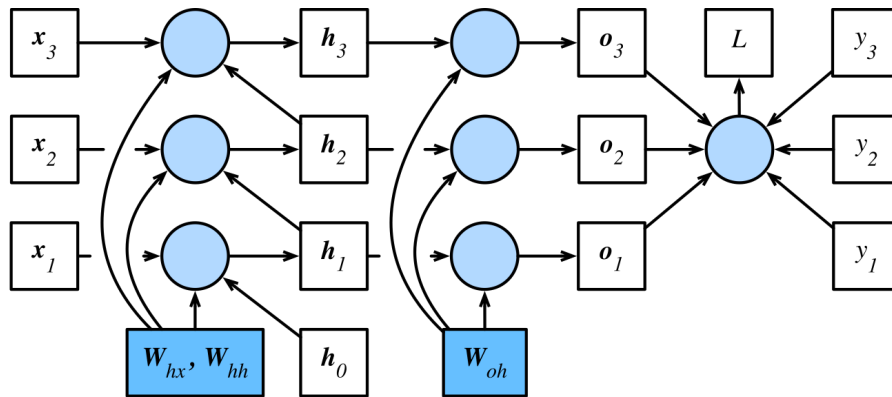
$$\frac{\partial L}{\partial \mathbf{w}_{qh}} = \sum_{t=1}^T \text{prod} \left(\frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{w}_{qh}} \right)$$

Toy Model

$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{y}_t, \mathbf{o}_t)$$



$$\frac{\partial L}{\partial \mathbf{h}_T} = \text{prod} \left(\frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_T} \right)$$

$$\frac{\partial L}{\partial \mathbf{h}_t} = \text{prod} \left(\frac{\partial L}{\partial \mathbf{h}_{t+1}}, \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_T} \right) + \text{prod} \left(\frac{\partial L}{\partial \mathbf{o}_t}, \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_T} \right)$$

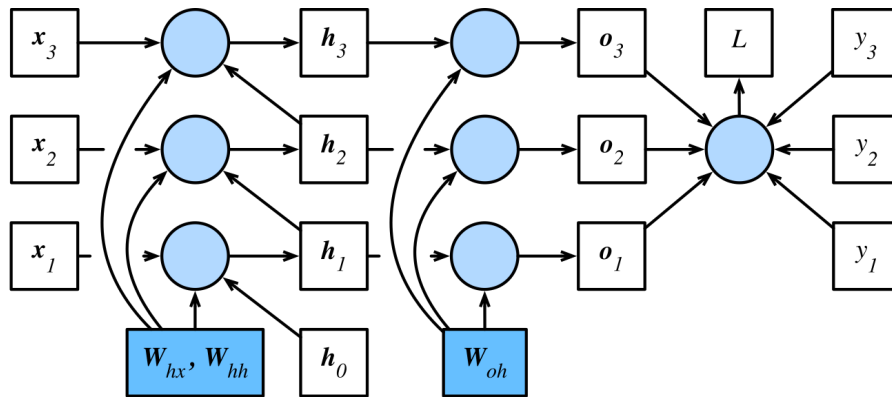
$$\frac{\partial L}{\partial \mathbf{h}_t} = \sum_{i=t}^T (\mathbf{W}_{hh}^T)^{T-i} \mathbf{W}_{qh}^T \frac{\partial L}{\partial \mathbf{o}_{T+t-i}} \quad 1 \leq t \leq T$$

