



Deep Learning for Healthcare

**Recurrent Neural
Networks (RNN)**

Jimeng Sun

Recurrent Neural Networks

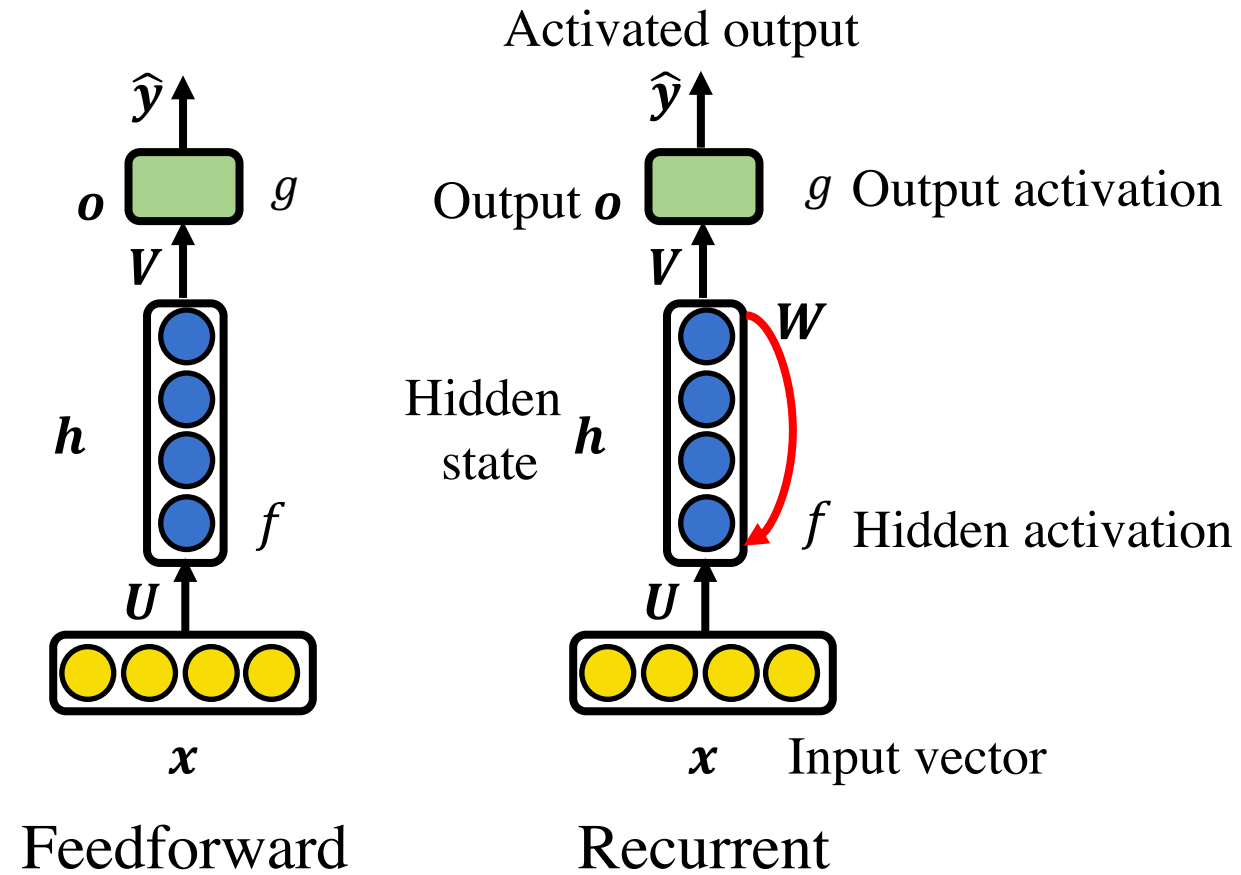
Agenda

- RNN Basics
- Learning RNN with Backpropagation Through Time (BPTT)
- Long-Short Term Memory Networks (LSTM)
- Gated Recurrent Unit (GRU)
- Bidirectional RNN
- Sequence-to-Sequence RNN
- Healthcare Applications

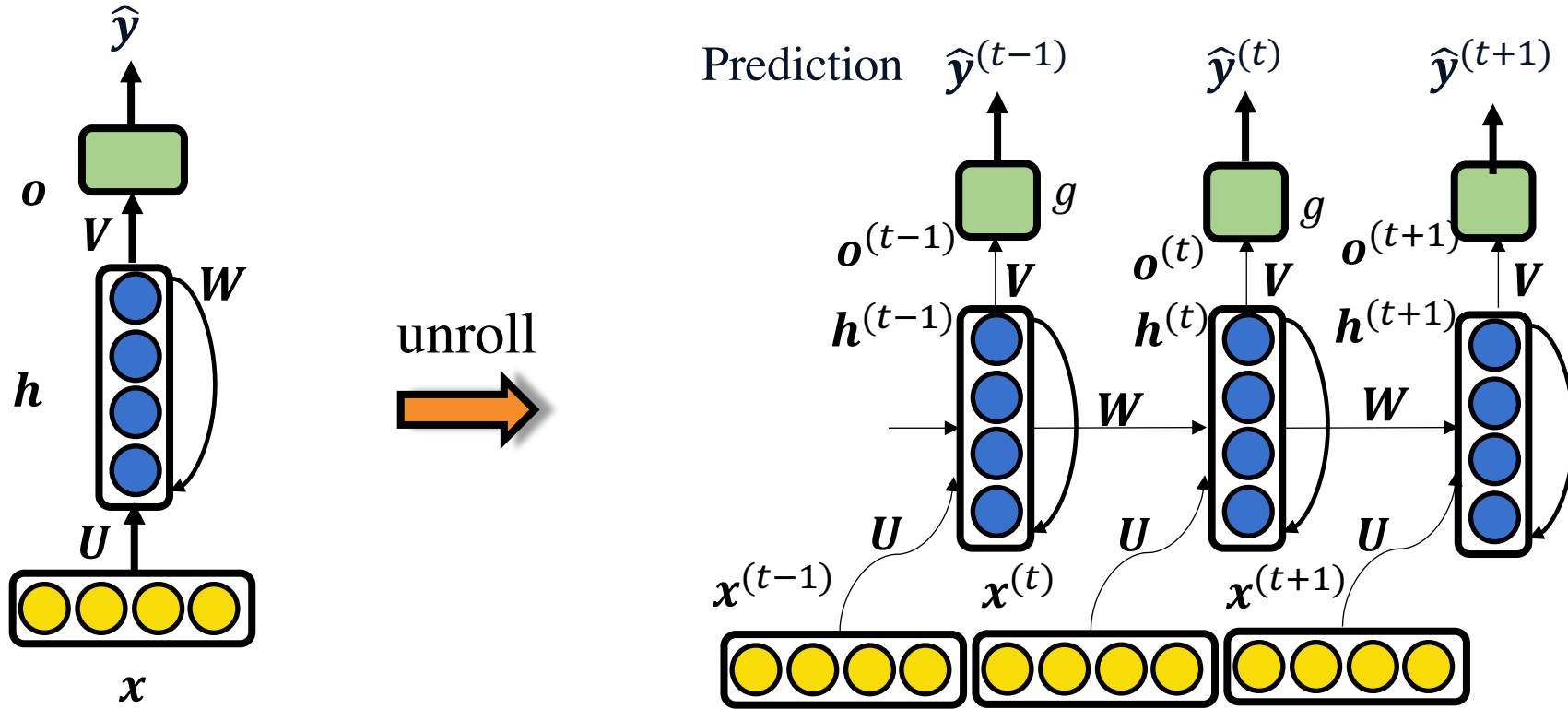
Agenda

- **RNN Basics**
- Learning RNN with Backpropagation Through Time (BPTT)
- Long-Short Term Memory Networks (LSTM)
- Gated Recurrent Unit (GRU)
- Bidirectional RNN
- Sequence-to-Sequence RNN
- Healthcare Applications

Basic Concepts of RNN



Basic RNN Structure

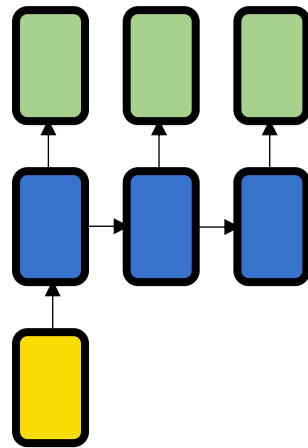


$$h^{(t)} = f(Ux^{(t)} + \mathbf{W}h^{(t-1)} + b_1)$$

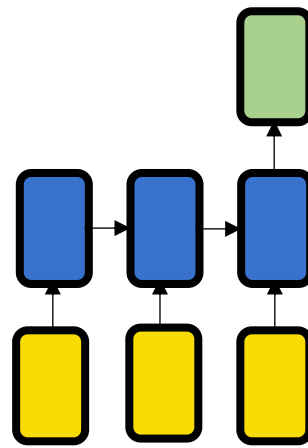
$$o^{(t)} = Vh^{(t)} + b_2$$

$$\hat{y}^{(t)} = g(o^{(t)})$$

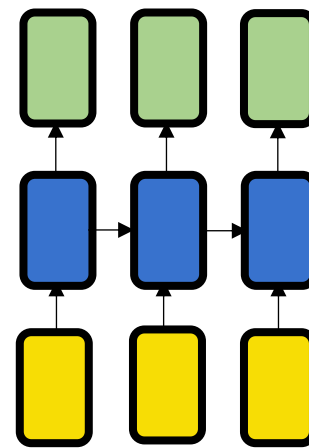
Basic RNN Structure



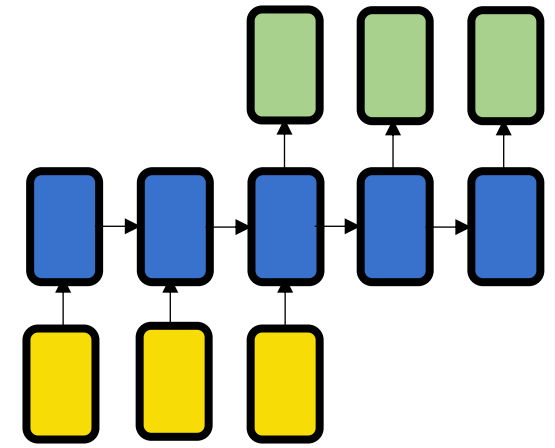
One-to-Many
(e.g., image
to text)



Many-to-One
(e.g., sequence
classification)



Many-to-Many
(e.g., sequential
prediction)



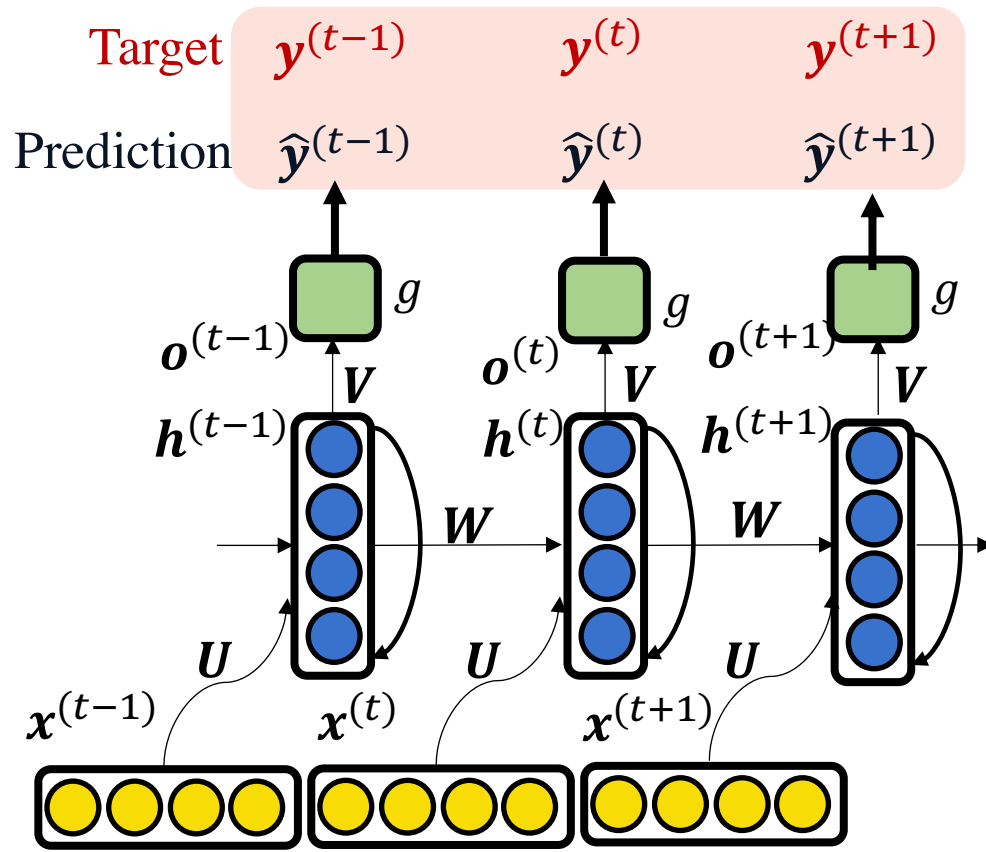
Many-to-Many
(e.g., seq2seq)

The figure is inspired by page 12 on http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf

Agenda

- RNN Basics
- **Learning RNN with Backpropagation Through Time (BPTT)**
- Long-Short Term Memory Networks (LSTM)
- Gated Recurrent Unit (GRU)
- Bidirectional RNN
- Sequence-to-Sequence RNN
- Healthcare Applications

Forward Computation



$$\mathbf{z}^{(t)} = \mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{b}_1$$

$$\mathbf{h}^{(t)} = f(\mathbf{z}^{(t)})$$

$$\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)} + \mathbf{b}_2$$

$$\hat{\mathbf{y}}^{(t)} = g(\mathbf{o}^{(t)})$$

$$L = \sum_t L^{(t)} = - \sum_t \log p(\mathbf{y}^{(t)} | \{\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(T)}\})$$

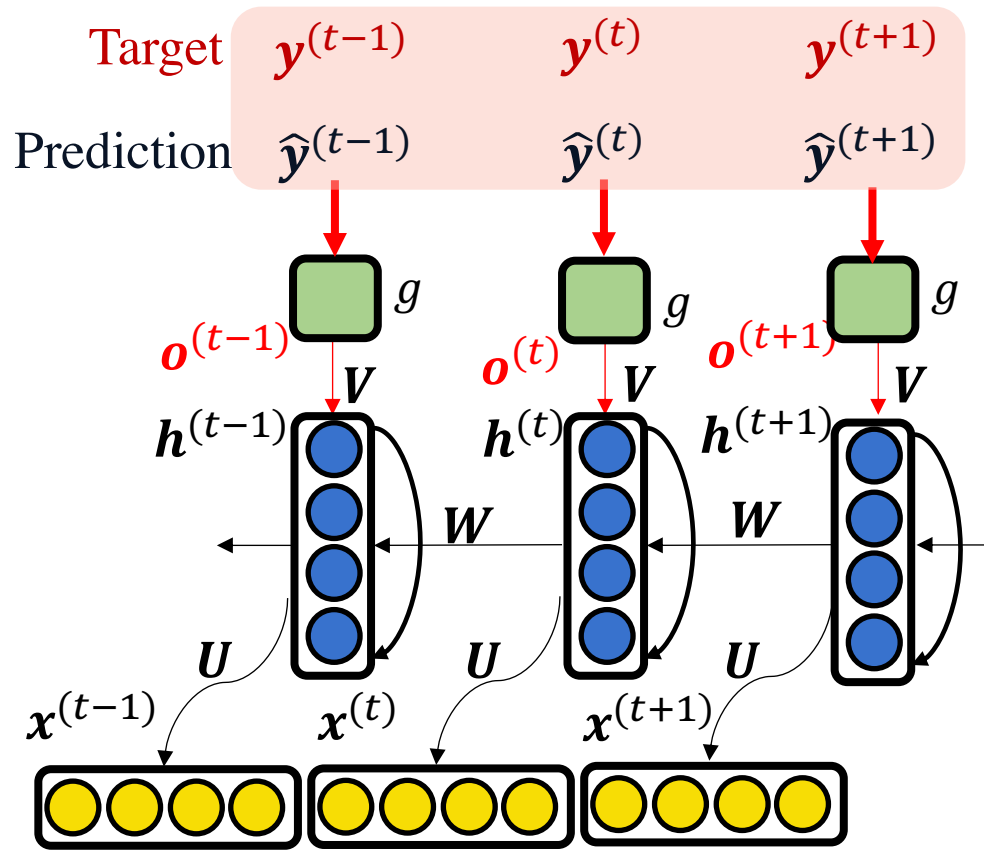
e. g.

