

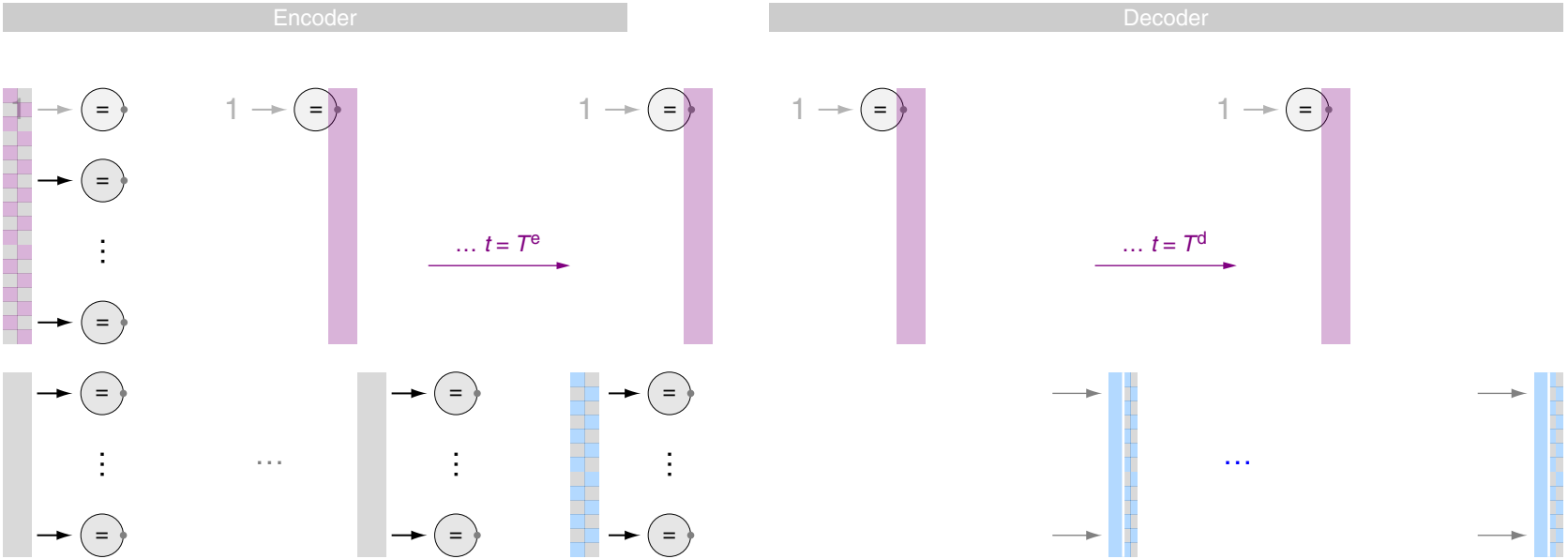
Chapter ML:IX

IX. Deep Learning

- ❑ Elements of Deep Learning
- ❑ Convolutional Neural Networks
- ❑ Autoencoder Networks
- ❑ Recurrent Neural Networks
- ❑ Long-Term Dependencies
- ❑ RNNs for Machine Translation
- ❑ Attention Mechanism
- ❑ Transformer
- ❑ Transformer Language Models
- ❑ Pretraining and Finetuning

Recurrent Neural Networks

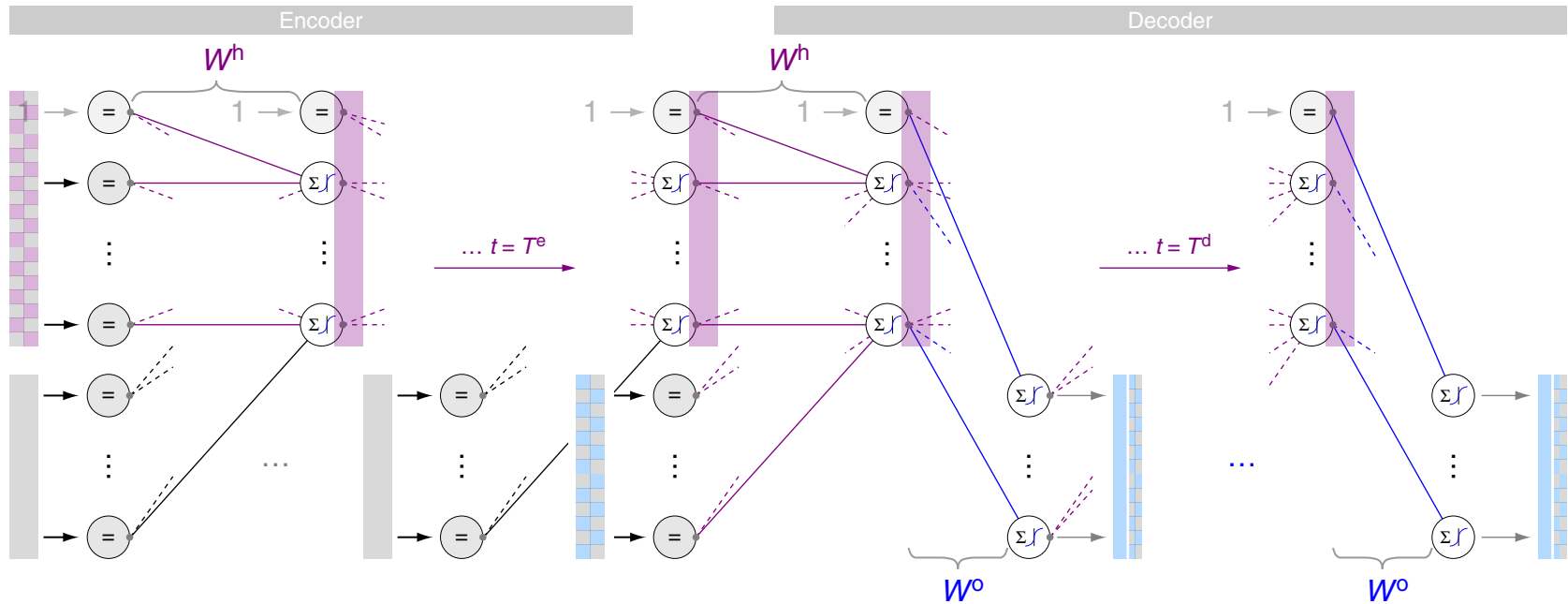
Notation I (MLP matrices) [notation: MLP matrices, computational graph, language model]



Recurrent Neural Networks

Notation I (MLP matrices)

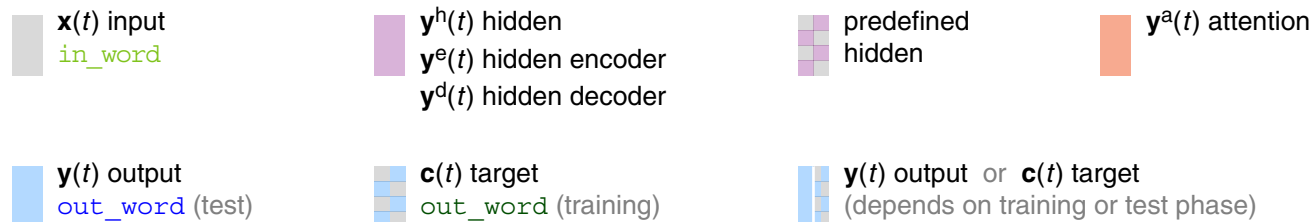
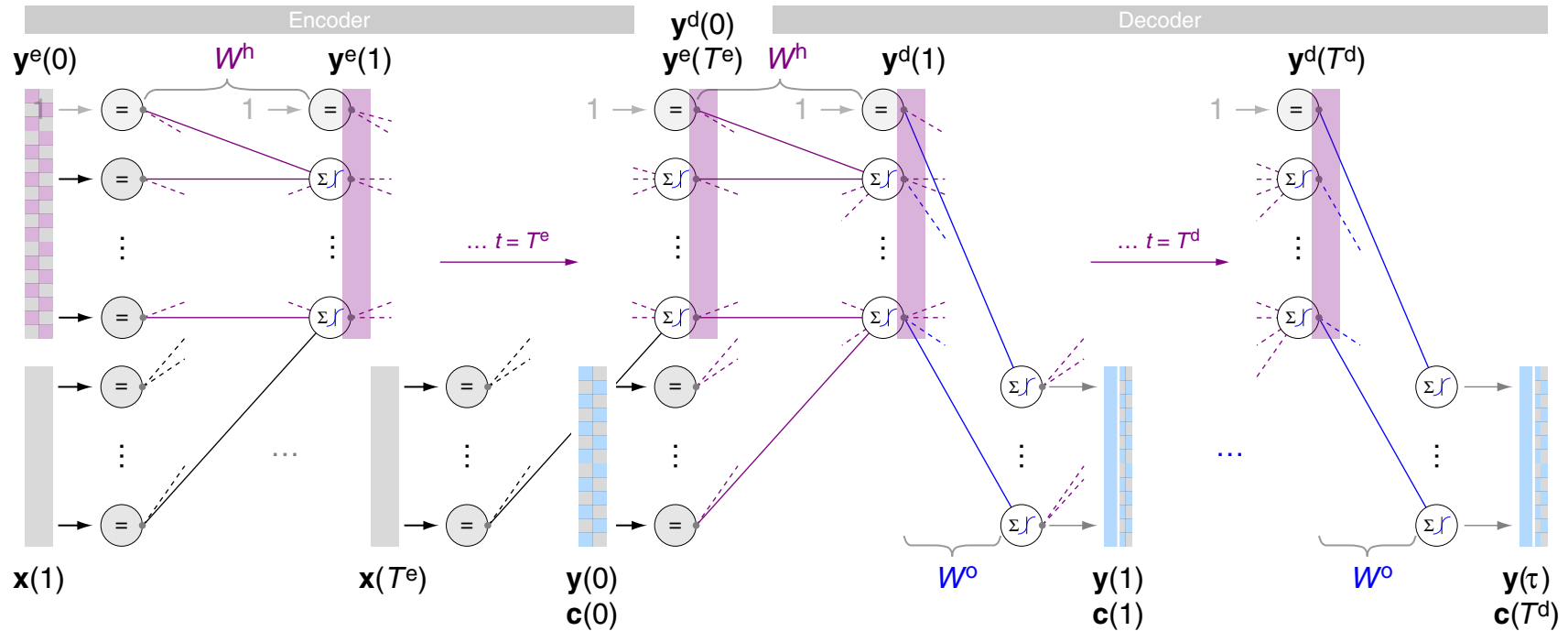
[notation: MLP matrices, computational graph, language model]



Recurrent Neural Networks

Notation I (MLP matrices)

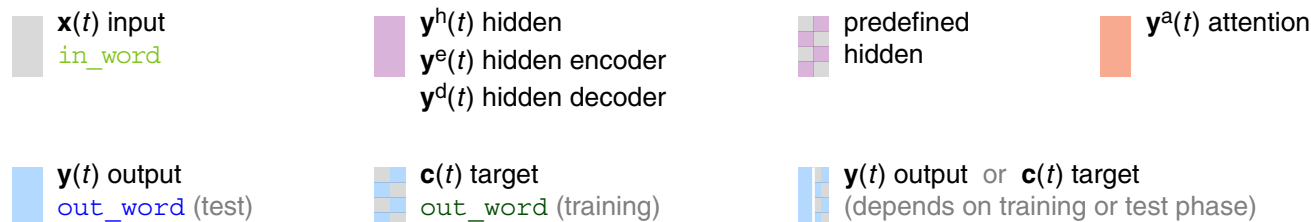
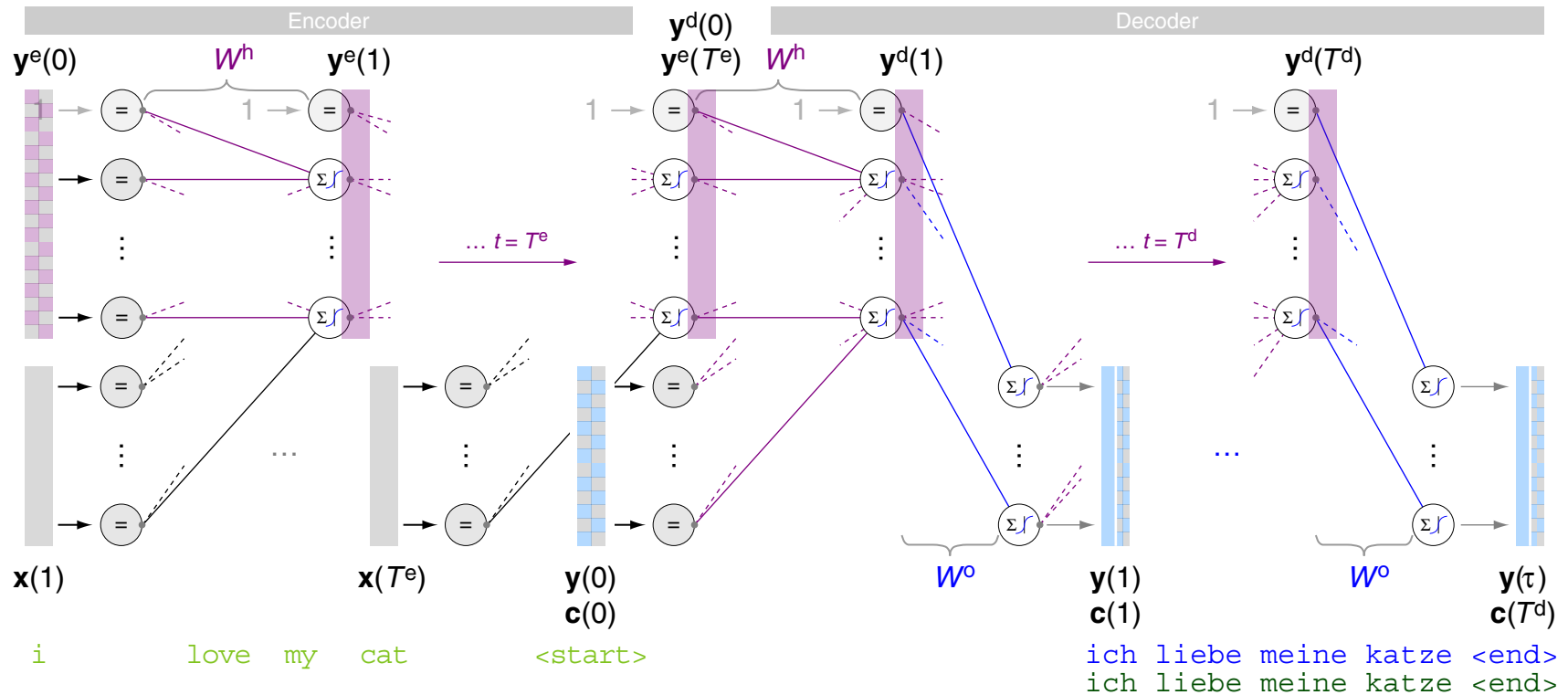
[notation: MLP matrices, computational graph, language model]



Recurrent Neural Networks

Notation I (MLP matrices)

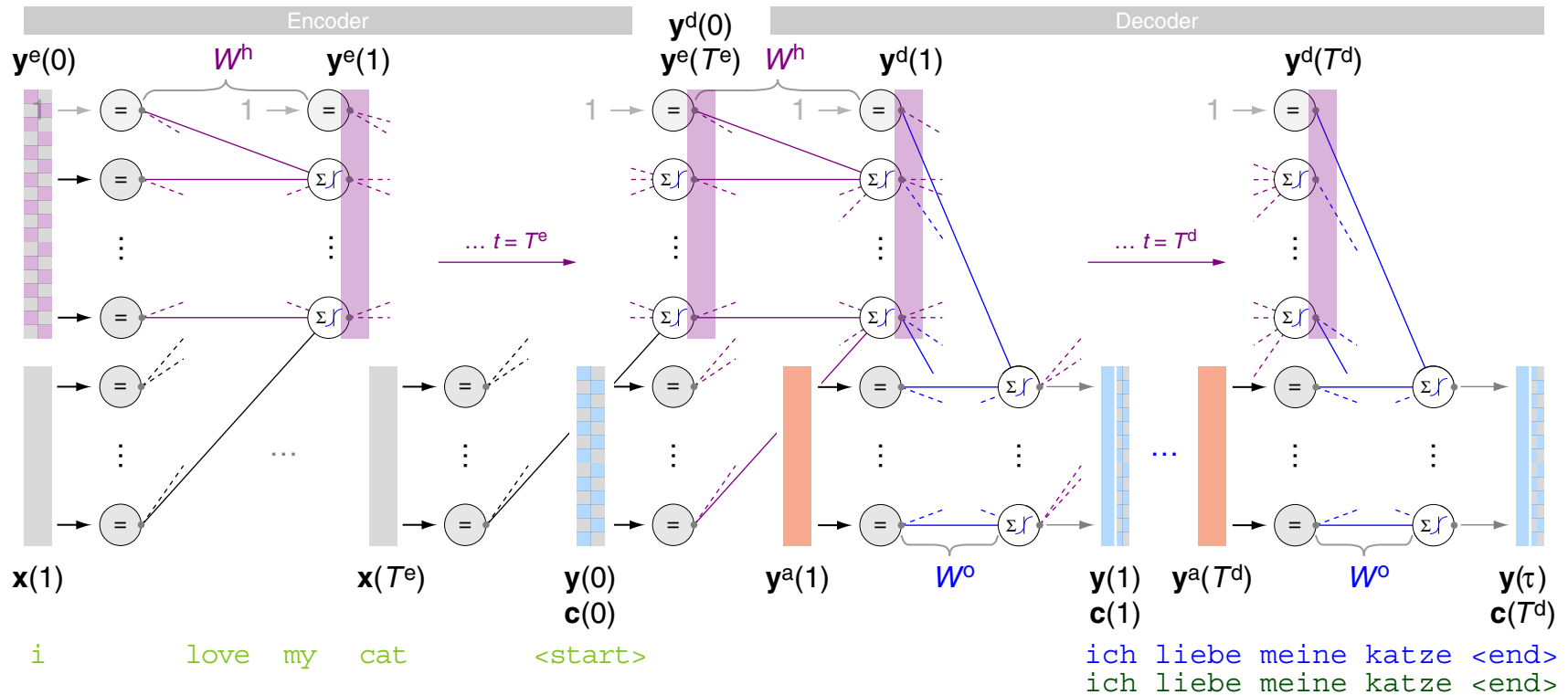
[notation: MLP matrices, computational graph, language model]



Recurrent Neural Networks

Notation I (MLP matrices)

[notation: MLP matrices, computational graph, language model]



Remarks:

- ❑ A hidden vector is the result of an intermediate computation in a multilayer network. If sequences are processed, hidden vectors may be distinguished as hidden *encoder* vectors (which consider the input at a certain time step) and hidden *decoder* vectors (which generate the output at a certain time step).
- ❑ A predefined hidden vector is used to for initialization purposes for the first hidden layer in a sequence-processing multilayer network. Typically, it is a constant vector of zeroes.
- ❑ Attention vectors combine the information in hidden vectors in a way that is specific to a certain time step. Keyword: vanishing gradient problem
- ❑ A target vector (or target vector sequence) encodes the desired output.
Keywords: supervised learning, ground truth

Recurrent Neural Networks

Types of Learning Tasks [Recap]

(S1) **sequence \rightarrow class**

sentence $\rightarrow \{\oplus, \ominus\}$

i love my cat $\rightarrow \oplus$

(S2) class \rightarrow sequence

$\{\oplus, \ominus\} \rightarrow$ sentence

$\oplus \rightarrow$ i love my cat

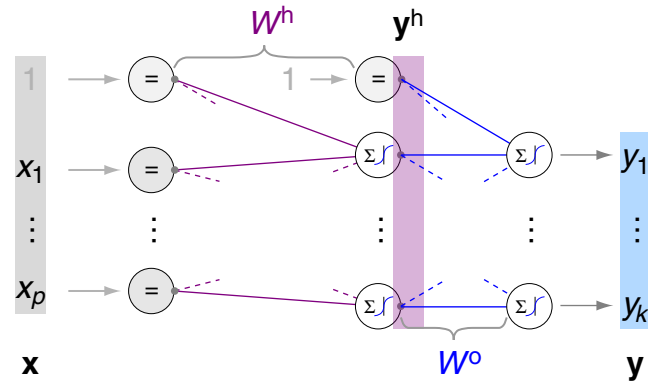
(S3) sequence \rightarrow sequence

English sentence \rightarrow German sentence

i love my cat \rightarrow ich liebe meine katze

Recurrent Neural Networks

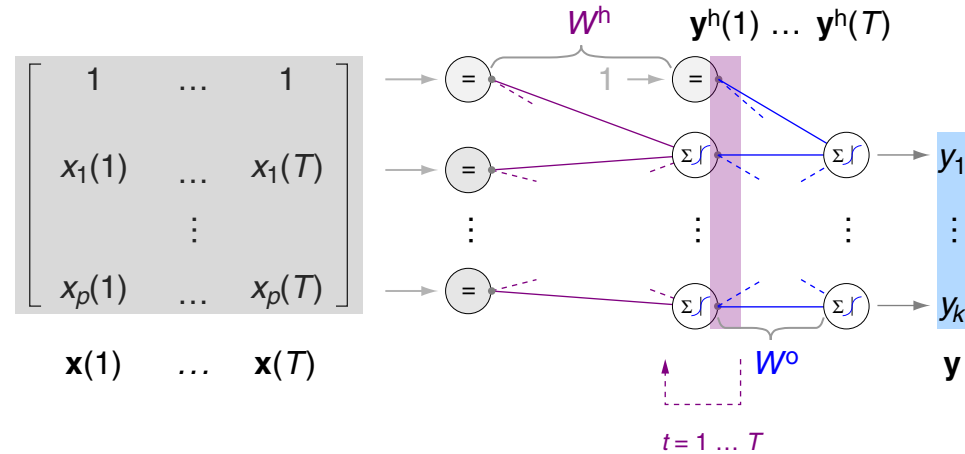
RNN Sequence Encoding



- ❑ One p -dimensional input vector \mathbf{x} .
- ❑ One hidden layer (general: $d-1$ hidden layers, i.e., d active layers).
- ❑ One k -dimensional output vector $\mathbf{y}(\mathbf{x})$.

