

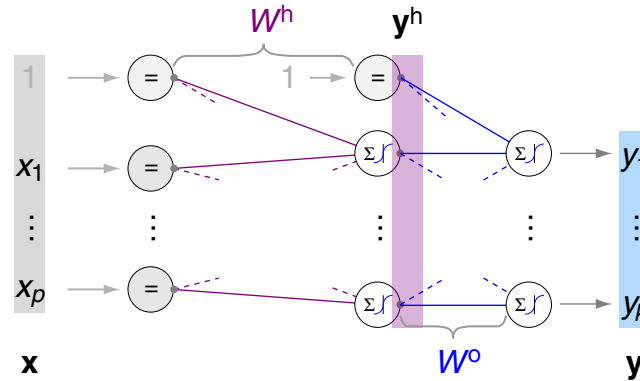
Chapter ML:IX

IX. Deep Learning for Sequence Processing

- Introduction to Deep Learning
- RNN Sequence Encoding
- RNN Sequence Decoding
- RNNs for Machine Translation
- Long-Term Dependencies
- Attention Mechanism
- Transformer Architecture
- Transformer Advanced
- Language Models
- Language Model Training

RNN Sequence Decoding

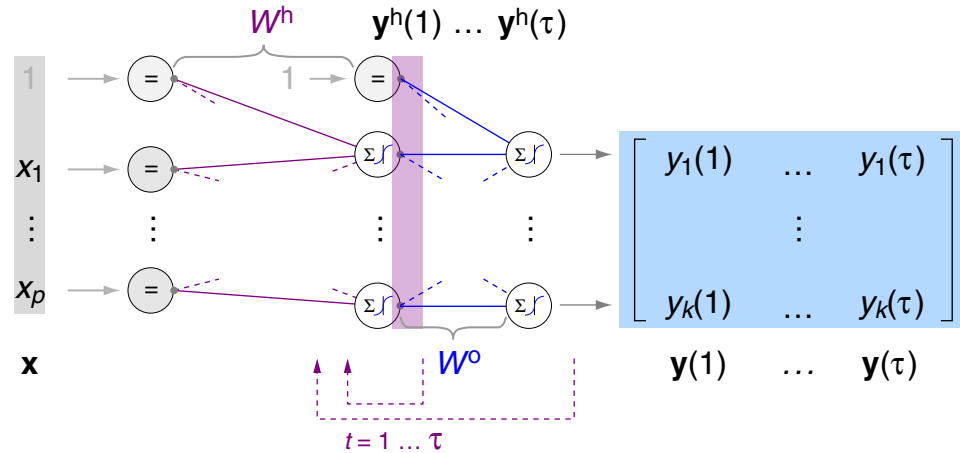
From Vector Output to Sequence Output



- 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 $\mathbf{y}(\mathbf{x})$.

RNN Sequence Decoding

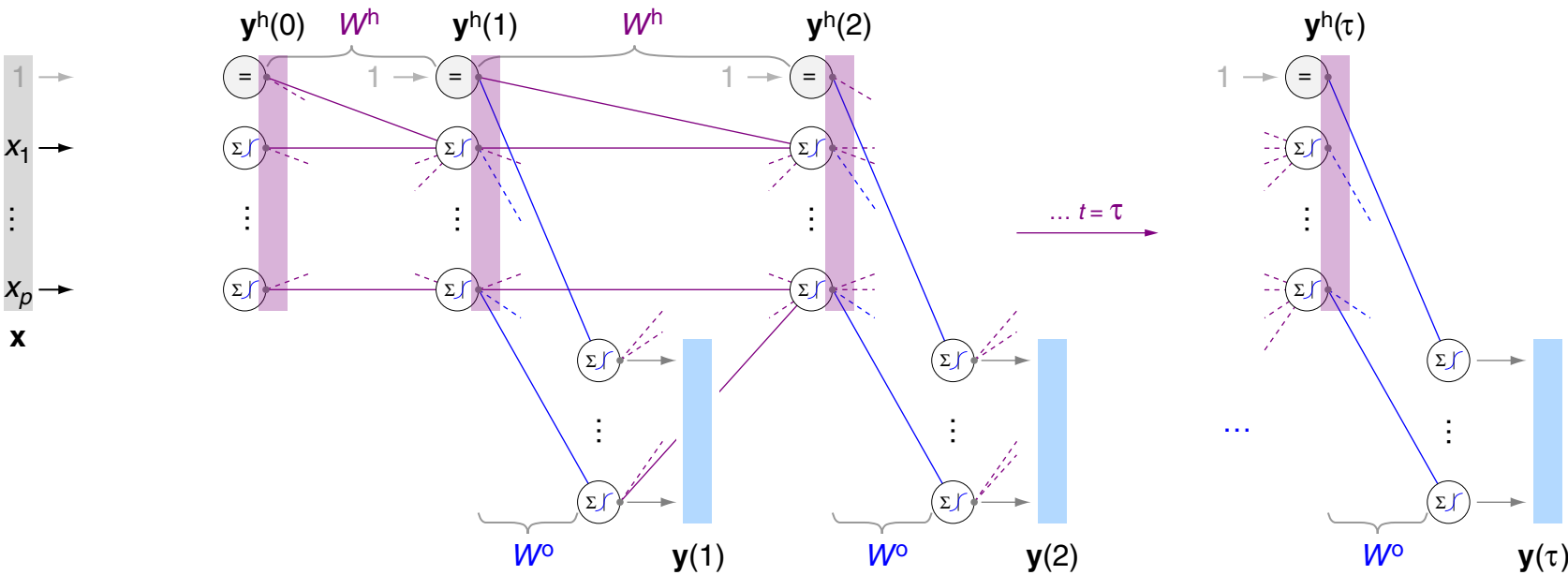
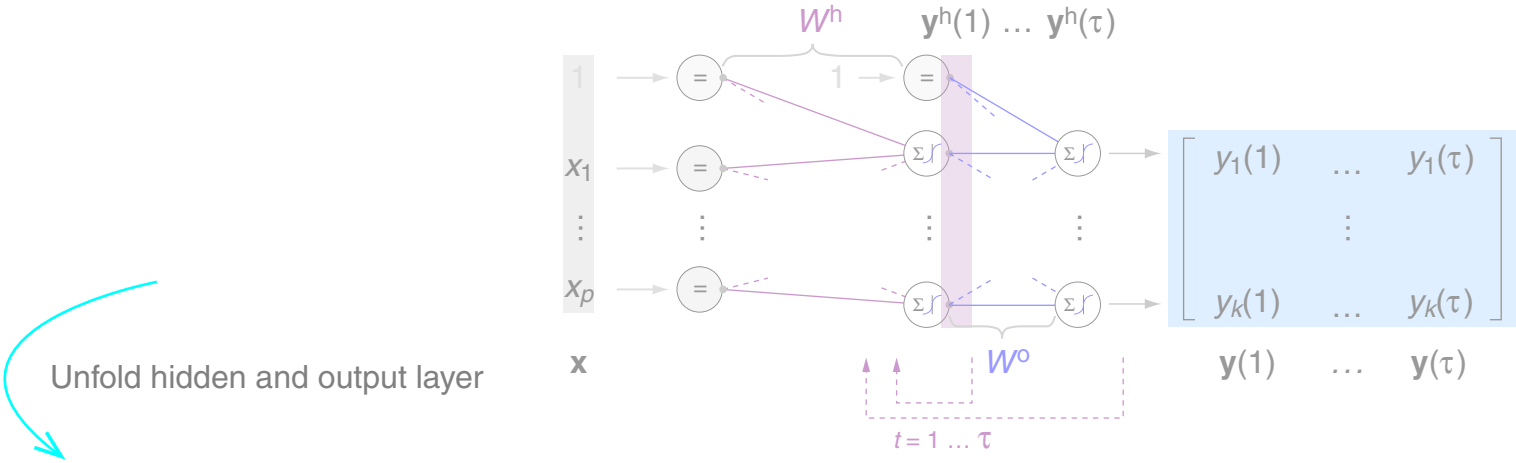
From Vector Output to Sequence Output



