

Course number: 80240743

Deep Learning

Xiaolin Hu (胡晓林) & Jun Zhu (朱军)

Dept. of Computer Science and Technology

Tsinghua University

Topic 6: Recurrent Neural Networks

Xiaolin Hu

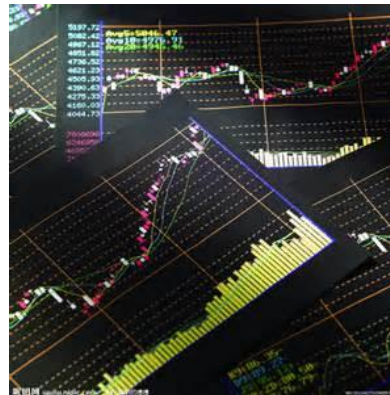
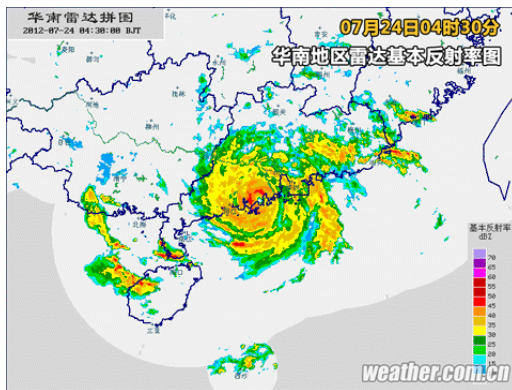
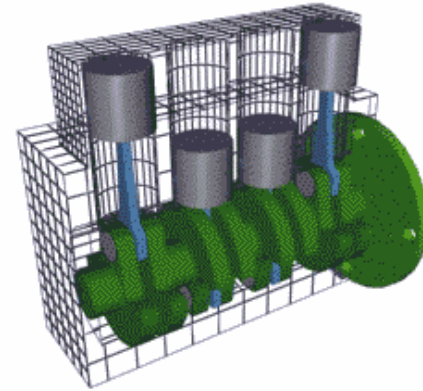
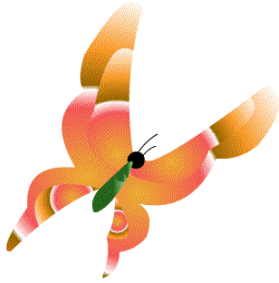
Dept. of Computer Science and
Technology

Tsinghua University

Outline

- Dynamic systems
- Simple RNNs
- Recurrent CNN
- Gated RNNs

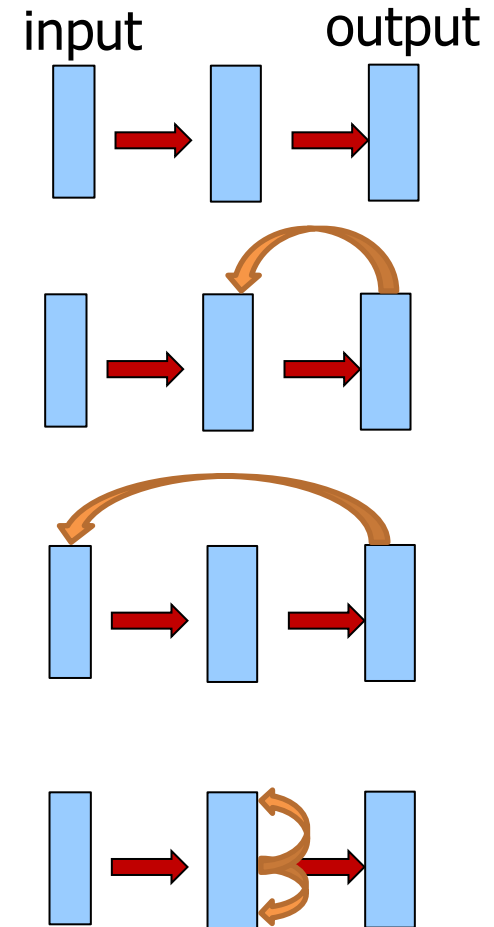
Dynamic systems



Feedback connections

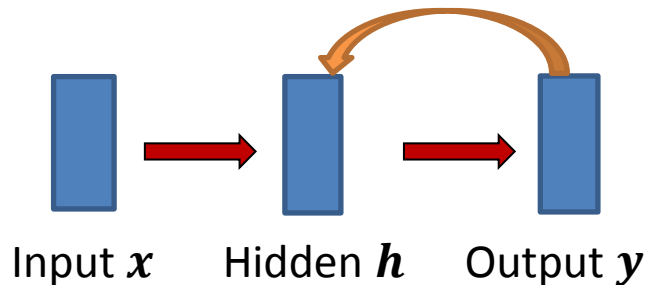
- Feedforward networks
 - No feedback
- Recurrent networks
 - **Between-layer** feedback: from output layer to input layer, or from hidden layer to input layer
 - **Within-layer** feedback

With feedback connections, the state (and therefore output) of neurons will change over time because their current input contains output of some neurons at the last time step which will change their current output



RNNs are dynamic systems

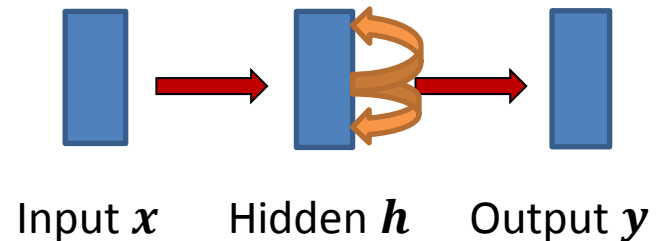
Jordan network



$$\begin{aligned} \mathbf{h}(t) &= \sigma_h(\mathbf{W}_h \mathbf{x}(t) + \mathbf{U}_h \mathbf{y}(t-1) + \mathbf{b}_h) \\ \mathbf{y}(t) &= \sigma_y(\mathbf{W}_y \mathbf{h}(t) + \mathbf{b}_y) \end{aligned}$$

where \mathbf{x} is input, \mathbf{h} is hidden state and \mathbf{y} is output, σ_h and σ_y are activation functions, \mathbf{W} , \mathbf{U} , \mathbf{b} are parameters

Elman network



$$\begin{aligned} \mathbf{h}(t) &= \sigma_h(\mathbf{W}_h \mathbf{x}(t) + \mathbf{U}_h \mathbf{h}(t-1) + \mathbf{b}_h) \\ \mathbf{y}(t) &= \sigma_y(\mathbf{W}_y \mathbf{h}(t) + \mathbf{b}_y) \end{aligned}$$

where \mathbf{x} is input, \mathbf{h} is hidden state and \mathbf{y} is output, σ_h and σ_y are activation functions, \mathbf{W} , \mathbf{U} , \mathbf{b} are parameters

RNNs in general

- The states of the neurons in RNN can be written as

$$\mathbf{h}(t + 1) = f(\mathbf{h}(t), \mathbf{x}(t))$$

where \mathbf{h} denotes the states of *all neurons*, \mathbf{x} denotes input to the network, and f is a nonlinear function (not activation function)

- Often, the output neurons \mathbf{y} are separated from the above equation

$$\mathbf{y}(t) = g(\mathbf{h}(t), \mathbf{x}(t))$$

where g denotes the output function

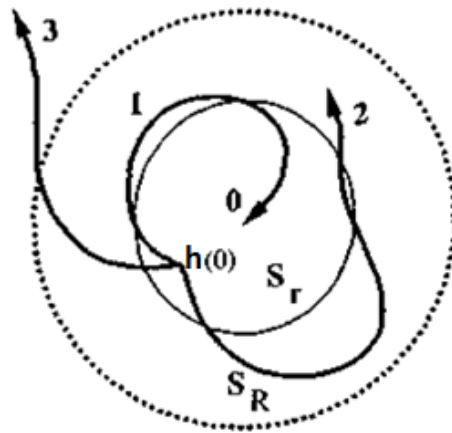
- Such systems are termed **(discrete) dynamic systems**
- **Linear VS nonlinear**: f is linear or nonlinear
- **Autonomous VS non-autonomous**: f doesn't depend on t explicitly or does

Properties of dynamic systems

- For autonomous system

$$\mathbf{h}(t + 1) = f(\mathbf{h}(t))$$

- A point \mathbf{h}^* satisfying $\mathbf{h}^* = f(\mathbf{h}^*)$ is called a **fixed point or equilibrium point**
- Stability



curve 1 - asymptotically stable
curve 2 - marginally stable
curve 3 - unstable

- If $\mathbf{h}^*(t + T) = \mathbf{h}^*(t)$ for some T , $\mathbf{h}^*(t)$, $\mathbf{h}^*(t + 1)$, ..., $\mathbf{h}^*(t + T - 1)$ is called **periodic points**

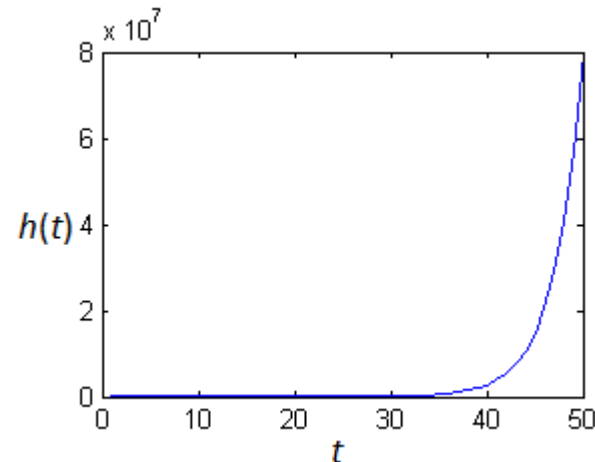
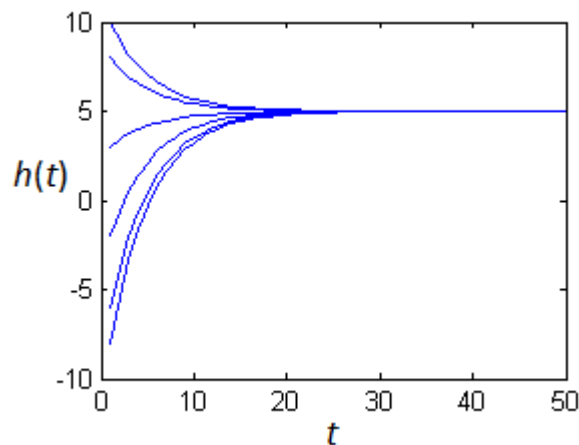
Examples about stability

- For 1D problem, by approximating f with a linear function, we get that a fixed point h^* is **stable** whenever $|f'(h^*)| < 1$
- Consider the linear system