Toy Model

$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1}$$

$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t$$

$$L = \frac{1}{T} \sum_{t=1}^T l(\mathbf{y}_t, \mathbf{o}_t)$$

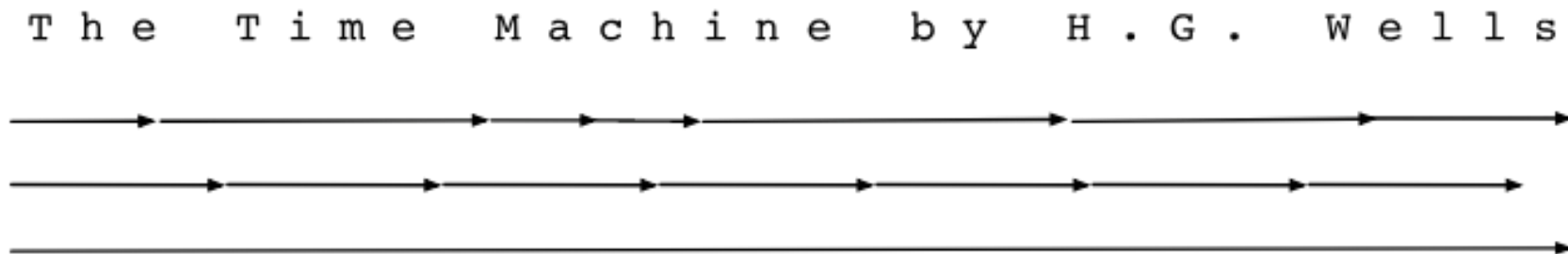


$$\frac{\partial L}{\partial W_{hx}} = \sum_{i=t}^T \text{prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hx}} \right)$$

$$\frac{\partial L}{\partial W_{hh}} = \sum_{i=t}^T \text{prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hh}} \right)$$

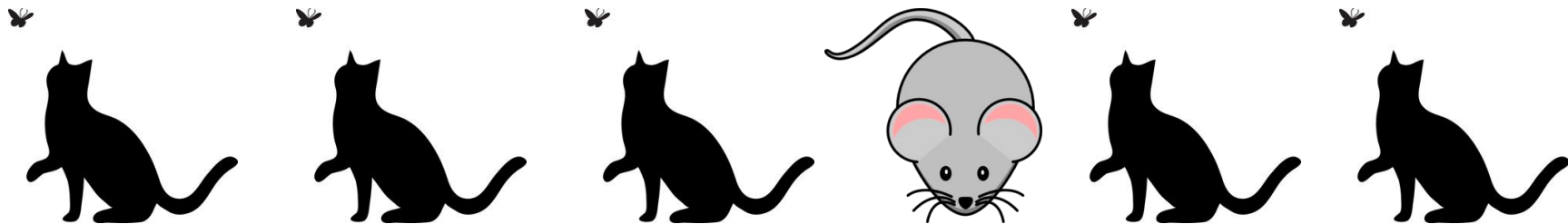
Truncated BPTT

- Don't truncate (naive strategy, costly and divergent)
- Truncate at fixed intervals
(standard approach, is approximation but works well)
- Variable length (Tallec and Olivier, 2015)



Paying attention to a sequence

- Not all observations are equally relevant

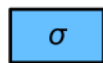
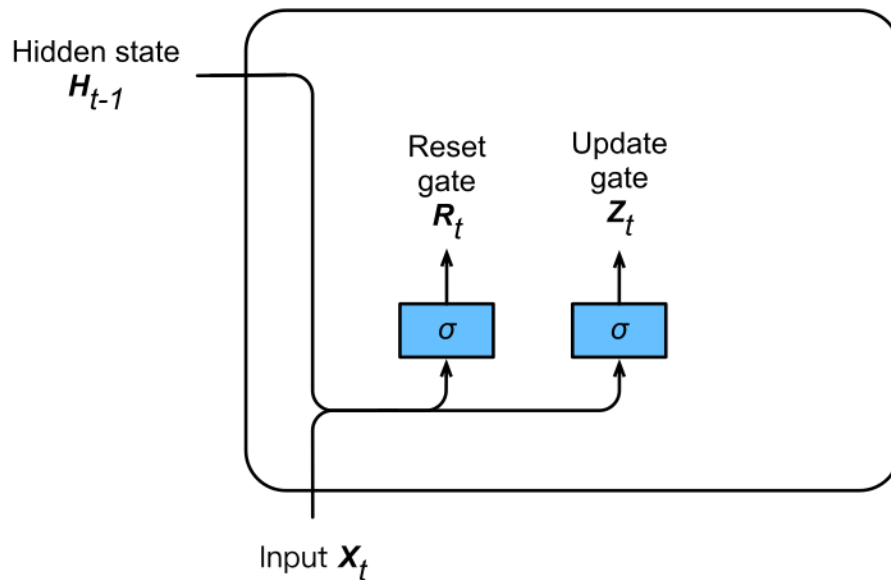


