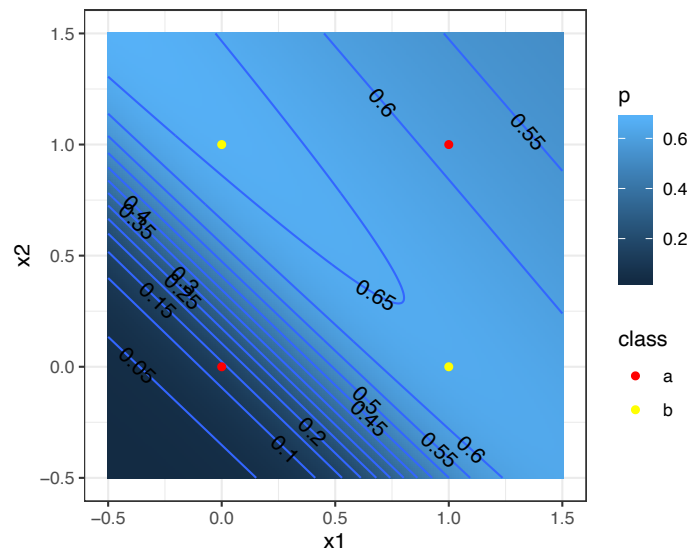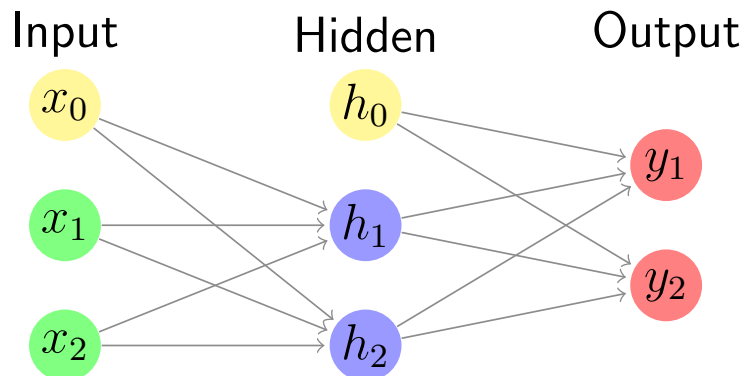# Neural networks for natural language

Roger Levy
9.19: Computational Psycholinguistics
2 November 2023

# Agenda for the day

- Last time: with a hidden layer, a NN can learn XOR...



- ...but language isn't just 2D input+2-class output! So, **today:**

- Dealing with language in neural networks

- Recurrent neural networks (RNNs)

  - Simple recurrent networks (SRNs)

  - Gated recurrent units (GRUs)

  - Long short-term memory networks (LSTMs)

- Examining RNN behavior

# Dealing with language inputs

`Adam adores zebras ...`

# Dealing with language inputs

For language, input $\{x_i\}$ and output prediction $y$ seem discrete:

**Adam adores zebras ...**

# Dealing with language inputs

For language, input $\{x_i\}$ and output prediction $y$ seem discrete:

**Adam adores zebras ...**

Simplest approach is *localist* or *one-hot* representations:

$$\text{Adam} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \text{adores} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \text{zebras} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

# Dealing with language inputs

For language, input $\{x_i\}$ and output prediction $y$ seem discrete:

**`Adam adores zebras ...`**

Simplest approach is *localist* or *one-hot* representations:

$$\texttt{Adam} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \texttt{adores} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \quad \texttt{zebras} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

But lower-dimensional *embeddings* capture word similarities:

$$\texttt{Adam} \rightarrow \begin{bmatrix} 0.6 \\ 0.3 \end{bmatrix} \quad \texttt{adores} \rightarrow \begin{bmatrix} -0.3 \\ 0.4 \end{bmatrix} \quad \texttt{zebras} \rightarrow \begin{bmatrix} 0.7 \\ -0.1 \end{bmatrix}$$

# Example feed-forward+embedding LM

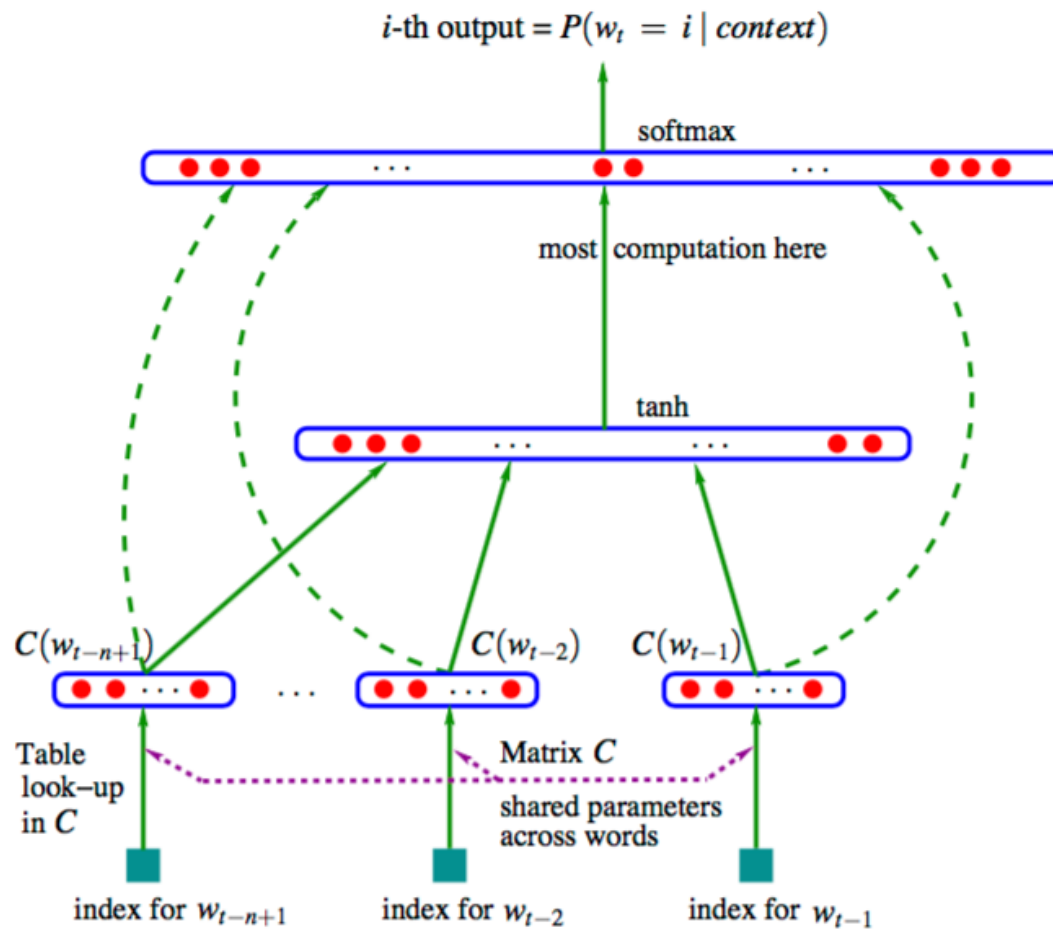Bengio et al., 2003: Neural *n*-gram language model



Figure 1: Neural architecture: $f(i, w_{t-1}, \cdots, w_{t-n+1}) = g(i, C(w_{t-1}), \cdots, C(w_{t-n+1}))$ where $g$ is the neural network and $C(i)$ is the $i$-th word feature vector.
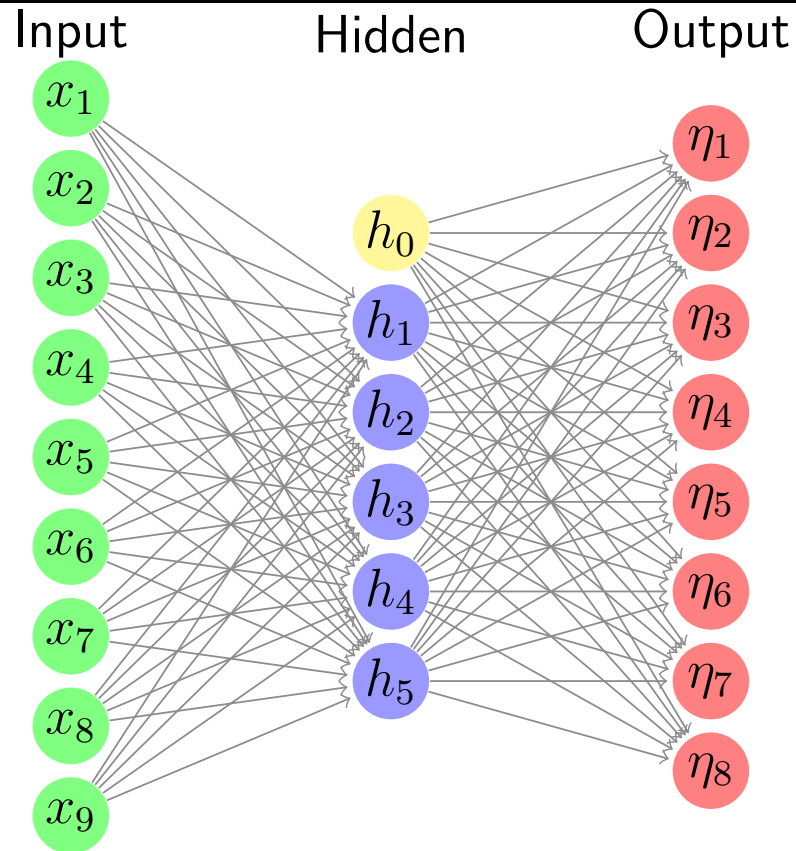
# Old (2003!) perplexity results on Brown corpus

| | n | c | h | m | direct | mix | train. | valid. | test. |
|---|---|---|---|---|---|---|---|---|---|
| MLP1 | 5 | | 50 | 60 | yes | no | 182 | 284 | 268 |
| MLP2 | 5 | | 50 | 60 | yes | yes | | 275 | 257 |
| MLP3 | 5 | | 0 | 60 | yes | no | 201 | 327 | 310 |
| MLP4 | 5 | | 0 | 60 | yes | yes | | 286 | 272 |
| MLP5 | 5 | | 50 | 30 | yes | no | 209 | 296 | 279 |
| MLP6 | 5 | | 50 | 30 | yes | yes | | 273 | 259 |
| MLP7 | 3 | | 50 | 30 | yes | no | 210 | 309 | 293 |
| MLP8 | 3 | | 50 | 30 | yes | yes | | 284 | 270 |
| MLP9 | 5 | | 100 | 30 | no | no | 175 | 280 | 276 |
| MLP10 | 5 | | 100 | 30 | no | yes | | 265 | **252** |
| Del. Int. | 3 | | | | | | 31 | 352 | 336 |
| Kneser-Ney back-off | 3 | | | | | | | 334 | 323 |
| Kneser-Ney back-off | 4 | | | | | | | 332 | 321 |
| Kneser-Ney back-off | 5 | | | | | | | 332 | 321 |
| class-based back-off | 3 | 150 | | | | | | 348 | 334 |
| class-based back-off | 3 | 200 | | | | | | 354 | 340 |
| class-based back-off | 3 | 500 | | | | | | 326 | **312** |
| class-based back-off | 3 | 1000 | | | | | | 335 | 319 |
| class-based back-off | 3 | 2000 | | | | | | 343 | 326 |
| class-based back-off | 4 | 500 | | | | | | 327 | 312 |
| class-based back-off | 5 | 500 | | | | | | 327 | 312 |