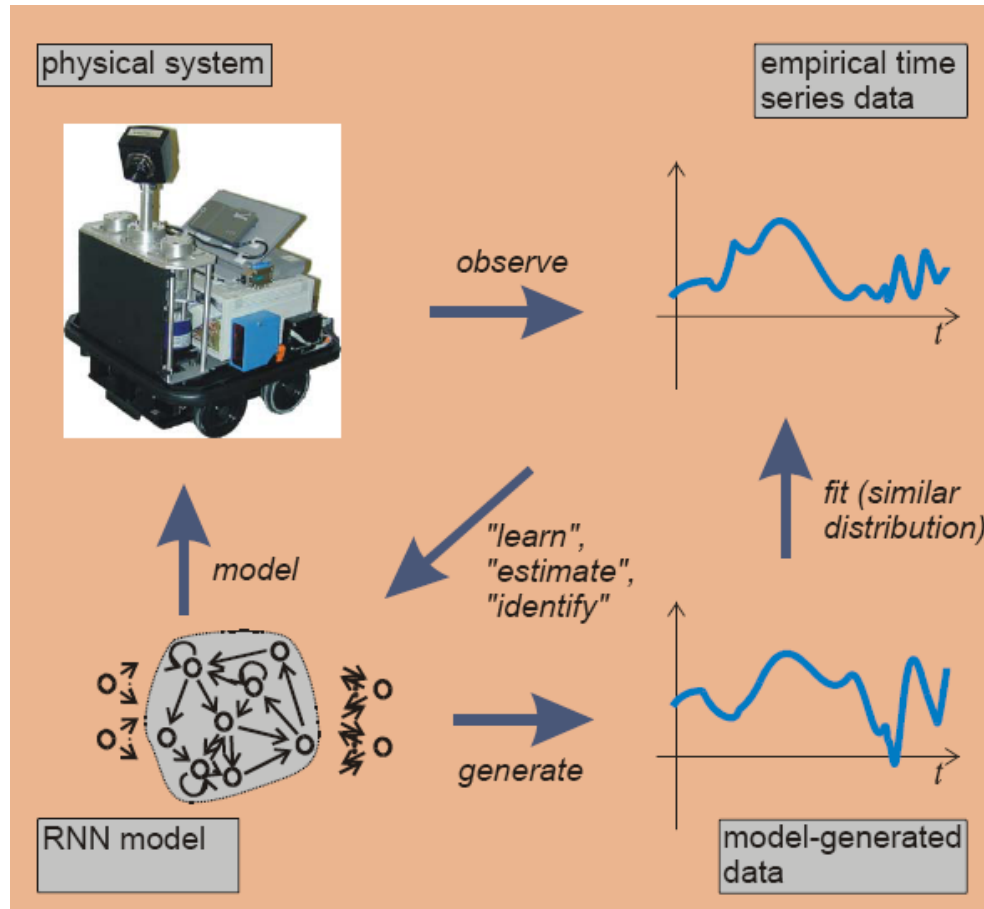
$$h(t+1) = 0.8h(t) + 1$$

It is stable (globally converges to $h^* = 5$)

- But $h(t+1) = 1.4h(t) + 1$ is **unstable**



Modeling dynamic systems with RNN



Why do we need RNN?

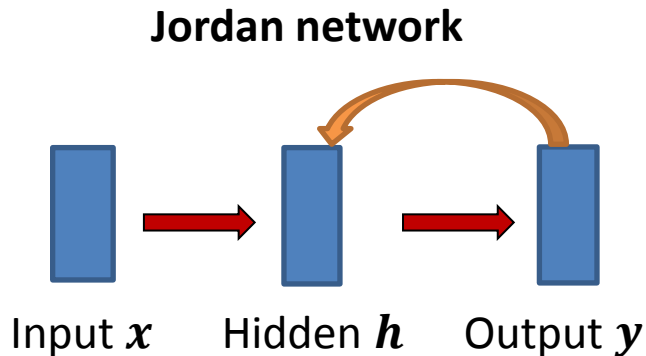
- If at every time step t you always have an input $\mathbf{x}(t)$ and the desired output $\mathbf{r}(t)$
 - Then you have a set of labeled data $\{\mathbf{x}(t), \mathbf{r}(t)\}$ mapping from the input space to the output space
 - You can construct a feedforward model (e.g., MLP) to learn this mapping
 - But you cannot expect that the approach works well as it ignores the dynamic nature of the system under study
 - Different **orders** of inputs may induce totally different outputs (not only differing in the order)!
- Sometimes, you only have input at the beginning, but the desired output at different time is different

We use another dynamic system (RNN) to approximate the real dynamic system!

Outline

- Dynamic systems
- Simple RNNs
- Recurrent CNN
- Gated RNNs

Jordan network



- Dynamic system:

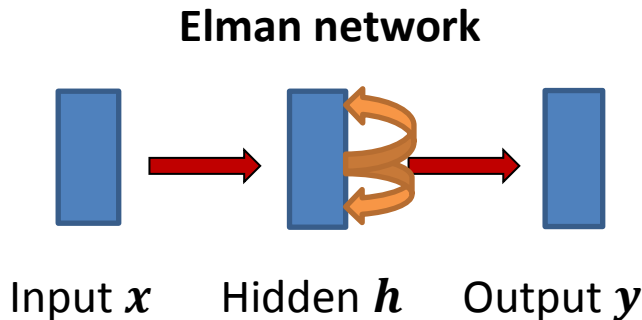
$$\mathbf{h}(t) = \sigma_h(\mathbf{W}_h \mathbf{x}(t) + \mathbf{U}_h \mathbf{y}(t-1) + \mathbf{b}_h)$$

$$\mathbf{y}(t) = \sigma_y(\mathbf{W}_y \mathbf{h}(t) + \mathbf{b}_y)$$

where \mathbf{x} is input, \mathbf{h} is hidden state and \mathbf{y} is output, σ_h and σ_y are activation functions, \mathbf{W} , \mathbf{U} , \mathbf{b} are parameters

- Michael I Jordan
 - BS in Psychology in 1978 from the Louisiana State University,
 - MS in Mathematics in 1980 from Arizona State University
 - PhD in Cognitive Science in 1985 from the UCSD
- At UCSD, Jordan was a student of [David Rumelhart](#)
- Now at UC Berkeley

Elman network



January 22, 1948 – June 28, 2018

- Dynamic system:

$$\mathbf{h}(t) = \sigma_h(\mathbf{W}_h \mathbf{x}(t) + \mathbf{U}_h \mathbf{h}(t-1) + \mathbf{b}_h)$$

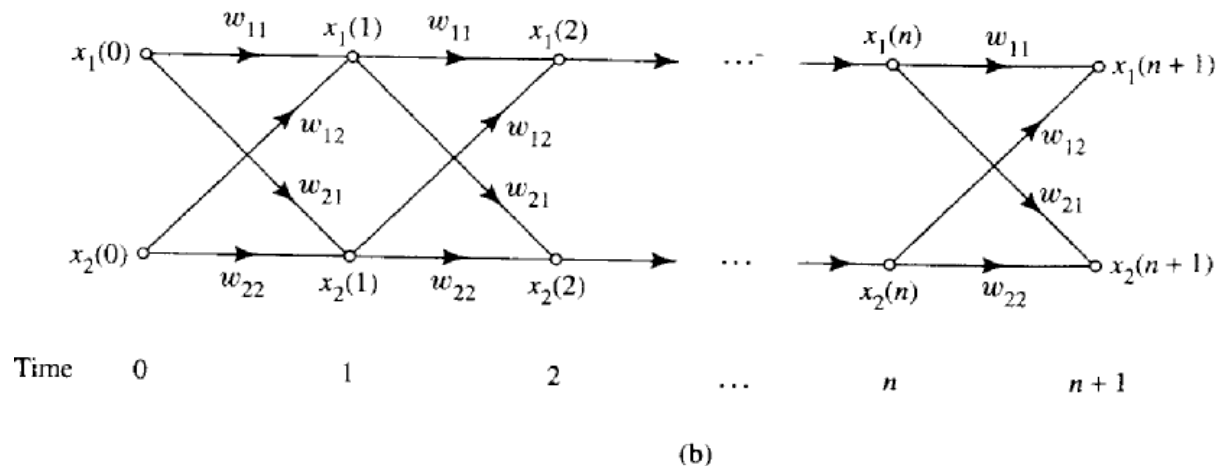
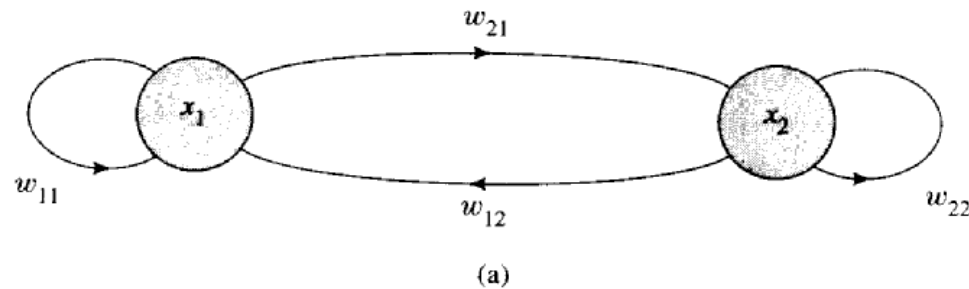
$$\mathbf{y}(t) = \sigma_y(\mathbf{W}_y \mathbf{h}(t) + \mathbf{b}_y)$$

where \mathbf{x} is input, \mathbf{h} is hidden state and \mathbf{y} is output, σ_h and σ_y are activation functions, \mathbf{W} , \mathbf{U} , \mathbf{b} are parameters

- Jeffrey Elman
 - BS in 1969 from Harvard University
 - Ph.D. in 1977 from the University of Texas at Austin
- Professor of cognitive science at the UCSD

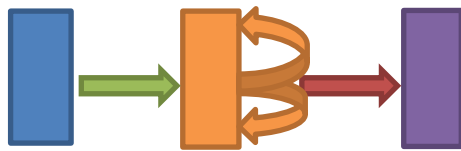
Back-propagation through time (BPTT)

Unfold the temporal operation of the network into a layered feedforward network, the topology of which grows by one layer at every time step.



