

CS2109S: Introduction to AI and Machine Learning

Lecture 11: Neural Networks on Sequential Data

8 April 2025

PollEv.com/conghuihu365

Outline

- Recurrent Neural Networks
 - Motivation
 - Recurrent Neural Networks
 - Applications
- Self-Attention
 - Self-Attention Layer
 - Positional Encoding
- Issues with Deep Learning
 - Overfitting
 - Vanishing/Exploding Gradient

Outline

- **Recurrent Neural Networks**
 - Motivation
 - Recurrent Neural Networks
 - Applications
- Self-Attention
 - Self-Attention Layer
 - Positional Encoding
- Issues with Deep Learning
 - Overfitting
 - Vanishing/Exploding Gradient

Sequential Data

Sequential data refers to data where the order of elements matters, and each element depends on its position in the sequence.

Text data:

“we saw this saw”

Audio data:



Video data:

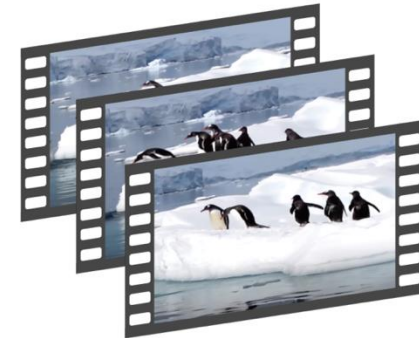


Image credit:

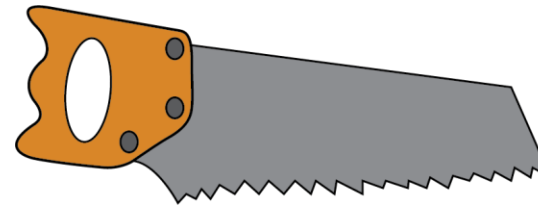
<https://towardsdatascience.com/all-you-need-to-know-to-start-speech-processing-with-deep-learning-102c916edf62> 4
https://openaccess.thecvf.com/content/ICCV2021/papers/Arnab_ViViT_A_Video_Vision_Transformer_ICCV_2021_paper.pdf

Sequential Data

Text data: “we saw this saw”

Verb? Verb? Verb? Verb?

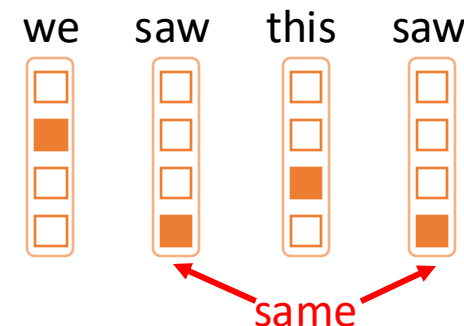
True Label: No Yes No No



One-hot Encoding:

- For each word, create a vector with length equal to the vocabulary size.
- Each word is assigned a unique index, and its corresponding vector is all zeros except for a 1 in the position of that index.

Word	Index	One-hot encoded vector
apple	0	[1, 0, 0, 0]
we	1	[0, 1, 0, 0]
this	2	[0, 0, 1, 0]
saw	3	[0, 0, 0, 1]

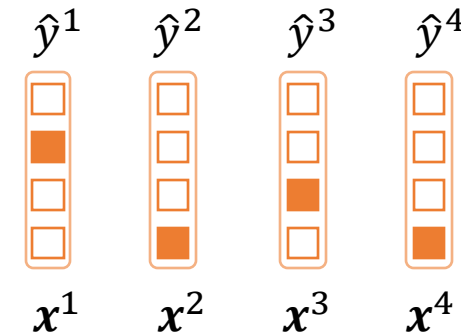
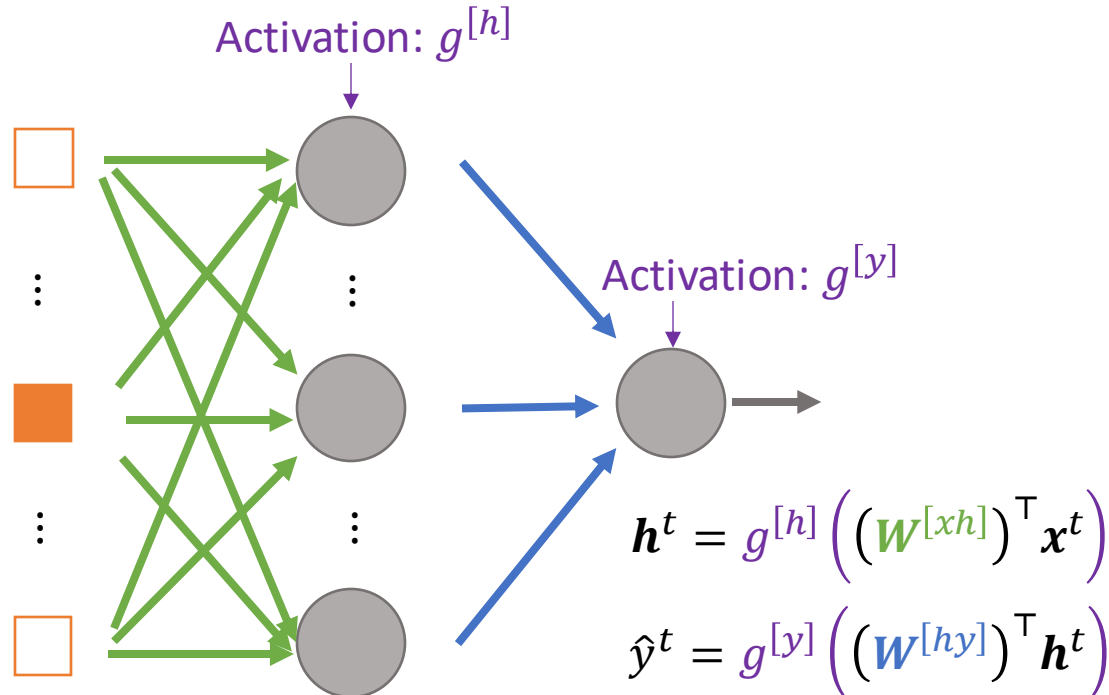


Sequential Data: A Naïve Attempt

Text data: “we saw this saw”

Verb? Verb? Verb? Verb?

True Label: No Yes No No

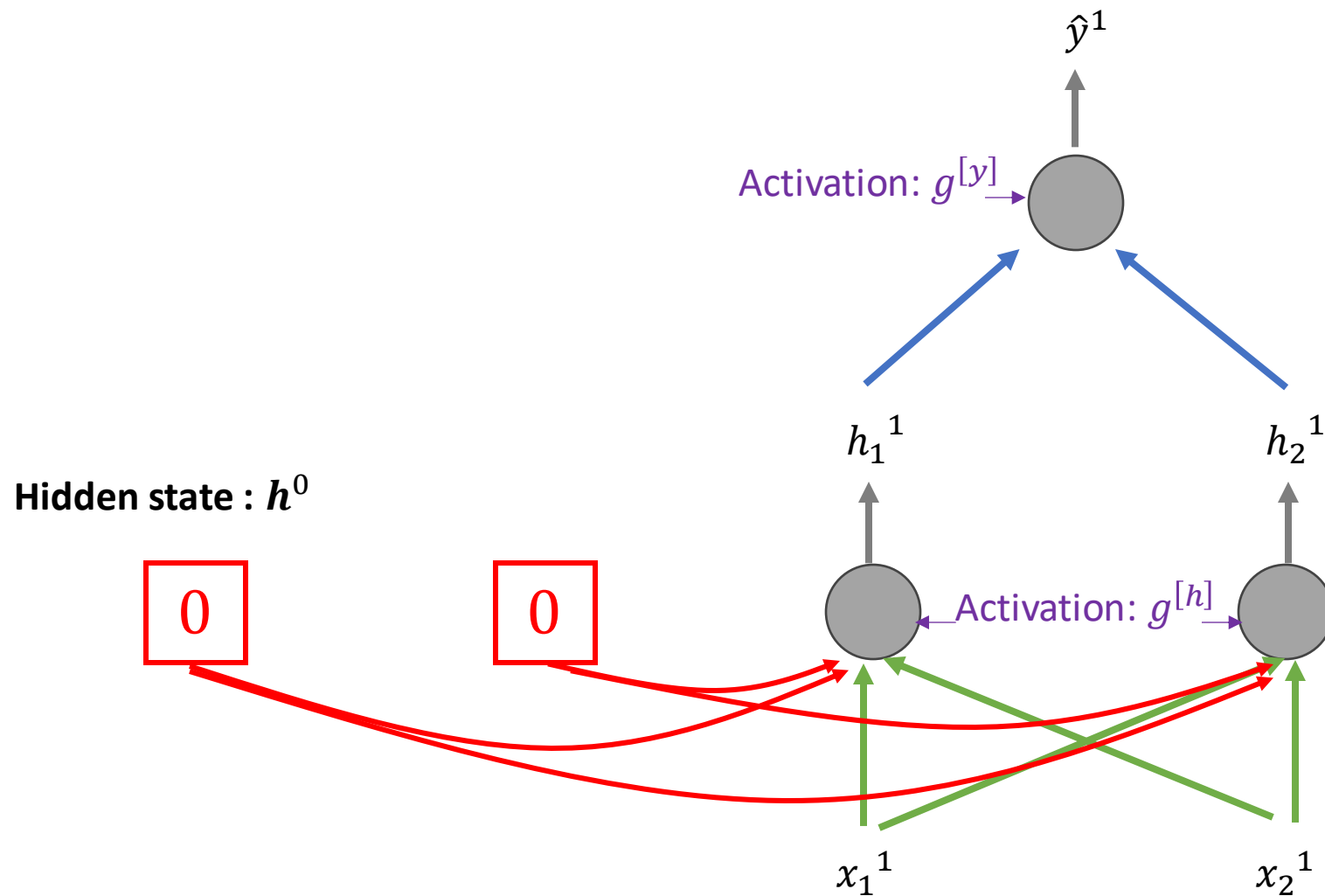


\mathbf{x}^2 and \mathbf{x}^4 are both one-hot vectors representing “saw”.

- ➡ $\mathbf{x}^2 = \mathbf{x}^4$
- ➡ $\hat{\mathbf{y}}^2 = \hat{\mathbf{y}}^4$
- ➡ For the two “saw”, one of the predicted labels must be incorrect.
- ➡ We should not predict each label independently.

Obtain contextual information by using RNN!

Recurrent Neural Network



A sequence with 3 elements

$$\mathbf{x}^1 = \begin{bmatrix} x_1^1 \\ x_2^1 \end{bmatrix}$$

$$\mathbf{x}^2 = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$

$$\mathbf{x}^3 = \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}$$

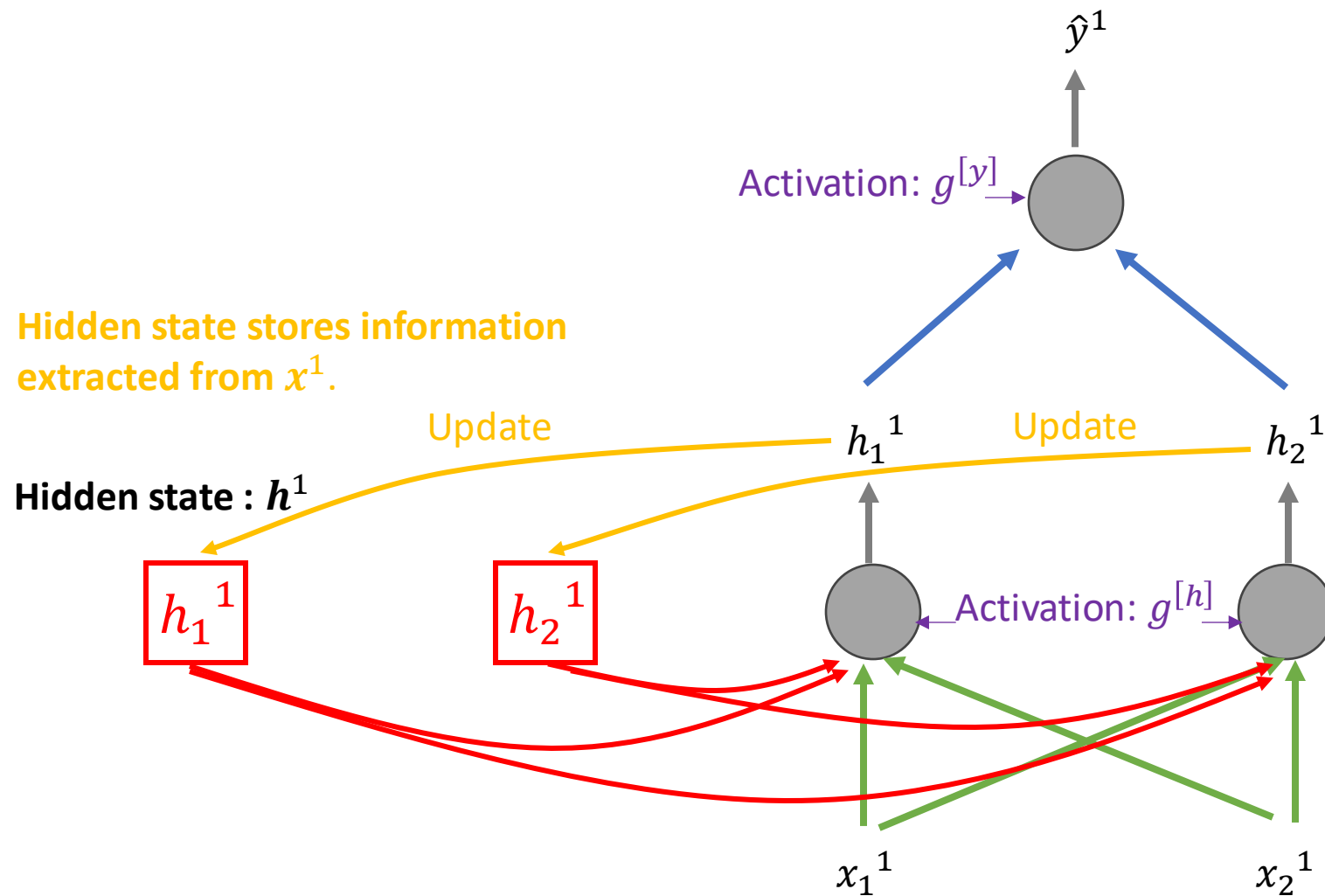
A time step is a single iteration in an RNN where it processes one element of a sequence.

Time step 1:

Hidden state: $\mathbf{h}^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\mathbf{h}^1 = g[h] \left((\mathbf{w}^{[xh]})^\top \mathbf{x}^1 + (\mathbf{w}^{[hh]})^\top \mathbf{h}^0 \right)$$

Recurrent Neural Network



$$x^1 = \begin{bmatrix} x_1^1 \\ x_2^1 \end{bmatrix}$$

$$x^2 = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$

$$x^3 = \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}$$

Time step 1:

Hidden state: $h^1 = \begin{bmatrix} h_1^1 \\ h_2^1 \end{bmatrix}$

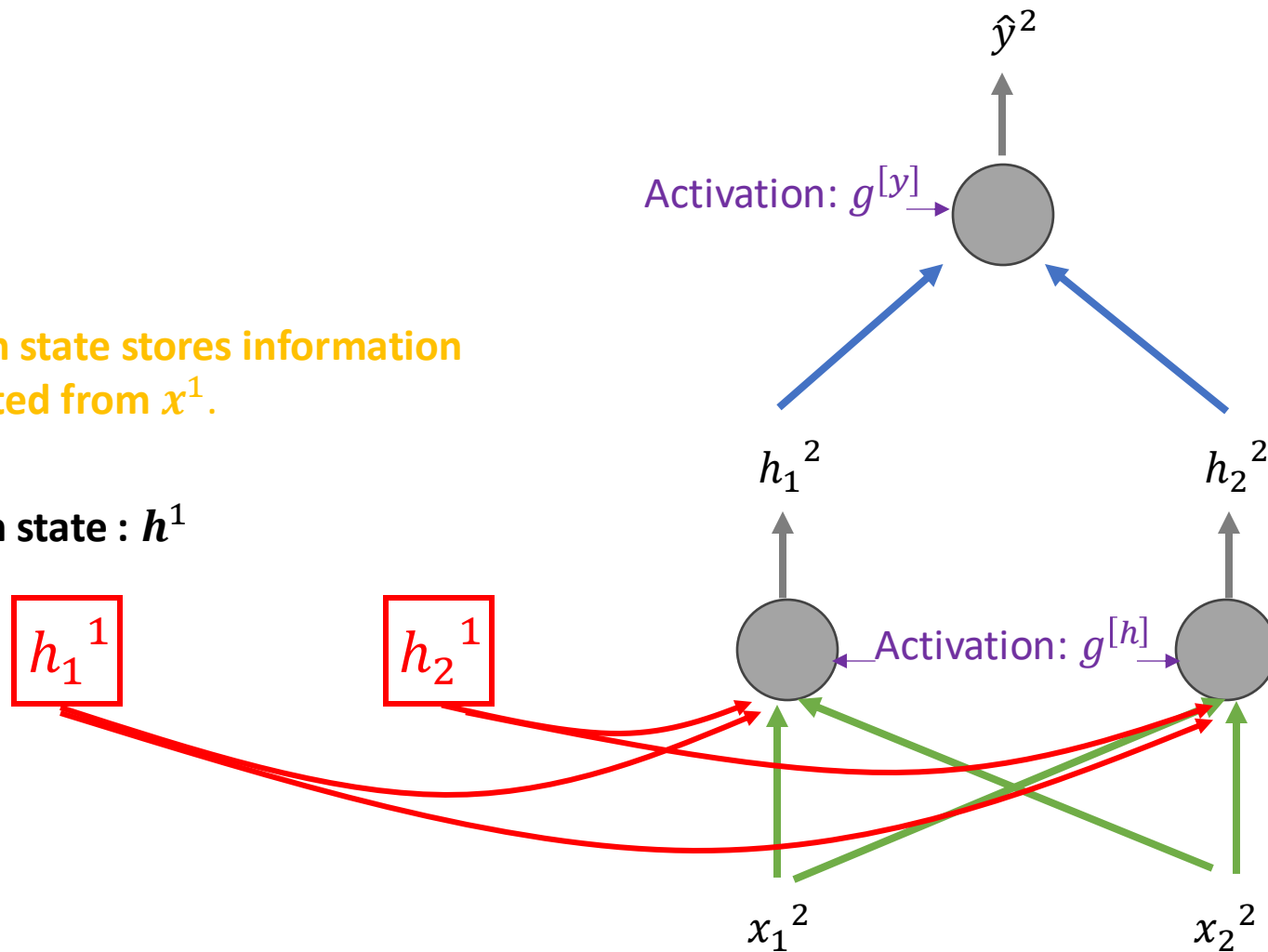
$$h^1 = g^{[h]} \left((W^{[xh]})^\top x^1 + (W^{[hh]})^\top h^0 \right)$$

$$\hat{y}^1 = g^{[y]} \left((W^{[hy]})^\top h^1 \right)$$

Recurrent Neural Network

Hidden state stores information extracted from x^1 .

Hidden state : h^1



$$x^1 = \begin{bmatrix} x_1^1 \\ x_2^1 \end{bmatrix}$$

$$x^2 = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$

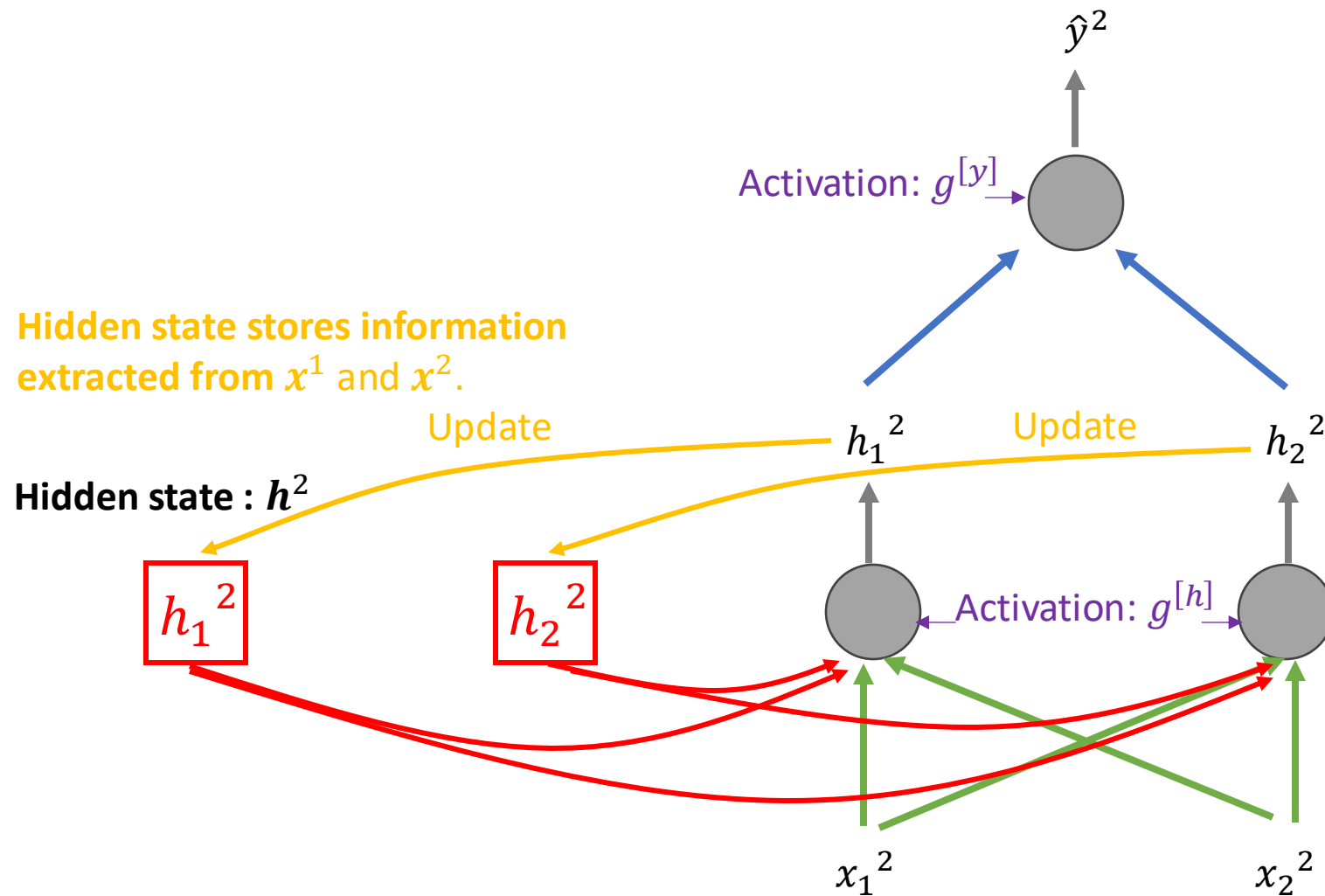
$$x^3 = \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}$$

Time step 2:

Hidden state: $h^1 = \begin{bmatrix} h_1^1 \\ h_2^1 \end{bmatrix}$

$$h^2 = g[h] \left((W^{[xh]})^\top x^2 + (W^{[hh]})^\top h^1 \right)$$

Recurrent Neural Network



$$x^1 = \begin{bmatrix} x_1^1 \\ x_2^1 \end{bmatrix}$$

$$x^2 = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$

$$x^3 = \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}$$

Time step 2:

Hidden state: $h^2 = \begin{bmatrix} h_1^2 \\ h_2^2 \end{bmatrix}$

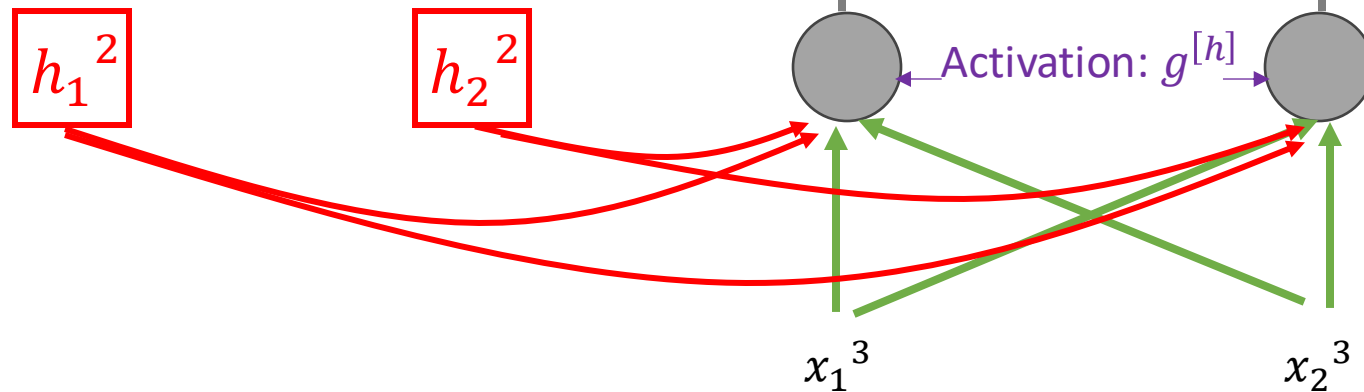
$$h^2 = g^{[h]} \left((W^{[xh]})^\top x^2 + (W^{[hh]})^\top h^1 \right)$$

$$\hat{y}^2 = g^{[y]} \left((W^{[hy]})^\top h^2 \right)$$

Recurrent Neural Network

Hidden state stores information extracted from x^1 and x^2 .