*neural language models*

*n-gram language models*

*(Bengio et al., 2003)*

5

# The neural *n*-gram model

Input       Hidden       Output

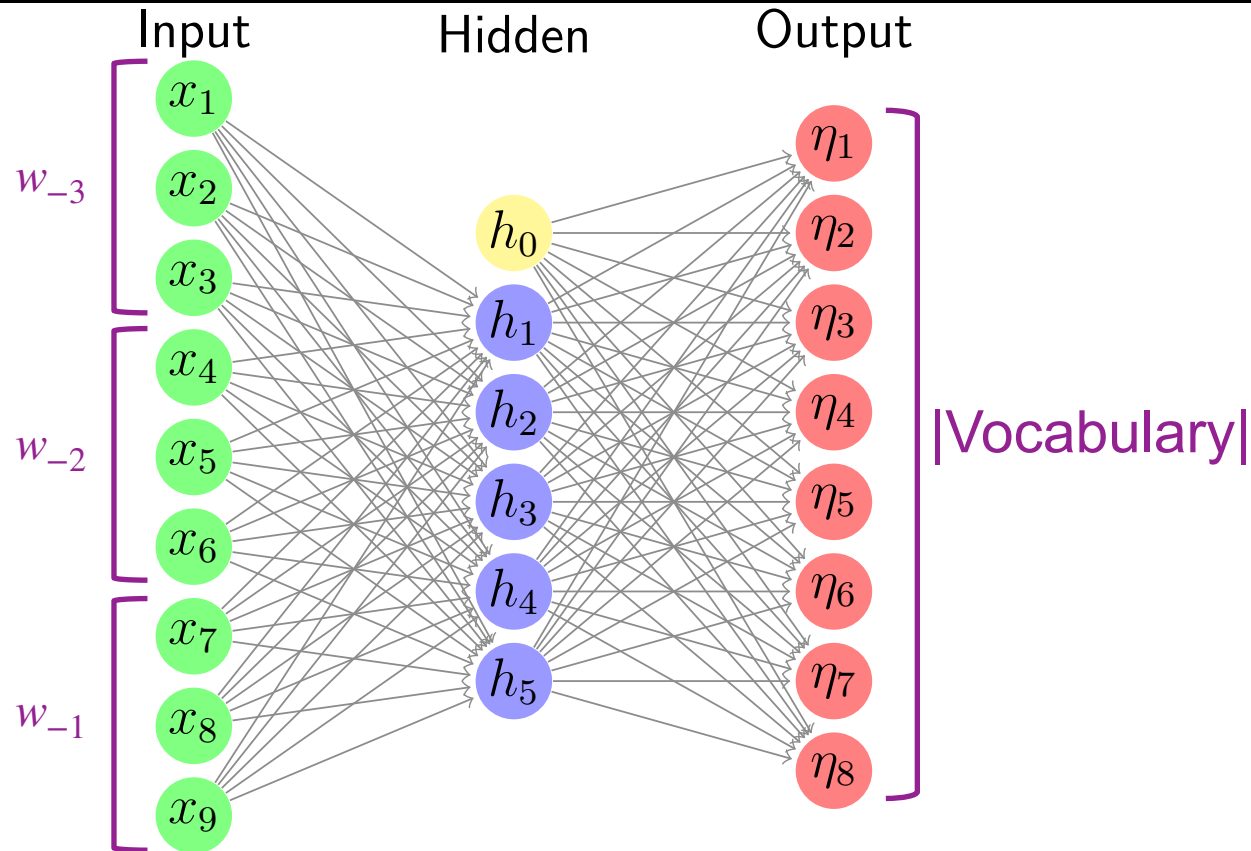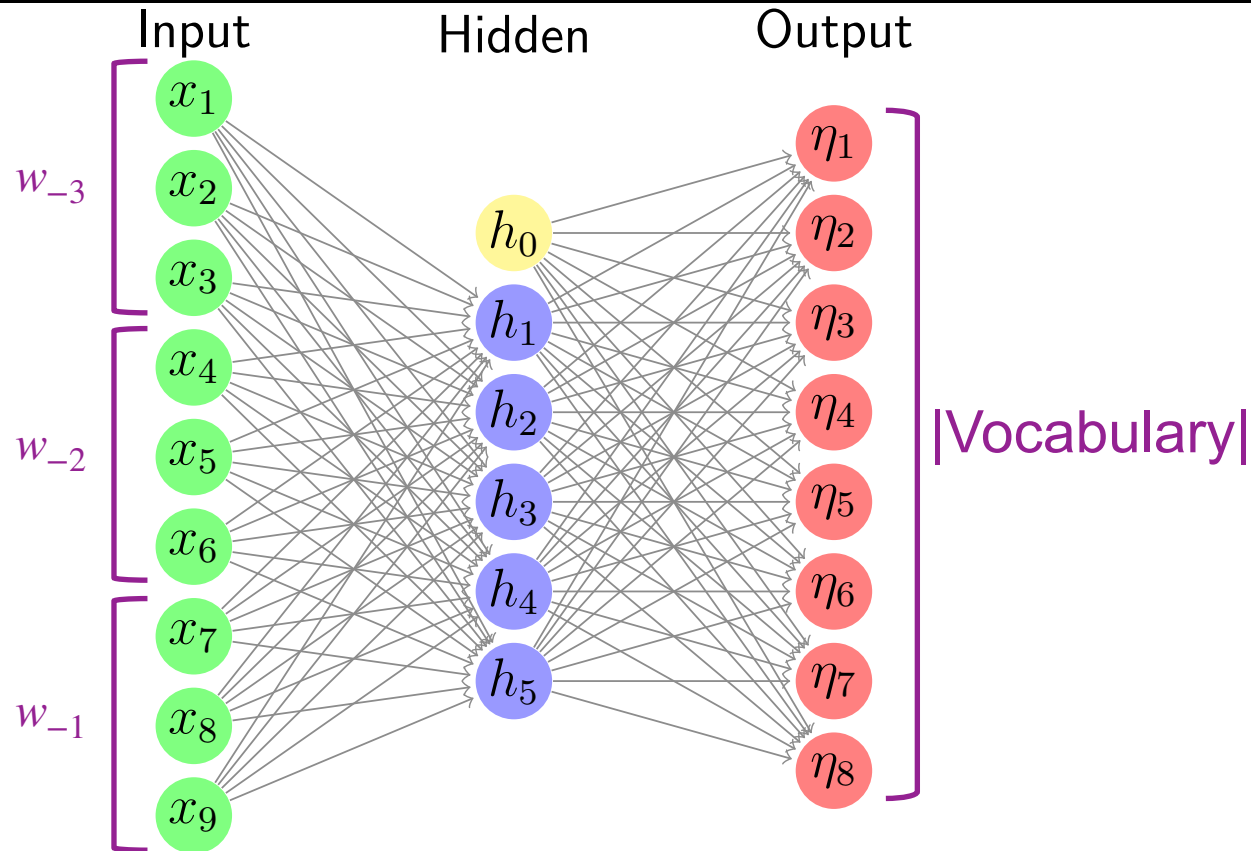# The neural *n*-gram model