Unfold the Elman network

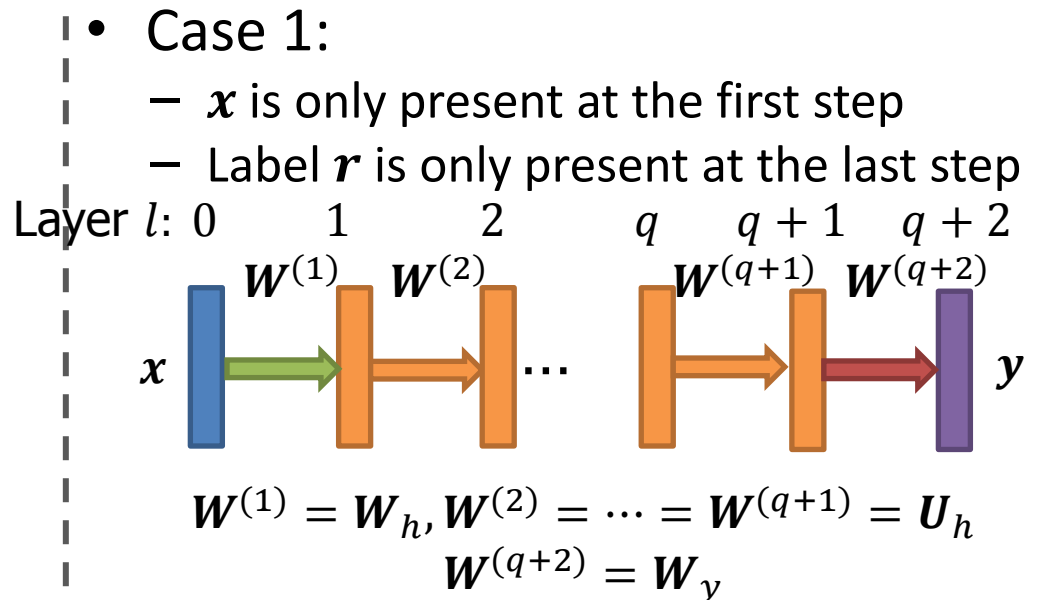
Unfold for q times



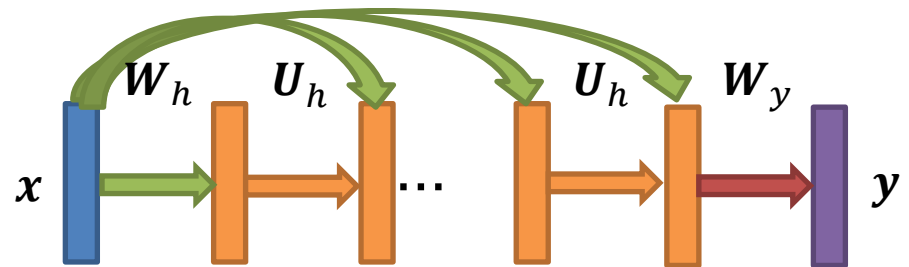
Input x Hidden h Output y

$$h(t) = \sigma_h(W_h x(t) + U_h h(t-1) + b_h)$$

$$y(t) = \sigma_y(W_y h(t) + b_y)$$



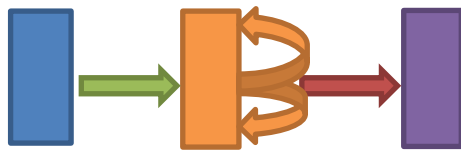
- Case 2:
 - x is fixed but present at all steps
 - Label r is only present at the last step
 - E.g., image classification (Liang, Hu, CVPR 2015)



(Arrows in the same color share weights)₁₆

Unfold the Elman network

Unfold for q times

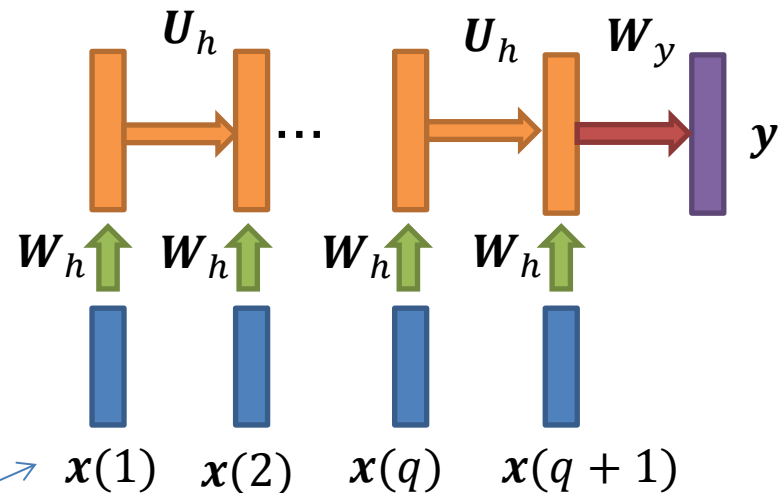


Input x Hidden h Output y

$$h(t) = \sigma_h(W_h x(t) + U_h h(t-1) + b_h)$$

$$y(t) = \sigma_y(W_y h(t) + b_y)$$

- Case 3:
 - x is time-varying
 - Label r is only present at the last step
 - E.g., sentence classification

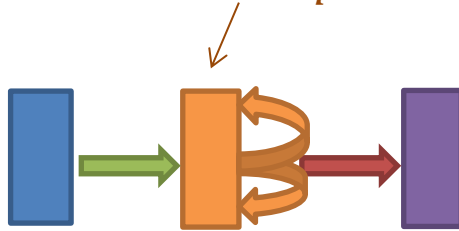


This can be viewed as layer 0
attached to the orange backbone

(Arrows in the same color share weights)

Unfold the Elman network

Unfold for q times



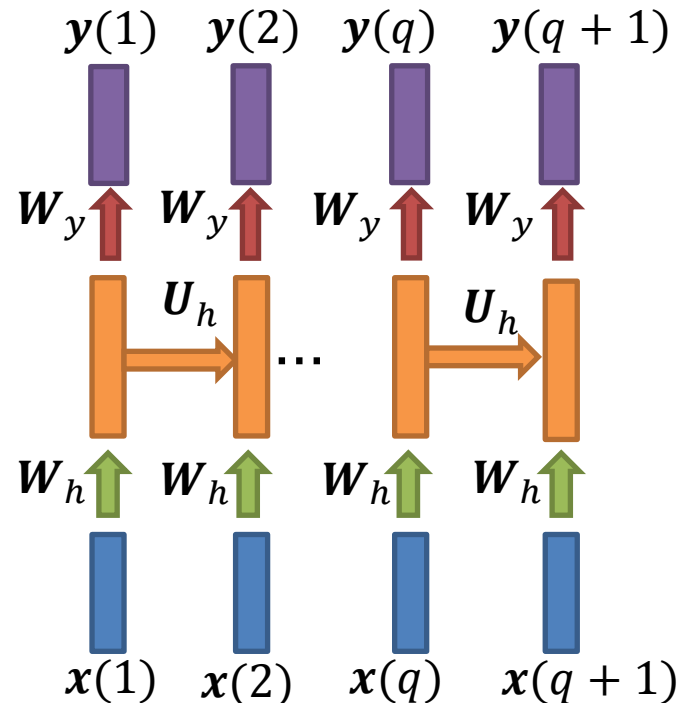
Input x Hidden h Output y

$$h(t) = \sigma_h(W_h x(t) + U_h h(t-1) + b_h)$$

$$y(t) = \sigma_y(W_y h(t) + b_y)$$

Do you know
other cases?

- Case 4:
 - x is time-varying
 - Label r is present at all steps
 - E.g., speech recognition

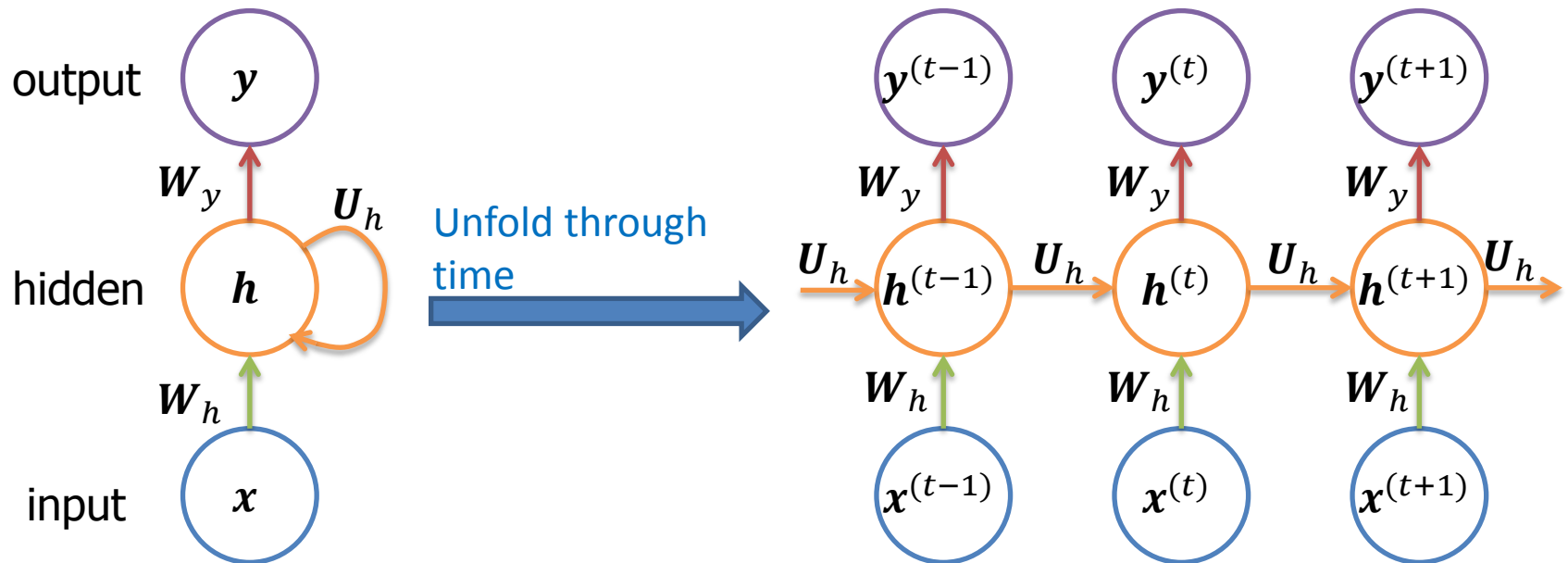


(Arrows in the same color share weights)

Simplified illustration (Elman network)

- Use circles to represent vectors (one circle one layer)
- Put time step into superscript

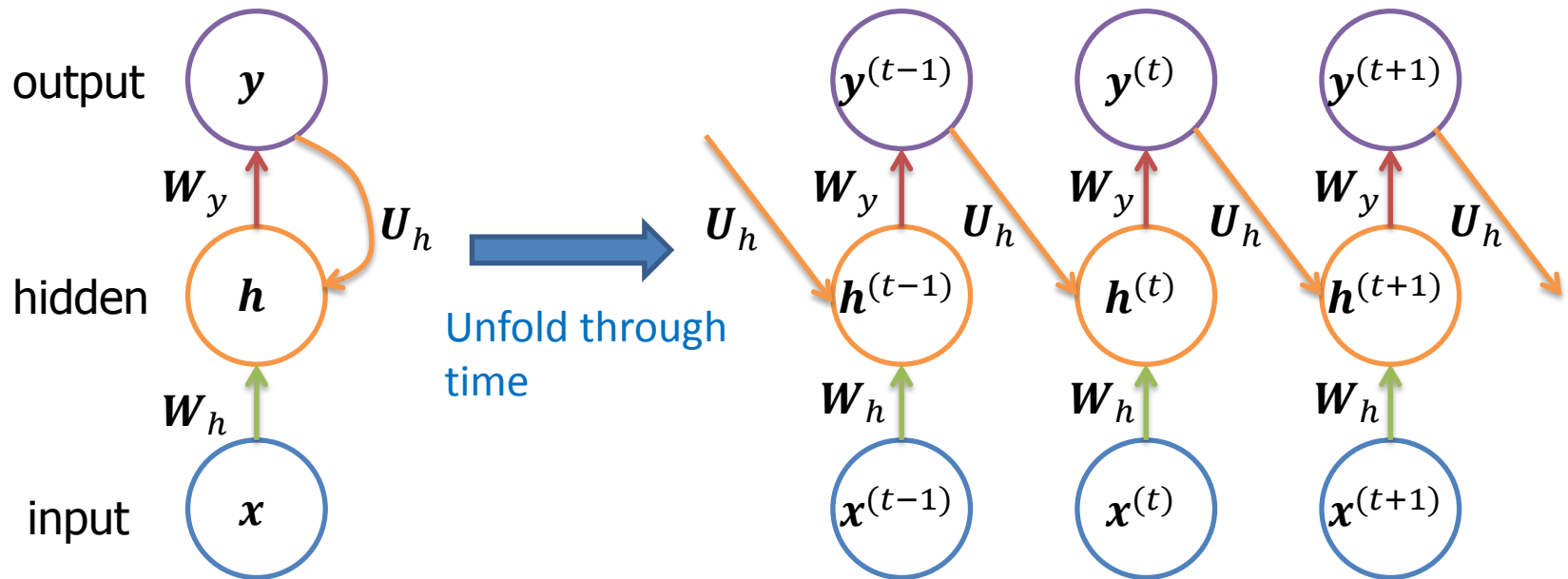
There are reference $r(t)$ at each time step



