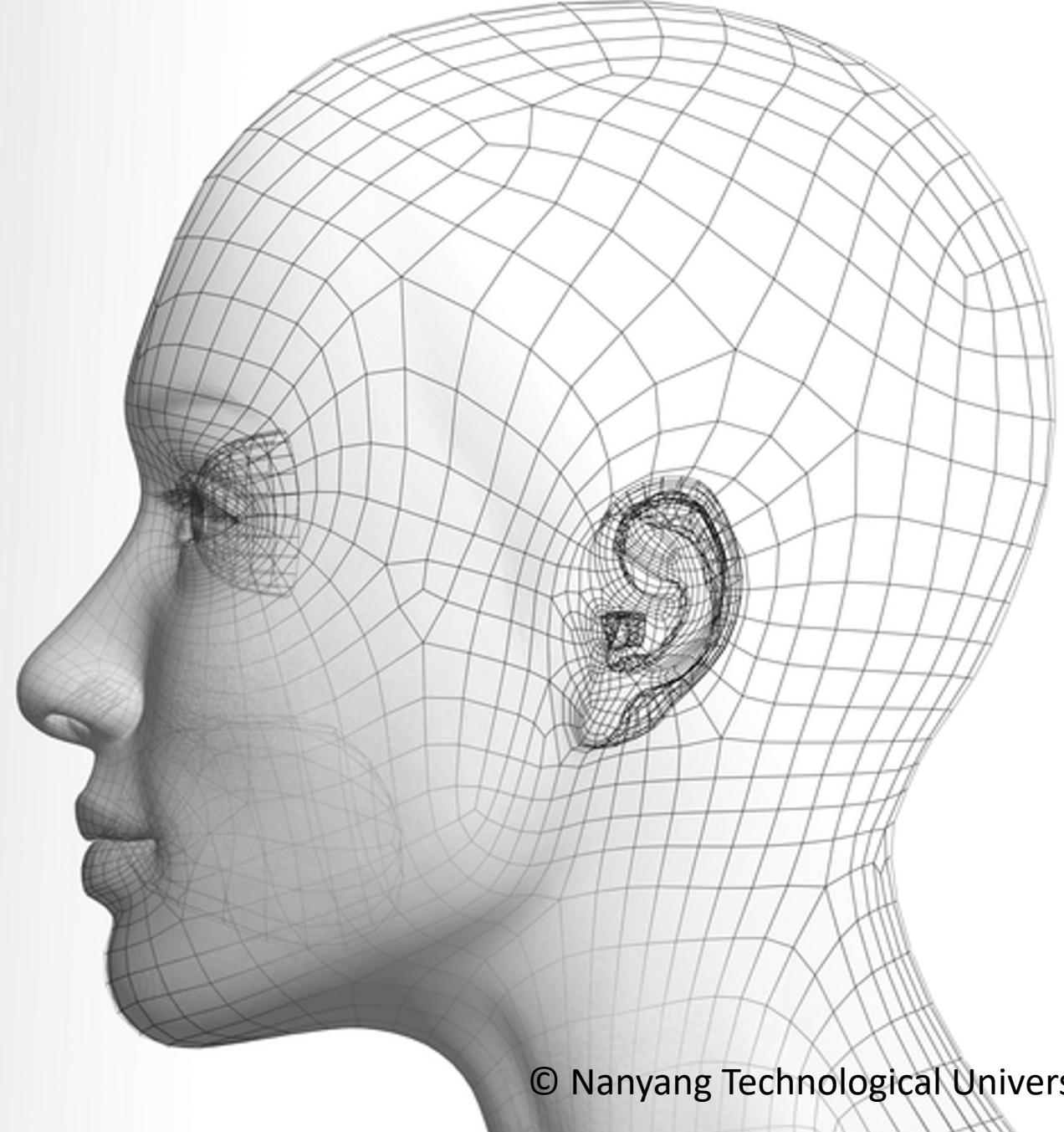
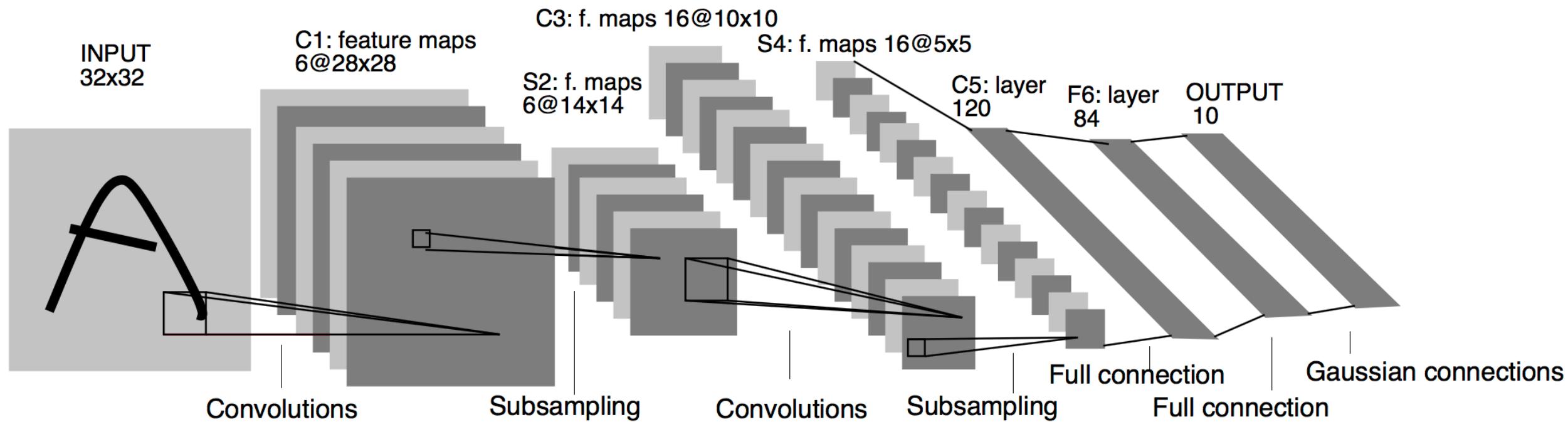


Recurrent Neural Networks (RNN)

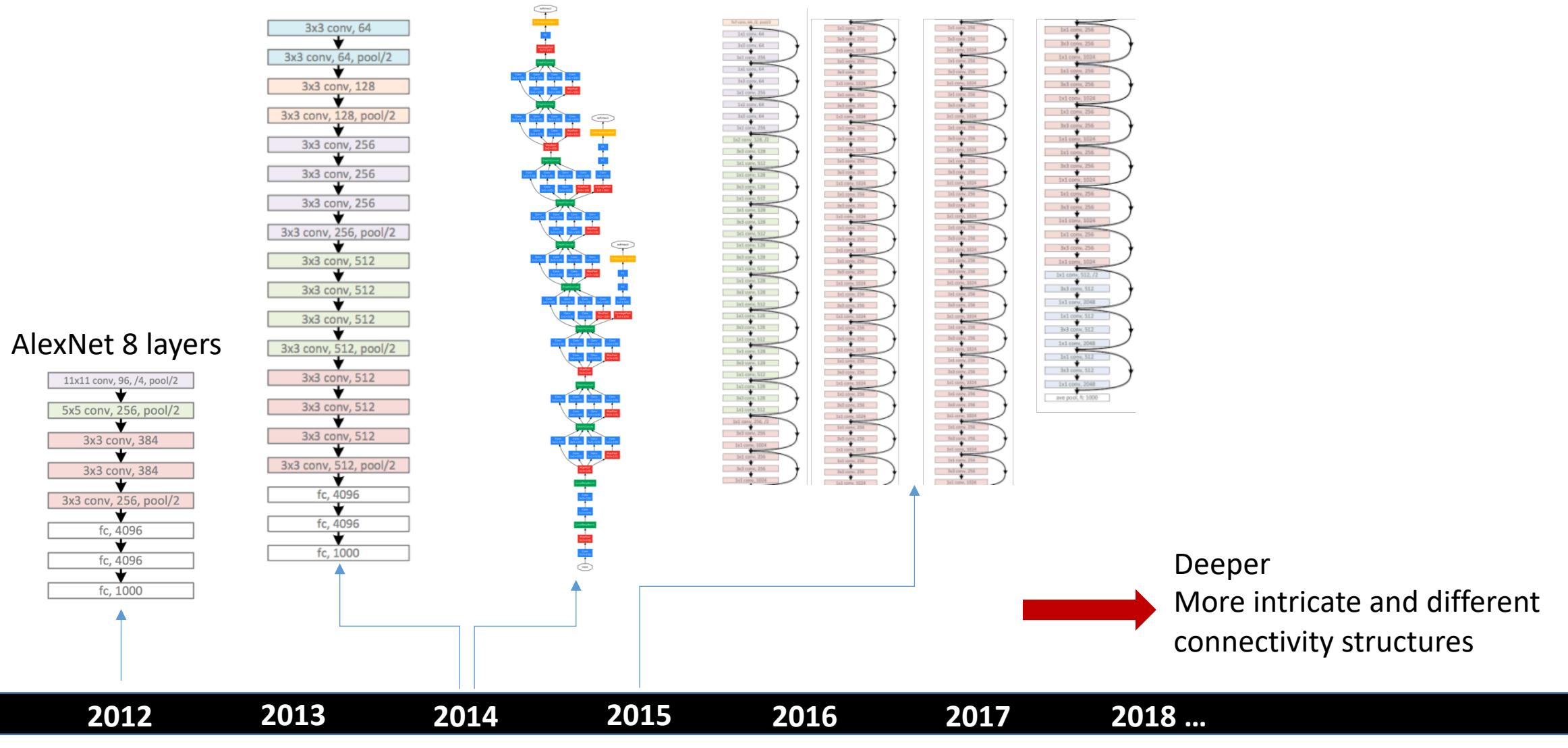
Neural networks and deep learning



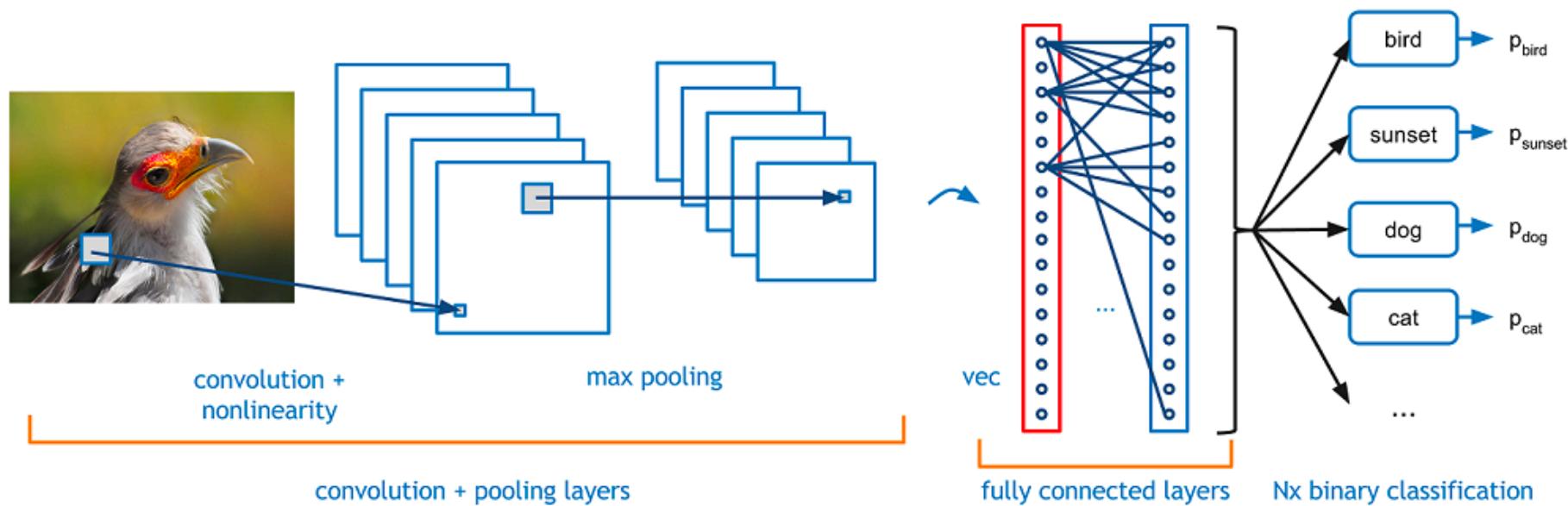
Previous: Convolutional Neural Networks (CNN)



Previous: Convolutional Neural Networks (CNN)



Previous: Convolutional Neural Networks (CNN)



How about sequential information?

- Turn a sequence of sound pressures into a sequence of word identities?
- Speech recognition?
- Video prediction?

Recurrent Neural Networks (RNN)

Recurrent neural networks (RNN) are designed to process **sequential information**. That is, the data presented in a sequence.

The next data point in the sequence is usually **dependent** on the current data point.

RNN attempts to **capture dependency** among the data points in the sequence.

Examples:

- Natural language processing (spoken words and written sentences). The next word in a sentence depends on the word which comes before it.
- Genomic sequences: a nucleotide in a DNA sequence is dependent on its neighbors.

Text Generation with an RNN

QUEENE:

*I had thought thou hadst a Roman; for the oracle,
Thus by All bids the man against the word,
Which are so weak of care, by old care done;
Your children were in your holy love,
And the precipitation through the bleeding throne.*

BISHOP OF ELY:

*Marry, and will, my lord, to weep in such a one were prettiest;
Yet now I was adopted heir
Of the world's lamentable day,
To watch the next way with his father with his face?*

Text Generation with an RNN

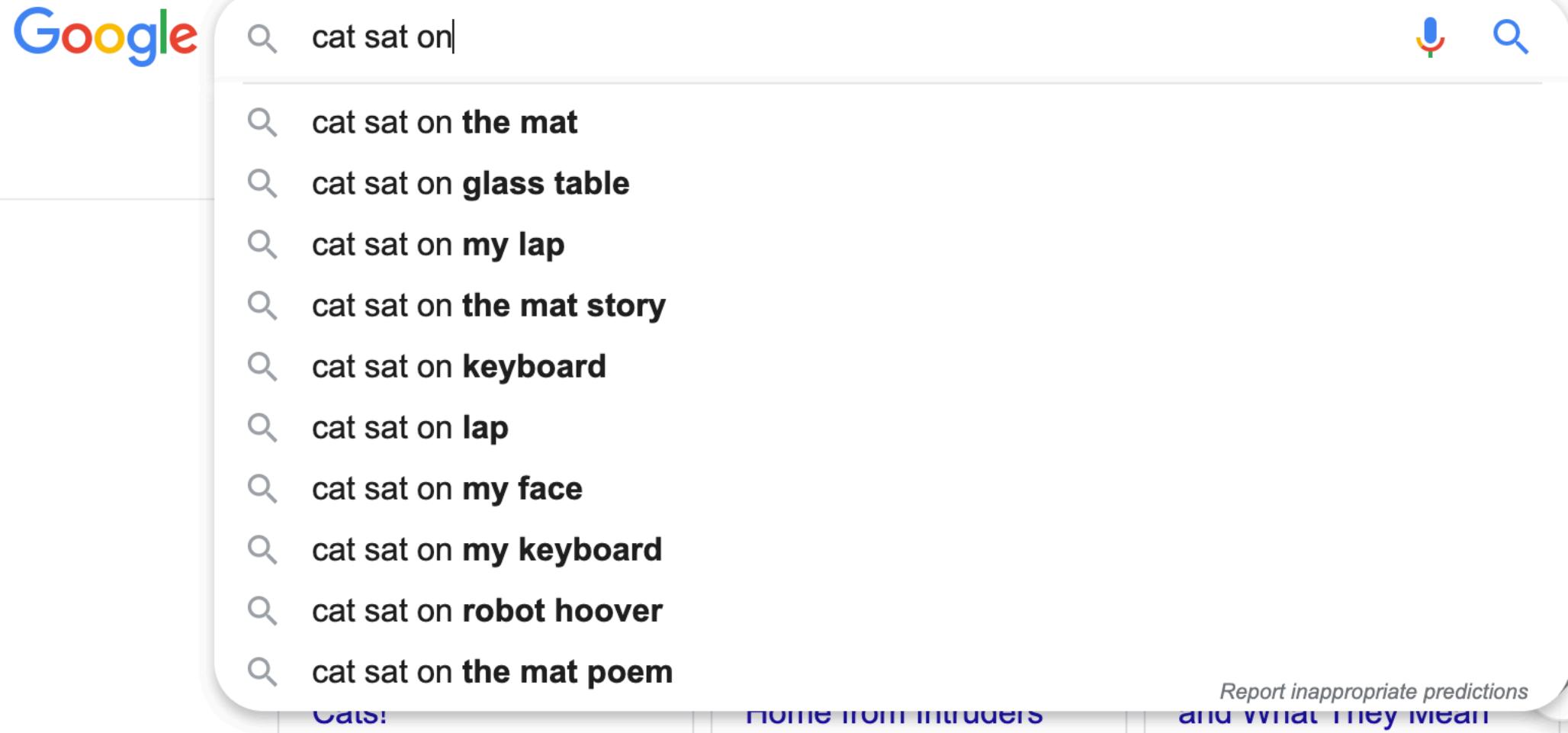


Image Captioning



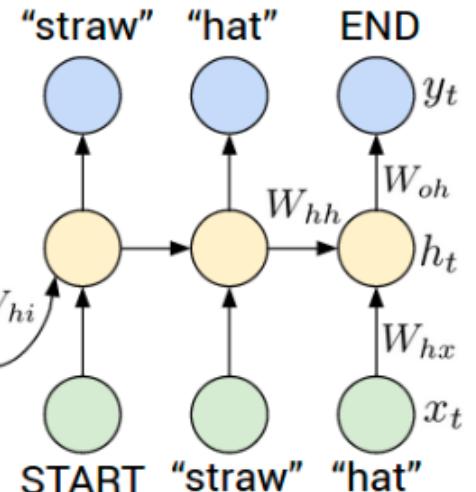
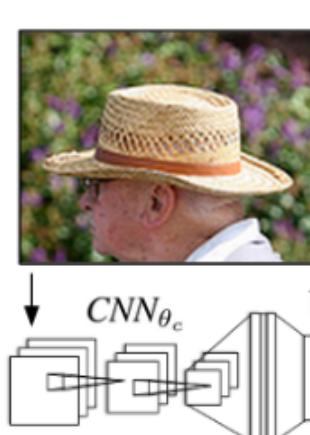
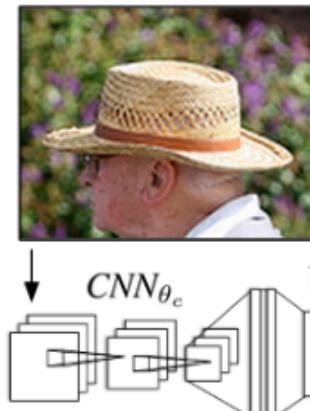
"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