Recurrent Neural Networks

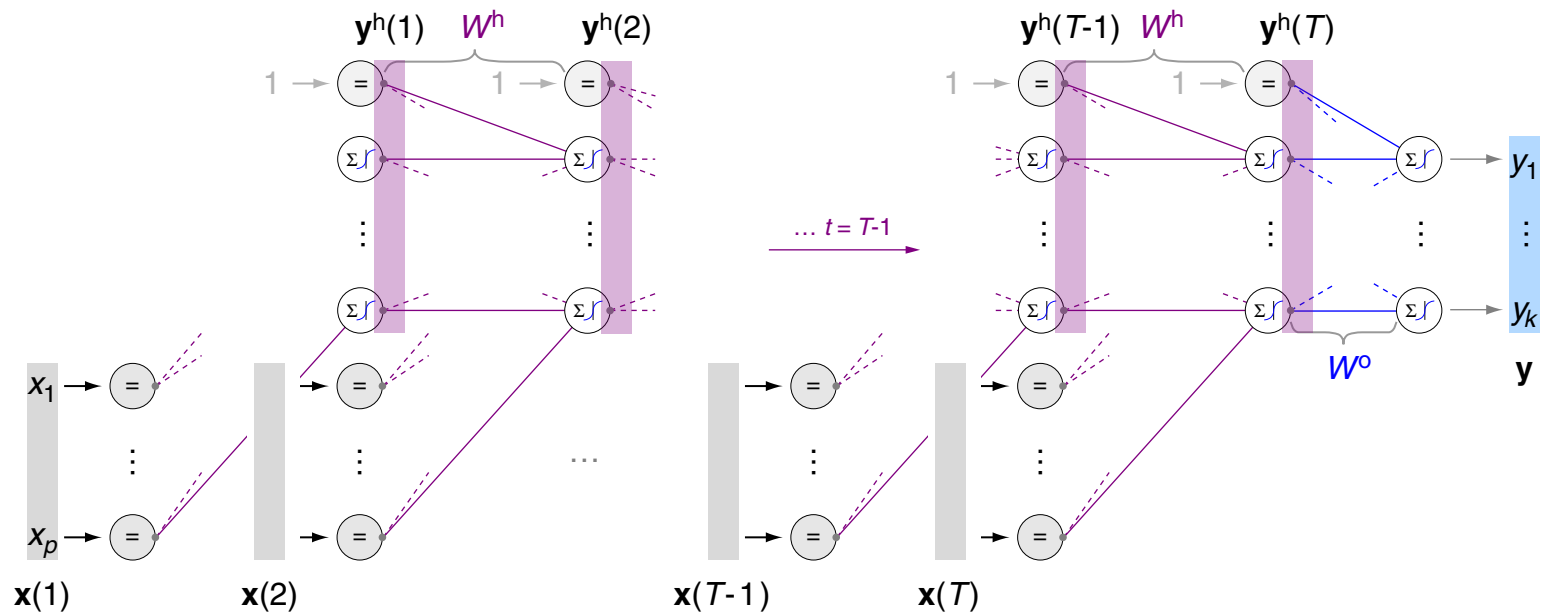
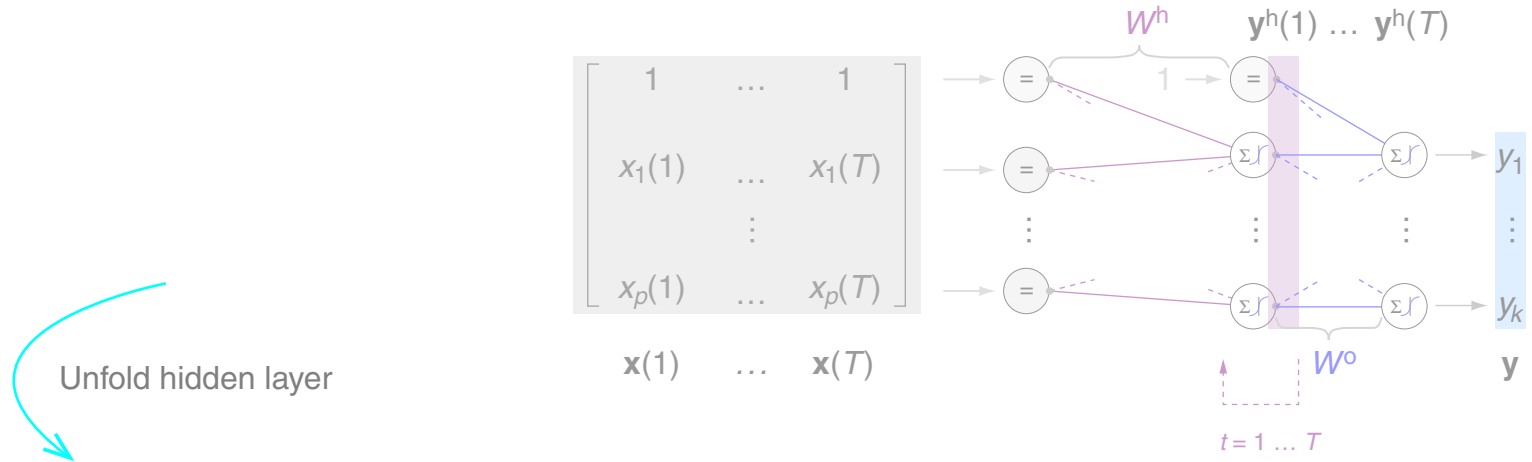
RNN Sequence Encoding (continued)



- ❑ Sequence of p -dimensional input vectors $[\mathbf{x}(1), \dots, \mathbf{x}(T)]$.
- ❑ One hidden layer that is recurrently updated.
- ❑ One k -dimensional output vector $\mathbf{y}([\mathbf{x}(1), \dots, \mathbf{x}(T)])$ or \mathbf{y} .

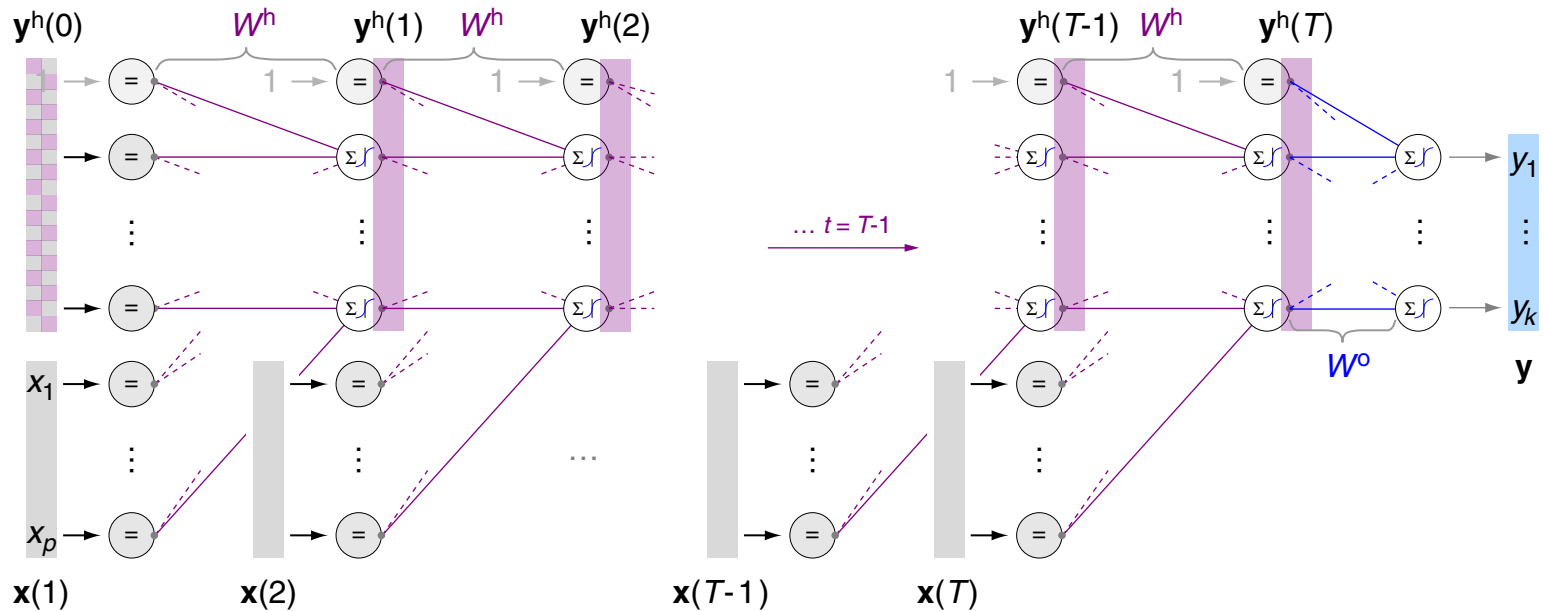
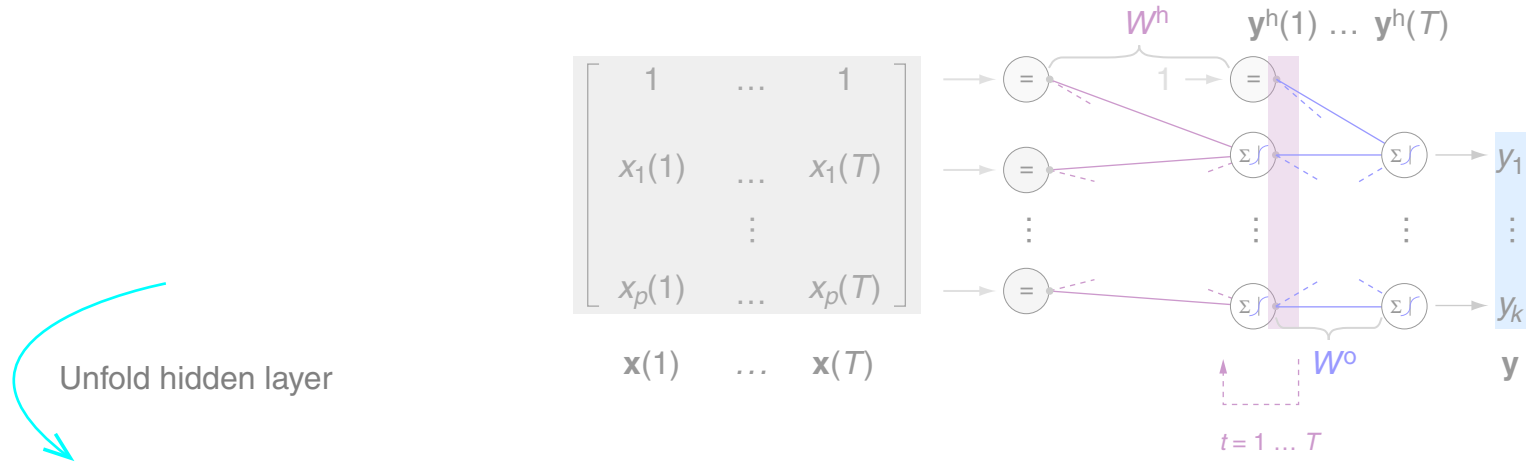
Recurrent Neural Networks

RNN Sequence Encoding (continued)



Recurrent Neural Networks

RNN Sequence Encoding (continued)



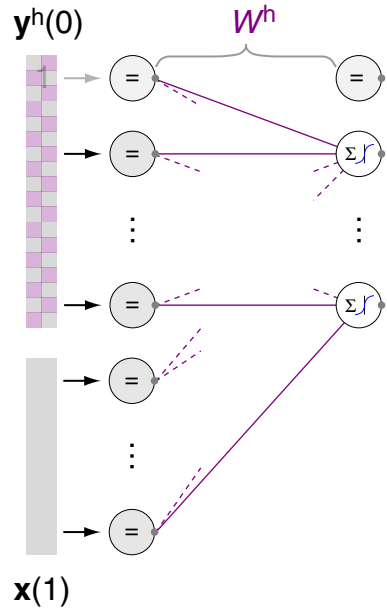
Remarks:

- ❑ An input sequence is written in brackets, $[\mathbf{x}(1), \dots, \mathbf{x}(T)]$, where $\mathbf{x}(t), t = 1, \dots, T$, denotes the input vector at time step t .
- ❑ The words in the input sequence are usually one-hot-encoded, i.e., by a p -dimensional input vector with a “1” whose position indicates the word, and zeros elsewhere.

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 0 \rightarrow 1$

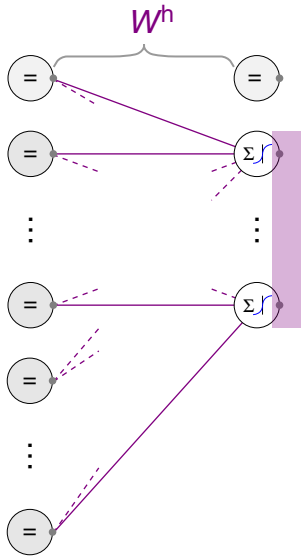


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: 0 \rightarrow 1

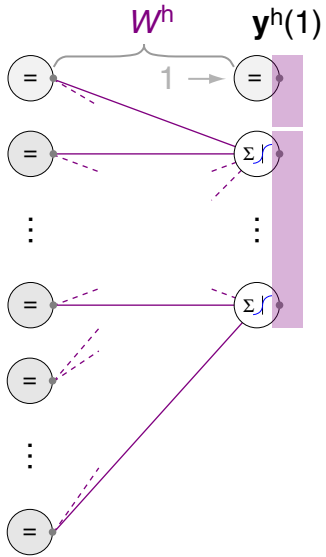


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: 0 \rightarrow 1

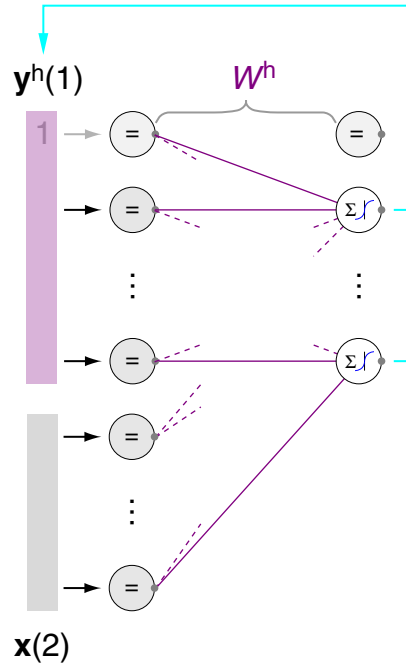


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 0 \rightarrow 1$

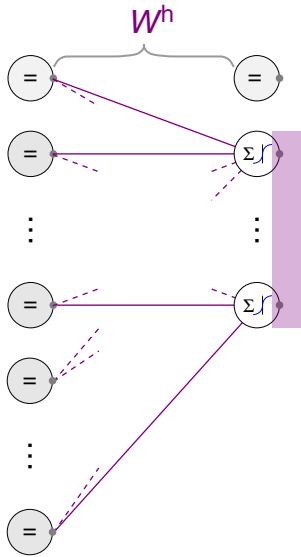


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: 1 \rightarrow 2

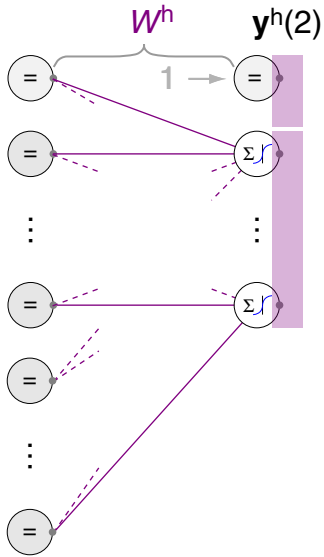


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: 1 \rightarrow 2

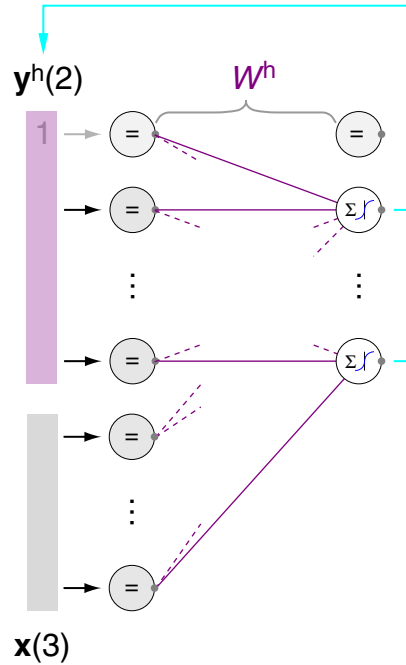


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 1 \rightarrow 2$

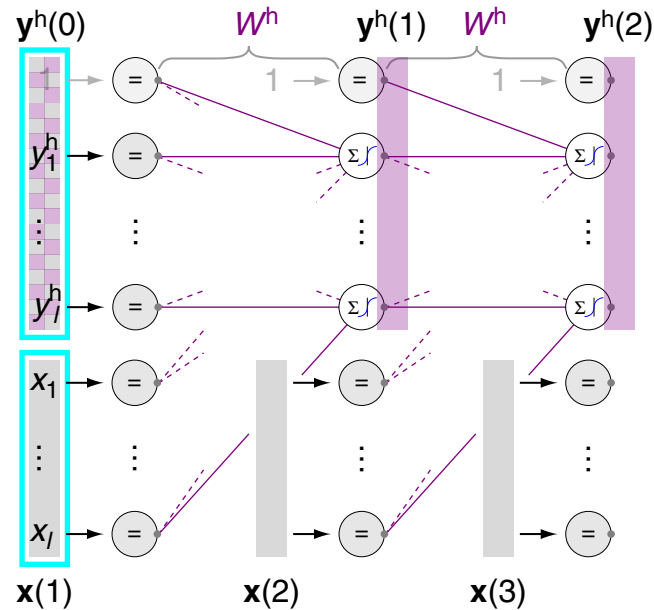
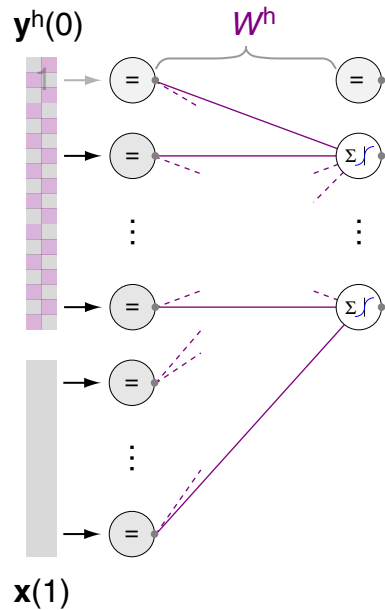


Input encoding over t .

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 0 \rightarrow 1$



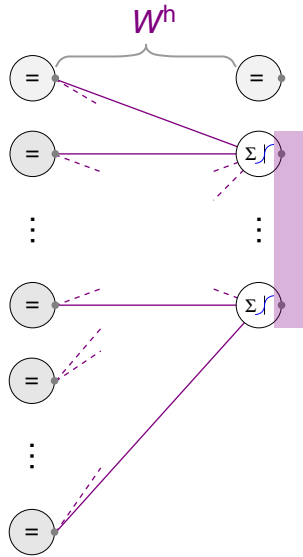
Input encoding over t .

The hidden layer at subsequent time steps.

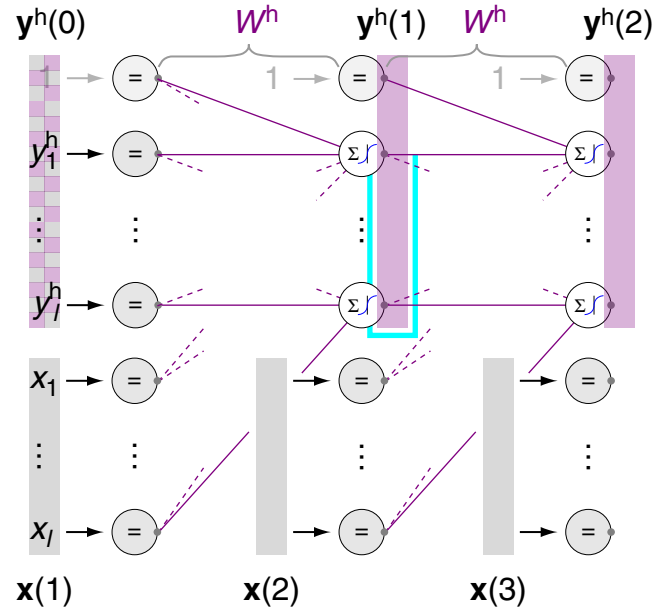
Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: 0 \rightarrow 1



Input encoding over t .

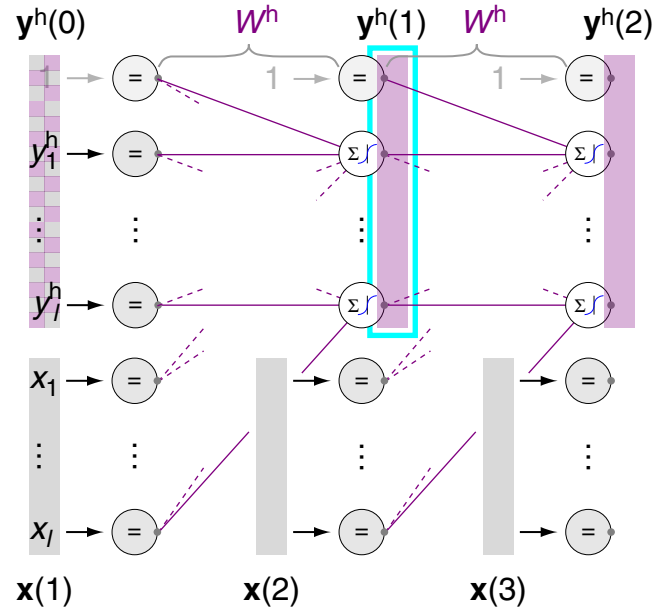
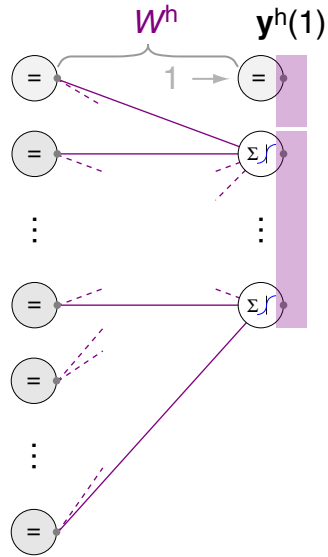


The hidden layer at subsequent time steps.