- The forward propagation runtime is $O(q)$ and cannot be reduced
- Equivalent to a **Turing machine** (due to the hidden-to-hidden) connections

Unfold the Jordan network

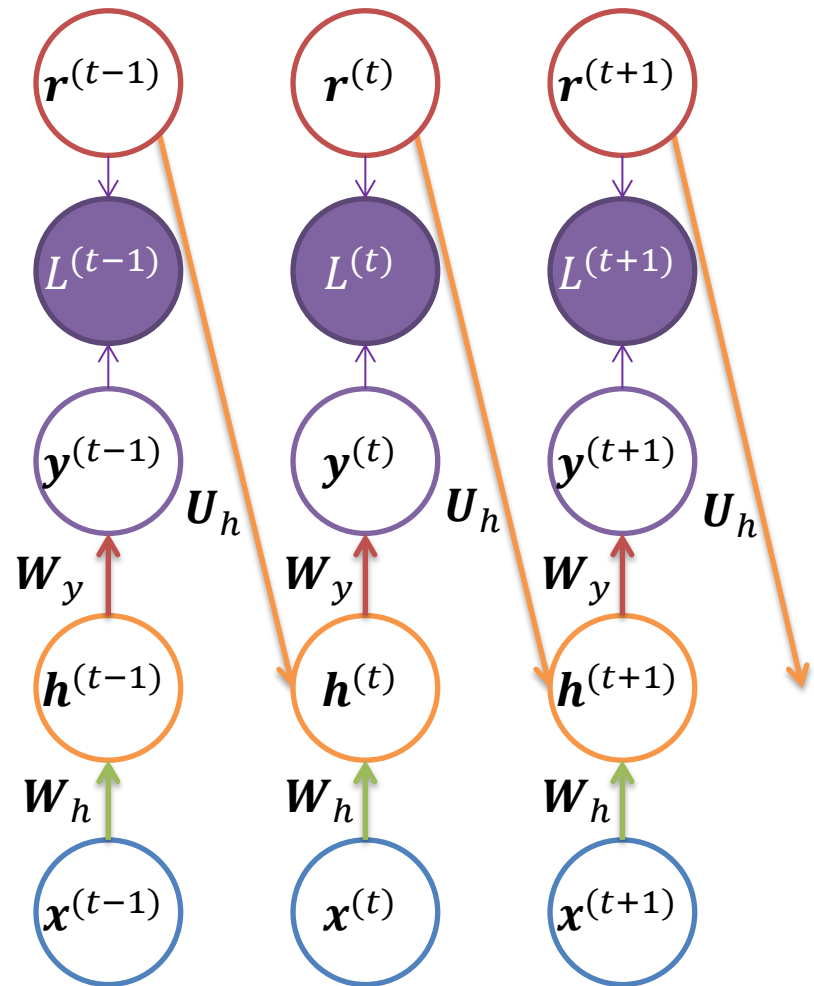
There are reference $\mathbf{r}(t)$ at each time step



- It is strictly **less powerful** and it cannot simulate a universal Turing machine
- If the loss is based on comparing $\mathbf{y}(t)$ and $\mathbf{r}(t)$, all time steps are decoupled and training can be parallelized

Teacher forcing

- Some networks such as Jordan network, have connections from the **output** at one time step to values computed in the next time step
- Then, what should be input to the next time step to represent the **output**?
- **Teacher forcing**: in training, we use the **reference signal**



Teacher forcing

- However, **in testing**, there is **no reference signal** and we have to use the network's output at time t
- The kind of inputs that the network sees during training could be quite different from the kind of inputs that it will see at test time
- To mitigate this problem:
 - Alternate use of teacher-forced inputs and free-running inputs for a number of time steps
 - Randomly choose between the teacher-forced input and free-running input at every time step

Bidirectional RNN

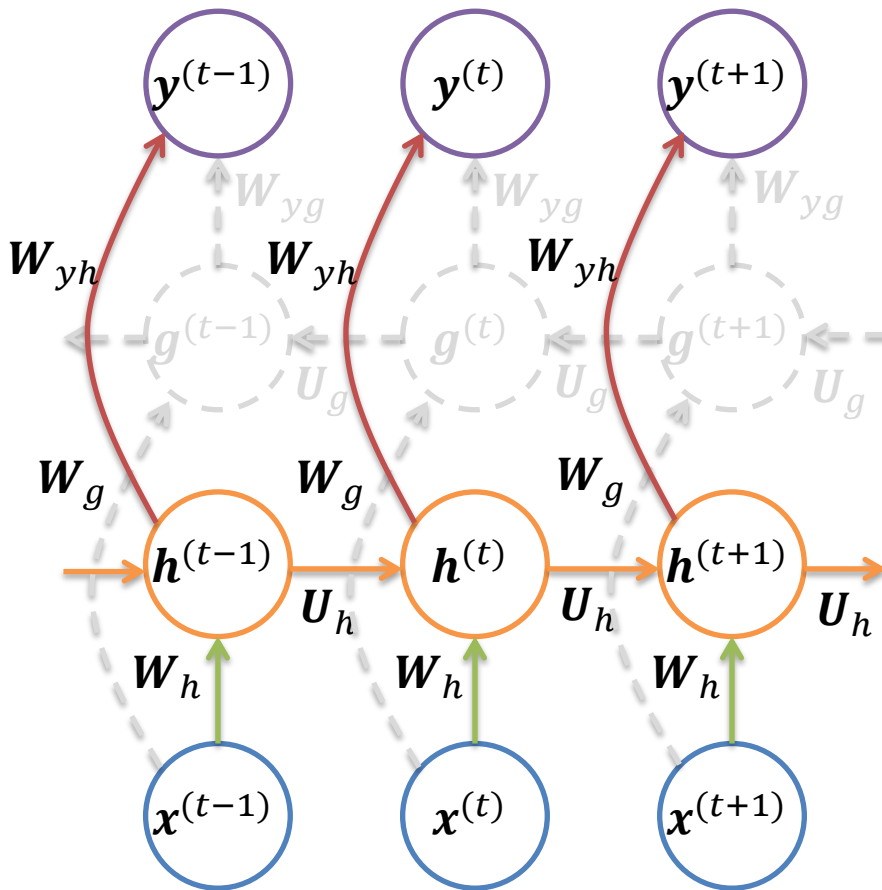
- In many applications, the prediction at the current time step may depend on **the whole input sequence** (both the past and the future)

Bank is the side of a river.

Bank provides various financial services.

- Bidirection RNNs combine an RNN that moves forward through time with another RNN that moves backward through time
- The output of the entire network at every time step then receives **two inputs**