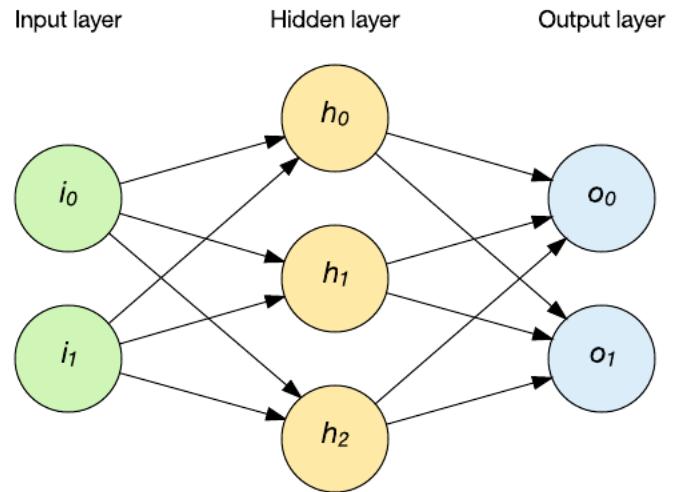
Recurrent Neural Networks (RNN)

RNNs are called recurrent because they perform **the same task** for every data element (frame) of a sequence, with the output depending on the previous computations.

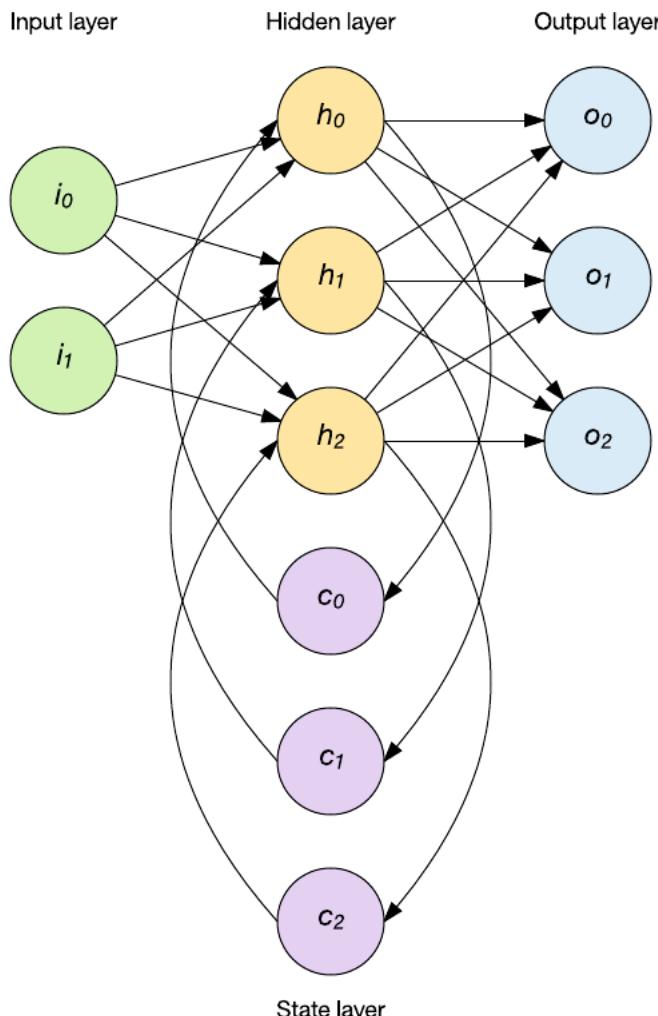
In other words, it has a **memory** that captures information about what has been processed so far.

RNN achieves its function by using **feedback connections** that enables learning of sequential (temporal) information of sequences.

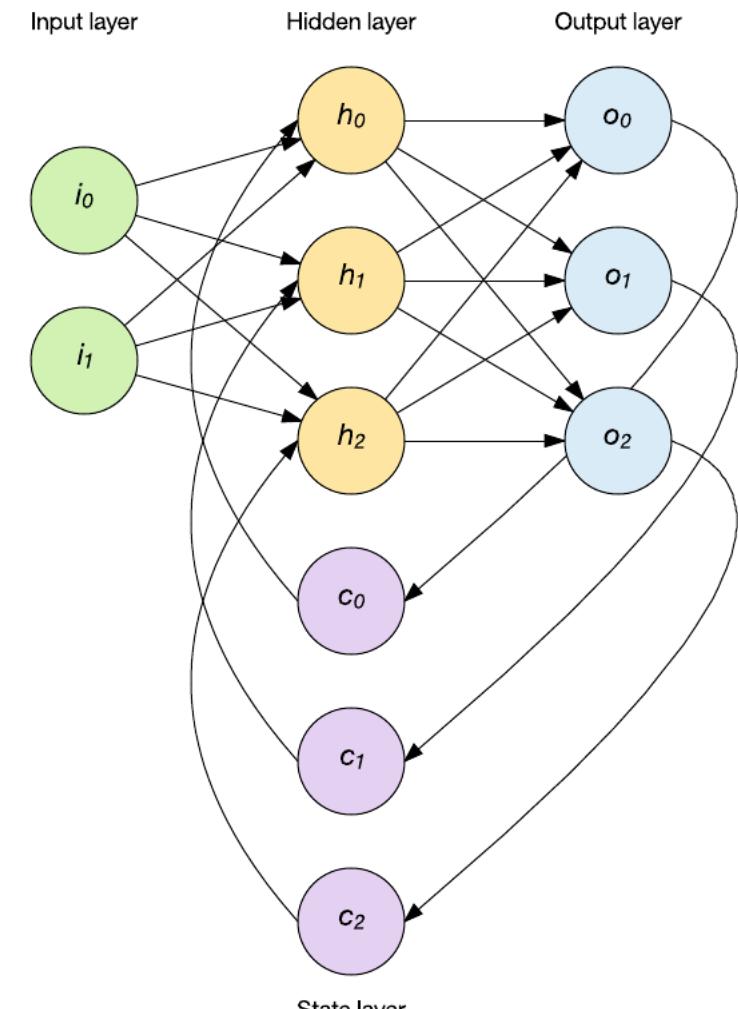
Types of RNN



Feedforward NN



RNN with hidden recurrence
(Elman-type)



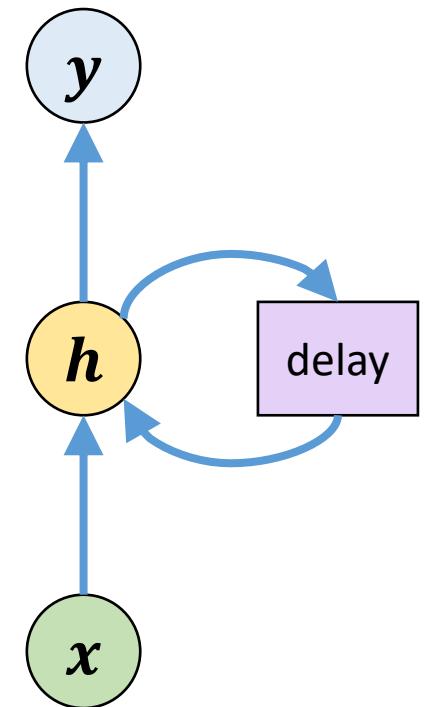
RNN with top-down recurrence
(Jordan-type)

RNN with hidden recurrence (Elman type)

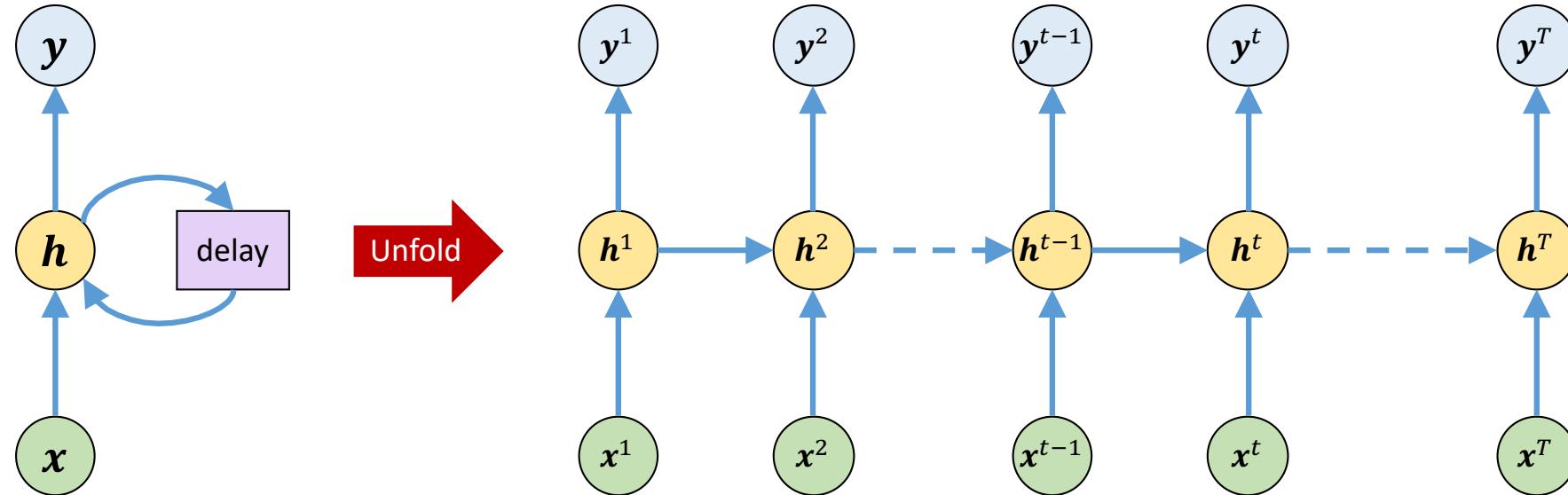
Elman type networks (sometimes called a “Vanilla RNN”) that produce an output at each time step and have recurrent connections between hidden units.

Data is presented as a sequence of instance t (time) that is discretized and activations are updated at each time instant t . The hidden-layer activation at time $t - 1$ is kept by the delay unit and fed to the hidden layer at time t together with the raw input $x(t)$.

The **delay unit** represents that the activation is held for one time unit until the next time instance. Here, one unit represents the time between two adjacent data points in the sequence .



RNN with hidden recurrence



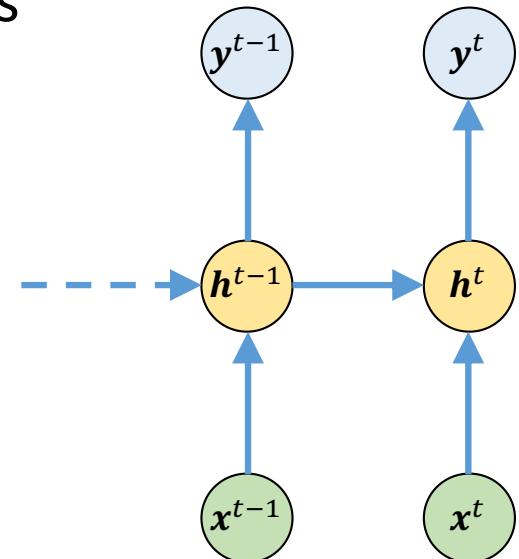
- Note that the purple square in the computational graph indicates that a feedback takes place with a delay of a unit time, from state at time t to the state at time $t + 1$.
- The recurrent connections in the hidden-layer can be unfolded and represented in time for processing a time series $(x(t))_{t=1}^T$

RNN with hidden recurrence

By considering the unfolded structure, $\mathbf{h}(t)$ is dependent on all the inputs before time t :

$$h(t) = f^t(x(t), x(t-1), \dots, x(2), x(1))$$

By using recurrence relationship, we can write:



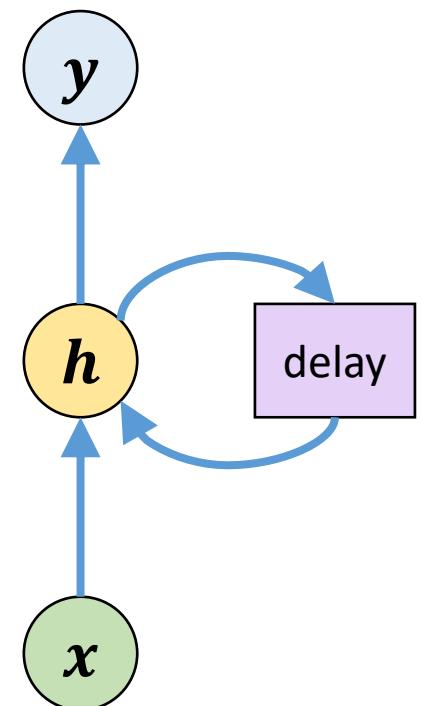
The function f^t takes the whole past sequence $(x(t), x(t-1), \dots, x(2), x(1))$ as input and produce the hidden layer activation.

The past inputs are represented by the hidden-layer activations in the previous instant. The recurrent structure allows us to factorize f^t into repeated applications of a function f .

RNN with hidden recurrence

The folded structure introduces two major advantages:

1. Regardless of the sequence length, **the learned model always has the same size**, rather than specified in terms of a variable-length history of states.
2. It is possible to use **same transition function f** with the **same parameters** at every time step.