- f is tanh
- g is softmax, which produces a normalized probability over output classes
- $L^{(t)}$ is negative log-likelihood loss

e.g., in binary classification

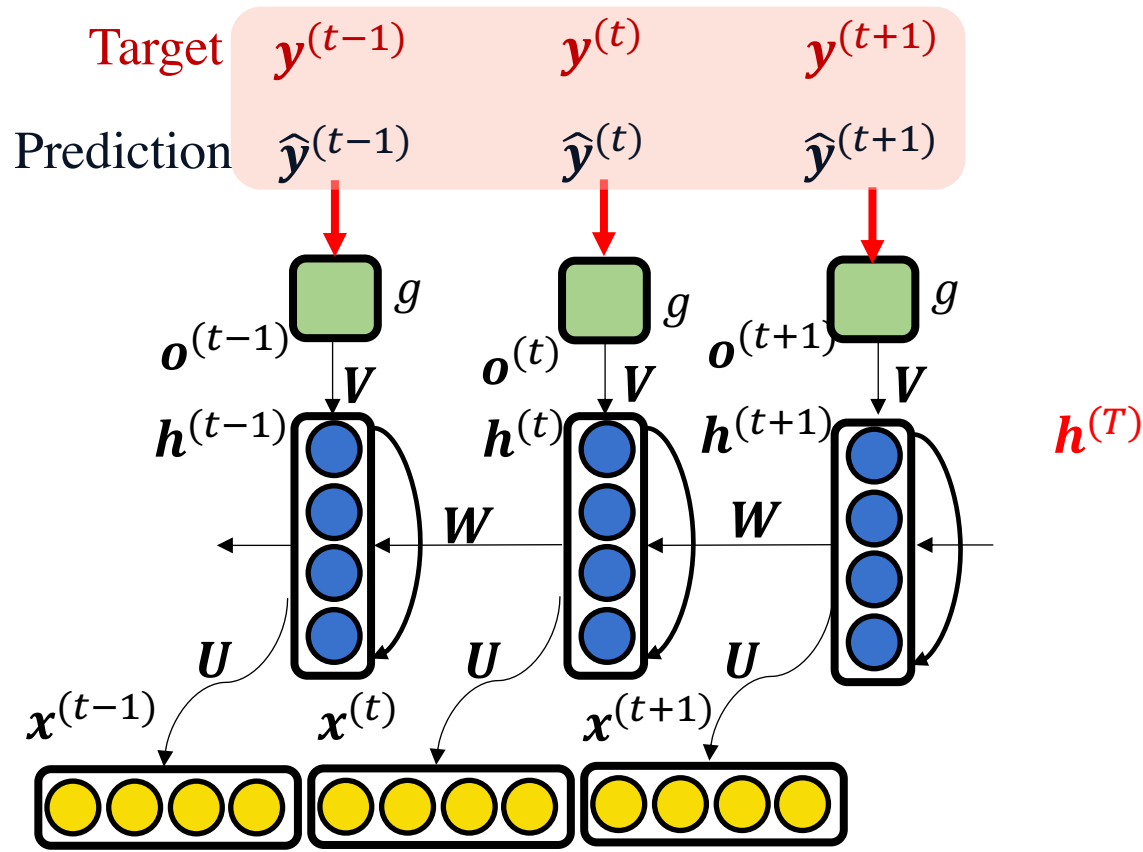
$$L = - \sum_t y^{(t)} \log \hat{y}^{(t)} + (1 - y^{(t)}) \log(1 - \hat{y}^{(t)})$$

Backpropagation through time (BPTT): $\nabla_{\mathbf{o}^{(t)}} L$



$$(\nabla_{\mathbf{o}^{(t)}} L)_i = \frac{\partial L}{\partial \mathbf{o}_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}}$$

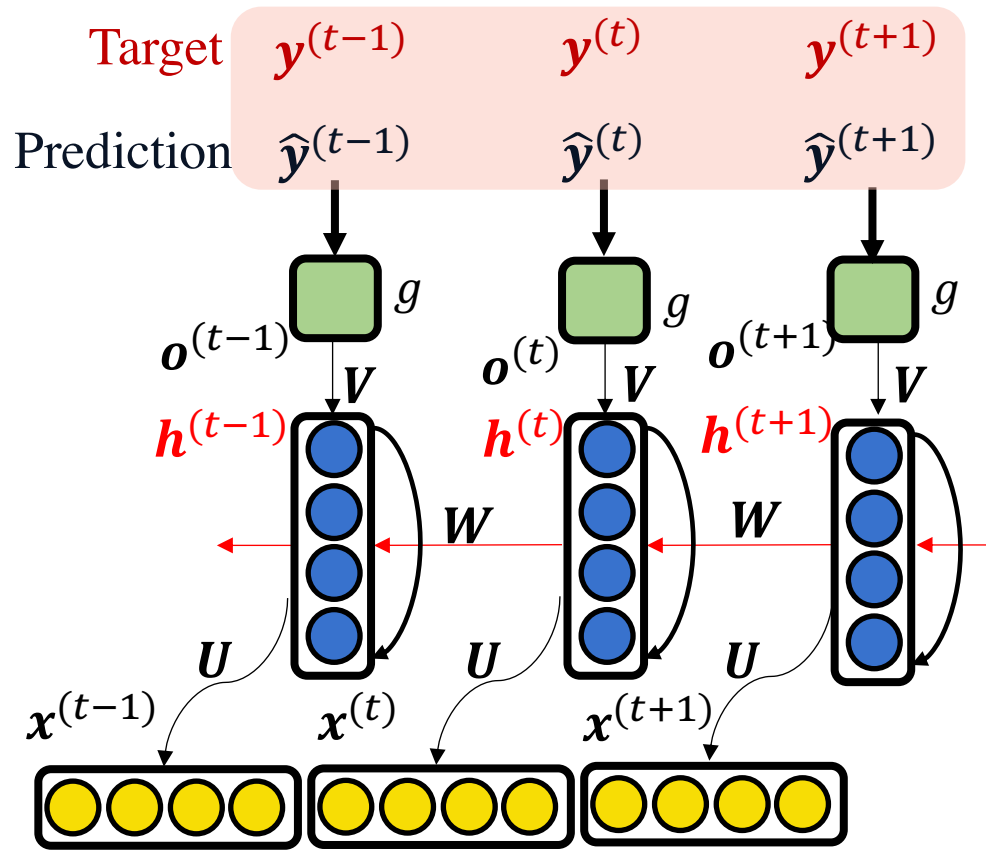
Backpropagation through time (BPTT): $\nabla_{h^{(T)}} L$



Last time stamp

$$\nabla_{h^{(T)}} L = V^T \nabla_{o^{(T)}} L$$

Backpropagation through time (BPTT): $\nabla_{\mathbf{h}^{(t)}} L$



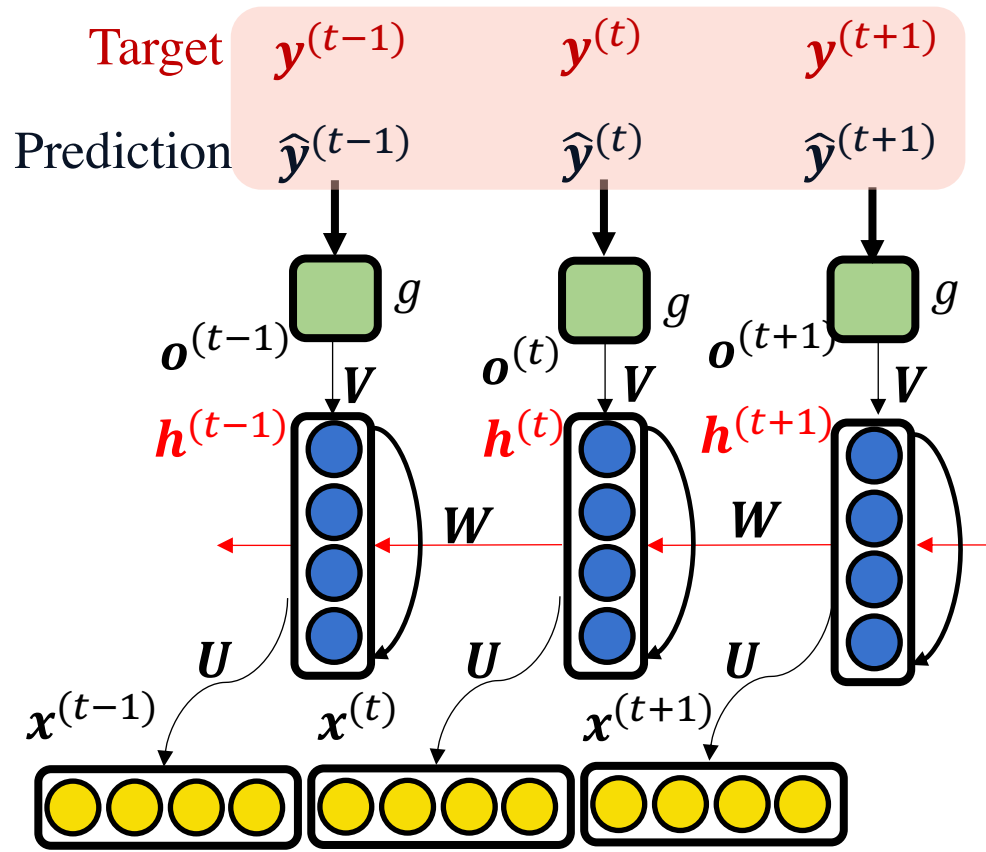
Last time stamp

$$\nabla_{\mathbf{h}^{(T)}} L = V^T \nabla_{\mathbf{o}^{(T)}} L$$

Other time stamp

$$\nabla_{\mathbf{h}^{(t)}} L = \left(\frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^T (\nabla_{\mathbf{h}^{(t+1)}} L) + \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^T (\nabla_{\mathbf{o}^{(t)}} L)$$

Backpropagation through time (BPTT)



$$\nabla_{\mathbf{V}} \mathcal{L} = \sum_t \sum_k \left(\frac{\partial \mathcal{L}}{\partial o_k^{(t)}} \right) \nabla_{\mathbf{V}} o_k^{(t)} = \sum_t (\nabla_{\mathbf{o}^{(t)}} \mathcal{L}) \mathbf{h}^{(t)T}$$

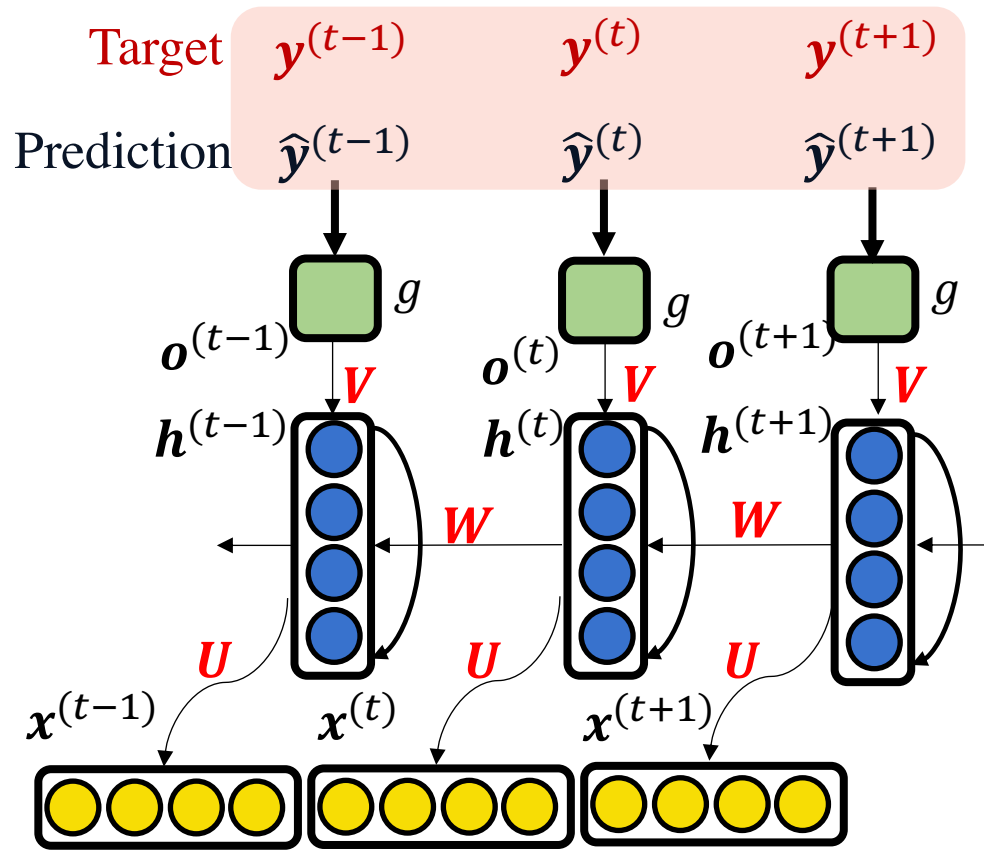
$$\nabla_{\mathbf{W}} \mathcal{L} = \sum_t \sum_j \left(\frac{\partial \mathcal{L}}{\partial h_j^{(t)}} \right) \nabla_{\mathbf{W}} h_j^{(t)} = \sum_t \text{diag}(1 - (\mathbf{h}^{(t)})^2) (\nabla_{\mathbf{h}^{(t)}} \mathcal{L}) \mathbf{h}^{(t-1)T}$$

$$\nabla_{\mathbf{U}} \mathcal{L} = \sum_t \sum_j \left(\frac{\partial \mathcal{L}}{\partial h_j^{(t)}} \right) \nabla_{\mathbf{U}} h_j^{(t)} = \sum_t \text{diag}(1 - (\mathbf{h}^{(t)})^2) (\nabla_{\mathbf{h}^{(t)}} \mathcal{L}) \mathbf{x}^{(t)T}$$

$$\nabla_{b_1} \mathcal{L} = \sum_t \left(\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^T \nabla_{\mathbf{h}^{(t)}} \mathcal{L} = \sum_t \text{diag}(1 - (\mathbf{h}^{(t)})^2) \nabla_{\mathbf{h}^{(t)}} \mathcal{L}$$

$$\nabla_{b_2} \mathcal{L} = \sum_t \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{b}_2} \right)^T \nabla_{\mathbf{o}^{(t)}} \mathcal{L} = \sum_t \nabla_{\mathbf{o}^{(t)}} \mathcal{L}$$