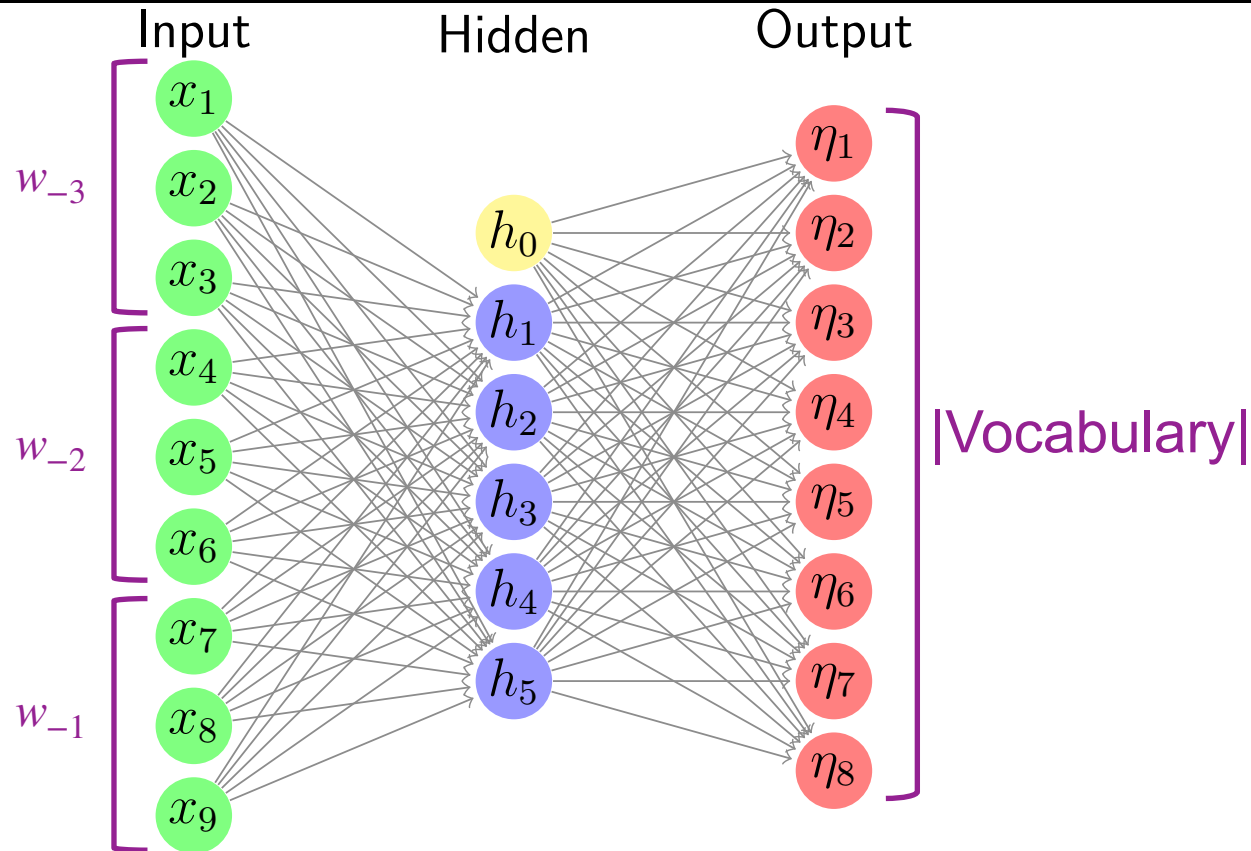# The neural *n*-gram model

# The neural *n*-gram model

- **Advantages:** generalizes over *n*-gram contexts

# The neural *n*-gram model
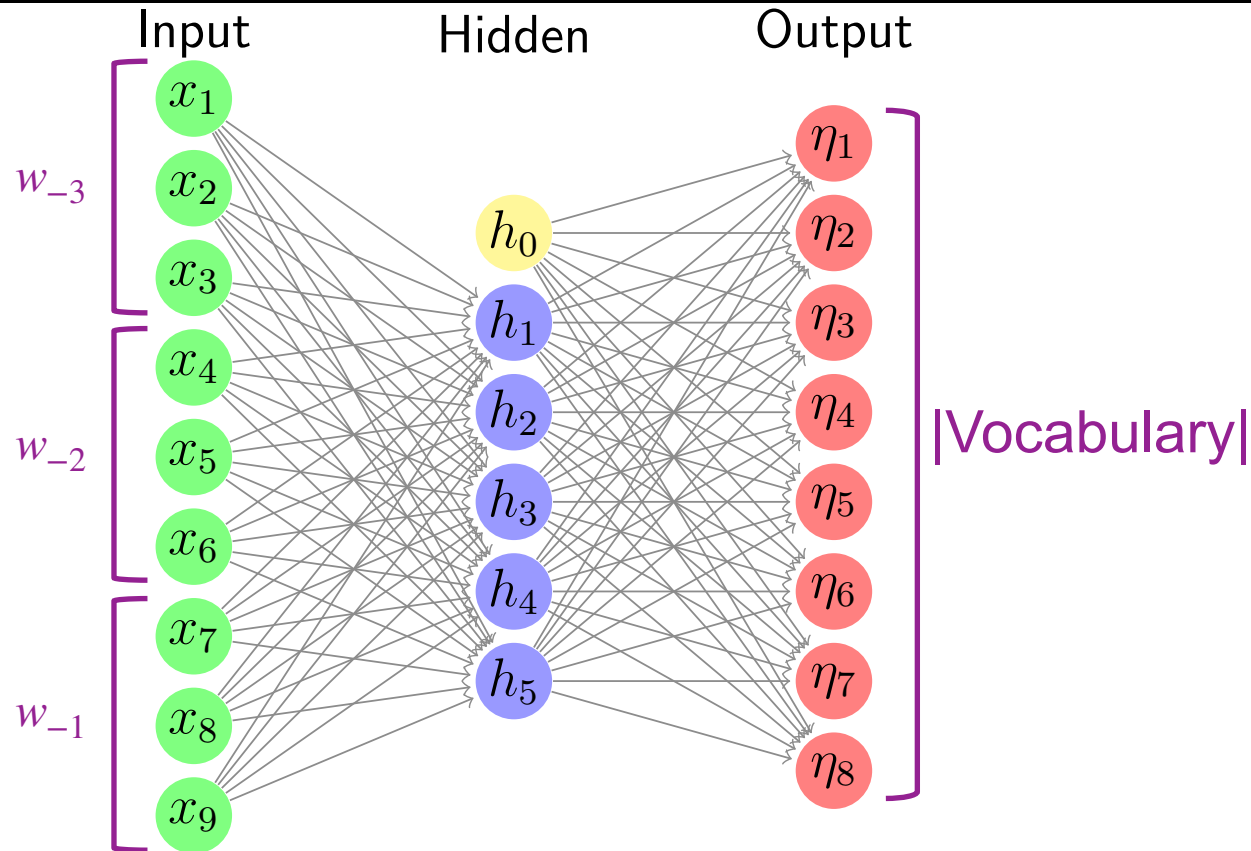


- **Advantages:** generalizes over *n*-gram contexts
- **Limitations:**

# The neural *n*-gram model



- **Advantages:** generalizes over *n*-gram contexts
- **Limitations:**
  - this is for a fixed dimensionality input context

# The neural *n*-gram model



- **Advantages:** generalizes over *n*-gram contexts
- **Limitations:**
  - this is for a fixed dimensionality input context
  - how to model variable-length context, like sentences?

# Recurrent neural networks for language

*(Jordan, 1986; Elman, 1990)*

# Recurrent neural networks for language

- Draw inspiration from real-time nature of human language processing

*(Jordan, 1986; Elman, 1990)*
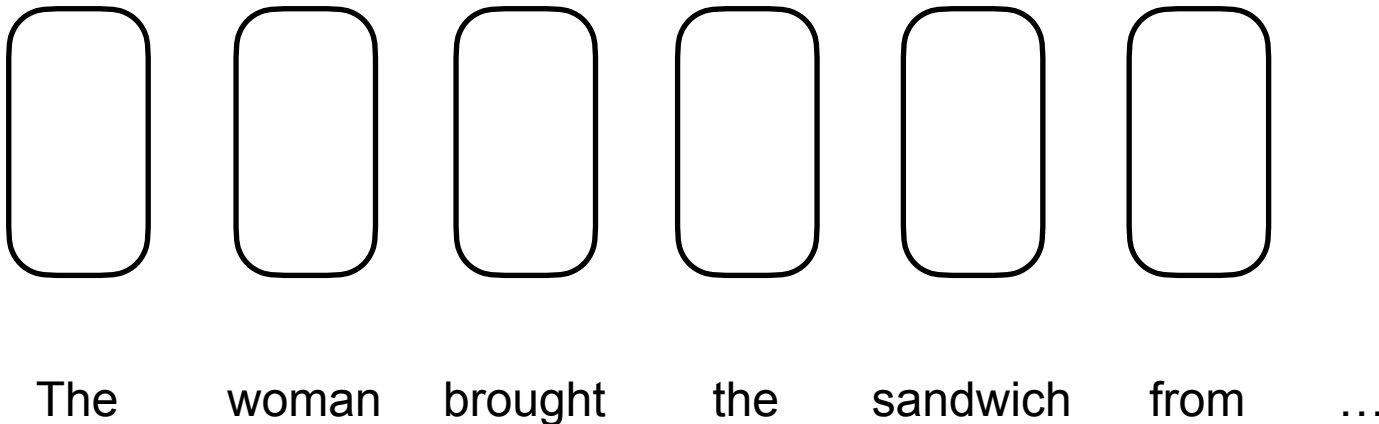
# Recurrent neural networks for language

- Draw inspiration from real-time nature of human language processing

- Previous inputs must be integrated and remembered all together in a uniform representational space

*(Jordan, 1986; Elman, 1990)*

# Recurrent neural networks for language

- Draw inspiration from real-time nature of human language processing
- Previous inputs must be integrated and remembered all together in a uniform representational space

The      woman      brought      the      sandwich      from      …

*(Jordan, 1986; Elman, 1990)*

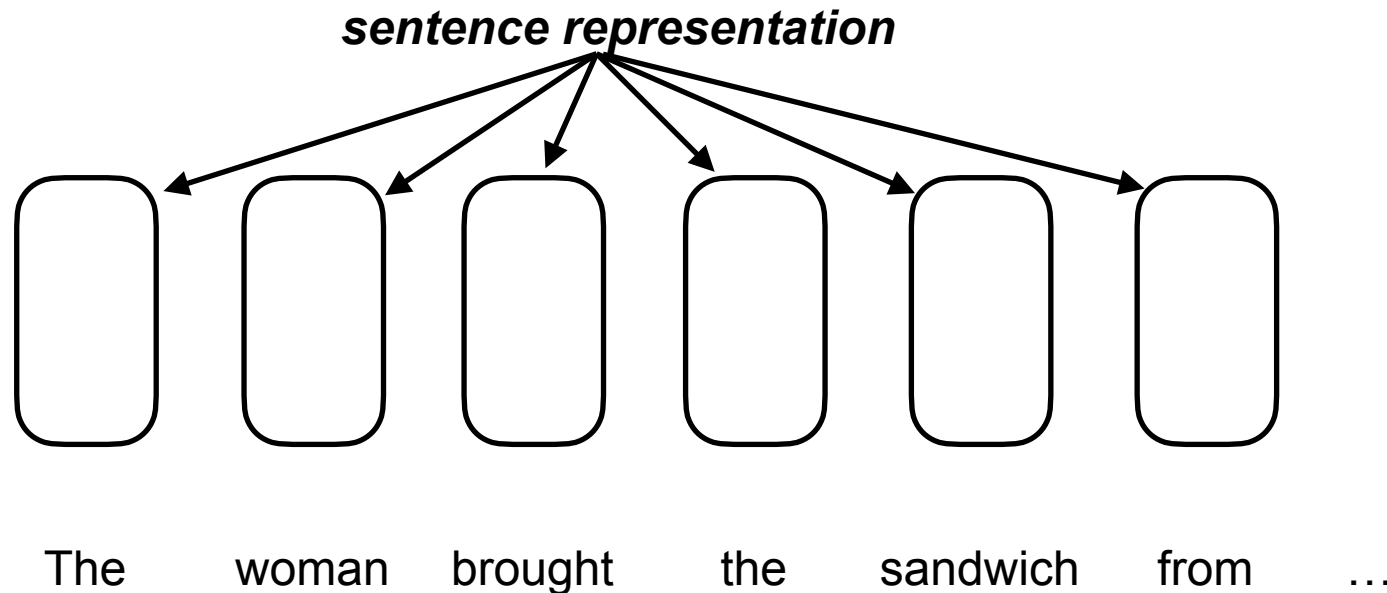# Recurrent neural networks for language

- Draw inspiration from real-time nature of human language processing
- Previous inputs must be integrated and remembered all together in a uniform representational space



The    woman    brought    the    sandwich    from    …

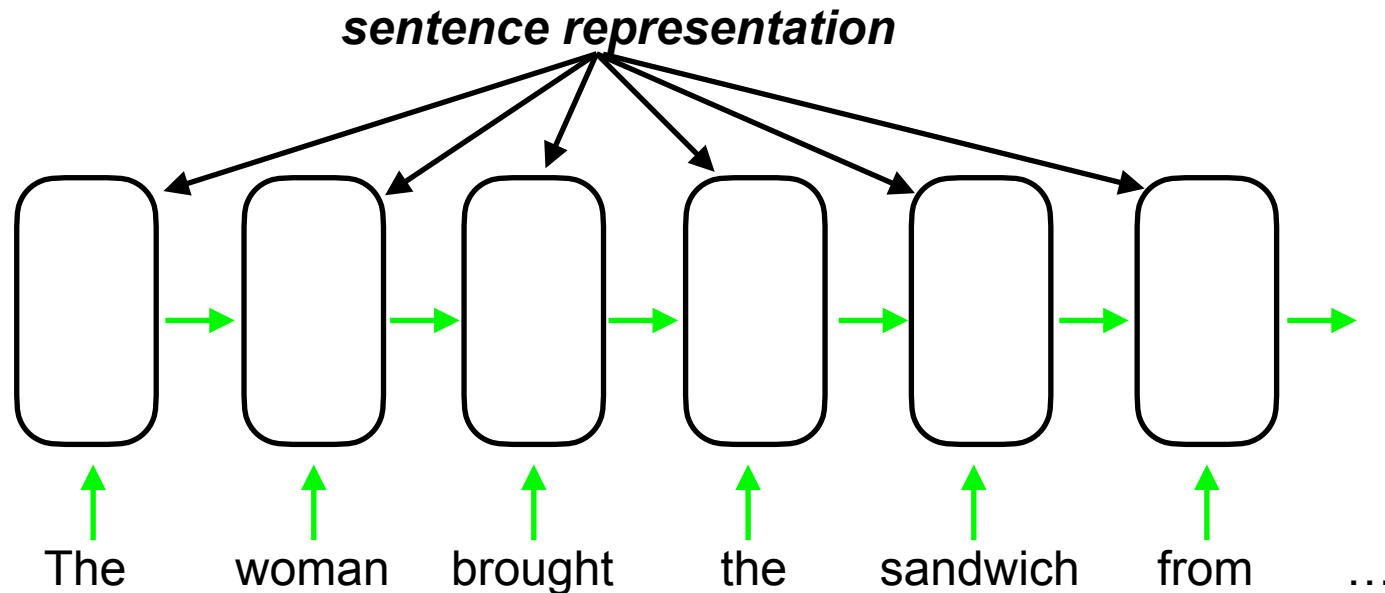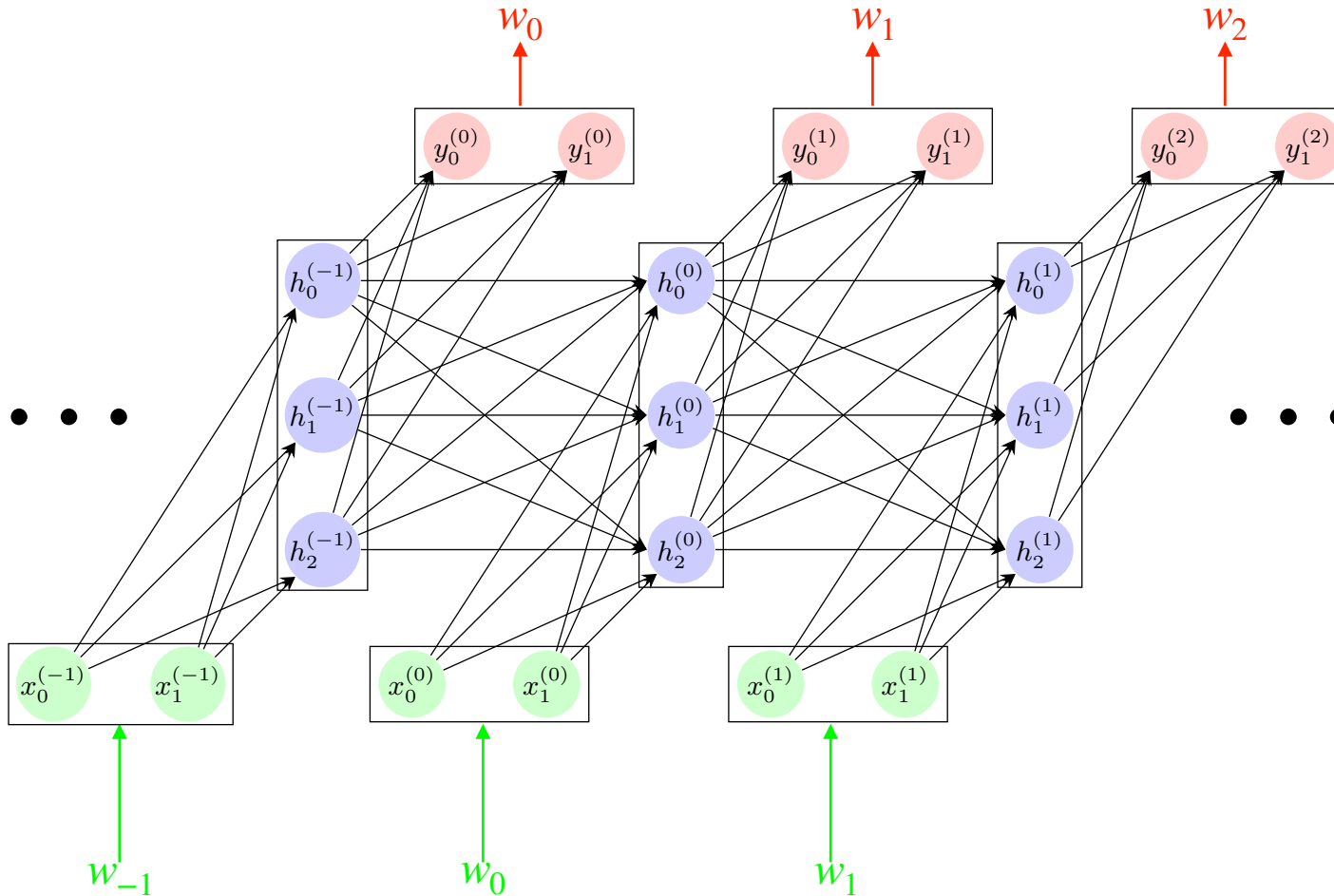*(Jordan, 1986; Elman, 1990)*