RNN with hidden recurrence

U : weight vector that transforms raw inputs to the hidden-layer

W : recurrent weight vector connecting previous hidden-layer output to hidden input

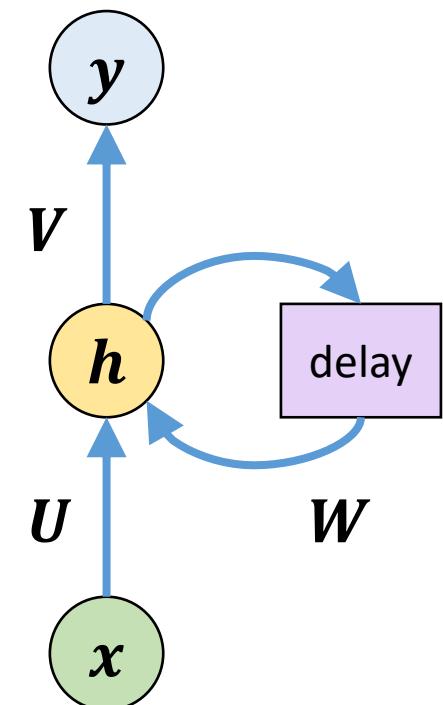
V : weight vector of the output layer

b : bias connected to hidden layer

c : bias connected to the output layer

ϕ : the tanh hidden-layer activation function

σ : the linear/softmax output-layer activation function



RNN with hidden recurrence

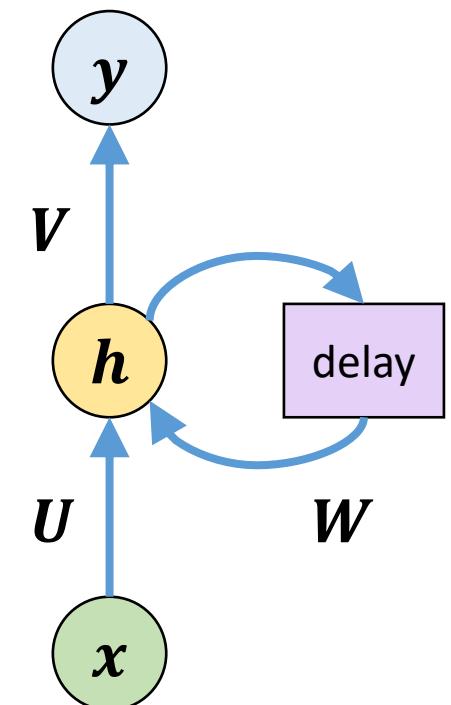
Let $x(t)$, $y(t)$, and $h(t)$ be the input, output, and hidden output of the network at time t .

Activation of the Elman-type RNN with one hidden-layer is given by:

$$h(t) = \phi(U^T x(t) + W^T h(t-1) + b)$$

$$y(t) = \sigma(V^T h(t) + c)$$

σ is a *softmax* function for classification and a *linear* function for regression.



RNN with hidden recurrence: batch processing

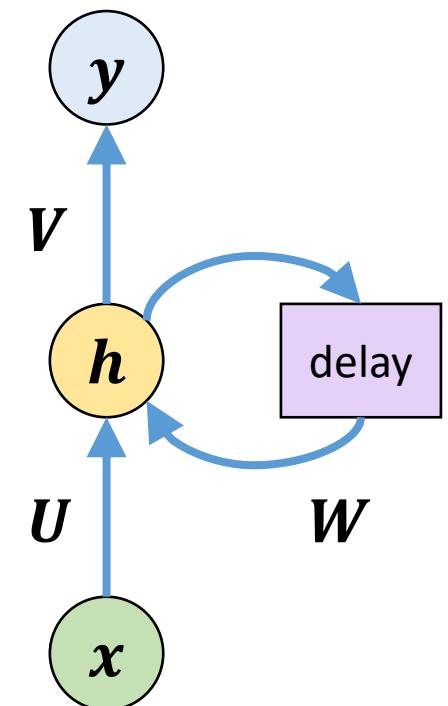
Given P patterns $\{x_p\}_{p=1}^P$ where $x_p = (x_p(t))_{t=1}^T$,

$$X(t) = \begin{pmatrix} x_1(t)^T \\ x_2(t)^T \\ \vdots \\ x_P(t)^T \end{pmatrix}$$

Let $X(t)$, $Y(t)$, and $H(t)$ be batch input, output, and hidden output of the network at time t

Activation of the three-layer Elman-type RNN is given by:

$$\begin{aligned} H(t) &= \phi(X(t)U + H(t-1)W + B) \\ Y(t) &= \sigma(H(t)V + C) \end{aligned}$$



Example 1

A recurrent neural network with hidden recurrence has two input neurons, three hidden neurons, and two output neurons. The parameters of the

network are initialized as $\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix}$, $\mathbf{W} = \begin{pmatrix} 2.0 & 1.3 & -1.0 \\ 1.5 & 0.0 & -0.5 \\ -0.2 & 1.5 & 0.4 \end{pmatrix}$ and $\mathbf{V} = \begin{pmatrix} 2.0 & -1.0 \\ -1.5 & 0.5 \\ 0.2 & 0.8 \end{pmatrix}$.

Bias to the hidden layer $\mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \\ 0.2 \end{pmatrix}$ and to the output layer $\mathbf{c} = \begin{pmatrix} 0.5 \\ 0.1 \end{pmatrix}$.

For a sequence of inputs $(\mathbf{x}(1), \mathbf{x}(2), \mathbf{x}(3), \mathbf{x}(4))$ where

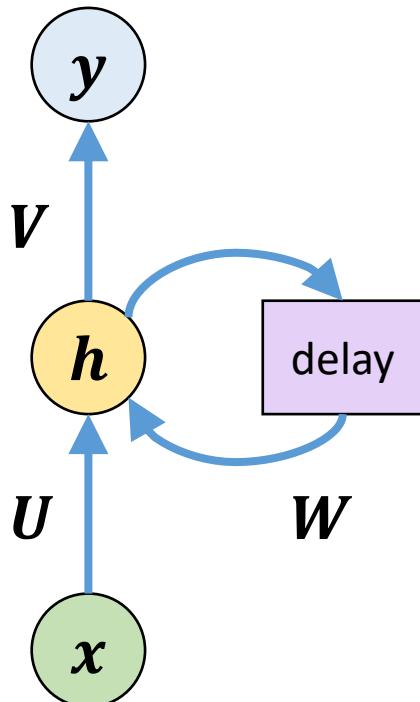
$$\mathbf{x}(1) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \mathbf{x}(2) = \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \mathbf{x}(3) = \begin{pmatrix} 0 \\ 3 \end{pmatrix}, \text{ and } \mathbf{x}(4) = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$$

find the output of RNN.

Assume that hidden layer activations are initialized to zero and \tanh and sigmoid functions for the hidden and output layer activation functions, respectively.

Example 1 (con't)

$$\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix}, \quad \mathbf{W} = \begin{pmatrix} 2.0 & 1.3 & -1.0 \\ 1.5 & 0.0 & -0.5 \\ -0.2 & 1.5 & 0.4 \end{pmatrix} \text{ and } \mathbf{V} = \begin{pmatrix} 2.0 & -1.0 \\ -1.5 & 0.5 \\ 0.2 & 0.8 \end{pmatrix}$$



$$\mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \\ 0.2 \end{pmatrix} \text{ and } \mathbf{c} = \begin{pmatrix} 0.5 \\ 0.1 \end{pmatrix}$$

Three hidden neurons and two output neurons

$$\begin{aligned}\mathbf{h}(t) &= \phi(\mathbf{U}^T \mathbf{x}(t) + \mathbf{W}^T \mathbf{h}(t-1) + \mathbf{b}) \\ \mathbf{y}(t) &= \sigma(\mathbf{V}^T \mathbf{h}(t) + \mathbf{c})\end{aligned}$$

$$\begin{aligned}\phi(u) &= \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}} \\ \sigma(u) &= \text{sigmoid}(u) = \frac{1}{1 + e^{-u}}.\end{aligned}$$

Assume $\mathbf{h}(0) = (0 \quad 0 \quad 0)^T$.

Example 1 (con't)

At $t=1$, $\mathbf{x}(1) = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$:

$$\begin{aligned}\mathbf{h}(1) &= \phi(\mathbf{U}^T \mathbf{x}(1) + \mathbf{W}^T \mathbf{h}(0) + \mathbf{b}) \\ &= \tanh\left(\begin{pmatrix} -1.0 & 0.5 \\ 0.5 & 0.1 \\ 0.2 & -2.0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} + \begin{pmatrix} 2.0 & 1.5 & -0.2 \\ 1.3 & 0.0 & 1.5 \\ -1.0 & -0.5 & 0.4 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \\ 0.2 \end{pmatrix}\right) \\ &= \begin{pmatrix} 0.2 \\ 0.72 \\ -1.0 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{y}(1) &= \sigma(\mathbf{V}^T \mathbf{h}(1) + \mathbf{c}) \\ &= \text{sigmoid}\left(\begin{pmatrix} 2.0 & -1.5 & 0.2 \\ -1.0 & 0.5 & 0.8 \end{pmatrix} \begin{pmatrix} 0.2 \\ 0.72 \\ -1.0 \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.1 \end{pmatrix}\right) \\ &= \begin{pmatrix} 0.41 \\ 0.37 \end{pmatrix}\end{aligned}$$

Example 1 (con't)

At $t=2$, $\mathbf{x}(2) = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$:

$$\begin{aligned}\mathbf{h}(2) &= \phi(\mathbf{U}^T \mathbf{x}(2) + \mathbf{W}^T \mathbf{h}(1) + \mathbf{b}) \\ &= \tanh\left(\begin{pmatrix} -1.0 & 0.5 \\ 0.5 & 0.1 \\ 0.2 & -2.0 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} + \begin{pmatrix} 2.0 & 1.5 & -0.2 \\ 1.3 & 0.0 & 1.5 \\ -1.0 & -0.5 & 0.4 \end{pmatrix} \begin{pmatrix} 0.2 \\ 0.72 \\ -1.0 \end{pmatrix} + \begin{pmatrix} 0.2 \\ 0.2 \\ 0.2 \end{pmatrix}\right) \\ &= \begin{pmatrix} 1.0 \\ -0.89 \\ -0.99 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{y}(2) &= \sigma(\mathbf{V}^T \mathbf{h}(2) + \mathbf{c}) \\ &= \text{sigmoid}\left(\begin{pmatrix} 2.0 & -1.5 & 0.2 \\ -1.0 & 0.5 & 0.8 \end{pmatrix} \begin{pmatrix} 1.0 \\ -0.89 \\ -0.99 \end{pmatrix} + \begin{pmatrix} 0.5 \\ 0.1 \end{pmatrix}\right) = \begin{pmatrix} 0.97 \\ 0.11 \end{pmatrix}\end{aligned}$$

Example 1 (con't)

Similarly,

$$\text{at } t=3, \mathbf{x}(3) = \begin{pmatrix} 0 \\ 3 \end{pmatrix}; \mathbf{h}(3) \begin{pmatrix} 0.99 \\ 0.3 \\ -1.0 \end{pmatrix} \text{ and } \mathbf{y}(3) = \begin{pmatrix} 0.86 \\ 0.18 \end{pmatrix}$$

$$\text{at } t=4, \mathbf{x}(4) = \begin{pmatrix} 2 \\ -1 \end{pmatrix}; \mathbf{h}(4) \begin{pmatrix} 0.31 \\ 0.71 \\ 0.79 \end{pmatrix} \text{ and } \mathbf{y}(4) = \begin{pmatrix} 0.55 \\ 0.68 \end{pmatrix}$$

The output is $(\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \mathbf{y}(4))$ where

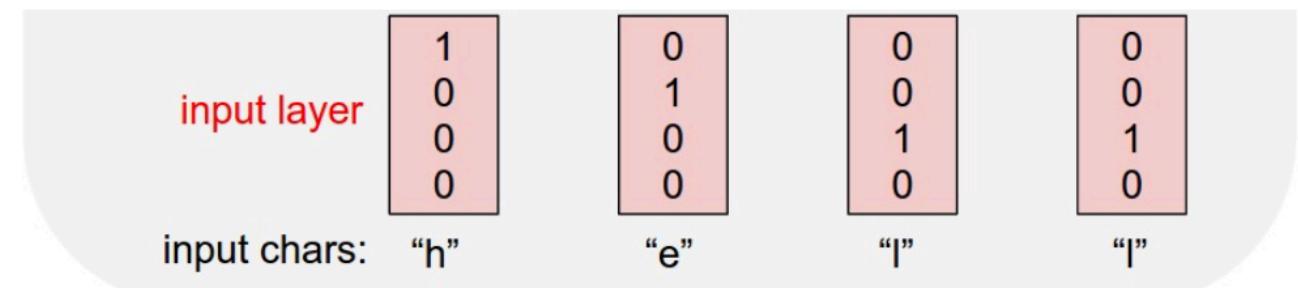
$$\mathbf{y}(1) = \begin{pmatrix} 0.41 \\ 0.37 \end{pmatrix}, \mathbf{y}(2) = \begin{pmatrix} 0.97 \\ 0.11 \end{pmatrix}, \mathbf{y}(3) = \begin{pmatrix} 0.86 \\ 0.18 \end{pmatrix}, \mathbf{y}(4) = \begin{pmatrix} 0.55 \\ 0.68 \end{pmatrix}$$

An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

Training sequence: "hello"

Vocabulary: [h, e, l, o]



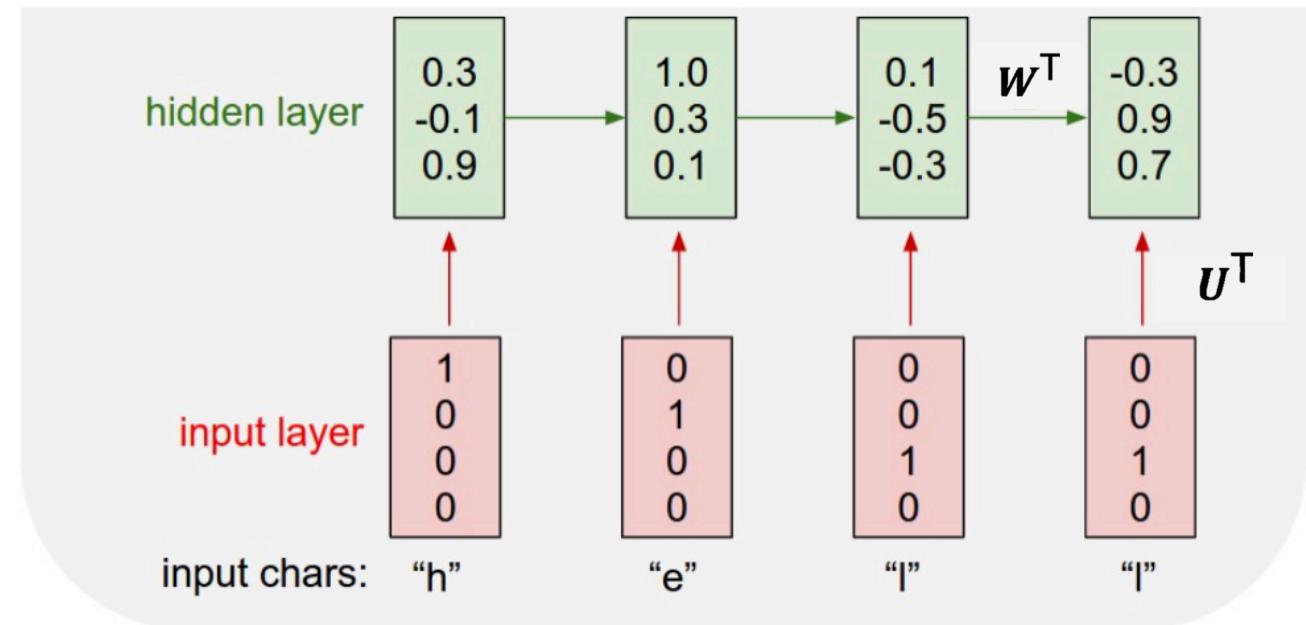
An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