BPTT: Vanishing gradient problem

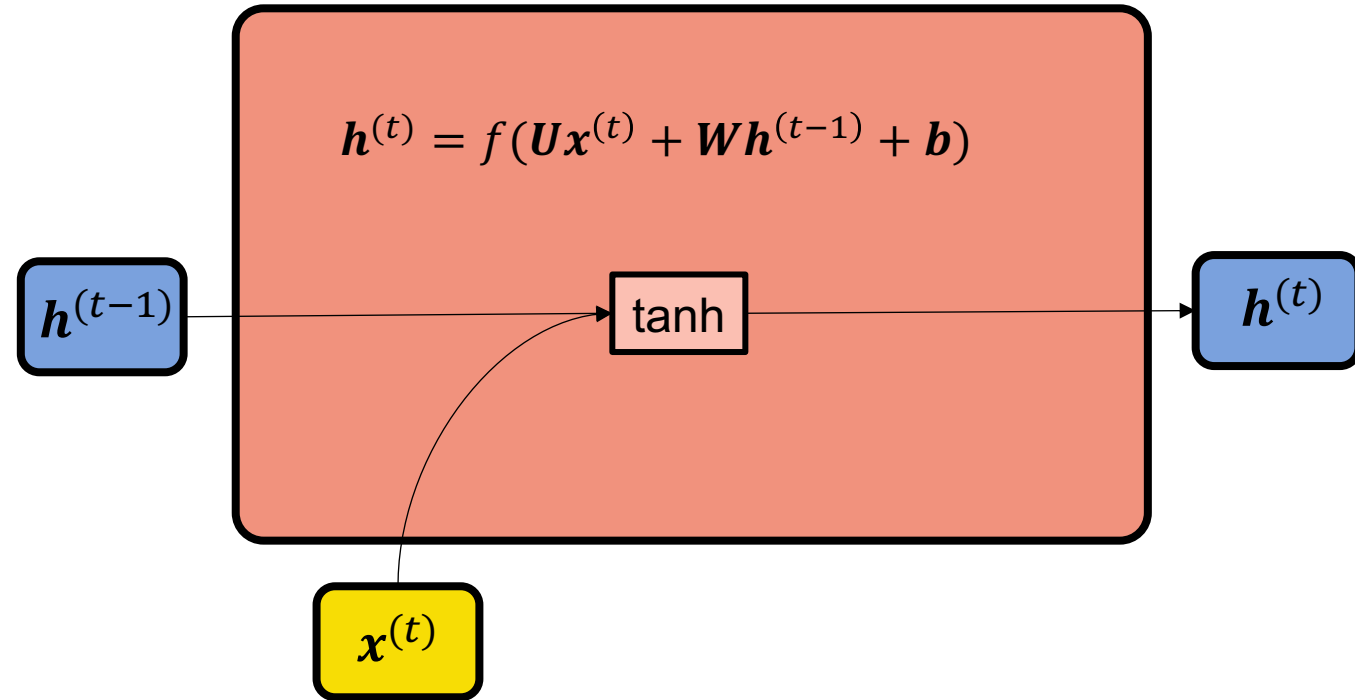
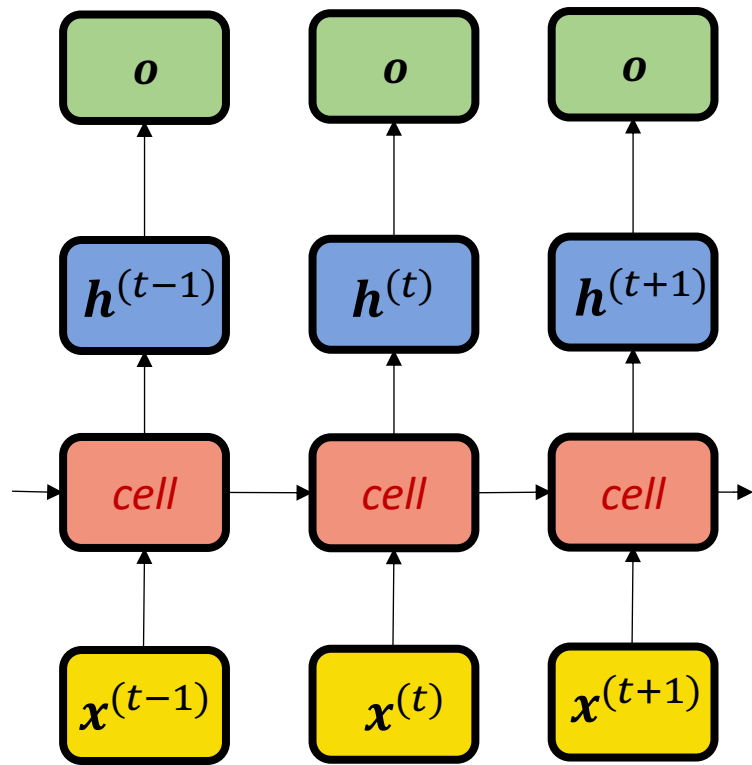


- Gradient can become very small over a long sequence
- Standard RNN will have difficulty to “remember” state from early part of the input sequence

Agenda

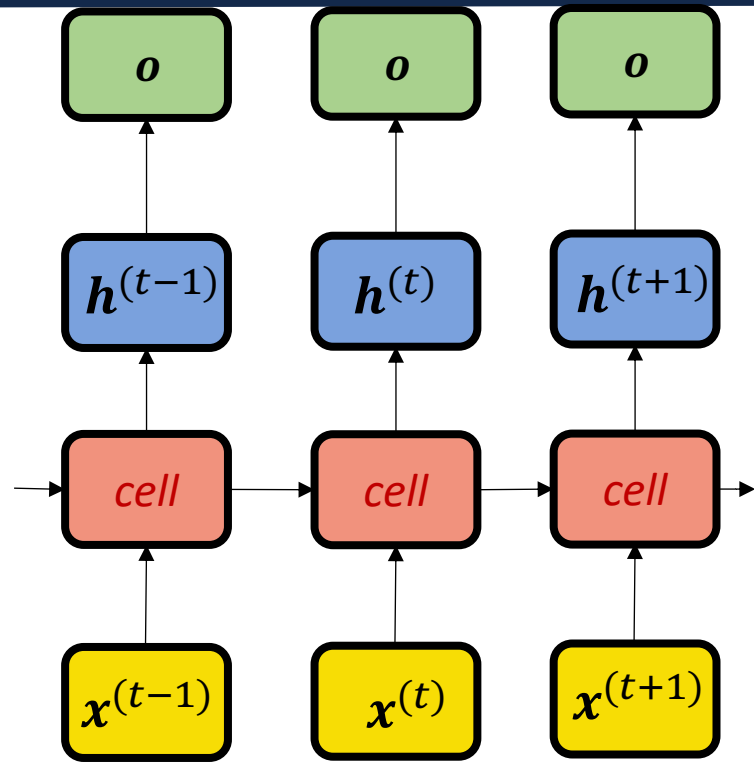
- RNN Basics
- Learning RNN with Backpropagation Through Time (BPTT)
- **Long-Short Term Memory Networks (LSTM)**
- Gated Recurrent Unit (GRU)
- Bidirectional RNN
- Sequence-to-Sequence RNN
- Healthcare Applications

Standard RNN

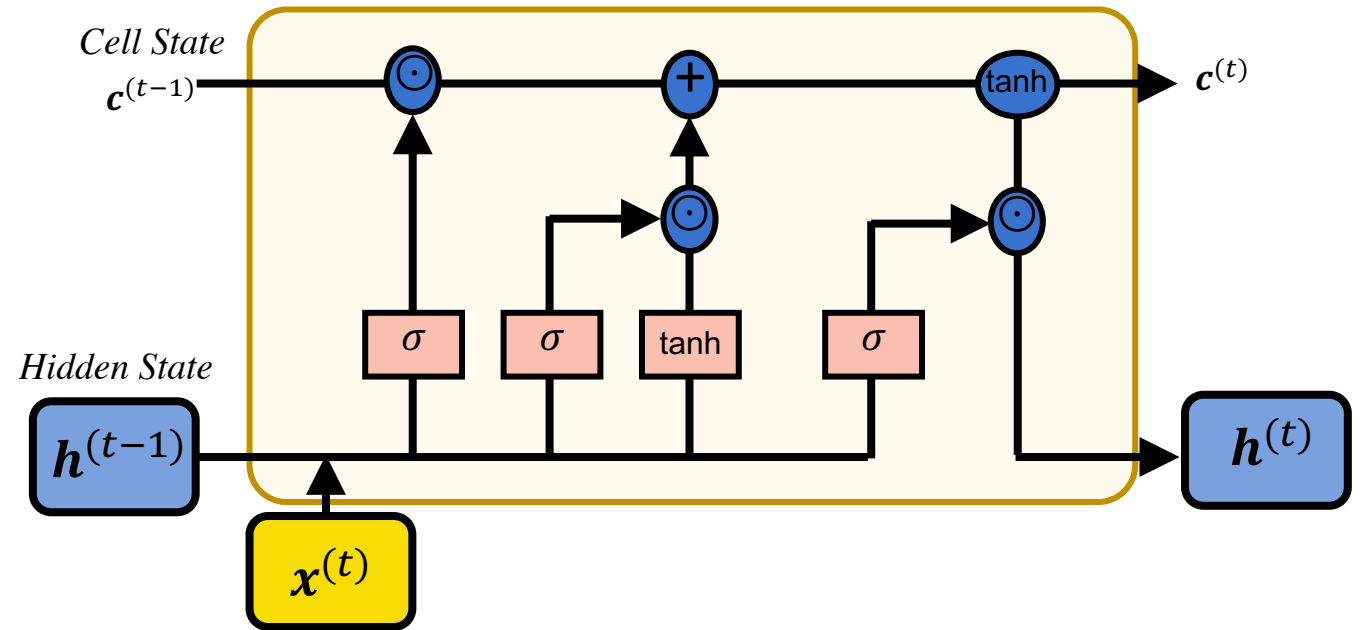


- Standard RNN has a simple computation *cell* from input x to latent state h

Long short term memory networks (LSTM)