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

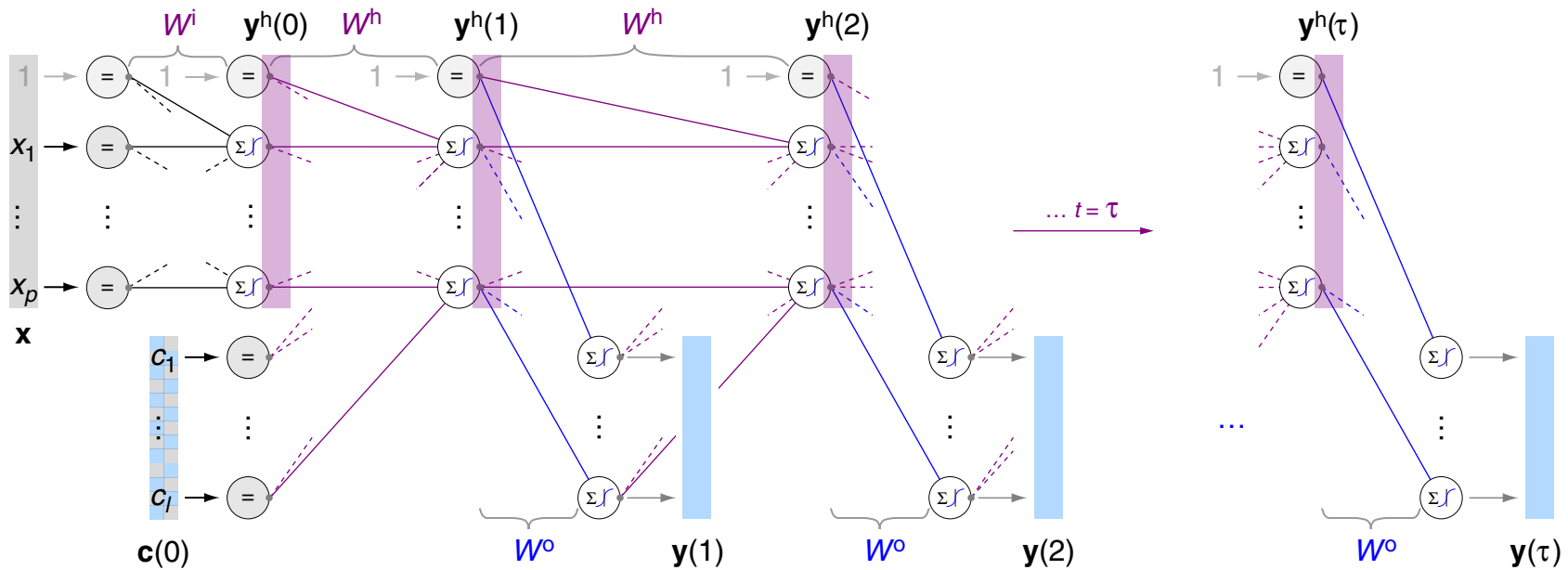
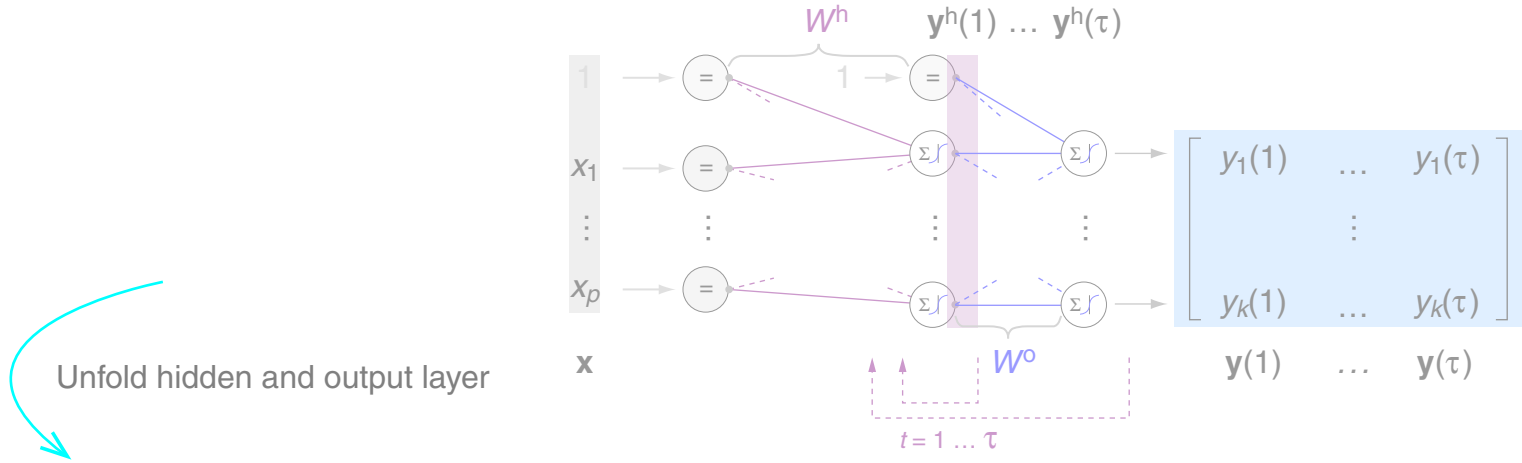
RNN Sequence Decoding

From Vector Output to Sequence Output



RNN Sequence Decoding

From Vector Output to Sequence Output



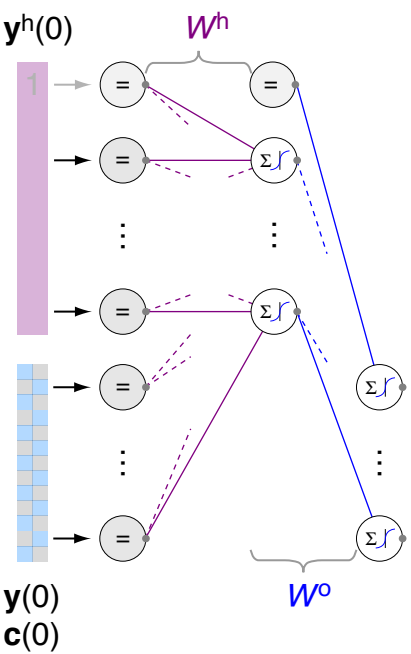
Remarks:

- ❑ 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 $\mathbf{y}(\tau) = (0, 0, \dots, 0, 1)^T$ ($\hat{=}$ `<end>`).
 τ (the length of the computed output) does not have to be equal to T (the length of the target or ground truth).

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 0 → 1

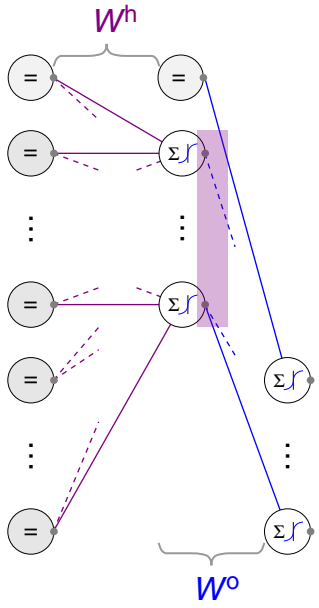


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 0 → 1

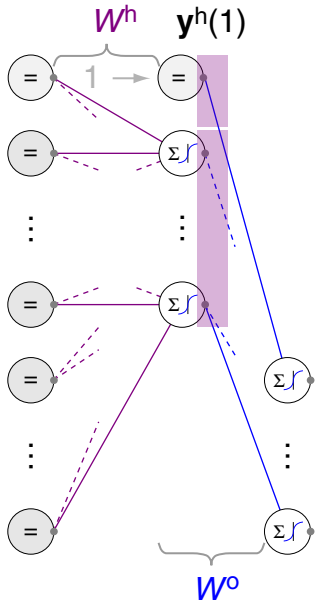


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 0 → 1

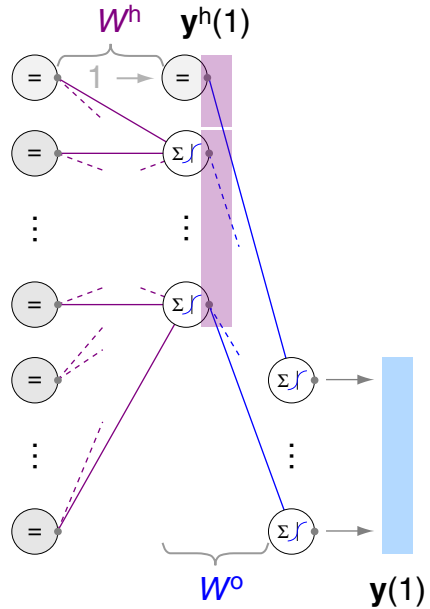


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 0 \rightarrow 1

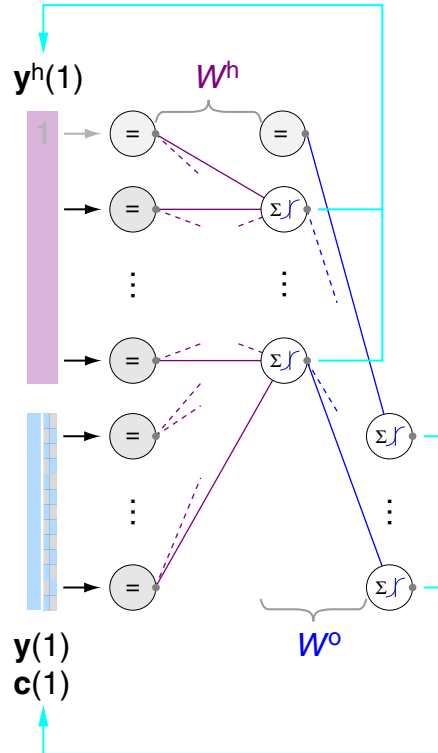


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 0 → 1

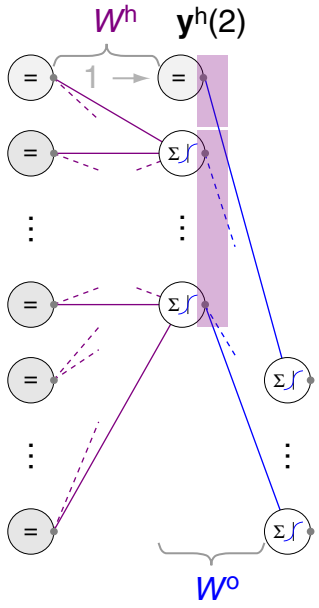


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 1 → 2

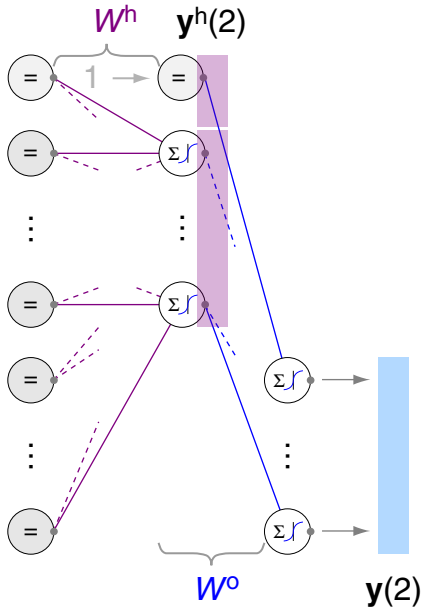


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 1 → 2

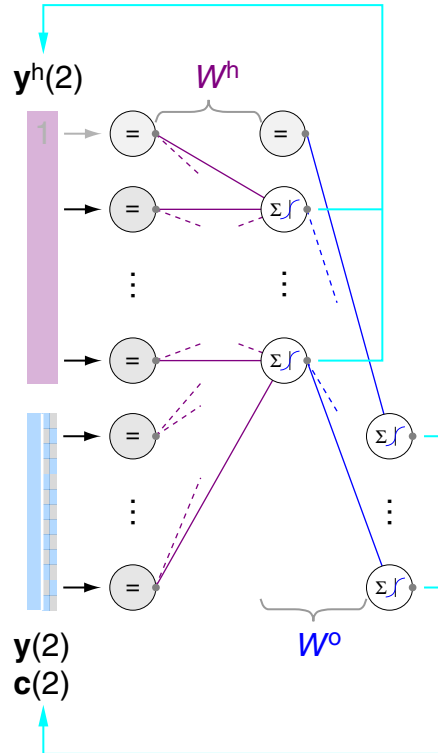


Output generation over t .