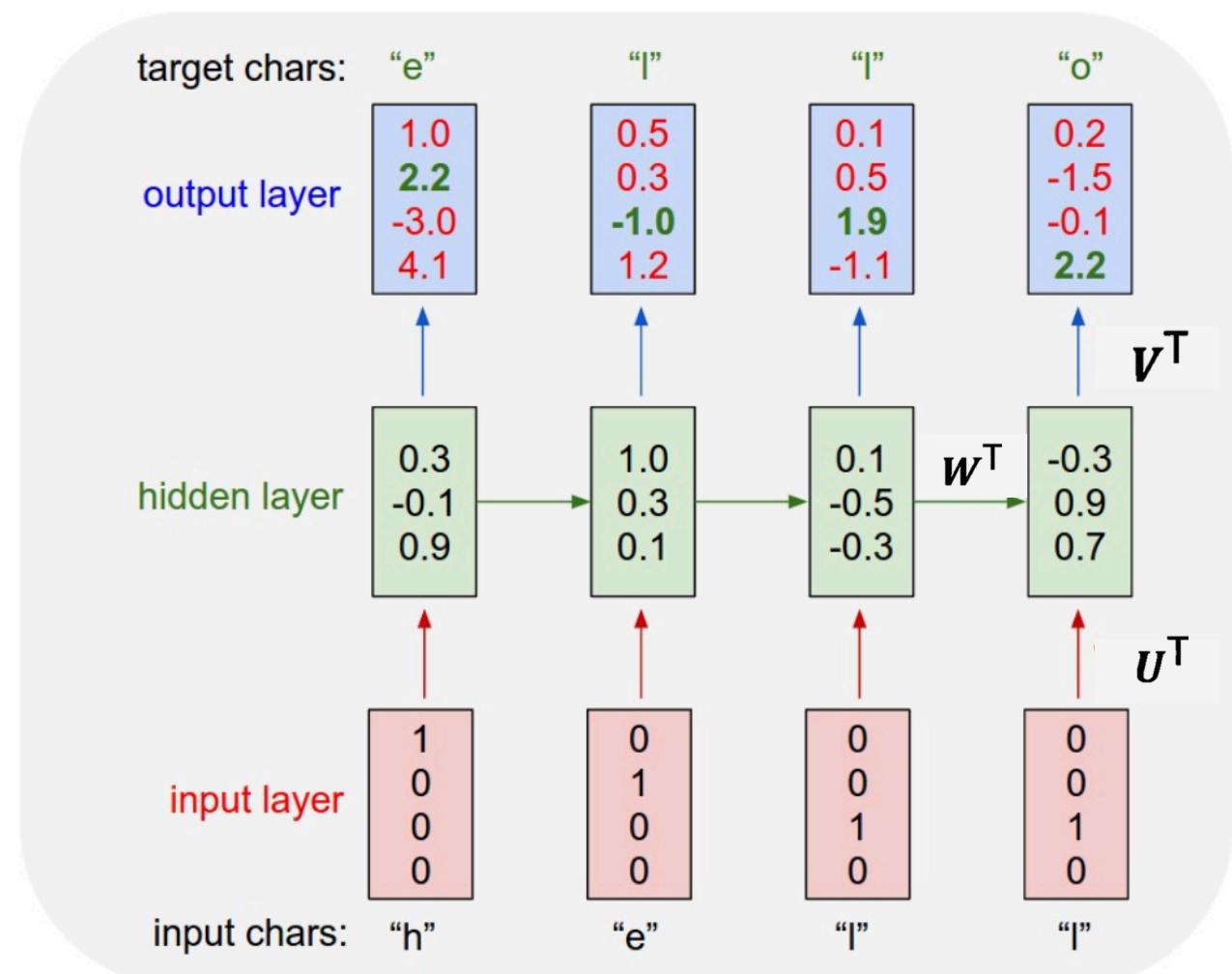
An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



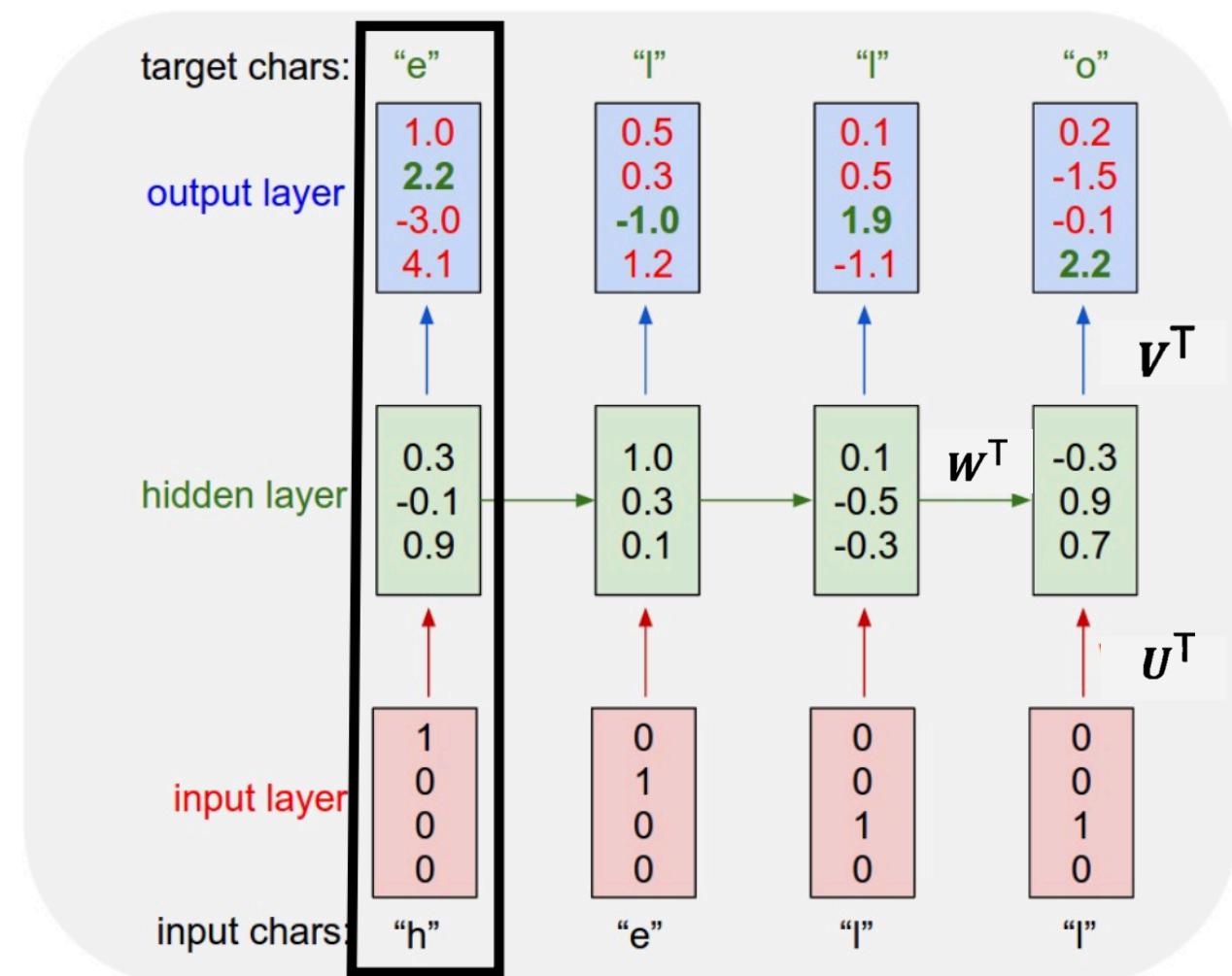
An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



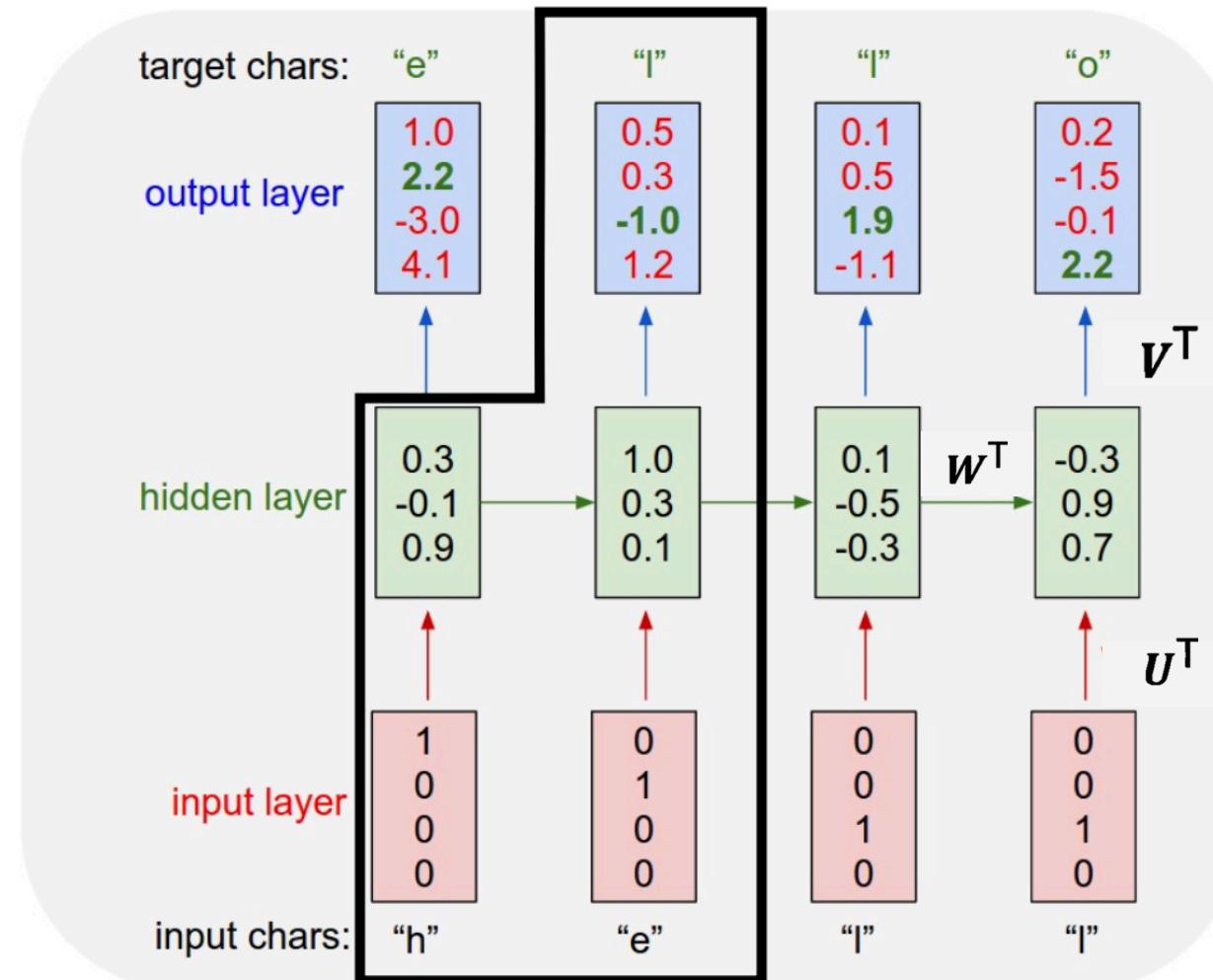
An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



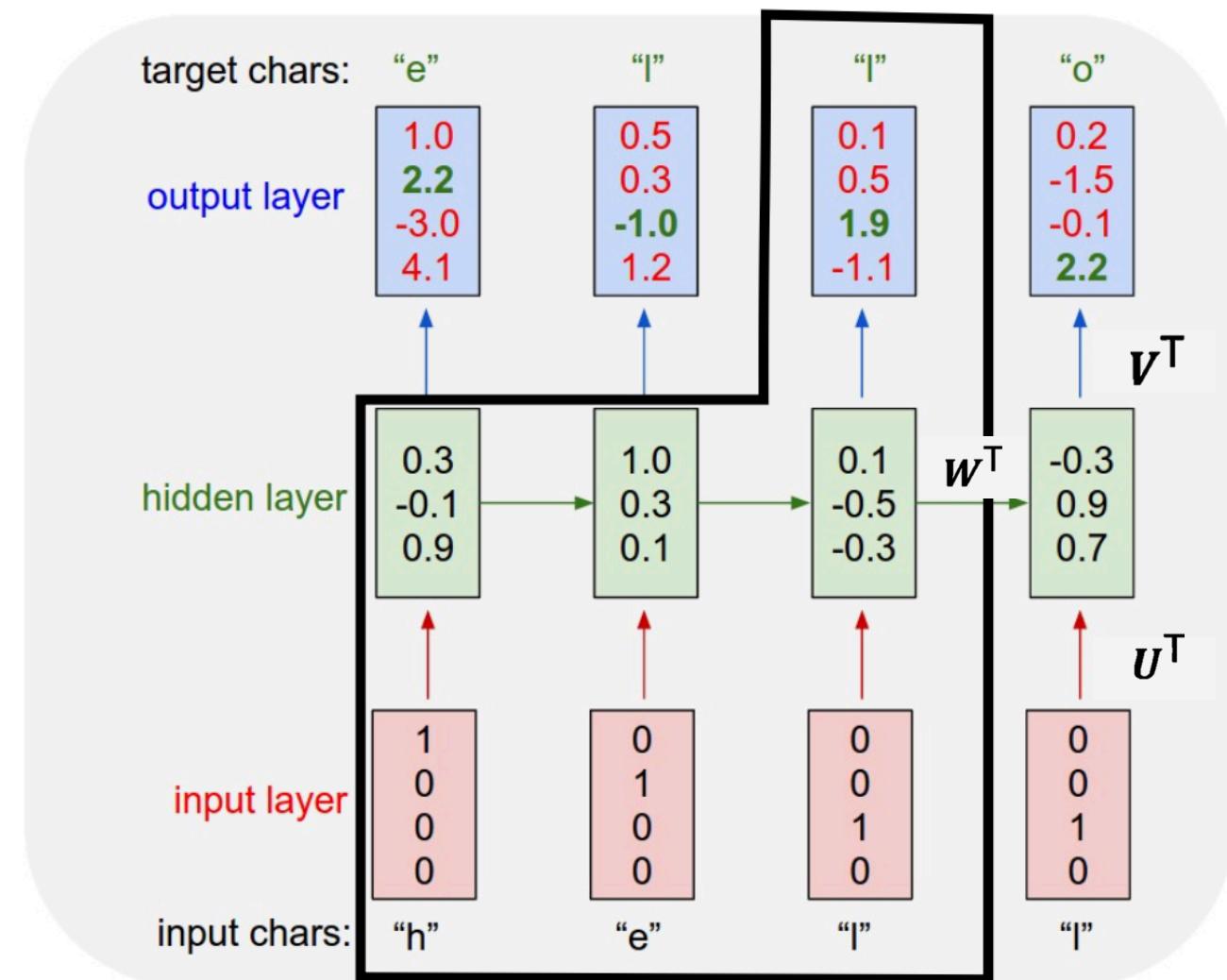
An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



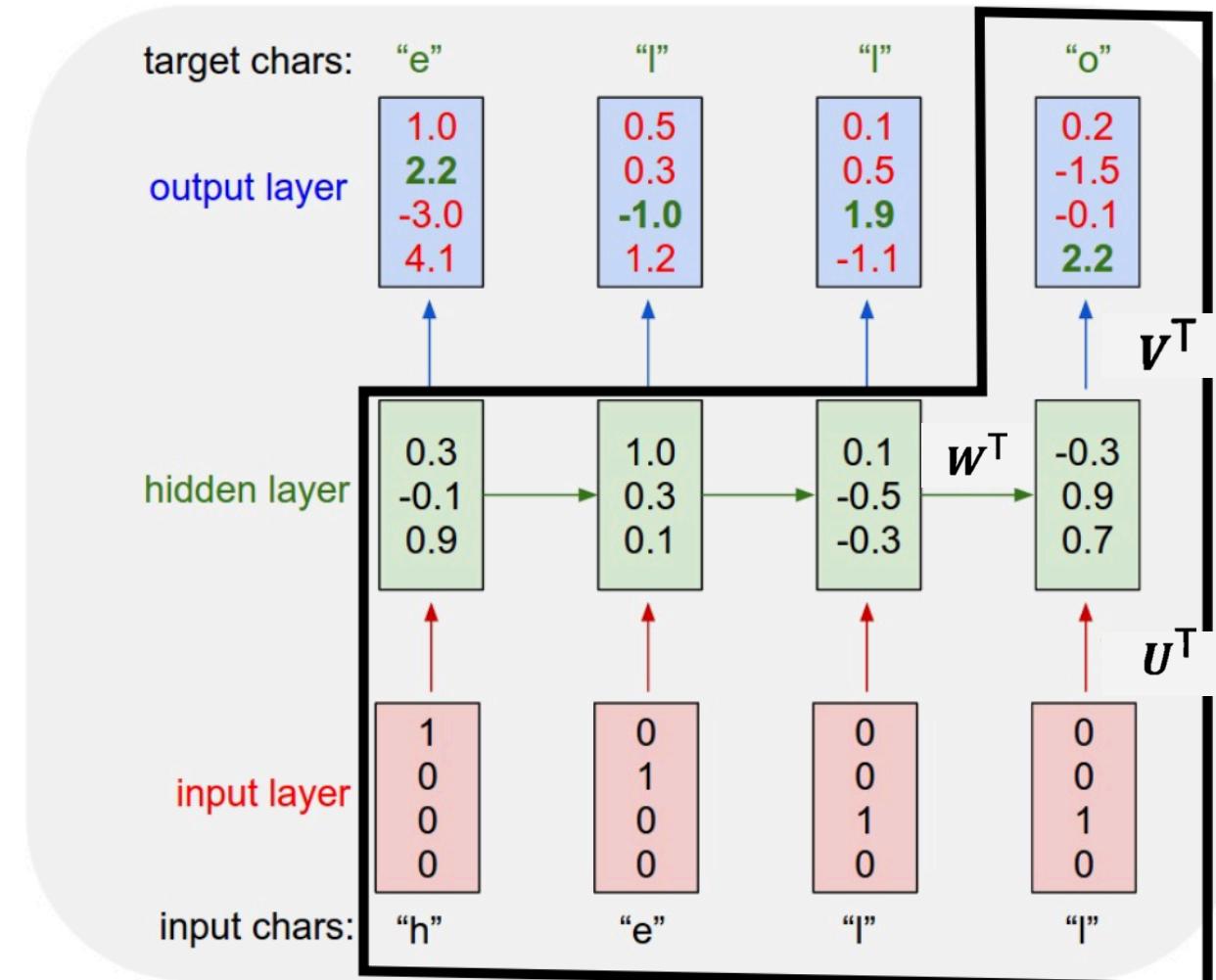
An example – Language modelling

Given characters 1, 2, ..., t,
model predicts character t

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

Training sequence: "hello"