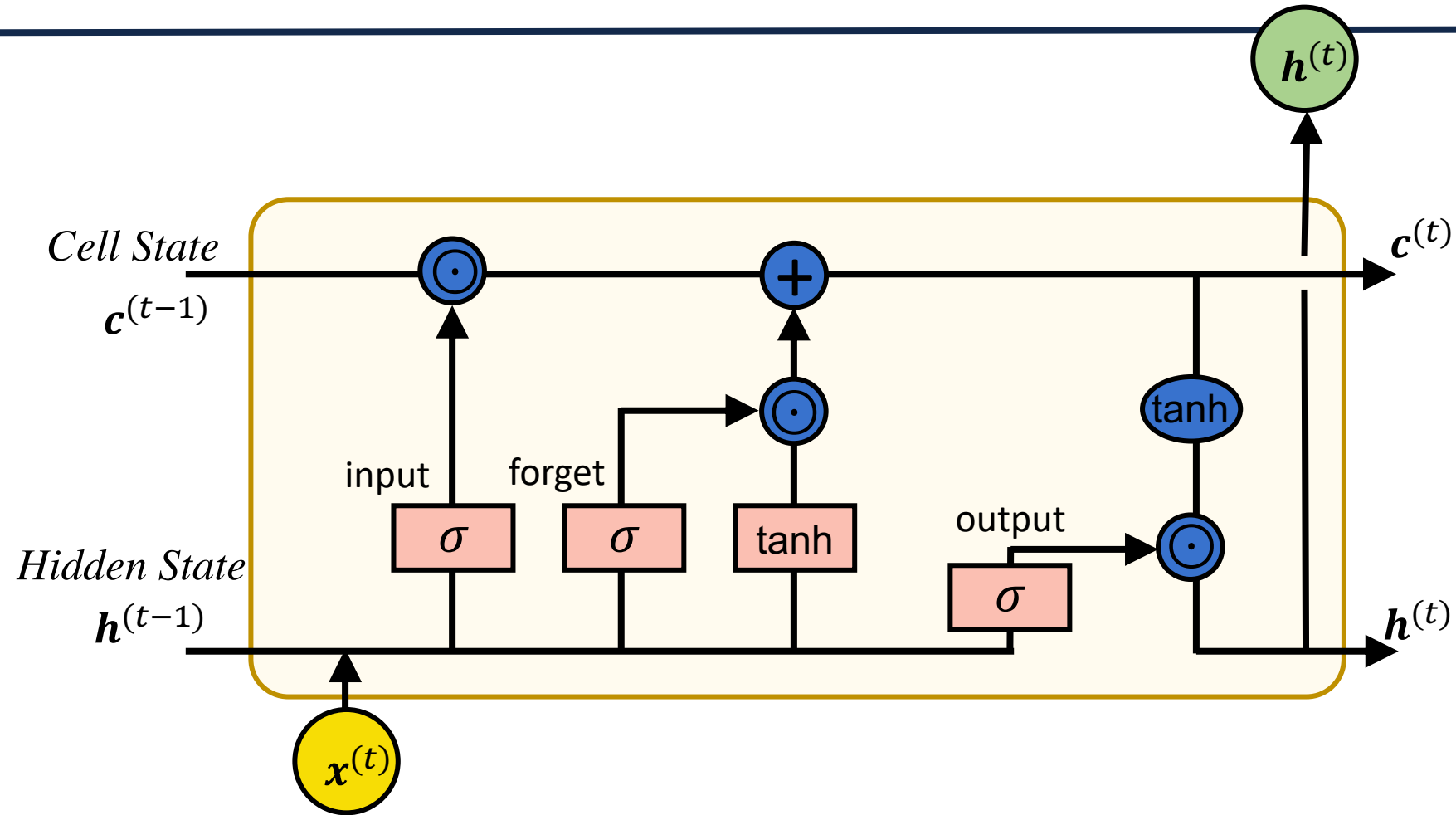


$$h^{(t)} = f(Ux^{(t)} + Wh^{(t-1)} + b)$$

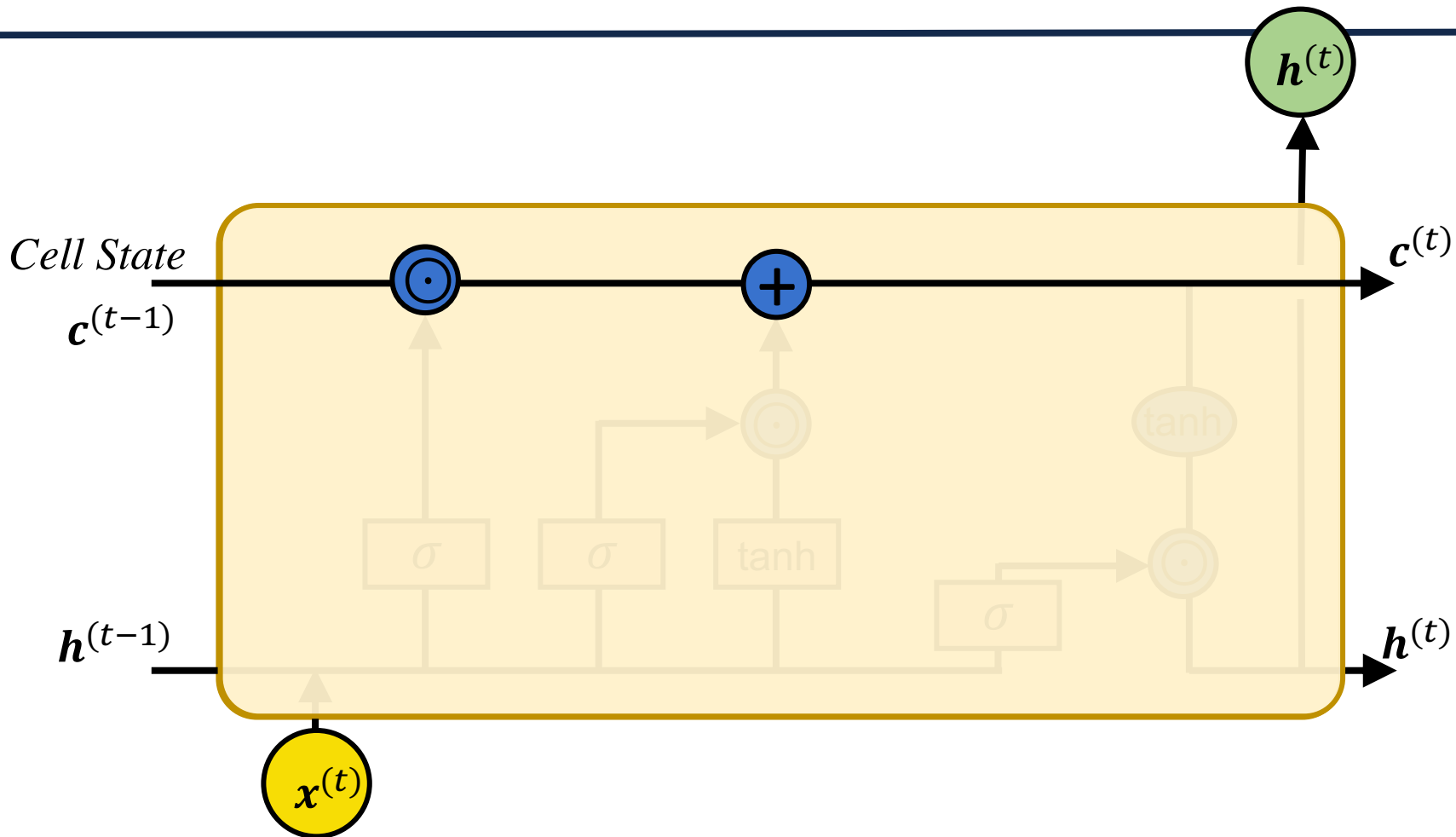


- Standard RNN has a simple computation *cell* from input x to latent state h
- LSTM provides a more sophisticated *cell*