RNN Sequence Decoding

Recurrent Output Generation [\[decoding overview\]](#)

t: 1 → 2

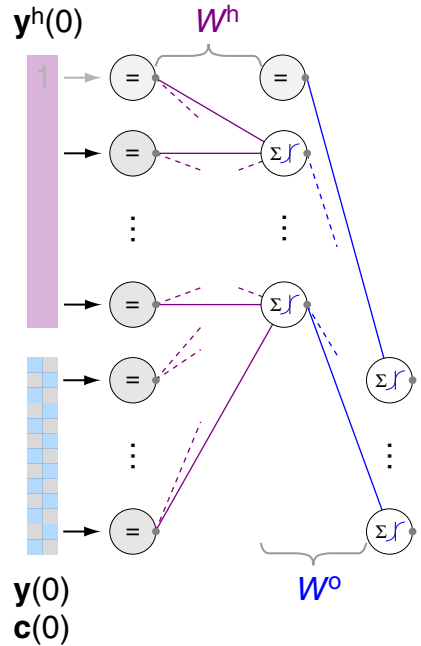


Output generation over t .

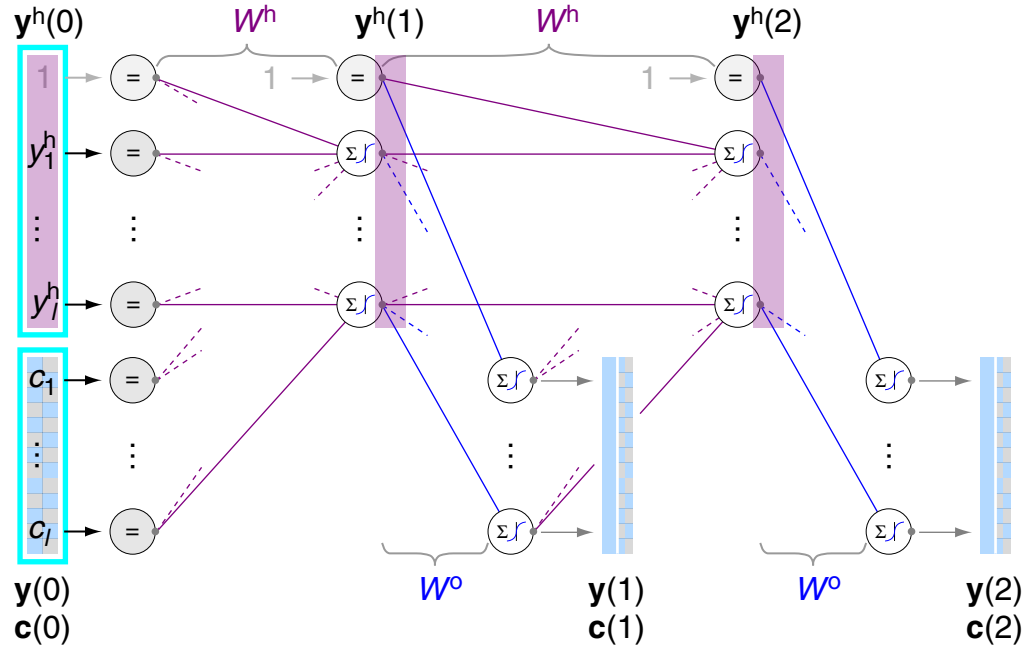
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 0 \rightarrow 1



Output generation over t .

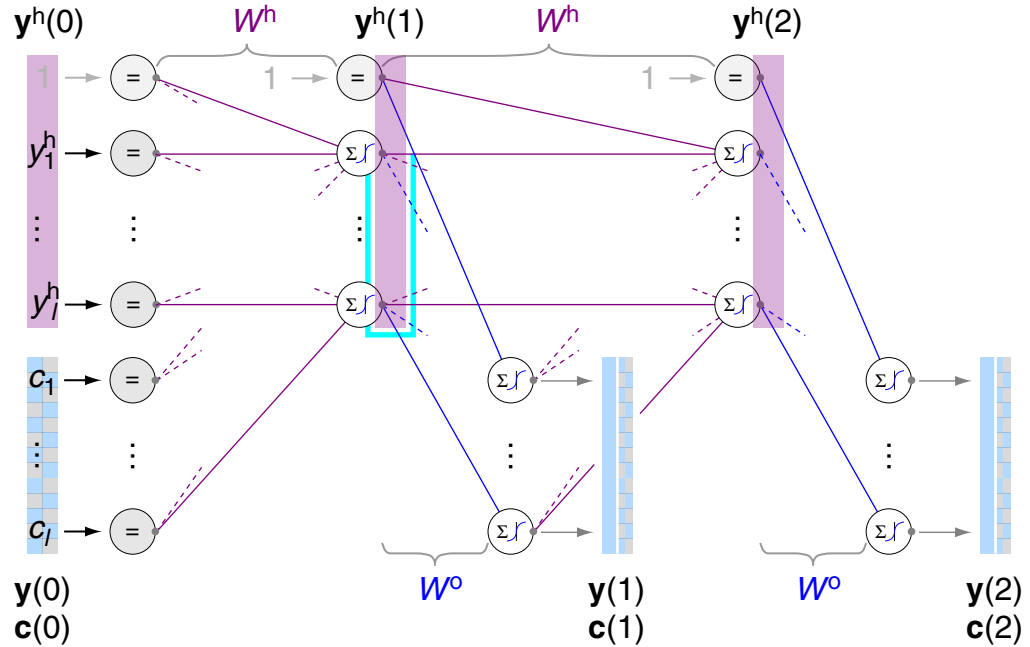
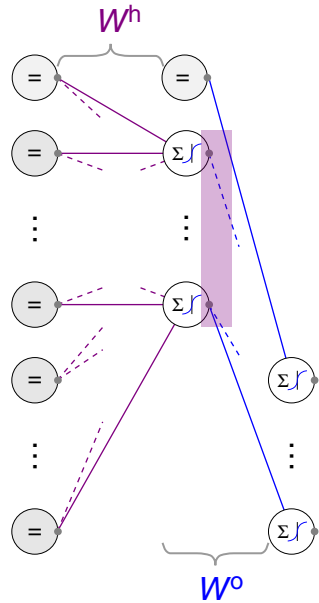


Hidden and output layer at subsequent time steps.

RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 0 → 1



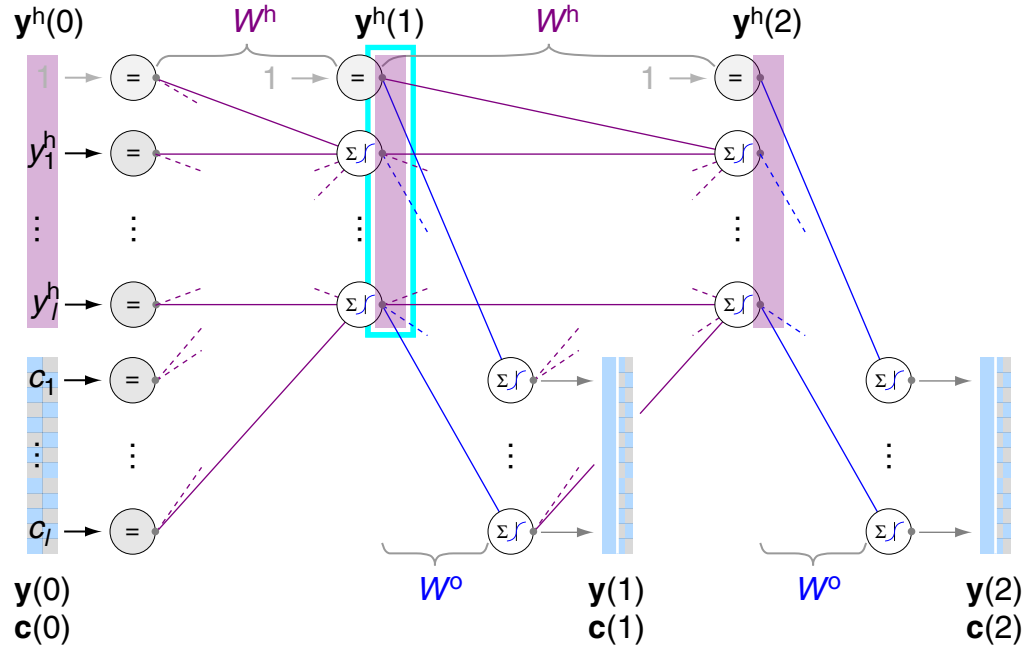
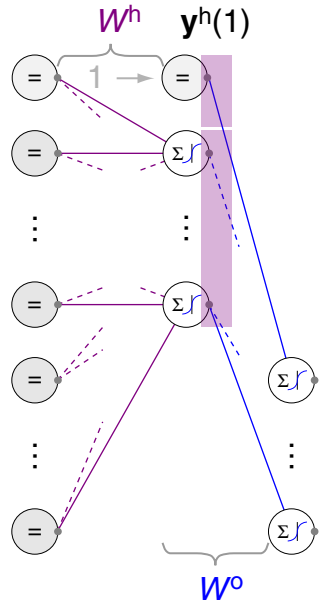
Output generation over t .

Hidden and output layer at subsequent time steps.

RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 0 → 1



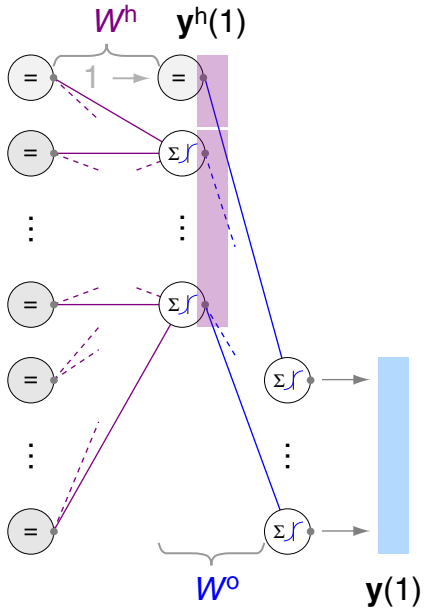
Output generation over t .

Hidden and output layer at subsequent time steps.