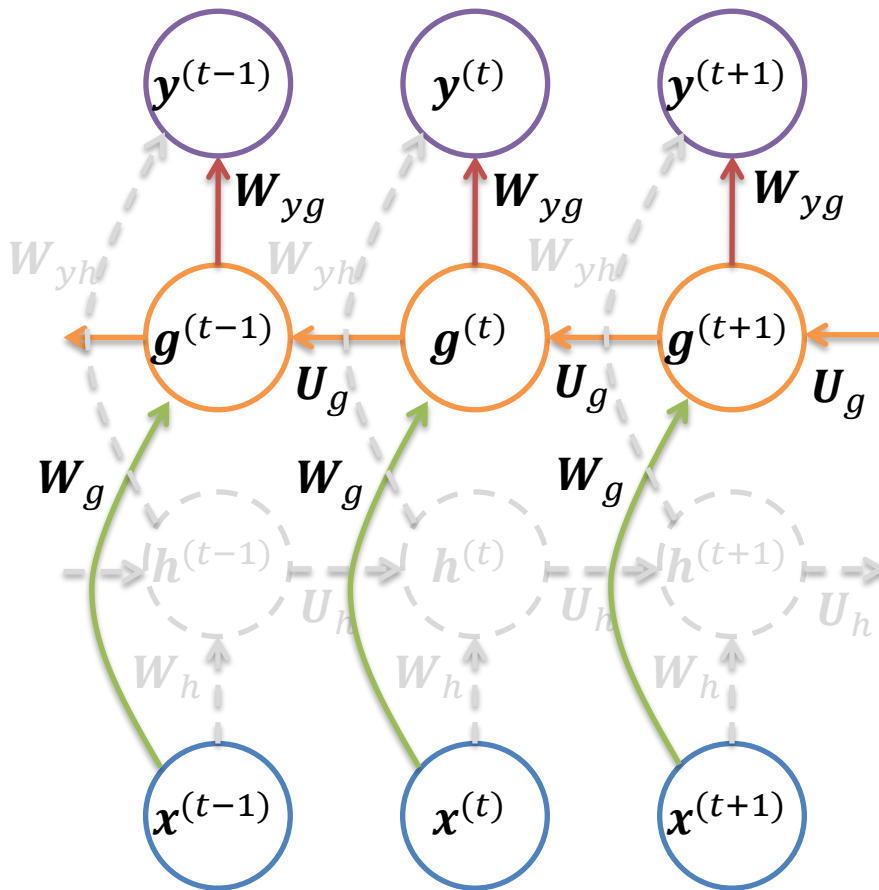
Bidirectional RNN



W_{yh} : connection weights from the sub-RNN g to output

One sub-RNN moves forward, same as before

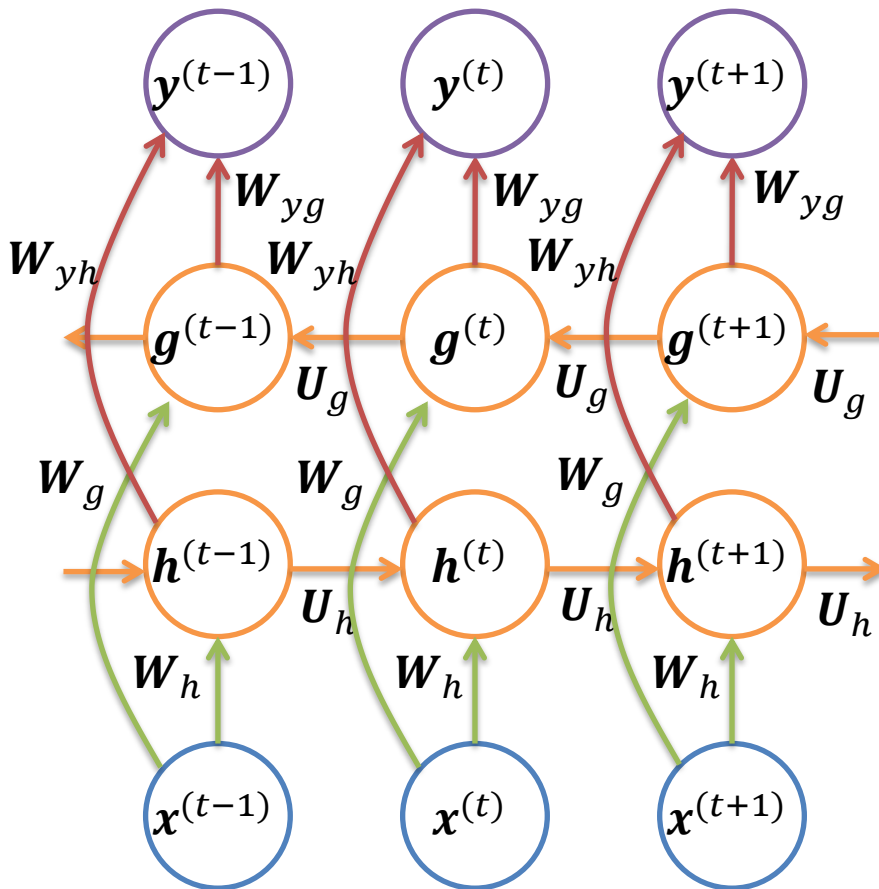
Bidirectional RNN



W_{yg} : connection weights from the sub-RNN g to output

One sub-RNN moves backward

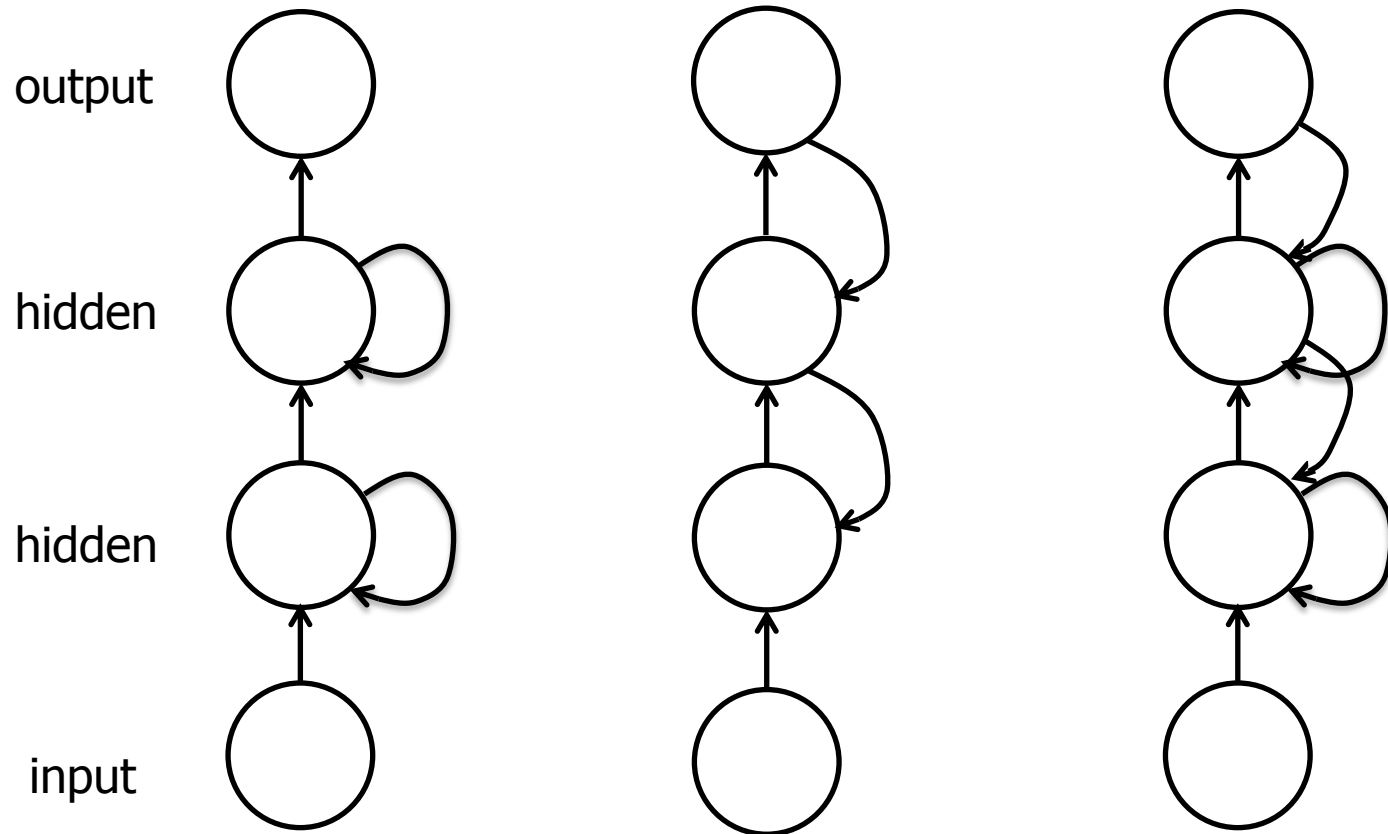
Bidirectional RNN



The output $y^{(t)}$ receives input from both sub-RNNs via W_{yh} and W_{yg}

Deep RNNs

- Many ways to construct deep RNNs



Challenges

- Simple RNNs tend to converge to **fixed points** or **explode**
- Recall the Elman network:


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

- Suppose σ_h is an identity mapping, $\mathbf{b} = \mathbf{0}$ and the input is only present at the beginning. After t steps from zero

$$\mathbf{h}(t) = \mathbf{U}^t \mathbf{h}(0)$$

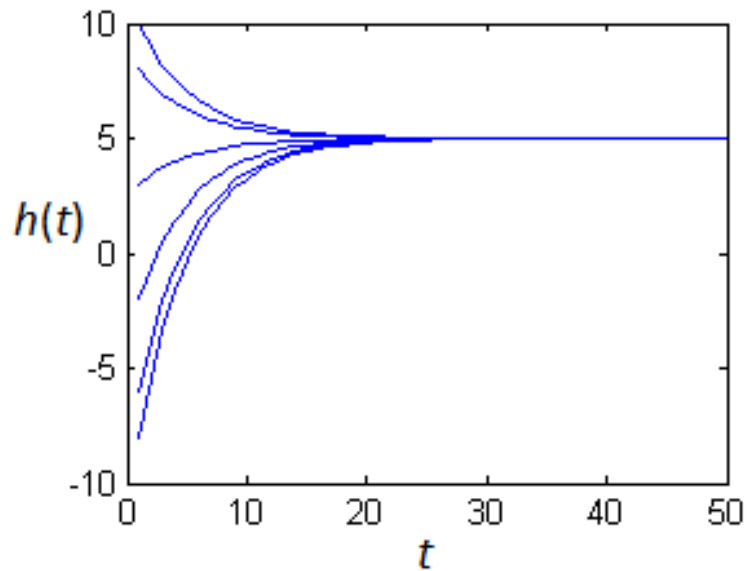
- Suppose \mathbf{U}^t has an eigen-decomposition $\mathbf{U} = \mathbf{V}\text{diag}(\boldsymbol{\lambda})\mathbf{V}^{-1}$, then

$$\mathbf{U}^t = (\mathbf{V}\text{diag}(\boldsymbol{\lambda})\mathbf{V}^{-1})^t = \mathbf{V}\text{diag}(\boldsymbol{\lambda})^t\mathbf{V}^{-1}$$

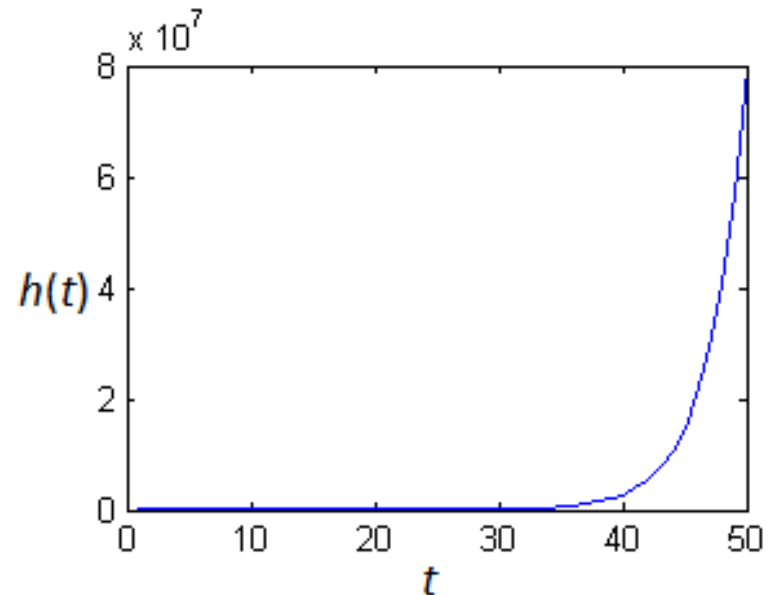
- Any eigenvalue λ_i that is not near an **absolute value of 1** will either explode or vanish
- Then the components of $\mathbf{h}(t)$ would
 - either **converge to zero**  If \mathbf{b} is not zero, then converge to nonzero values
 - or **explode**

Recall the 1D case

$$h(t+1) = 0.8h(t) + 1$$



$$h(t+1) = 0.8h(t) + 1$$



Challenges (backward pass)

- Recall the BP algorithm for MLP

$$\delta^{(l)} = (\mathbf{W}^{(l+1)})^\top \delta^{(l+1)} \odot \mathbf{f}'(\mathbf{u}^{(l)})$$

- For the Elman network, \mathbf{W} changes to \mathbf{U}
- With the same assumption as in the previous slides, we have

$$\delta^{(T-t)} = (\mathbf{U}^\top)^t \delta^{(T)}$$

where T denotes the last time step

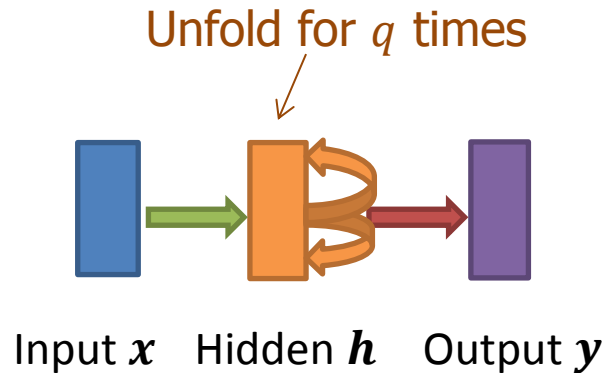
- Suppose \mathbf{U}^\top has an eigen-decomposition $\mathbf{U}^\top = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}$, then
$$(\mathbf{U}^\top)^t = (\mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1})^t = \mathbf{V} \text{diag}(\boldsymbol{\lambda})^t \mathbf{V}^{-1}$$
- Any eigenvalue λ_i that is not near an **absolute value of 1** will either explode or vanish
- Then the gradient either **vanishes** or **explodes**