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: $0 \rightarrow 1$



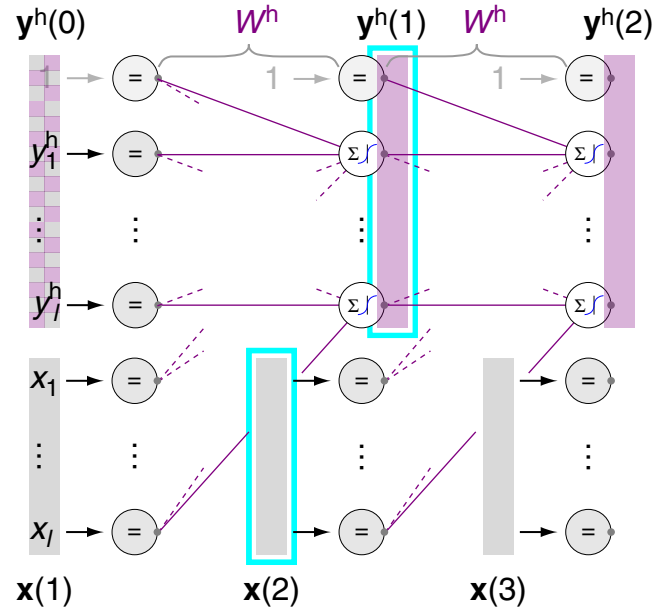
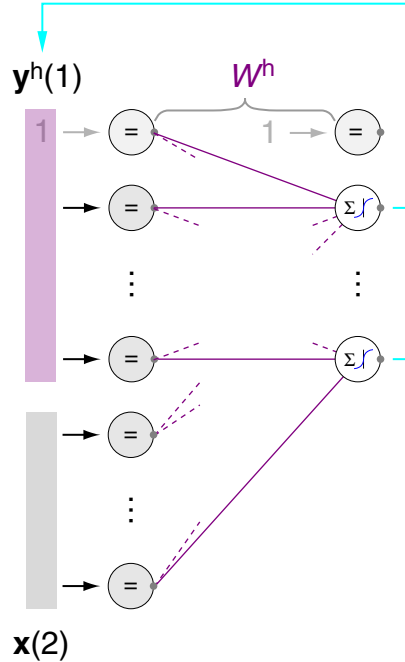
Input encoding over t .

The hidden layer at subsequent time steps.

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 0 \rightarrow 1$



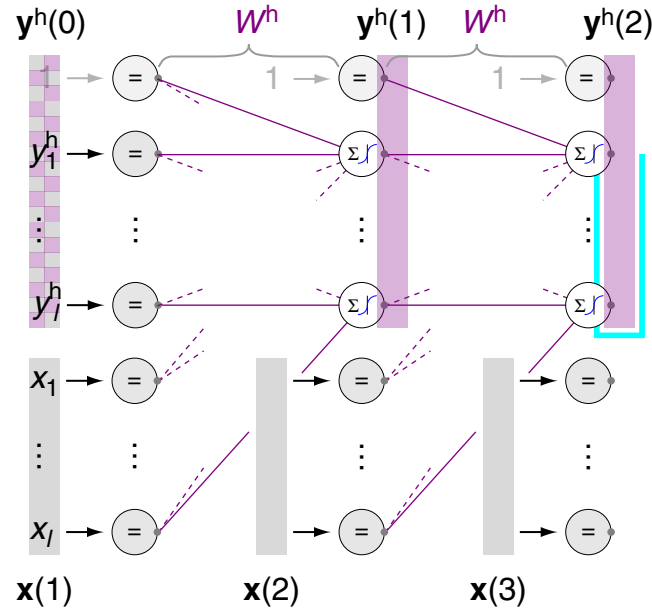
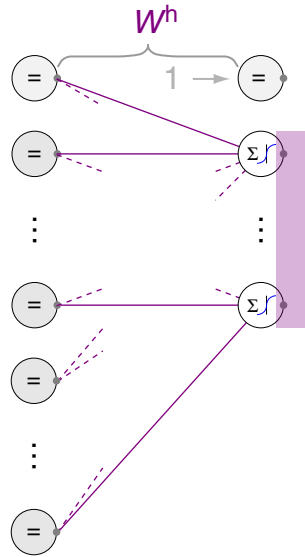
Input encoding over t .

The hidden layer at subsequent time steps.

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

t: $1 \rightarrow 2$



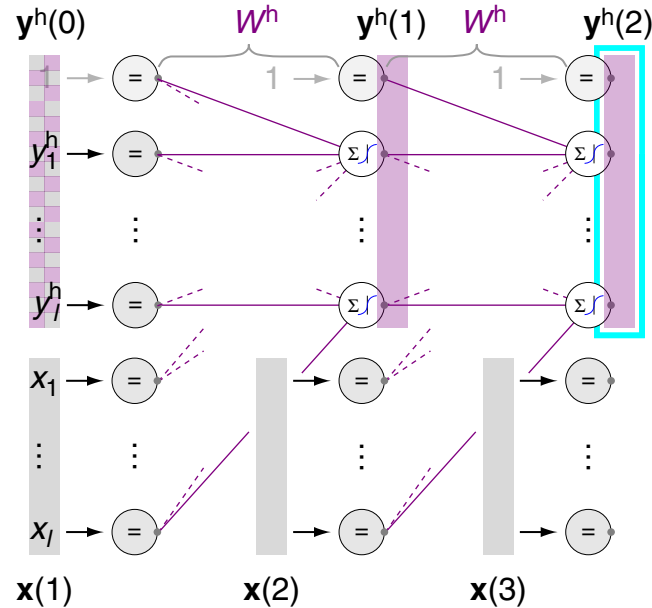
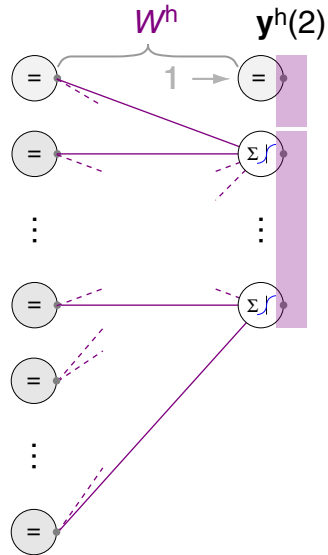
Input encoding over t .

The hidden layer at subsequent time steps.

Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 1 \rightarrow 2$



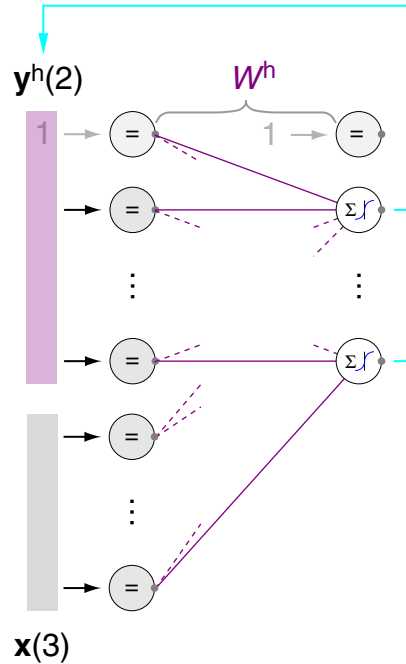
Input encoding over t .

The hidden layer at subsequent time steps.

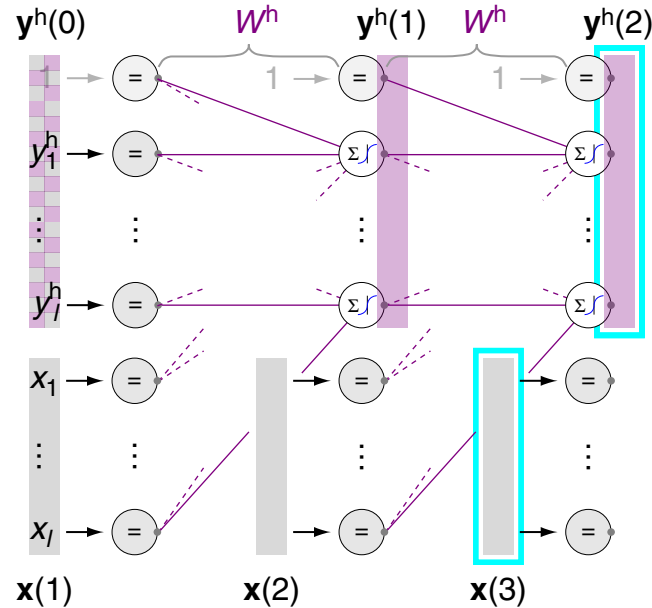
Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 1 \rightarrow 2$



Input encoding over t .

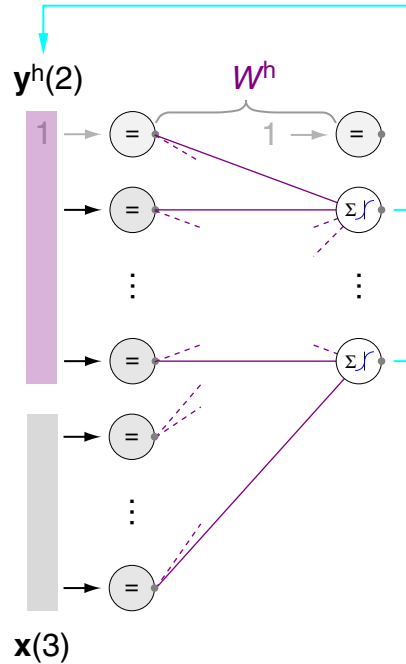


The hidden layer at subsequent time steps.

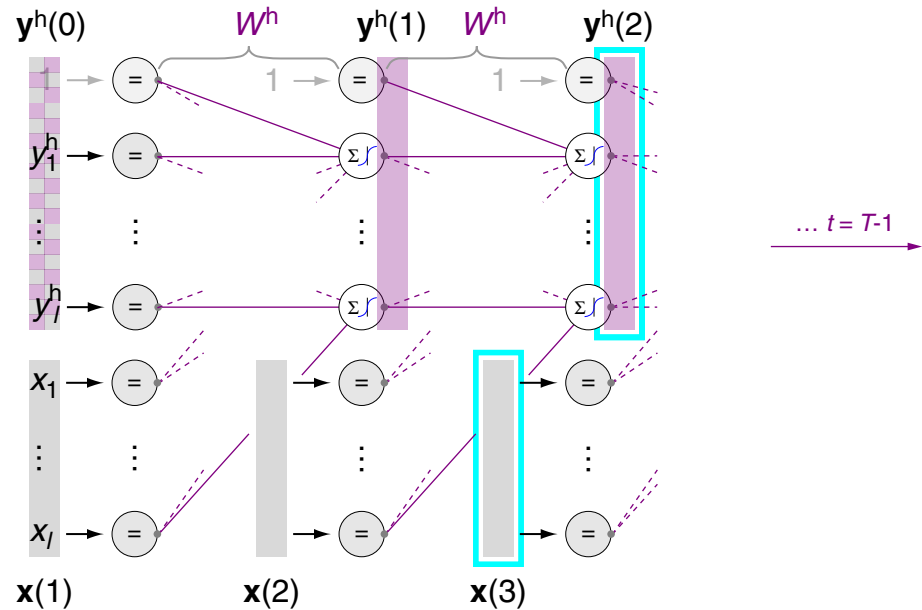
Recurrent Neural Networks

RNN Sequence Encoding (continued) [\[encoding overview\]](#)

$t: 1 \rightarrow 2$



Input encoding over t .



The hidden layer at subsequent time steps.

Recurrent Neural Networks

(S1) Sequence-to-Class: Sentiment Classification

- ❑ I love my cat. $\rightarrow \oplus$
- ❑ Cats and dogs lap water. $\rightarrow \oplus$
- ❑ It is raining cats and dogs. $\rightarrow \ominus$
- ❑ Cats and dogs are not allowed. $\rightarrow \ominus$
- ❑ Cats and dogs have always been natural enemies. $\rightarrow \ominus$

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water)

Recurrent Neural Networks

(S1) Sequence-to-Class: Sentiment Classification

- ❑ I love my cat. $\rightarrow \oplus$
- ❑ Cats and dogs lap water. $\rightarrow \oplus$
- ❑ It is raining cats and dogs. $\rightarrow \ominus$
- ❑ Cats and dogs are not allowed. $\rightarrow \ominus$
- ❑ Cats and dogs have always been natural enemies. $\rightarrow \ominus$

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water)

Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(4)] = \left[\begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \right]$$
$$\hat{=} [\text{word}_{11}, \text{word}_{15}, \text{word}_{16}, \text{word}_6]$$
$$\hat{=} \text{I love my cat}$$

Recurrent Neural Networks

(S1) Sequence-to-Class: Sentiment Classification

- ❑ I love my cat. $\rightarrow \oplus$
- ❑ Cats and dogs lap water. $\rightarrow \oplus$
- ❑ It is raining cats and dogs. $\rightarrow \ominus$
- ❑ Cats and dogs are not allowed. $\rightarrow \ominus$
- ❑ Cats and dogs have always been natural enemies. $\rightarrow \ominus$

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water)

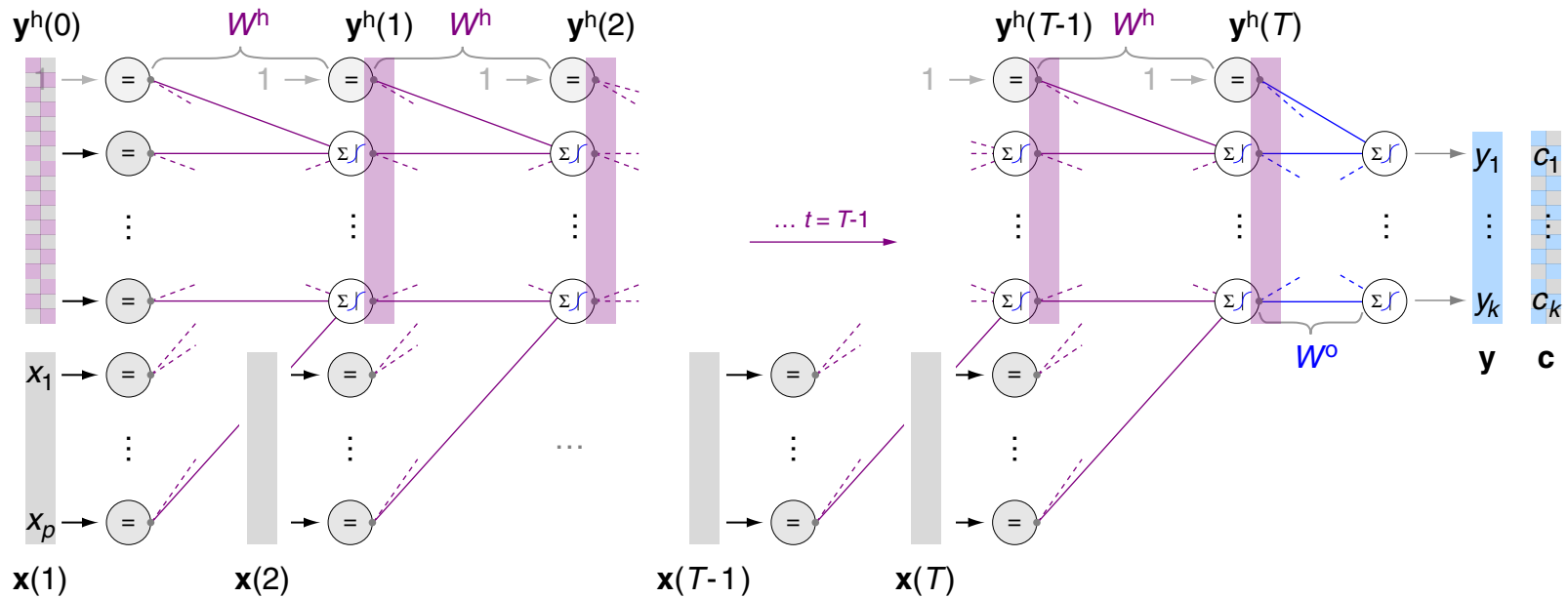
Input:
$$[\mathbf{x}(1), \dots, \mathbf{x}(4)] = \left[\begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \right]$$
$$\hat{=} [\text{word}_{11}, \text{word}_{15}, \text{word}_{16}, \text{word}_6]$$
$$\hat{=} \text{I love my cat}$$

Output:
$$\mathbf{y}([\mathbf{x}(1), \dots, \mathbf{x}(4)]) = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

Target:
$$\mathbf{c} = \begin{pmatrix} \oplus \\ \cdot \end{pmatrix}$$

Recurrent Neural Networks

(S1) Sequence-to-Class Mapping with RNNs



Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(T)]$$

Hidden:

$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right), t = 1, \dots, T$$

Target:

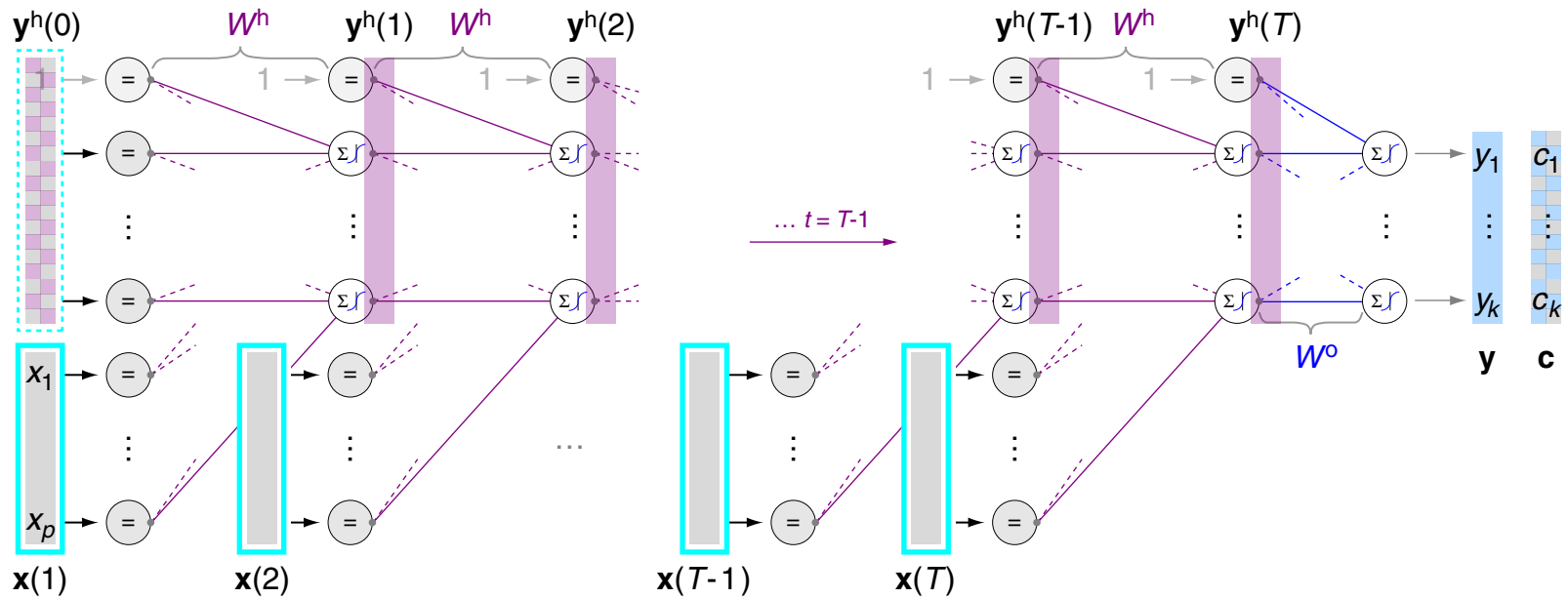
$$\mathbf{c}$$

Output:

$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(T))$$

Recurrent Neural Networks

(S1) Sequence-to-Class Mapping with RNNs



Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(T)]$$

Output:

$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(T))$$

Hidden:

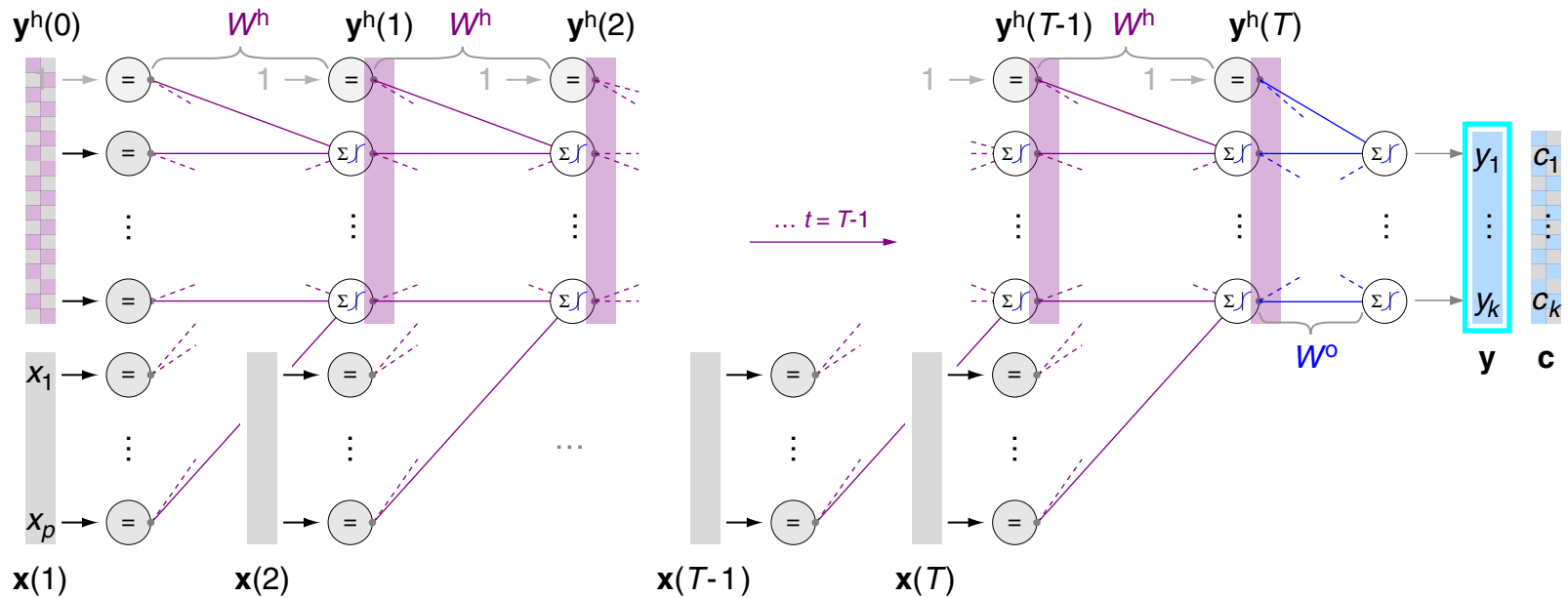
$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right), t = 1, \dots, T$$