Hidden state : h^2



$$x^1 = \begin{bmatrix} x_1^1 \\ x_2^1 \end{bmatrix}$$

$$x^2 = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$

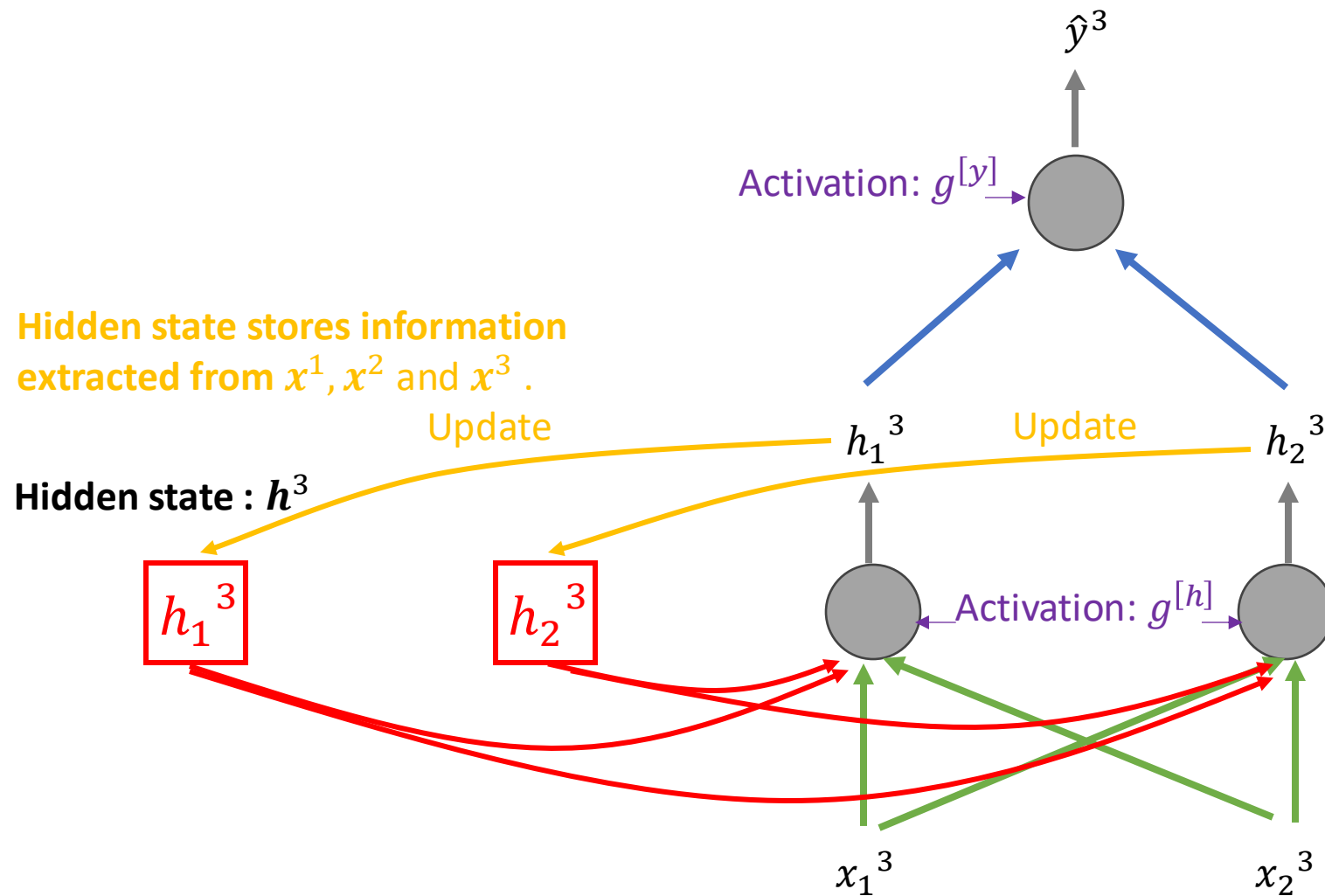
$$x^3 = \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}$$

Time step 3:

Hidden state: $h^2 = \begin{bmatrix} h_1^2 \\ h_2^2 \end{bmatrix}$

$$h^3 = g[h] \left((W^{[xh]})^\top x^3 + (W^{[hh]})^\top h^2 \right)$$

Recurrent Neural Network



$$x^1 = \begin{bmatrix} x_1^1 \\ x_2^1 \end{bmatrix}$$

$$x^2 = \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$

$$x^3 = \begin{bmatrix} x_1^3 \\ x_2^3 \end{bmatrix}$$

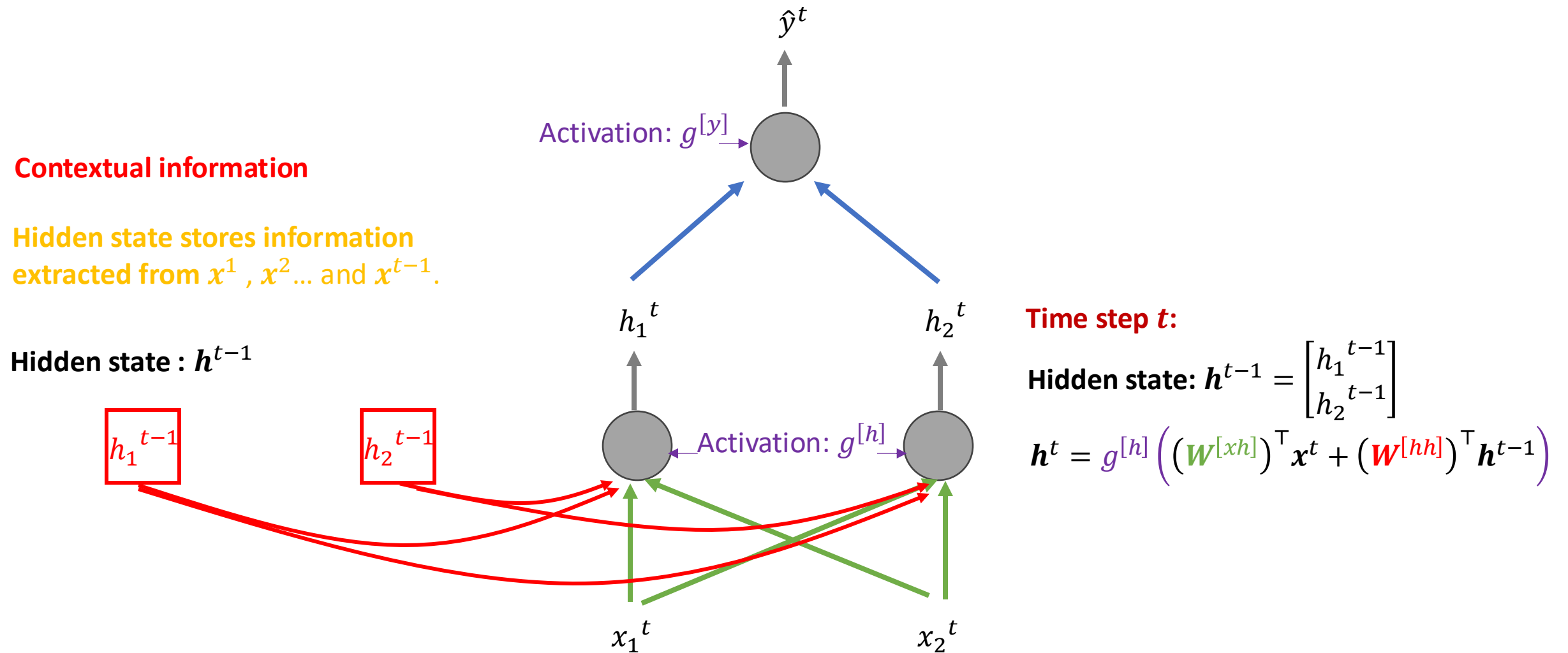
Time step 3:

Hidden state: $h^3 = \begin{bmatrix} h_1^3 \\ h_2^3 \end{bmatrix}$

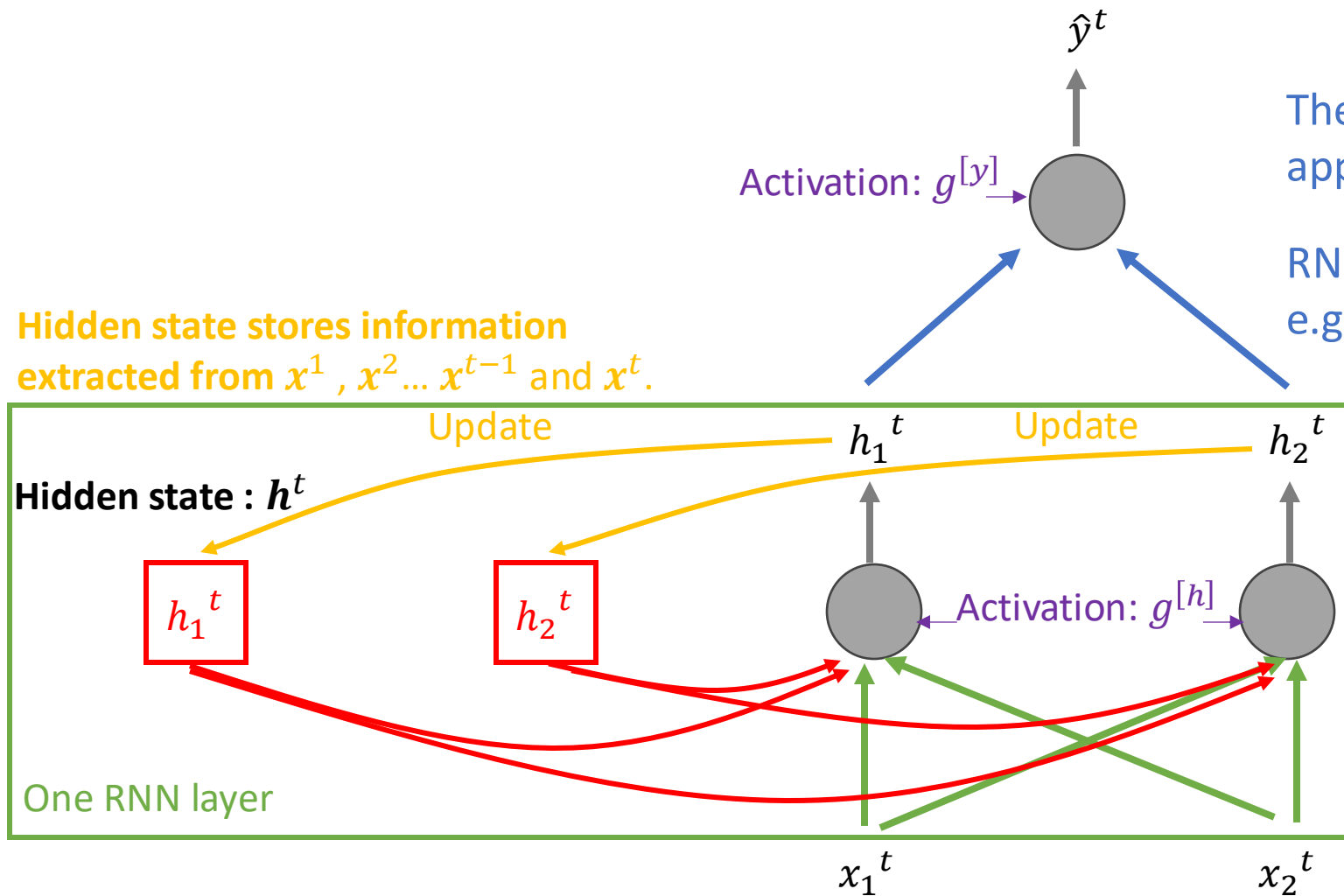
$$h^3 = g^{[h]} \left((W^{[xh]})^T x^3 + (W^{[hh]})^T h^2 \right)$$

$$\hat{y}^3 = g^{[y]} \left((W^{[hy]})^T h^3 \right)$$

Recurrent Neural Network



Recurrent Neural Network



The same weights ($W^{[xh]}$, $W^{[hh]}$, $W^{[hy]}$) are applied at each time step.

RNN can handle sentences of varying lengths, e.g., we saw this saw v.s. we saw this old saw.

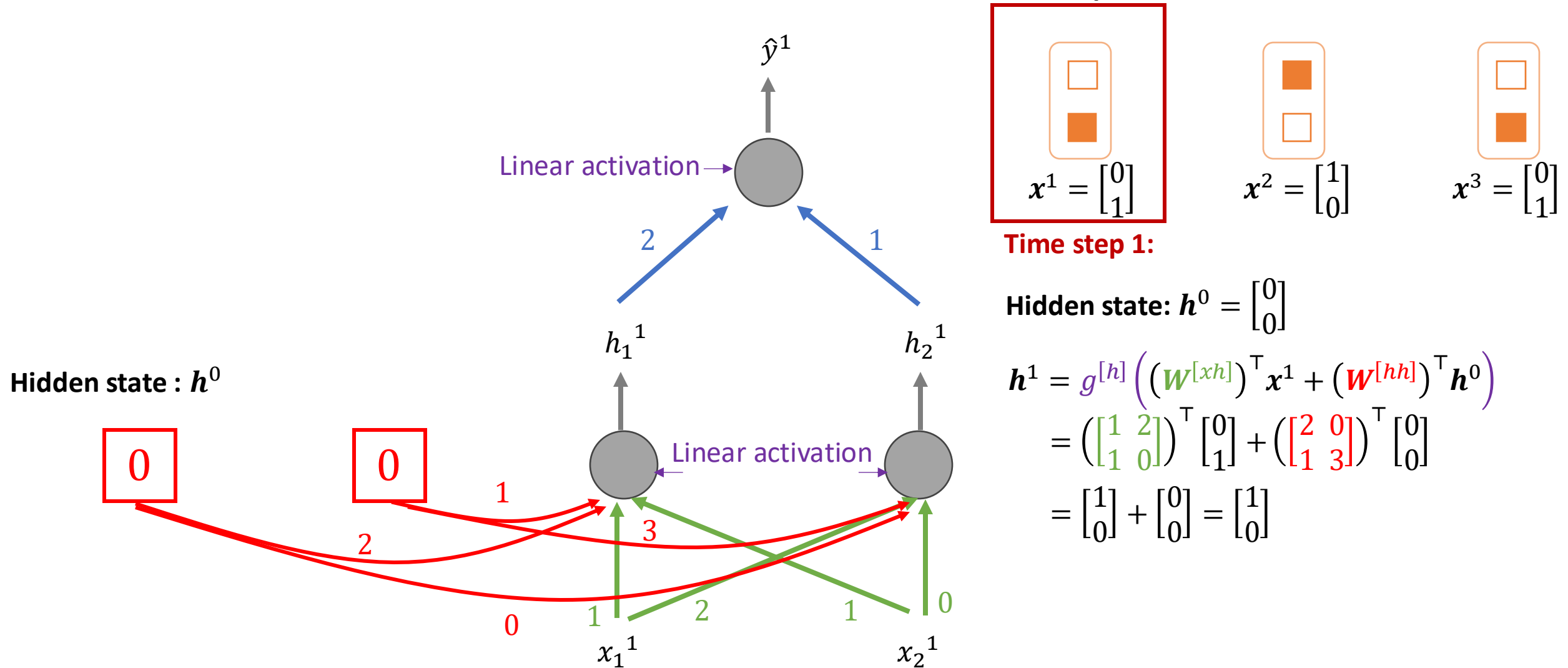
Time step t :

Hidden state: $\mathbf{h}^t = \begin{bmatrix} h_1^t \\ h_2^t \end{bmatrix}$

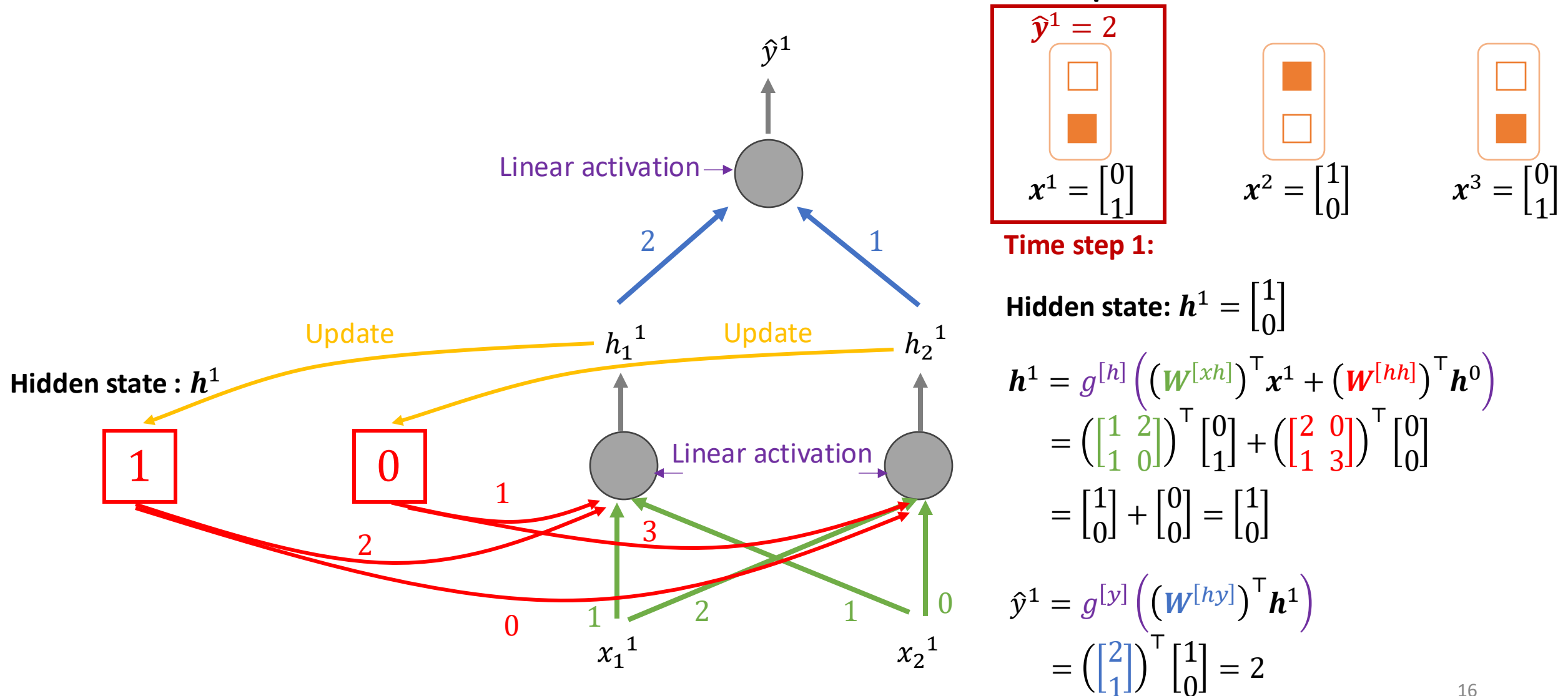
$$\mathbf{h}^t = g^{[h]} \left((W^{[xh]})^\top \mathbf{x}^t + (W^{[hh]})^\top \mathbf{h}^{t-1} \right)$$