Vocabulary: [h, e, l, o]



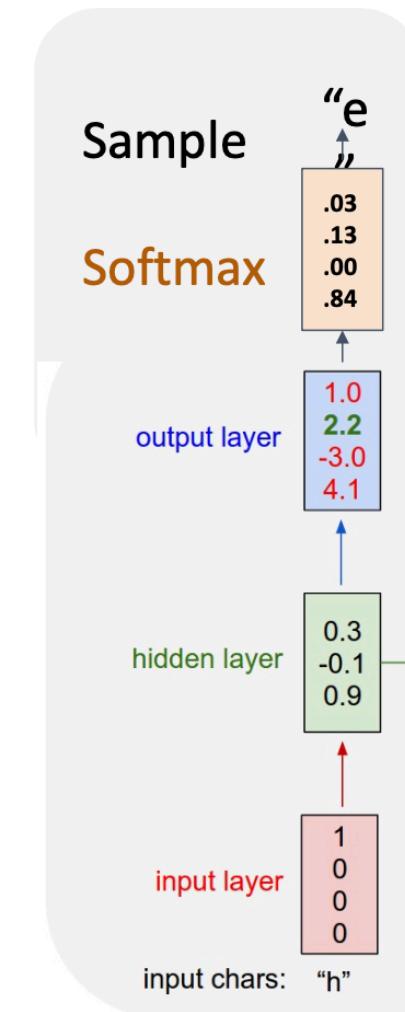
An example – Language modelling

Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]



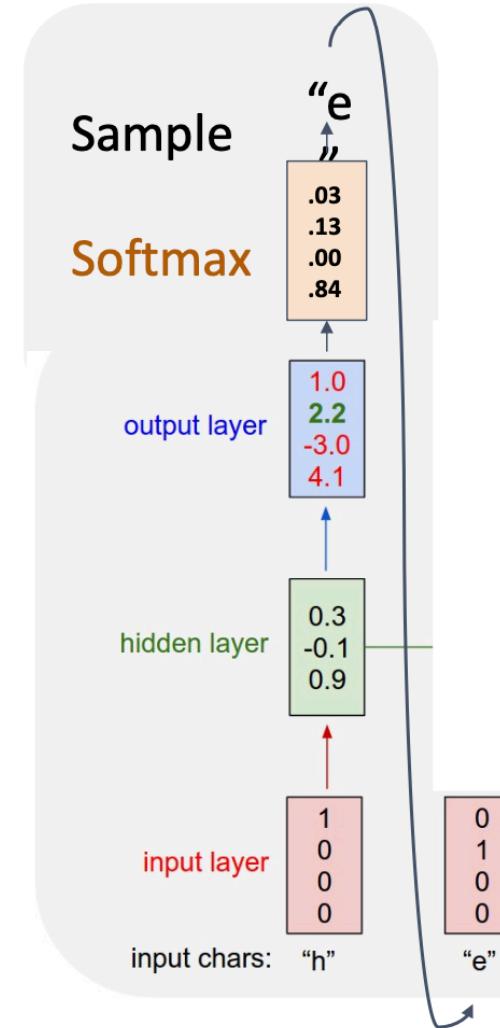
An example – Language modelling

Example: Language Modeling

At test-time, **generate new text**: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]



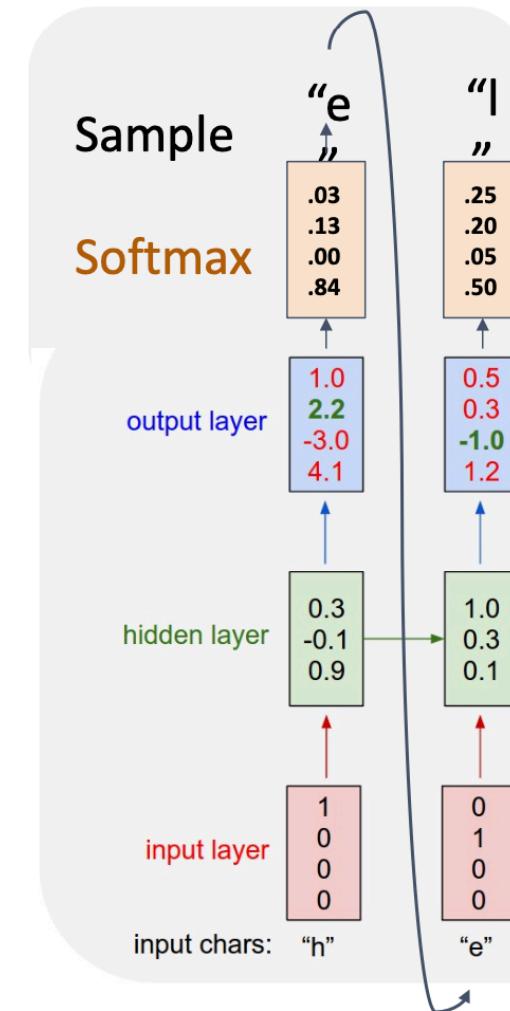
An example – Language modelling

Example: Language Modeling

At test-time, **generate** new text: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]



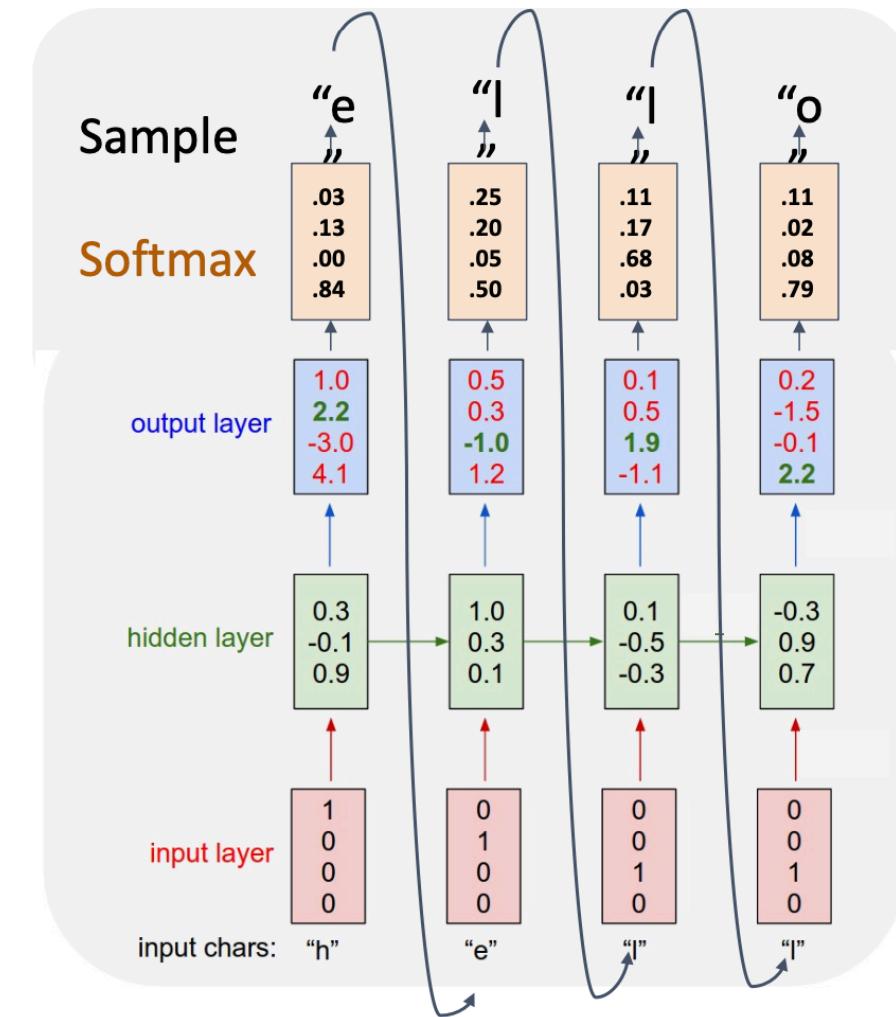
An example – Language modelling

Example: Language Modeling

At test-time, **generate new text**: sample characters one at a time, feed back to model

Training sequence: "hello"

Vocabulary: [h, e, l, o]



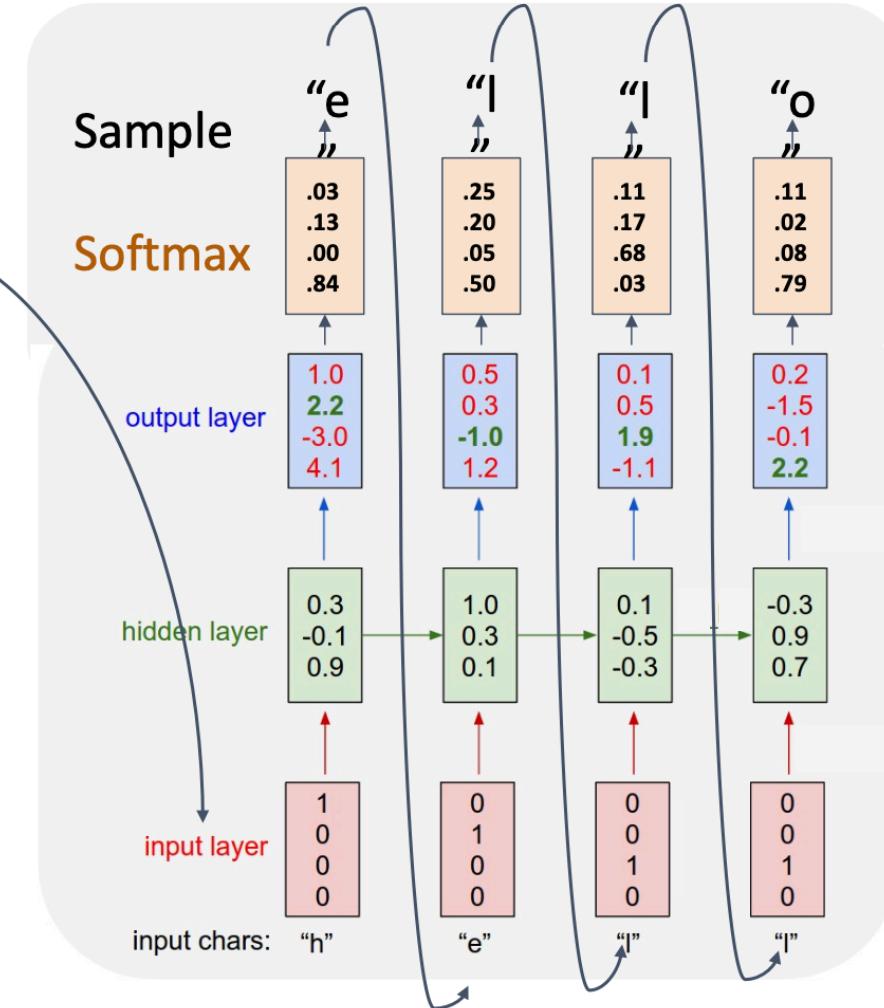
An example – Language modelling

Example: Language Modeling

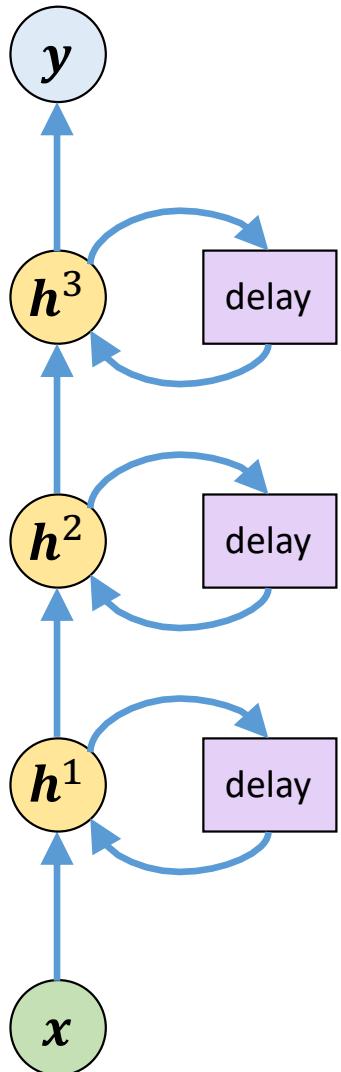
So far: encode inputs
as **one-hot-vector**

$$\begin{aligned} [w_{11} \ w_{12} \ w_{13} \ w_{14}] [1] &= [w_{11}] \\ [w_{21} \ w_{22} \ w_{23} \ w_{14}] [0] &= [w_{21}] \\ [w_{31} \ w_{32} \ w_{33} \ w_{14}] [0] &= [w_{31}] \\ &\quad [0] \end{aligned}$$

Matrix multiply with a one-hot vector just extracts a column from the weight matrix.
Often extract this into a separate
embedding layer



Deep RNN with hidden recurrence



\mathbf{U}^l : weight matrix connecting output of hidden-layer $l - 1$ to hidden-layer l .

\mathbf{W}^l : weight matrix connecting hidden-layer output of the previous instant to the current output of the hidden-layer.

\mathbf{b}^l : the bias vector connected to the layer l

First hidden-layer $l = 1$:

$$\mathbf{h}^1(t) = \phi \left(\mathbf{U}^{1\top} \mathbf{x}(t) + \mathbf{W}^{1\top} \mathbf{h}^1(t-1) + \mathbf{b}^1 \right)$$

For layer $l = 2, \dots, L-1$:

$$\mathbf{h}^l(t) = \phi \left(\mathbf{U}^{l\top} \mathbf{h}^{l-1}(t) + \mathbf{W}^{l\top} \mathbf{h}^l(t-1) + \mathbf{b}^l \right)$$

First output layer $l = L$:

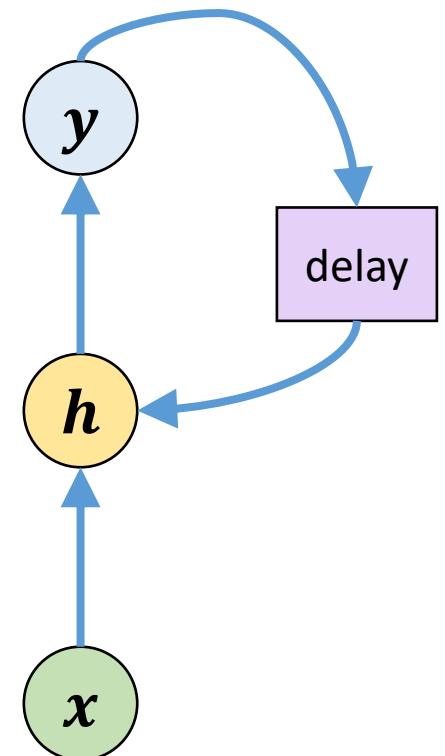
$$\mathbf{y}(t) = \sigma \left(\mathbf{U}^{L\top} \mathbf{h}^{L-1}(t) + \mathbf{b}^L \right)$$

RNN with top-down recurrence (Jordan type)

In Jordan-type RNN (that is, RNN with top-down recurrence), the output of the output-layer at time $t - 1$ is kept and fed to the hidden layer at time t together with the raw input $x(t)$ for time t .

Note that time t is assumed to be discretized and activations are updated at each time instant t .

The **delay unit** is to indicate that the output is held until the next time instance and fed to the hidden layer.



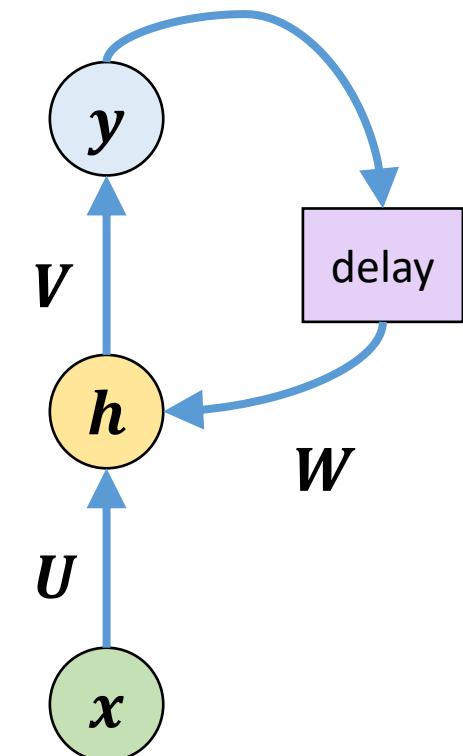
RNN with top-down recurrence

Let $x(t)$, $y(t)$, and $h(t)$ be the input, output, and hidden output of the network at time t .