Target:

$$\mathbf{c}$$

Recurrent Neural Networks

(S1) Sequence-to-Class Mapping with RNNs



Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(T)]$$

Hidden:

$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right), t = 1, \dots, T$$

Target:

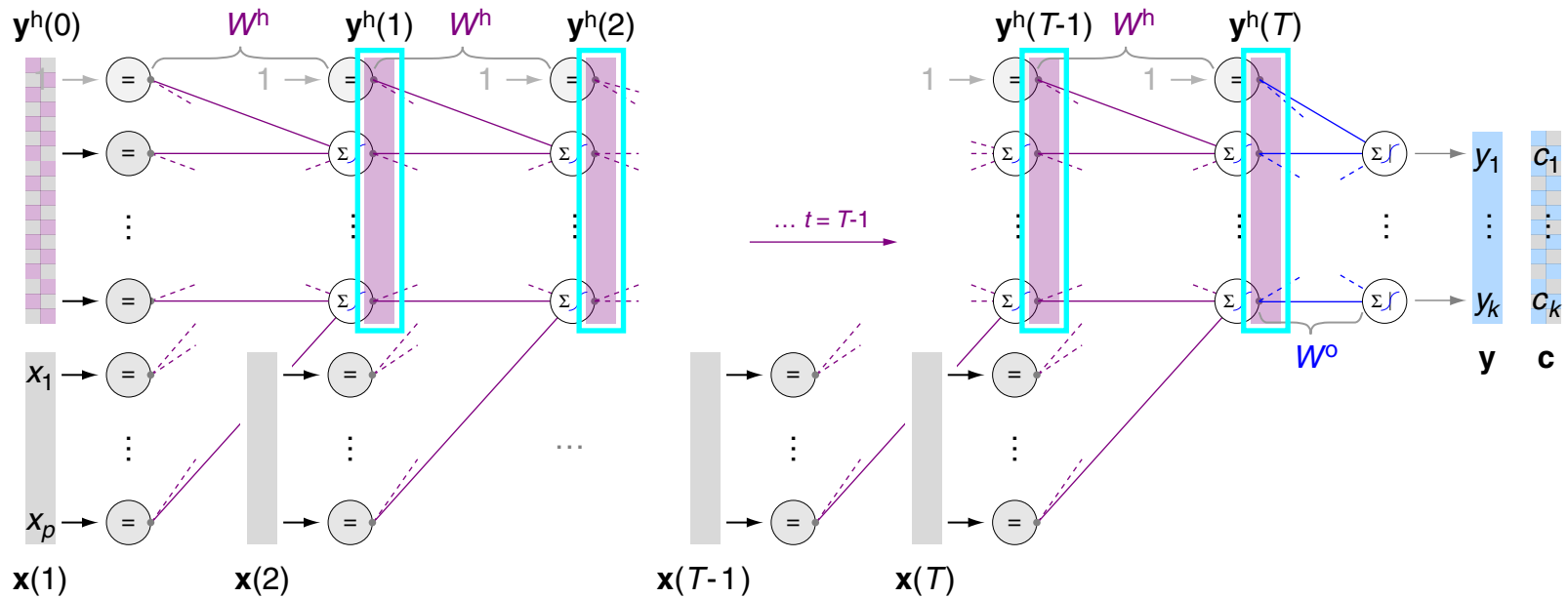
$$\mathbf{c}$$

Output:

$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(T))$$

Recurrent Neural Networks

(S1) Sequence-to-Class Mapping with RNNs



Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(T)]$$

Output:

$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(T))$$

Hidden:

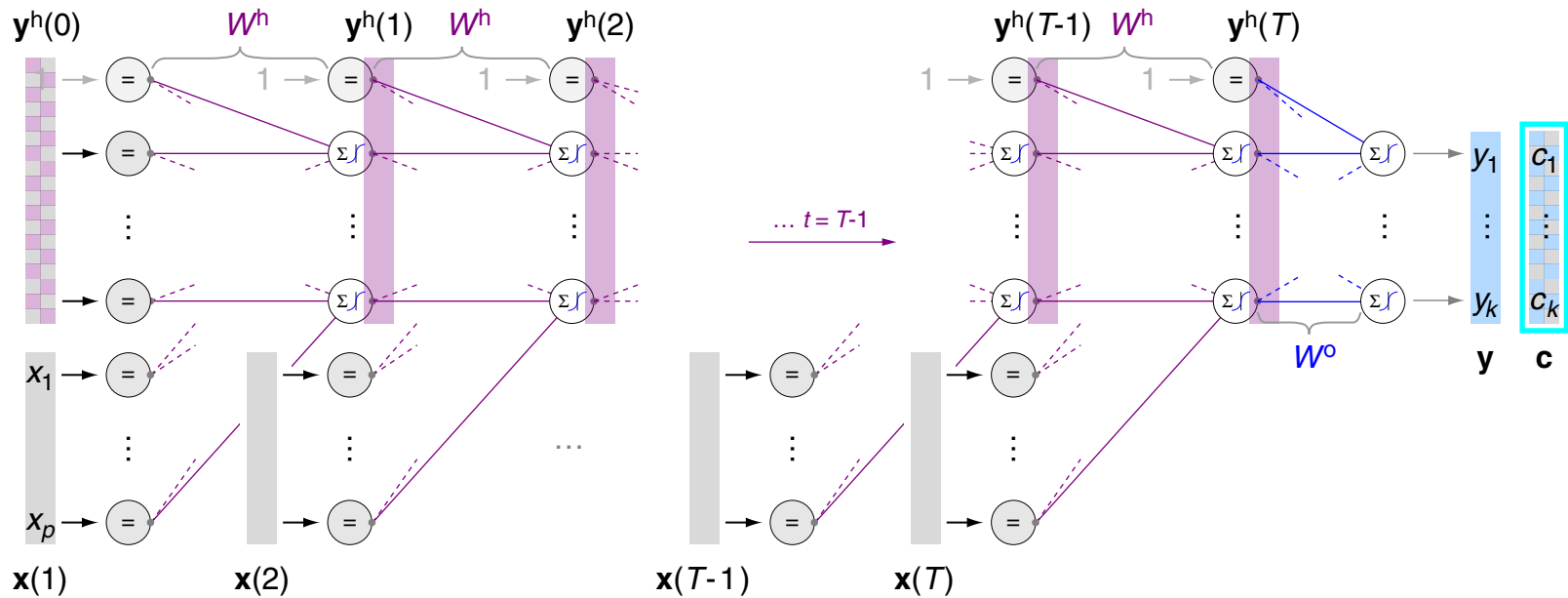
$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right), t = 1, \dots, T$$

Target:

\mathbf{c}

Recurrent Neural Networks

(S1) Sequence-to-Class Mapping with RNNs



Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(T)]$$

Output:

$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(T))$$

Hidden:

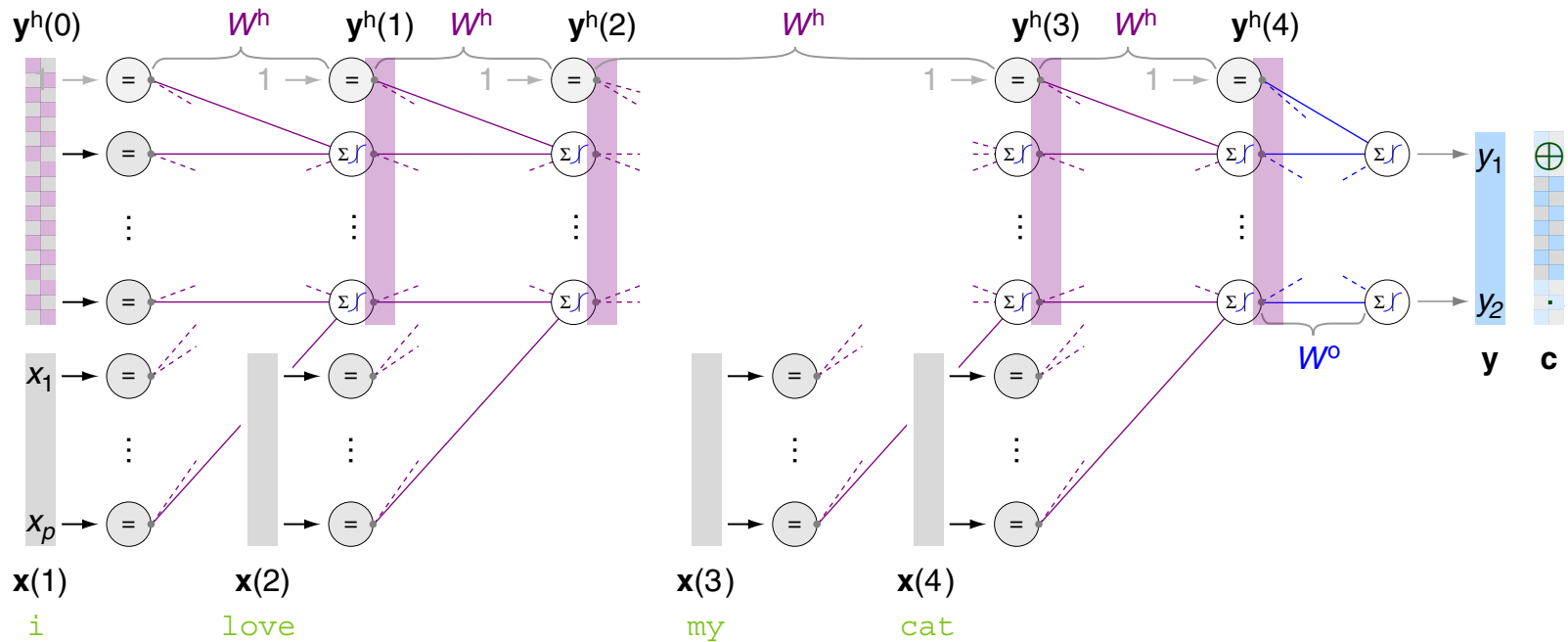
$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right), t = 1, \dots, T$$

Target:

$$\mathbf{c}$$

Recurrent Neural Networks

(S1) Sequence-to-Class Mapping with RNNs



Input:

$$[\mathbf{x}(1), \dots, \mathbf{x}(4)]$$

Output:

$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(4))$$

Hidden:

$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right), t = 1, \dots, 4$$

Target:

\mathbf{c}

Remarks:

- ❑ We denote $\mathbf{y}^h(0)$ not as input since this kind of predefined hidden vector does not contain any “actual knowledge”, but is usually initialized as vector of zeros.
- ❑ To keep the illustrations clear we use the bag-of-words model for representing (= embedding) the words as vectors $\mathbf{x}(t)$.

In practice, however, one considers semantically stronger (language-model-based) embeddings, which also encode information about neighborhoods and occurrence probabilities. In this regard, either a previously computed embedding can be used, or the embedding can be learned along with the task, end-to-end.

- ❑ Recap. $\sigma_{\Delta}()$ denotes the softmax function. $\sigma_{\Delta} : \mathbf{R}^k \rightarrow \Delta^{k-1}$, generalizes the logistic (sigmoid) function to k dimensions (to k exclusive classes), where $\sigma_{\Delta}(\mathbf{z})|_i = e^{z_i} / \sum_{j=1}^k e^{z_j}$. [\[Wikipedia\]](#)

$\Delta^{k-1} \subset \mathbf{R}^k$ denotes the standard $k-1$ -simplex, which contains all k -tuples with non-negative elements that sum to 1. [\[Wikipedia\]](#)

Recurrent Neural Networks

The IGD Algorithm for Sequence-to-Class Tasks [IGD_{c2seq}]

Algorithm: IGD_{seq2c} Incremental Gradient Descent for RNNs at seq2class tasks.

Input: D Multiset of examples $([\mathbf{x}(1), \dots, \mathbf{x}(T)], \mathbf{c})$ with $\mathbf{x}(t) \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.

η Learning rate, a small positive constant.

Output: W^h, W^o Weight matrices. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $\mathbf{y}^h(0) = (0, \dots, 0)^T$, $t_{\text{epoch}} = 0$
2. **REPEAT**
3. $t_{\text{epoch}} = t_{\text{epoch}} + 1$
4. **FOREACH** $([\mathbf{x}(1), \dots, \mathbf{x}(T)], \mathbf{c}) \in D$ **DO**
5. Model function evaluation.
6. Calculation of residual vector.
- 7a.
- 7b. Calculation of derivative of the loss.
8. Parameter update $\hat{=}$ one gradient step down.
9. **ENDDO**
10. **UNTIL**(*convergence*($D, \mathbf{y}(\cdot), t_{\text{epoch}}$))
11. *return*(W^h, W^o)

Recurrent Neural Networks

The IGD Algorithm for Sequence-to-Class Tasks [IGD_{c2seq}]

Algorithm: IGD_{seq2c} Incremental Gradient Descent for RNNs at seq2class tasks.
Input: D Multiset of examples $([\mathbf{x}(1), \dots, \mathbf{x}(T)], \mathbf{c})$ with $\mathbf{x}(t) \in \mathbb{R}^p$, $\mathbf{c} \in \{0, 1\}^k$.
 η Learning rate, a small positive constant.
Output: W^h, W^o Weight matrices. (= hypothesis)

1. *initialize_random_weights*(W^h, W^o), $\mathbf{y}^h(0) = (0, \dots, 0)^T$, $t_{\text{epoch}} = 0$
2. **REPEAT**
3. $t_{\text{epoch}} = t_{\text{epoch}} + 1$
4. **FOREACH** $([\mathbf{x}(1), \dots, \mathbf{x}(T)], \mathbf{c}) \in D$ **DO**
5. **FOR** $t = 1$ **TO** T **DO** // forward propagation
$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{x}(t) \end{pmatrix} \right)$$
ENDDO
$$\mathbf{y} = \sigma_{\Delta} (W^o \mathbf{y}^h(T))$$
6. $\boldsymbol{\delta} = \mathbf{c} - \mathbf{y}([\mathbf{x}(1), \dots, \mathbf{x}(T)])$
- 7a. $\ell(\mathbf{w}) = l(\boldsymbol{\delta}) + \frac{\lambda}{n} R(\mathbf{w})$, $\nabla \ell(\mathbf{w}) = \text{autodiff}(\ell(), \mathbf{w})$ // backpropagation (7a+7b)
- 7b. $\Delta W^h = \eta \cdot \nabla^h \ell(\mathbf{w})$, $\Delta W^o = \eta \cdot \nabla^o \ell(\mathbf{w})$
8. $W^h = W^h + \Delta W^h$, $W^o = W^o + \Delta W^o$
9. **ENDDO**
10. **UNTIL** (*convergence*($D, \mathbf{y}(\cdot), t_{\text{epoch}}$))
11. *return*(W^h, W^o)

Recurrent Neural Networks

Types of Learning Tasks [Recap]

(S1) sequence \rightarrow class

sentence $\rightarrow \{\oplus, \ominus\}$

i love my cat $\rightarrow \oplus$

(S2) **class \rightarrow sequence**

$\{\oplus, \ominus\} \rightarrow$ **sentence**

$\oplus \rightarrow$ i love my cat

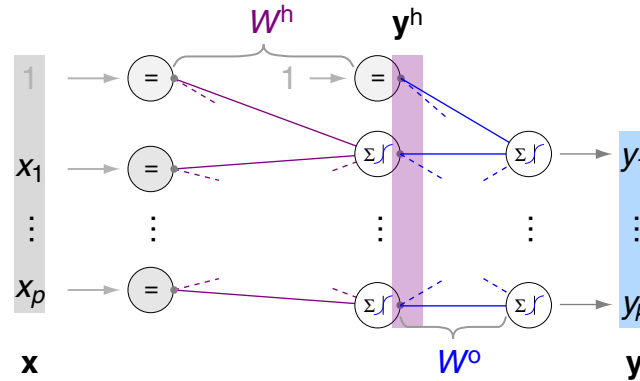
(S3) sequence \rightarrow sequence

English sentence \rightarrow German sentence

i love my cat \rightarrow ich liebe meine katze

Recurrent Neural Networks

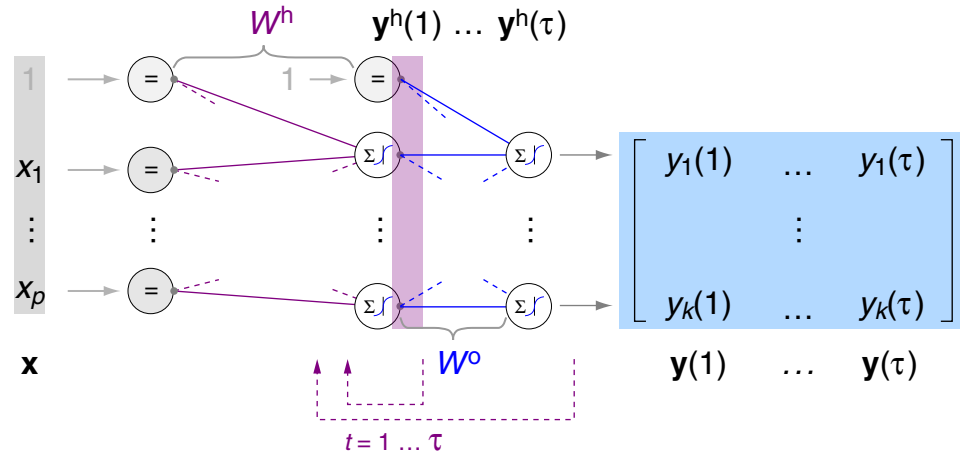
RNN Sequence Decoding



- ❑ One p -dimensional input vector \mathbf{x} .
- ❑ One hidden layer (general: $d-1$ hidden layers, i.e., d active layers).
- ❑ One k -dimensional output vector $\mathbf{y}(\mathbf{x})$.

Recurrent Neural Networks

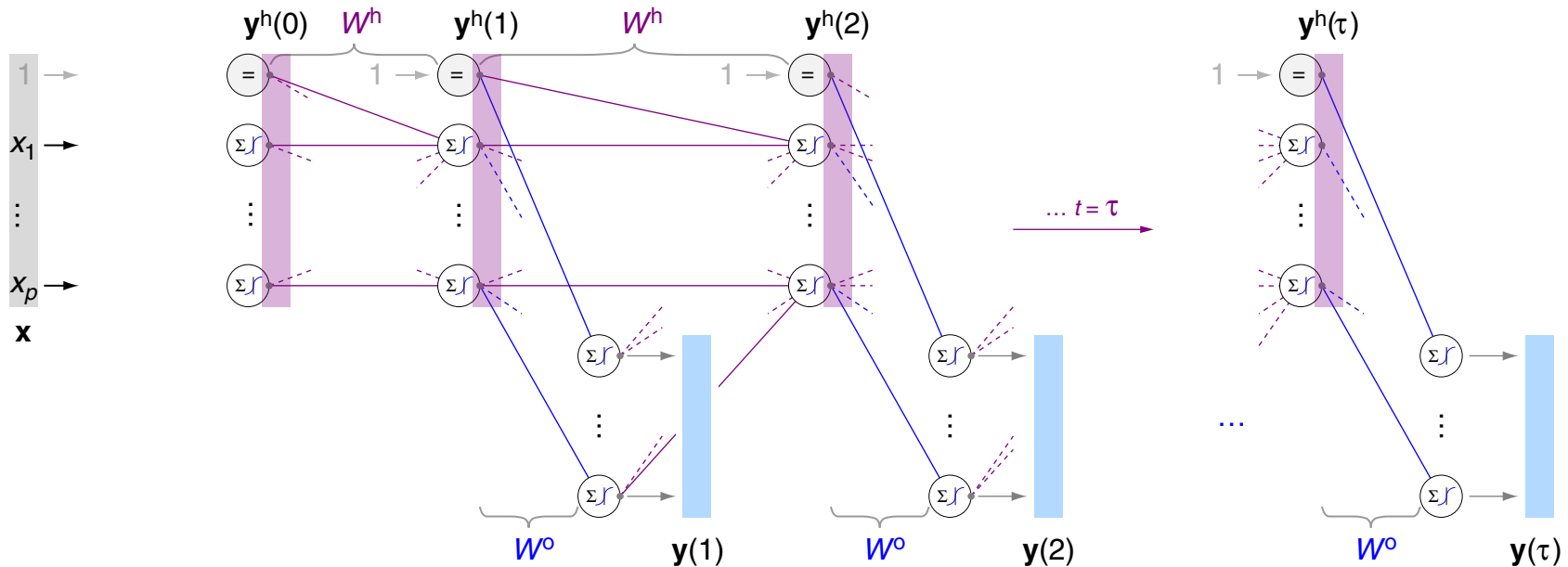
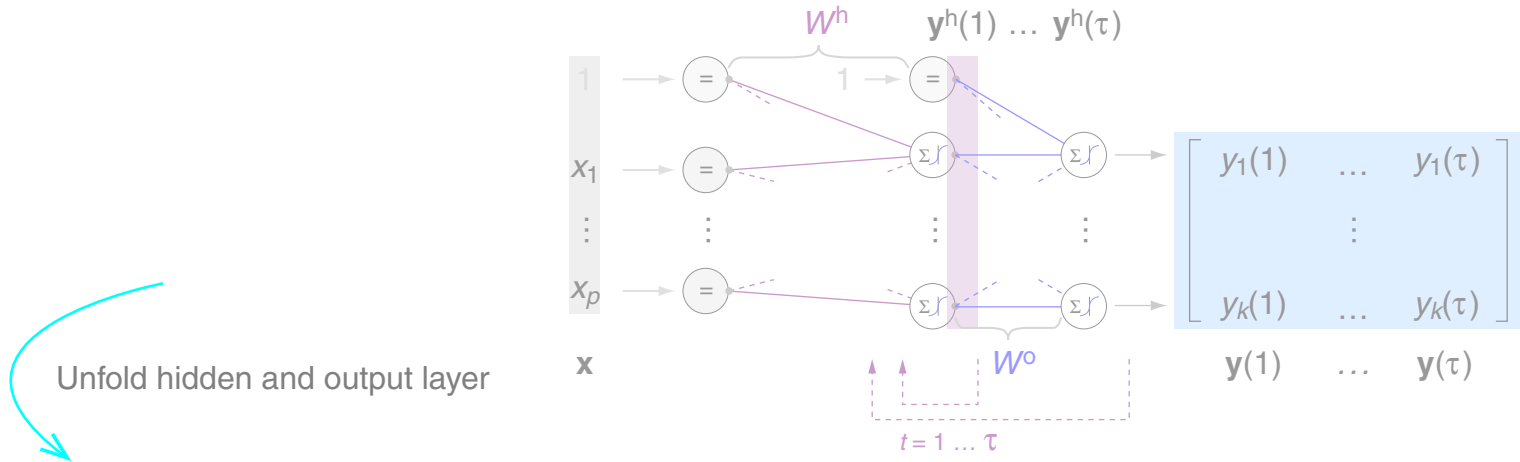
RNN Sequence Decoding (continued)



- ❑ One p -dimensional input vector \mathbf{x} .
- ❑ One hidden and one output layer, which are recurrently updated.
- ❑ Sequence of k -dimen. output vectors $[\mathbf{y}(\mathbf{x}, 1), \dots, \mathbf{y}(\mathbf{x}, \tau)]$ or $[\mathbf{y}(1), \dots, \mathbf{y}(\tau)]$.