Even if the input is always present, the challenges exist

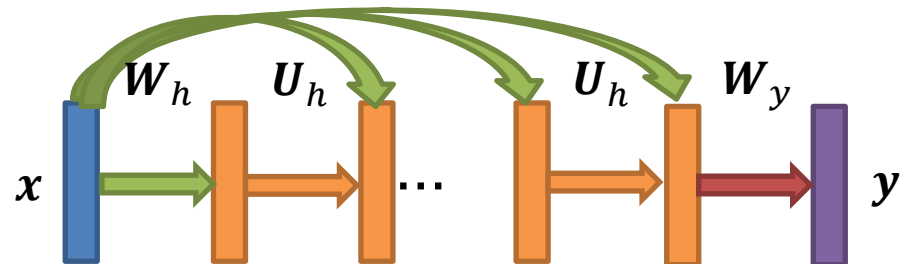
Outline

- Dynamic systems
- Simple RNNs
- Recurrent CNN
- Gated RNNs

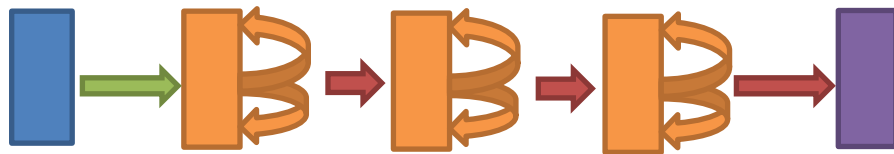
Recall one case of unfolding the Elman network



- Case 2:
 - x is fixed but present at all steps
 - Label r is only present at the last step



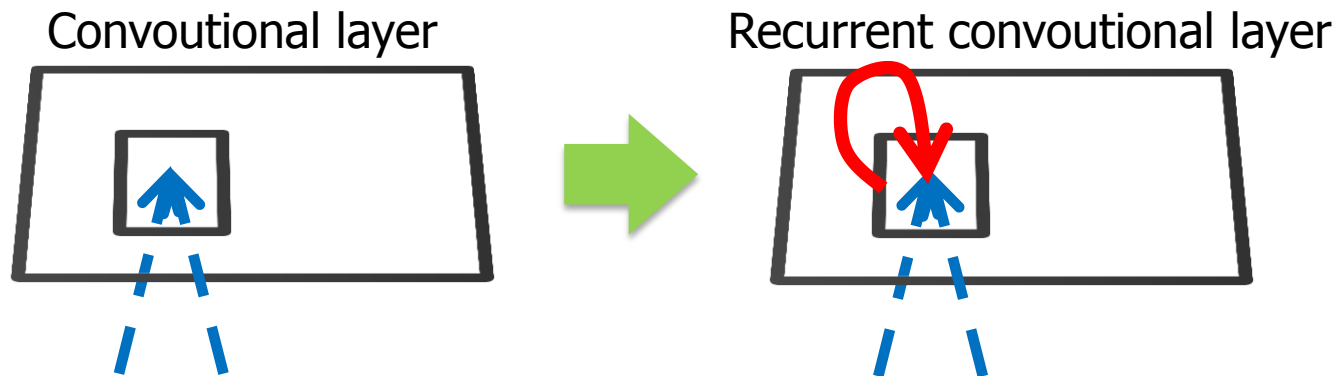
- If you have multiple such layers, it is called a **recurrent MLP**



- If this idea is applied to CNN, then you obtain a **recurrent CNN** (Liang, Hu, CVPR 2015; Liang, Hu, Zhang, NIPS 2015; Zhao et al., ICASSP 2017)

Recurrent convolutional layer (RCL)

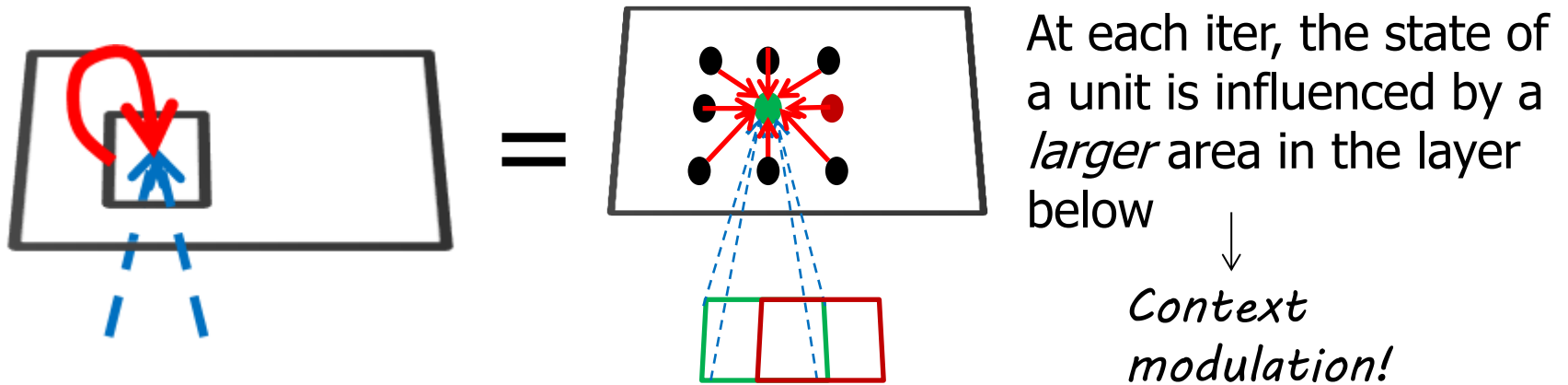
- Incorporate intra-layer recurrent connections into each convolutional layer



Blue: feed-forward connections

Red: recurrent connections

RCL as 2D RNN



It becomes a dynamic system since the states change over time (here discrete iterations)

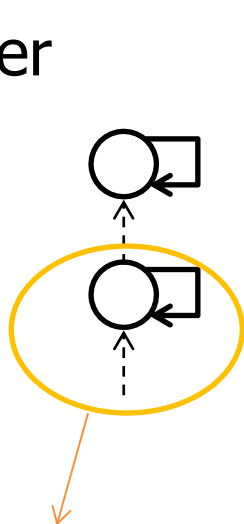
State (time-varying) Input (static)

$$x_{ijk}(t) = \sigma \left((\mathbf{w}_k^{in})^T \mathbf{u}^{(i,j)} + (\mathbf{w}_k^{rec})^T \mathbf{x}^{(i,j)}(t-1) + b_k \right)$$

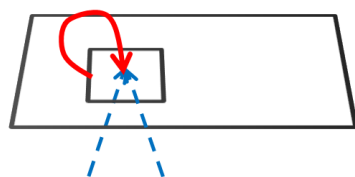
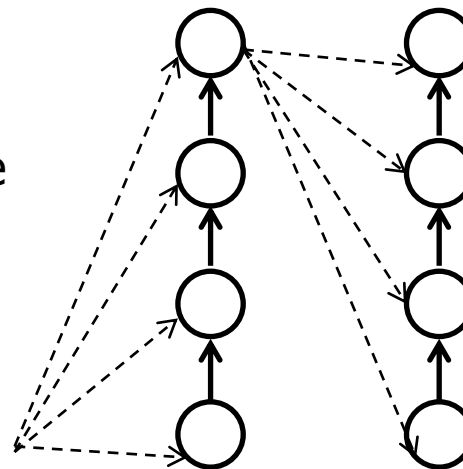
In RCL, both feed-forward and recurrent computation take the form of convolution (local connections)

Stacking RCL's to a deep structure

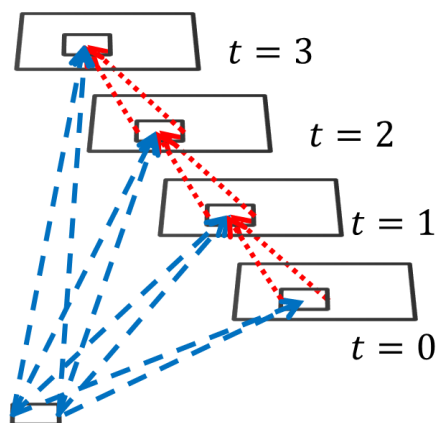
A two-layer RCNN



Unfolding through time



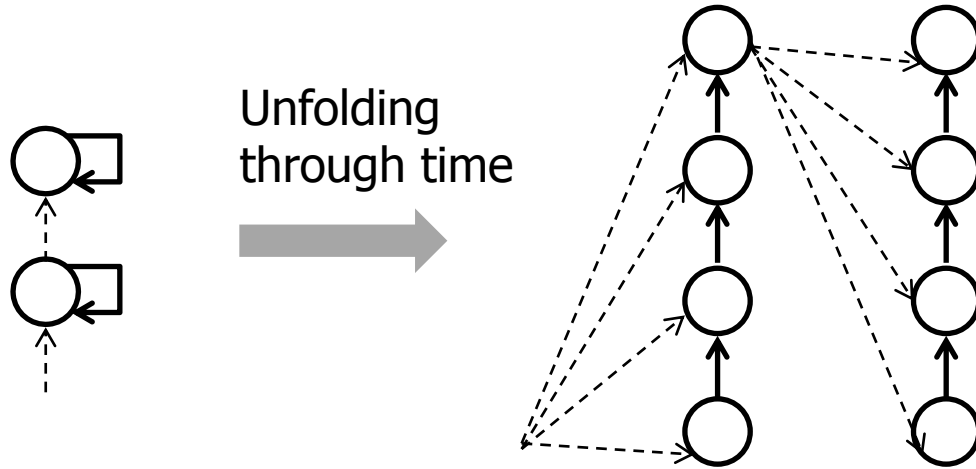
Unfolding RCL through time



$$x_{ijk}(t) = \sigma \left((\mathbf{w}_k^{in})^T \mathbf{u}^{(i,j)} + (\mathbf{w}_k^{rec})^T \mathbf{x}^{(i,j)}(t-1) + b_k \right)$$

RCNN properties

A two-layer
RCNN

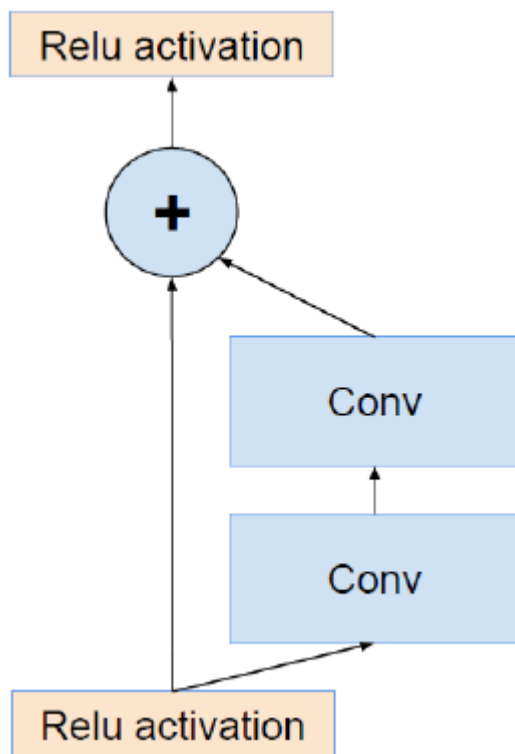


- A deep model with limited number of parameters
- Neurons in each layer can capture global context
- Multi-path structure