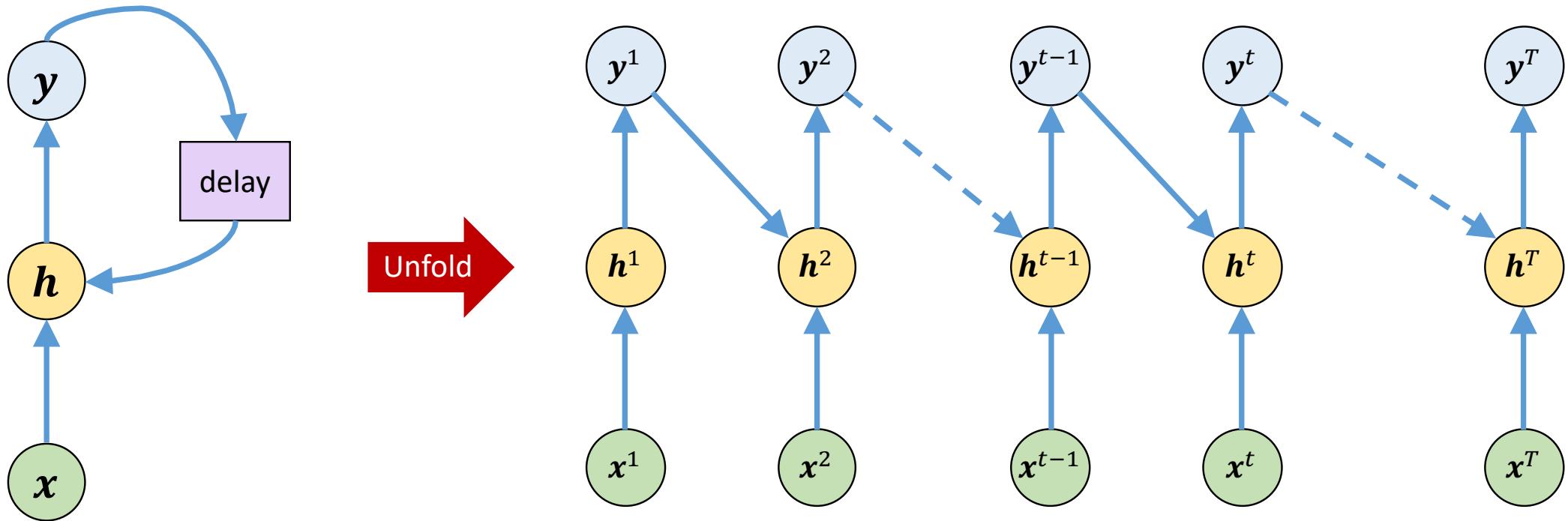
Activation of the Jordan-type RNN with one hidden-layer is given by:

$$\begin{aligned} h(t) &= \phi(U^T x(t) + W^T y(t-1) + b) \\ y(t) &= \sigma(V^T h(t) + c) \end{aligned}$$

Note that output of the previous time instant is fed back to the hidden layer and W represents the recurrent weight matrix connecting previous output to the current hidden input

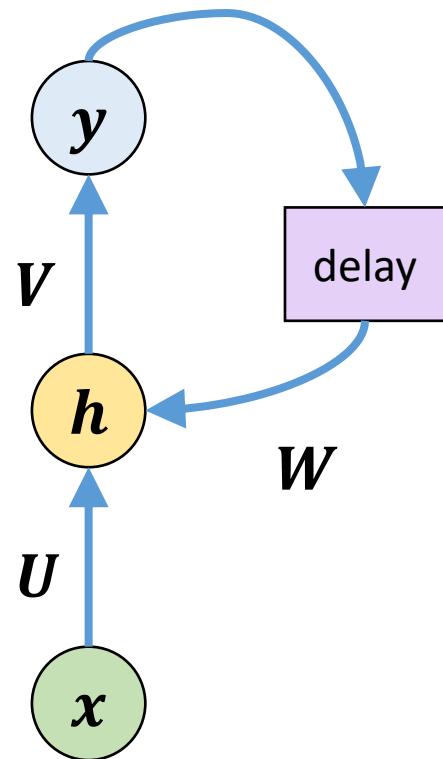


RNN with top-down recurrence



The recurrent connections in the hidden-layer is unfolded and represented in time for processing a time series $(x(t))_{t=1}^T$

RNN with top-down recurrence: batch processing



Given P patterns $\{x_p\}_{p=1}^P$ where $x_p = (x_p(t))_{t=1}^T$,

$$X(t) = \begin{pmatrix} x_1(t)^T \\ x_2(t)^T \\ \vdots \\ x_P(t)^T \end{pmatrix}$$

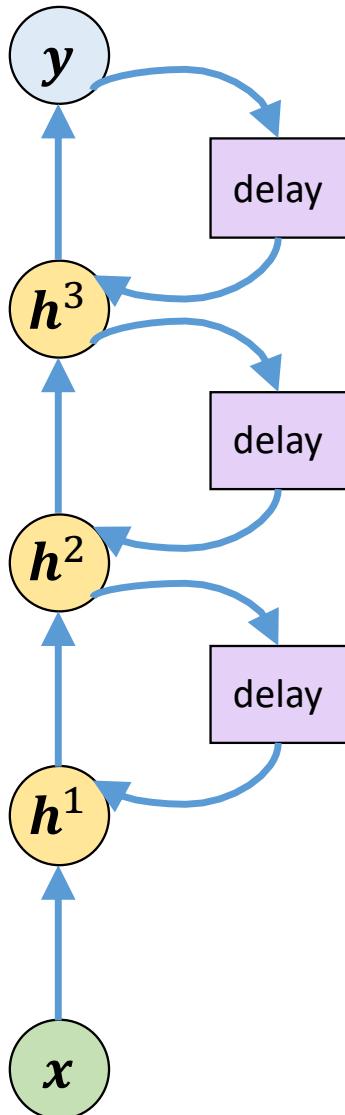
Let $X(t)$, $Y(t)$, and $H(t)$ be batch input, output, and hidden output of the network at time t

Activation of the three-layer Jordan-type RNN is given by:

$$H(t) = \phi(X(t)U + Y(t-1)W + B)$$

$$Y(t) = \sigma(H(t)V + C)$$

Deep RNN with top-down recurrence



\mathbf{U}^l : weight matrix connecting output of hidden-layer $l - 1$ to hidden-layer l .

\mathbf{W}^l : weight matrix connecting hidden-layer output of the previous instant to the current output of the hidden-layer.

\mathbf{b}^l : the bias vector connected to the layer l

First hidden-layer $l = 1$:

$$\mathbf{h}^1(t) = \phi(\mathbf{U}^{1\top} \mathbf{x}(t) + \mathbf{W}^{1\top} \mathbf{h}^2(t-1) + \mathbf{b}^1)$$

For hidden-layer $l = 2, \dots, L-2$:

$$\mathbf{h}^l(t) = \phi(\mathbf{U}^{l\top} \mathbf{h}^{l-1}(t) + \mathbf{W}^{l\top} \mathbf{h}^{l+1}(t-1) + \mathbf{b}^l)$$

For hidden-layer $l = L-1$:

$$\mathbf{h}^{L-1}(t) = \phi(\mathbf{U}^{L-1\top} \mathbf{h}^{L-2}(t) + \mathbf{W}^{L-1\top} \mathbf{y}(t-1) + \mathbf{b}^{L-1})$$

First output layer $l = L$:

$$\mathbf{y}(t) = \sigma(\mathbf{U}^{L\top} \mathbf{h}^{L-1}(t) + \mathbf{b}^L)$$

Example 2

A recurrent neural network with top-down recurrence receives 2-dimensional input sequences and produce 1-dimensional output sequences. It has three hidden neurons and following weight matrices and biases:

Weight matrix connecting input to the hidden layer $\mathbf{U} = \begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix}$

Recurrence weight matrix connecting previous output to the hidden layer $\mathbf{W} = \begin{pmatrix} 2.0 & 1.3 & -1.0 \end{pmatrix}$

Weight matrix connecting hidden layer to the output $\mathbf{V} = \begin{pmatrix} 2.0 \\ -1.5 \\ 0.2 \end{pmatrix}$

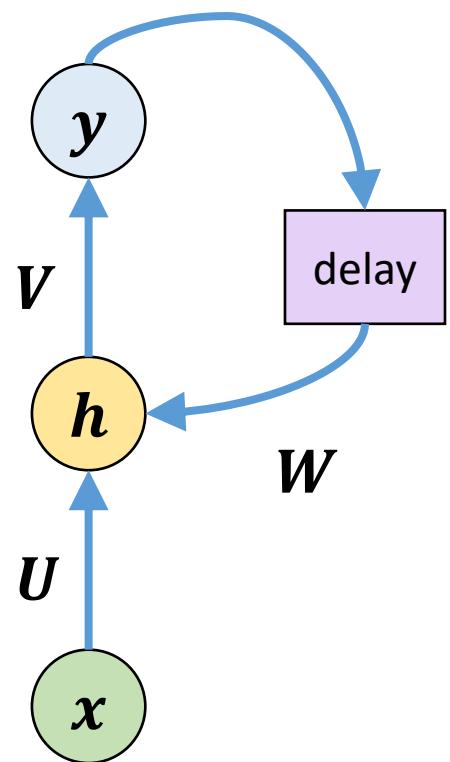
Bias to the hidden layer $\mathbf{b} = \begin{pmatrix} 0.2 \\ 0.2 \\ 0.2 \end{pmatrix}$ and bias to the output layer $\mathbf{c} = 0.1$.

Given the following two input sequences:

$$\mathbf{x}_1 = \left(\begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \begin{pmatrix} -1 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 3 \end{pmatrix} \right)$$
$$\mathbf{x}_2 = \left(\begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 2 \\ -1 \end{pmatrix} \quad \begin{pmatrix} 3 \\ -1 \end{pmatrix} \right)$$

Using batch processing, find output sequences. Assume outputs are initialized to zero at the beginning. Assume tanh and sigmoid activations for hidden and output layer, respectively.

Example 2 (con't)



$$\begin{aligned} \mathbf{x}_1 &= \begin{pmatrix} 1 \\ 2 \\ -1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} -1 \\ 1 \\ 3 \\ 2 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 3 \\ -1 \end{pmatrix} \\ \mathbf{x}_2 &= \begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 2 \\ -1 \end{pmatrix} \quad \begin{pmatrix} 3 \\ -1 \end{pmatrix} \end{aligned}$$

Two sequences as batch of sequences:

$$\mathbf{x} = \begin{pmatrix} \begin{pmatrix} 1 & 2 \\ -1 & 0 \end{pmatrix} & \begin{pmatrix} -1 & 1 \\ 2 & -1 \end{pmatrix} & \begin{pmatrix} 0 & 3 \\ 3 & -1 \end{pmatrix} \end{pmatrix}$$

Three hidden neurons and one output neuron.

$$\begin{aligned} \mathbf{H}(t) &= \phi(\mathbf{X}(t)\mathbf{U} + \mathbf{Y}(t-1)\mathbf{W} + \mathbf{B}) \\ \mathbf{Y}(t) &= \sigma(\mathbf{H}(t)\mathbf{V} + \mathbf{C}) \end{aligned}$$

Initially,

$$\mathbf{Y}(0) = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix}$$

Example 2 (con't)

At $t = 1$: $\mathbf{X}(1) = \begin{pmatrix} 1 & 2 \\ -1 & 0 \end{pmatrix}$

$$\begin{aligned}\mathbf{H}(1) &= \tanh(\mathbf{X}(1)\mathbf{U} + \mathbf{Y}(0)\mathbf{W} + \mathbf{B}) \\ &= \tanh \left(\begin{pmatrix} 1 & 2 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} + \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix} (2.0 \quad 1.3 \quad -1.0) + \begin{pmatrix} 0.2 & 0.2 & 0.2 \\ 0.2 & 0.2 & 0.2 \end{pmatrix} \right) \\ &= \begin{pmatrix} 0.2 & 0.72 & -1.0 \\ 0.83 & -0.29 & 0.0 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{Y}(1) &= \text{sigmoid}(\mathbf{H}(1)\mathbf{V} + \mathbf{C}) \\ &= \text{sigmoid} \left(\begin{pmatrix} 0.2 & 0.72 & -1.0 \\ 0.83 & -0.29 & 0.0 \end{pmatrix} \begin{pmatrix} 2.0 \\ -1.5 \\ 0.2 \end{pmatrix} + \begin{pmatrix} 0.1 \\ 0.1 \end{pmatrix} \right) \\ &= \begin{pmatrix} 0.31 \\ 0.90 \end{pmatrix}\end{aligned}$$

Example 2 (con't)

At $t = 2$: $\mathbf{X}(2) = \begin{pmatrix} -1 & 1 \\ 2 & -1 \end{pmatrix}$

$$\begin{aligned}\mathbf{H}(2) &= \tanh(\mathbf{X}(2)\mathbf{U} + \mathbf{Y}(1)\mathbf{W} + \mathbf{B}) \\ &= \tanh \left(\begin{pmatrix} -1 & 1 \\ 2 & -1 \end{pmatrix} \begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} + \begin{pmatrix} 0.31 \\ 0.9 \end{pmatrix} \begin{pmatrix} 2.0 & 1.3 & -1.0 \end{pmatrix} + \begin{pmatrix} 0.2 & 0.2 & 0.2 \\ 0.2 & 0.2 & 0.2 \end{pmatrix} \right) \\ &= \begin{pmatrix} 0.98 & 0.21 & -0.98 \\ -0.46 & 0.98 & 0.94 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{Y}(2) &= \text{sigmoid}(\mathbf{H}(2)\mathbf{V} + \mathbf{C}) \\ &= \text{sigmoid} \left(\begin{pmatrix} 0.98 & 0.21 & -0.98 \\ -0.46 & 0.98 & 0.94 \end{pmatrix} \begin{pmatrix} 2.0 \\ -1.5 \\ 0.2 \end{pmatrix} + \begin{pmatrix} 0.1 \\ 0.1 \end{pmatrix} \right) \\ &= \begin{pmatrix} 0.83 \\ 0.11 \end{pmatrix}\end{aligned}$$

Example 2 (con't)

At $t = 3$: $\mathbf{X}(3) = \begin{pmatrix} 0 & 3 \\ 3 & -1 \end{pmatrix}$

$$\begin{aligned}\mathbf{H}(3) &= \tanh(\mathbf{X}(3)\mathbf{U} + \mathbf{Y}(1)\mathbf{W} + \mathbf{B}) \\ &= \tanh \left(\begin{pmatrix} 0 & 3 \\ 3 & -1 \end{pmatrix} \begin{pmatrix} -1.0 & 0.5 & 0.2 \\ 0.5 & 0.1 & -2.0 \end{pmatrix} + \begin{pmatrix} 0.83 \\ 0.11 \end{pmatrix} (2.0 \quad 1.3 \quad -1.0) + \begin{pmatrix} 0.2 & 0.2 & 0.2 \\ 0.2 & 0.2 & 0.2 \end{pmatrix} \right) \\ &= \begin{pmatrix} 1.0 & 0.92 & -1.0 \\ -1.00 & 0.94 & 0.99 \end{pmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{Y}(3) &= \text{sigmoid}(\mathbf{H}(3)\mathbf{V} + \mathbf{C}) \\ &= \text{sigmoid} \left(\begin{pmatrix} 1.0 & 0.92 & -1.0 \\ -1.0 & 0.94 & 0.99 \end{pmatrix} \begin{pmatrix} 2.0 \\ -1.5 \\ 0.2 \end{pmatrix} + \begin{pmatrix} 0.1 \\ 0.1 \end{pmatrix} \right) \\ &= \begin{pmatrix} 0.63 \\ 0.04 \end{pmatrix}\end{aligned}$$

Example 2 (con't)

Output (batch):

$$\begin{aligned}\mathbf{Y} &= (\mathbf{Y}(1) \quad \mathbf{Y}(2) \quad \mathbf{Y}(3)) \\ &= \left(\begin{pmatrix} 0.31 \\ 0.9 \end{pmatrix} \quad \begin{pmatrix} 0.83 \\ 0.11 \end{pmatrix} \quad \begin{pmatrix} 0.63 \\ 0.04 \end{pmatrix} \right)\end{aligned}$$