# Recurrent neural networks for language

- Draw inspiration from real-time nature of human language processing
- Previous inputs must be integrated and remembered all together in a uniform representational space

*sentence representation*

The     woman     brought     the     sandwich     from     …

*(Jordan, 1986; Elman, 1990)*

# Recurrent neural networks for language

- Draw inspiration from real-time nature of human language processing
- Previous inputs must be integrated and remembered all together in a uniform representational space
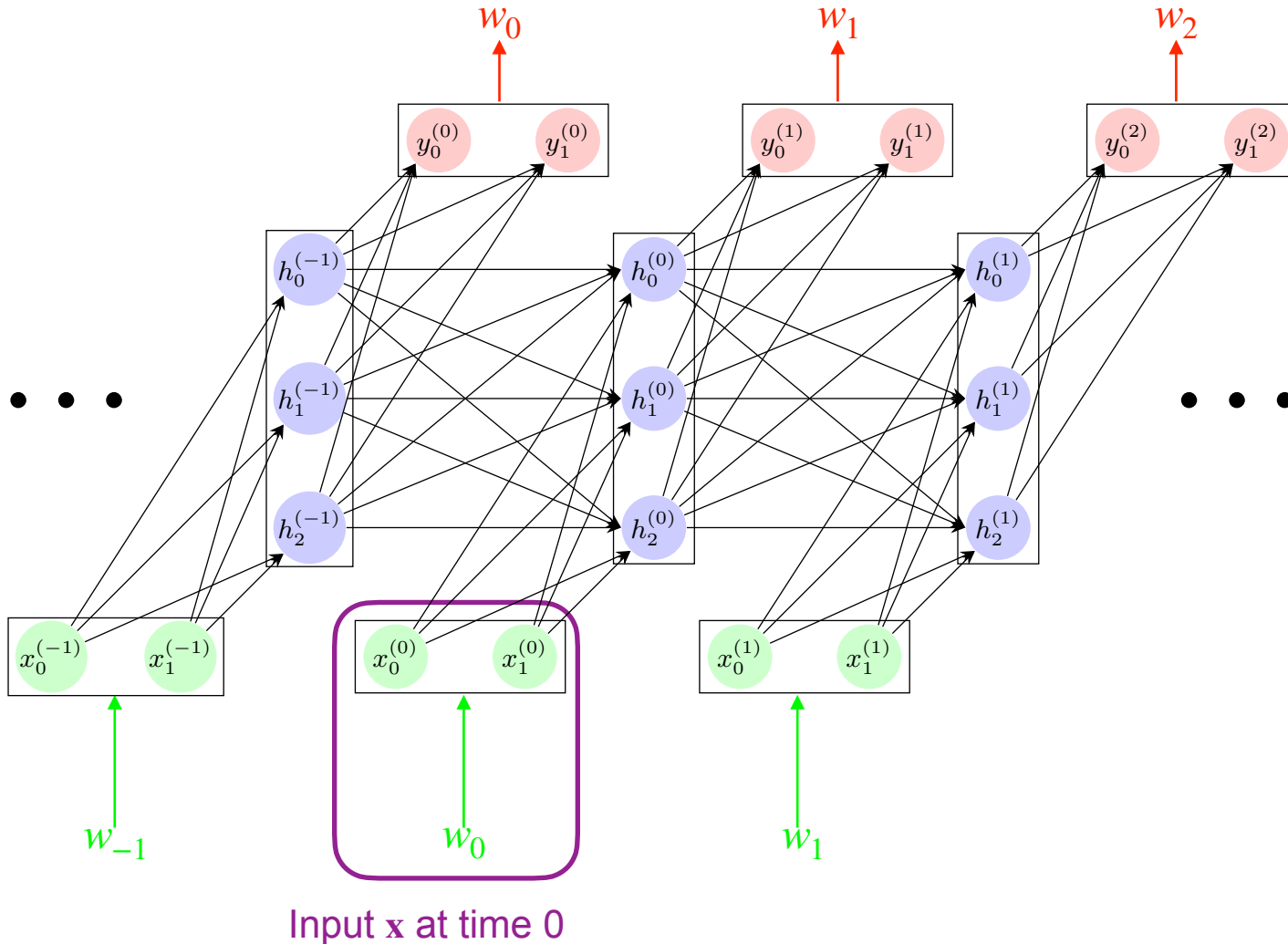
*sentence representation*



The    woman    brought    the    sandwich    from    …

*(Jordan, 1986; Elman, 1990)*

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



(Elman, 1990)

# The Simple Recurrent Network (SRN)

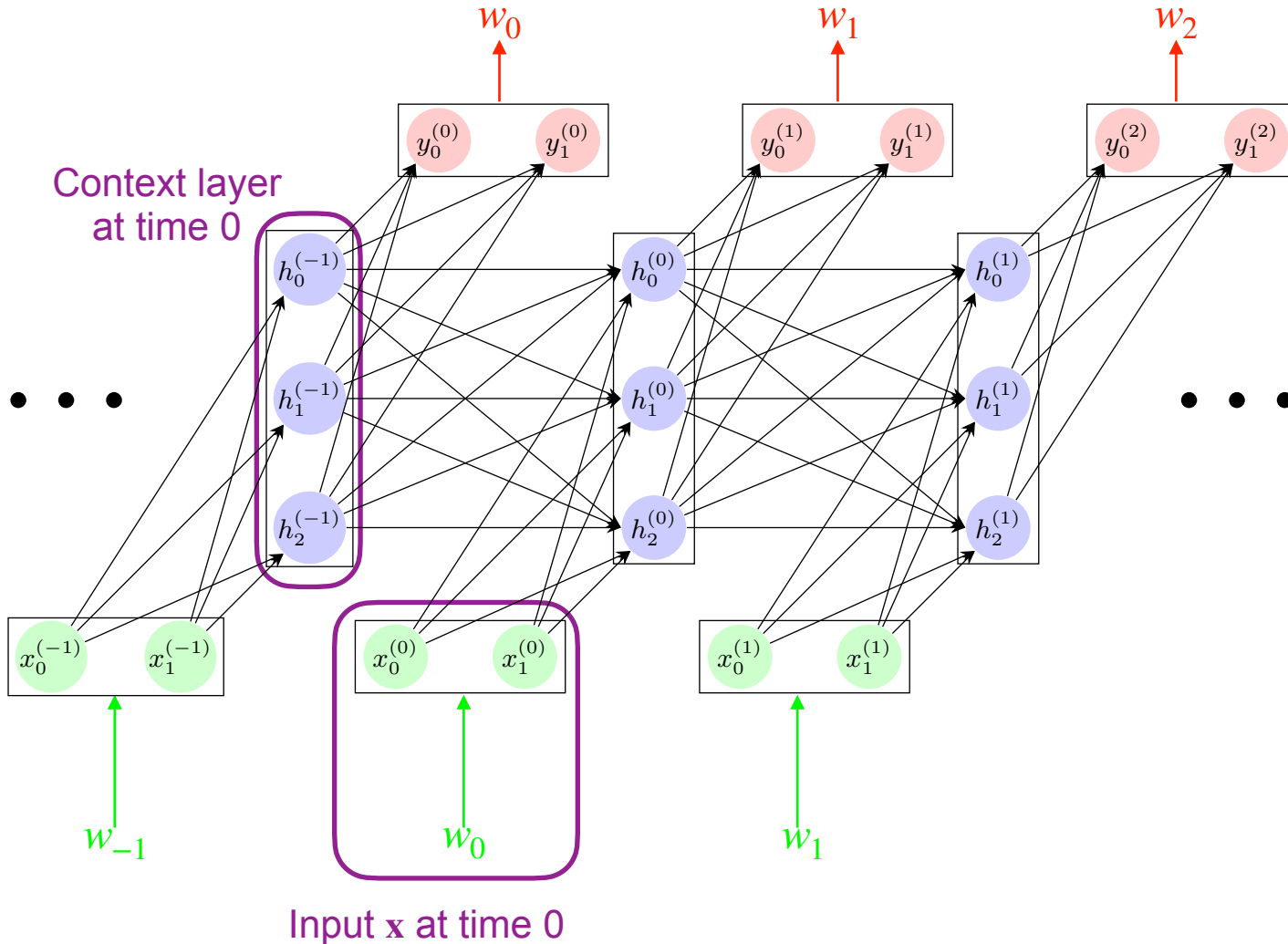(bias nodes not shown)