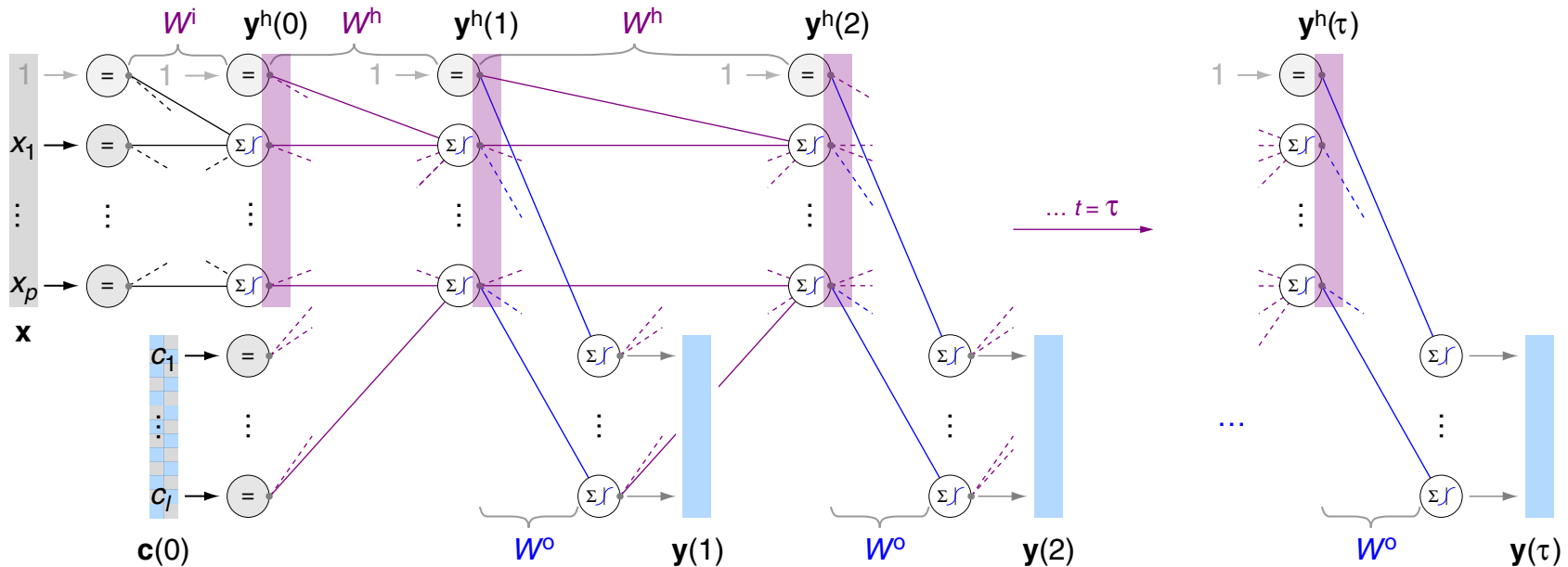
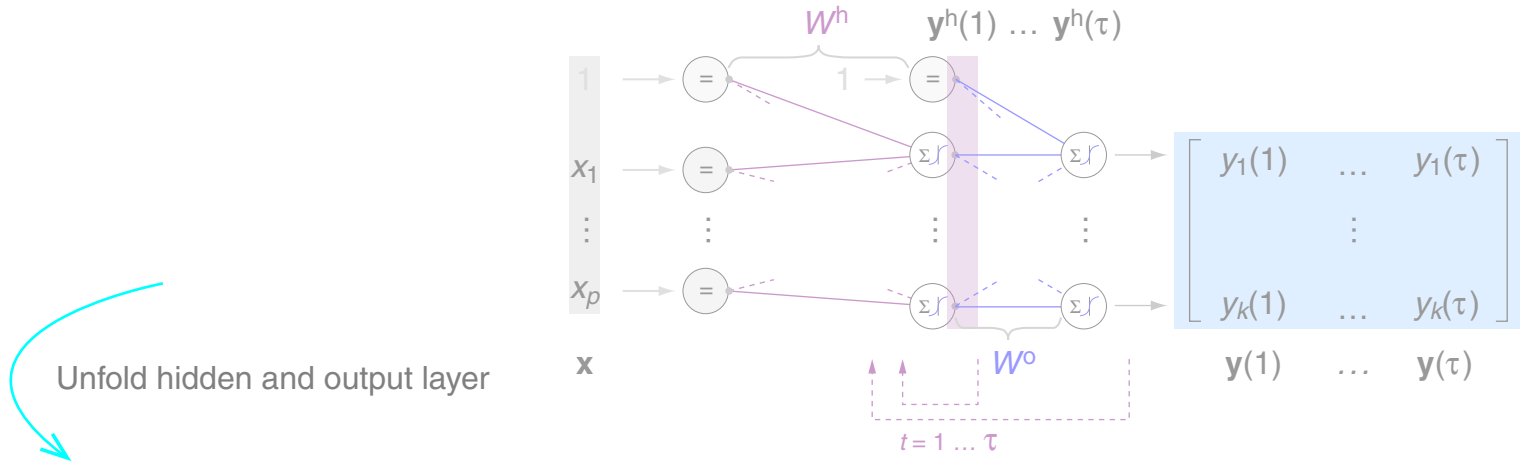
Recurrent Neural Networks

RNN Sequence Decoding (continued)



Recurrent Neural Networks

RNN Sequence Decoding (continued)



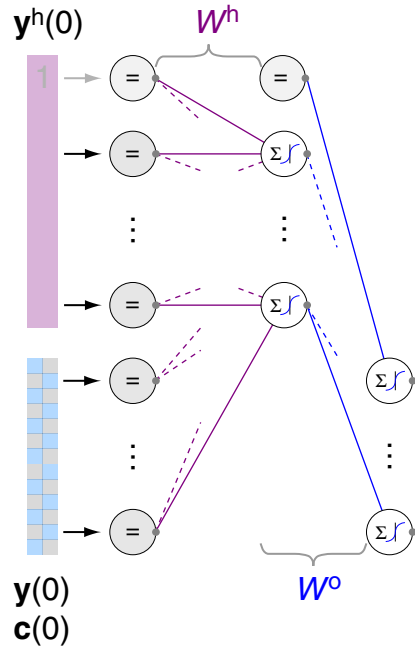
Remarks:

- ❑ An output sequence is written in brackets, $[\mathbf{y}(1), \dots, \mathbf{y}(\tau)]$, where $\mathbf{y}(t), t = 1, \dots, \tau$, denotes the output vector at time step t .
- ❑ The words in the output sequence are usually one-hot-encoded, i.e., by a k -dimensional output vector with a “1” whose position indicates the word, and zeros elsewhere.
- ❑ If the input, \mathbf{x} , is clear from the context, we usually note $\mathbf{y}(\mathbf{x}, t)$ as $\mathbf{y}(t)$.
- ❑ The matrix W^i is necessary to embed the typically low-dimensional input vector \mathbf{x} regarding the high-dimensional hidden vectors \mathbf{y}^h : $\mathbf{y}^h(0) = \sigma(W^i \mathbf{x})$.
- ❑ The parameter τ in $\mathbf{y}(\tau)$ is unknown. More specifically, the generation process terminates at that time step τ for which $\mathbf{y}(\tau) = (0, 0, \dots, 0, 1)^T (\hat{=} \text{<end>})$.
 τ (the length of the computed output) does not have to be equal to T (the length of the target or ground truth).

Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 0 \rightarrow 1$

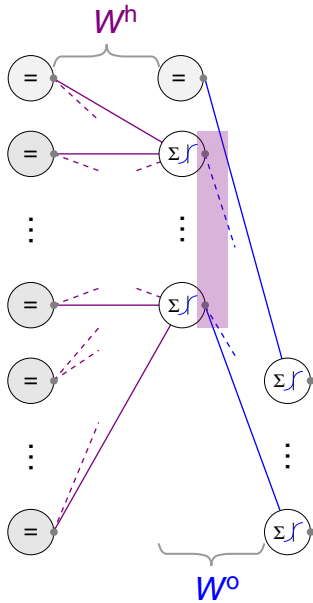


Output decoding over t .

Recurrent Neural Networks

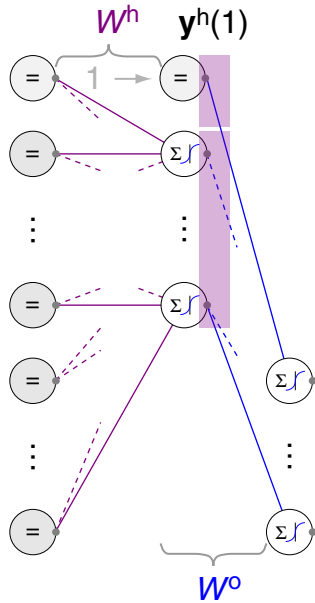
RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: 0 \rightarrow 1



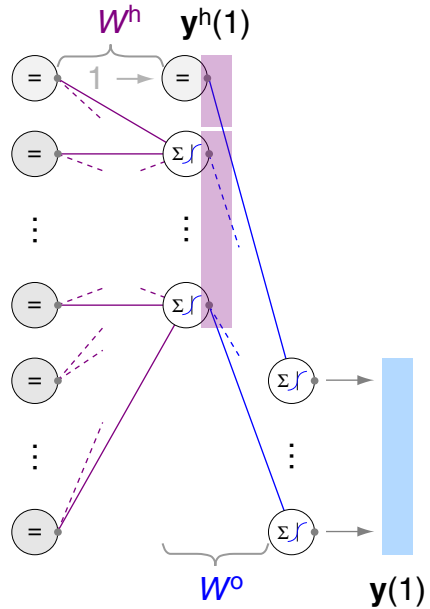
Output decoding over t .

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$$t: 0 \rightarrow 1$$


Output decoding over t .

Recurrent Neural Networks

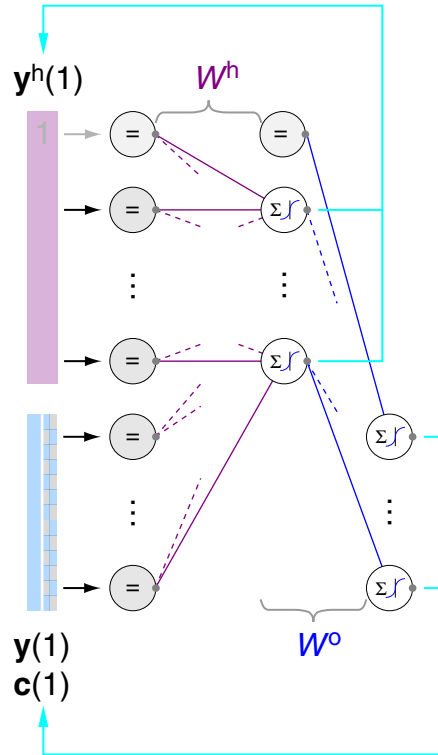
$$t: 0 \rightarrow 1$$


Output decoding over t .

Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 0 \rightarrow 1$

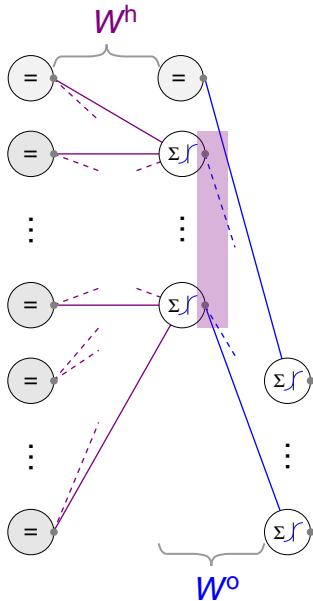


Output decoding over t .

Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: $1 \rightarrow 2$

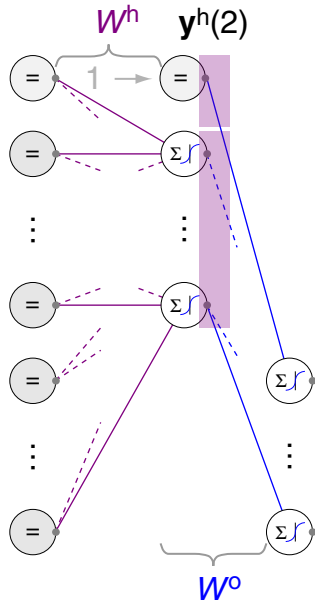


Output decoding over t .

Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: 1 \rightarrow 2



Output decoding over t .

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

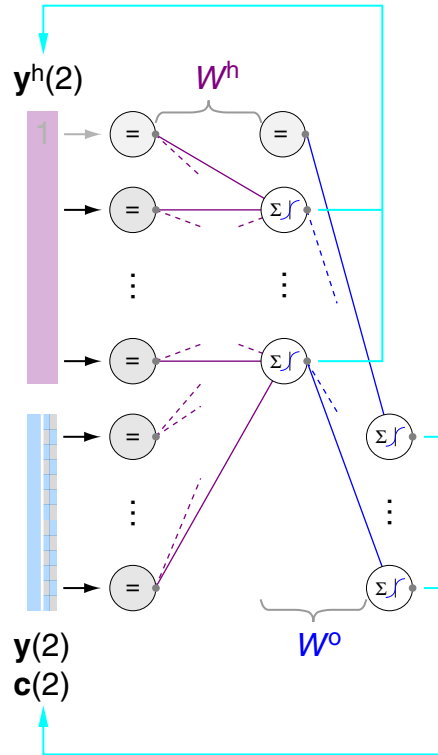
[illegible]

©STEIN/VÖLSKE 2024

Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 1 \rightarrow 2$

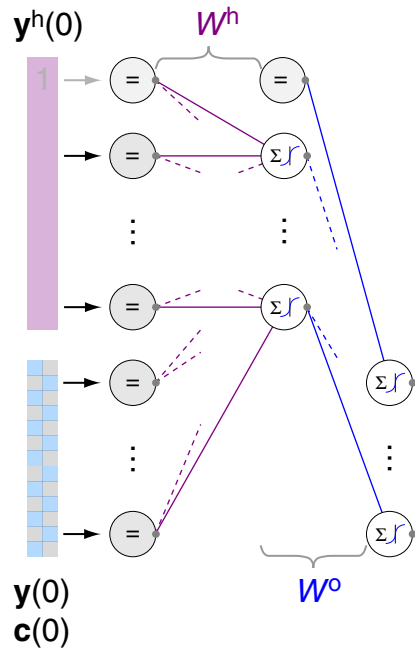


Output decoding over t .

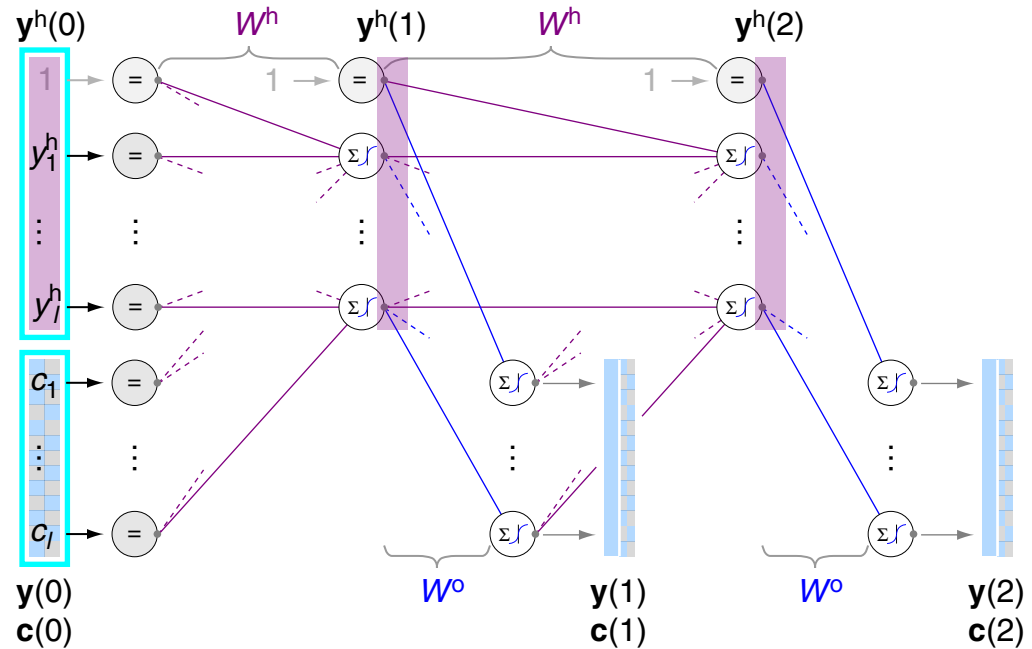
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 0 \rightarrow 1$



Output decoding over t .

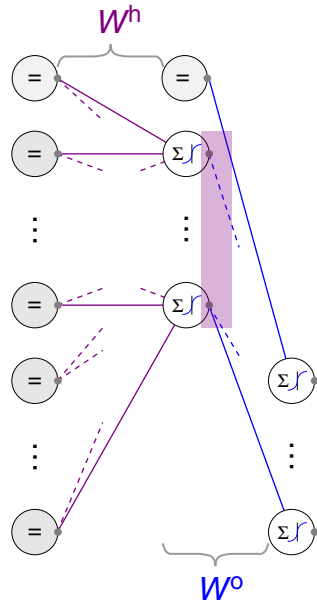


Hidden and output layer at subsequent time steps.

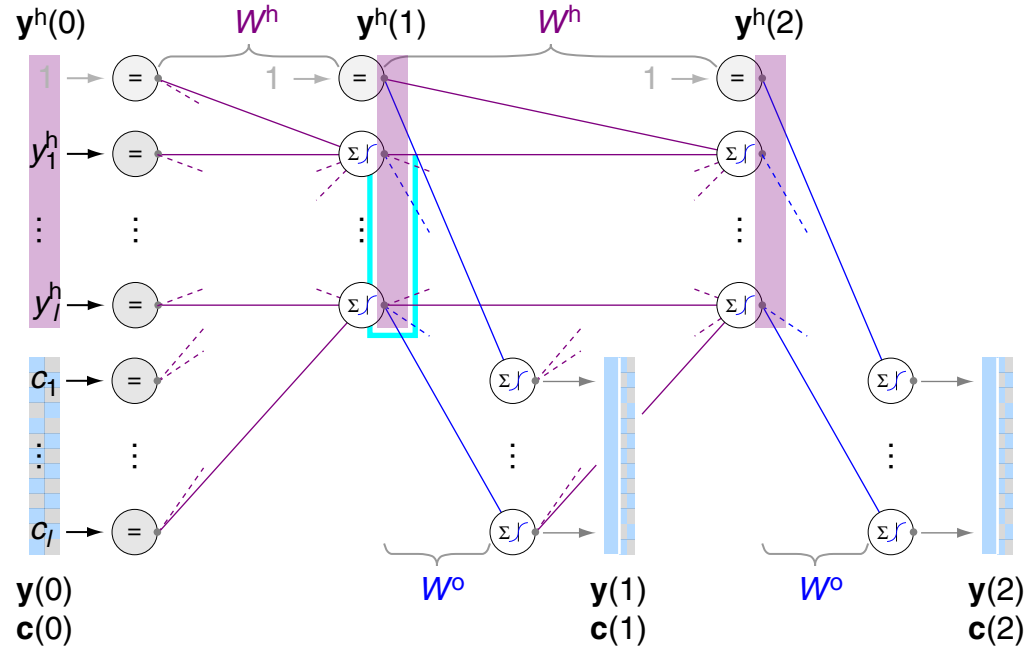
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 0 \rightarrow 1$



Output decoding over t .

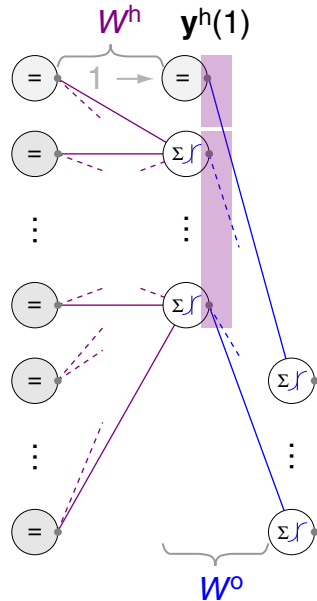


Hidden and output layer at subsequent time steps.