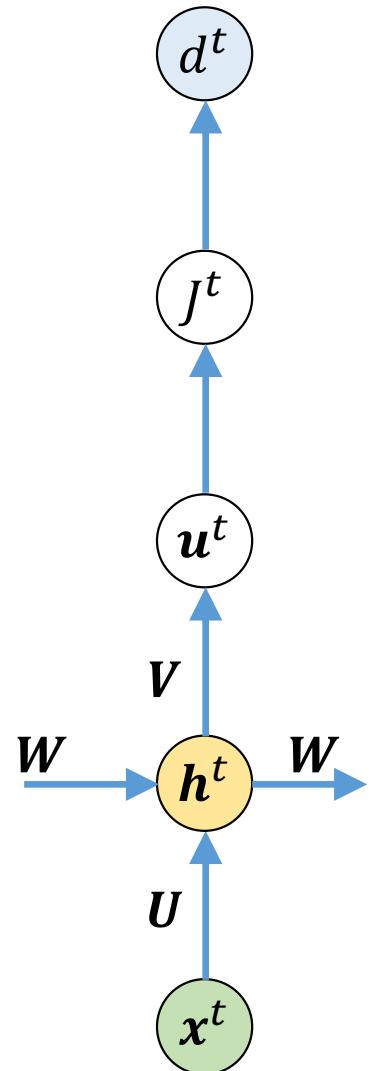
Outputs for inputs

$$\begin{aligned}x_1 &= \left(\begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad \begin{pmatrix} -1 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 3 \end{pmatrix} \right) \\ x_2 &= \left(\begin{pmatrix} -1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 2 \\ -1 \end{pmatrix} \quad \begin{pmatrix} 3 \\ -1 \end{pmatrix} \right)\end{aligned}$$

are

$$\begin{aligned}y_1 &= (0.31, 0.83, 0.63) \\ y_2 &= (0.9, 0.11, 0.04)\end{aligned}$$

Backpropagation through time (BPTT)



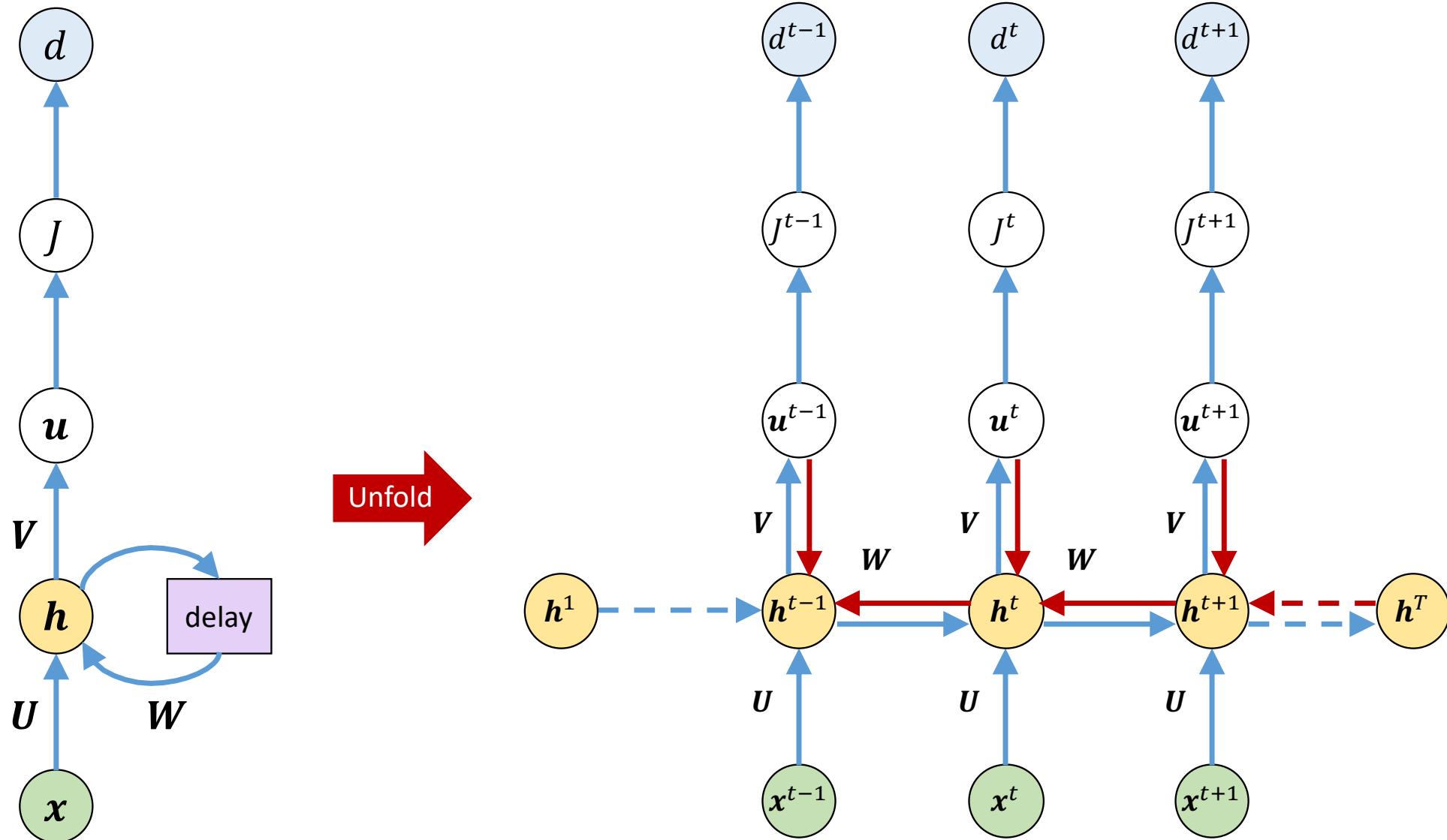
Let's consider vanilla-RNN with hidden recurrence.
Forward propagation equations:

$$\begin{aligned} h(t) &= \tanh(\mathbf{U}^\top \mathbf{x}(t) + \mathbf{W}^\top \mathbf{h}(t-1) + \mathbf{b}) \\ u(t) &= \mathbf{V}^\top \mathbf{h}(t) + \mathbf{c} \end{aligned}$$

For classification,
 $y(t) = \text{softmax}(\mathbf{u}(t))$

For regression,
 $y(t) = \mathbf{u}(t)$

Backpropagation through time (BPTT)



Backpropagation through time (BPTT)

Total cost:

$$J = \sum_{t=1}^T J(t)$$

where T is the length of the sequence and $J(t)$ represent the cost at the instant t .

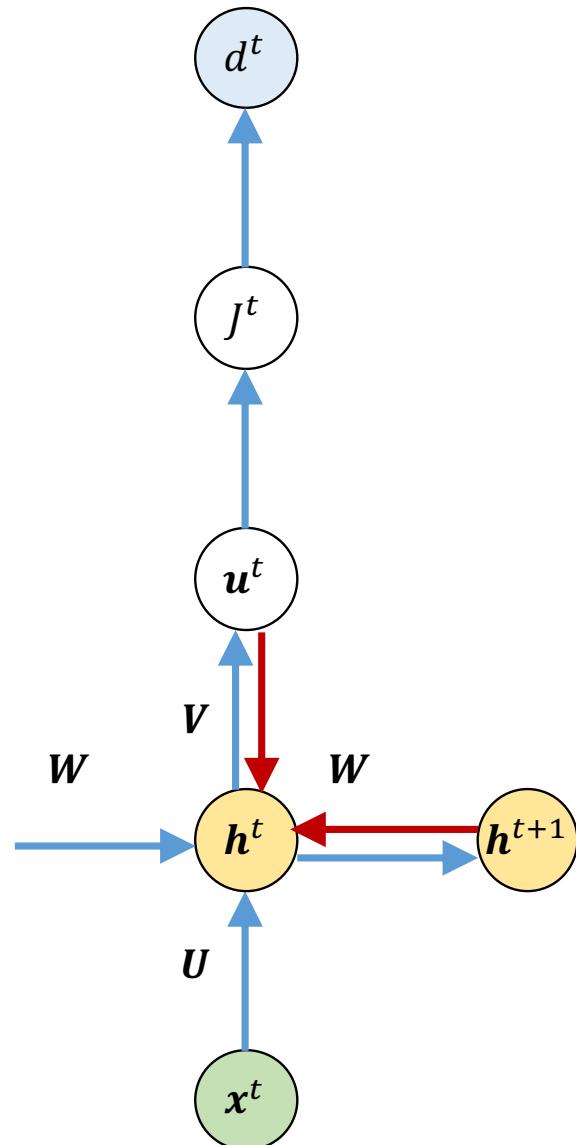
Note that

$$\frac{\partial J}{\partial J(t)} = 1$$

Error gradients at the output:

$$\nabla_{\mathbf{u}(t)} J = \nabla_{\mathbf{u}(t)} J(t) = \begin{cases} -\left(1(\mathbf{k} = d(t)) - \sigma(\mathbf{u}(t))\right), & \text{for softmax} \\ -(d(t) - \mathbf{u}(t)), & \text{for linear regression} \end{cases}$$

Backpropagation through time (BPTT)



The gradients propagate backward, starting from the end of the sequence from two directions: top-down direction and reverse sequence direction.

At the final time step T , $\mathbf{h}(t)$ has only one descendant.

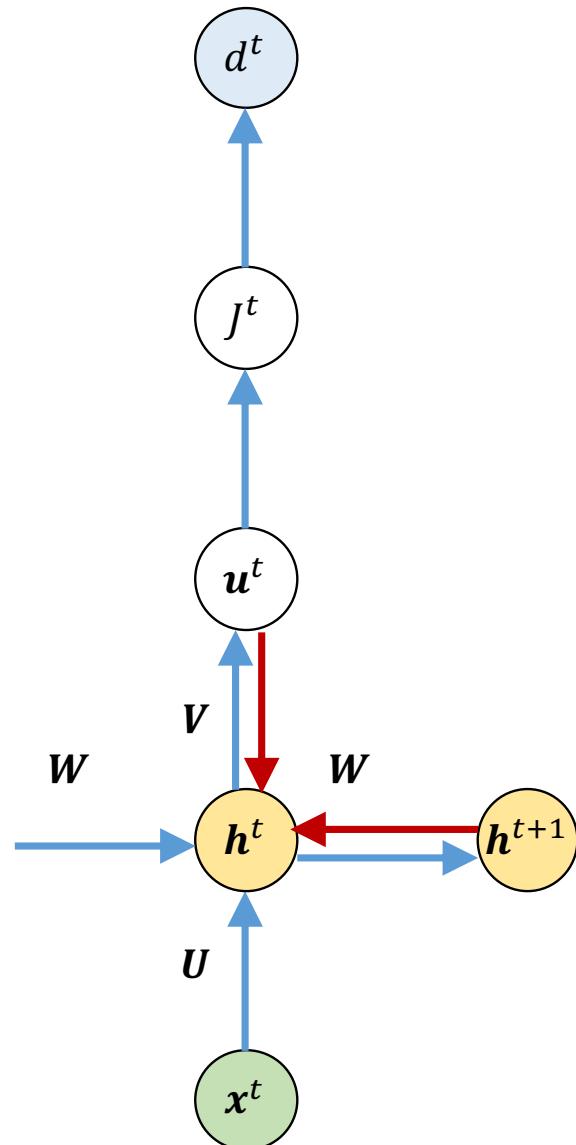
For $t = T$:

$$\nabla_{\mathbf{h}(T)} J = \mathbf{V}^T \nabla_{\mathbf{u}(T)} J \quad (\text{A})$$

For $t < T$, $\mathbf{h}(t)$ has descendants of gradients from both $\mathbf{h}(t)$ and $\mathbf{u}(t)$. That is, for $t = 1, 2, \dots, T-1$, from chain rule for multi-dimensions,

$$\nabla_{\mathbf{h}(t)} J = \left(\frac{\partial \mathbf{h}(t+1)}{\partial \mathbf{h}(t)} \right)^T \nabla_{\mathbf{h}(t+1)} J + \left(\frac{\partial \mathbf{u}(t)}{\partial \mathbf{h}(t)} \right)^T \nabla_{\mathbf{u}(t)} J \quad (\text{B})$$

Backpropagation through time (BPTT)



For $t = 1, 2, \dots, T - 1$

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

where:

$$\frac{\partial \mathbf{h}(t+1)}{\partial \mathbf{h}(t)} = \frac{\partial \mathbf{h}(t+1)}{\partial \mathbf{z}(t+1)} \frac{\partial \mathbf{z}(t+1)}{\partial \mathbf{h}(t)}$$

and $\mathbf{z}(t)$ is the synaptic input to the hidden-layer at instant t and the hidden-layer has a \tanh activation function.

$$\frac{\partial \mathbf{h}(t+1)}{\partial \mathbf{z}(t+1)} = \phi'(\mathbf{z}(t+1)) = \text{diag}(1 - \mathbf{h}^2(t+1)) \quad (\text{C})$$

And

$$\frac{\partial \mathbf{z}(t+1)}{\partial \mathbf{h}(t)} = \mathbf{W}^T \quad \text{and} \quad \frac{\partial \mathbf{u}(t)}{\partial \mathbf{h}(t)} = \mathbf{V}^T \quad (\text{D})$$

Backpropagation through time (BPTT)

Substituting (C) and (D) in (B);

$$\nabla_{\mathbf{h}(t)} J = \mathbf{W} \text{diag}(1 - \mathbf{h}^2(t+1)) \nabla_{\mathbf{h}(t+1)} J + \mathbf{V} \nabla_{\mathbf{u}(t)} J \quad (\text{E})$$

Starting from (A):

$$\nabla_{\mathbf{h}(T)} J = \mathbf{V} \nabla_{\mathbf{u}(T)} J$$

The above equation (E) can be recursively used to compute $\nabla_{\mathbf{h}(t)} J$ backwards.

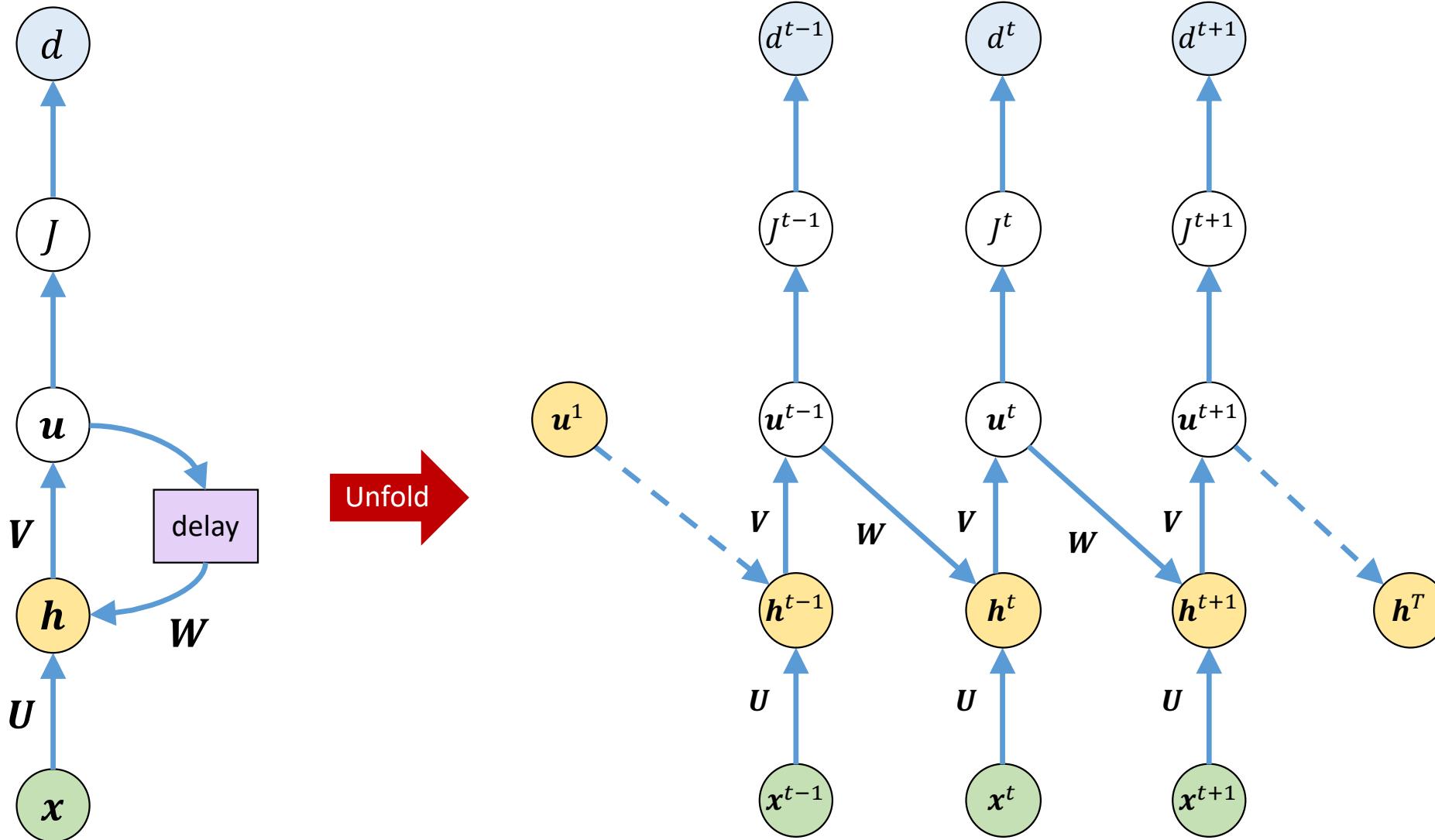
Note that each backward propagation through time, the $\nabla_{\mathbf{h}(t)} J$ is multiplied by \mathbf{W}^T .