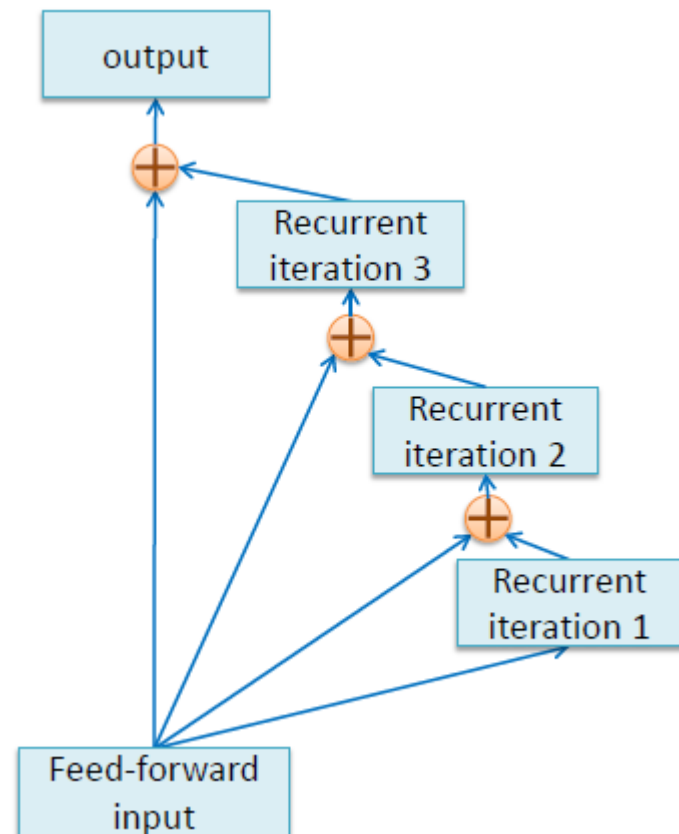
Multi-path in Residual Net

- 152 layers for ImageNet; 100 or 1000 layers for CIFAR-10 (He et al., 2016)

Residual Net



RCL



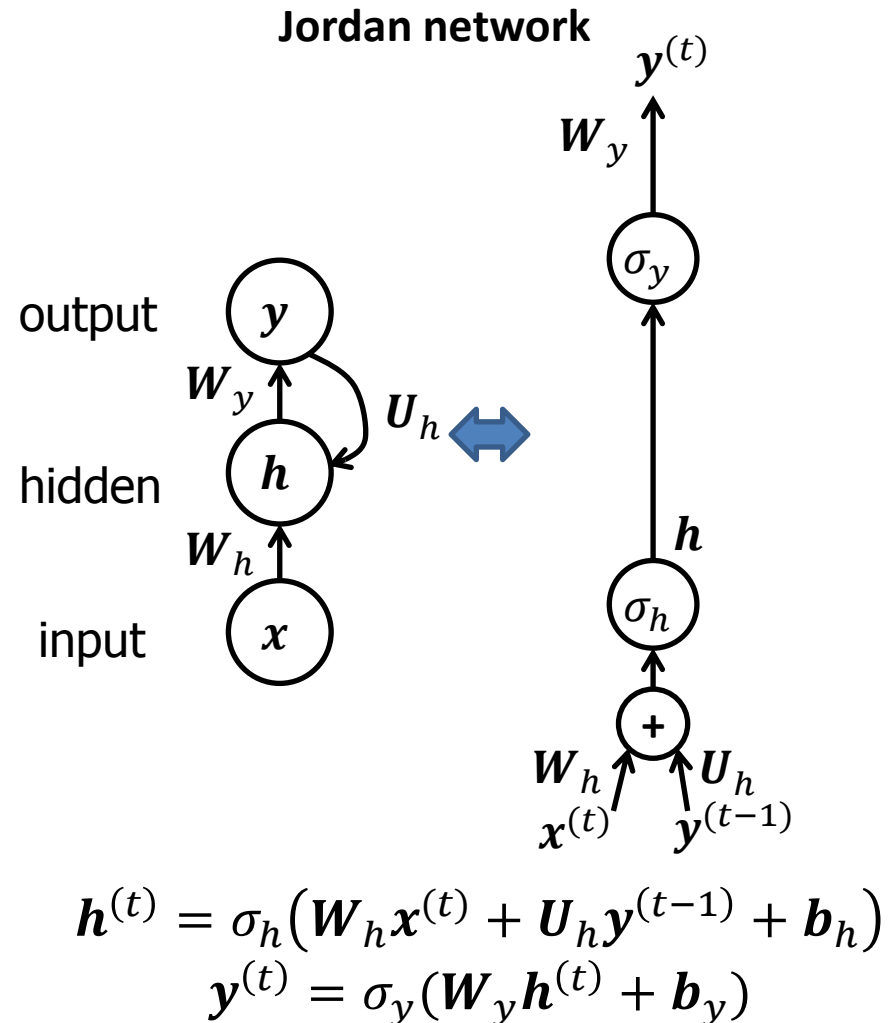
Adapted from Szegedy et al. 2016

Outline

- Dynamic systems
- Simple RNNs
- Recurrent CNN
- Gated RNNs

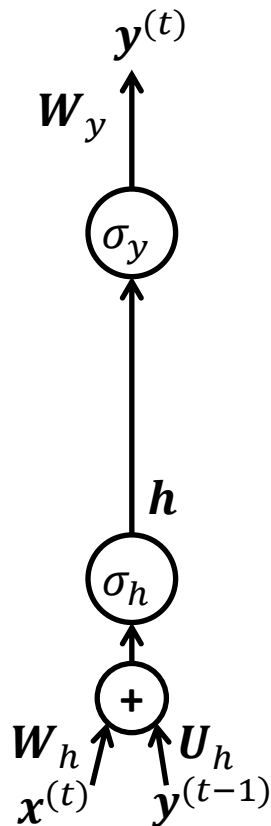
Long short-term memory (LSTM) cell

- It can be viewed as a combination of the Jordan network and the Elman network
 - The output is connect to the input
 - A self-loop is used to capture the information about the past
- Redraw the Jordan network
 - Use circles to denote operations
 - Variables are indicated on arrows



Step 1: add a self-loop

Jordan network



$$h^{(t)} = \sigma_h(W_h x^{(t)} + U_h y^{(t-1)} + b_h)$$

$$y^{(t)} = \sigma_y(W_y h^{(t)} + b_y)$$



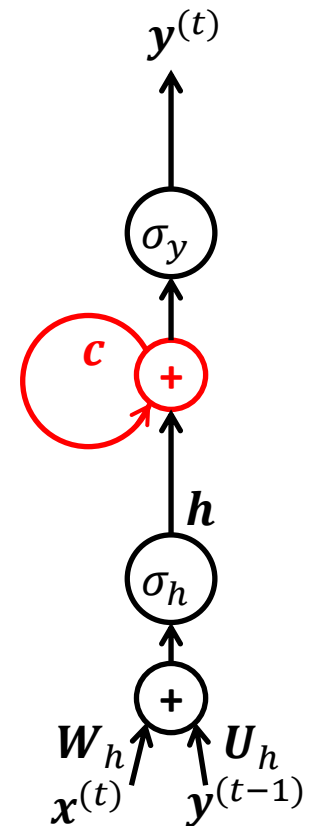
$$h^{(t)} = \sigma_h(W_h x^{(t)} + U_h y^{(t-1)} + b_h)$$

$$\mathbf{c}^{(t)} = \mathbf{c}^{(t-1)} + \mathbf{h}^{(t)}$$

$$y^{(t)} = \sigma_y(\mathbf{c}^{(t)})$$

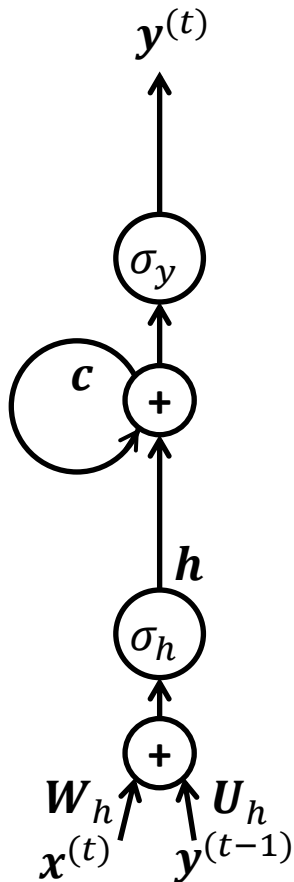
(We eliminate the linear transformation in the output)

- σ_h is either the logistic sigmoid function or tanh function
- σ_y is often tanh function



Step 2: add three gates

Gates are introduced to adaptively control the flow of information

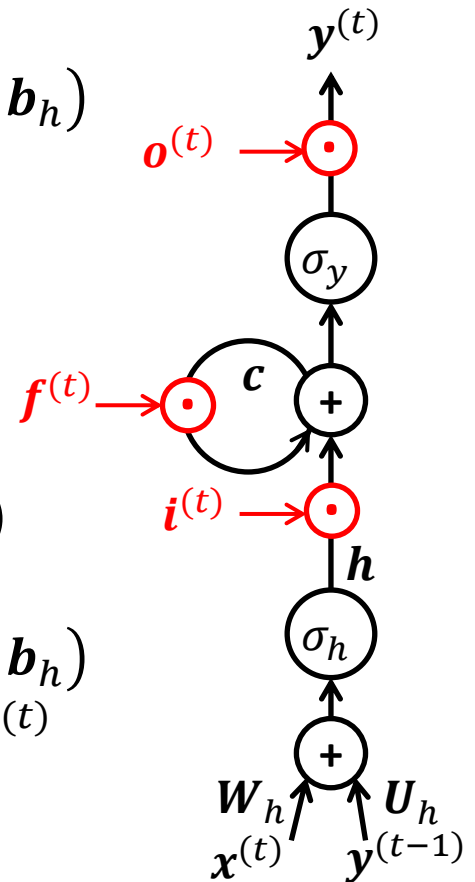


$$\begin{aligned} h^{(t)} &= \sigma_h(W_h x^{(t)} + U_h y^{(t-1)} + b_h) \\ c^{(t)} &= c^{(t-1)} + h^{(t)} \\ y^{(t)} &= \sigma_y(c^{(t)}) \end{aligned}$$