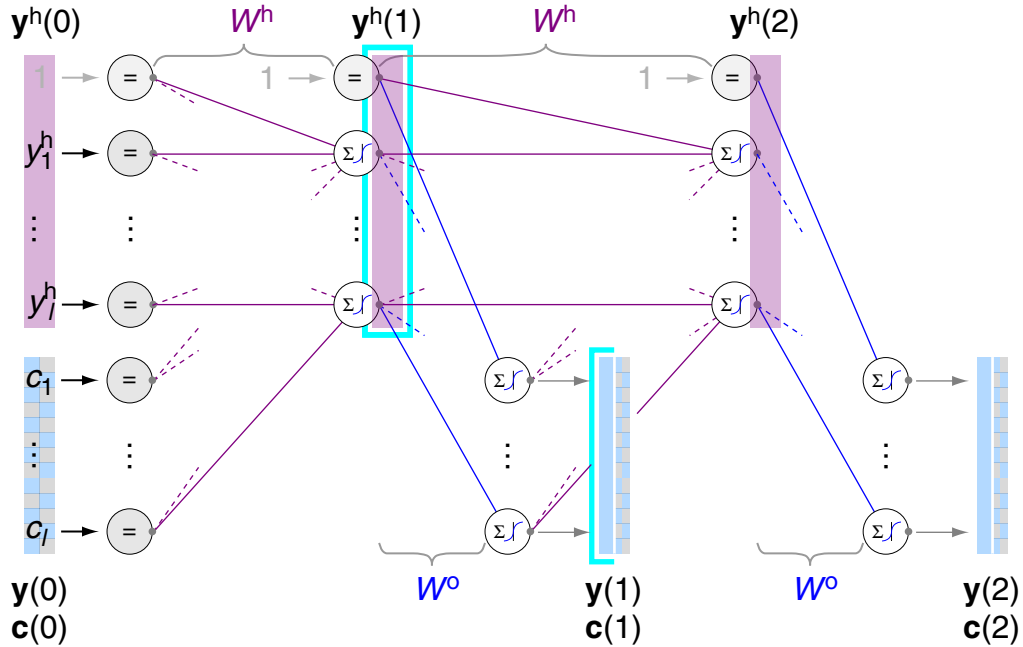
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 0 → 1



Output generation over t .

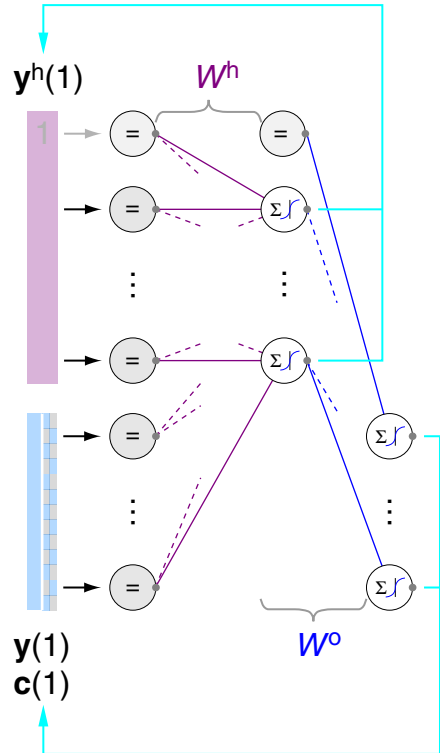


Hidden and output layer at subsequent time steps.

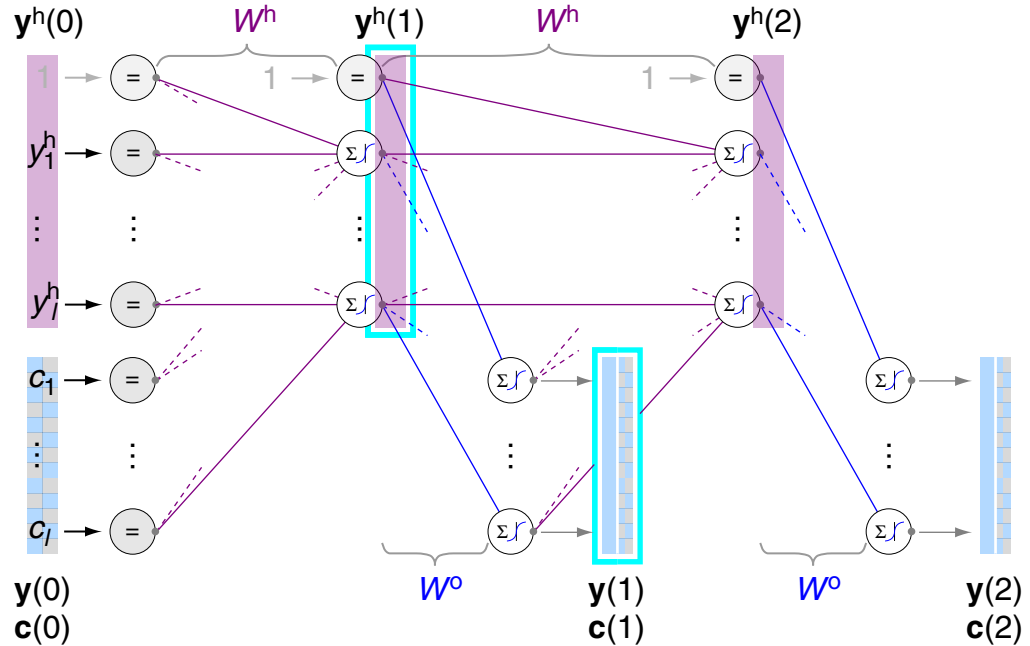
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 0 → 1



Output generation over t .

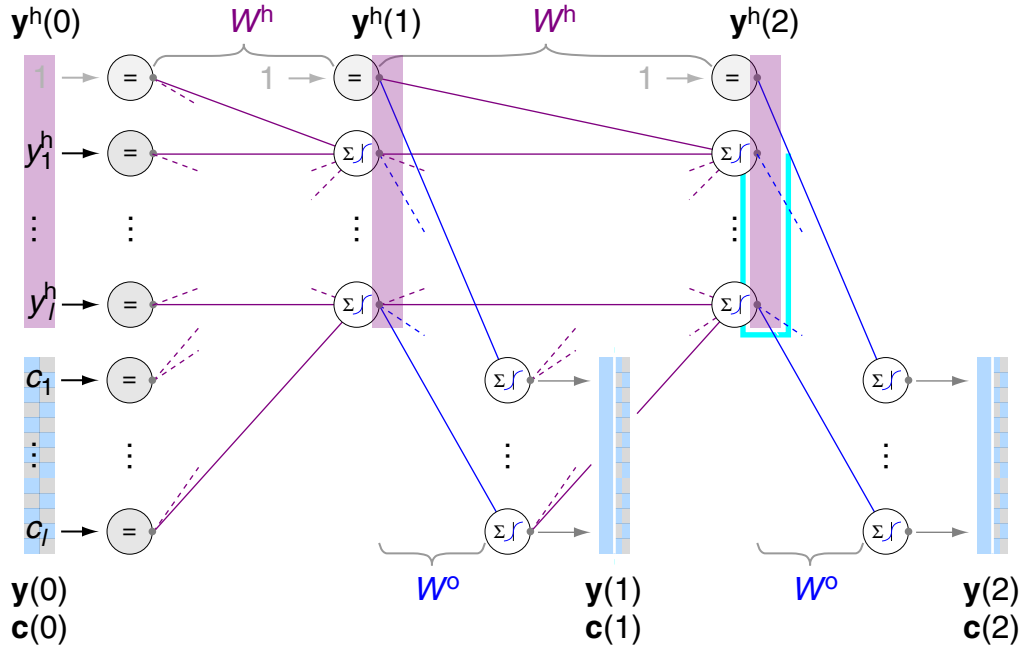
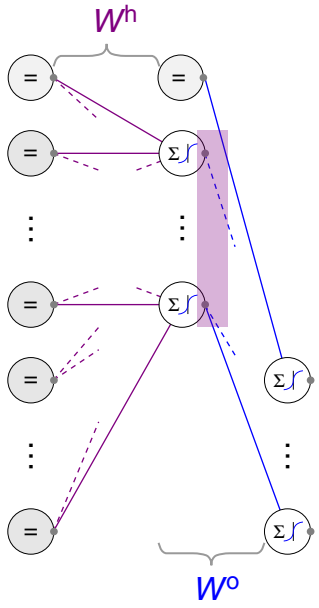


Hidden and output layer at subsequent time steps.

RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 1 → 2



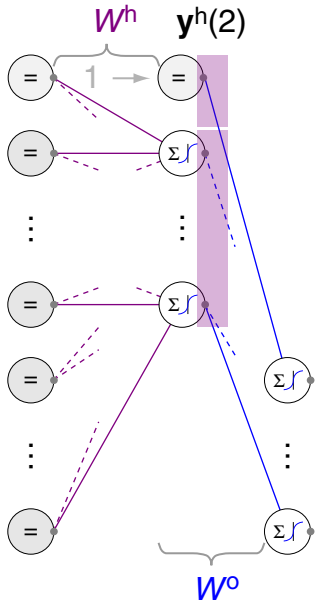
Output generation over t .

Hidden and output layer at subsequent time steps.

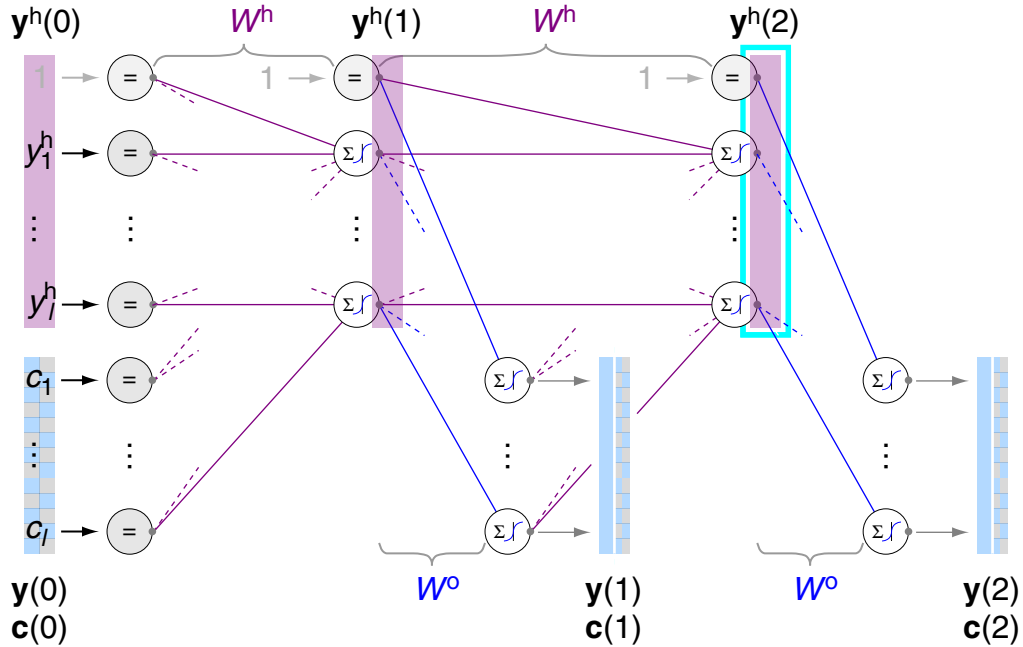
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 1 → 2



Output generation over t .