Backpropagation through time (BPTT)

Using gradients $\nabla_{\mathbf{u}(t)}J$ and $\nabla_{\mathbf{h}(t)}J$, the error gradients with respect to the parameters $\{\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{b}, \mathbf{c}\}$ can be obtained:

$$\begin{aligned}\nabla_{\mathbf{c}}J &= \sum_t \frac{\partial J(t)}{\partial \mathbf{c}} = \sum_t \left(\frac{\partial \mathbf{u}(t)}{\partial \mathbf{c}} \right)^T \nabla_{\mathbf{u}(t)}J = \sum_t \nabla_{\mathbf{u}(t)}J \\ \nabla_{\mathbf{b}}J &= \sum_t \left(\frac{\partial \mathbf{h}(t)}{\partial \mathbf{b}} \right)^T \nabla_{\mathbf{h}(t)}J = \sum_t \left(\frac{\partial \mathbf{h}(t)}{\partial \mathbf{z}(t)} \frac{\partial \mathbf{z}(t)}{\partial \mathbf{b}} \right)^T \nabla_{\mathbf{h}(t)}J = \sum_t \text{diag}(1 - \mathbf{h}(t)^2) \nabla_{\mathbf{h}(t)}J \\ \nabla_{\mathbf{V}}J &= \sum_t \mathbf{h}(t) (\nabla_{\mathbf{u}(t)}J)^T \\ \nabla_{\mathbf{W}}J &= \sum_t \mathbf{h}(t-1) (\nabla_{\mathbf{z}(t)}J)^T = \sum_t \mathbf{h}(t-1) (\nabla_{\mathbf{h}(t)}J)^T \text{diag}(1 - \mathbf{h}(t)^2) \\ \nabla_{\mathbf{U}}J &= \sum_t \mathbf{x}(t) (\nabla_{\mathbf{z}(t)}J)^T = \sum_t \mathbf{x}(t) (\nabla_{\mathbf{h}(t)}J)^T \text{diag}(1 - \mathbf{h}(t)^2)\end{aligned}$$

BPTT for RNN with top-down recurrence



Backpropagation through time (BPTT)

The gradient computation involves performing a **forward propagation pass moving left to right** through the unfolded graph, followed by a **backpropagation pass moving right to left** through the graph. The gradients are propagated from the final time point to the initial time points .

The runtime is $O(T)$ where T is the length of input sequence and cannot be reduced by parallelization because the forward propagation is inherently **sequential**. The back-propagation algorithm applied to the unrolled graph with $O(T)$ cost is called *back-propagation through time (BPTT)*.

The network with recurrence between hidden units is thus very powerful but also **expensive** to train.

Example 3

Generate 8-dimensional 16 input sequences of 64 time-steps with each input is a random number between [0.0, 1.0].

Generate the corresponding 1-dimensional labels by randomly generating a number from [0, 1, 2].

Create an RNN with one hidden layer of 5 neurons.

Plot the learning curves and predicted labels.

Example 3 (con't)

```
# Initialize hyper-parameters
n_in = 8
n_hidden = 5
n_out = 3
n_steps = 64
n_seqs = 16
n_iters = 25000
lr = 0.0001
```

```
# Generate training data
x_train = tf.constant(np.random.rand(n_seqs, n_steps, n_in), dtype=tf.float32)
x_train = tf.split(x_train, n_steps, axis=1)
y_train = np.random.randint(size=(n_seqs, n_steps), low=0, high=n_out)
y_train = tf.reshape(tf.constant(y_train, dtype=tf.float32), [n_seqs, n_steps, 1])
```

Example 3 (con't)

```
class RNN(Model):
    def __init__(self, n_in, n_hidden, n_out):
        super(RNN, self).__init__()
        self.n_hidden = n_hidden

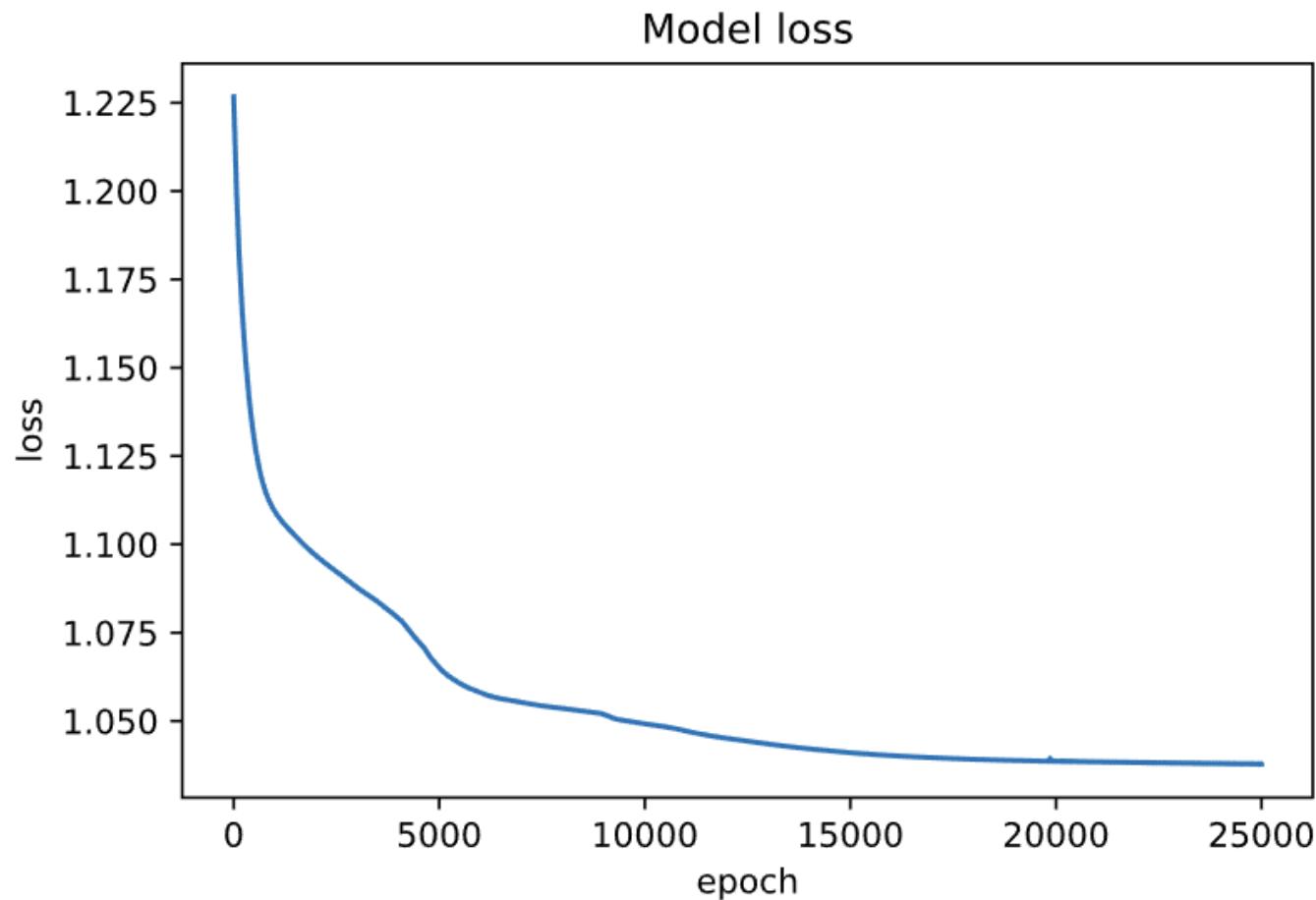
        # Specify the weights and biases
        self.u = tf.Variable(tf.random.truncated_normal([n_in, n_hidden], stddev=1 / np.sqrt(n_in)))
        self.w = tf.Variable(tf.random.truncated_normal([n_hidden, n_hidden], stddev=1 / np.sqrt(n_hidden)))
        self.v = tf.Variable(tf.random.truncated_normal([n_hidden, n_out], stddev=1 / np.sqrt(n_hidden)))
        self.b = tf.Variable(tf.zeros([n_hidden]))
        self.c = tf.Variable(tf.zeros([n_out]))

    def call(self, x):
        # The initial hidden state, which is a zero vector
        h = tf.constant(np.array([0] * self.n_hidden), dtype=tf.float32)
        h = tf.tile(tf.reshape(h, [1, h.shape[0]]), [x[0].shape[0], 1])

        ys = []
        for i in range(0, len(x)):
            h = tf.tanh(tf.matmul(tf.squeeze(x[i]), self.u) + tf.matmul(h, self.w) + self.b)
            u = tf.matmul(h, self.v) + self.c
            ys.append(u)

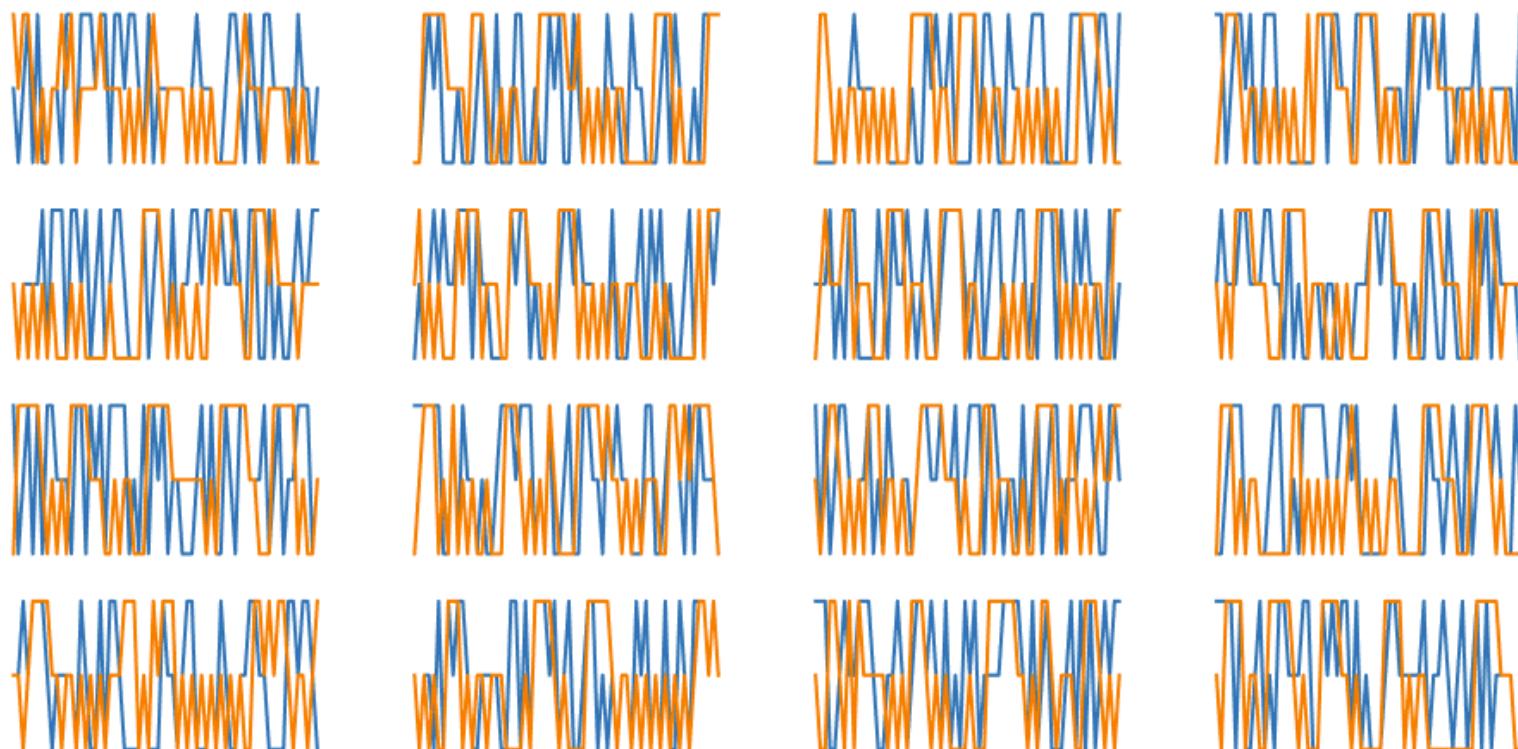
        return tf.stack(ys, axis=1)
```

Example 3 (con't)

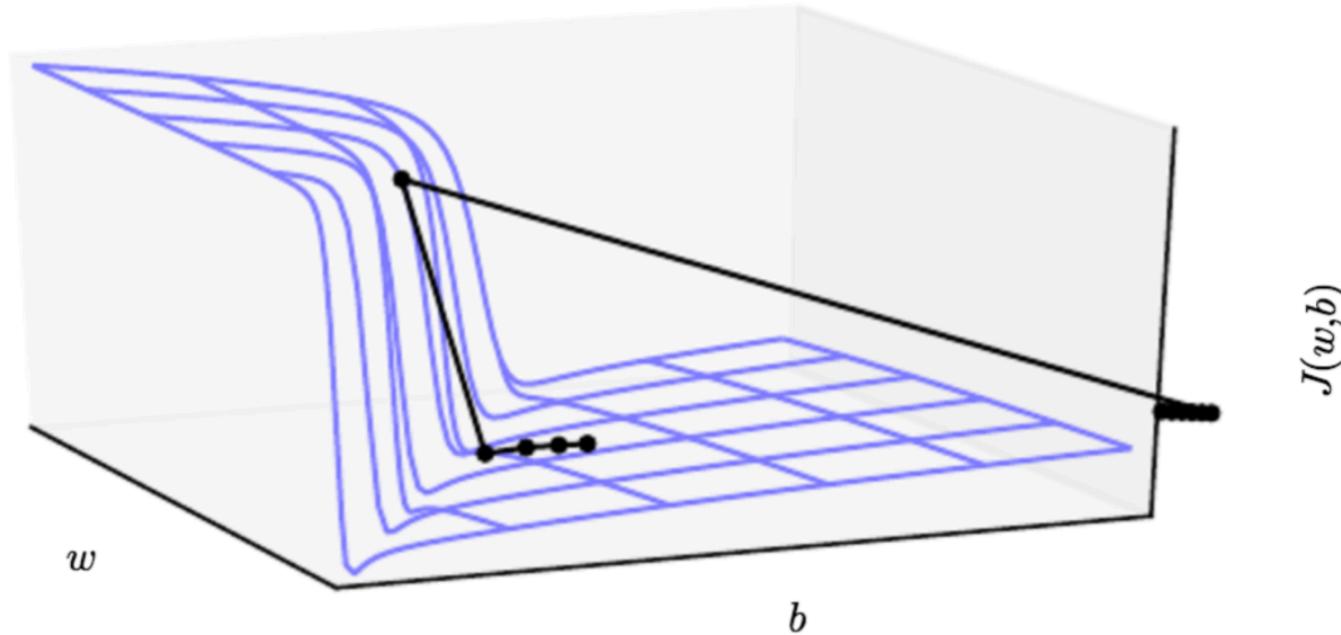


Example 3 (con't)

Prediction (orange) vs Groundtruth (blue)



Gradient Clipping



Due to long term dependencies, RNN tend to have gradients having very large or very small magnitudes. The large gradients resemble cliffs in the error landscape and when the gradient descent encounters the gradient updates can move the parameters away from true minimum.

A gradient clipping is employed to avoid the gradients to become too large.

Gradient Clipping

Commonly, two methods are used:

1. Clip the gradient g when it exceeds a threshold:

```
if ||g|| > ν:  
    ||g|| ← ν
```

2. Normalize the gradient when it exceeds a threshold:

```
if ||g|| > ν:  
    ||g|| ←  $\frac{g}{||g||} \nu$ 
```

Image Credits

- <https://cs.stanford.edu/people/karpathy/sfmltalk.pdf>
- <https://developer.ibm.com/articles/cc-cognitive-recurrent-neural-networks/>