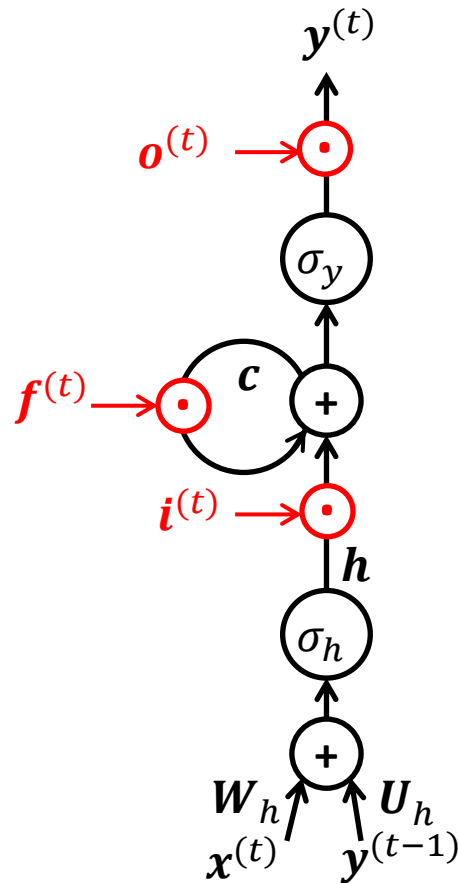


Forget gate $f^{(t)}$, input gate $i^{(t)}$,
output gate $o^{(t)}$: between (0,1)

$$\begin{aligned} h^{(t)} &= \sigma_h(W_h x^{(t)} + U_h y^{(t-1)} + b_h) \\ c^{(t)} &= f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot h^{(t)} \\ y^{(t)} &= o^{(t)} \odot \sigma_y(c^{(t)}) \end{aligned}$$



What determine these gates?



$$h^{(t)} = \sigma_h(W_h x^{(t)} + U_h y^{(t-1)} + b_h)$$

$$c^{(t)} = f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot h^{(t)}$$

$$y^{(t)} = o^{(t)} \odot \sigma_y(c^{(t)})$$

- All of the gates are determined by the input $x^{(t)}$ and $y^{(t-1)}$

$$f^{(t)} = \sigma(W_f x^{(t)} + U_f y^{(t-1)} + b_f)$$

$$i^{(t)} = \sigma(W_i x^{(t)} + U_i y^{(t-1)} + b_i)$$

$$o^{(t)} = \sigma(W_o x^{(t)} + U_o y^{(t-1)} + b_o)$$

where σ is the logistic sigmoid function

- Sometimes, they are also determined by $c^{(t)}$ and $c^{(t-1)}$: *peepholes*

Note: sometimes, the output y is also called *hidden state of LSTM*, especially when LSTM is integrated into a larger system.

Gated recurrent unit (GRU)

- In the Elman network, the **hidden units h** are used to capture the history information

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

- In an LSTM cell without gates, a **new vector c** is introduced for this purpose

$$\mathbf{h}^{(t)} = \sigma_h(\mathbf{W}_h \mathbf{x}^{(t)} + \mathbf{U}_h \mathbf{y}^{(t-1)} + \mathbf{b}_h)$$

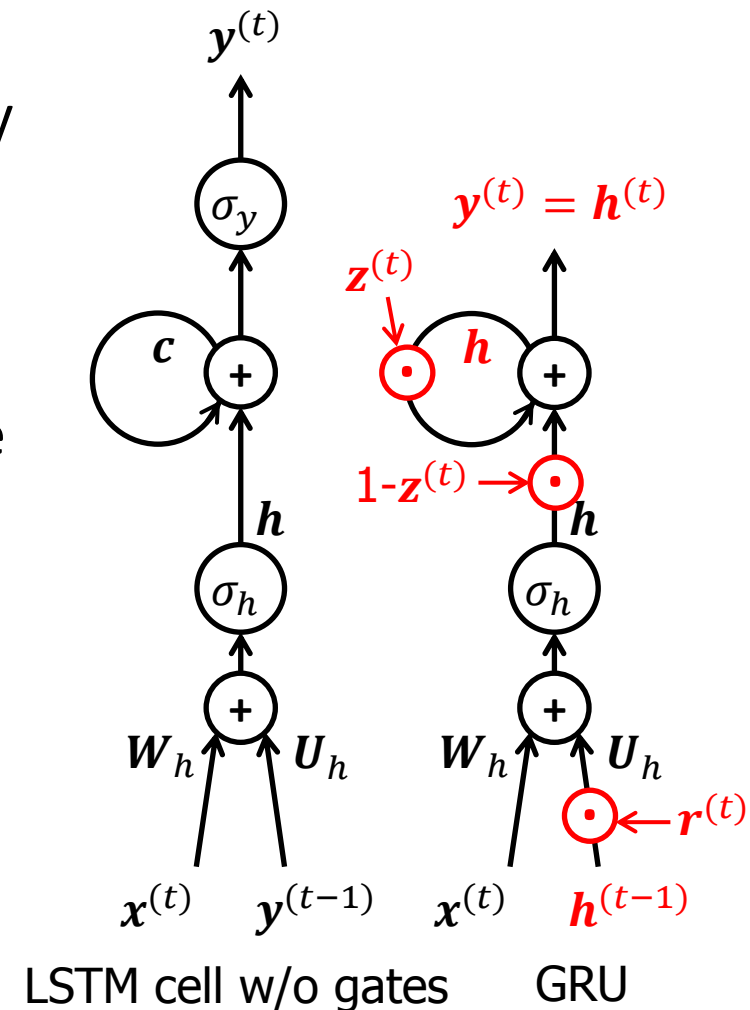
$$\mathbf{c}^{(t)} = \mathbf{c}^{(t-1)} + \mathbf{h}^{(t)}$$

- Why not **use h directly**?
- This is the **1st idea** of GRU

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

where $\mathbf{z}^{(t)} \in (0,1)$ and

$$\tilde{\mathbf{h}}^{(t)} = \sigma_h(\mathbf{W}_h \mathbf{x}^{(t)} + \mathbf{U}_h \mathbf{h}^{(t-1)} + \mathbf{b}_h)$$



Gated recurrent unit (GRU)

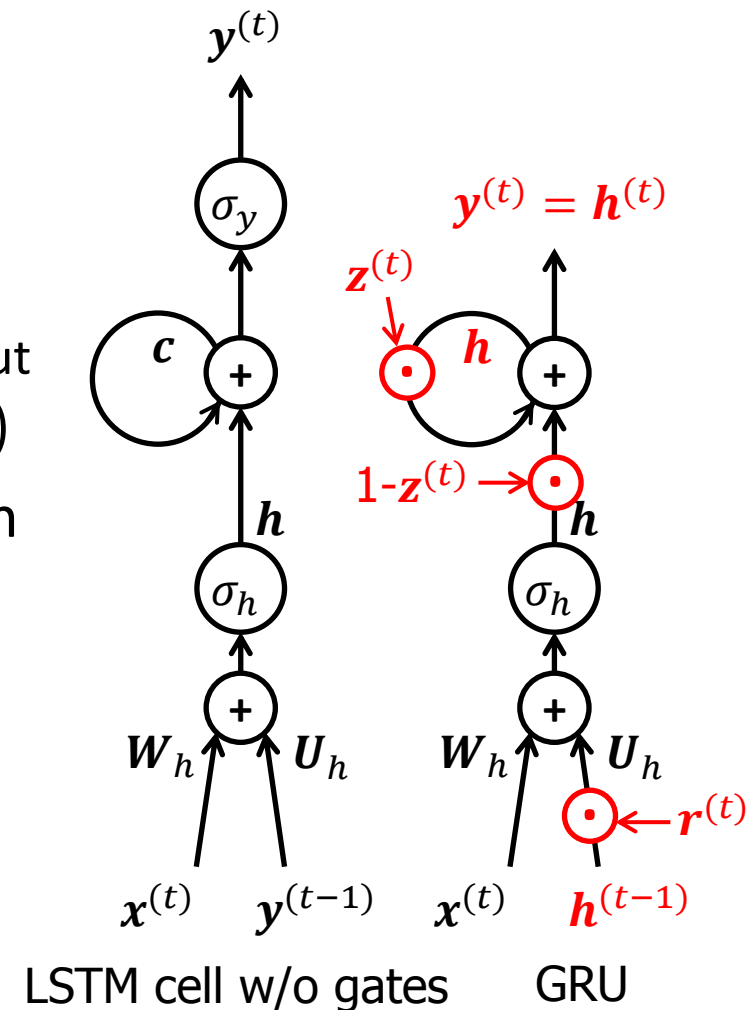
- The **2nd idea** of GRU
 - Let output be equal to hidden states:

$$\mathbf{y}^{(t)} = \mathbf{h}^{(t)}$$
- The **3rd idea** of GRU
 - Use a gate to modulate the recurrent input

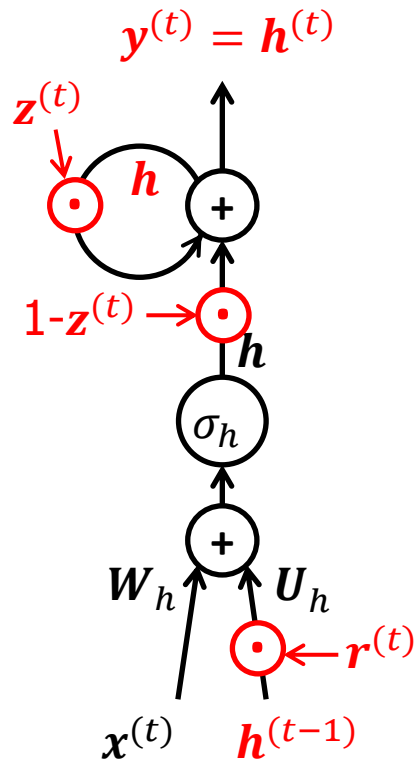
$$\tilde{\mathbf{h}}^{(t)} = \sigma_h(\mathbf{W}_h \mathbf{x}^{(t)} + \mathbf{U}_h(\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)})) + \mathbf{b}_h$$
- The gates depend on input and hidden states
 - **Update gate**

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

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



GRU in summary



- Dynamic equations

$$h^{(t)} = z^{(t)} \odot h^{(t-1)} + (1 - z^{(t)}) \odot \tilde{h}^{(t)}$$
$$\tilde{h}^{(t)} = \sigma_h(W_h x^{(t)} + U_h(r^{(t)} \odot h^{(t-1)}) + b_h)$$

where σ_h is either the logistic sigmoid function or tanh function

- The gates

$$z^{(t)} = \sigma(W_z x^{(t)} + U_z h^{(t-1)} + b_z)$$

$$r^{(t)} = \sigma(W_r x^{(t)} + U_r h^{(t-1)} + b_r)$$

σ is the logistic sigmoid function

There are many other variants of LSTM and GRU, but none of them would clearly beat both of these across a wide range of tasks (Greff et al., 2015)

Summary

- Dynamic systems
 - Autonomous and non-autonomous
- Simple RNNs
 - Jordan network and Elman network
 - BPTT and teacher forcing
 - Bidirectional RNN
 - Deep RNN
- Recurrent CNN
 - CNN + recurrent connections
- Gated RNNs
 - LSTM and GRU

Further reading

- Goodfellow, Bengio and Courville, 2016
Deep Learning, MIT Press, Chapters 10
- Understanding LSTM networks
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Prepare for the next lecture

- Form groups of 2 and every group prepares a 5-minute presentation with slides for the following paper
 - Santurkar, Tsipras, Ilyas, Madry (2018) How does batch normalization help optimization? NeurIPS

Nobody prepared for this paper today. So I'll randomly select a group to present the paper in the next week (April 25)!

April 18