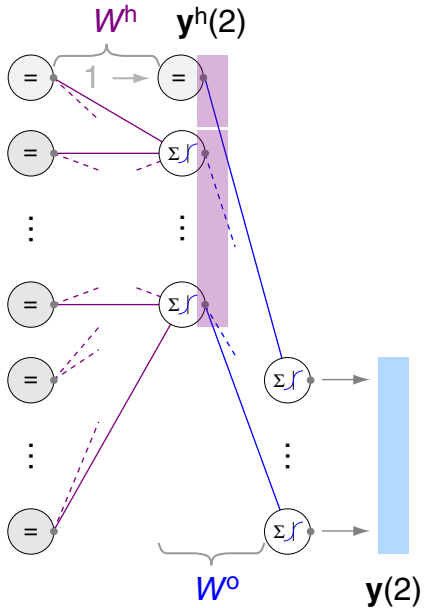


Hidden and output layer at subsequent time steps.

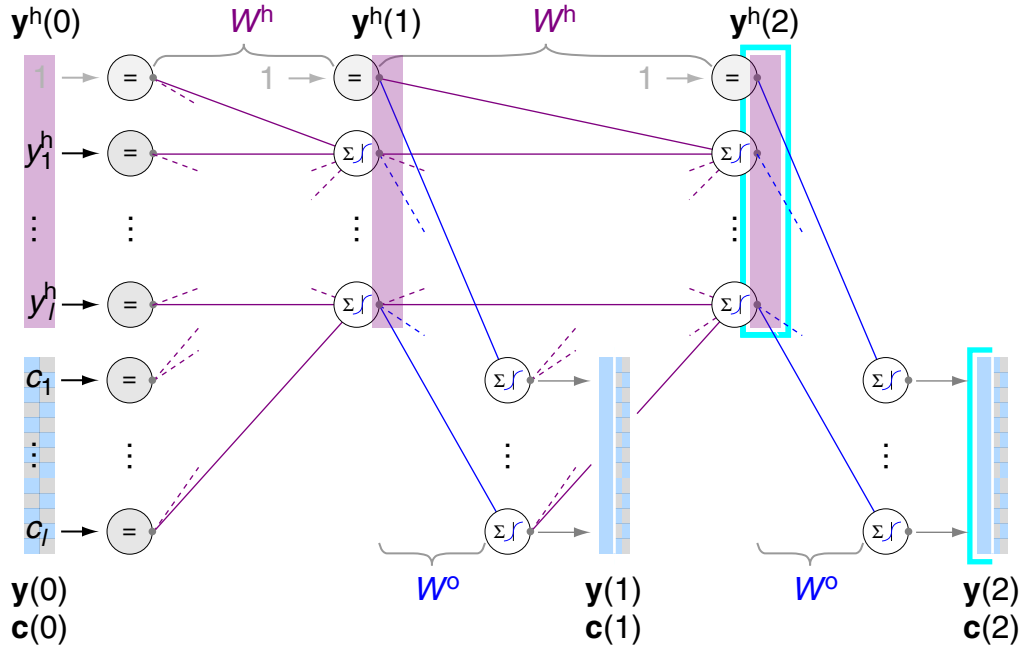
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 1 → 2



Output generation over t .

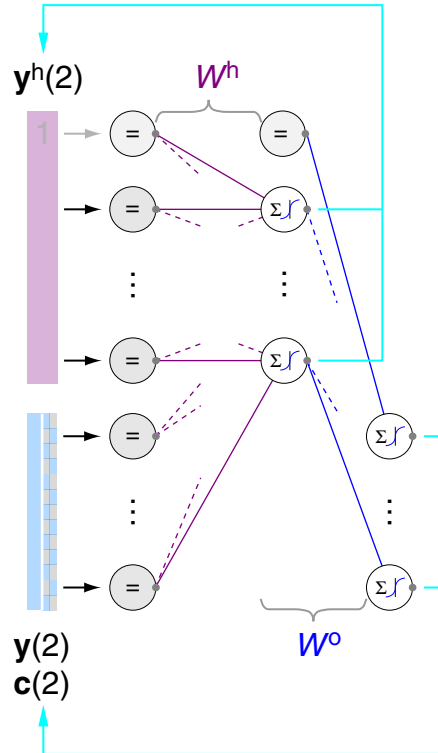


Hidden and output layer at subsequent time steps.

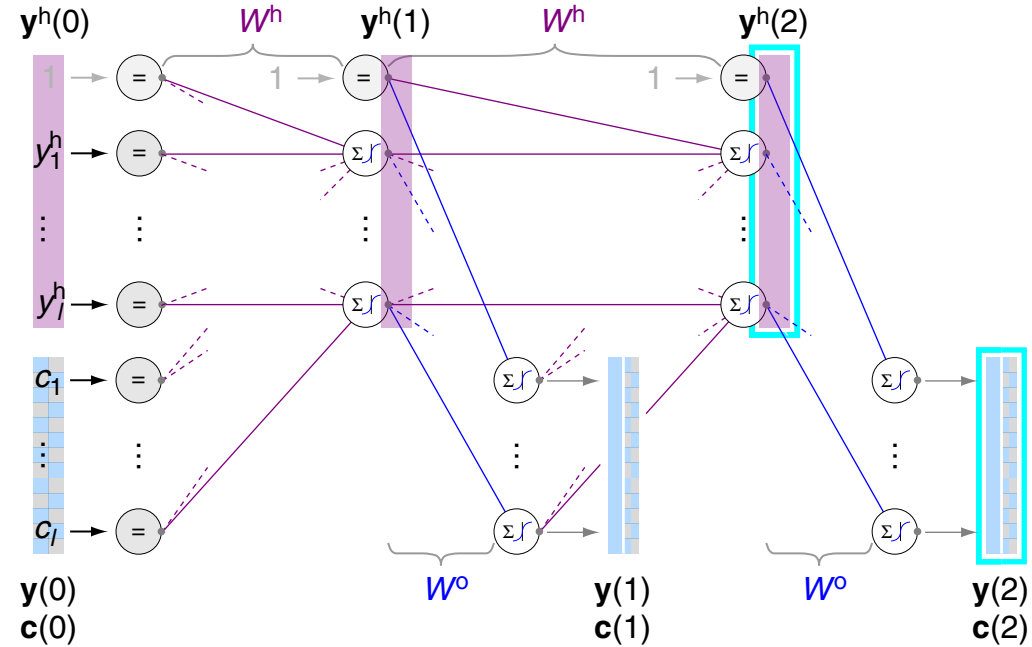
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 1 → 2



Output generation over t .

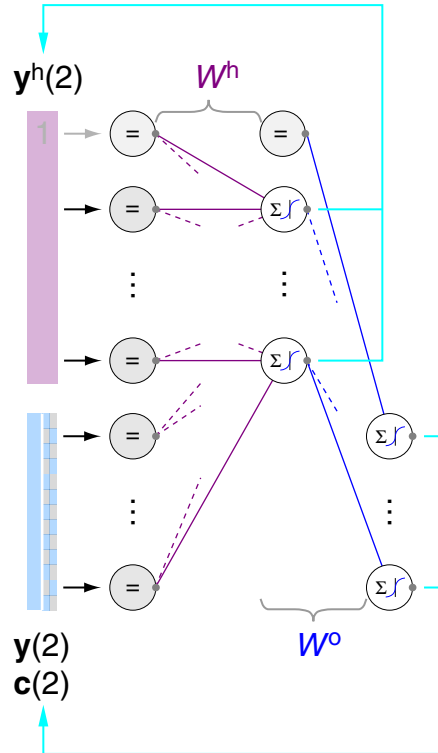


Hidden and output layer at subsequent time steps.

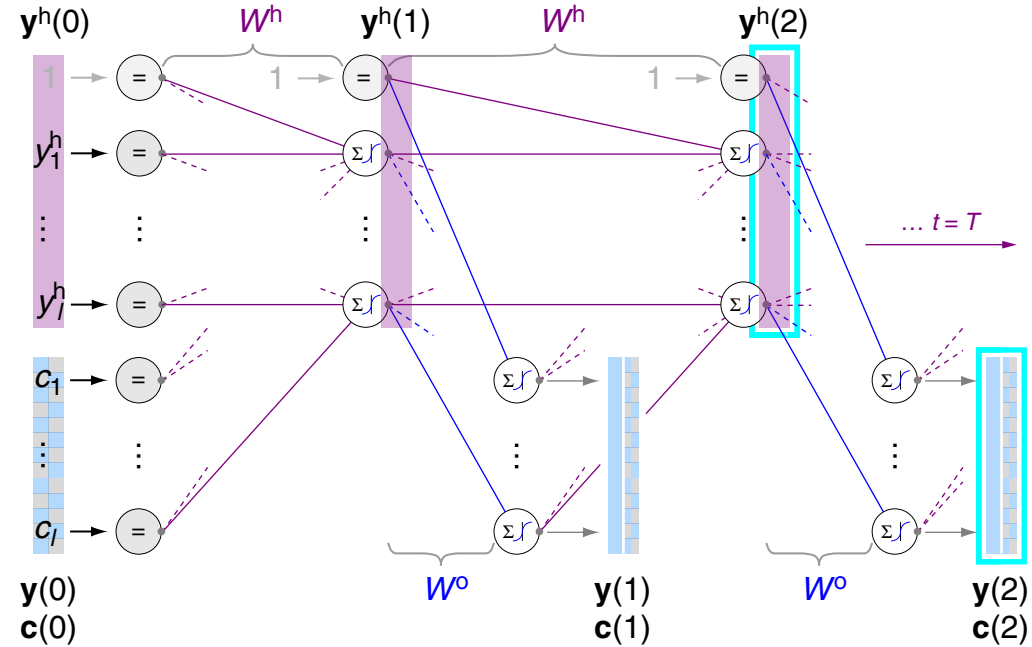
RNN Sequence Decoding

Recurrent Output Generation Unfolded [\[decoding overview\]](#)

t: 1 → 2



Output generation over t .



Hidden and output layer at subsequent time steps.