$$\hat{y}^t = g^{[y]} \left((W^{[hy]})^\top \mathbf{h}^t \right)$$

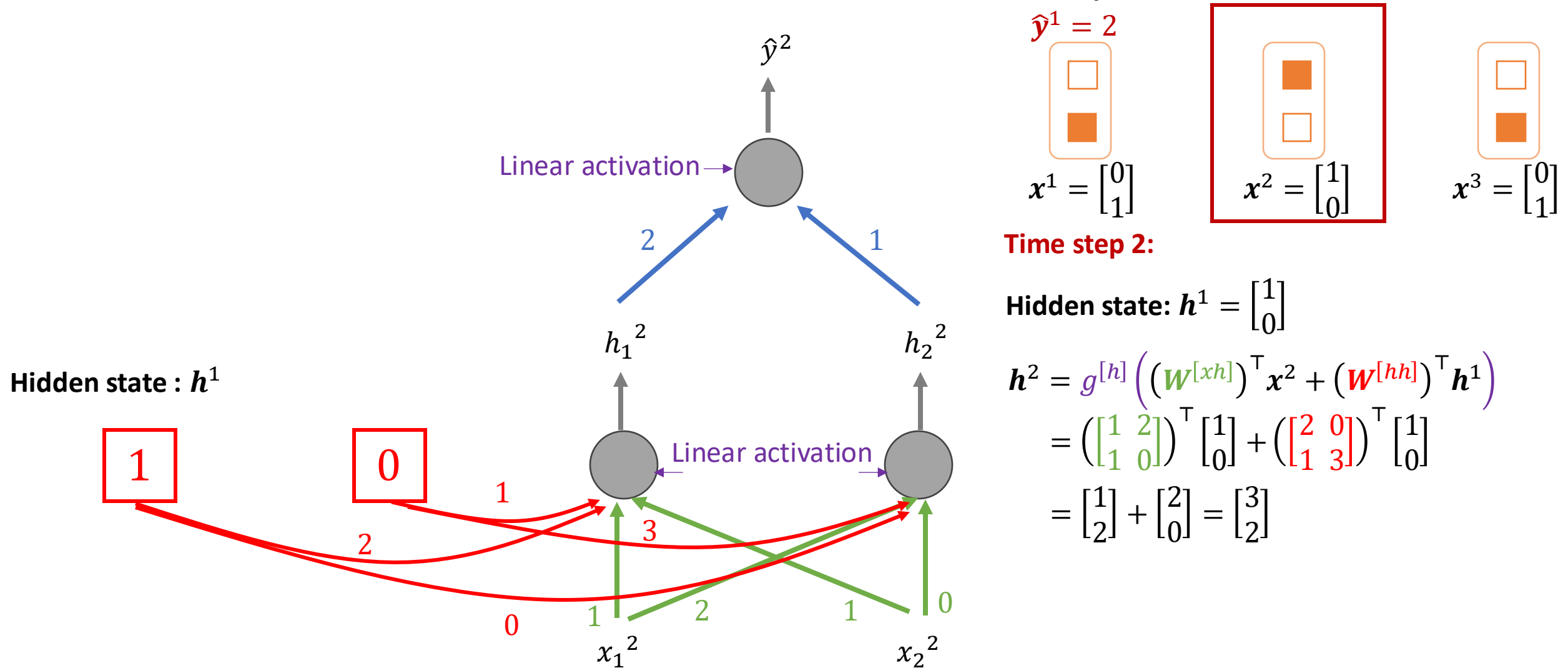
Recurrent Neural Network: Example



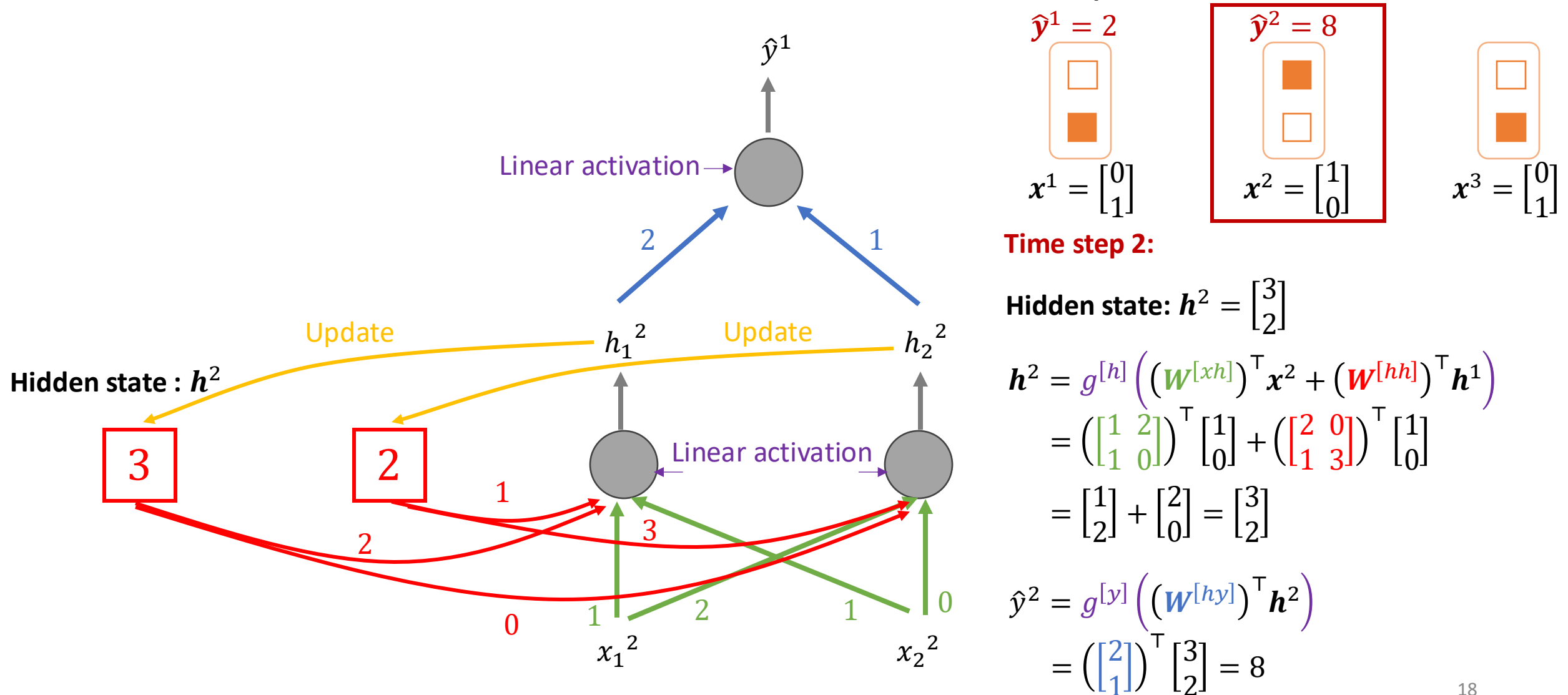
Recurrent Neural Network: Example



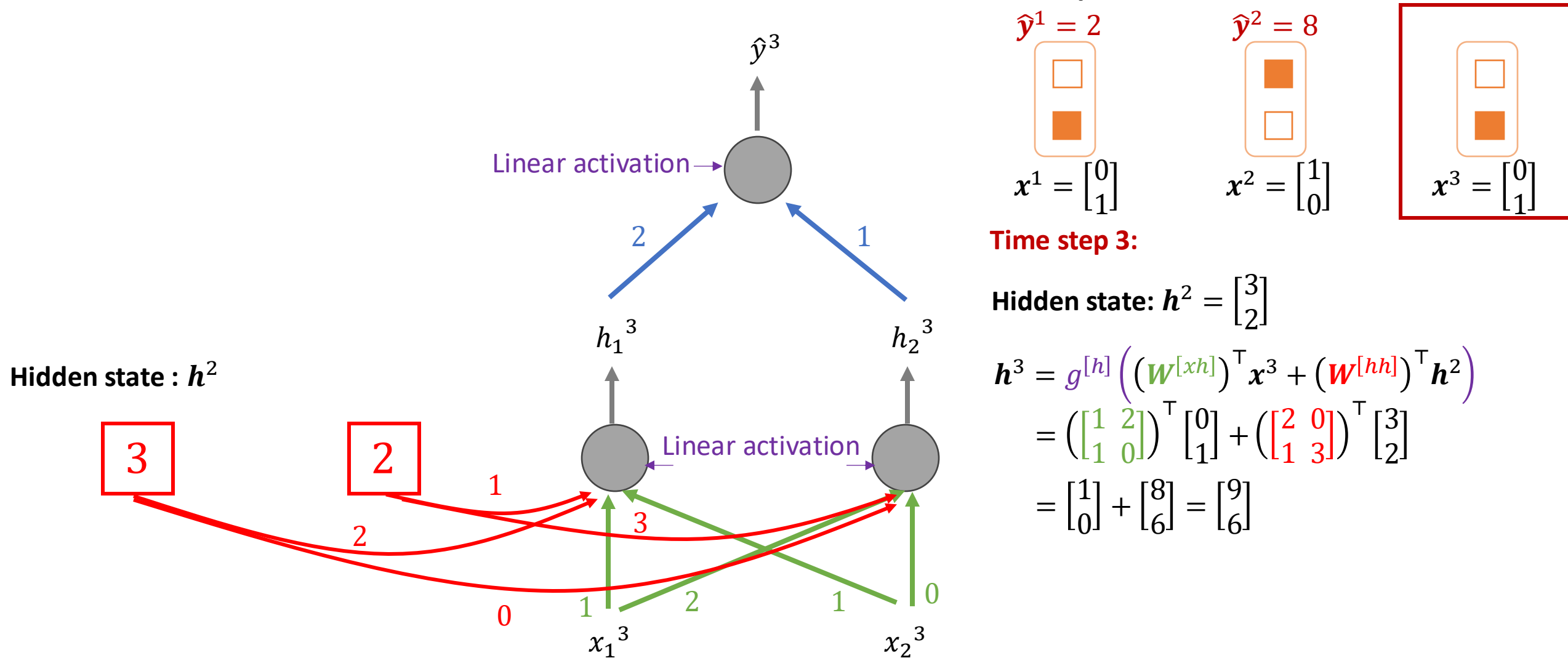
Recurrent Neural Network: Example



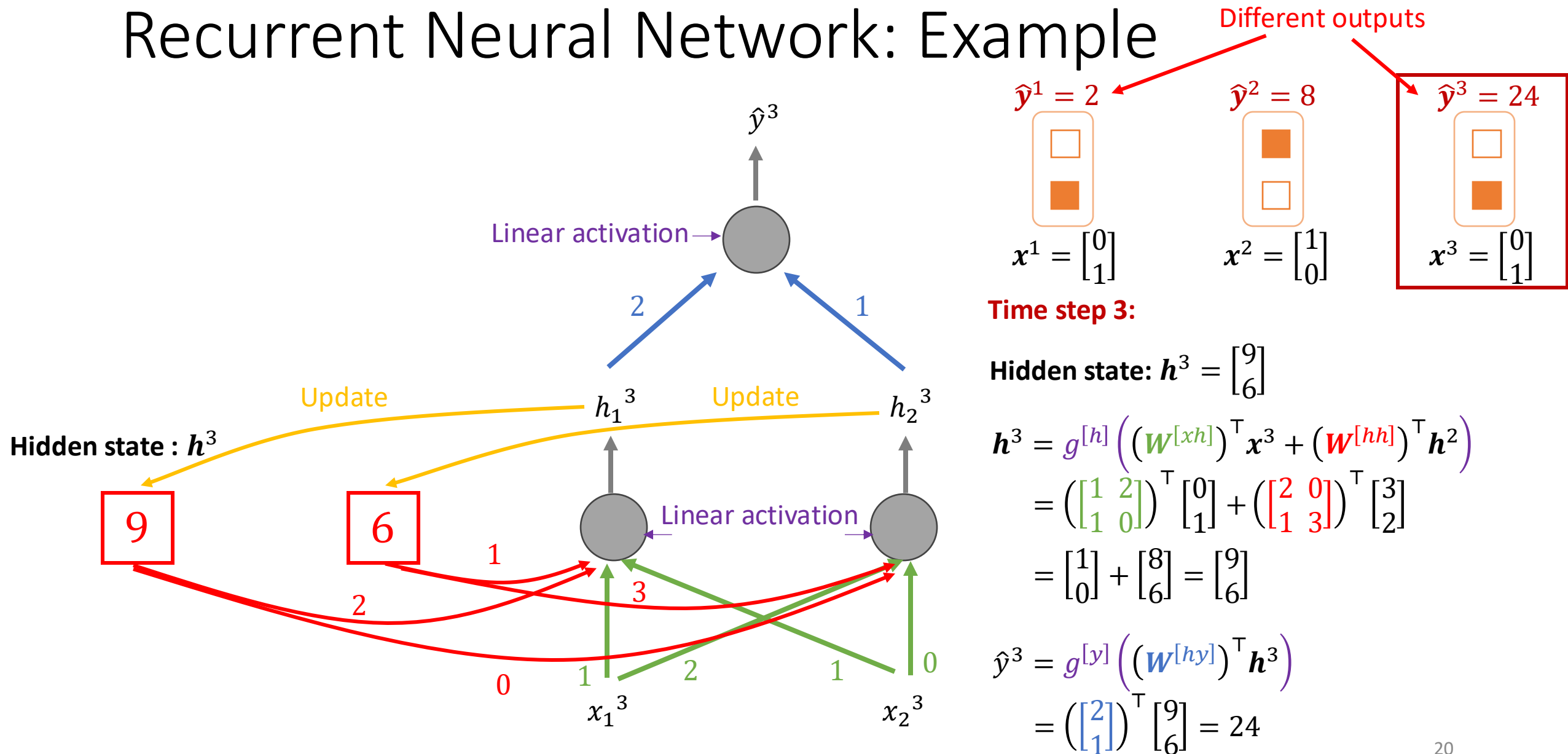
Recurrent Neural Network: Example



Recurrent Neural Network: Example

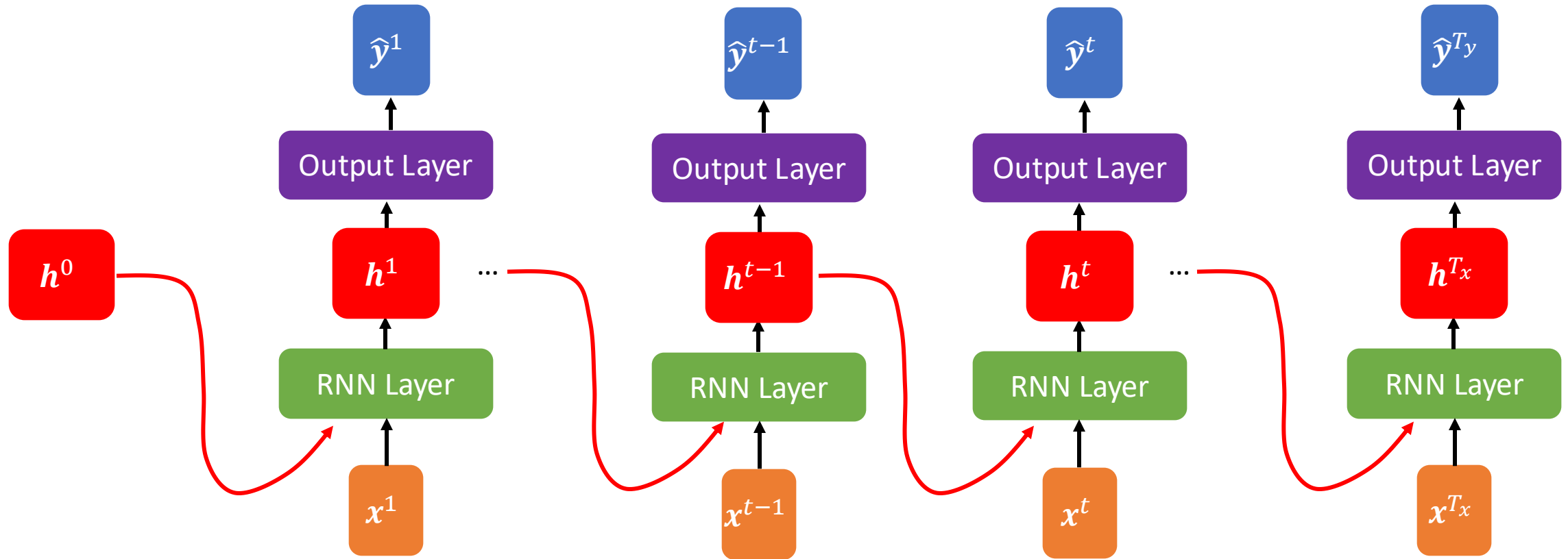


Recurrent Neural Network: Example



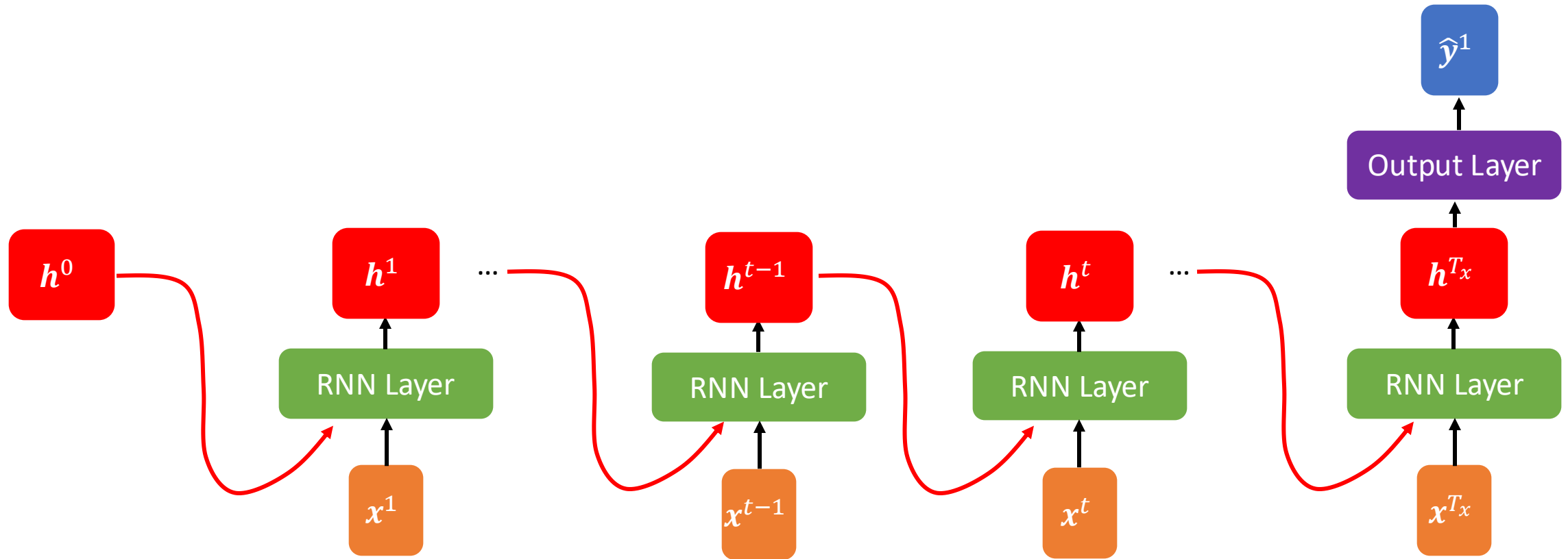
Recurrent Neural Network: Many-to-Many

$$T_x = T_y > 1$$



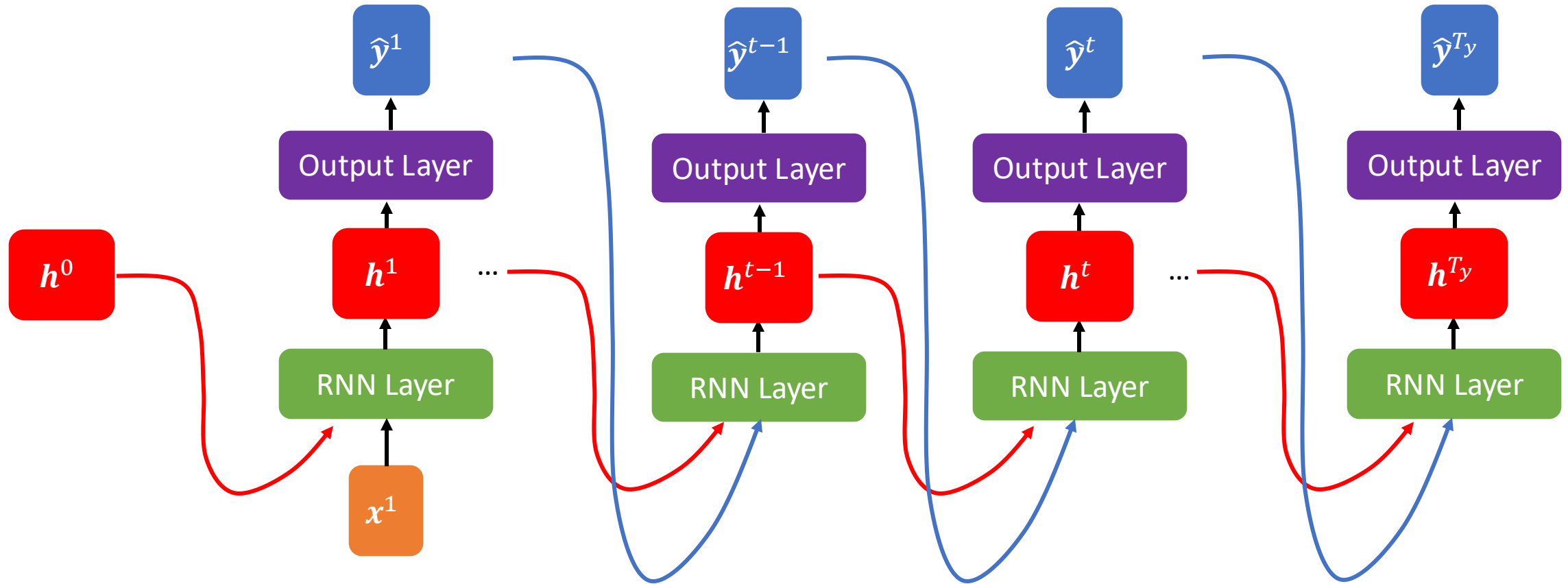
Recurrent Neural Network: Many-to-One

$$T_x > 1, T_y = 1$$



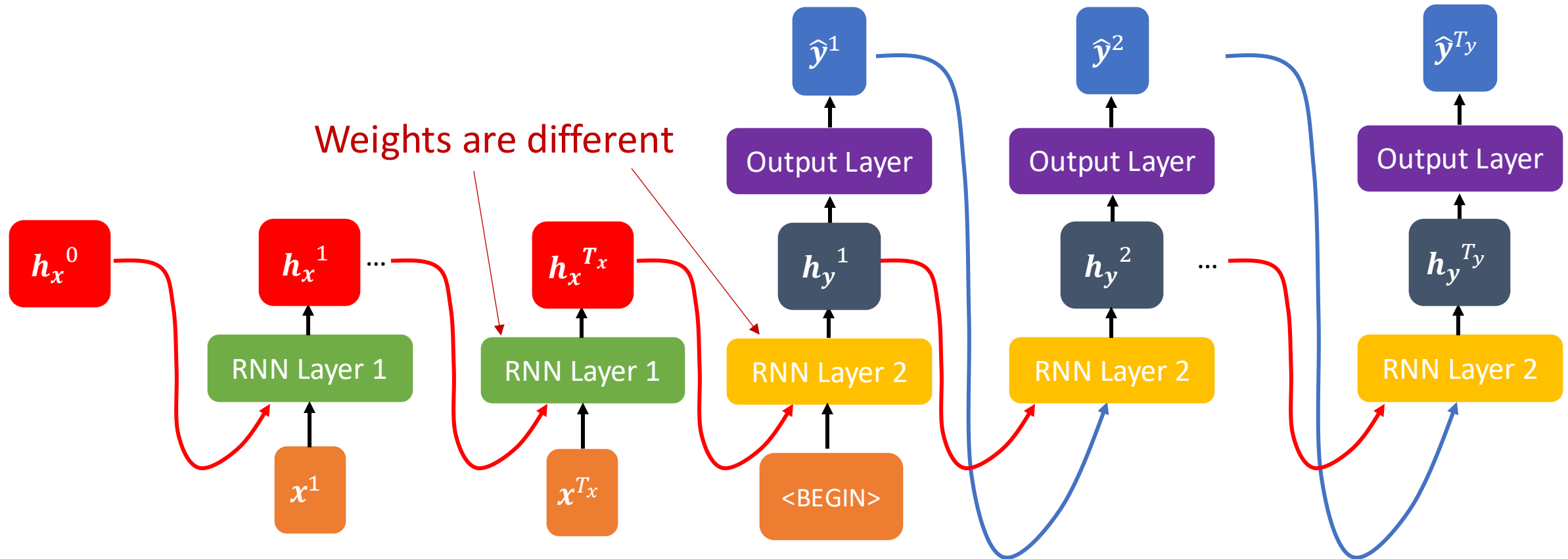
Recurrent Neural Network: One-to-Many

$$T_x = 1, T_y > 1$$

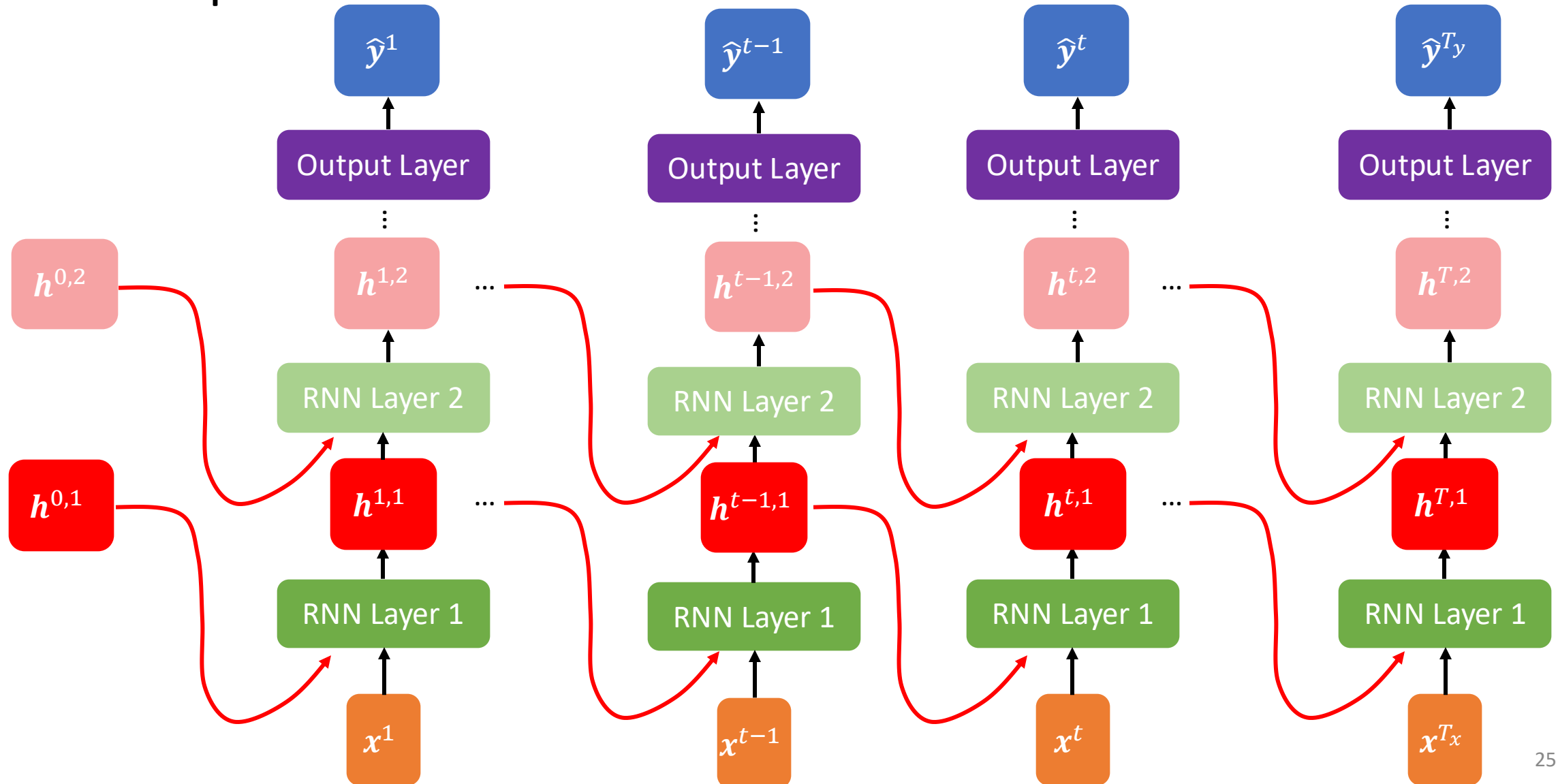


Recurrent Neural Network: Many-to-Many

$$T_x \neq T_y, T_x > 1, T_y > 1$$



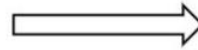
Deep Recurrent Neural Networks



Applications of RNN

Sentiment Analysis:

“Decent effort. The plot
could have been better.”

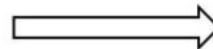
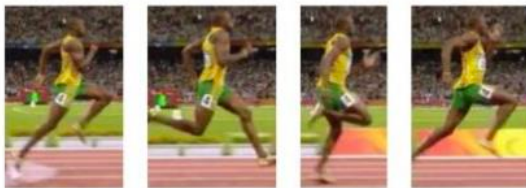


Speech Recognition:





“Fuzzy Wuzzy was a bear.
Fuzzy Wuzzy had no hair.”

Video Captioning:



“A man is running.”

Properties of RNN

- Capture contextual information 
- The prediction at time step t must wait until all previous steps have been completed.  **Not parallelism-friendly**

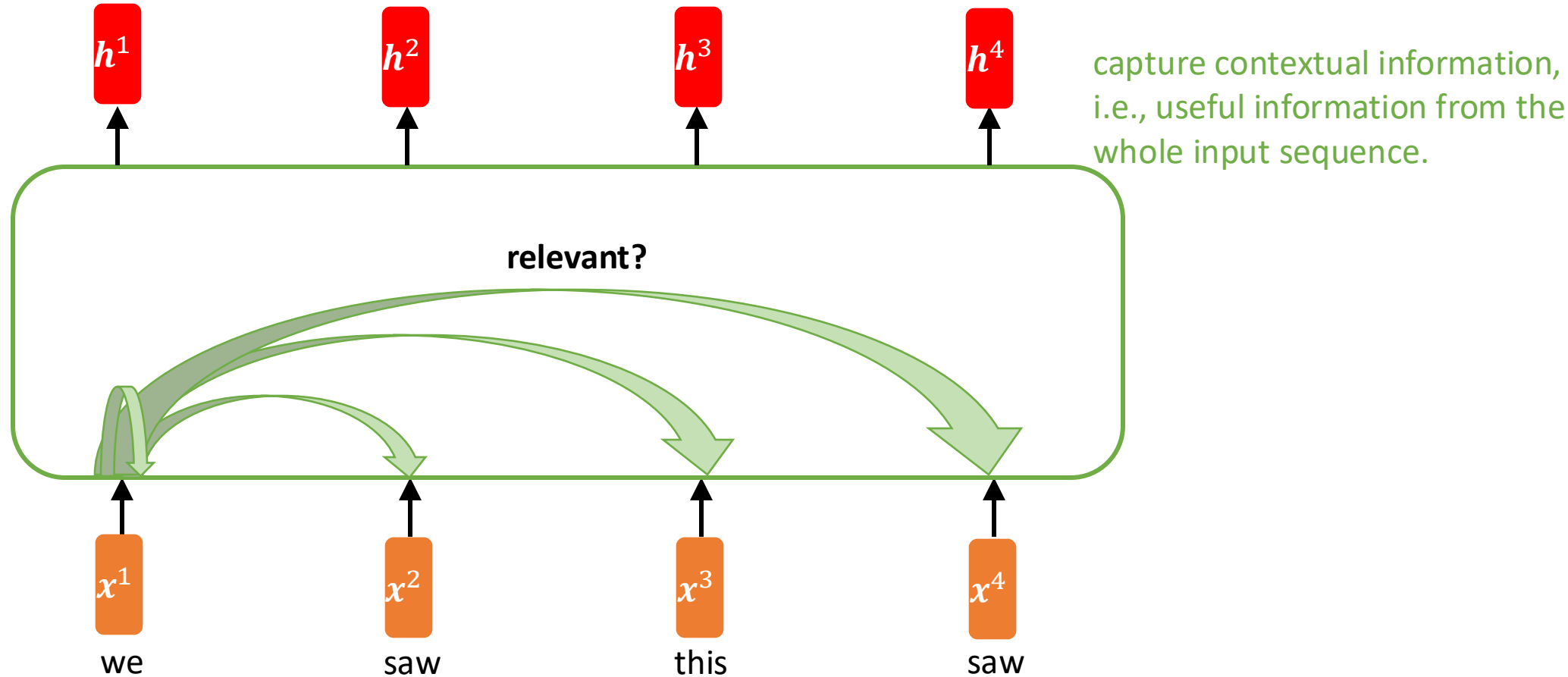
Is it possible to capture contextual information without waiting for previous time steps?

Self-attention layer!

Outline

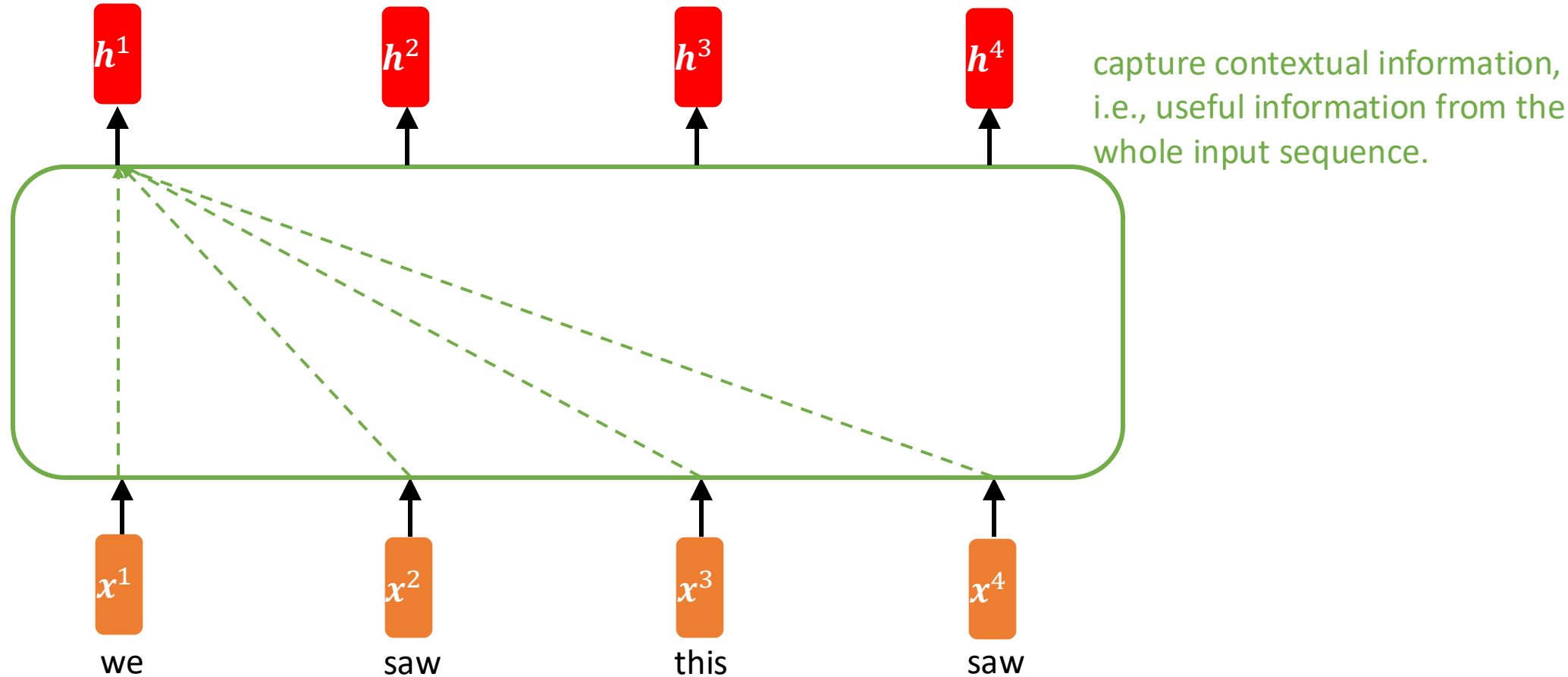
- Recurrent Neural Networks
 - Motivation
 - Recurrent Neural Networks
 - Applications
- **Self-Attention**
 - Self-Attention Layer
 - Positional Encoding
- Issues with Deep Learning
 - Overfitting
 - Vanishing/Exploding Gradient

Self-Attention Layer



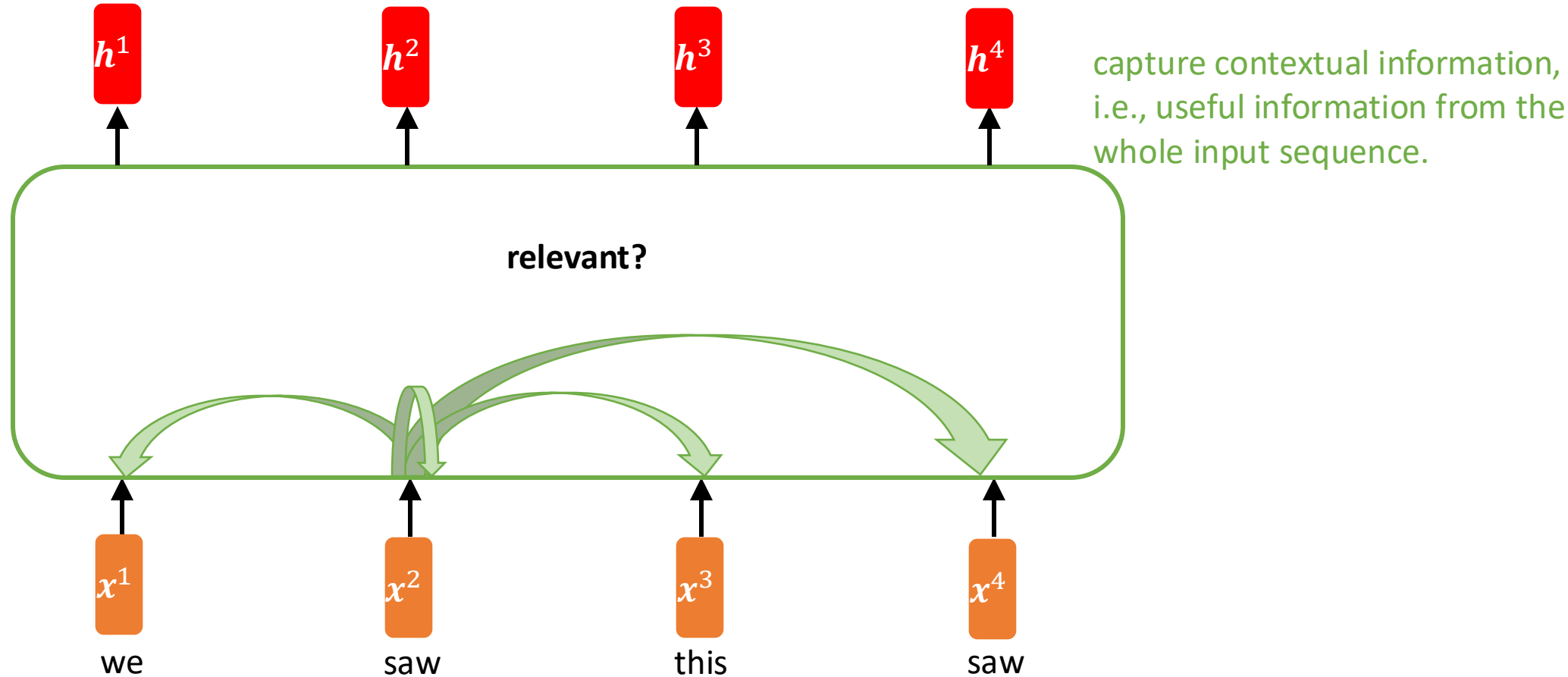
Identify whether the information from elements in the input sequence is relevant for generating the current output.

Self-Attention Layer



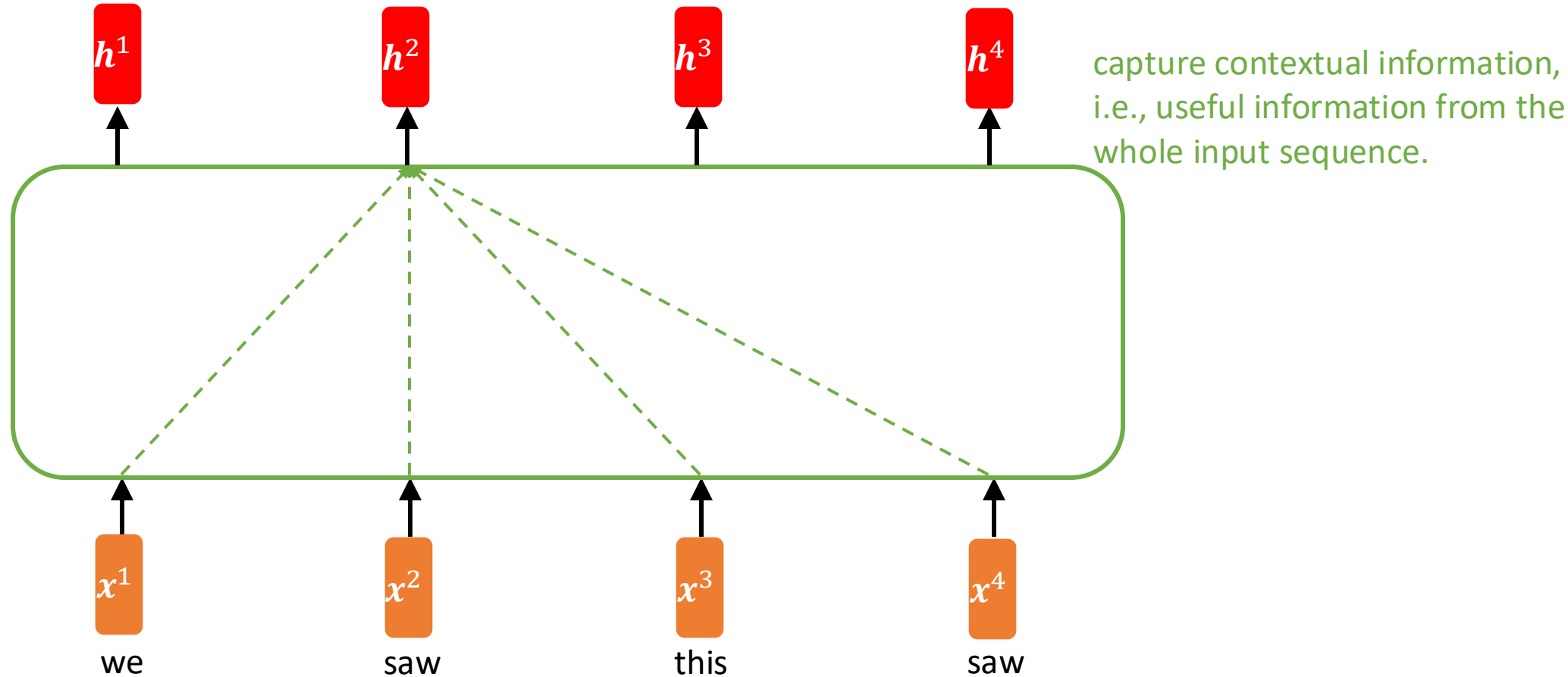
Aggregate useful information from all elements in the sequence to predict the current output.

Self-Attention Layer



Identify whether the information from elements in the input sequence is relevant for generating the current output.