Input **x** at time 0

(Elman, 1990)
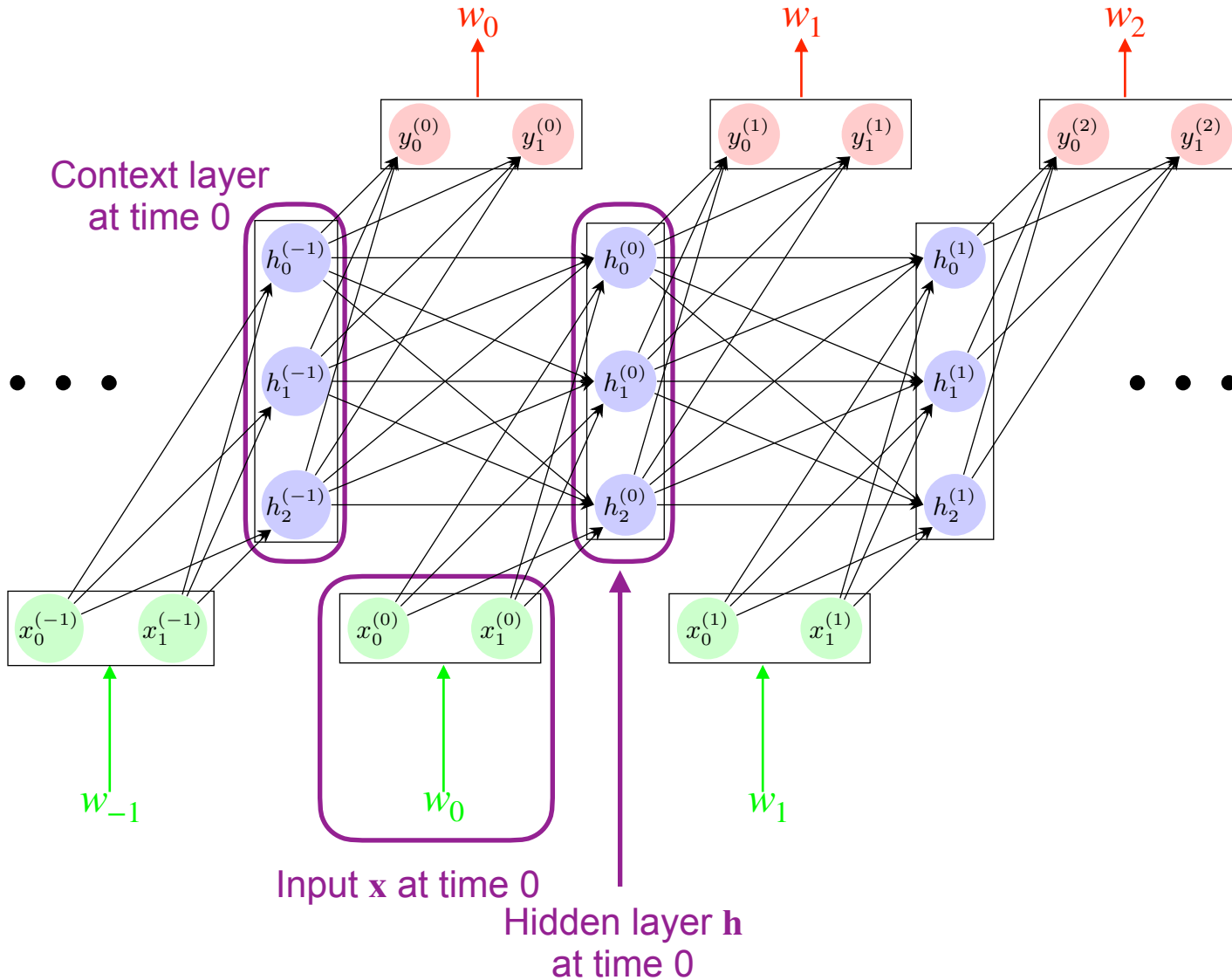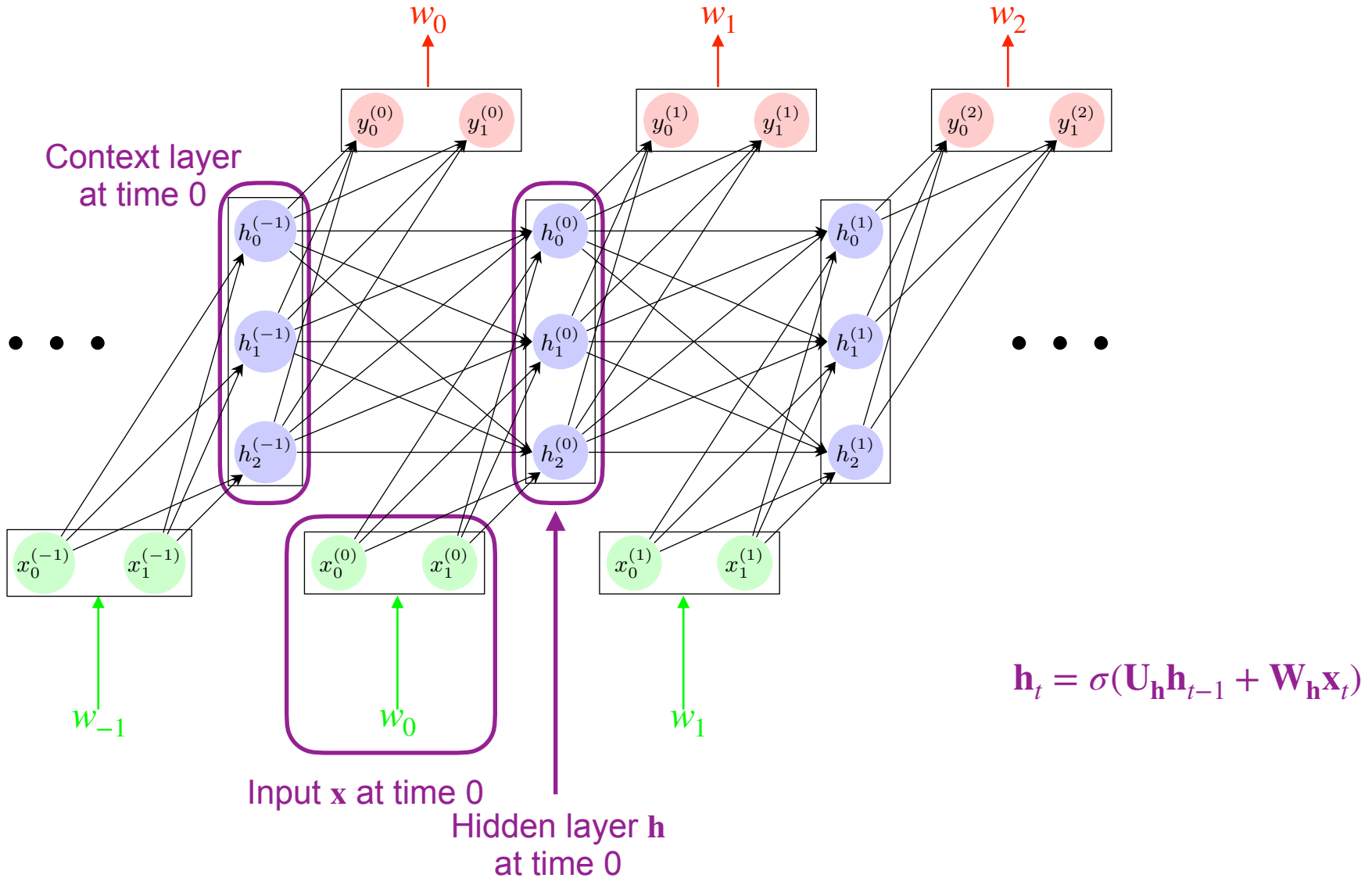
# The Simple Recurrent Network (SRN)

(bias nodes not shown)



(Elman, 1990)

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



Context layer at time 0

Input **x** at time 0

Hidden layer **h** at time 0

(Elman, 1990)