RNN Sequence Decoding

Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

RNN Sequence Decoding

Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

Input: $[[[[[\mathbf{x}, \mathbf{y}(0)], \mathbf{y}(1)], \mathbf{y}(2)], \dots], \mathbf{y}(\tau-1)]$, $\mathbf{x} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \hat{=} \text{word}_{11}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)]$, $\mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}$, $\mathbf{y}(\tau) \hat{=} \mathbf{c}(4) \hat{=} \text{<end>}$

RNN Sequence Decoding

Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

Input: [[[[**x**, **y(0)**], **y(1)**], **y(2)**], ...], **y(τ-1)**], $\mathbf{x} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \hat{=} \text{word}_{11}$

Output: **y(1)**, **y(2)**, **y(3)**, ..., **y(τ)**, $\mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}$, $\mathbf{y}(\tau) \hat{=} \mathbf{c}(4) \hat{=} \text{<end>}$

RNN Sequence Decoding

Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

Input: [[[[**x**, **y(0)**], **y(1)**], **y(2)**], ...], **y(τ-1)**], $\mathbf{x} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \hat{=} \text{word}_{11}$

Output: [**y(1)**, **y(2)**, **y(3)**, ..., **y(τ)**], $\mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}$, $\mathbf{y}(\tau) \hat{=} \mathbf{c}(4) \hat{=} \text{<end>}$

RNN Sequence Decoding

Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

Input: [[[[[**x**, **y(0)**], **y(1)**], **y(2)**], ...], **y(τ-1)**], $\mathbf{x} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \hat{=} \text{word}_{11}$

Output: [**y(1)**, **y(2)**, **y(3)**, ..., **y(τ)**], $\mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}$, $\mathbf{y}(\tau) \hat{=} \mathbf{c}(4) \hat{=} \text{<end>}$

RNN Sequence Decoding

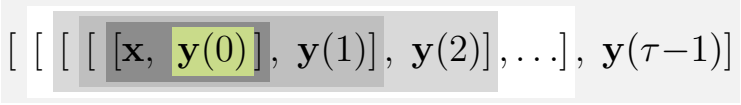
Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

Input:  $\mathbf{x} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \hat{=} \text{word}_{11}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \quad \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \quad \mathbf{y}(\tau) \hat{=} \mathbf{c}(4) \hat{=} \text{<end>}$

RNN Sequence Decoding

Class-to-Sequence: Text Generation

I → love my cat.

Cats → and dogs lap water.

It → is raining cats and dogs.

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

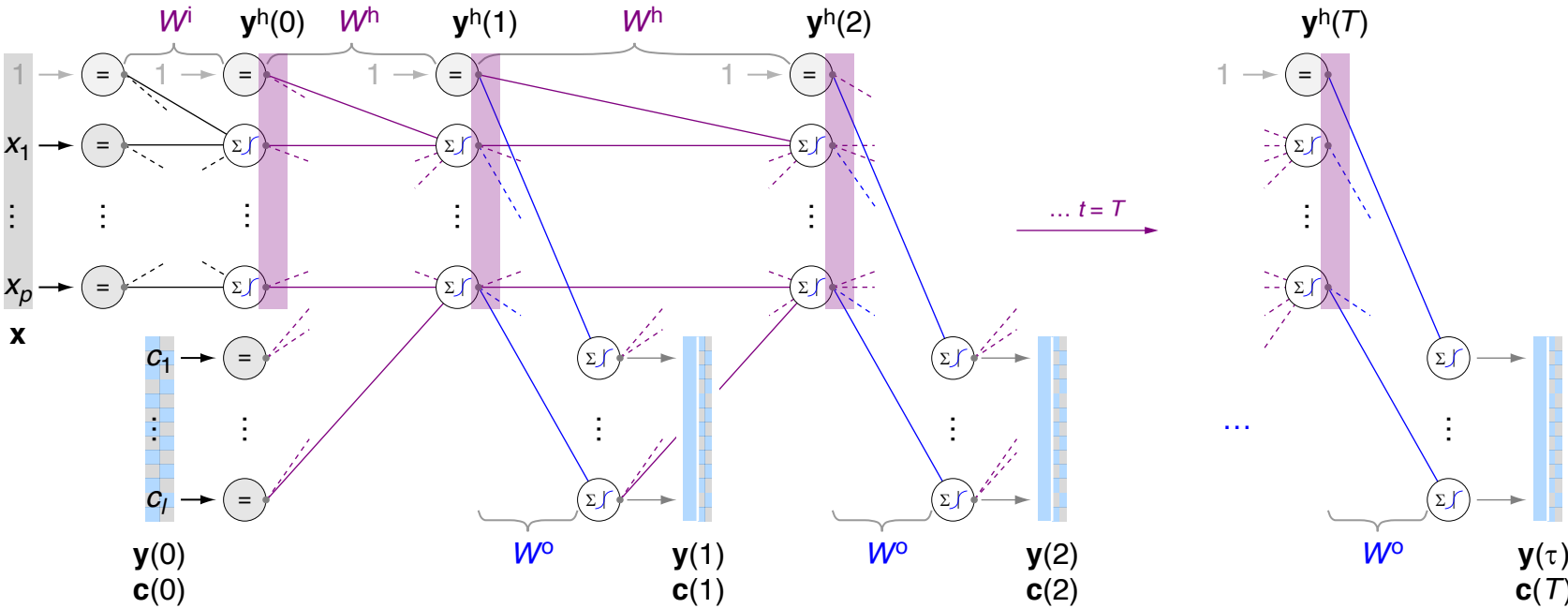
Input: $[[[[[[\mathbf{x}, \mathbf{y}(0)], \mathbf{y}(1)], \mathbf{y}(2)], \dots], \mathbf{y}(\tau-1)], \mathbf{x} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \hat{=} \text{word}_{11}$

Output: $[\mathbf{y}(1), \mathbf{y}(2), \mathbf{y}(3), \dots, \mathbf{y}(\tau)], \mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \text{<start>}, \mathbf{y}(\tau) \hat{=} \mathbf{c}(4) \hat{=} \text{<end>}$

Target: $[\mathbf{c}(1), \dots, \mathbf{c}(4)] = \left[\begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix}, \begin{pmatrix} 0 \\ \vdots \\ \vdots \\ 1 \end{pmatrix} \right]$
 $\hat{=} [\text{word}_{15}, \text{word}_{16}, \text{word}_6, \text{word}_{22}]$
 $\hat{=} \text{love my cat}$

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



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

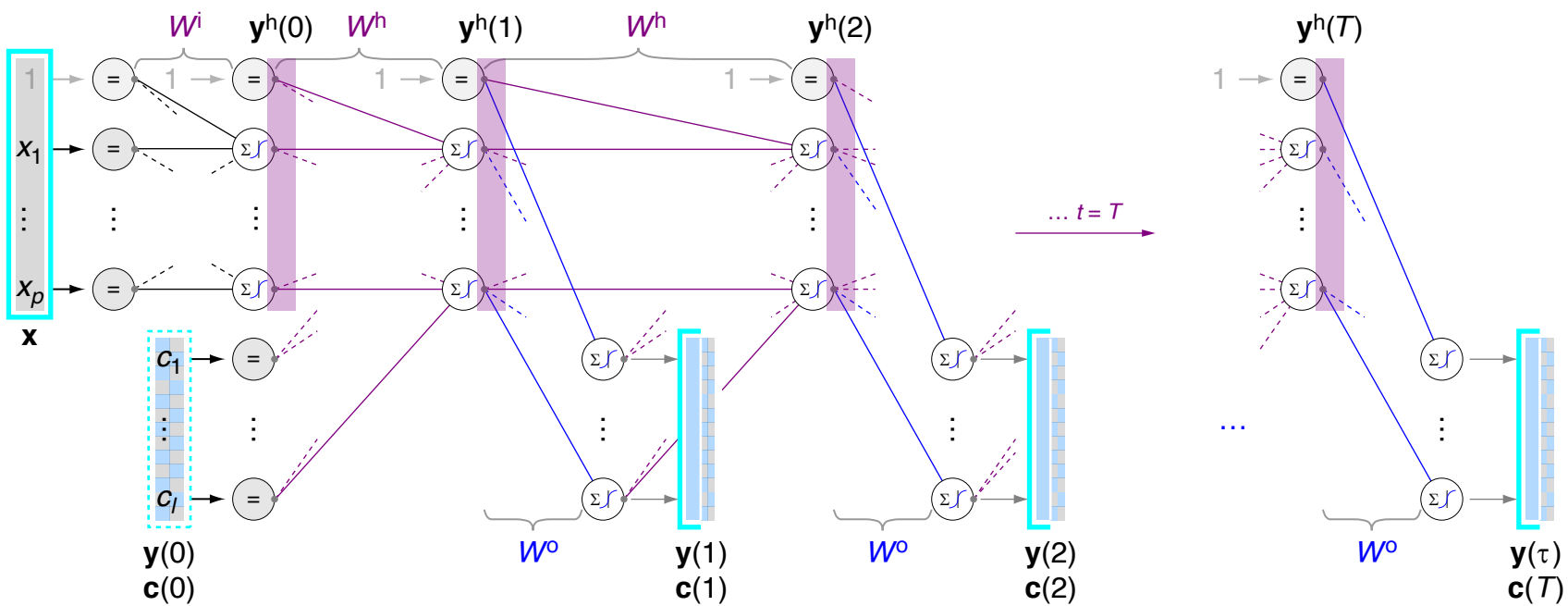
Hidden:
 $\mathbf{y}^h(0) = \sigma(W^i \mathbf{x})$
 $\mathbf{y}^h(t) = \sigma\left(W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix}\right), t = 1, \dots, \tau$

Target:
 $[\mathbf{c}(1), \dots, \mathbf{c}(T)]$
 $\mathbf{c}(T) \hat{=} \langle \text{end} \rangle$