- Only remember the relevant ones
 - Need mechanism to pay attention (update gate)
 - Need mechanism to forget (reset gate)

Gating

$$\mathbf{R}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r),$$

$$\mathbf{Z}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z)$$



FC layer with
activation fuction



Element-wise
Operator



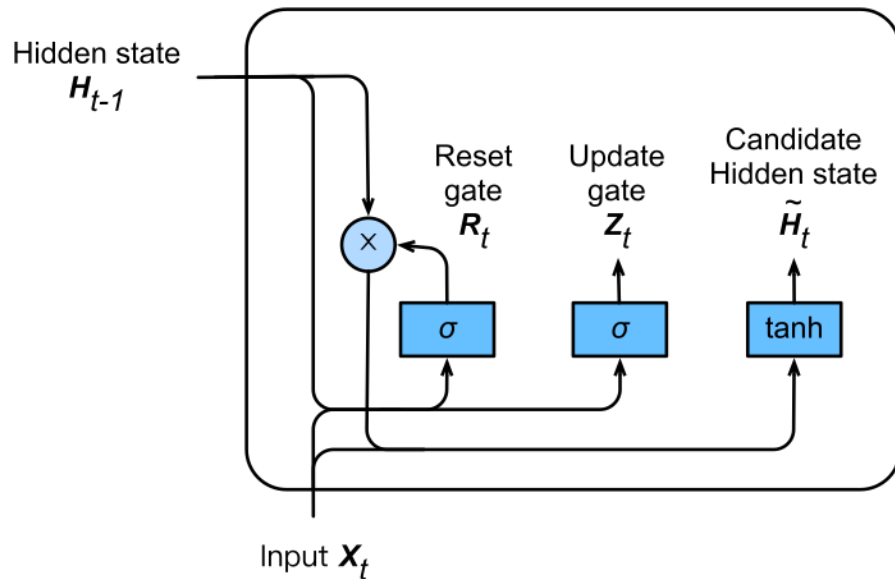
Copy



Concatenate

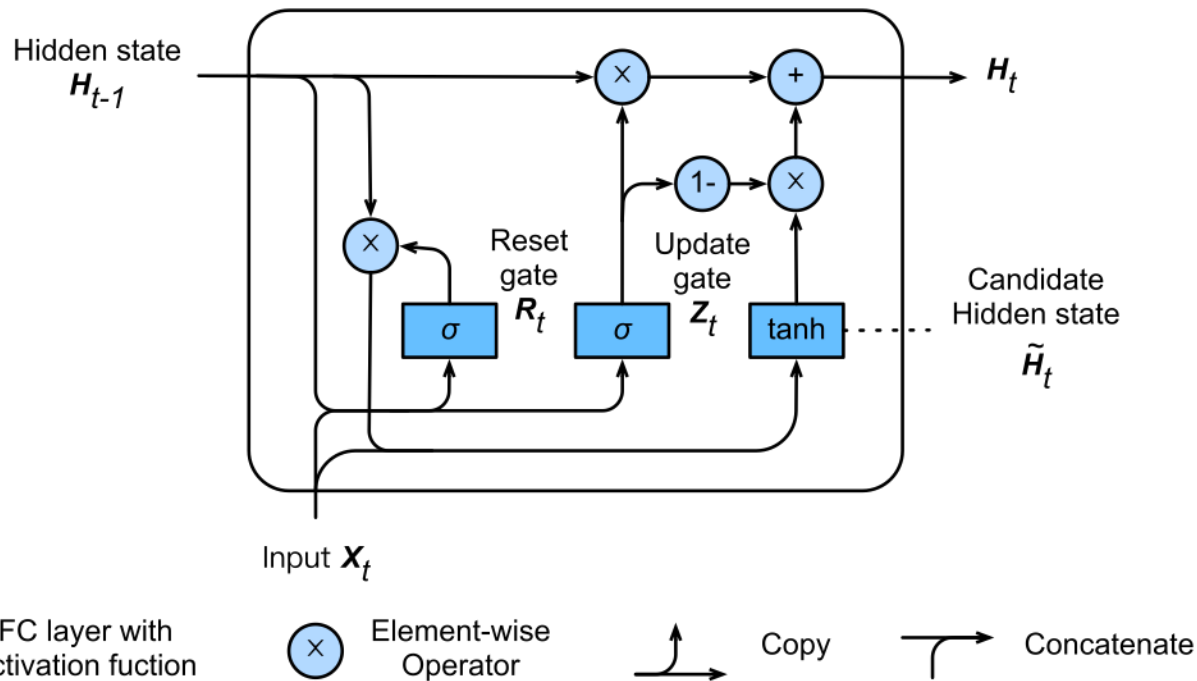
Candidate Hidden State

$$\tilde{H}_t = \tanh(X_t W_{xh} + (R_t \odot H_{t-1}) W_{hh} + b_h)$$



GRU-Hidden State