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

(Elman, 1990)
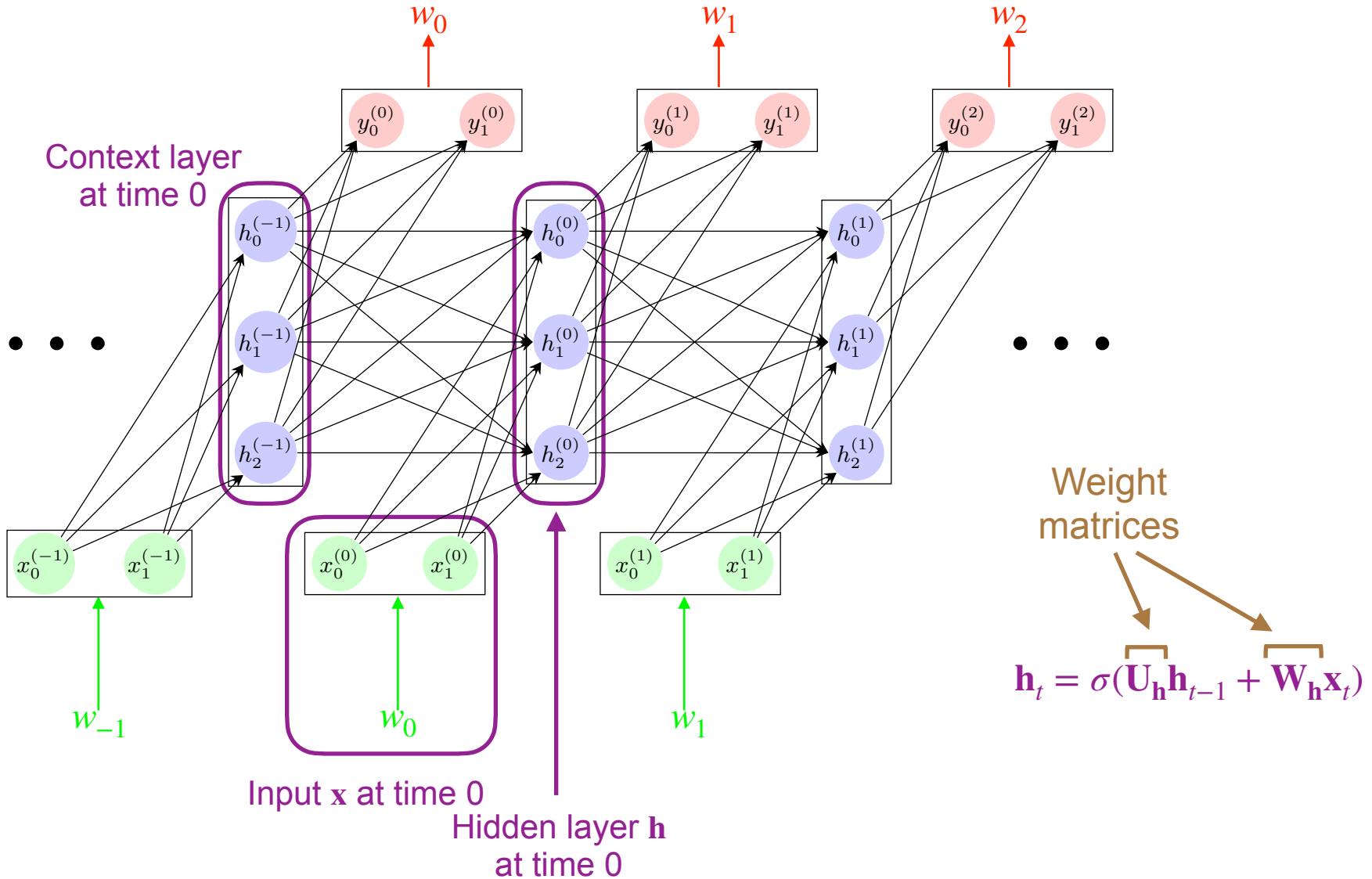
8

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

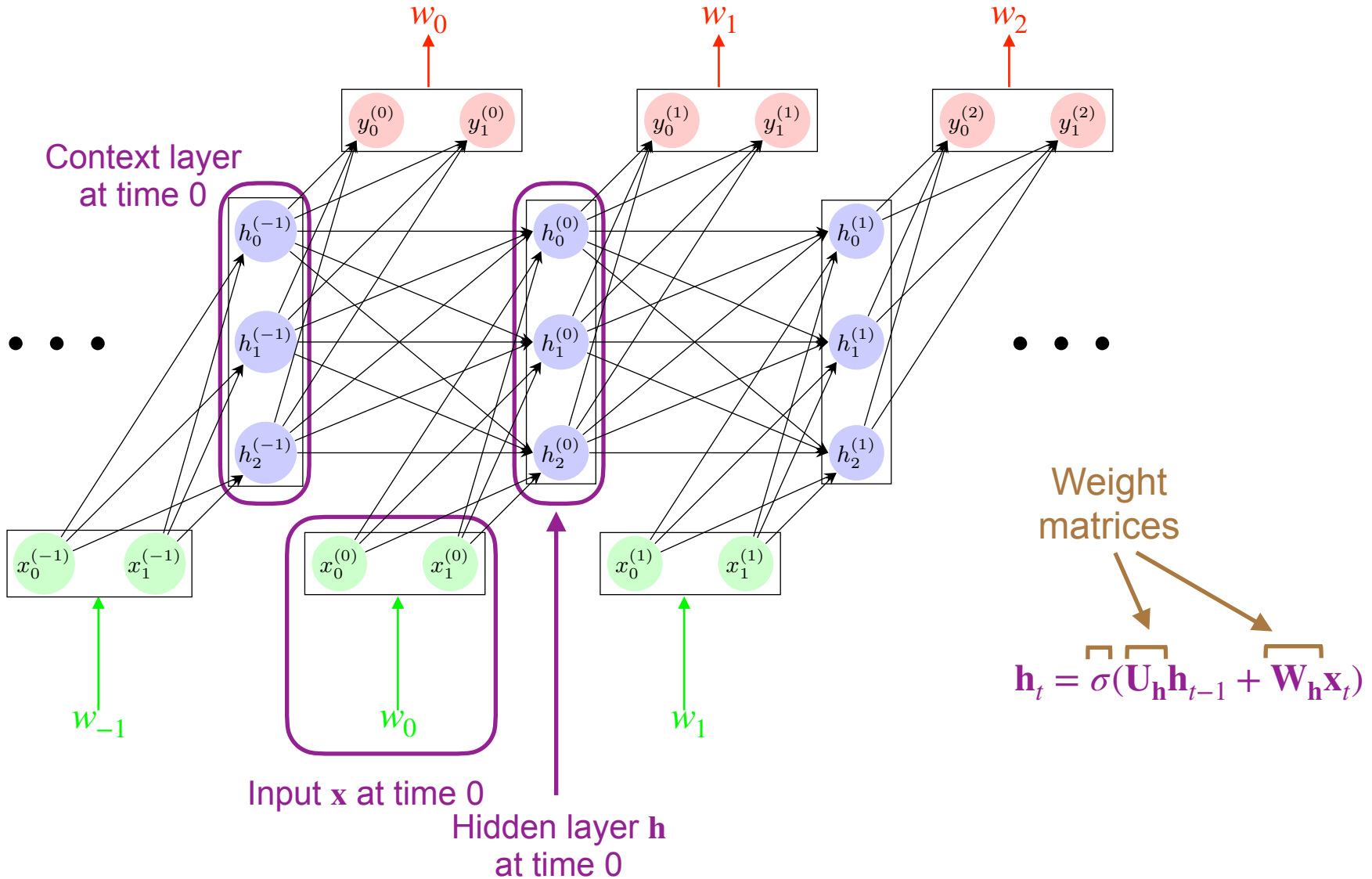(Elman, 1990)

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

Weight matrices

Context layer at time 0

Input **x** at time 0

Hidden layer **h** at time 0

(Elman, 1990)

8

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



Context layer at time 0

Weight matrices

$$\mathbf{h}_t = \sigma(\mathbf{U_h h}_{t-1} + \mathbf{W_h x}_t)$$

Input $\mathbf{x}$ at time 0

Hidden layer $\mathbf{h}$ at time 0

(Elman, 1990)

8

# The Simple Recurrent Network (SRN)

(bias nodes not shown)



$w_0$      $w_1$      $w_2$

Context layer at time 0

$h_0^{(-1)}$   $h_1^{(-1)}$   $h_2^{(-1)}$

$h_0^{(0)}$   $h_1^{(0)}$   $h_2^{(0)}$

$h_0^{(1)}$   $h_1^{(1)}$   $h_2^{(1)}$

$y_0^{(0)}$   $y_1^{(0)}$   $y_0^{(1)}$   $y_1^{(1)}$   $y_0^{(2)}$   $y_1^{(2)}$

$x_0^{(-1)}$   $x_1^{(-1)}$   $x_0^{(0)}$   $x_1^{(0)}$   $x_0^{(1)}$   $x_1^{(1)}$

$w_{-1}$      $w_0$      $w_1$

Input **x** at time 0

Hidden layer **h** at time 0

Weight matrices

Non-linear activation

$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

(Elman, 1990)

8

# The Simple Recurrent Network (SRN)

(bias nodes not shown)

Next-word prediction $\mathbf{y}$ at time 0

$w_0$ $\qquad$ $w_1$ $\qquad$ $w_2$

$$y_0^{(0)} \quad y_1^{(0)} \qquad y_0^{(1)} \quad y_1^{(1)} \qquad y_0^{(2)} \quad y_1^{(2)}$$

Context layer at time 0

$h_0^{(-1)} \quad h_0^{(0)} \quad h_0^{(1)}$

$h_1^{(-1)} \quad h_1^{(0)} \quad h_1^{(1)}$

$h_2^{(-1)} \quad h_2^{(0)} \quad h_2^{(1)}$

$x_0^{(-1)} \quad x_1^{(-1)} \qquad x_0^{(0)} \quad x_1^{(0)} \qquad x_0^{(1)} \quad x_1^{(1)}$

$w_{-1}$ $\qquad$ $w_0$ $\qquad$ $w_1$

Input $\mathbf{x}$ at time 0

Hidden layer $\mathbf{h}$ at time 0

Weight matrices

Non-linear activation

$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

(Elman, 1990)

8

# The Simple Recurrent Network (SRN)

(bias nodes not shown)

Next-word prediction $\mathbf{y}$ at time 0

$w_0$

$w_1$

$w_2$

$y_0^{(0)}$  $y_1^{(0)}$

$y_0^{(1)}$  $y_1^{(1)}$

$y_0^{(2)}$  $y_1^{(2)}$

Context layer at time 0

$h_0^{(-1)}$

$h_1^{(-1)}$

$h_2^{(-1)}$

$h_0^{(0)}$

$h_1^{(0)}$

$h_2^{(0)}$

$h_0^{(1)}$

$h_1^{(1)}$

$h_2^{(1)}$

$x_0^{(-1)}$  $x_1^{(-1)}$

$x_0^{(0)}$  $x_1^{(0)}$

$x_0^{(1)}$  $x_1^{(1)}$

$w_{-1}$

$w_0$

$w_1$

Input $\mathbf{x}$ at time 0

Hidden layer $\mathbf{h}$ at time 0

Weight matrices

Non-linear activation

$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

$$\mathbf{y}_{t+1} = \mathbf{W_y}\mathbf{h}_t$$

(Elman, 1990)

8

# The Simple Recurrent Network (SRN)

(bias nodes not shown)

Next-word prediction $\mathbf{y}$ at time 0

$w_0$

$w_1$

$w_2$

Context layer at time 0

$y_0^{(0)}$  $y_1^{(0)}$

$y_0^{(1)}$  $y_1^{(1)}$

$y_0^{(2)}$  $y_1^{(2)}$

$h_0^{(-1)}$

$h_1^{(-1)}$

$h_2^{(-1)}$

$h_0^{(0)}$

$h_1^{(0)}$

$h_2^{(0)}$

$h_0^{(1)}$

$h_1^{(1)}$

$h_2^{(1)}$

$x_0^{(-1)}$  $x_1^{(-1)}$

$x_0^{(0)}$  $x_1^{(0)}$

$x_0^{(1)}$  $x_1^{(1)}$

$w_{-1}$

$w_0$

$w_1$

Input $\mathbf{x}$ at time 0

Hidden layer $\mathbf{h}$ at time 0

Weight matrices

Non-linear activation

$$\mathbf{h}_t = \sigma(\mathbf{U_h}\mathbf{h}_{t-1} + \mathbf{W_h}\mathbf{x}_t)$$

$$\mathbf{y}_{t+1} = \mathbf{W_y}\mathbf{h}_t$$

$$P(w_t \mid \text{Context}) = \text{softmax}(\mathbf{y}_t)$$

(Elman, 1990)

8

# SRN "rolled up" and unrolled

- A "rolled-up" representation (Elman, 1990); and unrolled:



output units

hidden units

input units          context units

# SRN "rolled up" and unrolled

- A "rolled-up" representation (Elman, 1990); and unrolled:

output units

hidden units

input units          context units

# SRN "rolled up" and unrolled

- A "rolled-up" representation (Elman, 1990); and unrolled:

output units

hidden units

input units          context units

*sentence representation*

The      woman      brought      the      sandwich      from      …

# SRN "rolled up" and unrolled

- A "rolled-up" representation (Elman, 1990); and unrolled:

output units

hidden units

input units  context units

*sentence representation*

The   woman   brought   the   sandwich   from   …

*Predict!*

**Learning with artificial language input**

### TABLE 3
#### Categories of Lexical Items Used in Sentence Simulation

| Category | Examples |
|---|---|
| NOUN-HUM | man, woman |
| NOUN-ANIM | cat, mouse |
| NOUN-INANIM | book, rock |
| NOUN-AGRESS | dragon, monster |
| NOUN-FRAG | glass, plate |
| NOUN-FOOD | cookie, break |
| VERB-INTRAN | think, sleep |
| VERB-TRAN | see, chase |
| VERB-AGPAT | move, break |
| VERB-PERCEPT | smell, see |
| VERB-DESTROY | break, smash |
| VERB-EAT | eat |

### TABLE 4
#### Templates for Sentence Generator

| WORD 1 | WORD 2 | WORD 3 |
|---|---|---|
| NOUN-HUM | VERB-EAT | NOUN-FOOD |
| NOUN-HUM | VERB-PERCEPT | NOUN-INANIM |
| NOUN-HUM | VERB-DESTROY | NOUN-FRAG |
| NOUN-HUM | VERB-INTRAN | |
| NOUN-HUM | VERB-TRAN | NOUN-HUM |
| NOUN-HUM | VERB-AGPAT | NOUN-INANIM |
| NOUN-HUM | VERB-AGPAT | |
| NOUN-ANIM | VERB-EAT | NOUN-FOOD |
| NOUN-ANIM | VERB-TRAN | NOUN-ANIM |
| NOUN-ANIM | VERB-AGPAT | NOUN-INANIM |
| NOUN-ANIM | VERB-AGPAT | |
| NOUN-INANIM | VERB-AGPAT | |
| NOUN-AGRESS | VERB-DESTROY | NOUN-FRAG |
| NOUN-AGRESS | VERB-EAT | NOUN-HUM |
| NOUN-AGRESS | VERB-EAT | NOUN-ANIM |
| NOUN-AGRESS | VERB-EAT | NOUN-FOOD |