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

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



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

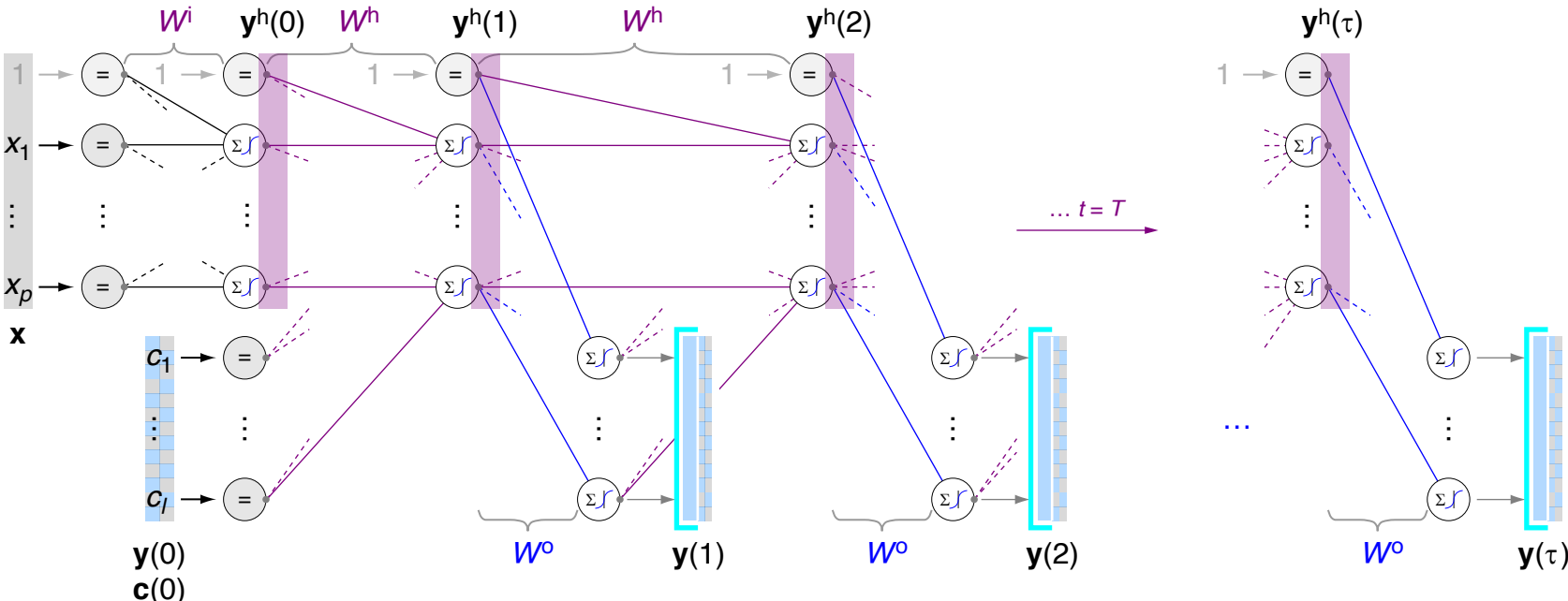
Output:
 $\mathbf{y}(t) = \sigma_{\Delta}(W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$

Hidden:
 $\mathbf{y}^h(0) = \sigma(W^i \mathbf{x})$
 $\mathbf{y}^h(t) = \sigma(W^h(\mathbf{y}^h(t-1), \mathbf{y}(t-1))), t = 1, \dots, \tau$

Target:
 $[\mathbf{c}(1), \dots, \mathbf{c}(T)]$
 $\mathbf{c}(T) \hat{=} \langle \text{end} \rangle$

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



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

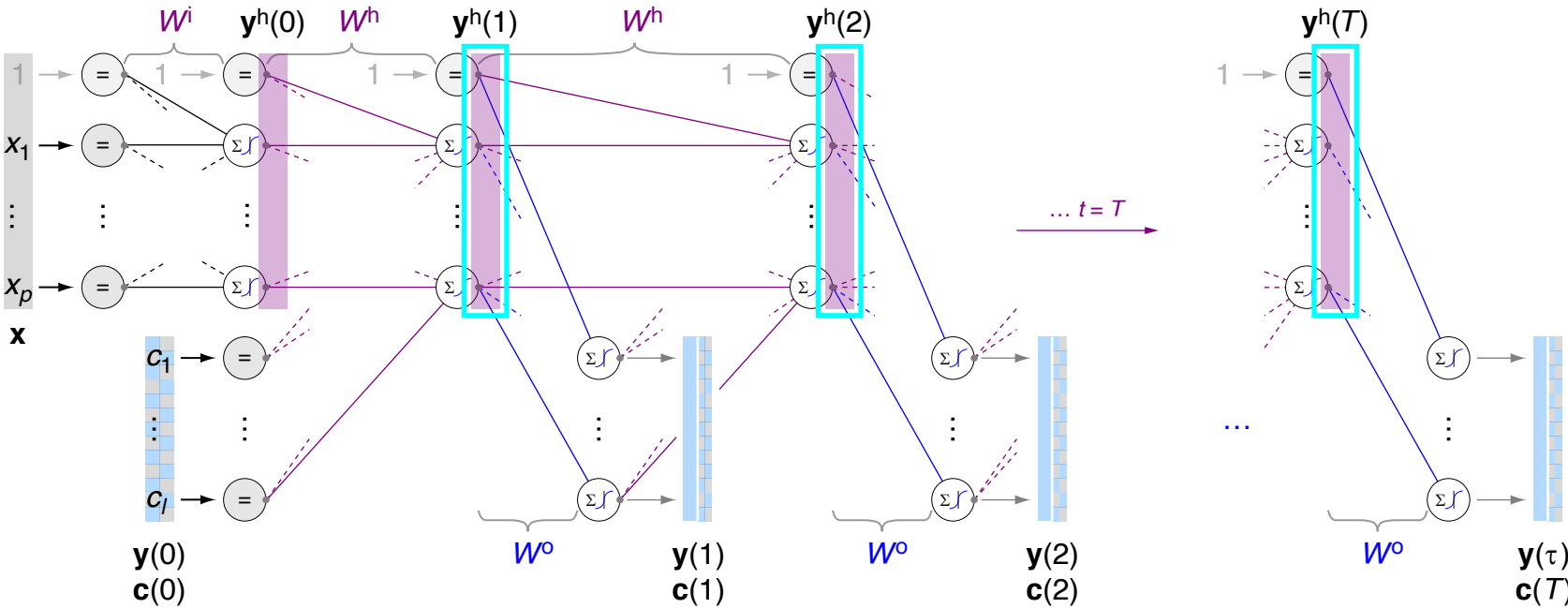
Output:
 $\mathbf{y}(t) = \sigma_{\Delta} (W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$

Hidden:
 $\mathbf{y}^h(0) = \sigma (W^i \mathbf{x})$
 $\mathbf{y}^h(t) = \sigma (W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix}), t = 1, \dots, \tau$

Target:
 $[\mathbf{c}(1), \dots, \mathbf{c}(T)]$
 $\mathbf{c}(T) \hat{=} \langle \text{end} \rangle$

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



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

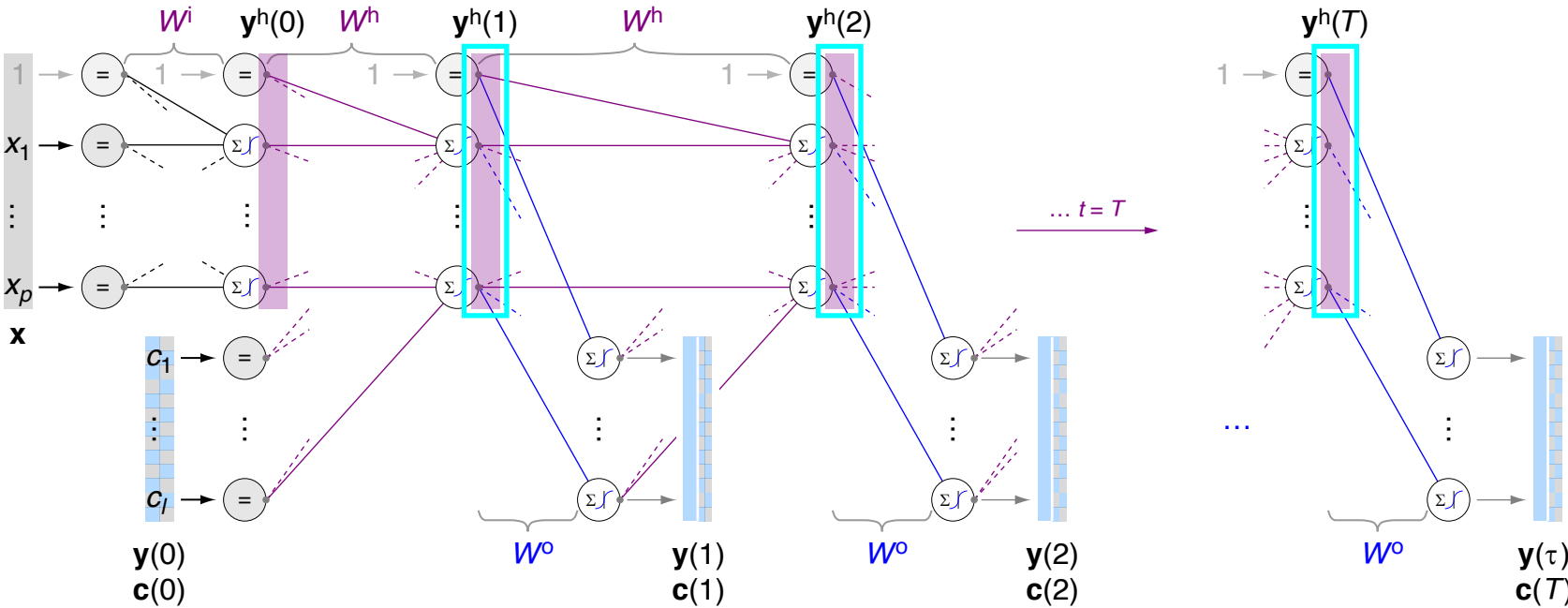
Output:
 $\mathbf{y}(t) = \sigma_{\Delta} (W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$

Hidden:
 $\mathbf{y}^h(0) = \sigma (W^i \mathbf{x})$
 $\mathbf{y}^h(t) = \sigma (W^h (\mathbf{y}^h(t-1))) , t = 1, \dots, T$

Target:
 $[\mathbf{c}(1), \dots, \mathbf{c}(T)]$
 $\mathbf{c}(T) \hat{=} \langle \text{end} \rangle$