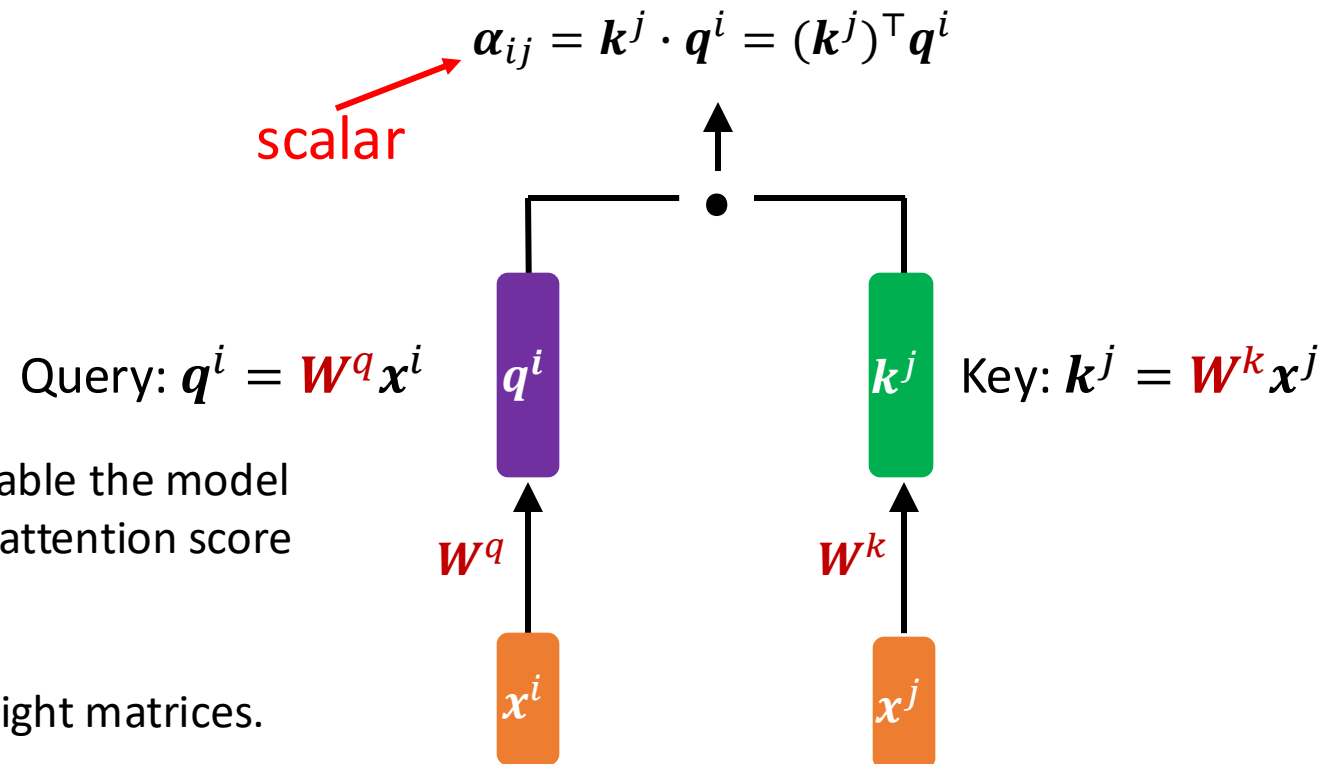
Self-Attention Layer



Aggregate useful information from all elements in the sequence to predict the current output.

Self-Attention Layer

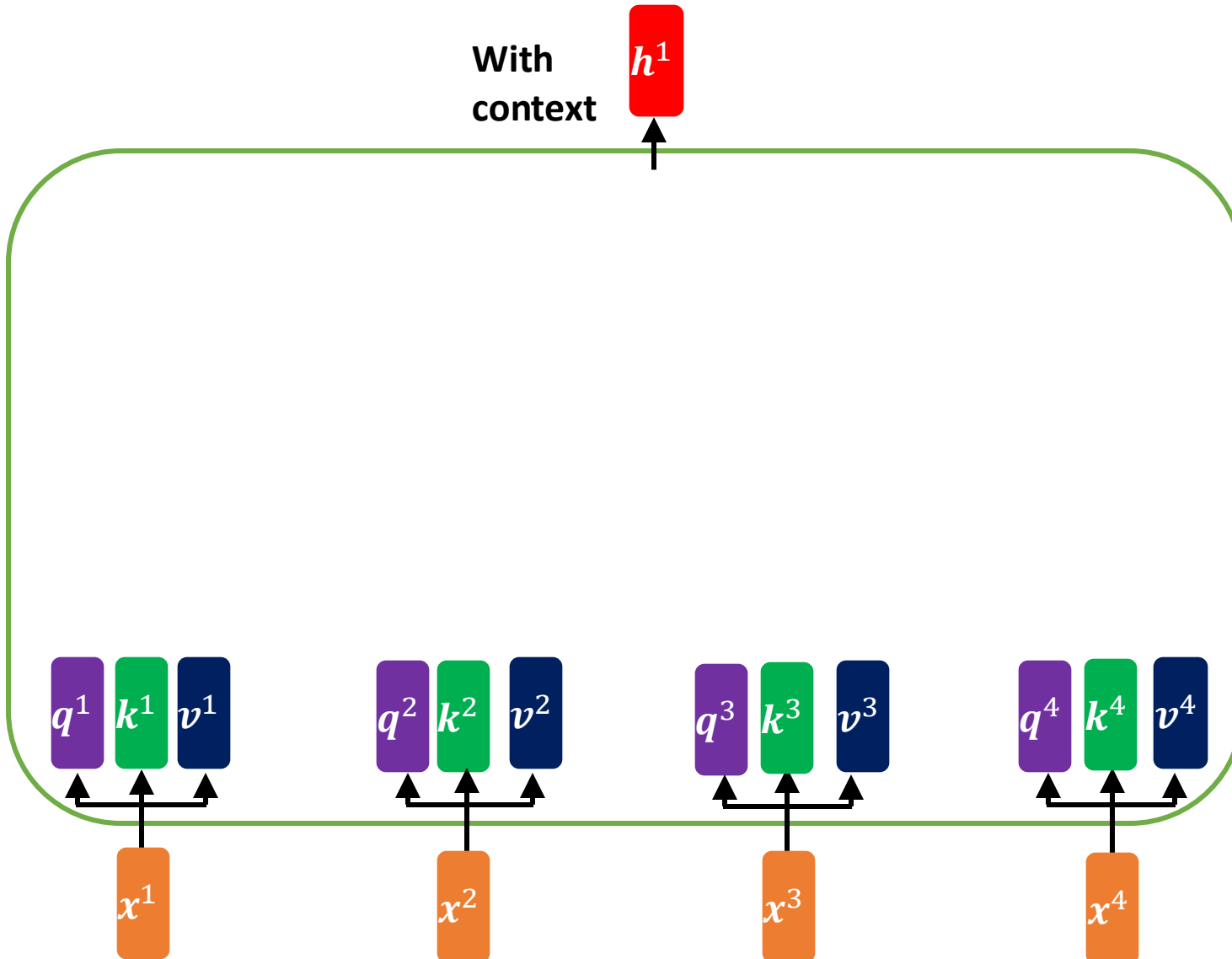
- **Attention score:** Determine how much focus (or "attention") each part of the input sequence (x^j) should receive when processing a specific element (x^i).



The matrices W^q and W^k enable the model to extract useful features for attention score calculation from the inputs

W^q and W^k are trainable weight matrices.

Self-Attention Layer



Step 1: Linear Projection

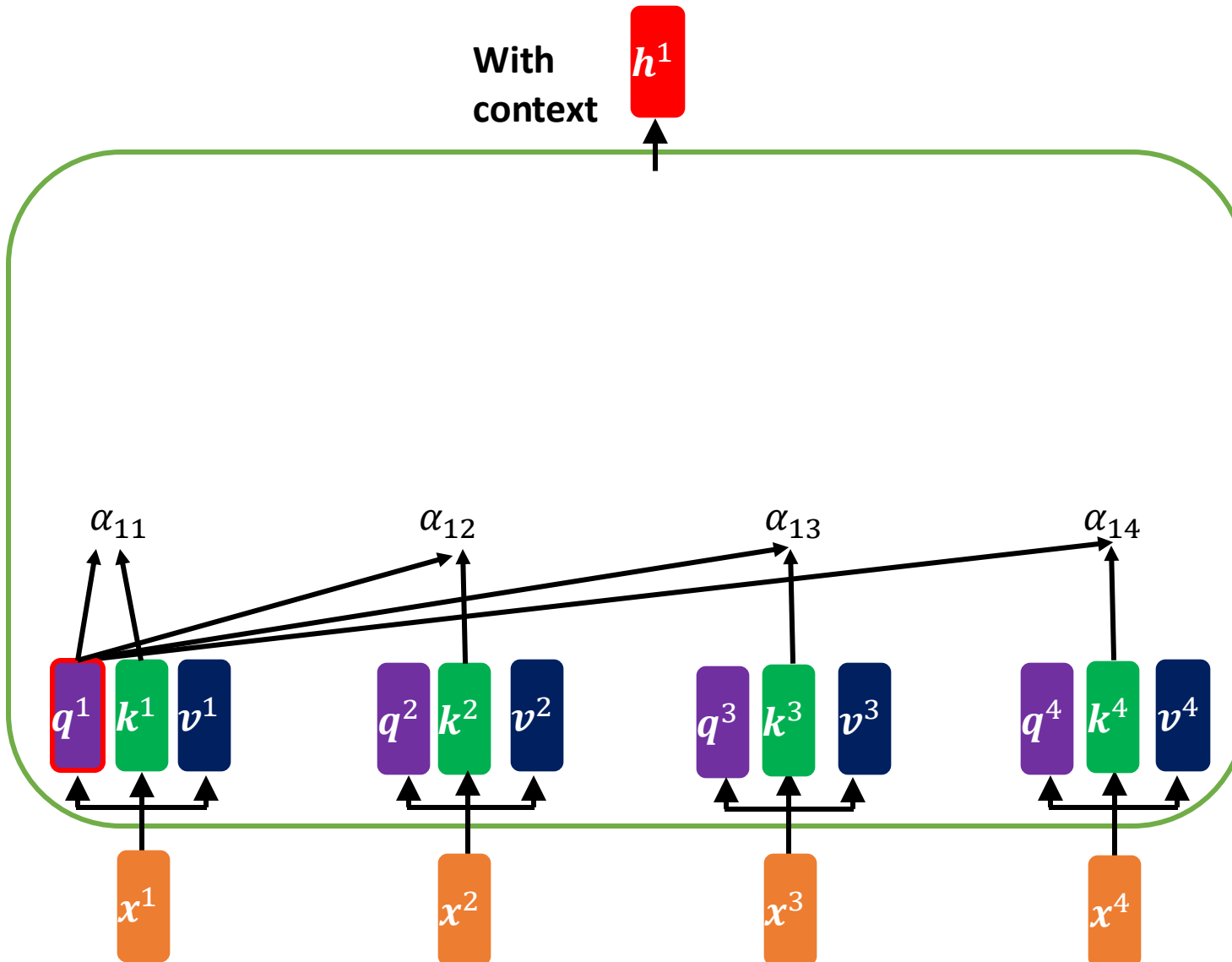
Transform input vectors into Query, Key, and Value using weights matrices (shared across inputs)

Query: $q^i = W^q x^i$

Key: $k^i = W^k x^i$

Value: $v^i = W^v x^i$

Self-Attention Layer



Step 2: Compute the attention scores:

$$\alpha_{1j} = k^j \cdot q^1 = (k^j)^\top q^1$$

Step 1: Linear Projection

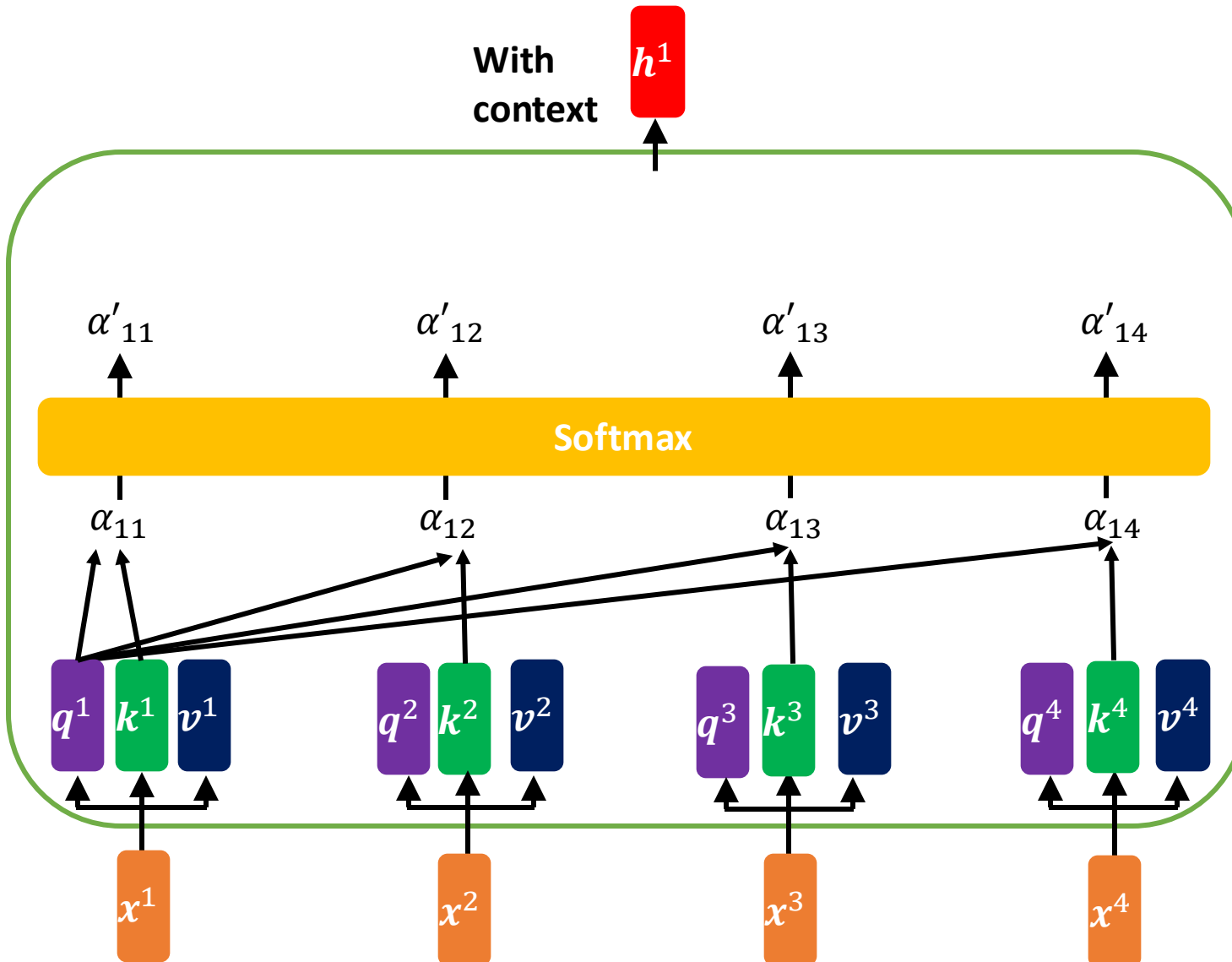
Transform input vectors into Query, Key, and Value using weights matrices (shared across inputs)

Query: $q^i = W^q x^i$

Key: $k^i = W^k x^i$

Value: $v^i = W^v x^i$

Self-Attention Layer



Step 3: Apply Softmax: $\alpha'_{1j} = \frac{e^{\alpha_{1j}}}{\sum_j e^{\alpha_{1j}}}$

Step 2: Compute the attention scores:
 $\alpha_{1j} = \mathbf{k}^j \cdot \mathbf{q}^1 = (\mathbf{k}^j)^\top \mathbf{q}^1$

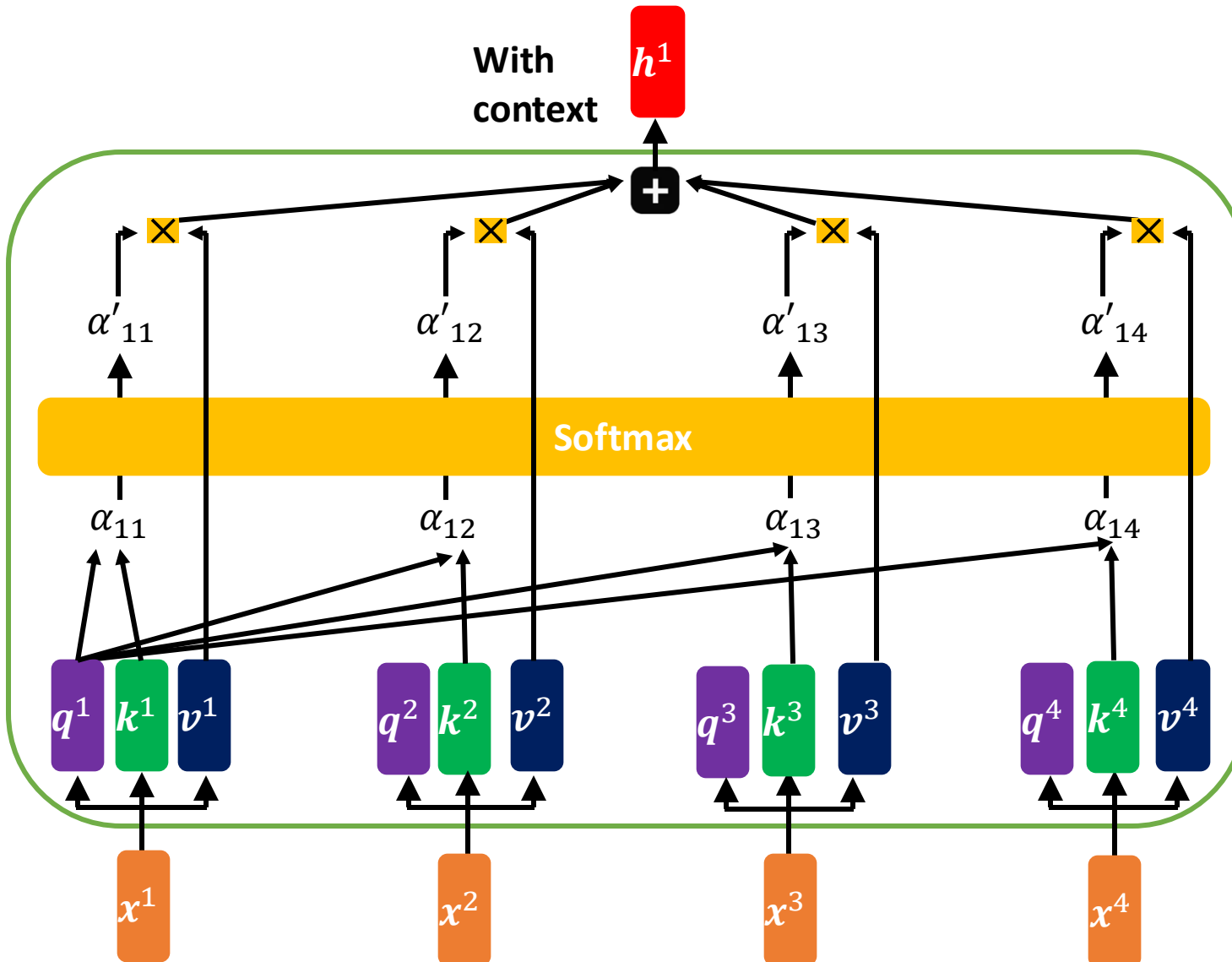
Step 1: Linear Projection
Transform input vectors into Query, Key, and Value using weights matrices (shared across inputs)

Query: $\mathbf{q}^i = \mathbf{W}^q \mathbf{x}^i$

Key: $\mathbf{k}^i = \mathbf{W}^k \mathbf{x}^i$

Value: $\mathbf{v}^i = \mathbf{W}^v \mathbf{x}^i$

Self-Attention Layer



The actual content to aggregate

Step 4: Aggregate information:
Multiply Values by attention
score (after Softmax)

$$h^1 = \sum_j \alpha'_{1j} v^j$$

Step 3: Apply Softmax: $\alpha'_{1j} = \frac{e^{\alpha_{1j}}}{\sum_j e^{\alpha_{1j}}}$

Step 2: Compute the attention scores:
 $\alpha_{1j} = k^j \cdot q^1 = (k^j)^\top q^1$

Step 1: Linear Projection

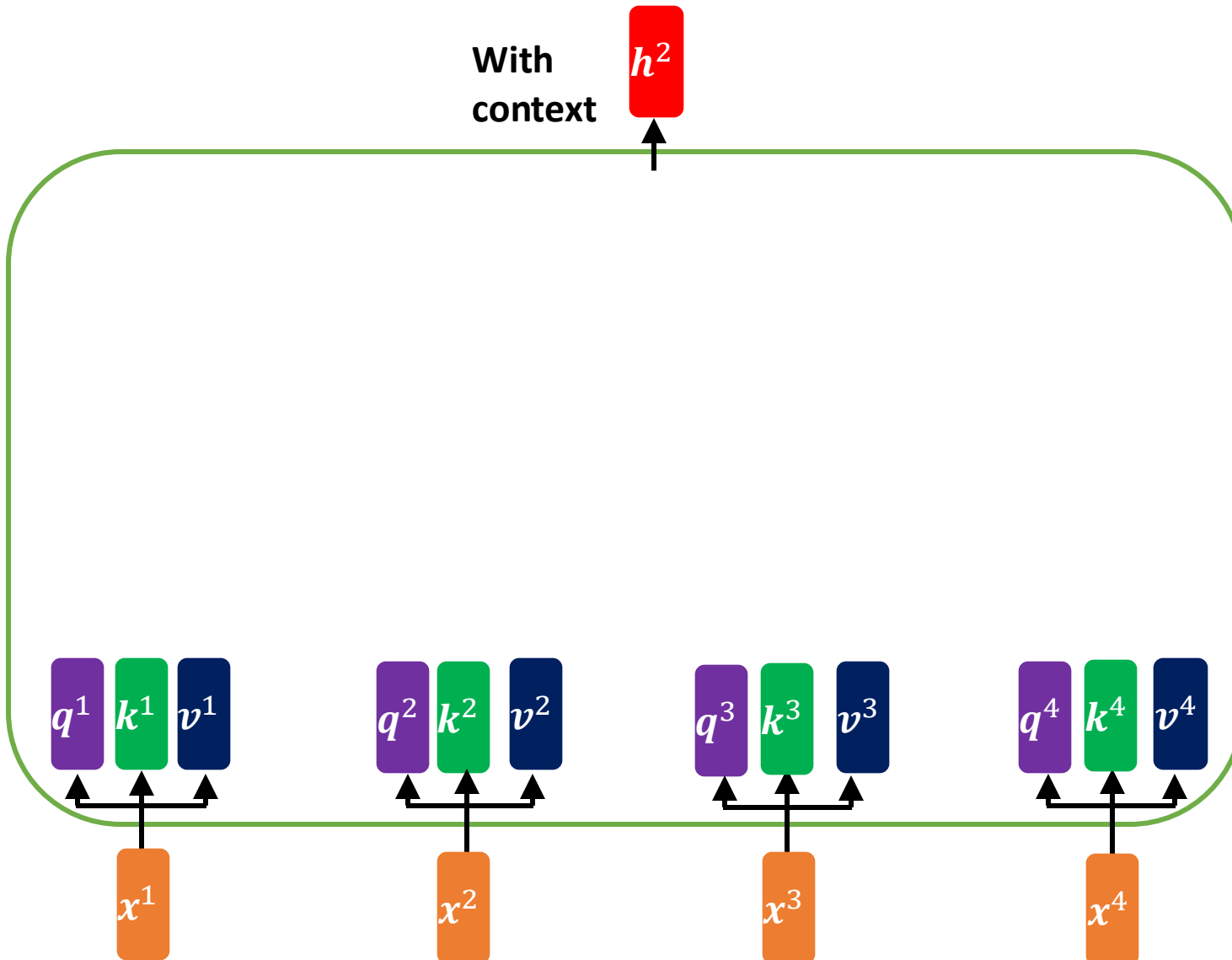
Transform input vectors into Query, Key, and Value
using weights matrices (shared across inputs)