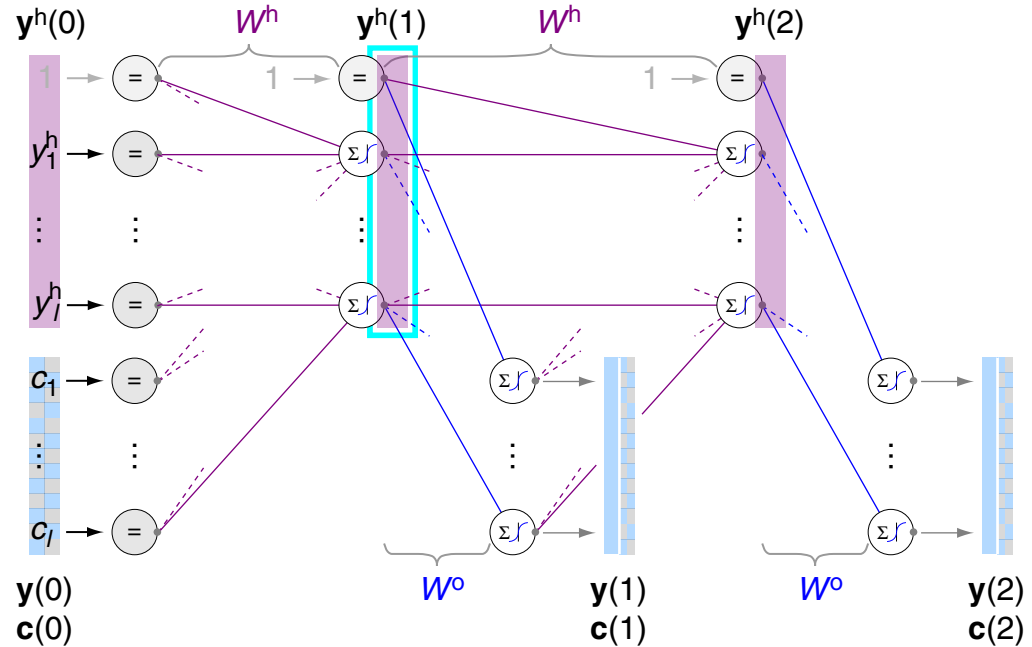
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 0 \rightarrow 1$



Output decoding over t .

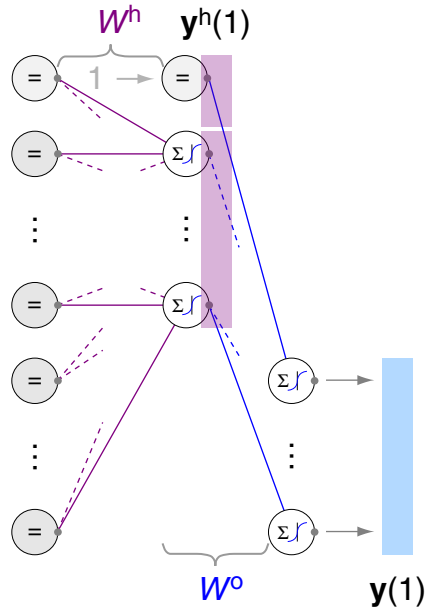


Hidden and output layer at subsequent time steps.

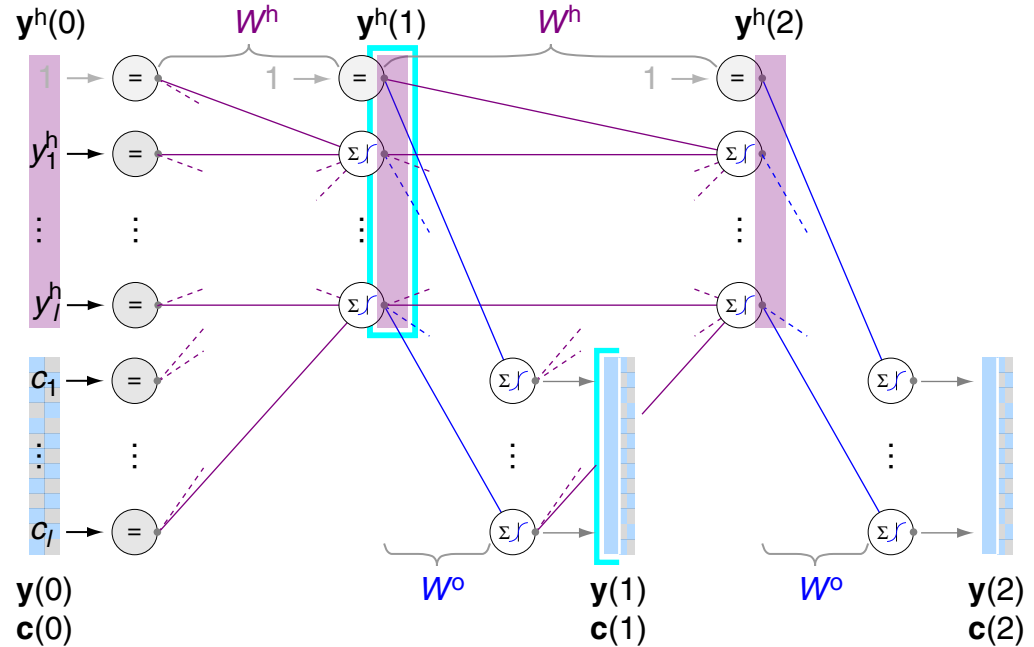
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: 0 \rightarrow 1



Output decoding over t .

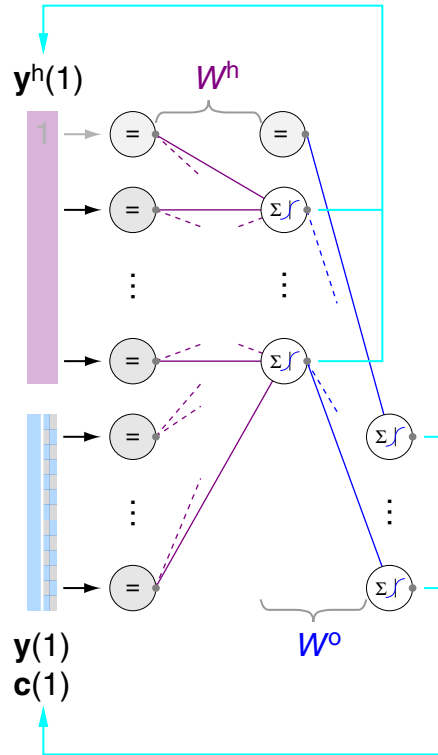


Hidden and output layer at subsequent time steps.

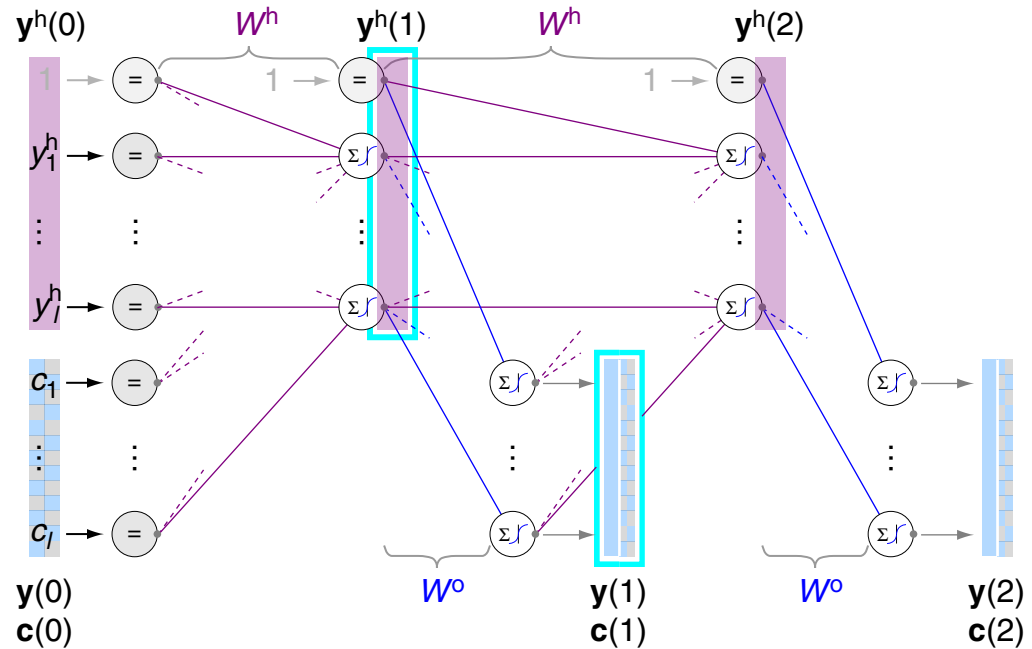
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: 0 \rightarrow 1



Output decoding over t .

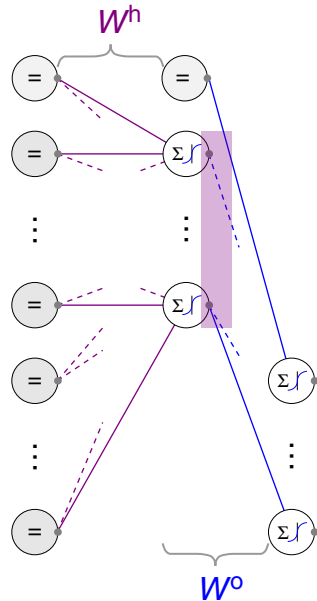


Hidden and output layer at subsequent time steps.

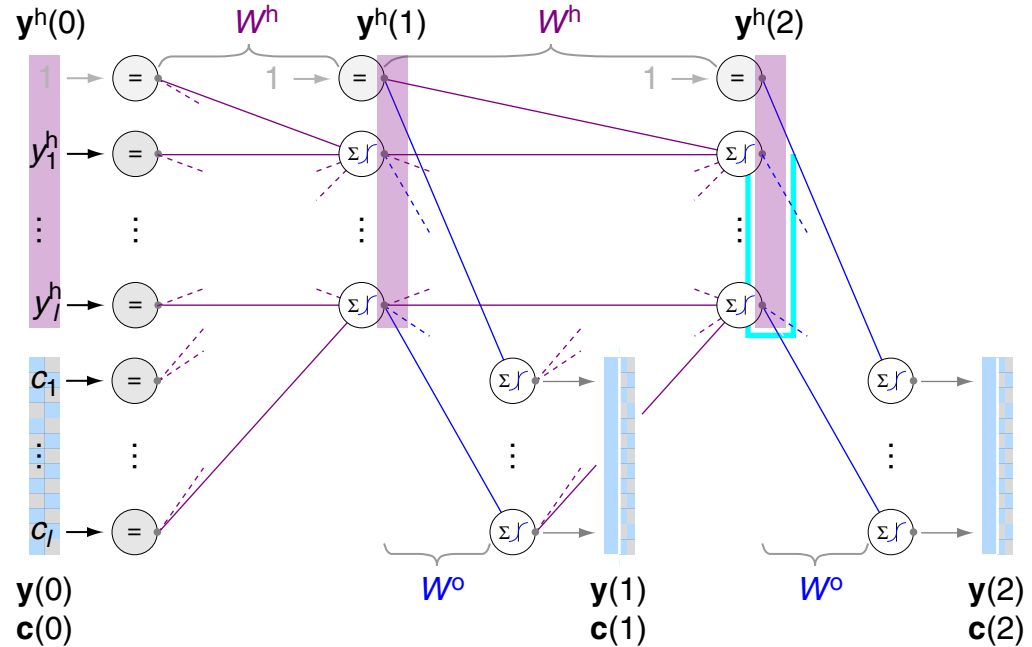
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: $1 \rightarrow 2$



Output decoding over t .

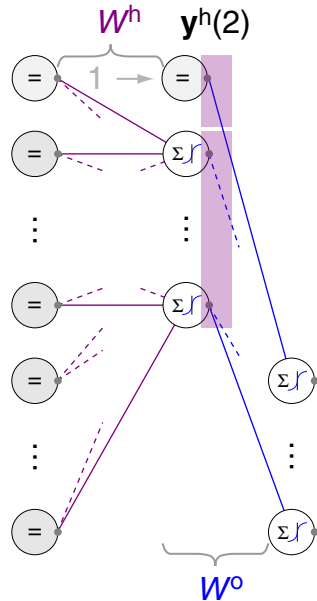


Hidden and output layer at subsequent time steps.

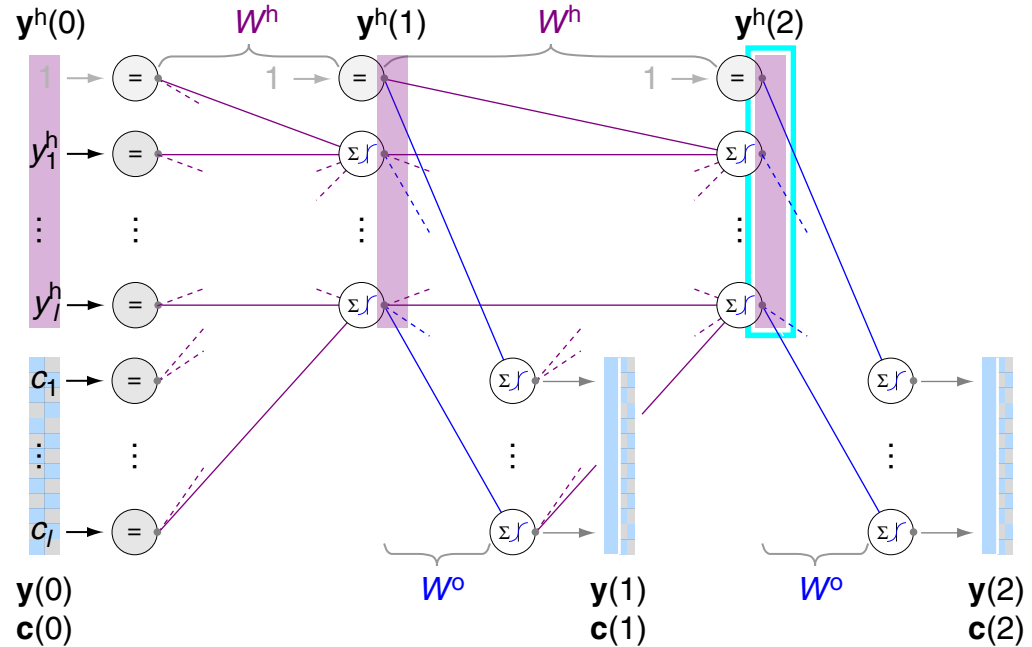
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: $1 \rightarrow 2$



Output decoding over t .

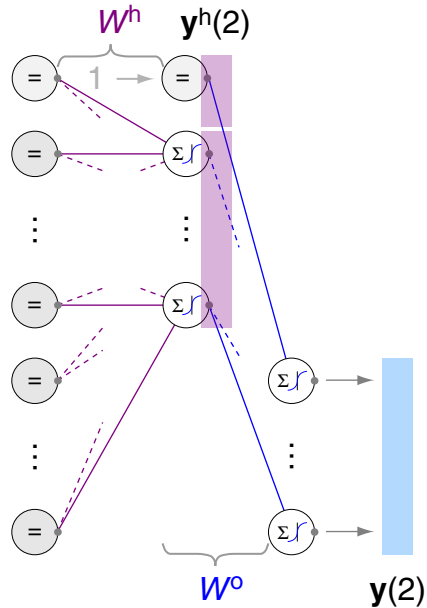


Hidden and output layer at subsequent time steps.

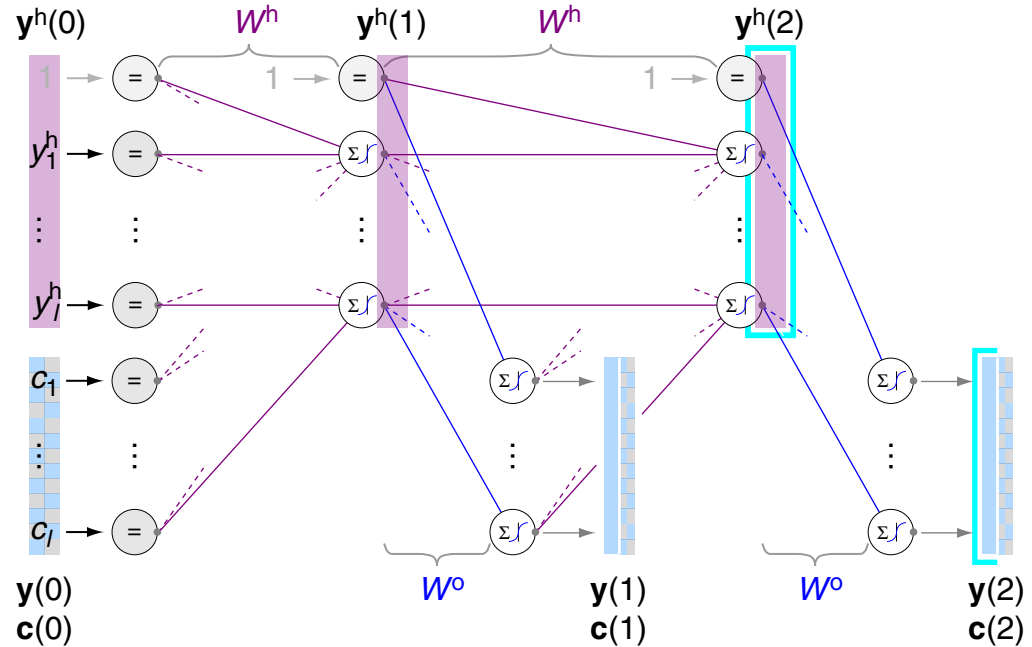
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

t: 1 \rightarrow 2



Output decoding over t .

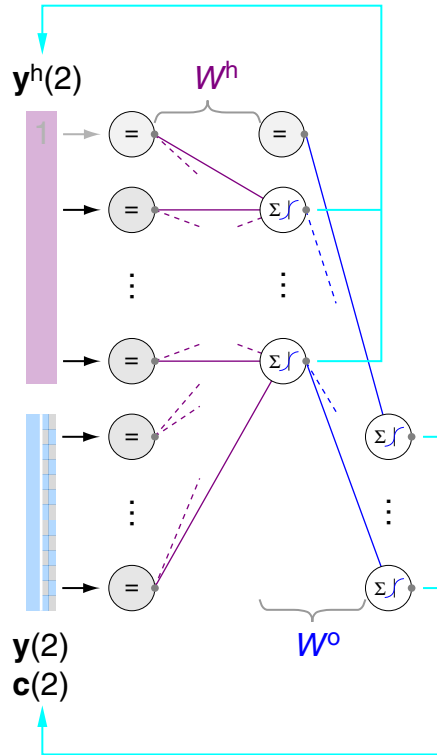


Hidden and output layer at subsequent time steps.

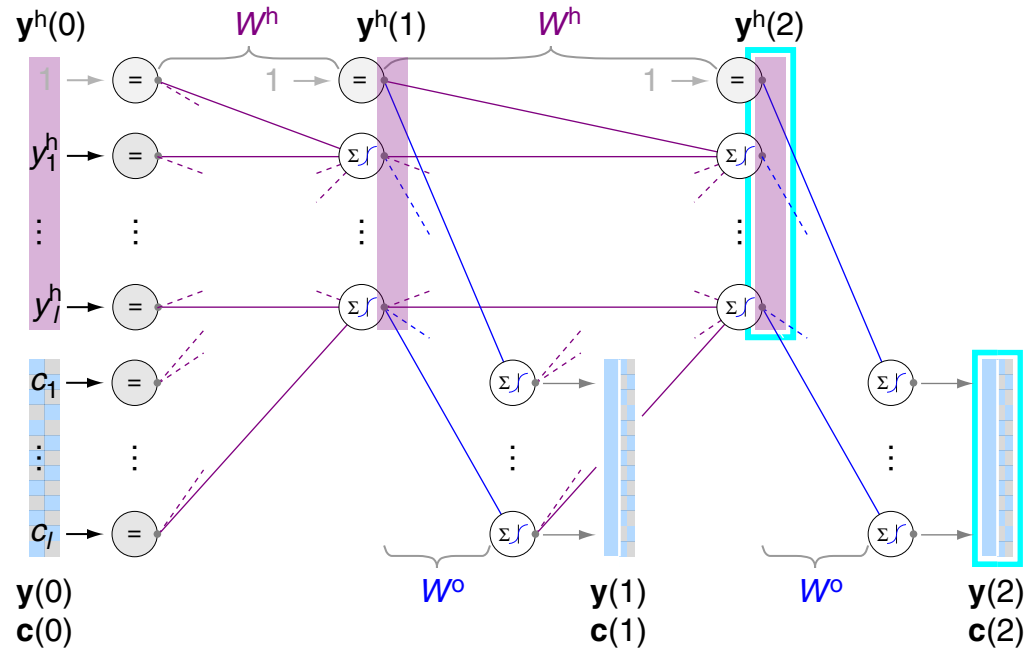
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 1 \rightarrow 2$



Output decoding over t .

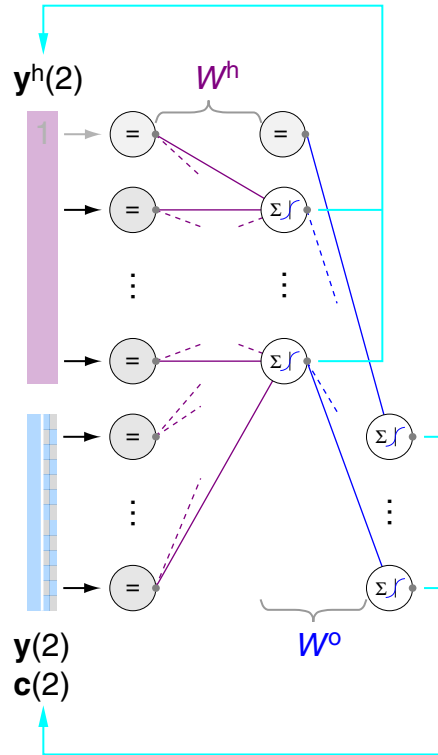


Hidden and output layer at subsequent time steps.