$$\mathbf{H}_t = \mathbf{Z}_t \odot \mathbf{H}_{t-1} + (1 - \mathbf{Z}_t) \odot \tilde{\mathbf{H}}_t$$



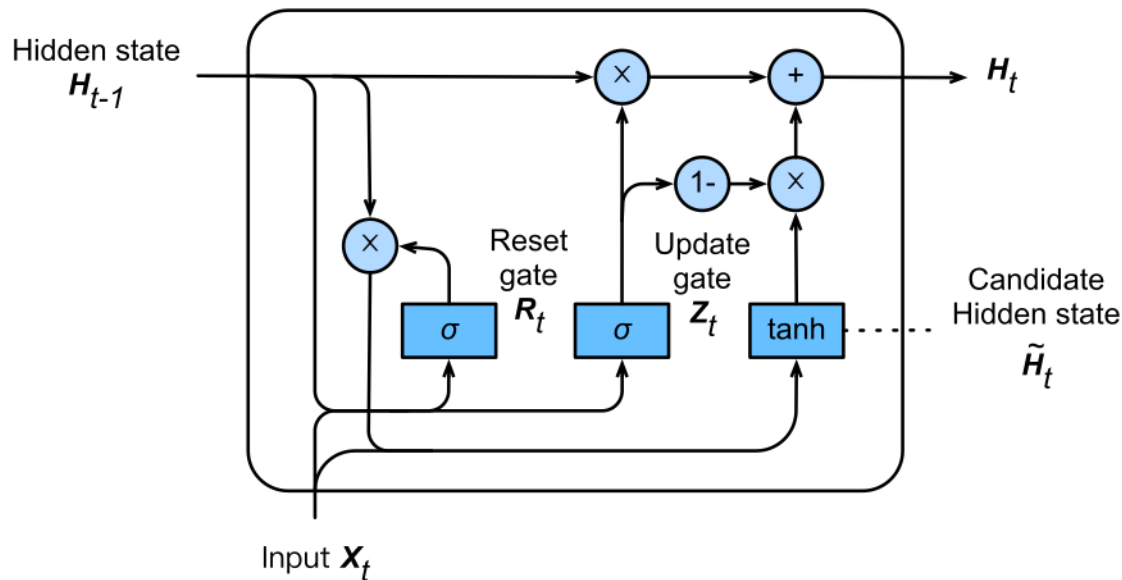
GRU-Summary

$$\mathbf{R}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r),$$

$$\mathbf{Z}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z)$$

$$\tilde{\mathbf{H}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xh} + (\mathbf{R}_t \odot \mathbf{H}_{t-1}) \mathbf{W}_{hh} + \mathbf{b}_h)$$

$$\mathbf{H}_t = \mathbf{Z}_t \odot \mathbf{H}_{t-1} + (1 - \mathbf{Z}_t) \odot \tilde{\mathbf{H}}_t$$