Query: $q^i = W^q x^i$

Key: $k^i = W^k x^i$

Value: $v^i = W^v x^i$

Self-Attention Layer



Step 1: Linear Projection

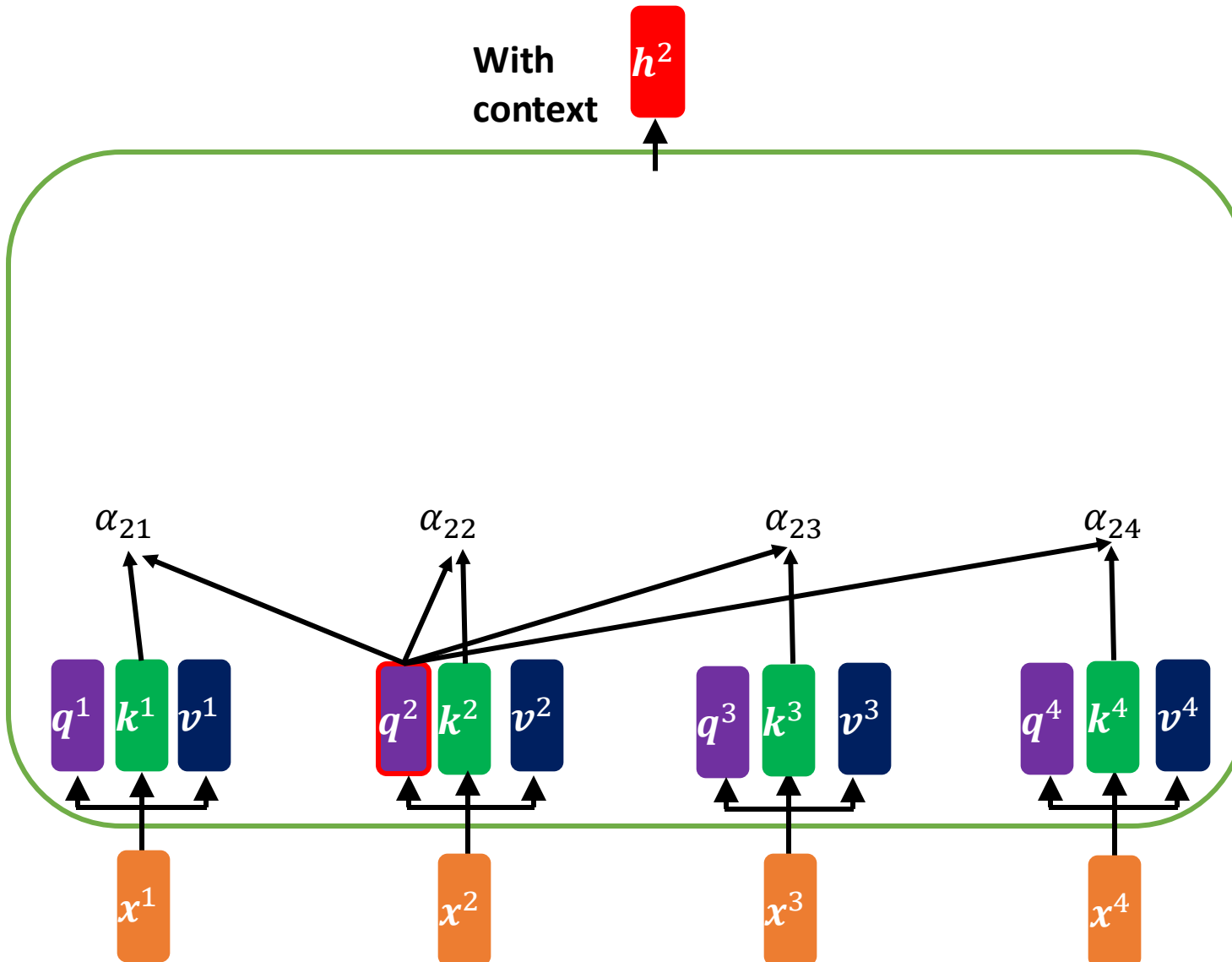
Transform input vectors into Query, Key, and Value using weights matrices (shared across inputs)

Query: $q^i = W^q x^i$

Key: $k^i = W^k x^i$

Value: $v^i = W^v x^i$

Self-Attention Layer



Step 2: Compute the attention scores:

$$\alpha_{2j} = \mathbf{k}^j \cdot \mathbf{q}^2 = (\mathbf{k}^j)^\top \mathbf{q}^2$$

Step 1: Linear Projection

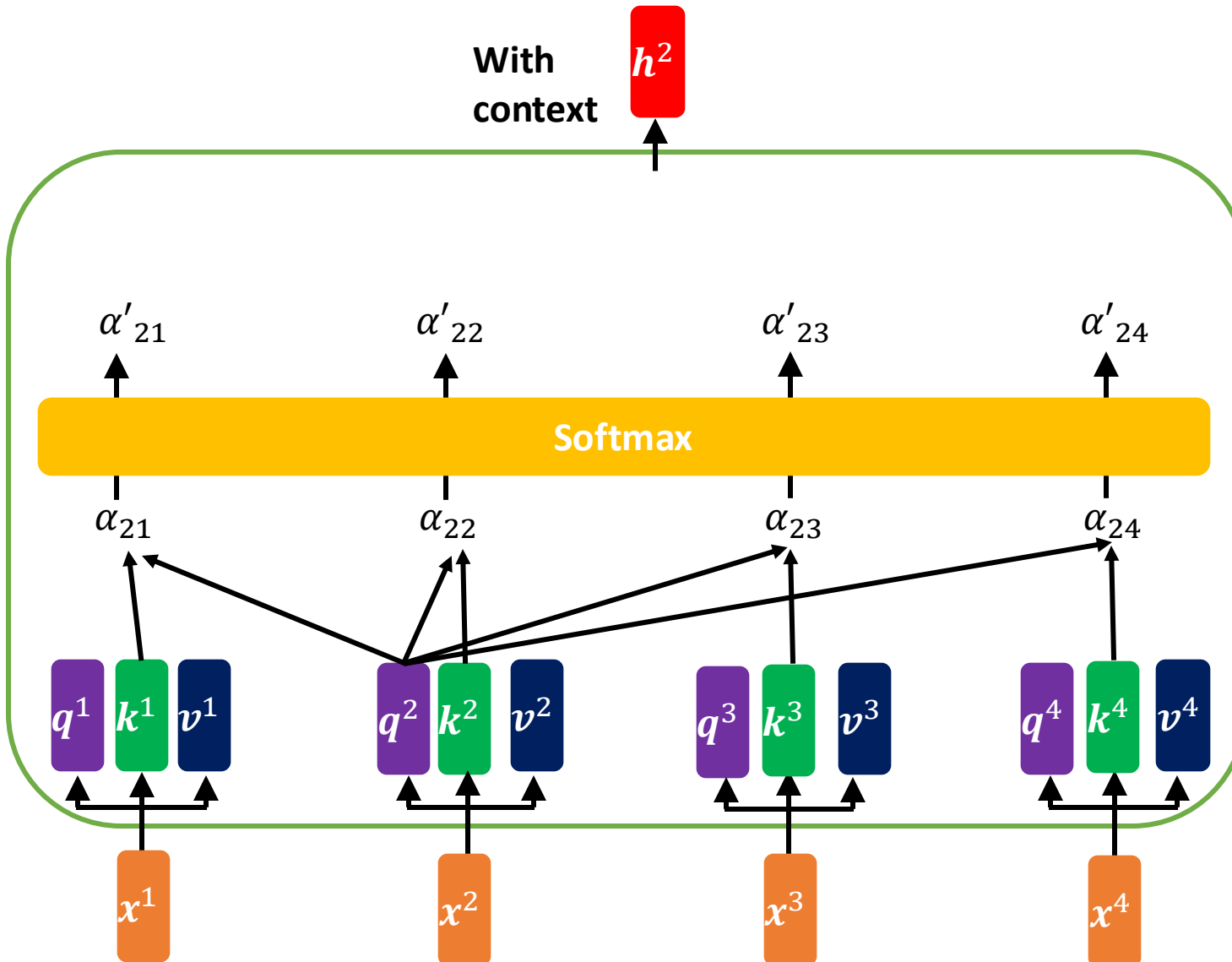
Transform input vectors into Query, Key, and Value using weights matrices (shared across inputs)

Query: $\mathbf{q}^i = \mathbf{W}^q \mathbf{x}^i$

Key: $\mathbf{k}^i = \mathbf{W}^k \mathbf{x}^i$

Value: $\mathbf{v}^i = \mathbf{W}^v \mathbf{x}^i$

Self-Attention Layer



Step 3: Apply Softmax: $\alpha'_{2j} = \frac{e^{\alpha_{2j}}}{\sum_j e^{\alpha_{2j}}}$

Step 2: Compute the attention scores:
 $\alpha_{2j} = \mathbf{k}^j \cdot \mathbf{q}^2 = (\mathbf{k}^j)^\top \mathbf{q}^2$

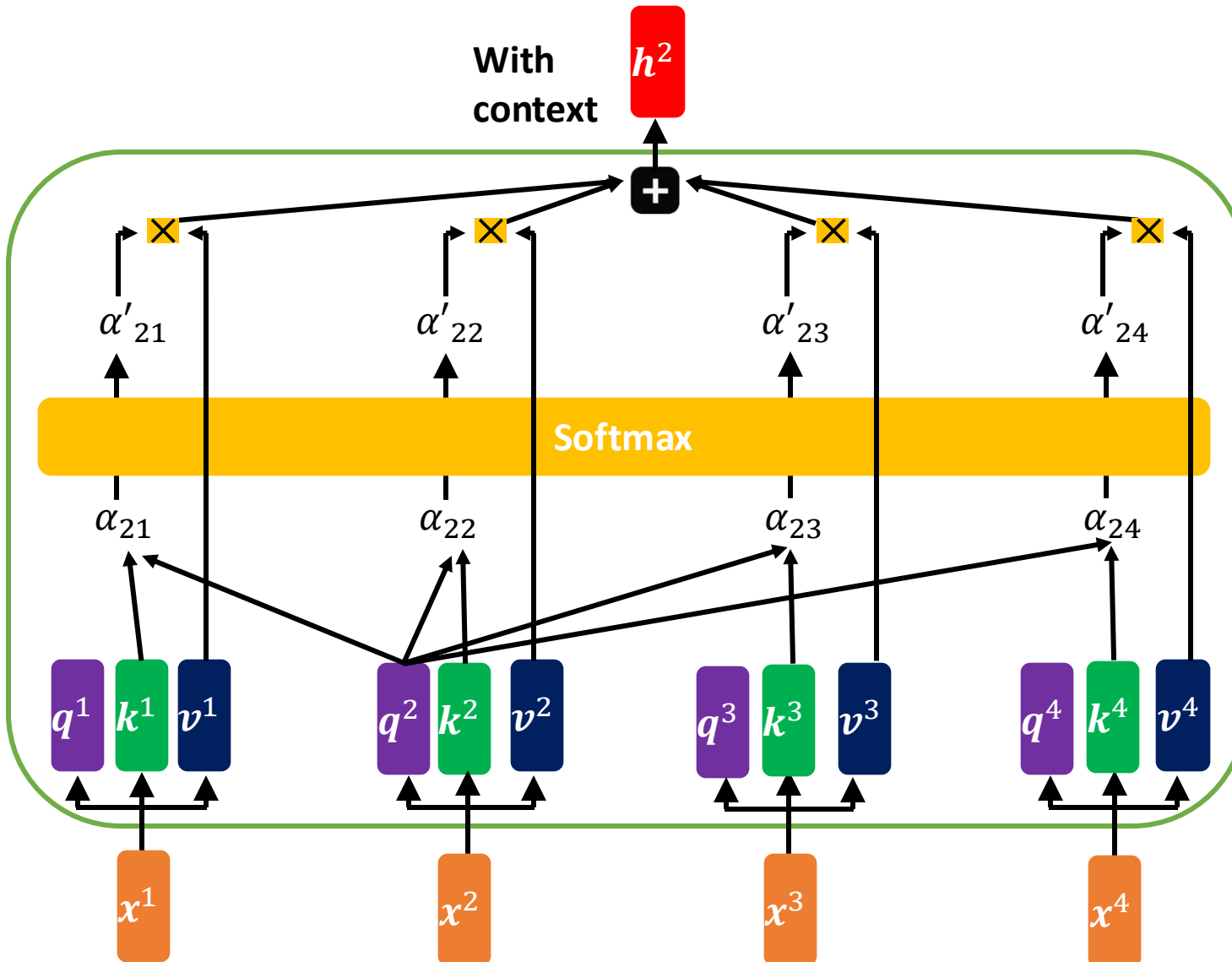
Step 1: Linear Projection
Transform input vectors into Query, Key, and Value using weights matrices (shared across inputs)

Query: $\mathbf{q}^i = \mathbf{W}^q \mathbf{x}^i$

Key: $\mathbf{k}^i = \mathbf{W}^k \mathbf{x}^i$

Value: $\mathbf{v}^i = \mathbf{W}^v \mathbf{x}^i$

Self-Attention Layer



Step 4: Aggregate information:
Multiply Values by attention
score (after Softmax)

$$h^2 = \sum_j \alpha'_{2j} v^j$$

Step 3: Apply Softmax: $\alpha'_{2j} = \frac{e^{\alpha_{2j}}}{\sum_j e^{\alpha_{2j}}}$

Step 2: Compute the attention scores:
 $\alpha_{2j} = k^j \cdot q^2 = (k^j)^\top q^2$

Step 1: Linear Projection

Transform input vectors into Query, Key, and Value
using weights matrices (shared across inputs)

Query: $q^i = W^q x^i$

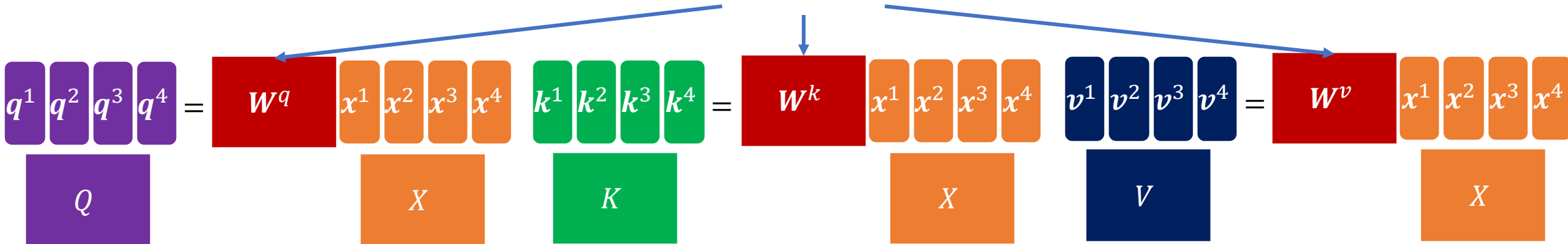
Key: $k^i = W^k x^i$

Value: $v^i = W^v x^i$

Self-Attention Layer

Step1: Linear Projection

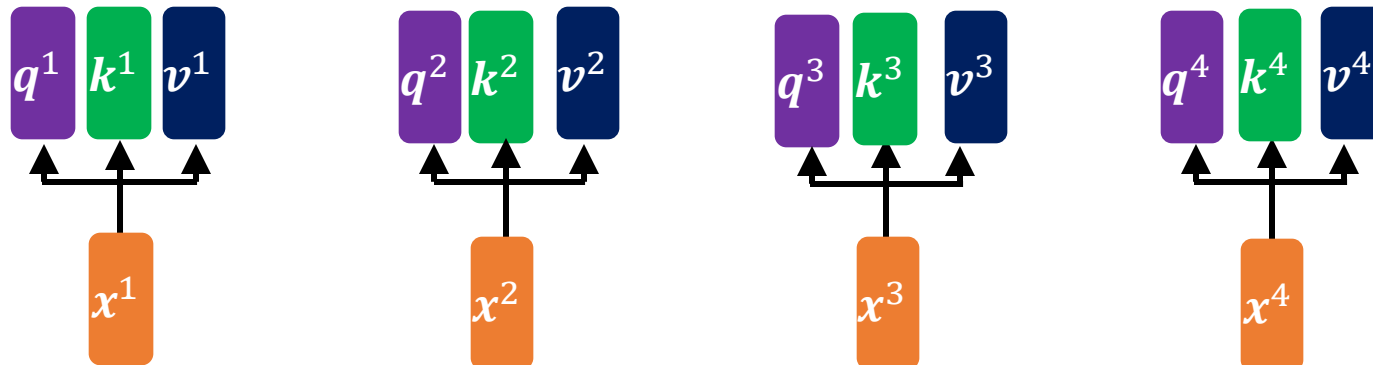
Trainable Weights/Parameters



Query $q^i = W^q x^i$

Key $k^i = W^k x^i$

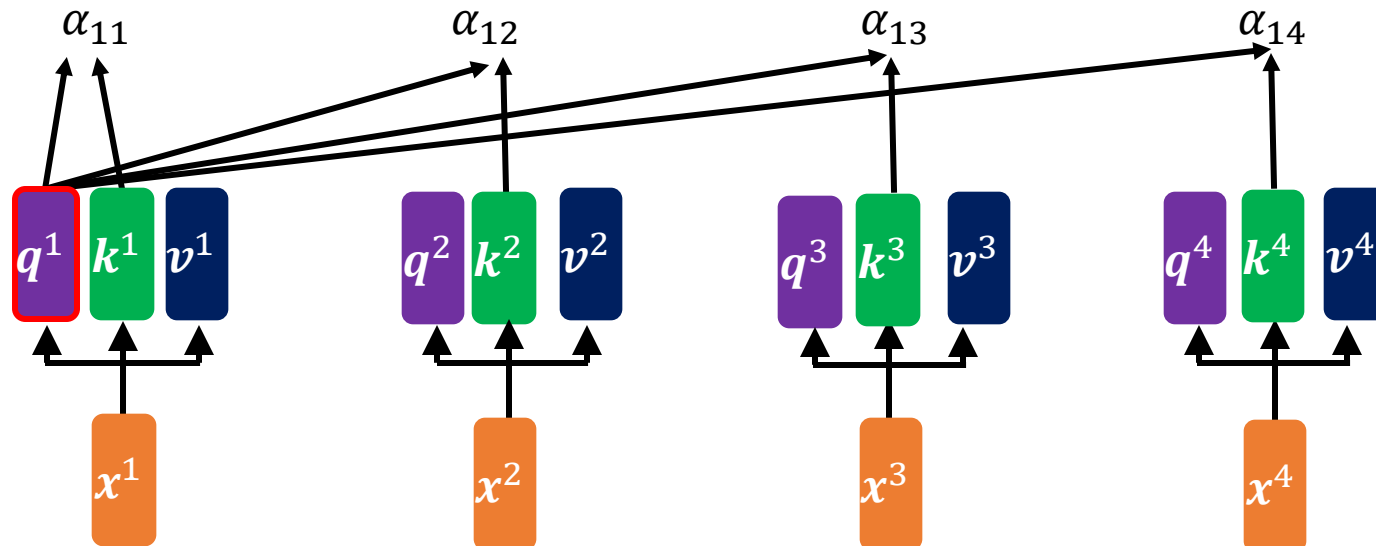
Value $v^i = W^v x^i$



Self-Attention Layer

Step2: Compute the attention scores

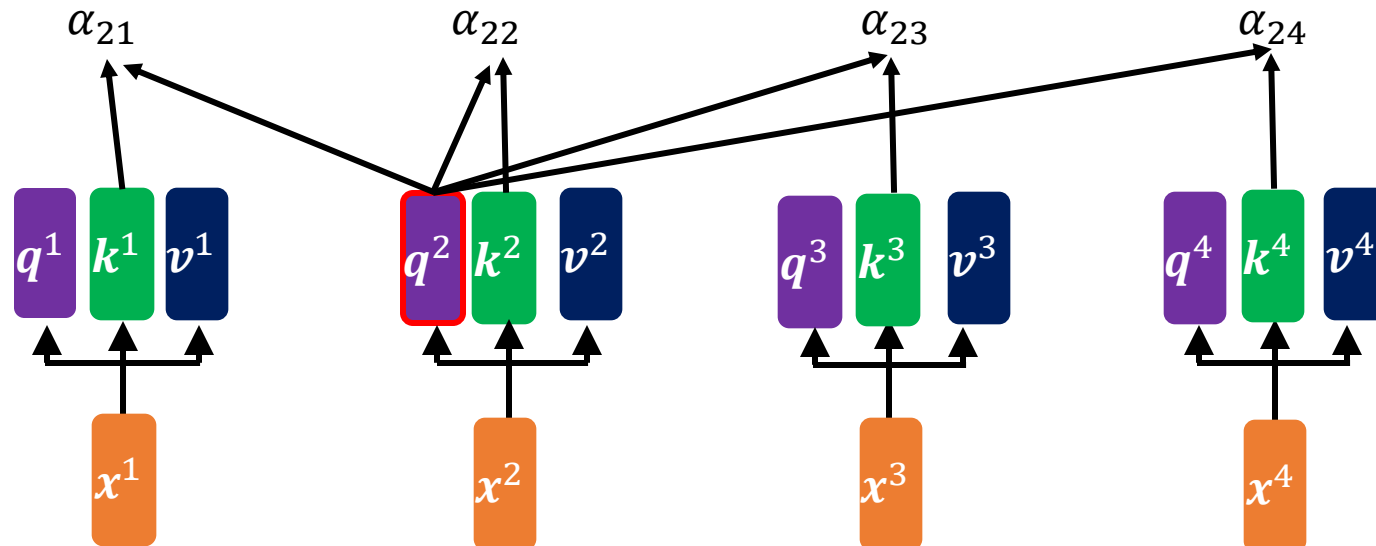
$$\begin{aligned}\alpha_{11} &= (k^1)^\top q^1 & \alpha_{12} &= (k^2)^\top q^1 \\ \alpha_{13} &= (k^3)^\top q^1 & \alpha_{14} &= (k^4)^\top q^1\end{aligned}$$
$$\begin{matrix} \alpha_{11} \\ \alpha_{12} \\ \alpha_{13} \\ \alpha_{14} \end{matrix} = \begin{matrix} (k^1)^\top \\ (k^2)^\top \\ (k^3)^\top \\ (k^4)^\top \end{matrix} q^1$$