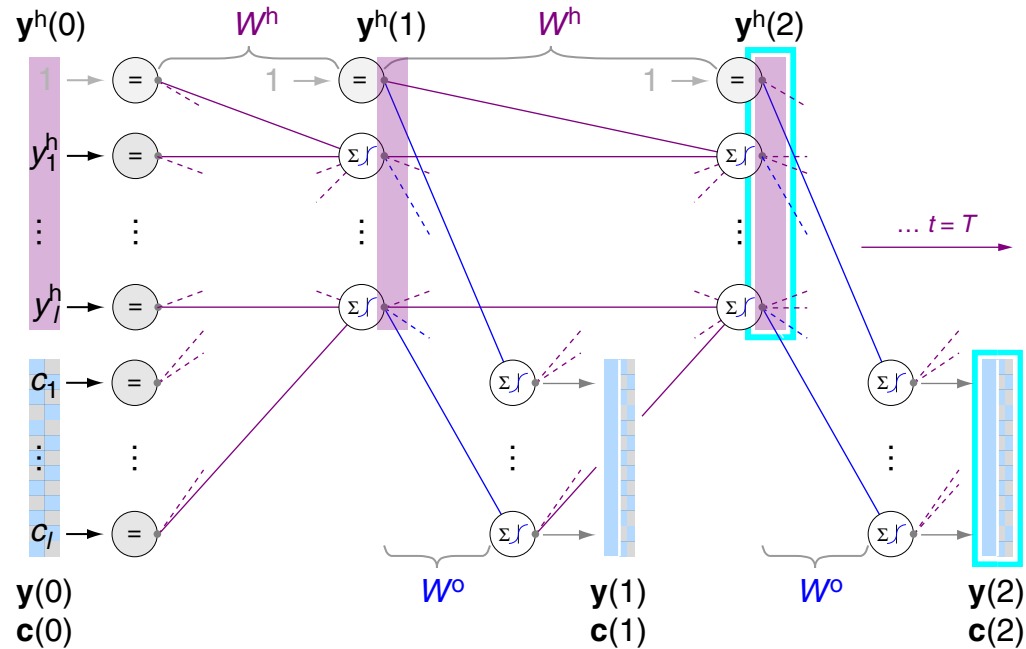
Recurrent Neural Networks

RNN Sequence Decoding (continued) [\[decoding overview\]](#)

$t: 1 \rightarrow 2$



Output decoding over t .



Hidden and output layer at subsequent time steps.

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- ⊕ → I love my cat.
- ⊕ → Cats and dogs lap water.
- ⊖ → It is raining cats and dogs.
- ⊖ → Cats and dogs are not allowed.
- ⊖ → Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water <start> <end>)

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- $\oplus \rightarrow$ I love my cat.
- $\oplus \rightarrow$ Cats and dogs lap water.
- $\ominus \rightarrow$ It is raining cats and dogs.
- $\ominus \rightarrow$ Cats and dogs are not allowed.
- $\ominus \rightarrow$ Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water <start> <end>)

Input: $[[[[\mathbf{x}, \mathbf{y}(0)], \mathbf{y}(1)], \mathbf{y}(2)], \dots], \mathbf{y}(\tau-1)], \quad \mathbf{x} = \begin{pmatrix} \oplus \\ . \end{pmatrix}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(5) \hat{=} \text{<end>}$

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- \oplus \rightarrow I love my cat.
- \oplus \rightarrow Cats and dogs lap water.
- \ominus \rightarrow It is raining cats and dogs.
- \ominus \rightarrow Cats and dogs are not allowed.
- \ominus \rightarrow Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water **<start>** <end>)

Input: $[[[[\mathbf{x}, \mathbf{y}(0)], \mathbf{y}(1)], \mathbf{y}(2)], \dots], \mathbf{y}(\tau-1)], \quad \mathbf{x} = \begin{pmatrix} \oplus \\ . \end{pmatrix}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(5) \hat{=} \text{<end>}$

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- \oplus \rightarrow I love my cat.
- \oplus \rightarrow Cats and dogs lap water.
- \ominus \rightarrow It is raining cats and dogs.
- \ominus \rightarrow Cats and dogs are not allowed.
- \ominus \rightarrow Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water **<start>** <end>)

Input:
$$\left[\left[\left[\left[\mathbf{x}, \mathbf{y}(0) \right], \mathbf{y}(1) \right], \mathbf{y}(2), \dots, \mathbf{y}(\tau-1) \right], \quad \mathbf{x} = \begin{pmatrix} \oplus \\ . \end{pmatrix} \right]$$

Output:
$$\left[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau) \right], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(5) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- \oplus \rightarrow I love my cat.
- \oplus \rightarrow Cats and dogs lap water.
- \ominus \rightarrow It is raining cats and dogs.
- \ominus \rightarrow Cats and dogs are not allowed.
- \ominus \rightarrow Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water **<start>** <end>)

Input:
$$[[[[[\mathbf{x}, \mathbf{y}(0)], \mathbf{y}(1)], \mathbf{y}(2)], \dots], \mathbf{y}(\tau-1)], \quad \mathbf{x} = \begin{pmatrix} \oplus \\ . \end{pmatrix}$$

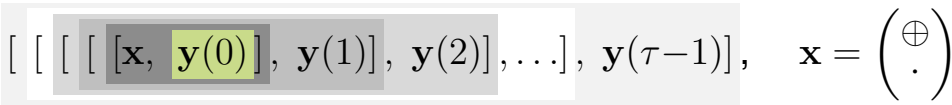
Output:
$$[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(5) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- \oplus \rightarrow I love my cat.
- \oplus \rightarrow Cats and dogs lap water.
- \ominus \rightarrow It is raining cats and dogs.
- \ominus \rightarrow Cats and dogs are not allowed.
- \ominus \rightarrow Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is
it lap love my natural not raining water `<start>` `<end>`)

Input:  $\mathbf{x} = \begin{pmatrix} \oplus \\ \cdot \end{pmatrix}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(5) \hat{=} \text{<end>}$

Recurrent Neural Networks

(S2) Class-to-Sequence: Text Generation

- \oplus \rightarrow I love my cat.
- \oplus \rightarrow Cats and dogs lap water.
- \ominus \rightarrow It is raining cats and dogs.
- \ominus \rightarrow Cats and dogs are not allowed.
- \ominus \rightarrow Cats and dogs have always been natural enemies.

Vocabulary: (allowed always and are been cat cats dogs enemies have i is it lap love my natural not raining water `<start>` `<end>`)

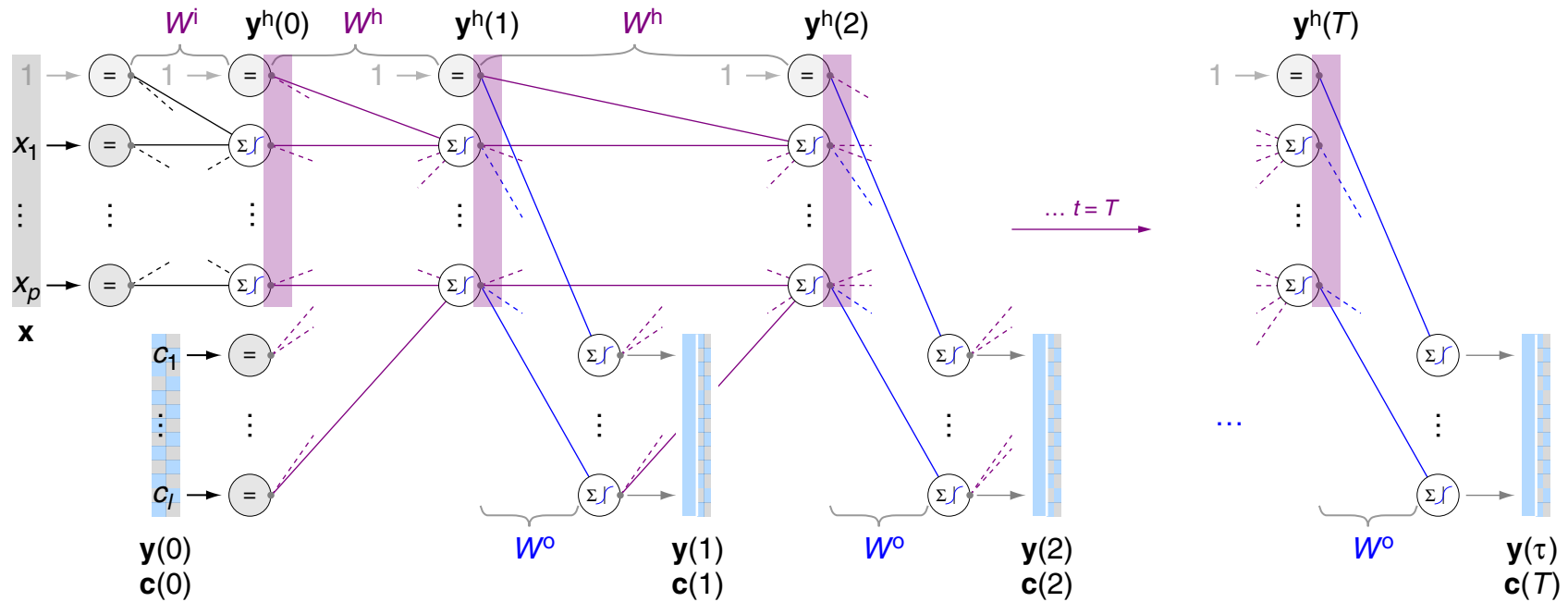
Input: $[[[[[\mathbf{x}, \mathbf{y}(0)], \mathbf{y}(1)], \mathbf{y}(2)], \dots], \mathbf{y}(\tau-1)], \quad \mathbf{x} = \begin{pmatrix} \oplus \\ . \end{pmatrix}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(5) \hat{=} \text{<end>}$

Target:
$$[\mathbf{c}(1), \dots, \mathbf{c}(5)] = \left[\begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \right]$$
$$\hat{=} [\text{word}_{11}, \text{word}_{15}, \text{word}_{16}, \text{word}_6, \text{word}_{22}]$$
$$\hat{=} \text{I love my cat}$$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(\tau-1)]$$

Output:

$$\mathbf{y}(t) = \sigma_{\Delta} (W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$$

Hidden:

$$\mathbf{y}^h(0) = \sigma (W^i \mathbf{x})$$

$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix} \right), t = 1, \dots, \tau$$

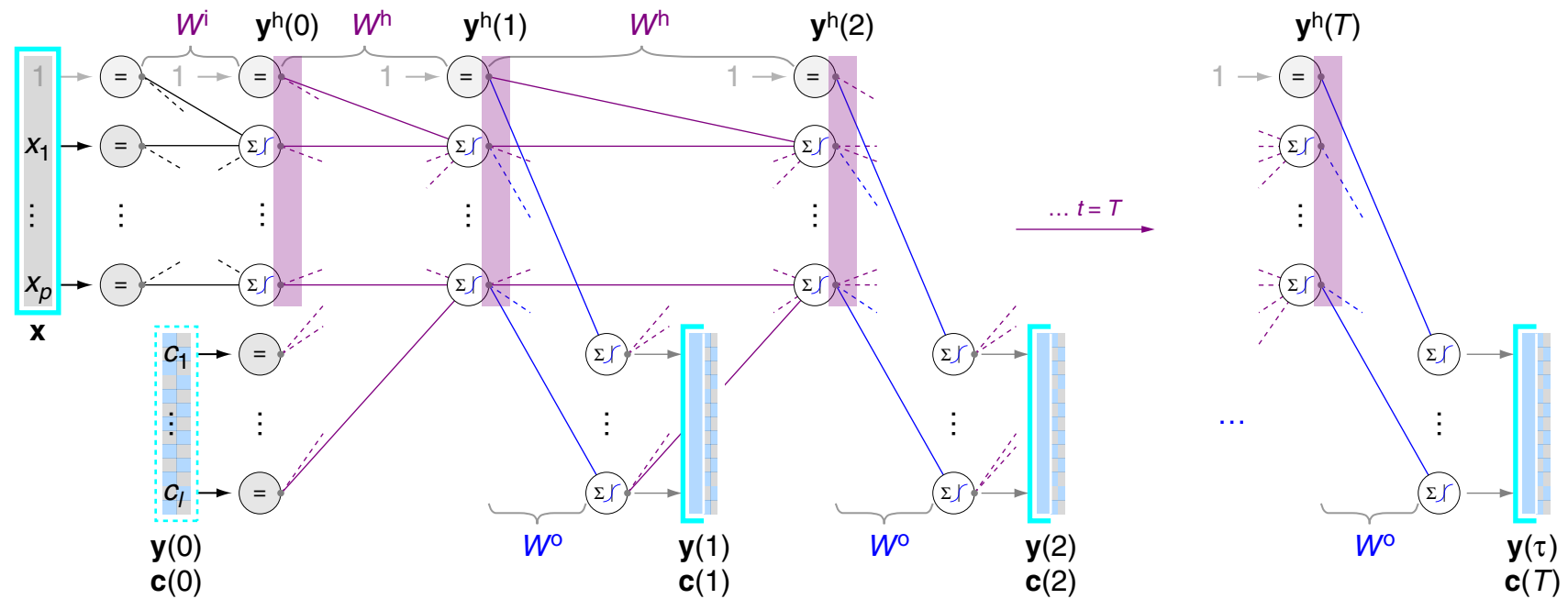
Target:

$$[\mathbf{c}(1), \dots, \mathbf{c}(T)]$$

$$\mathbf{c}(T) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(\tau-1)]$

Output:

$\mathbf{y}(t) = \sigma_\Delta(W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$

Hidden:

$\mathbf{y}^h(0) = \sigma(W^i \mathbf{x})$

$\mathbf{y}^h(t) = \sigma(W^h(\mathbf{y}^h(t-1))), t = 1, \dots, \tau$

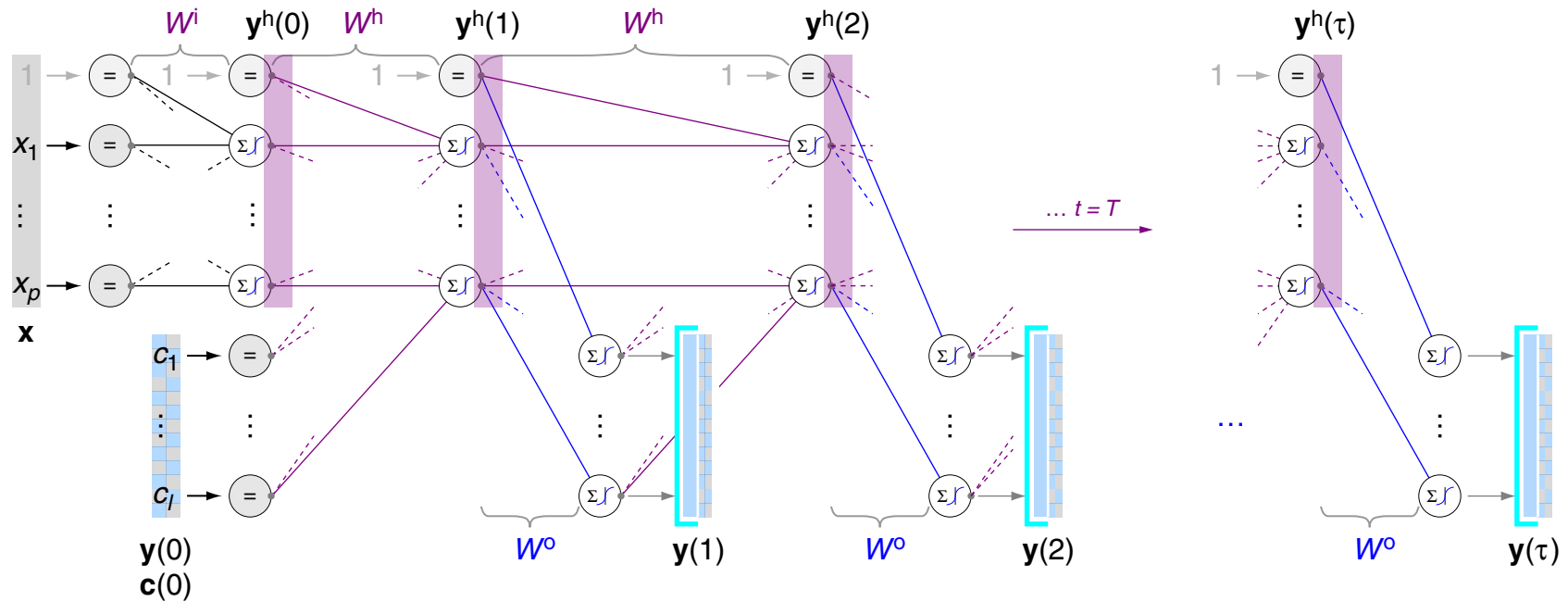
Target:

$[\mathbf{c}(1), \dots, \mathbf{c}(T)]$

$\mathbf{c}(T) \hat{=} \text{<end>}$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$$x, [y(1), \dots, y(\tau-1)]$$

Output:

$$y(t) = \sigma_{\Delta} (W^o y^h(t)), t = 1, \dots, \tau$$

Hidden:

$$y^h(0) = \sigma (W^i x)$$

$$y^h(t) = \sigma \left(W^h \begin{pmatrix} y^h(t-1) \\ y(t-1) \end{pmatrix} \right), t = 1, \dots, \tau$$

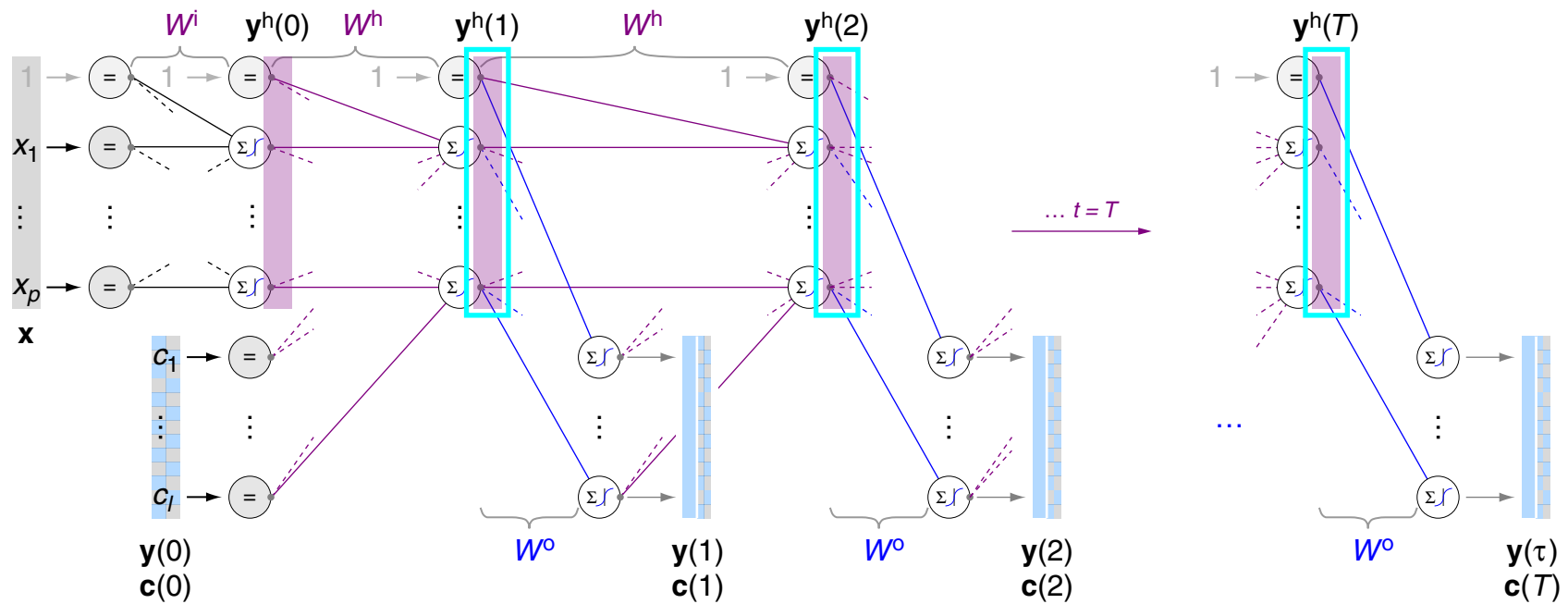
Target:

$$[c(1), \dots, c(T)]$$

$$c(T) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(\tau-1)]$$

Output:

$$\mathbf{y}(t) = \sigma_{\Delta} \left(W^o \mathbf{y}^h(t) \right), t = 1, \dots, \tau$$

Hidden:

$$\mathbf{y}^h(0) = \sigma \left(W^i \mathbf{x} \right)$$

$$\mathbf{y}^h(t) = \sigma \left(W^h \left(\mathbf{y}^h(t-1) \right) \right), t = 1, \dots, T$$