Long Short Term Memory-LSTM

- **Forget gate**

Shrink values towards zero

- **Input gate**

Decide whether we should ignore the input data

- **Output gate**

Decide whether the hidden state is used for the output generated by the LSTM

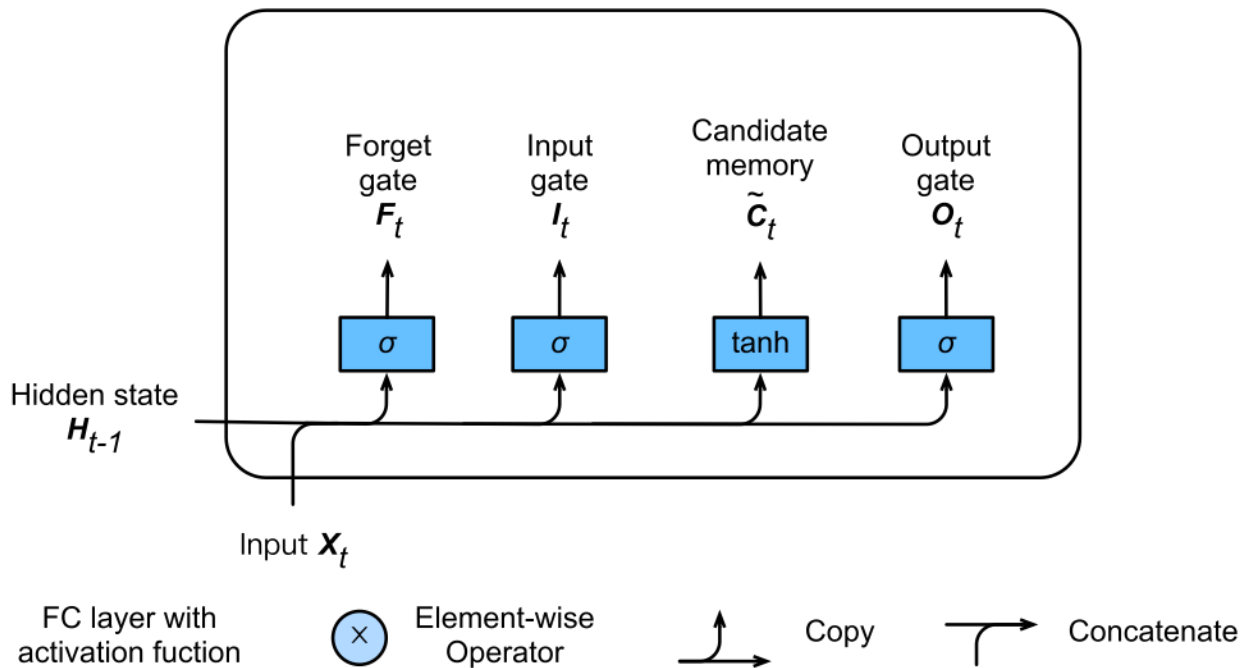
- **Hidden state and Memory cell**

Gates

$$I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i)$$

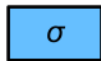
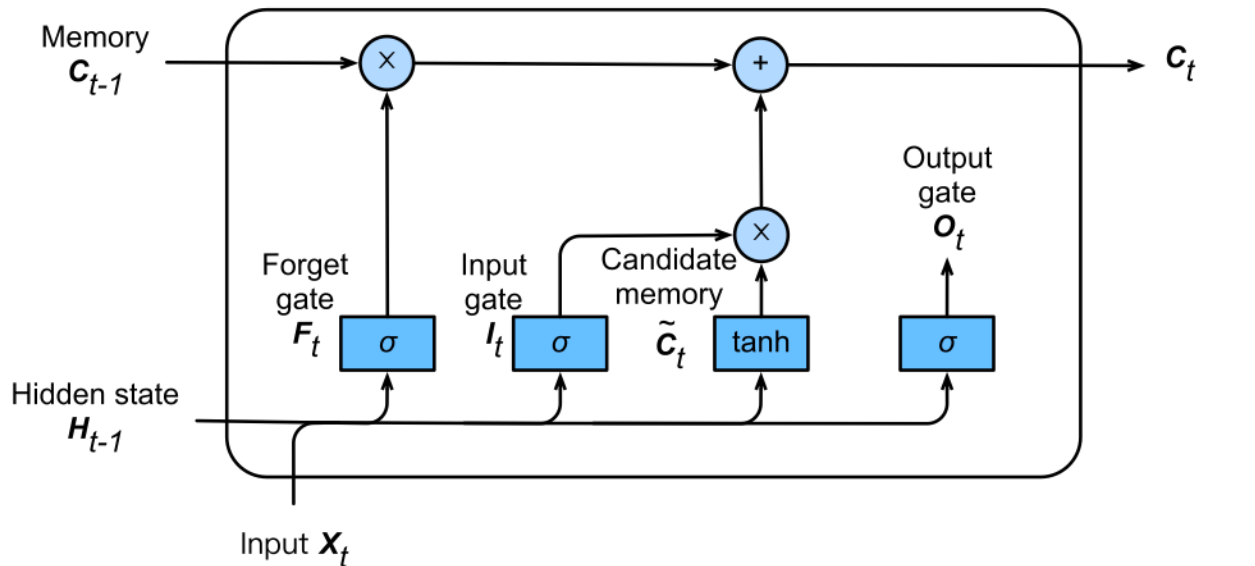
$$F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f)$$