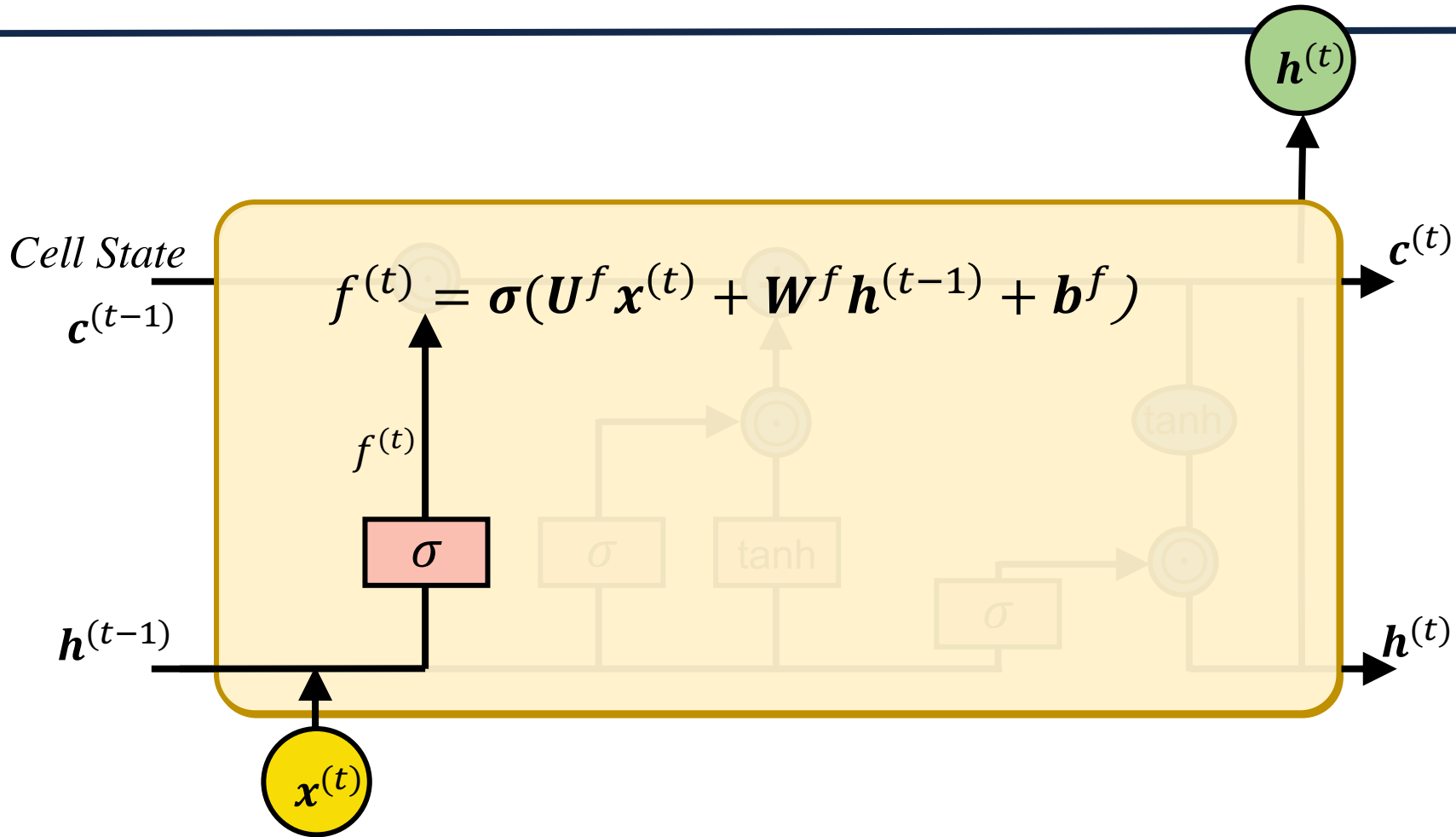
LSTM: Cell Structure



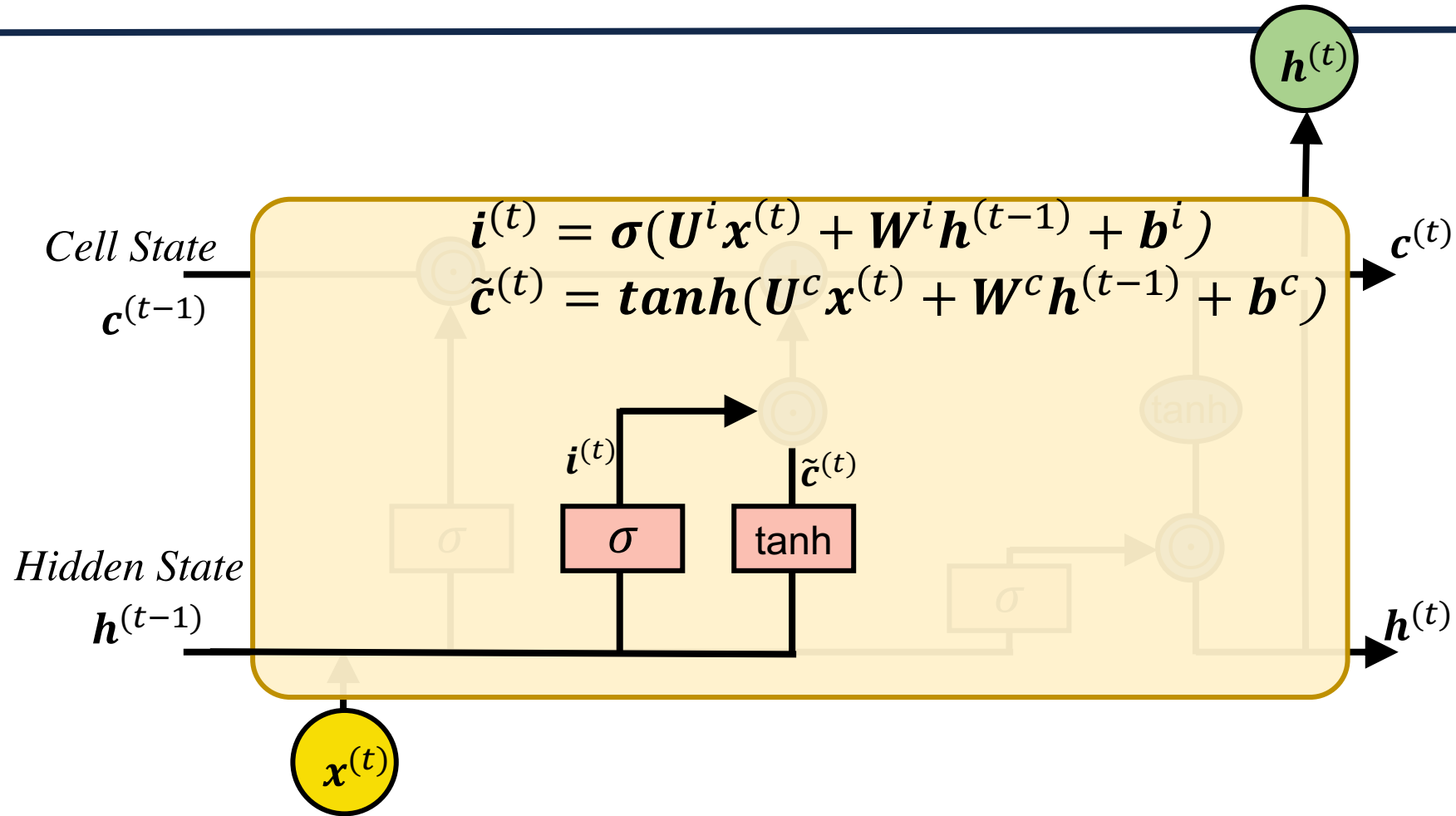
LSTM: Cell state



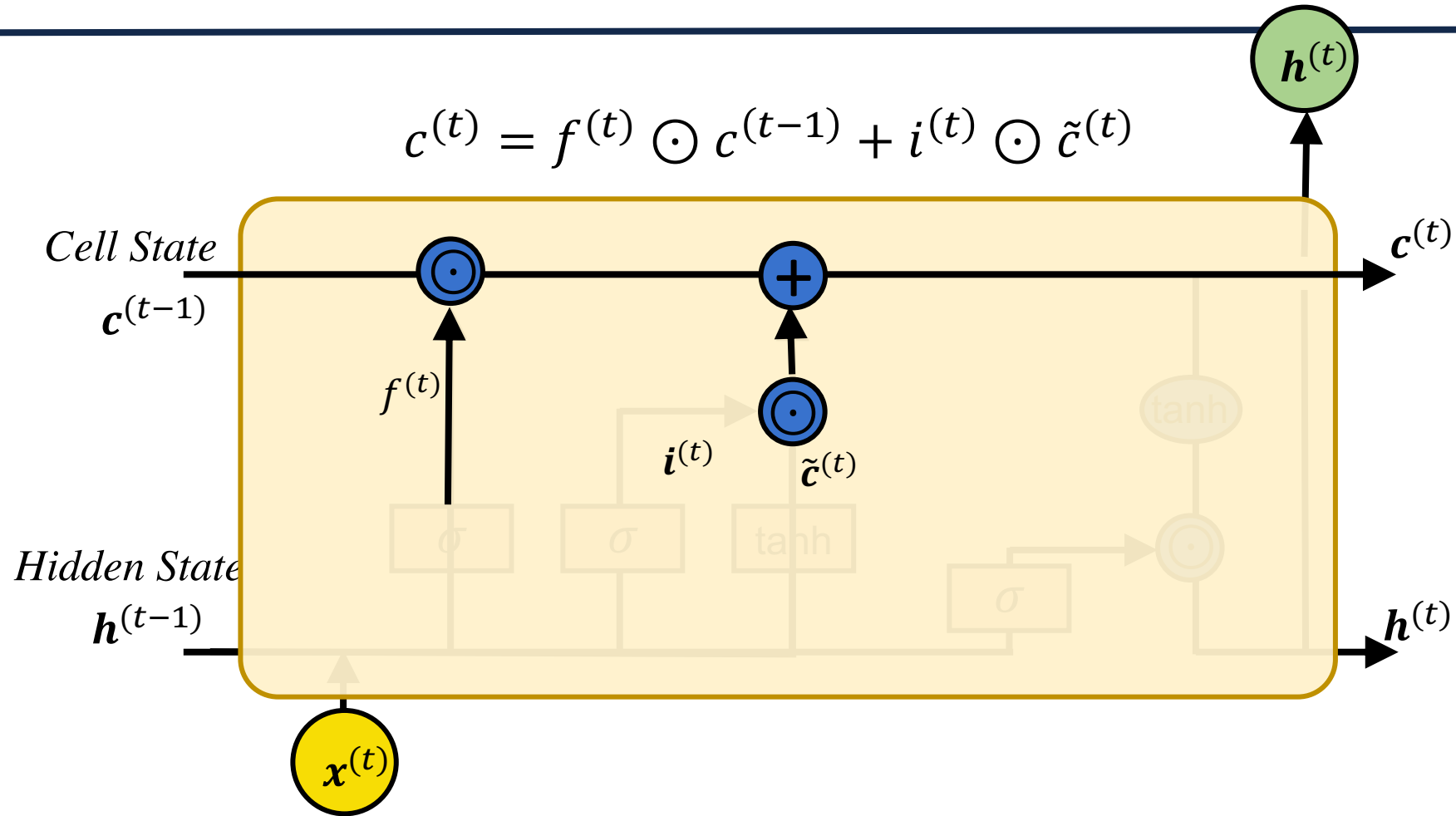
LSTM: Forget Gate



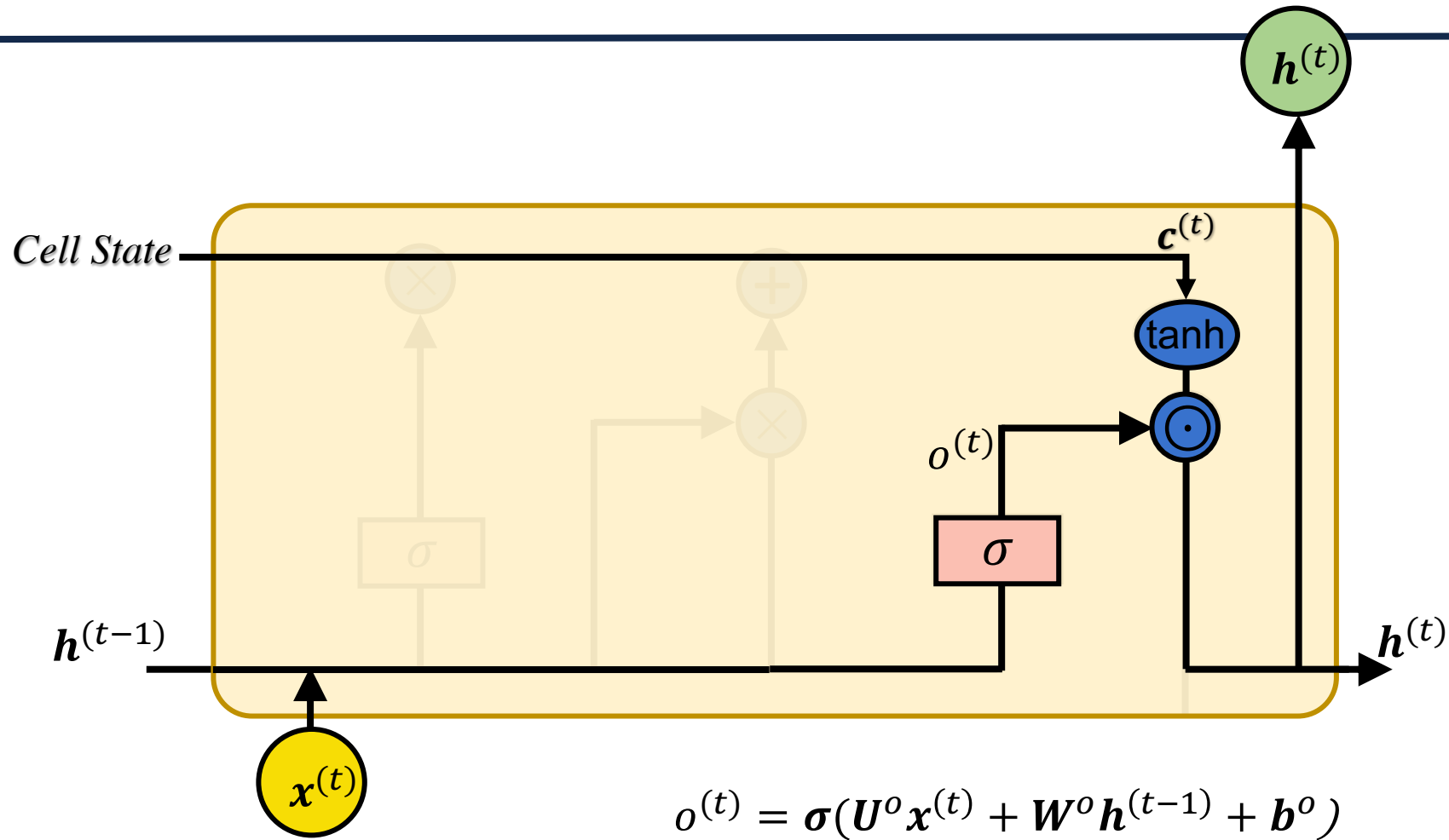
LSTM: Input Gate



LSTM: Update cell state

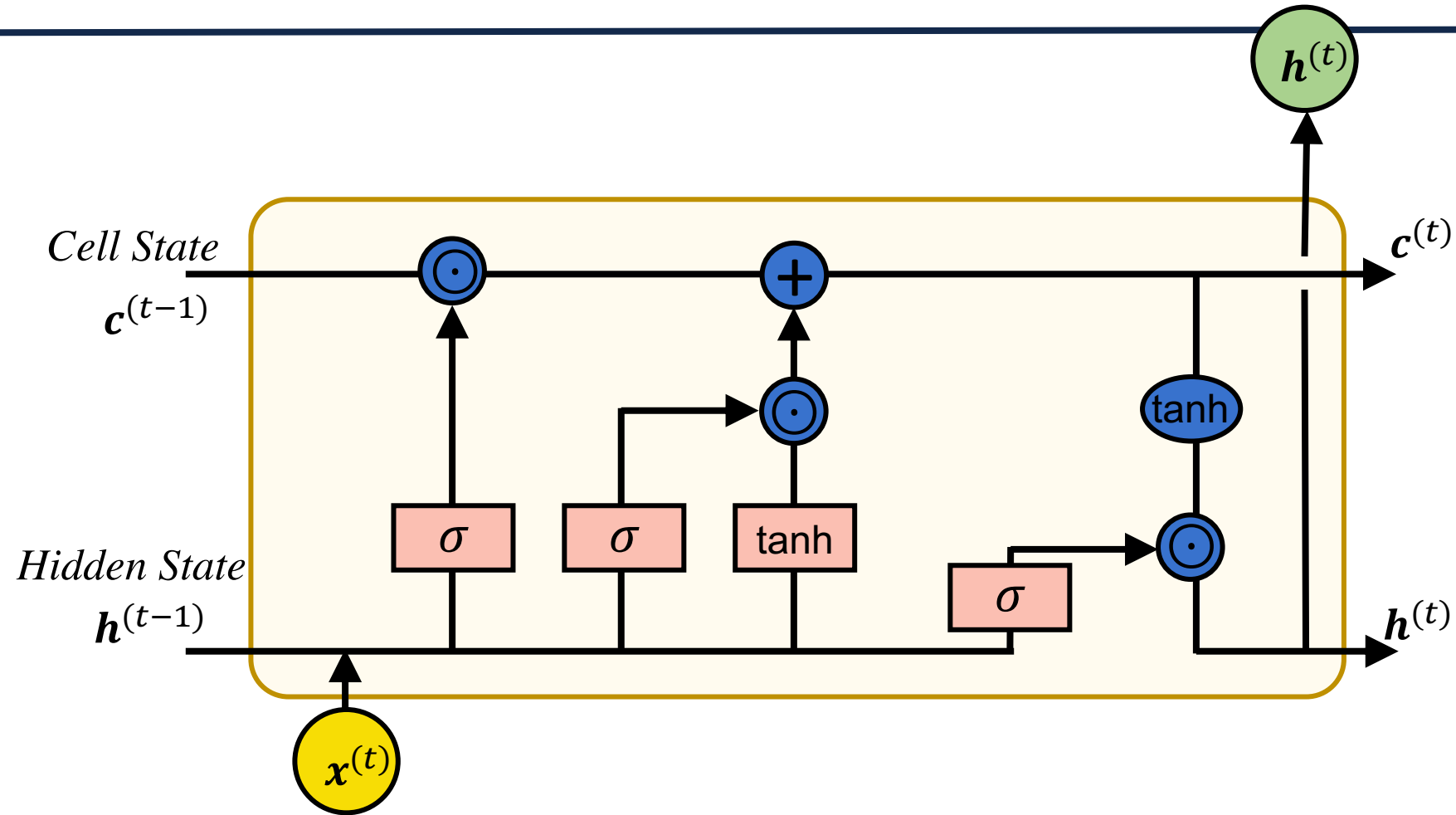


LSTM: Output Gate



$$o^{(t)} = \sigma(U^o x^{(t)} + W^o h^{(t-1)} + b^o)$$
$$h^{(t)} = o^{(t)} \odot \tanh(c^{(t)})$$

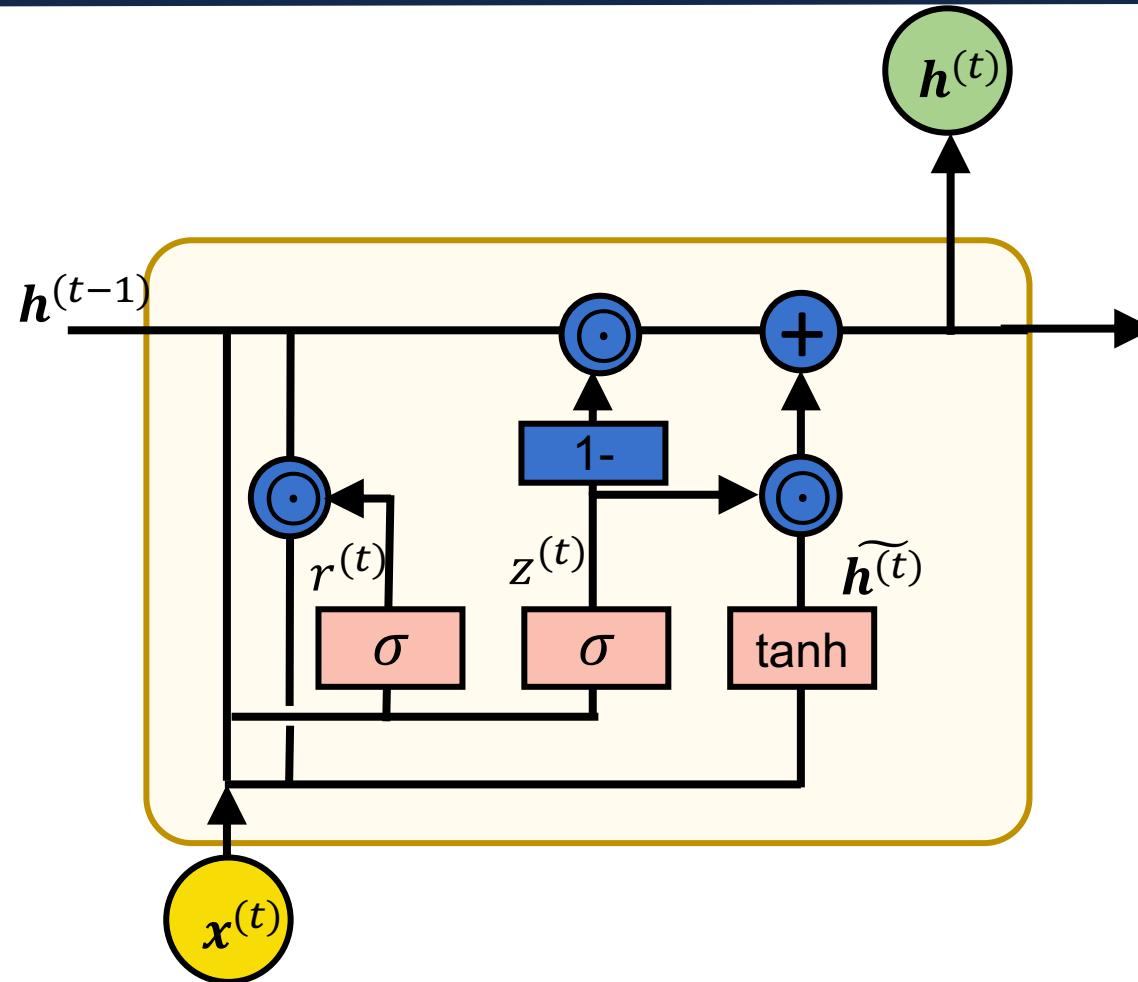
LSTM: Cell Structure



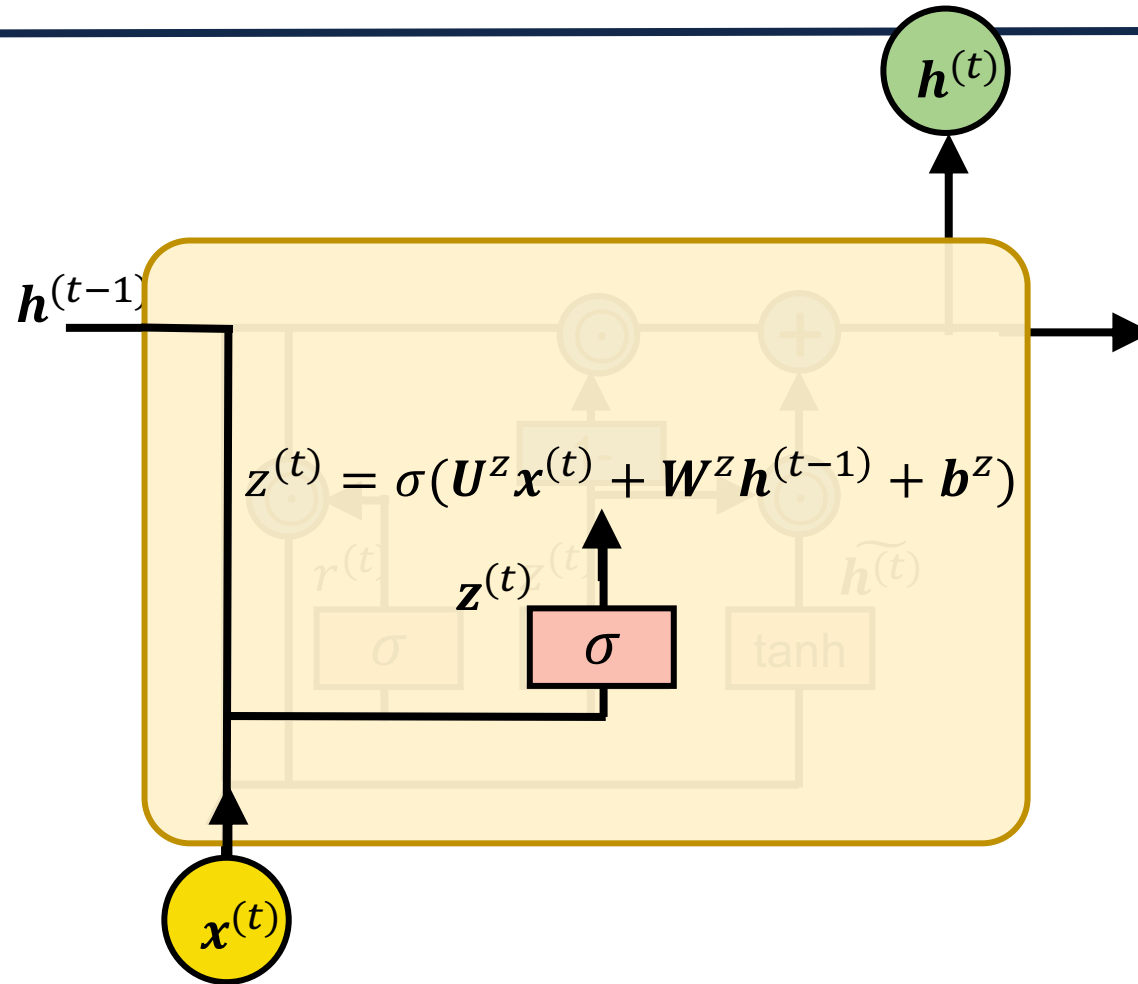
Agenda

- RNN Basics
- Learning RNN with Backpropagation Through Time (BPTT)
- Long-Short Term Memory Networks (LSTM)
- **Gated Recurrent Unit (GRU)**
- Bidirectional RNN
- Sequence-to-Sequence RNN
- Healthcare Applications