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



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

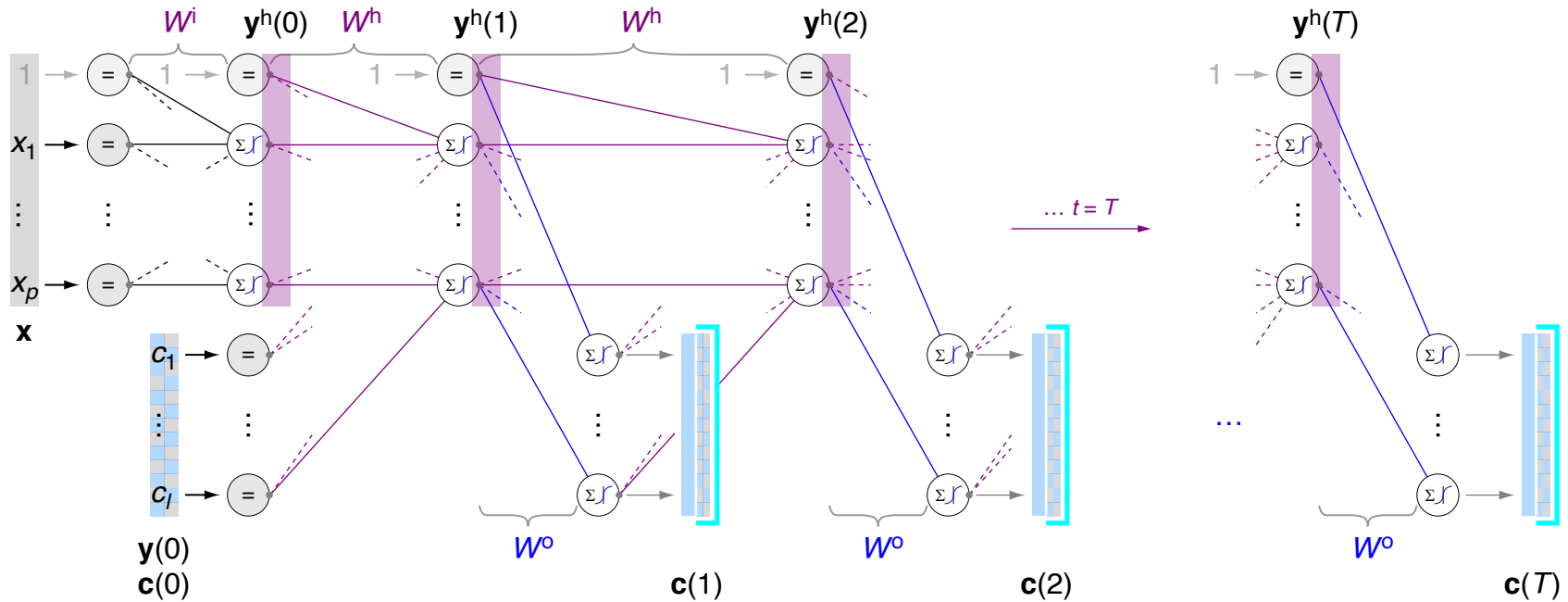
Output:
 $\mathbf{y}(t) = \sigma_{\Delta} (W^o \mathbf{y}^h(t)), t = 1, \dots, \tau$

Hidden:
 $\mathbf{y}^h(0) = \sigma (W^i \mathbf{x})$
 $\mathbf{y}^h(t) = \sigma (W^h \begin{pmatrix} \mathbf{y}^h(t-1) \\ \mathbf{y}(t-1) \end{pmatrix}), t = 1, \dots, \tau$

Target:
 $[\mathbf{c}(1), \dots, \mathbf{c}(T)]$
 $\mathbf{c}(T) \hat{=} \langle \text{end} \rangle$

RNN Sequence Decoding

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:

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

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

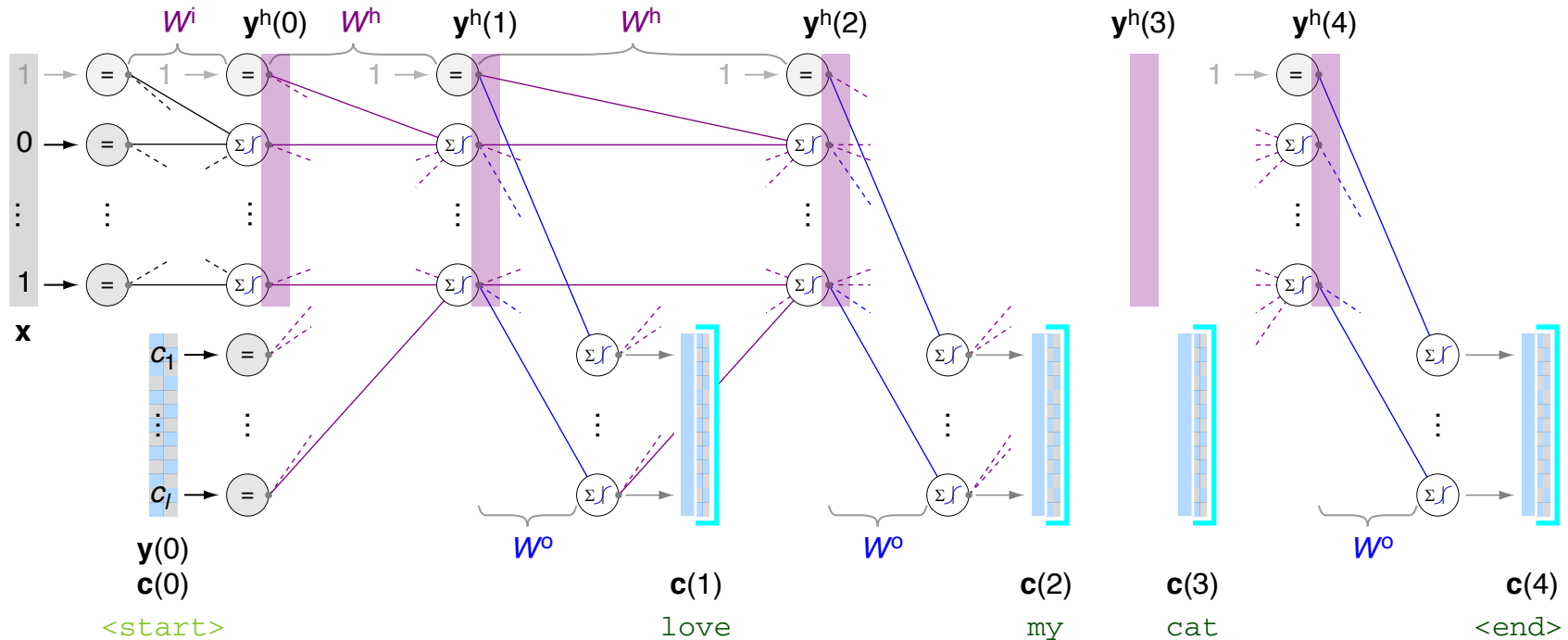
Target:

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

$$\mathbf{c}(T) \hat{=} \langle \text{end} \rangle$$

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



Input:

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

Output:

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

Hidden:

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

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

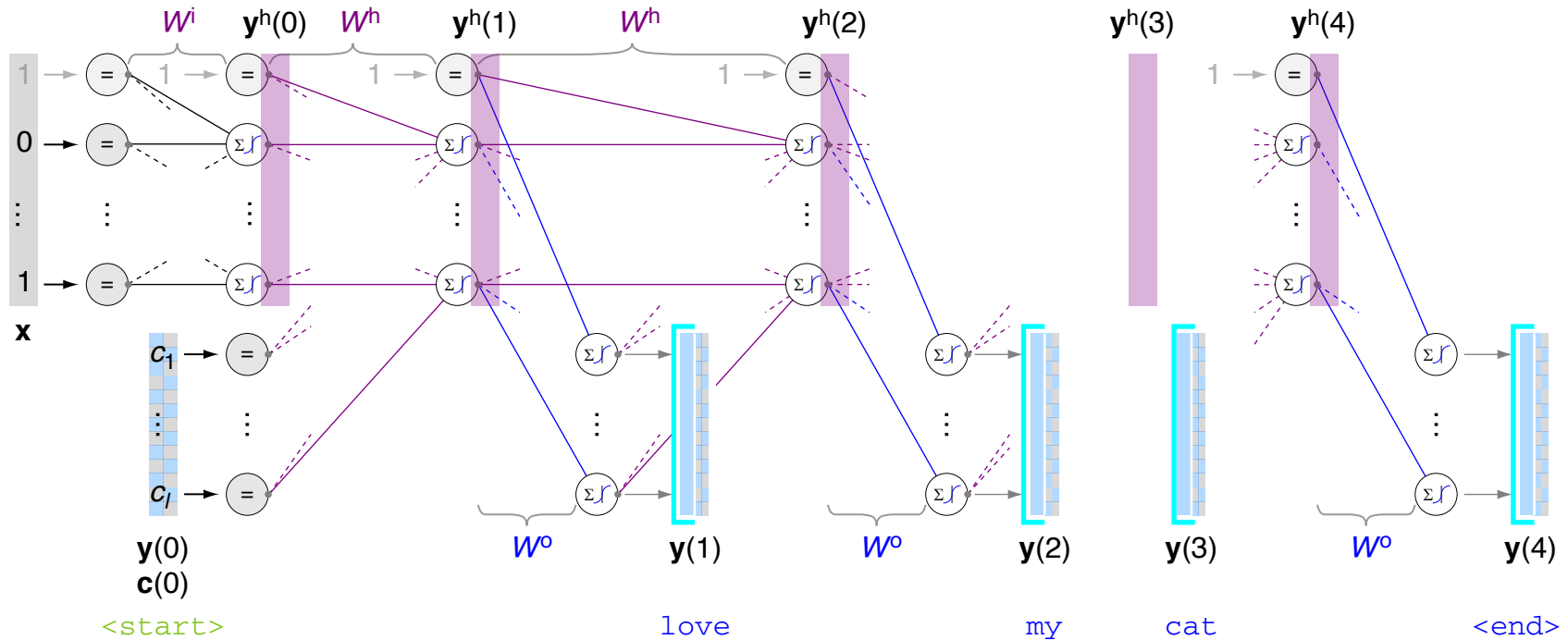
Target:

$$[c(1), \dots, c(4)]$$

$$c(4) \hat{=} \text{<end>}$$

RNN Sequence Decoding

Class-to-Sequence Mapping with RNNs



Input:

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

Output:

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

Hidden:

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

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

Target:

$$[c(1), \dots, c(4)]$$

$$c(4) \hat{=} \text{<end>}$$

Remarks:

- We denote $\mathbf{y}(0)$ not as input since it is predefined and does not contain any “actual knowledge”. In particular, $\mathbf{y}(0) \equiv \mathbf{c}(0) \hat{=} \langle \text{start} \rangle$.
- At training time the calculation of $\mathbf{y}^h(t)$ usually considers the ground truth $\mathbf{c}(t-1)$:

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

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

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

RNN Sequence Decoding

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 $\mathbf{x} \in \{0, 1\}^p$, $\mathbf{c}(t) \in \mathbf{R}^k$.
 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**(*convergence*($D, \mathbf{y}(\cdot), t_{\text{epoch}}$))
11. *return*(W^i, W^h, W^o)

RNN Sequence Decoding

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