Self-Attention Layer

Step2: Compute the attention scores

$$\begin{aligned}\alpha_{21} &= (k^1)^\top q^2 & \alpha_{22} &= (k^2)^\top q^2 \\ \alpha_{23} &= (k^3)^\top q^2 & \alpha_{24} &= (k^4)^\top q^2\end{aligned}$$
$$\begin{matrix} \alpha_{21} \\ \alpha_{22} \\ \alpha_{23} \\ \alpha_{24} \end{matrix} = \begin{matrix} (k^1)^\top \\ (k^2)^\top \\ (k^3)^\top \\ (k^4)^\top \end{matrix} q^2$$

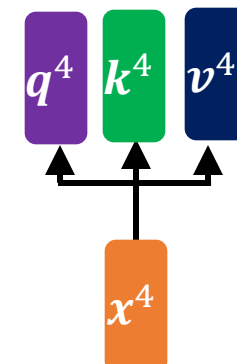
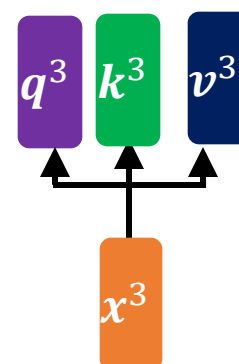
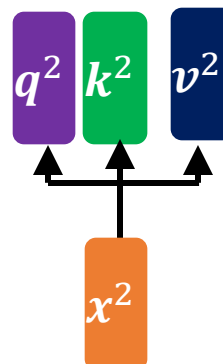
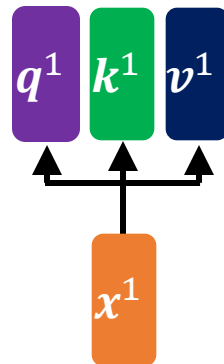


Self-Attention Layer

Step2: Compute the attention scores

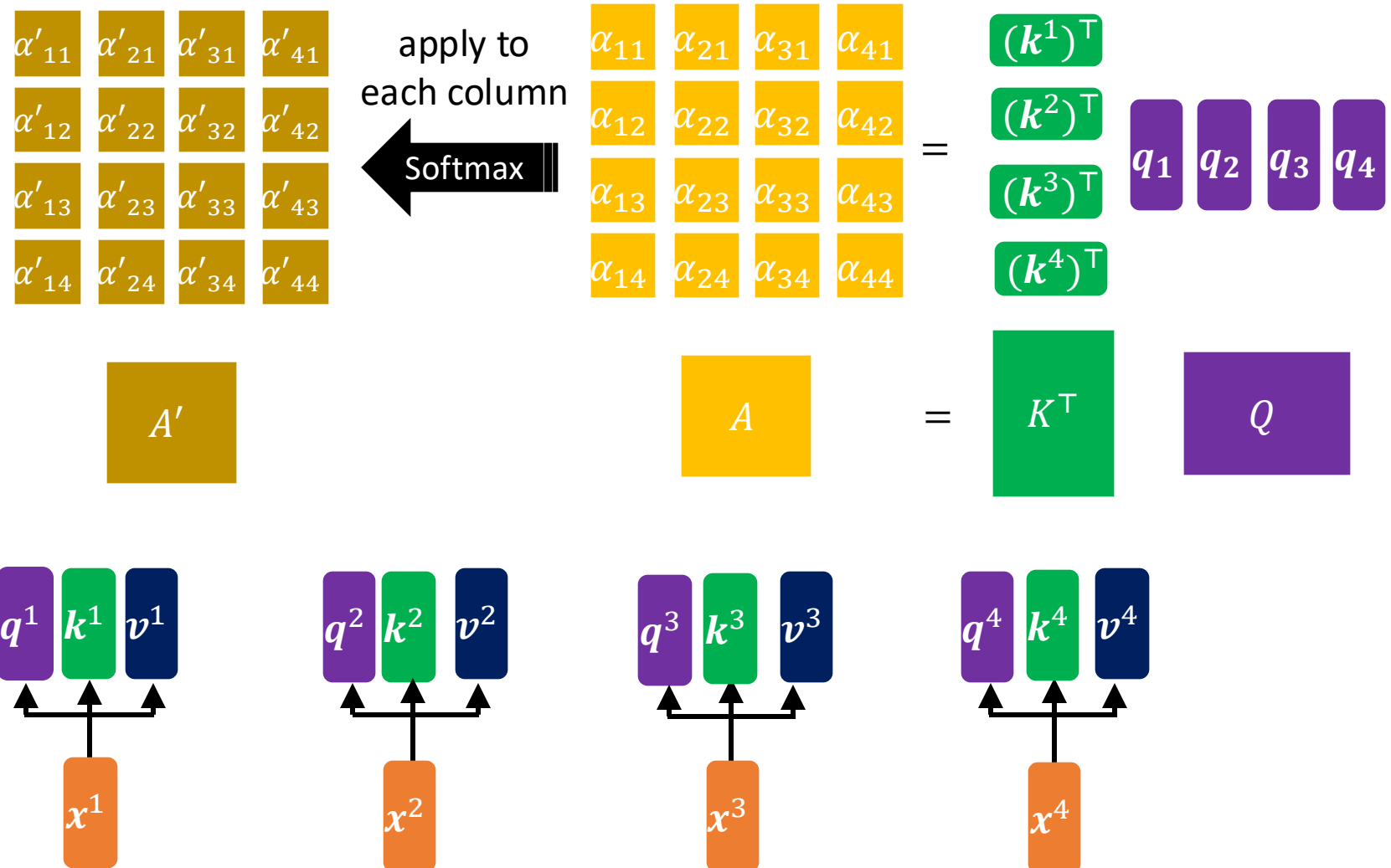
$$\begin{array}{cccc} \alpha_{11} & \alpha_{21} & \alpha_{31} & \alpha_{41} \\ \alpha_{12} & \alpha_{22} & \alpha_{32} & \alpha_{42} \\ \alpha_{13} & \alpha_{23} & \alpha_{33} & \alpha_{43} \\ \alpha_{14} & \alpha_{24} & \alpha_{34} & \alpha_{44} \end{array} = \begin{array}{c} (k^1)^\top \\ (k^2)^\top \\ (k^3)^\top \\ (k^4)^\top \end{array} \begin{array}{cccc} q_1 & q_2 & q_3 & q_4 \end{array}$$

$$A = K^\top Q$$



Self-Attention Layer

Step3: Apply Softmax



Background: Weighted Sum of Vectors

- Given vectors $\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^N$ and the corresponding weights for each vector a^1, a^2, \dots, a^N , we can compute the weighted sum of vectors:

$$\mathbf{s} = \sum_{j=1}^N a^j \mathbf{v}^j = \begin{bmatrix} | & | & \dots & | \\ \mathbf{v}^1 & \mathbf{v}^2 & \dots & \mathbf{v}^N \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} a^1 \\ a^2 \\ \vdots \\ a^N \end{bmatrix}$$

- For example:

Stack the vectors as columns in a matrix Arrange the weights as a column vector

$$\mathbf{v}^1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \mathbf{v}^2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \mathbf{v}^3 = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \text{and} \quad a^1 = 0.5, a^2 = 0.2, a^3 = 0.3$$

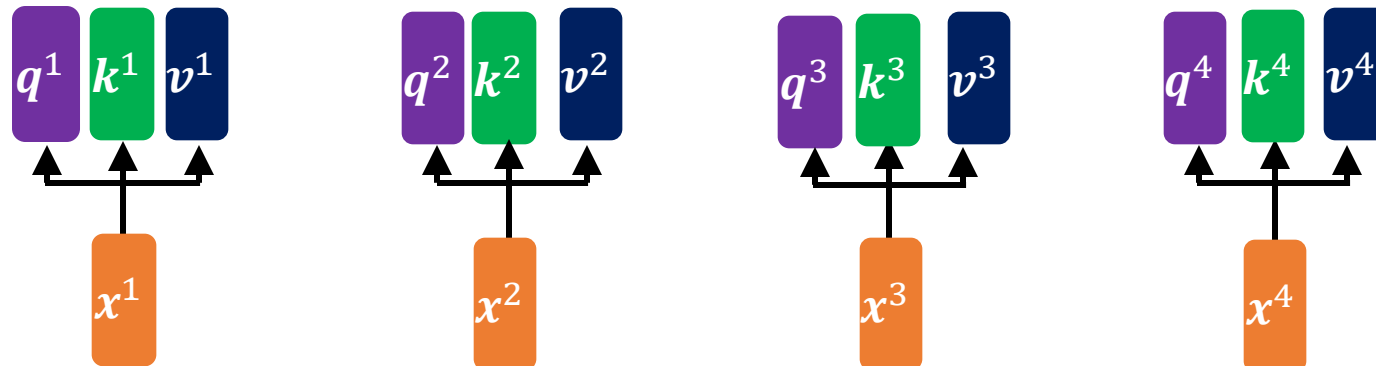
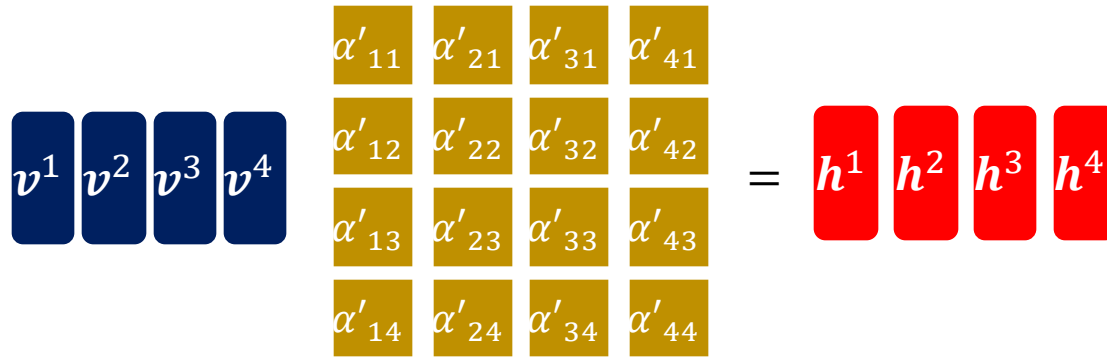
$$\mathbf{s} = 0.5 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 0.2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 0.3 \begin{bmatrix} 0 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \times 0.5 \\ 2 \times 0.5 \end{bmatrix} + \begin{bmatrix} 1 \times 0.2 \\ 0 \times 0.2 \end{bmatrix} + \begin{bmatrix} 0 \times 0.3 \\ 2 \times 0.3 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 1.6 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 & 0 \\ 2 & 0 & 2 \end{bmatrix} \begin{bmatrix} 0.5 \\ 0.2 \\ 0.3 \end{bmatrix} = \begin{bmatrix} 1 \times 0.5 + 1 \times 0.2 + 0 \times 0.3 \\ 2 \times 0.5 + 0 \times 0.2 + 2 \times 0.3 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 1.6 \end{bmatrix}$$

Self-Attention Layer

Step4: Aggregate Information: Multiply Values by attention score (after Softmax)

$$h^i = \sum_j \alpha'_{ij} v^j$$

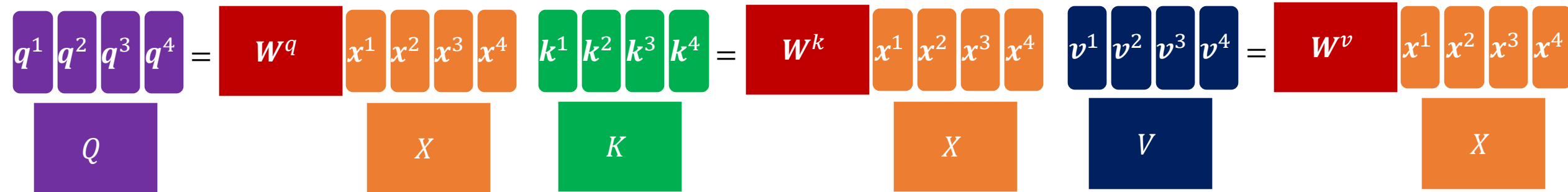


Self-Attention Layer

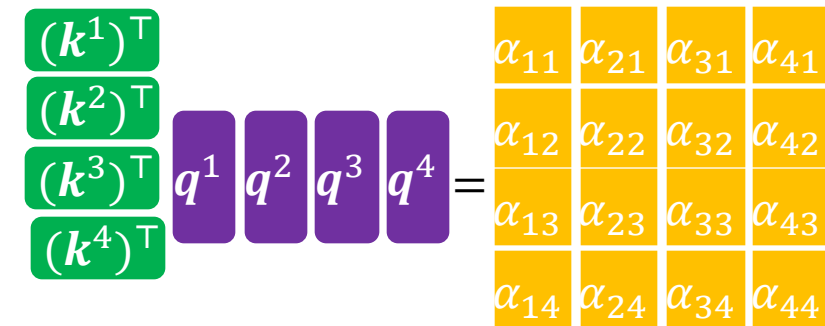
h_1, h_2, h_3, h_4

Can be generated together,
not need to wait for previous
time steps to complete first

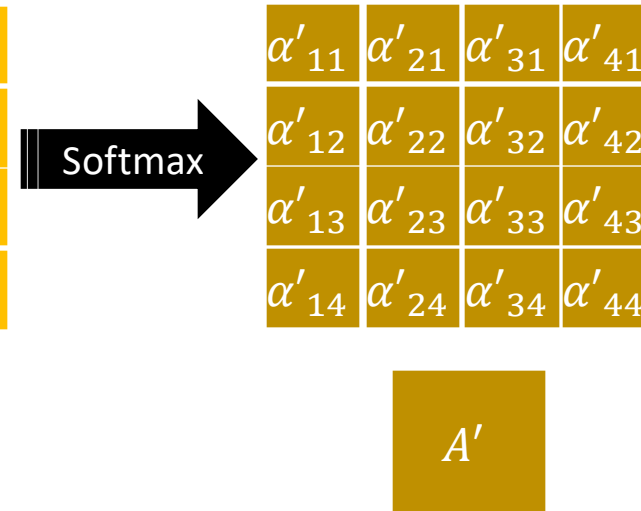
Step1: Linear Projection



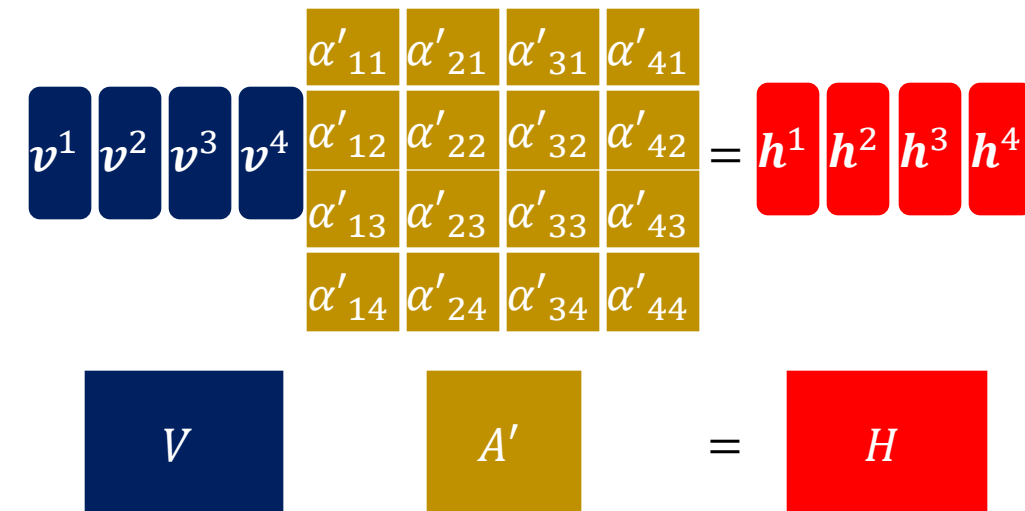
Step2: Compute the attention scores



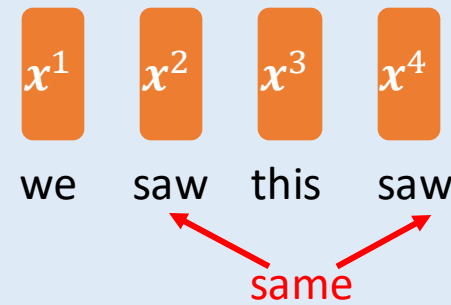
Step3: Apply Softmax



Step4: : Aggregate Information



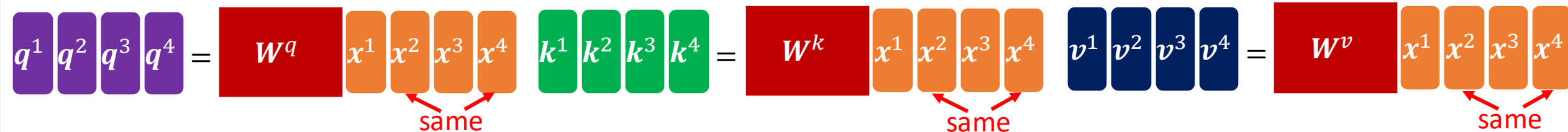
Poll Everywhere



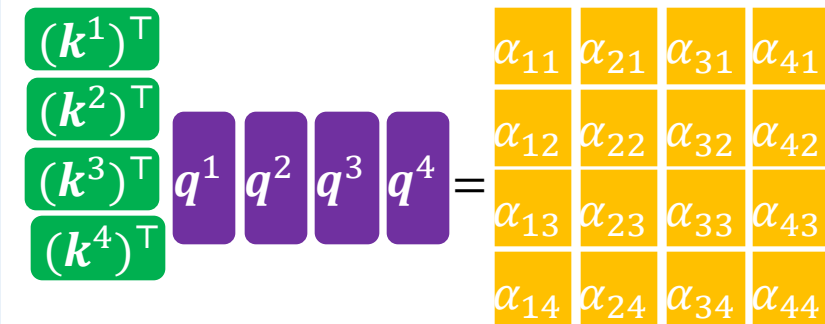
PollEv.com/conghuihu365

- Consider the sentence "we saw this saw". If we use one-hot encoding to represent the words and input these one-hot vectors into a self-attention layer, can the self-attention layer generate **different** outputs for the two instances of "saw"?