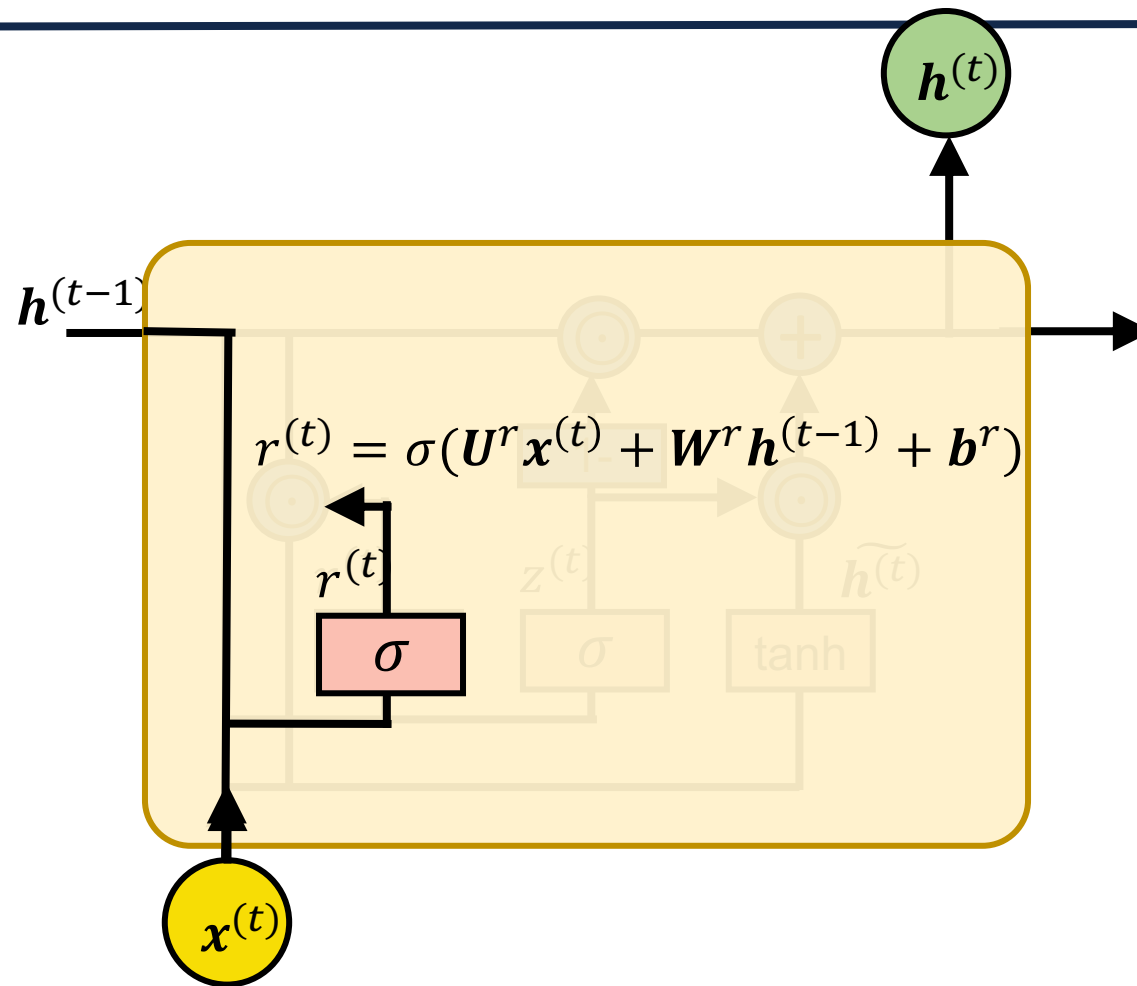
GRU : Gated Recurrent Unit



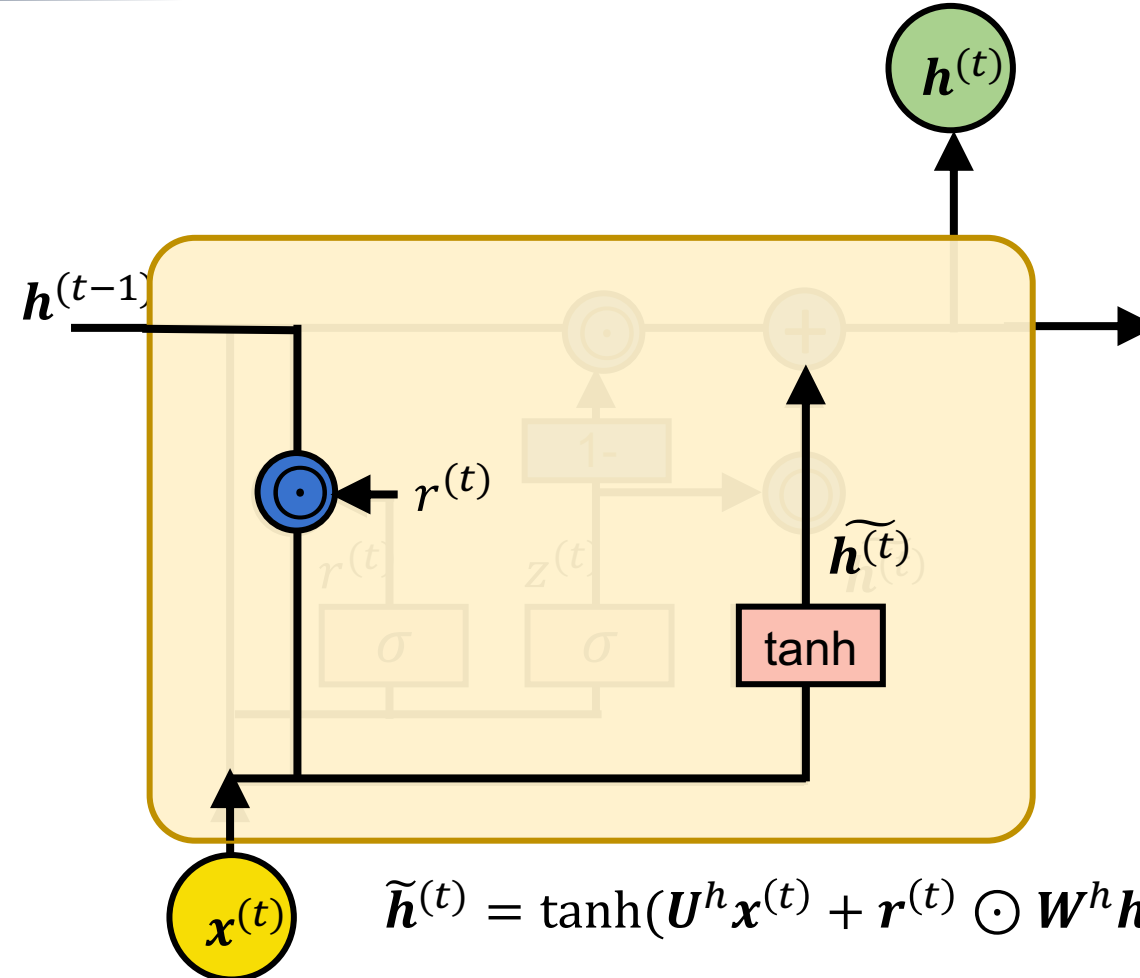
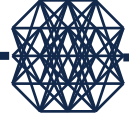
GRU: Update Gate



GRU: Reset Gate

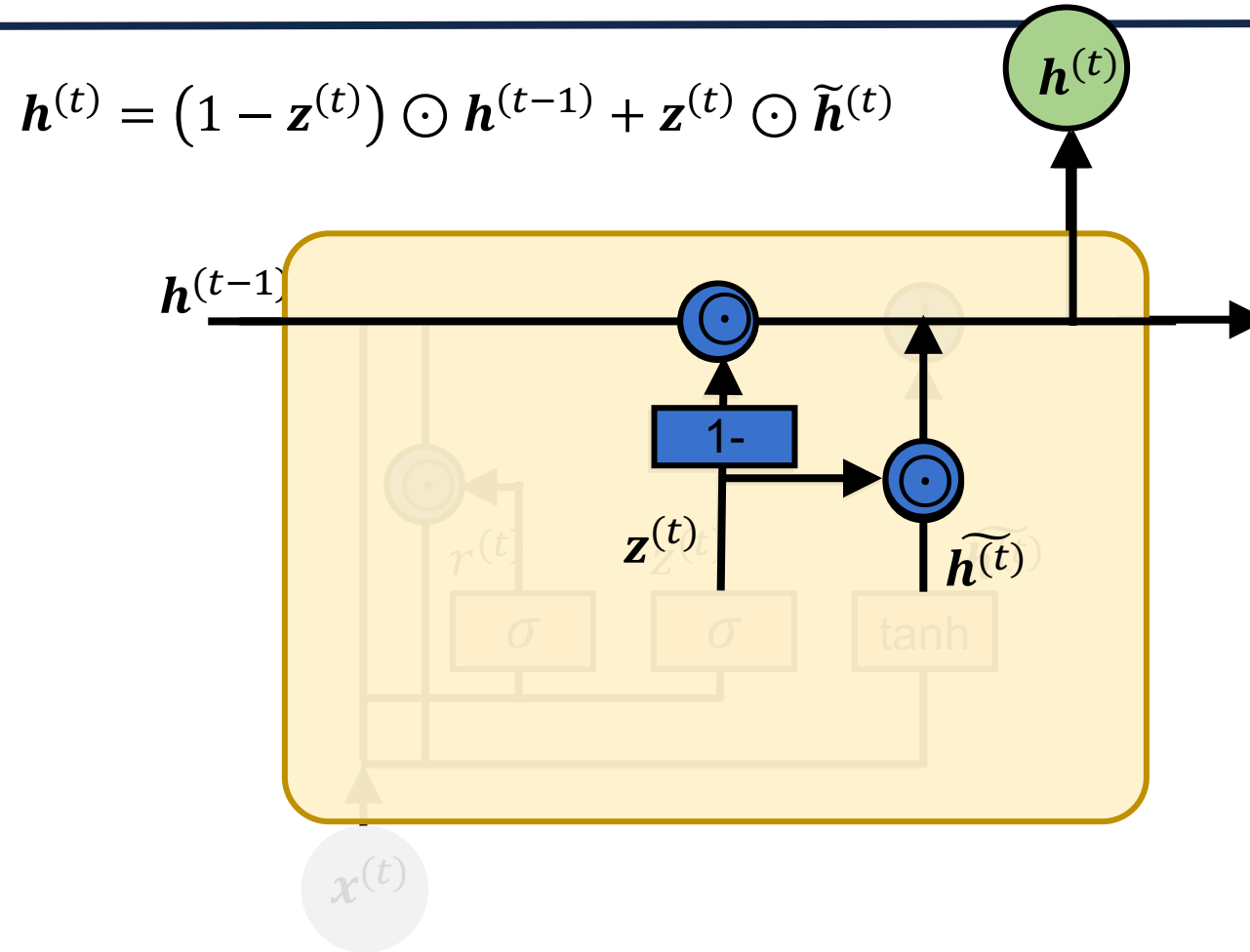


GRU: New information to the hidden state

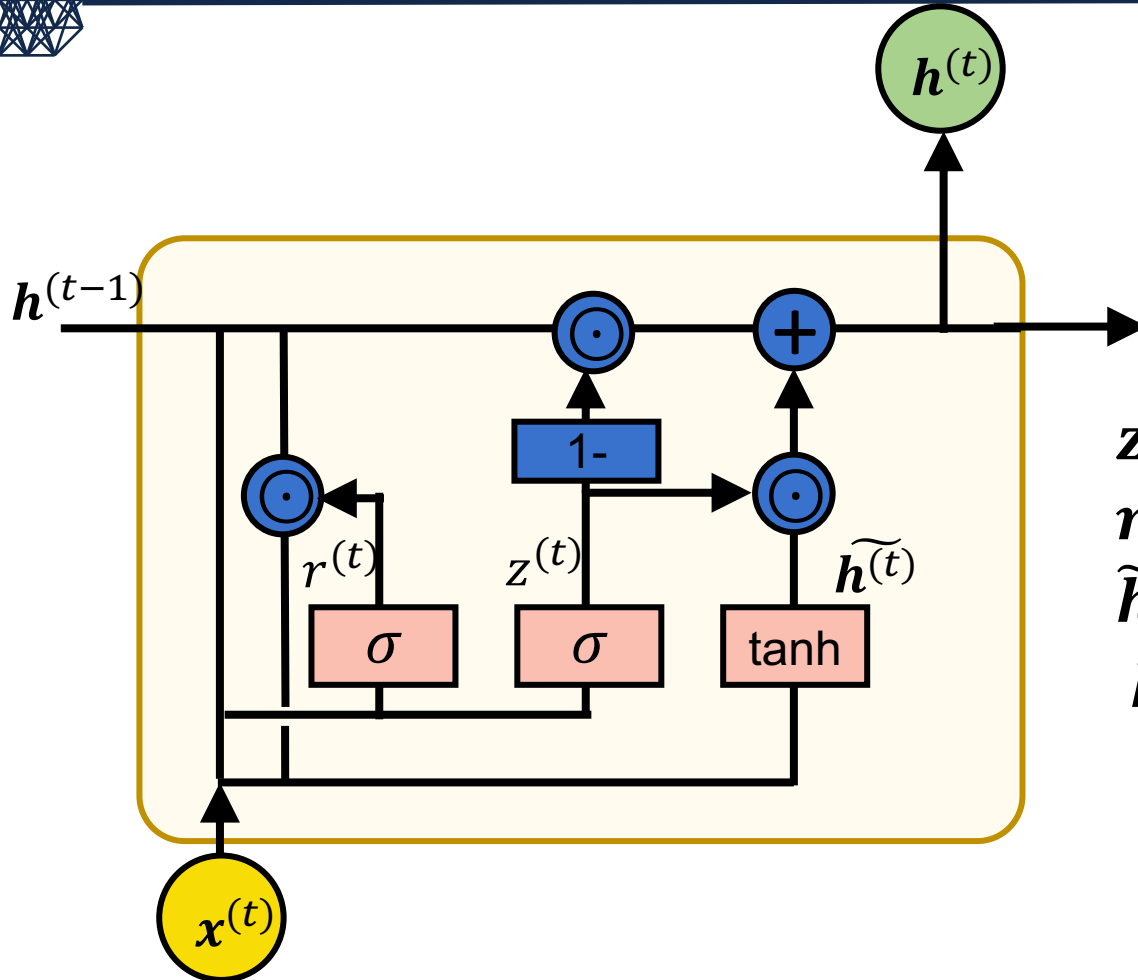


$$\tilde{h}^{(t)} = \tanh(U^h x^{(t)} + r^{(t)} \odot W^h h^{(t-1)} + b^h)$$

GRU: Final New Hidden State



GRU: Summary



$$\mathbf{z}^{(t)} = \sigma(\mathbf{U}^z \mathbf{x}^{(t)} + \mathbf{W}^z \mathbf{h}^{(t-1)} + \mathbf{b}^z)$$

$$\mathbf{r}^{(t)} = \sigma(\mathbf{U}^r \mathbf{x}^{(t)} + \mathbf{W}^r \mathbf{h}^{(t-1)} + \mathbf{b}^r)$$

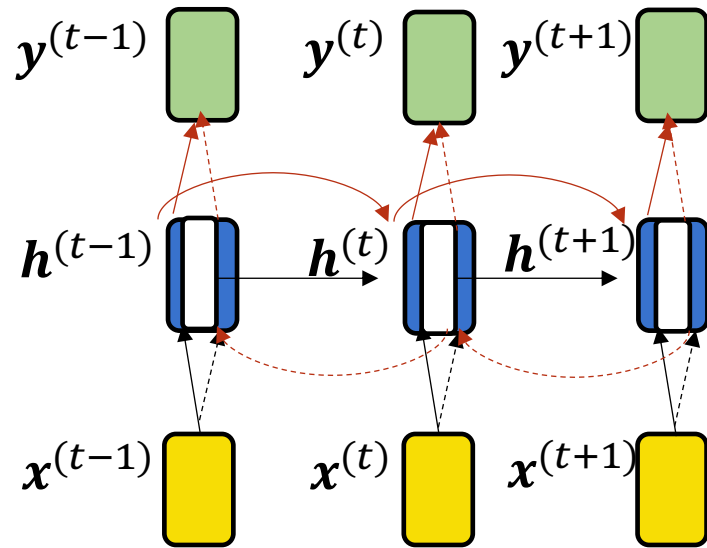
$$\tilde{\mathbf{h}}^{(t)} = \tanh(\mathbf{U}^h \mathbf{x}^{(t)} + \mathbf{r}^{(t)} \odot \mathbf{W}^h \mathbf{h}^{(t-1)} + \mathbf{b}^h)$$

$$\mathbf{h}^{(t)} = (1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \odot \tilde{\mathbf{h}}^{(t)}$$