$$O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o)$$



Candidate Memory Cell

$$\tilde{C}_t = \tanh(X_t W_{xc} + H_{t-1} W_{hc} + b_c)$$



FC layer with
activation fuction



Element-wise
Operator



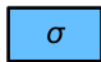
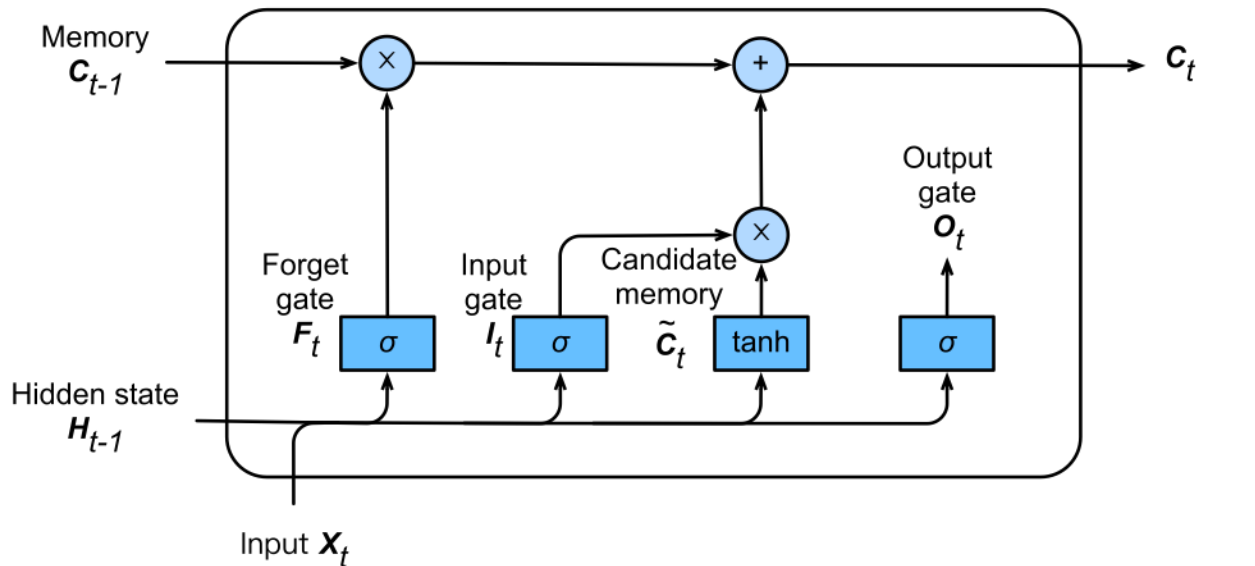
Copy