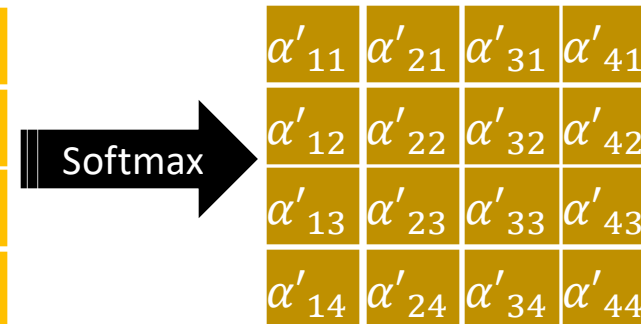
Step1: Linear Projection



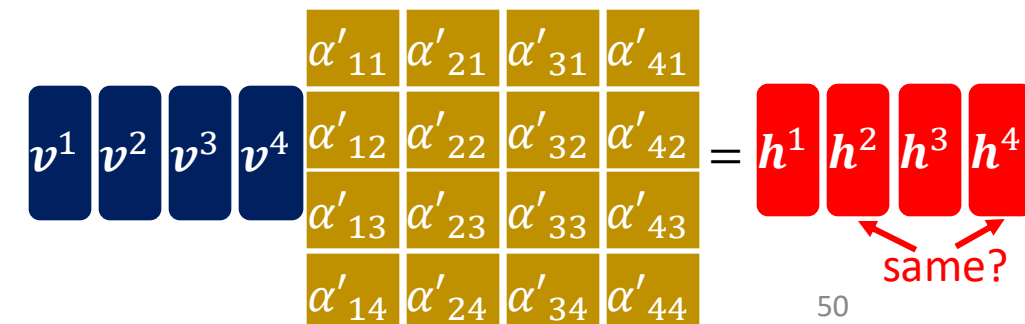
Step2: Compute the attention scores



Step3: Apply Softmax



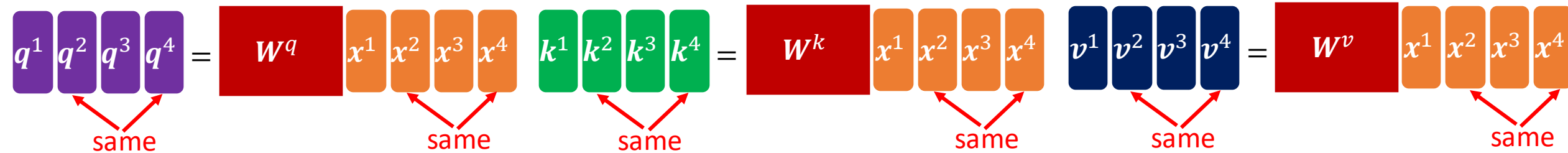
Step4: : Aggregate Information



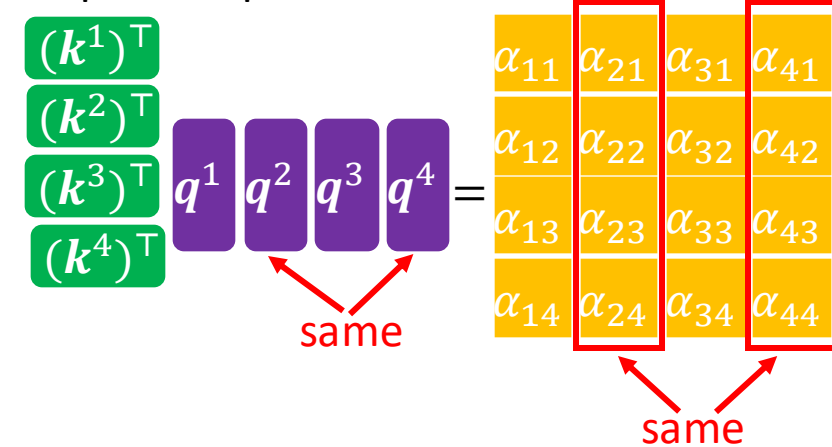
Self-Attention Layer: Example

x^1 x^2 x^3 x^4
 we saw this saw
 ↙ ↘
 same

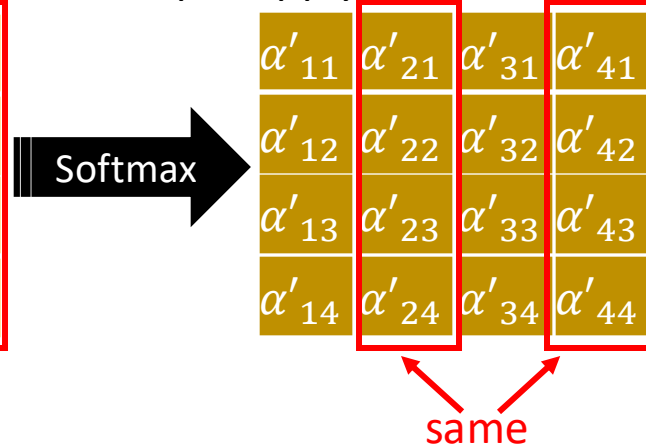
Step1: Linear Projection



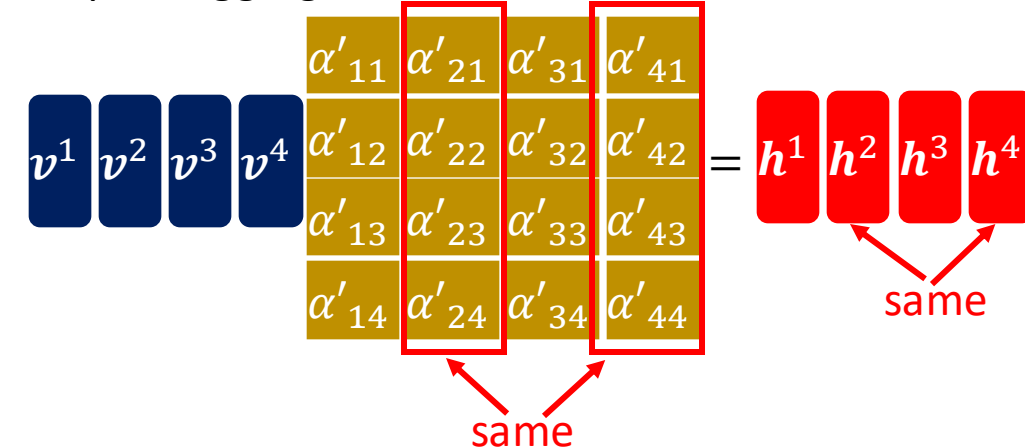
Step2: Compute the attention scores



Step3: Apply Softmax



Step4: : Aggregation Information



Cannot distinguish between the first “saw” (verb) and the second “saw” (noun)!

Need to add positional information to differentiate the two “saw”!

Positional Encoding

- Positional encoding explicitly injects positional information into the original input features (e.g., one-hot vectors):

$$\mathbf{x}'^i = \mathbf{x}^i + \mathbf{PE}^i$$

↑
Unique for each position!

\mathbf{x}^i : Original input feature vector for position i , where $\mathbf{x}^i \in \mathbb{R}^d$.

\mathbf{PE}^i : Positional encoding vector for position i , where $\mathbf{PE}^i \in \mathbb{R}^d$

\mathbf{x}'^i : Final input feature vector after adding positional encoding.

Positional Encoding: \sin and \cos Functions

- \sin and \cos functions can be used to generate the positional encoding as follows:

$$PE(i, 2k) = \sin\left(\frac{i}{10000^{\frac{2k}{d}}}\right), \quad PE(i, 2k + 1) = \cos\left(\frac{i}{10000^{\frac{2k}{d}}}\right)$$

i : The position index.

$2k$ and $2k + 1$: The dimension indices of the positional encoding vector.

k starts from 0 and goes up to $\frac{d}{2} - 1$, where d is the input feature dimension.

Positional Encoding: sin and cos Functions

- Suppose the input feature dimension is $d = 4$,
- PE^1 can be generate as follows:

- $i = 1, k = 0$
- $PE(1, 0) = \sin\left(\frac{1}{10000 \frac{2 \times 0}{4}}\right) = \sin\left(\frac{1}{10000 \frac{0}{4}}\right) = \sin(1).$
- $PE(1, 1) = \cos\left(\frac{1}{10000 \frac{2 \times 0}{4}}\right) = \cos\left(\frac{1}{10000 \frac{0}{4}}\right) = \cos(1)$
- $i = 1, k = 1$
- $PE(1, 2) = \sin\left(\frac{1}{10000 \frac{2 \times 1}{4}}\right) = \sin\left(\frac{1}{10000 \frac{1}{2}}\right) = \sin\left(\frac{1}{100}\right).$
- $PE(1, 3) = \cos\left(\frac{1}{10000 \frac{2 \times 1}{4}}\right) = \cos\left(\frac{1}{10000 \frac{1}{2}}\right) = \cos\left(\frac{1}{100}\right)$

$$PE(i, 2k) = \sin\left(\frac{i}{10000 \frac{2k}{d}}\right)$$

$$PE(i, 2k + 1) = \cos\left(\frac{i}{10000 \frac{2k}{d}}\right)$$

$$0 \leq k \leq \frac{d}{2} - 1$$

$$0 \leq k \leq \frac{4}{2} - 1 = 1$$

$$PE^1 = \begin{bmatrix} \sin(1) \\ \cos(1) \\ \sin\left(\frac{1}{100}\right) \\ \cos\left(\frac{1}{100}\right) \end{bmatrix}$$

Positional Encoding: sin and cos Functions

- Suppose the input feature dimension is $d = 4$,
- PE^2 can be generate as follows:

- $i = 2, k = 0$
- $PE(2, 0) = \sin\left(\frac{2}{10000 \frac{2 \times 0}{4}}\right) = \sin\left(\frac{2}{10000 \frac{0}{4}}\right) = \sin(2).$
- $PE(2, 1) = \cos\left(\frac{2}{10000 \frac{2 \times 0}{4}}\right) = \cos\left(\frac{2}{10000 \frac{0}{4}}\right) = \cos(2)$
- $i = 2, k = 1$
- $PE(2, 2) = \sin\left(\frac{2}{10000 \frac{2 \times 1}{4}}\right) = \sin\left(\frac{2}{10000 \frac{1}{2}}\right) = \sin\left(\frac{2}{100}\right) = \sin\left(\frac{1}{50}\right).$
- $PE(2, 3) = \cos\left(\frac{2}{10000 \frac{2 \times 1}{4}}\right) = \cos\left(\frac{2}{10000 \frac{1}{2}}\right) = \cos\left(\frac{2}{100}\right) = \cos\left(\frac{1}{50}\right)$

$$PE(i, 2k) = \sin\left(\frac{i}{10000 \frac{2k}{d}}\right)$$

$$PE(i, 2k + 1) = \cos\left(\frac{i}{10000 \frac{2k}{d}}\right)$$

$$0 \leq k \leq \frac{d}{2} - 1$$

$$0 \leq k \leq \frac{4}{2} - 1 = 1$$

$$PE^2 = \begin{bmatrix} \sin(2) \\ \cos(2) \\ \sin\left(\frac{1}{50}\right) \\ \cos\left(\frac{1}{50}\right) \end{bmatrix}$$

Outline

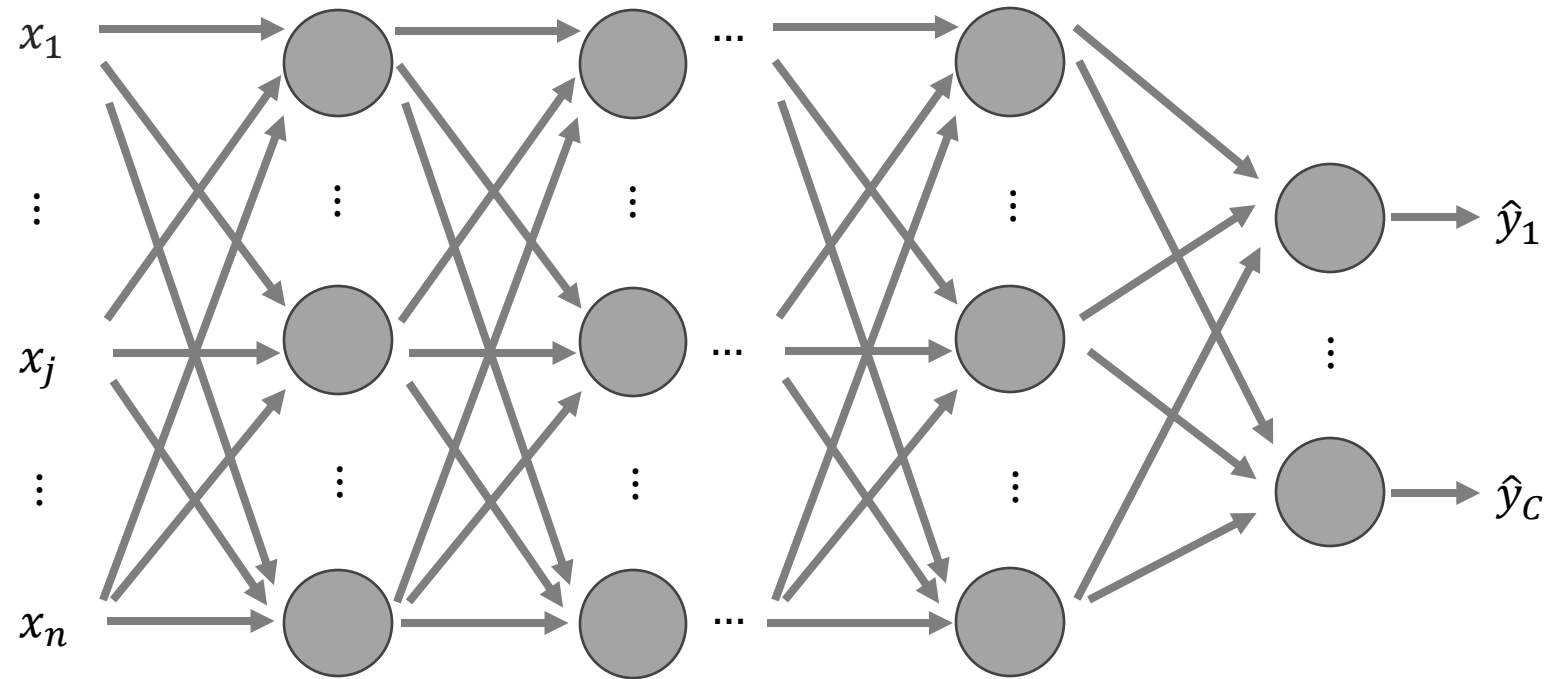
- Recurrent Neural Networks
 - Motivation
 - Recurrent Neural Networks
 - Applications
- Self-Attention
 - Self-Attention Layer
 - Positional Encoding
- **Issues with Deep Learning**
 - Overfitting
 - Vanishing/Exploding Gradient

Issues with Deep Learning

- Overfitting
- Gradient Vanishing/Exploding

Dropout

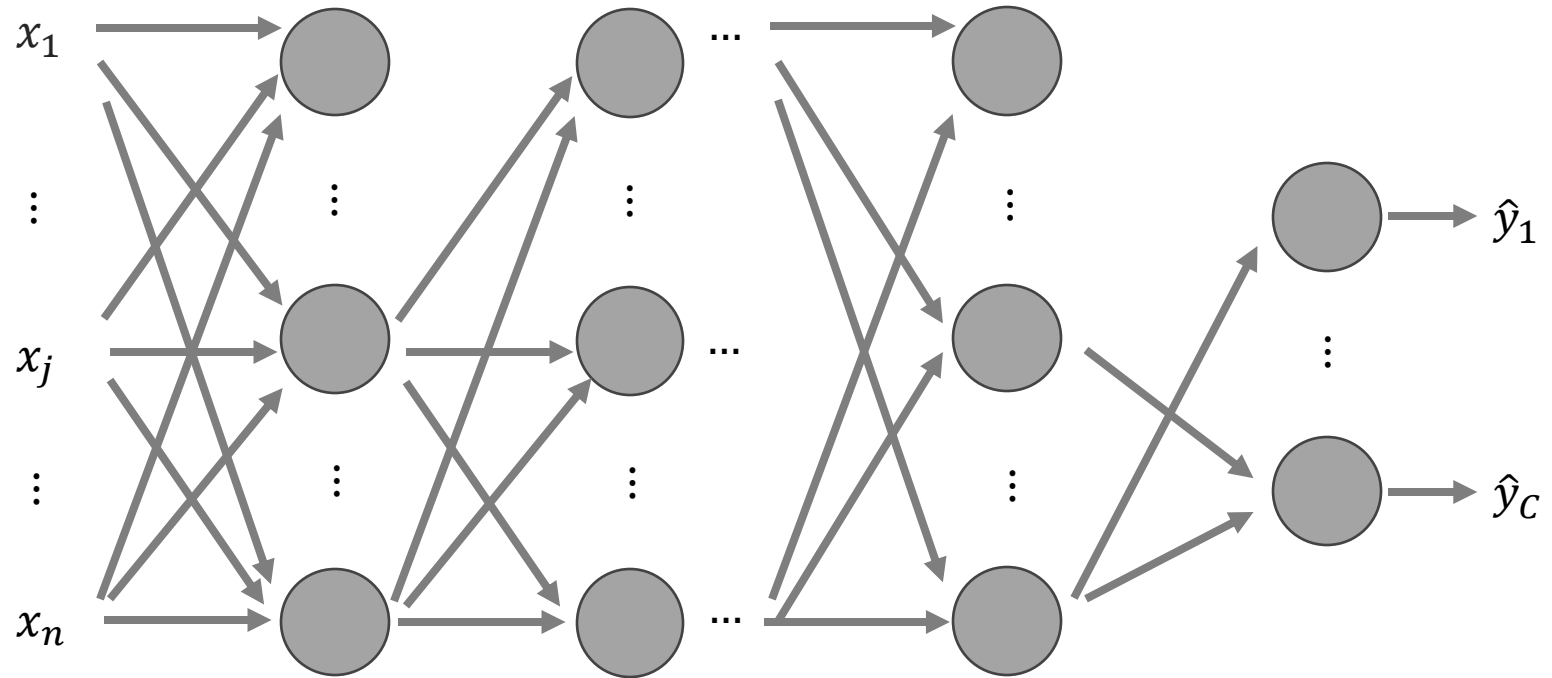
Prevent Overfitting



During training, **randomly** set some neurons' output to 0

Dropout

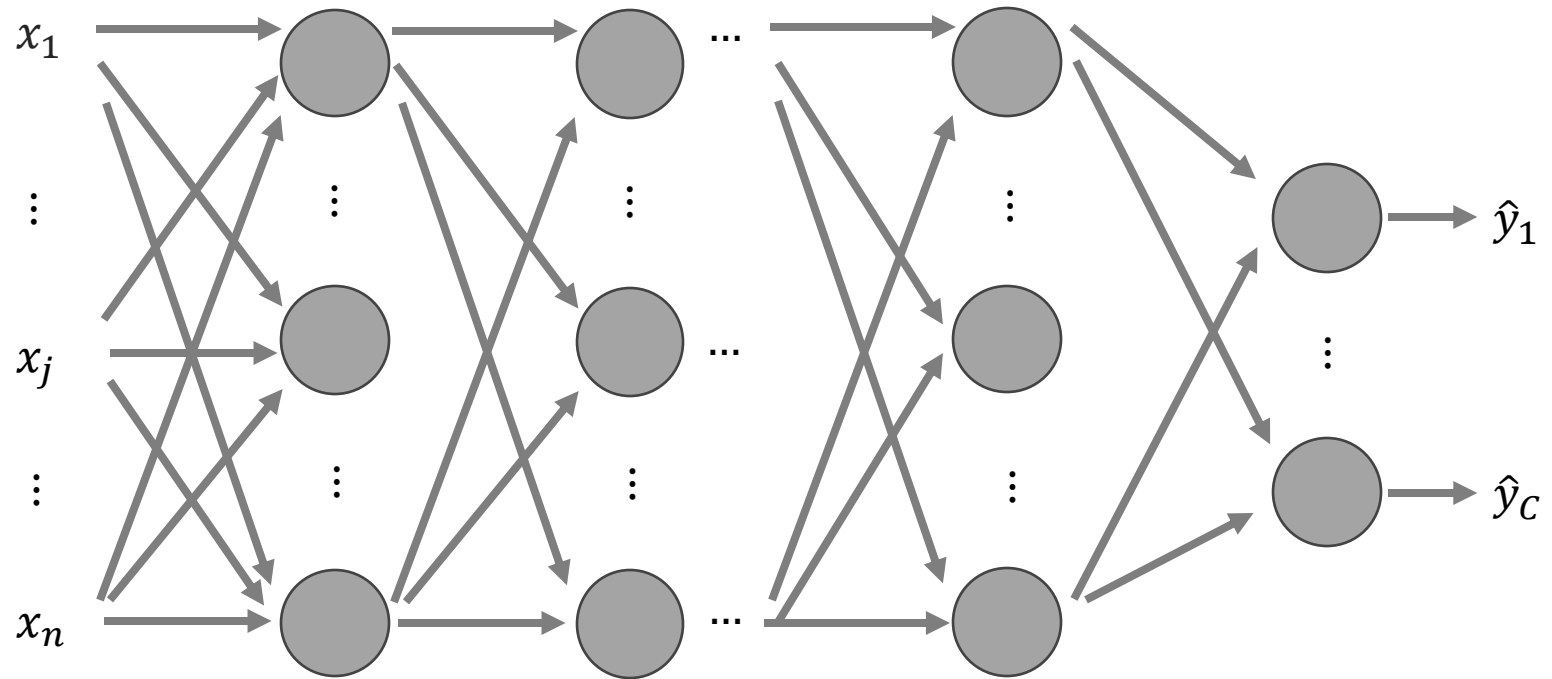
Prevent Overfitting



During training, **randomly** set some neurons' output to 0

Dropout

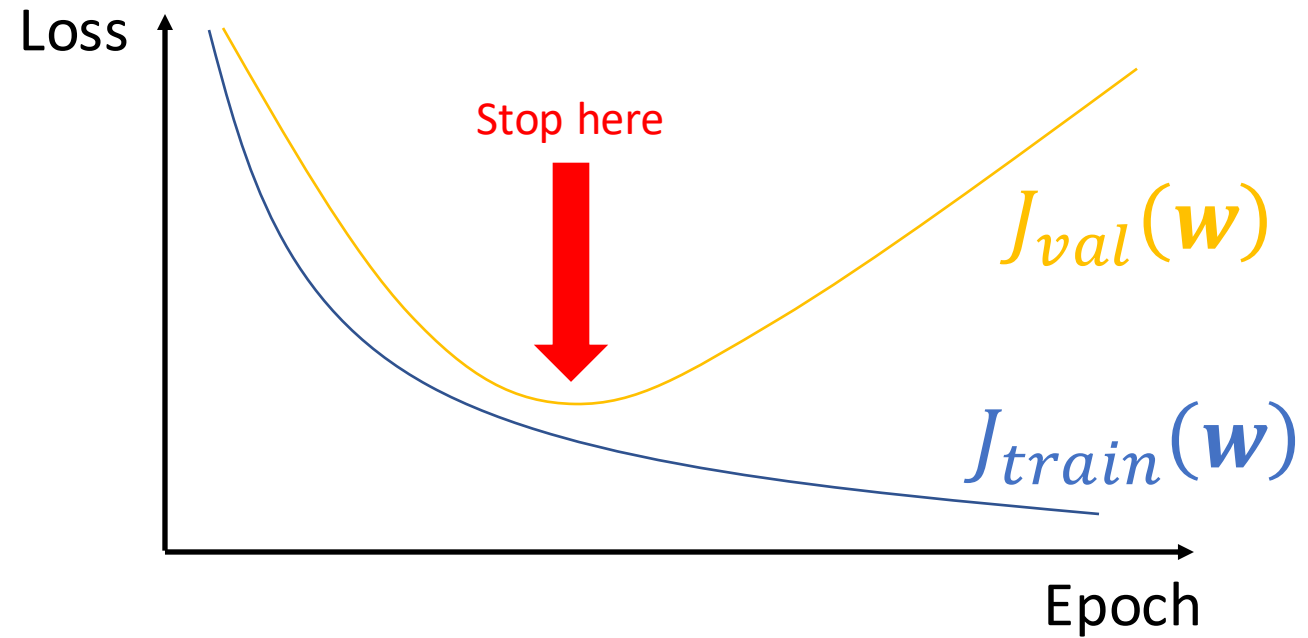
Prevent Overfitting



During training, **randomly** set some neurons' output to 0

Early Stopping

Prevent Overfitting



Vanishing/Exploding Gradient

- **Vanishing gradient:**

- **Small** gradients got multiplied again and again until it reaches almost zero.
- Mitigation: Change Activation Functions (Will be discussed during tutorial)

- **Exploding gradient:**

- **Large** gradients got multiplied again and again until it overflows.
- Mitigation: Gradient Clipping ~ clip gradient within range $[-clip_value, clip_value]$

Summary

- **Recurrent Neural Networks**

- Recurrent Neural Networks can capture contextual information.
- RNN types: Many-to-Many, Many-to-One, and One-to-Many
- Applications: sentiment analysis, speech recognition, etc

- **Self-Attention**

- Self-Attention layer can handle sequential data in parallel.
- Positional Encoding can be used to encode the positional information.

- **Issues with Deep Learning**

- Overfitting: Dropout, Early Stopping
- Vanishing Gradient: Change activation function
- Exploding Gradient: Gradient Clipping

Coming Up Next Week

- **Transformer**
- **AI Ethics**
- ...

To Do

- **Lecture Training 11**
 - +250 Free EXP
 - +100 Early bird bonus