Agenda

- RNN Basics
- Learning RNN with Backpropagation Through Time (BPTT)
- Long-Short Term Memory Networks (LSTM)
- Gated Recurrent Unit (GRU)
- **Bidirectional RNN**
- Sequence-to-Sequence RNN
- Healthcare Applications

Bidirectional RNN



$$\vec{h}^t = f(Ux^{(t)} + Wh^{(t-1)} + b_1)$$

$$\overleftarrow{h}^t = f(Ux^{(t)} + Wh^{(t+1)} + b_1)$$

$$y^{(t)} = g(V[\vec{h}^t; \overleftarrow{h}^t] + b_2)$$

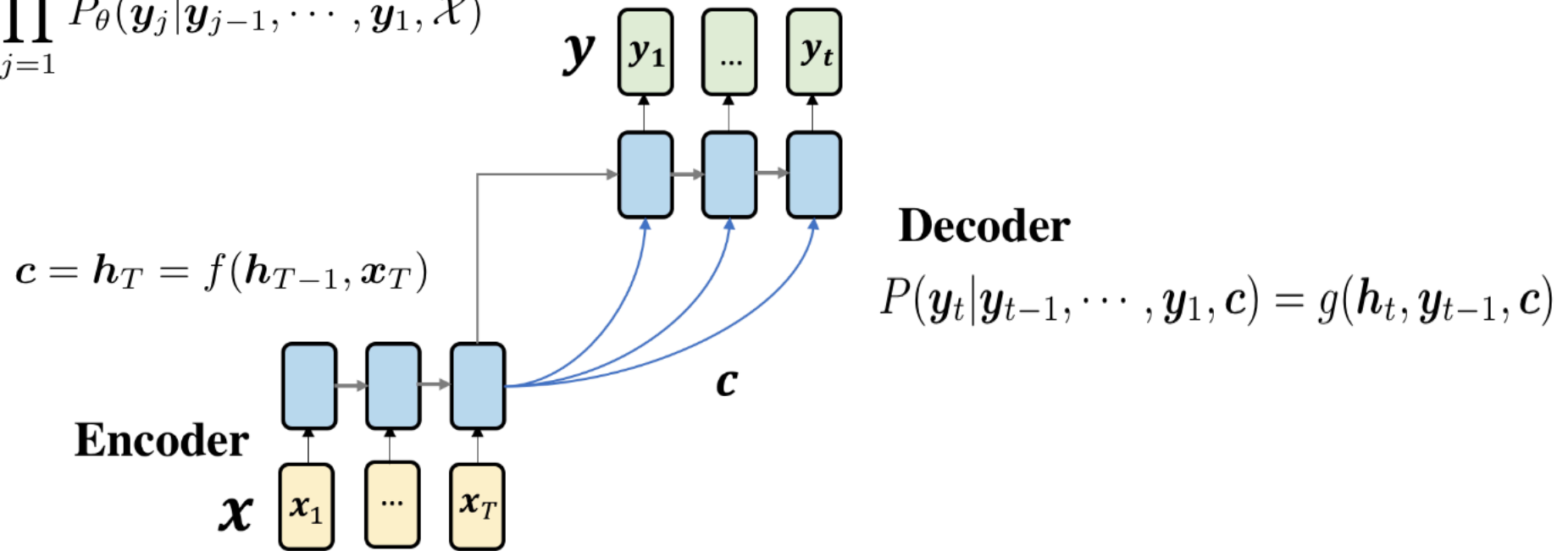
Agenda

- RNN Basics
- Learning RNN with Backpropagation Through Time (BPTT)
- Long-Short Term Memory Networks (LSTM)
- Gated Recurrent Unit (GRU)
- Bidirectional RNN
- **Sequence-to-Sequence RNN**
- Healthcare Applications

Encoder-Decoder Sequence-to-Sequence Model



$$P_{\theta}(\mathcal{Y}|\mathcal{X}) = \prod_{j=1}^{J+1} P_{\theta}(\mathbf{y}_j | \mathbf{y}_{j-1}, \dots, \mathbf{y}_1, \mathcal{X})$$



Agenda

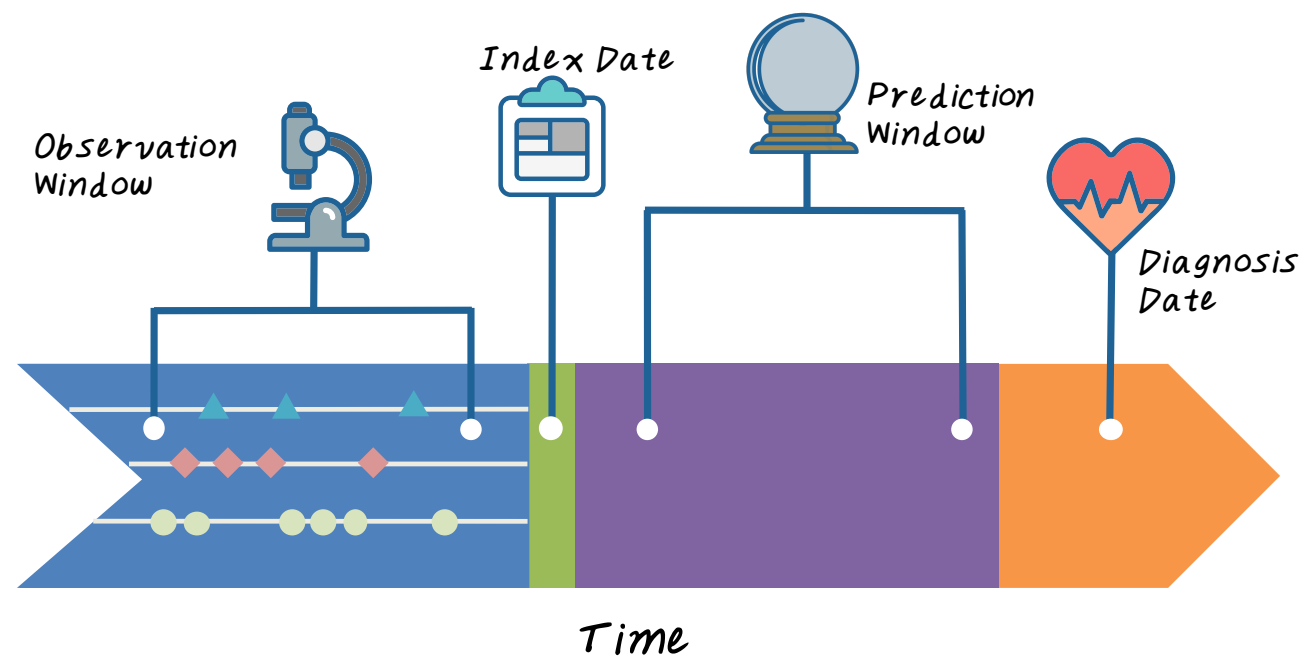
- RNN Basics
- Learning RNN with Backpropagation Through Time (BPTT)
- Long-Short Term Memory Networks (LSTM)
- Gated Recurrent Unit (GRU)
- Bidirectional RNN
- Sequence-to-Sequence RNN
- Healthcare Applications

Using Recurrent Neural Network Models For Early Detection Of Heart Failure Onset

How to model temporal relations in EHR

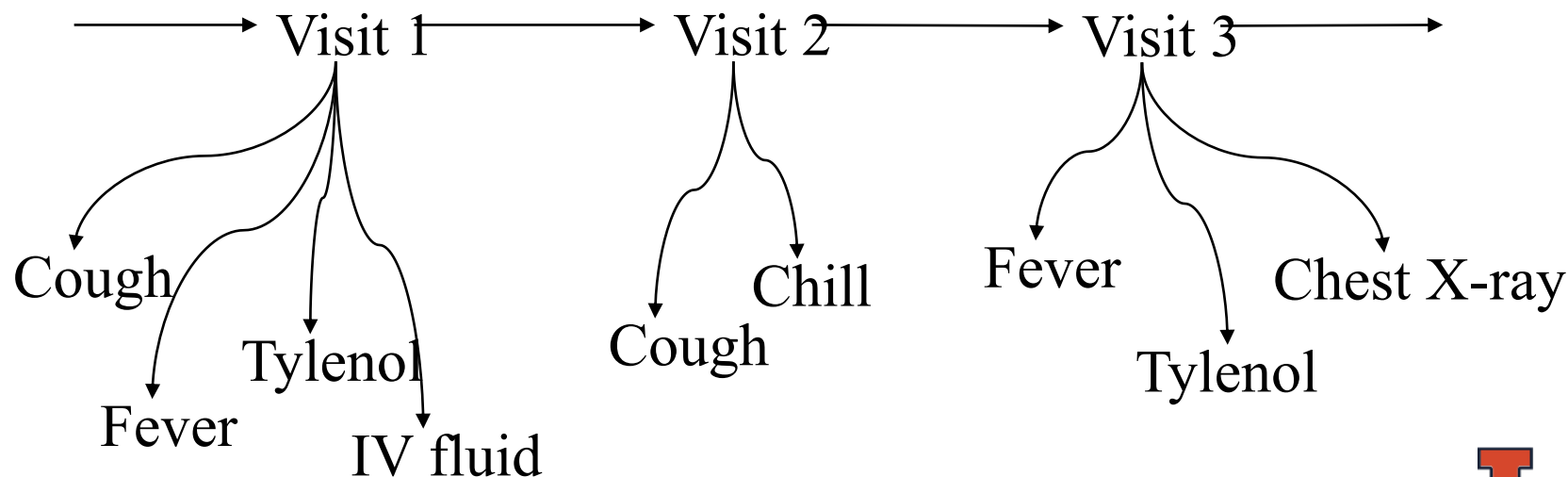
Heart Failure Prediction with RNN

- Given a patient record, predict if the patient will be diagnosed with heart failure in the future



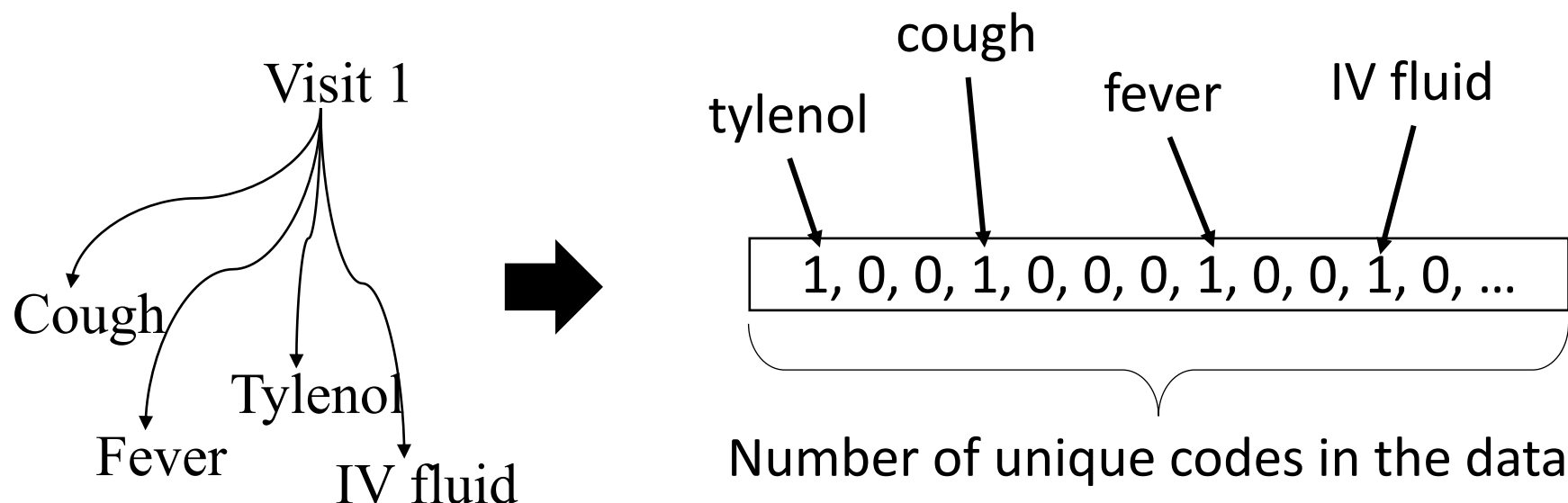
Heart Failure Prediction with RNN

- Input sample x
 - Patient record over time
 - Diagnosis codes, medication codes, procedure codes



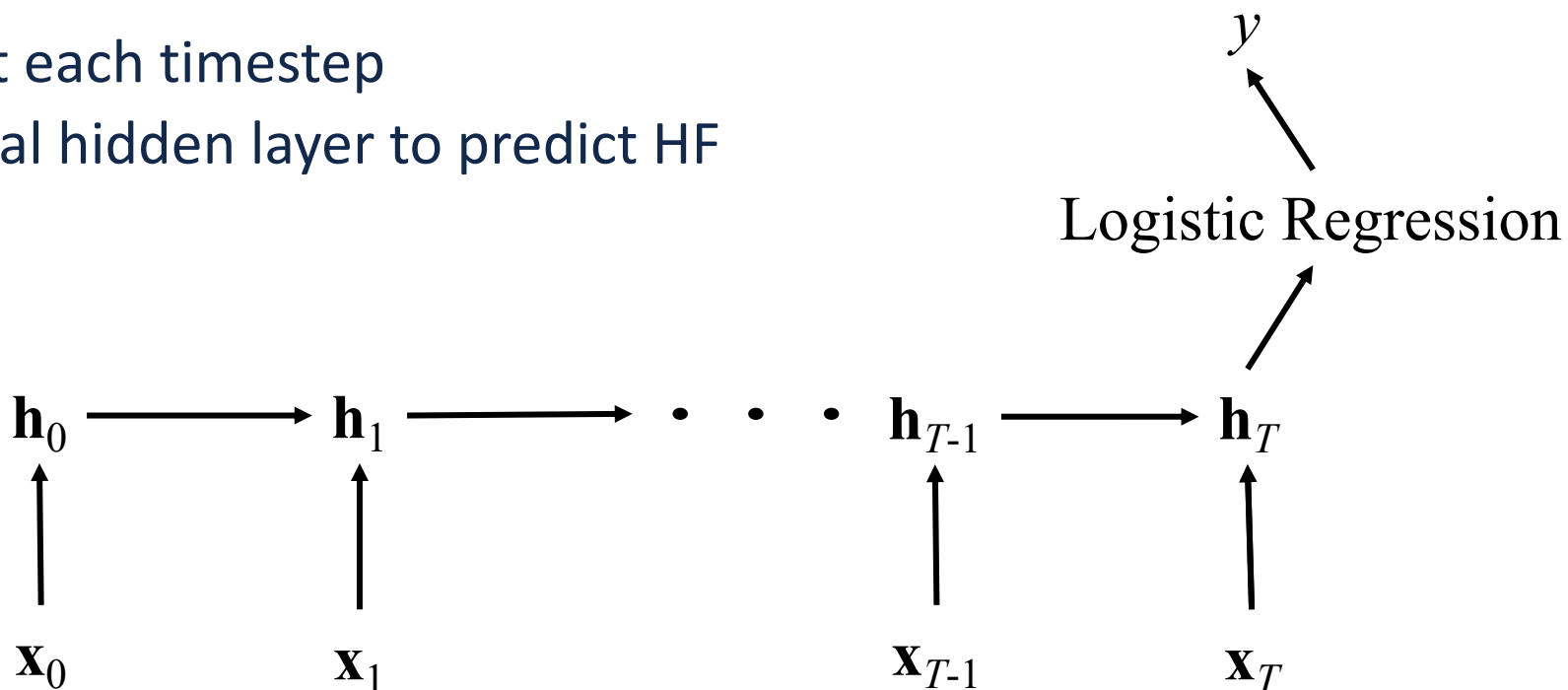
Heart Failure Prediction with RNN

- Input sample x
 - Patient record over time
 - Diagnosis codes, medication codes, procedure codes