Concatenate

Memory Cell

$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$



FC layer with
activation fuction



Element-wise
Operator



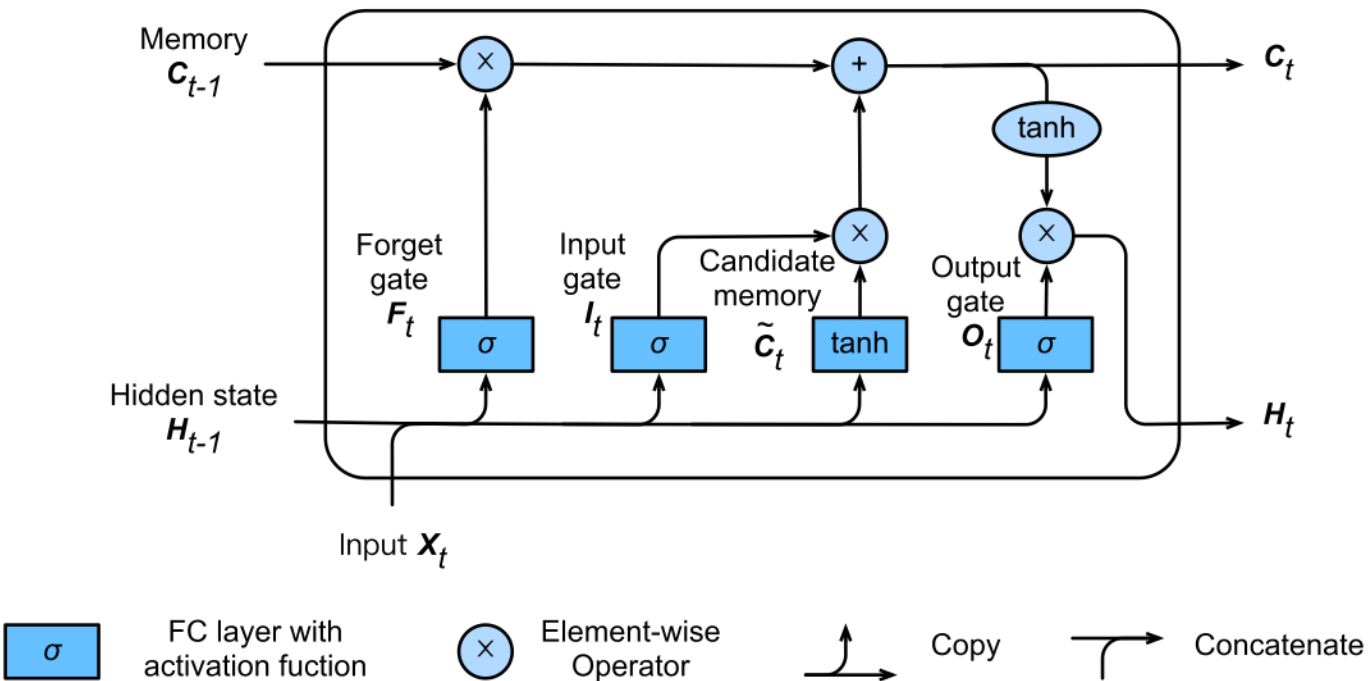
Copy



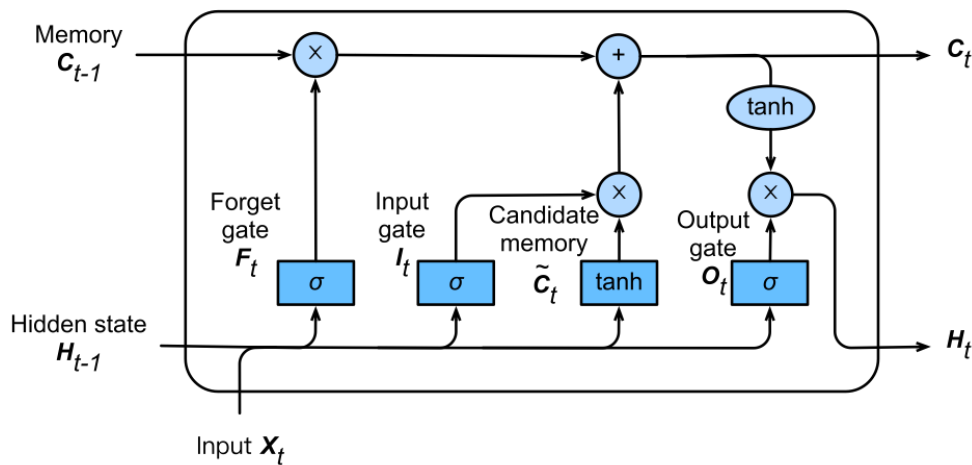
Concatenate

Hidden State / Output

$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t)$$



Hidden State / Output



$$I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i)$$

$$F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f)$$

$$O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o)$$

$$\tilde{C}_t = \tanh(X_t W_{xc} + H_{t-1} W_{hc} + b_c)$$

$$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$$

$$H_t = O_t \odot \tanh(C_t)$$