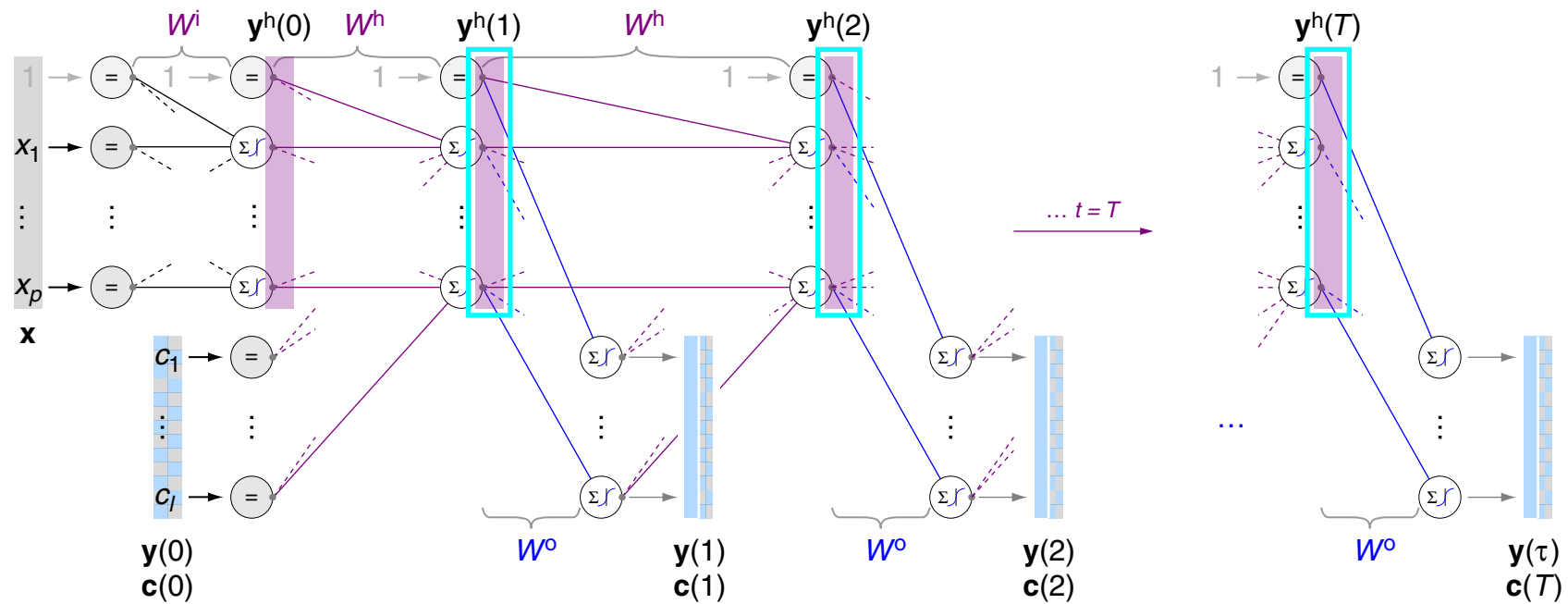
Target:

$$[\mathbf{c}(1), \dots, \mathbf{c}(T)]$$

$$\mathbf{c}(T) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(\tau-1)]$$

Output:

$$\mathbf{y}(t) = \sigma_{\Delta} (W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$$

Hidden:

$$\mathbf{y}^h(0) = \sigma (W^i \mathbf{x})$$

$$\mathbf{y}^h(t) = \sigma (W^h (\mathbf{y}^h(t-1))) , t = 1, \dots, \tau$$

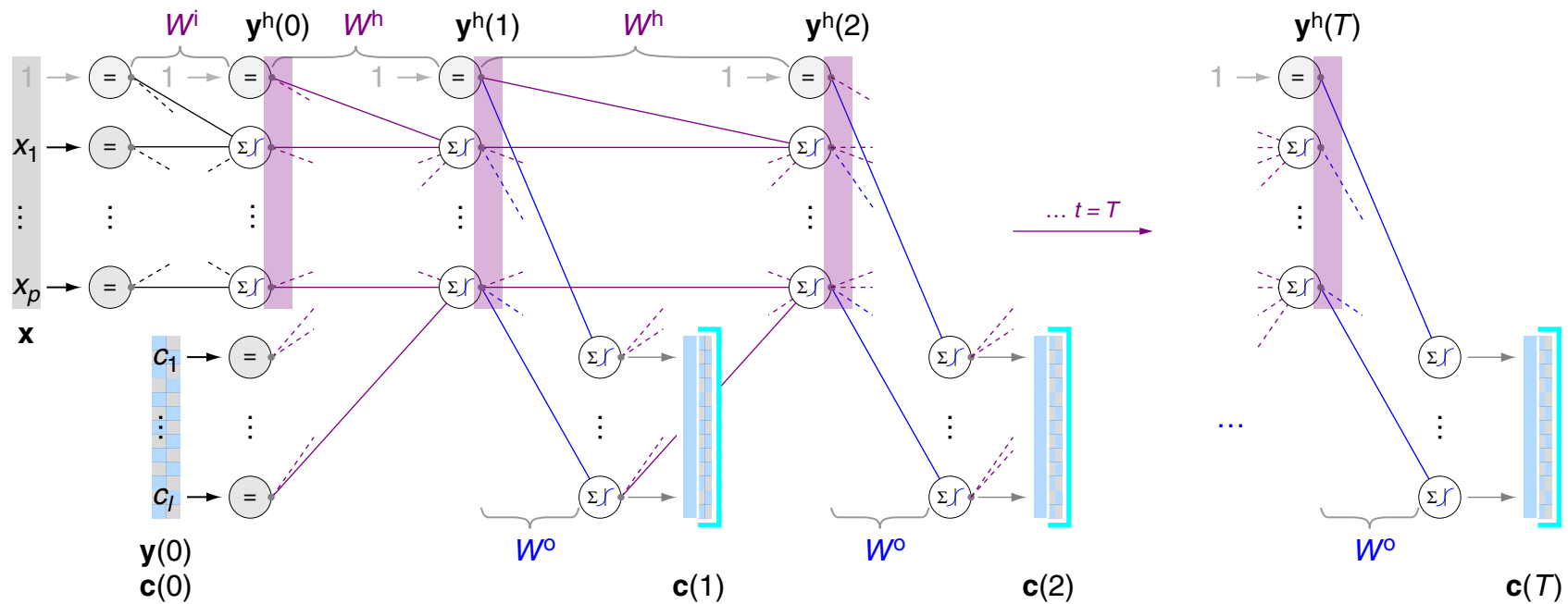
Target:

$$[\mathbf{c}(1), \dots, \mathbf{c}(T)]$$

$$\mathbf{c}(T) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(\tau-1)]$$

Output:

$$\mathbf{y}(t) = \sigma_{\Delta} \left(W^o \mathbf{y}^h(t) \right), t = 1, \dots, \tau$$

Hidden:

$$\mathbf{y}^h(0) = \sigma \left(W^i \mathbf{x} \right)$$

$$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix} \right), t = 1, \dots, \tau$$

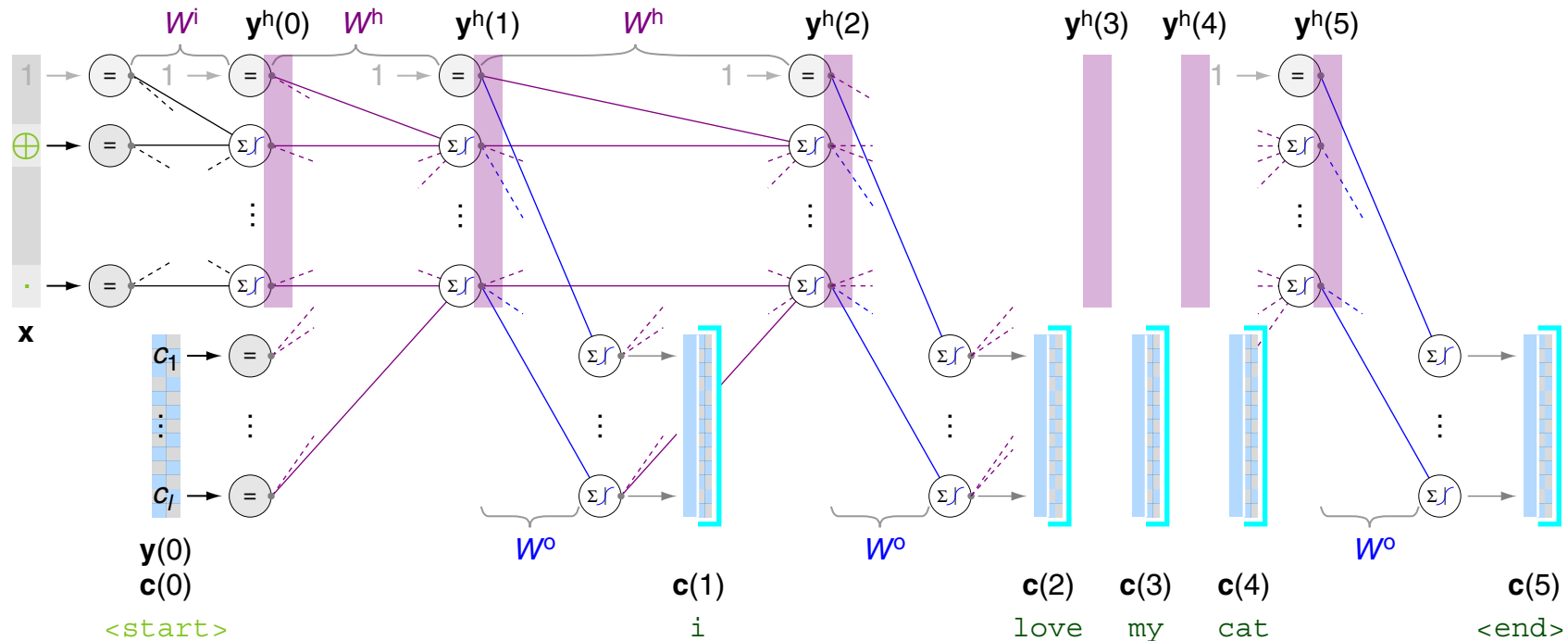
Target:

$$[\mathbf{c}(1), \dots, \mathbf{c}(T)]$$

$$\mathbf{c}(T) \hat{=} \text{<end>}$$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(4)]$

Output:

$\mathbf{y}(t) = \sigma_{\Delta} \left(W^o \mathbf{y}^h(t) \right), t = 1, \dots, 5$

Hidden:

$\mathbf{y}^h(0) = \sigma \left(W^i \mathbf{x} \right)$

$\mathbf{y}^h(t) = \sigma \left(W^h \left(\mathbf{y}^h(t-1) \right) \right), t = 1, \dots, 5$

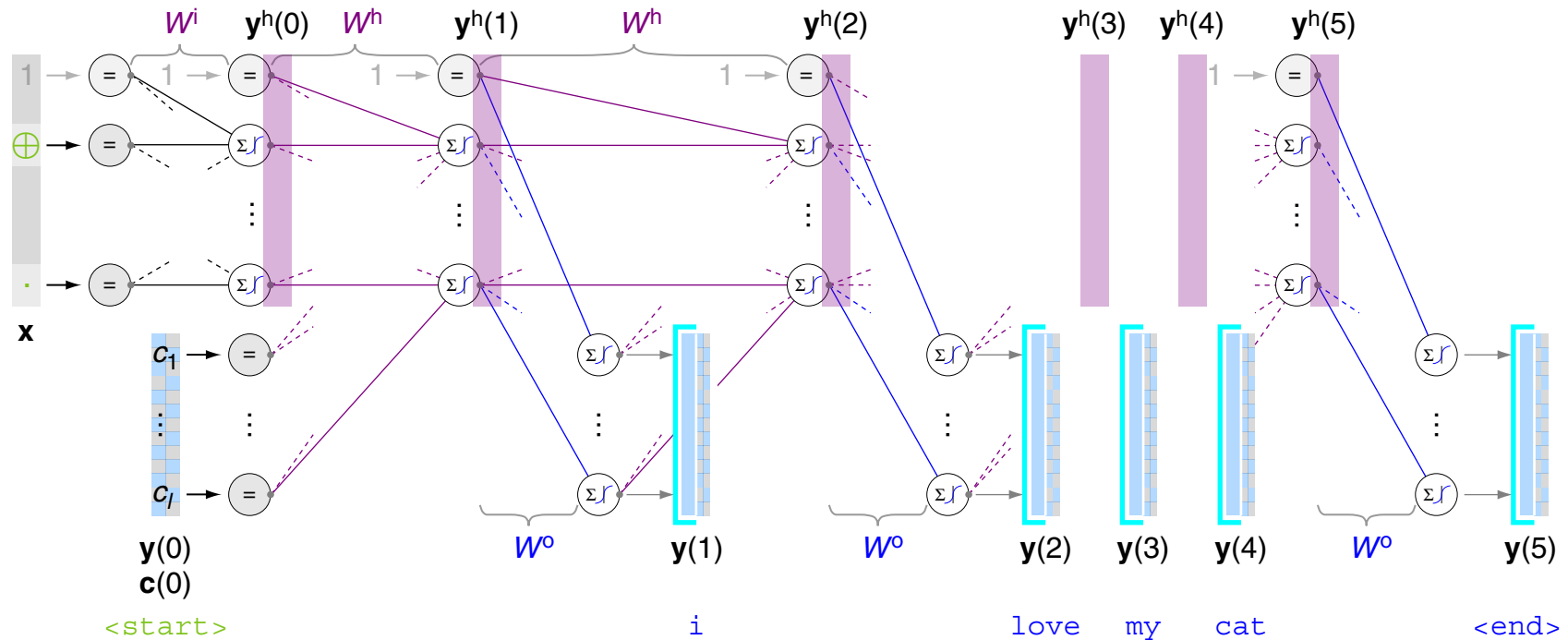
Target:

$[\mathbf{c}(1), \dots, \mathbf{c}(5)]$

$\mathbf{c}(5) \hat{=} \text{<end>}$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(4)]$

Output:

$\mathbf{y}(t) = \sigma_{\Delta} \left(W^o \mathbf{y}^h(t) \right), t = 1, \dots, 5$

Hidden:

$\mathbf{y}^h(0) = \sigma \left(W^i \mathbf{x} \right)$

$\mathbf{y}^h(t) = \sigma \left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix} \right), t = 1, \dots, 5$

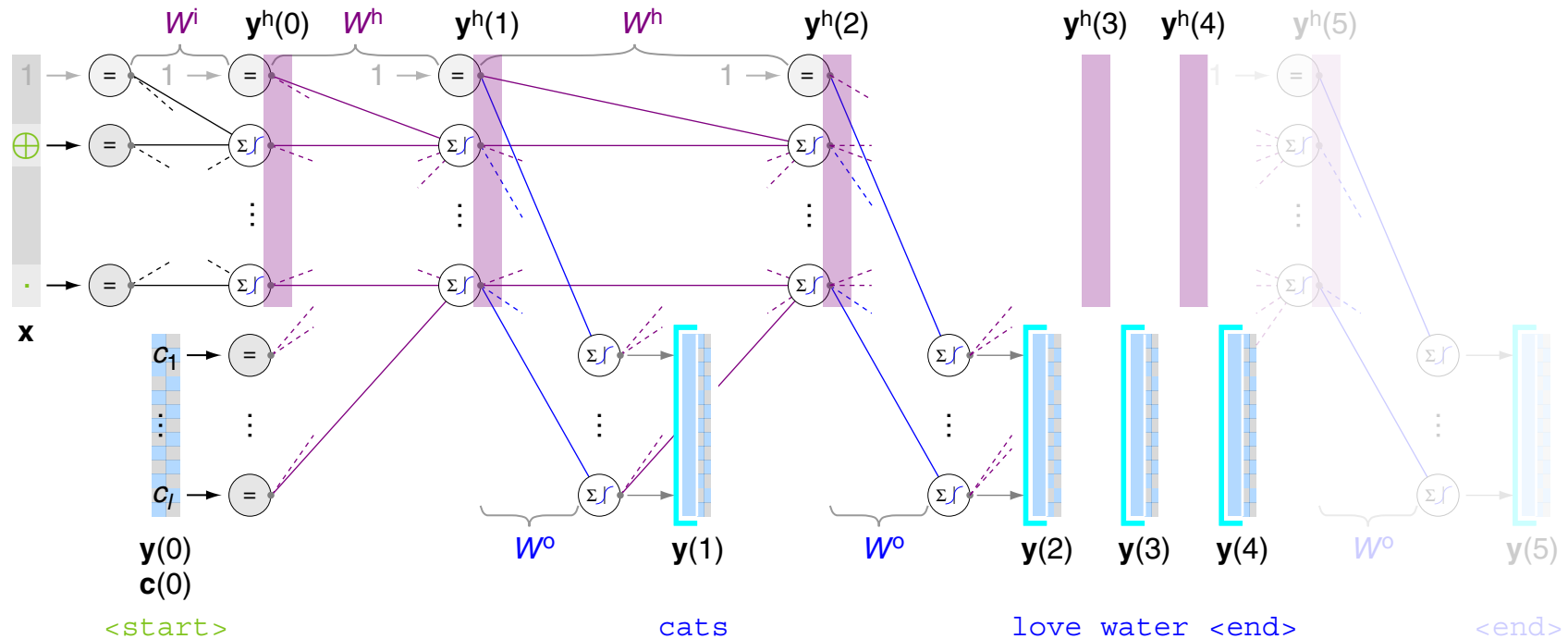
Target:

$[\mathbf{c}(1), \dots, \mathbf{c}(5)]$

$\mathbf{c}(5) \hat{=} \text{<end>}$

Recurrent Neural Networks

(S2) Class-to-Sequence Mapping with RNNs



Input:

$$\mathbf{x}, [\mathbf{y}(1), \dots, \mathbf{y}(3)]$$

Output:

$$\mathbf{y}(t) = \sigma_{\Delta} \left(W^o \mathbf{y}^h(t) \right), t = 1, \dots, 4$$

Hidden:

$$\mathbf{y}^h(0) = \sigma \left(W^i \mathbf{x} \right)$$

$$\mathbf{y}^h(t) = \sigma \left(W^h \left(\mathbf{y}^h(t-1) \right) \right), t = 1, \dots, 4$$

Target:

$$[\mathbf{c}(1), \dots, \mathbf{c}(5)]$$

$$\mathbf{c}(5) \hat{=} \text{<end>}$$

Remarks:

- We denote $\mathbf{y}(0)$ not as input since it is predefined and does not contain any “actual knowledge”. In particular, $\mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}$.

- At training time the calculation of $\mathbf{y}^h(t)$ usually considers the ground truth $\mathbf{c}(t-1)$:

$$\mathbf{y}^h(t) = \sigma \left(\mathbf{W}^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{c}(t-1) \end{pmatrix} \right)$$

- At test time (“production mode”) the calculation of $\mathbf{y}^h(t)$ has to consider the output $\mathbf{y}(t-1)$:

$$\mathbf{y}^h(t) = \sigma \left(\mathbf{W}^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix} \right)$$

Recurrent Neural Networks




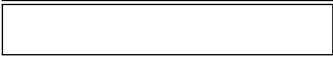
The IGD Algorithm for Class-to-Sequence Tasks [IGD_{seq2c}]

Algorithm: IGD_{c2seq} Incremental Gradient Descent for RNNs at class2seq tasks.

Input: D Multiset of examples $(\mathbf{x}, [\mathbf{c}(1), \dots, \mathbf{c}(T)])$ with $\mathbf{x} \in \{0, 1\}^p$, $\mathbf{c}(t) \in \mathbb{R}^k$.

η Learning rate, a small positive constant.

Output: W^i, W^h, W^o Weights matrices. (= hypothesis)

```
1. initialize_random_weights( $W^i, W^h, W^o$ ),  $t_{\text{epoch}} = 0$ 
2. REPEAT
3.    $t_{\text{epoch}} = t_{\text{epoch}} + 1$ 
4.   FOREACH  $(\mathbf{x}, [\mathbf{c}(1), \dots, \mathbf{c}(T)]) \in D$  DO
5.      Model function evaluation.
6.      Calculation of residual vectors.
7a.     Calculation of derivative of the loss.
7b.
8.      Parameter update  $\hat{=}$  one gradient step down.
9.   ENDDO
10.  UNTIL ( $\text{convergence}(D, \mathbf{y}(\cdot), t_{\text{epoch}})$ )
11.  return( $W^i, W^h, W^o$ )
```

Recurrent Neural Networks

The IGD Algorithm for Class-to-Sequence Tasks [IGD_{seq2c}]

Algorithm: IGD_{c2seq} Incremental Gradient Descent for RNNs at class2seq tasks.
Input: D Multiset of examples $(\mathbf{x}, [\mathbf{c}(1), \dots, \mathbf{c}(T)])$ with $\mathbf{x} \in \{0, 1\}^p$, $\mathbf{c}(t) \in \mathbb{R}^k$.
 η Learning rate, a small positive constant.
Output: W^i, W^h, W^o Weights matrices. (= hypothesis)

```
1. initialize_random_weights( $W^i, W^h, W^o$ ),  $t_{\text{epoch}} = 0$ 
2. REPEAT
3.    $t_{\text{epoch}} = t_{\text{epoch}} + 1$ 
4.   FOREACH  $(\mathbf{x}, [\mathbf{c}(1), \dots, \mathbf{c}(T)]) \in D$  DO
5.      $\mathbf{y}^h(0) = \sigma(W^i \mathbf{x})$ 
6.     FOR  $t = 1$  TO  $T$  DO // forward propagation
7.        $\mathbf{y}^h(t) = \sigma(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{c}(t-1) \end{pmatrix})$ ,  $\mathbf{y}(t) = \sigma_\Delta(W^o \mathbf{y}^h(t))$ 
8.     ENDDO
9.      $[\delta(1), \dots, \delta(T)] = [\mathbf{c}(1), \dots, \mathbf{c}(T)] \ominus [\mathbf{y}(1), \dots, \mathbf{y}(T)]$  // consider that  $T$  may  $\neq \tau$ 
10.     $\ell(\mathbf{w}) = \sum_t l(\delta(t)) + \frac{\lambda}{n} R(\mathbf{w})$ ,  $\nabla \ell(\mathbf{w}) = \text{autodiff}(\ell(), \mathbf{w})$  // backprop. (7a+7b)
11.     $\Delta W^i = \eta \cdot \nabla^i \ell(\mathbf{w})$ ,  $\Delta W^h = \eta \cdot \nabla^h \ell(\mathbf{w})$ ,  $\Delta W^o = \eta \cdot \nabla^o \ell(\mathbf{w})$ 
12.     $W^i = W^i + \Delta W^i$ ,  $W^h = W^h + \Delta W^h$ ,  $W^o = W^o + \Delta W^o$ 
13.  ENDDO
14. UNTIL (convergence( $D, \mathbf{y}(\cdot), t_{\text{epoch}}$ ))
15. return( $W^i, W^h, W^o$ )
```


Remarks:

- We use the operator $\gg\ominus\ll$ to compare the two sequences $[c(t)]$ and $[y(t)]$ of possibly different length. Note that the concrete semantics of $\gg\ominus\ll$ is left open here.