Heart Failure Prediction with RNN

- Feed visits into the RNN
 - One visit at each timestep
 - Use the final hidden layer to predict HF



Heart Failure Prediction with RNN



Data

34K patients from Sutter Health
4K cases, 30K controls
18-months observation window

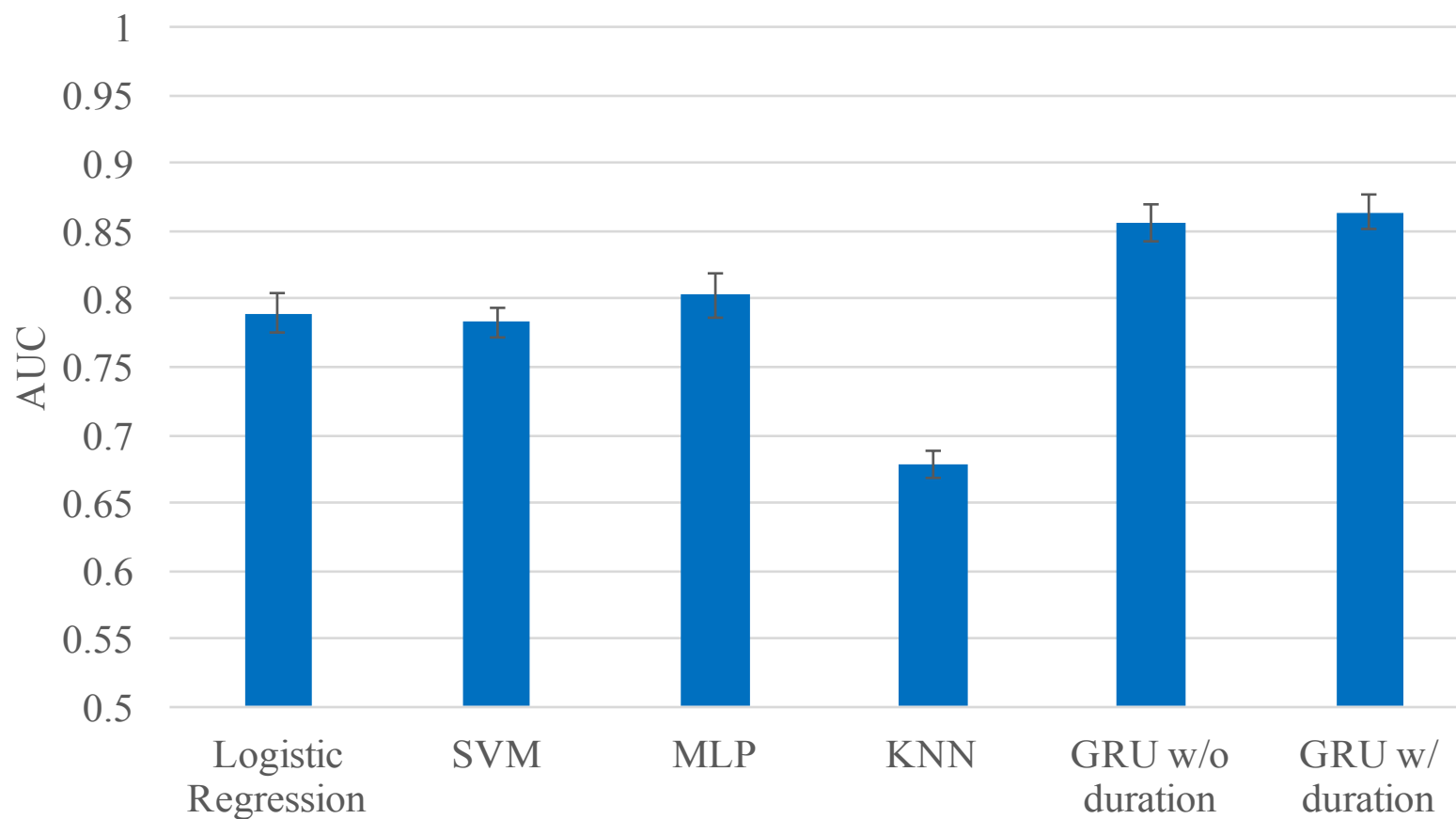


Case-control selection criteria

Age (40-85)
Types of diagnoses received
Number of hospital visits
Time span between diagnoses

Heart Failure Prediction with RNN

- Prediction performance



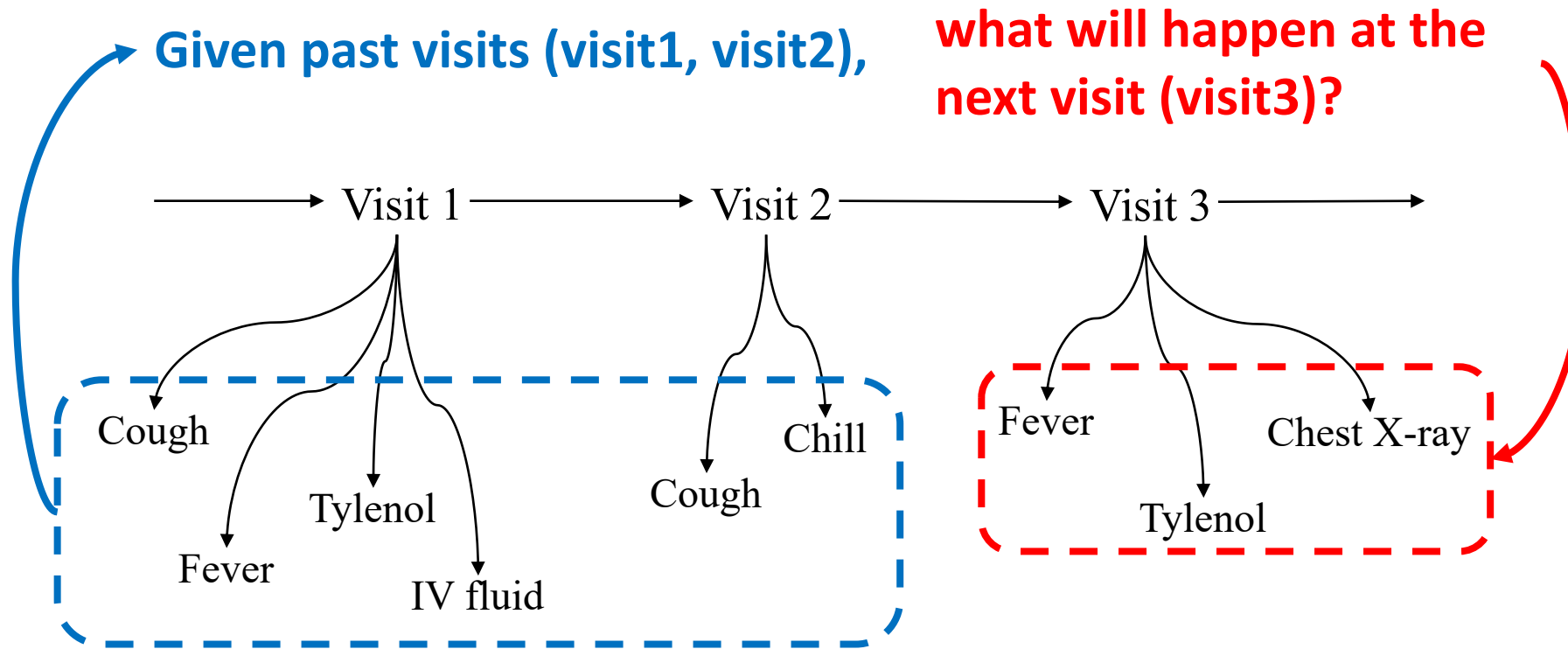
Doctor AI: Predicting Clinical Events via Recurrent Neural Networks

Edward Choi, Mohammad Taha Bahadori, Andy Schuetz, Walter F. Stewart, Jimeng Sun

Machine Learning for Healthcare Conference, 2016

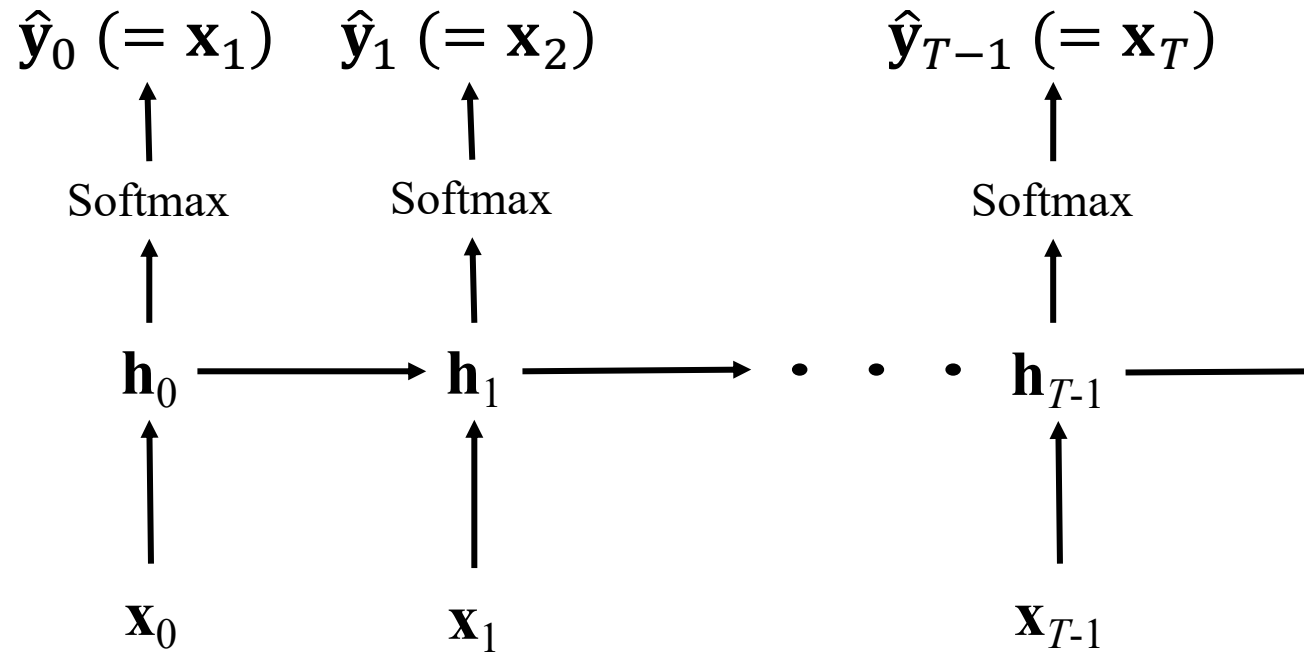
Doctor AI: Background

- Disease progression modeling



Doctor AI: Model

- Feed visits into the RNN
 - One visit at each timestep.
 - Predict next events at each timestep.



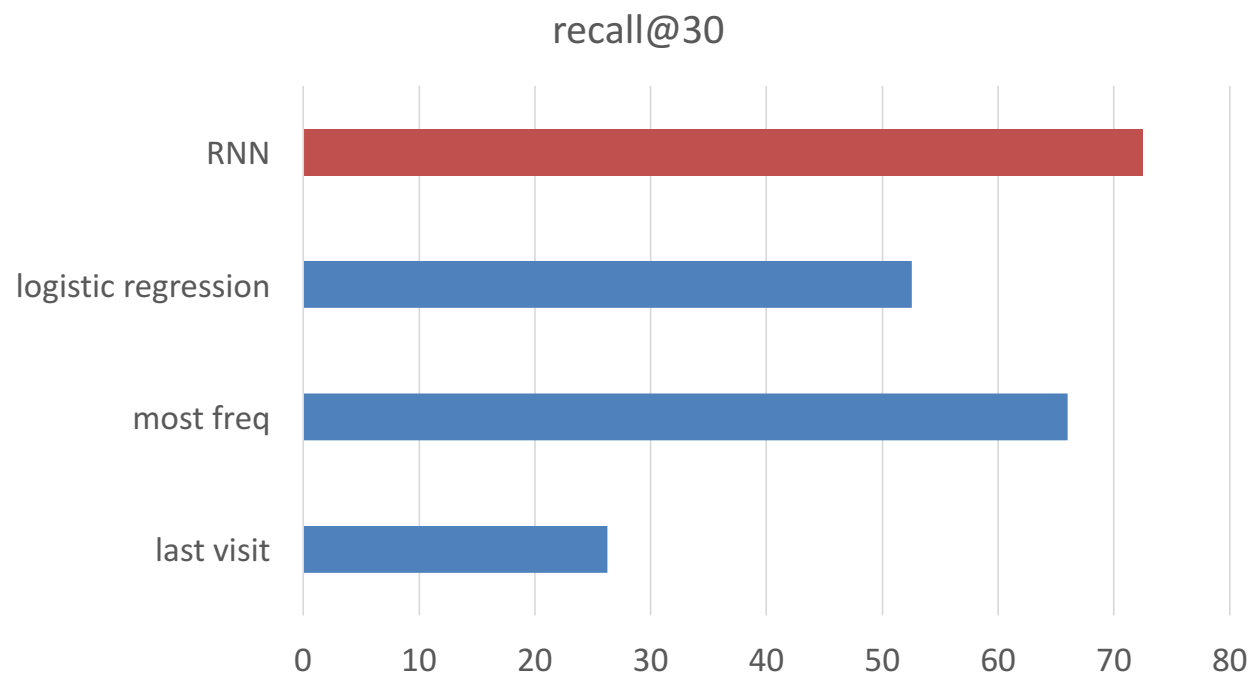
Doctor AI: Data



- 260K patients from Sutter Health
- Patient records over 10 years
- Input codes
 - Diagnosis codes, medication codes, procedure codes (38,000 codes)
- Output labels
 - 1,183 diagnosis codes

Doctor AI: Sequential Prediction

- Predicting diagnoses in the next visit



$$\text{top-}k \text{ recall} = \frac{\# \text{ of true positives in the top } k \text{ predictions}}{\# \text{ of true positives}}$$

Doctor AI: Knowledge Transfer

- Generalize RNN model from one hospital to another