*(Elman, 1990)*

10

# Used *localist* word representations

**Fragment of Training Sequences for Sentence Simulation**

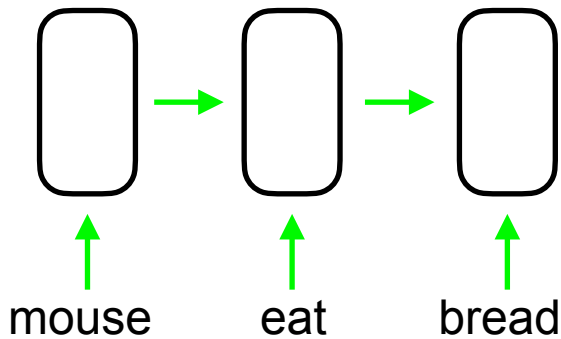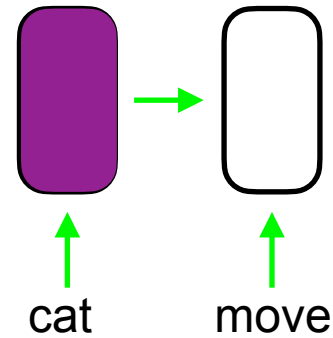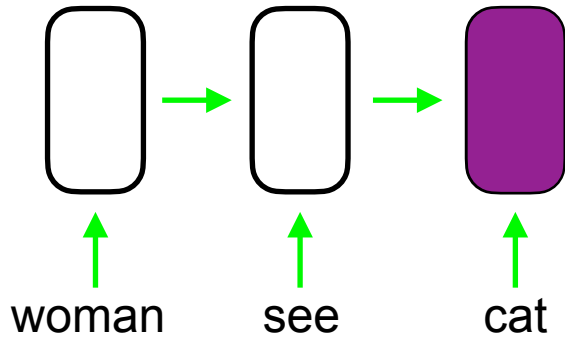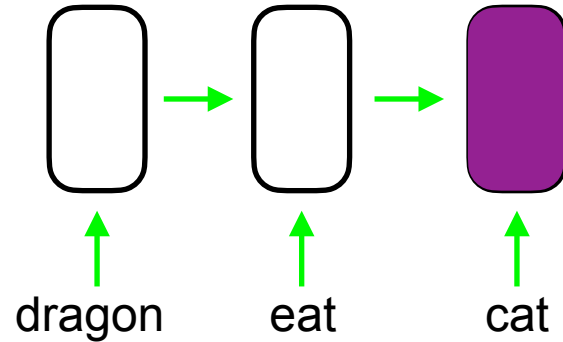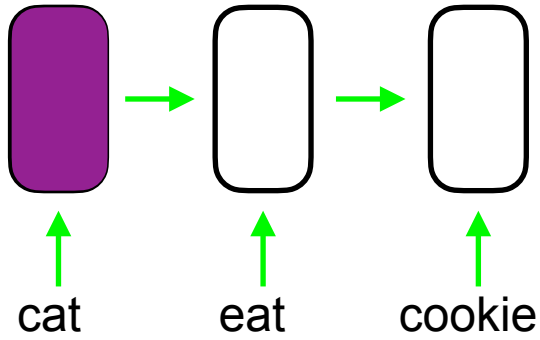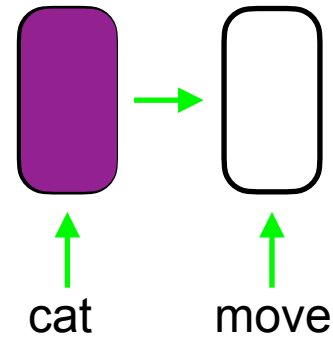| Input | Output |
|-------|--------|
| 00000000000000000000000000000010 (*woman*) | 00000000000000000000000000010000 (*smash*) |
| 00000000000000000000000000010000 (*smash*) | 00000000000000000000001000000000 (*plate*) |
| 00000000000000000000001000000000 (*plate*) | 00000100000000000000000000000000 (*cat*) |
| 00000100000000000000000000000000 (*cat*) | 00000000000000000010000000000000 (*move*) |
| 00000000000000000010000000000000 (*move*) | 00000000000000001000000000000000 (*man*) |
| 00000000000000001000000000000000 (*man*) | 00010000000000000000000000000000 (*break*) |
| 00010000000000000000000000000000 (*break*) | 00001000000000000000000000000000 (*car*) |
| 00001000000000000000000000000000 (*car*) | 01000000000000000000000000000000 (*boy*) |
| 01000000000000000000000000000000 (*boy*) | 00000000000000000010000000000000 (*move*) |
| 00000000000000000010000000000000 (*move*) | 00000000000010000000000000000000 (*girl*) |
| 00000000000010000000000000000000 (*girl*) | 00000000001000000000000000000000 (*eat*) |
| 00000000001000000000000000000000 (*eat*) | 00100000000000000000000000000000 (*bread*) |
| 00100000000000000000000000000000 (*bread*) | 00000001000000000000000000000000 (*dog*) |
| 00000001000000000000000000000000 (*dog*) | 00000000000000000010000000000000 (*move*) |
| 00000000000000000010000000000000 (*move*) | 00000000000000000001000000000000 (*mouse*) |
| 00000000000000000001000000000000 (*mouse*) | 00000000000000000001000000000000 (*mouse*) |
| 00000000000000000001000000000000 (*mouse*) | 00000000000000000010000000000000 (*move*) |
| 00000000000000000010000000000000 (*move*) | 10000000000000000000000000000000 (*book*) |
| 10000000000000000000000000000000 (*book*) | 00000000000000010000000000000000 (*lion*) |

*(Elman, 1990)*

# Learning word classes



cat → eat → cookie

dragon → eat → cat

woman → see → cat

cat → move

mouse → eat → bread

cat → chase → mouse

12

# Learning word classes



cat → eat → cookie

dragon → eat → cat

woman → see → cat

cat → move

mouse → eat → bread

cat → chase → mouse

12

# Learning word classes



cat     eat     cookie

dragon     eat     cat

woman     see     cat

cat     move

mouse     eat     bread

cat     chase     mouse

12

# Discovered similarity structure of words



*(Elman, 1990)*

13

# Beyond the simple recurrent network

- The SRN has a very strong *linear locality bias*
- But natural language syntax is characterized by *hierarchical structure*
- SRNs can learn hierarchy (Elman, 1991), but *it is hard*—their inductive bias disfavors it

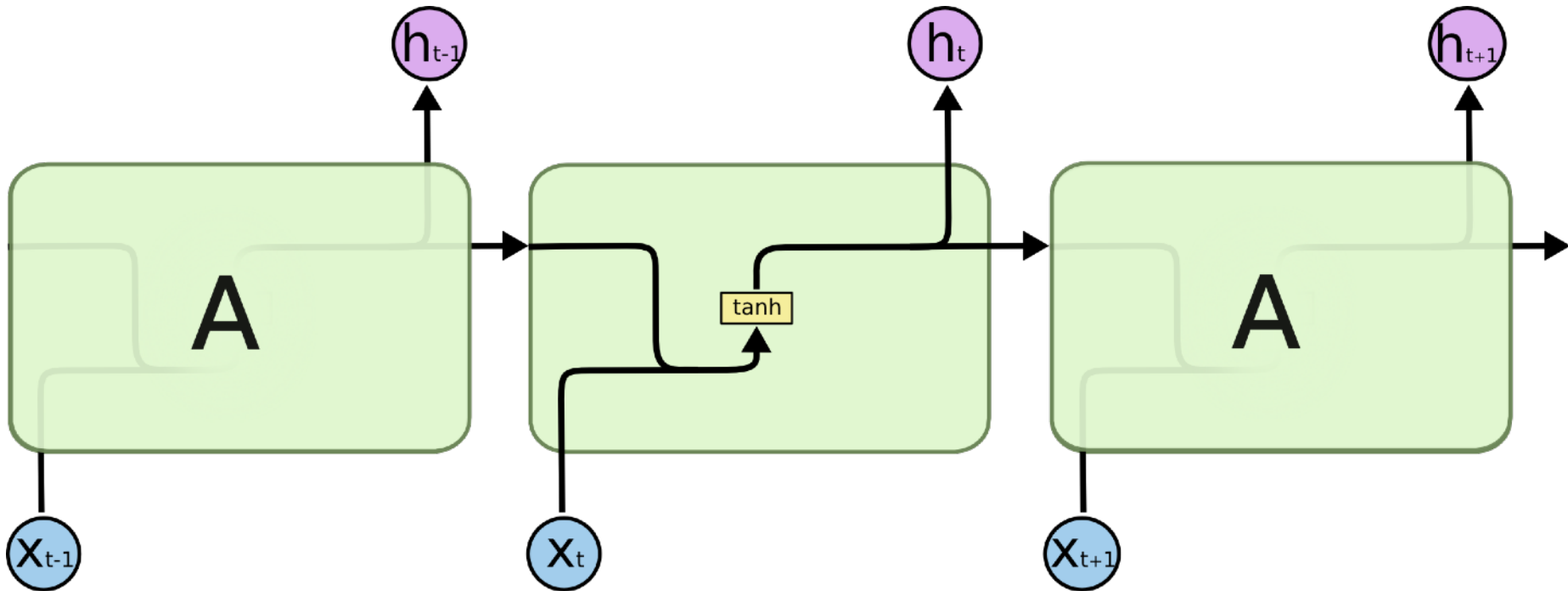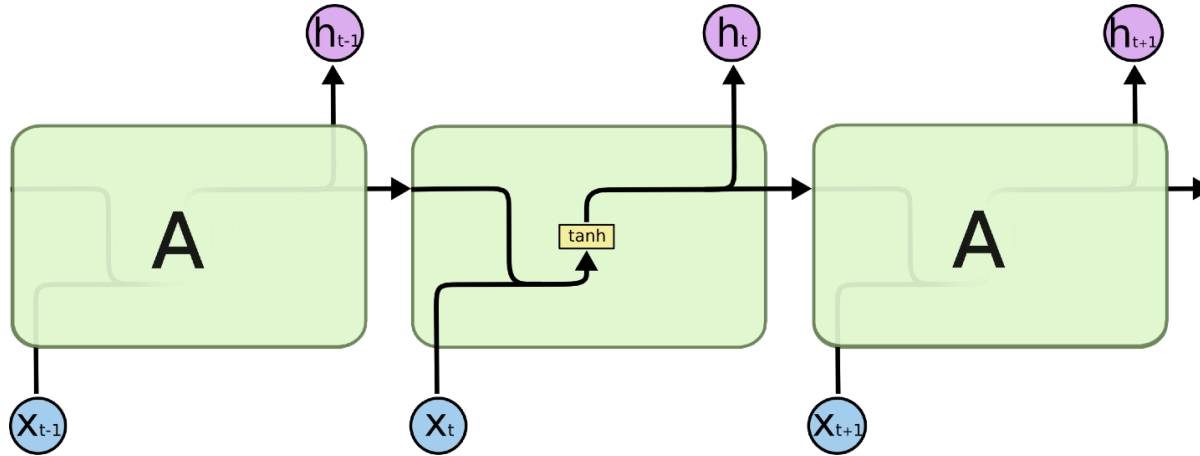*sentence representation*

The cat that the dog chased …

# More sophisticated recurrent units
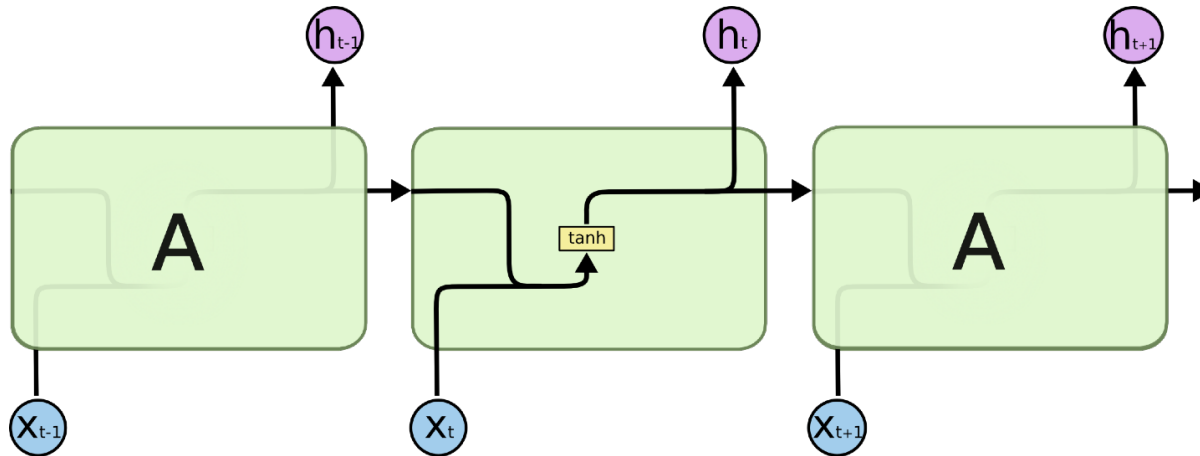
15

# More sophisticated recurrent units

- Another view of an unrolled SRN:

15

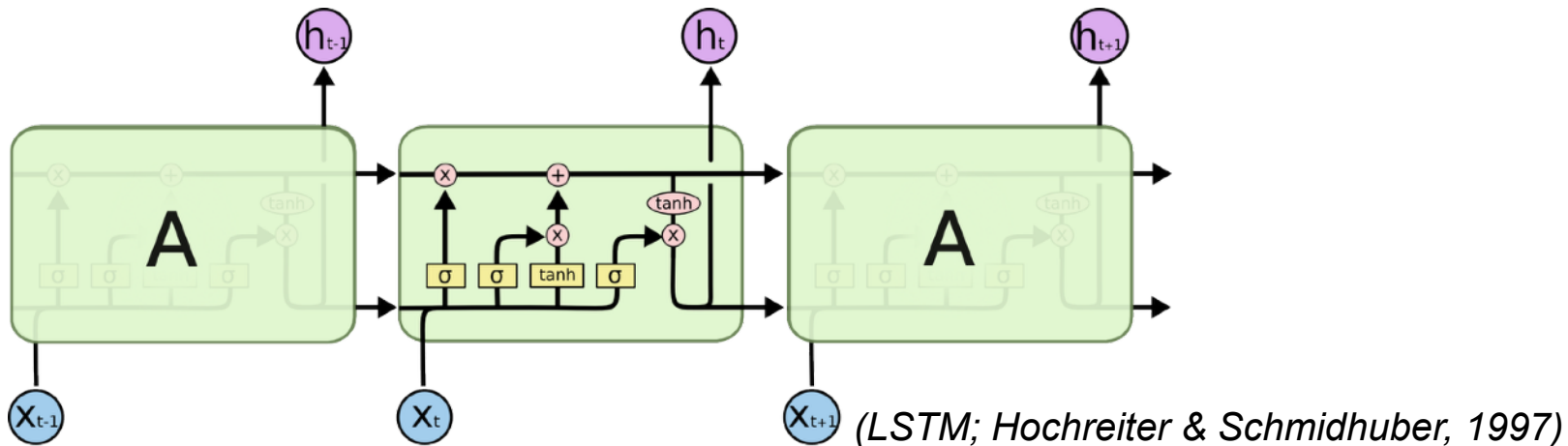# More sophisticated recurrent units

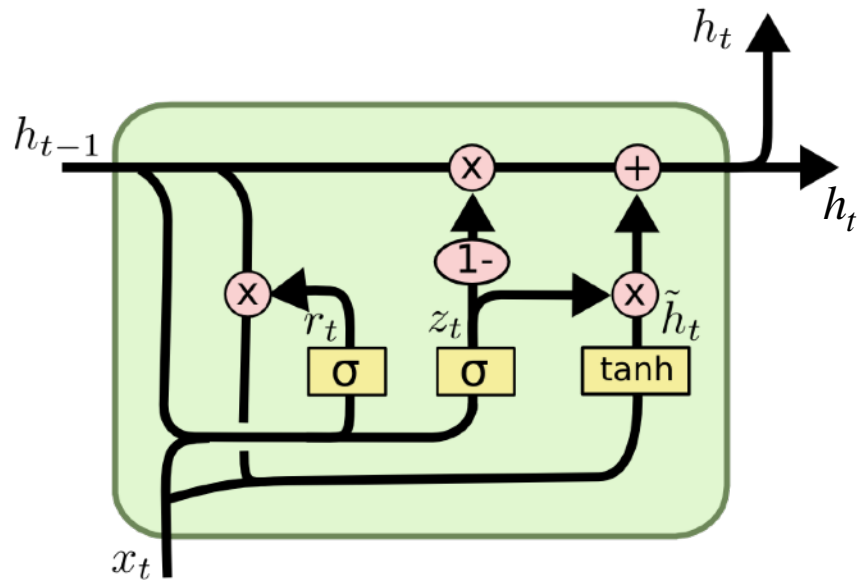- Another view of an unrolled SRN:

15

# More sophisticated recurrent units

- Another view of an unrolled SRN:



- Keep the recurrent structure and "swap in" a new unit:



*(LSTM; Hochreiter & Schmidhuber, 1997)*

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

15

# Gated Recurrent Unit (GRU) architecture

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

# Gated Recurrent Unit (GRU) architecture



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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W} \mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1})$$

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t$$
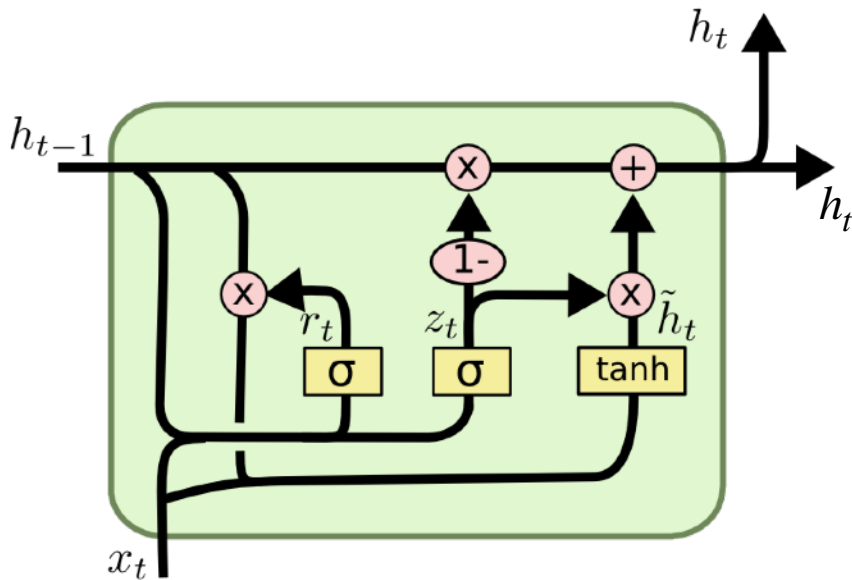
# Gated Recurrent Unit (GRU) architecture

*logistic/sigmoid activation function*

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1})$$

$$\mathbf{z}_t = \sigma(\mathbf{W}_z\mathbf{x}_t + \mathbf{U}_z\mathbf{h}_{t-1})$$

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

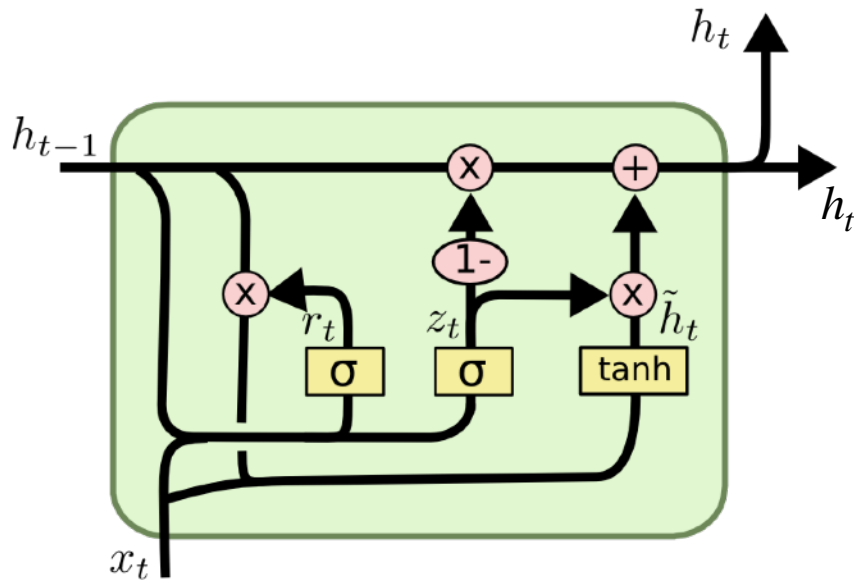# Gated Recurrent Unit (GRU) architecture

*logistic/sigmoid activation function*

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W} \mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1})$$



$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

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

# Gated Recurrent Unit (GRU) architecture

*logistic/sigmoid activation function*

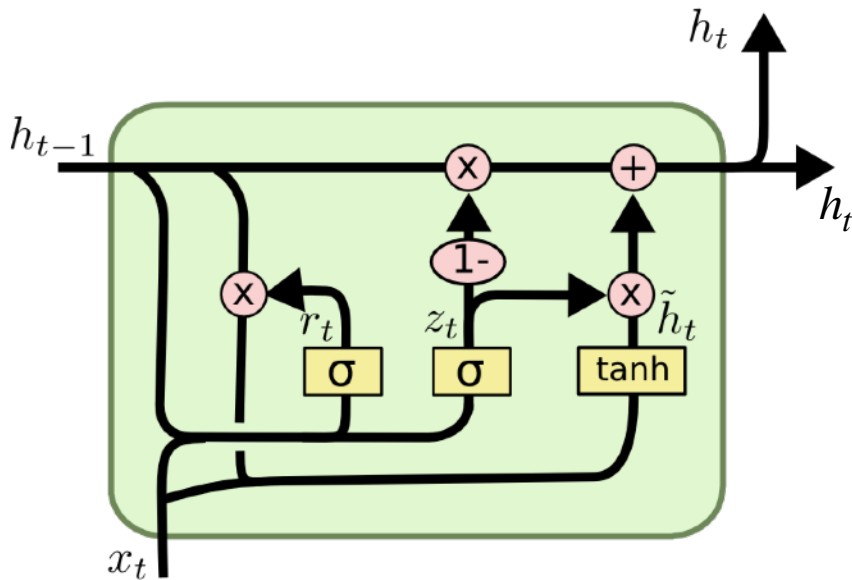$$\mathbf{r}_t = \sigma(\mathbf{W}_r \mathbf{x}_t + \mathbf{U}_r \mathbf{h}_{t-1})$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1}))$$
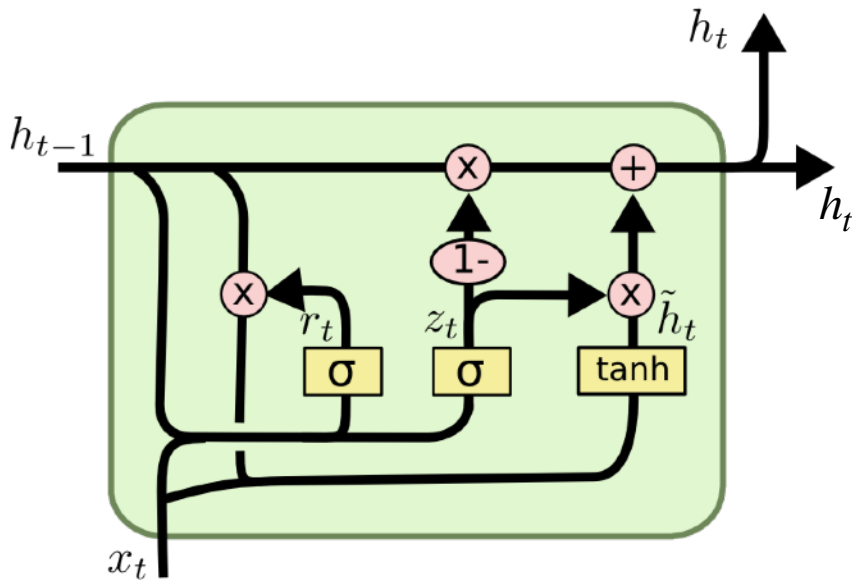
*element-wise multiplication*

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

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

(e.g., $\langle 1,2,3 \rangle \odot \langle 0.5,2,1 \rangle = \langle 0.5,4,3 \rangle$)

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

# Gated Recurrent Unit (GRU) architecture



*logistic/sigmoid activation function*

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1})$$
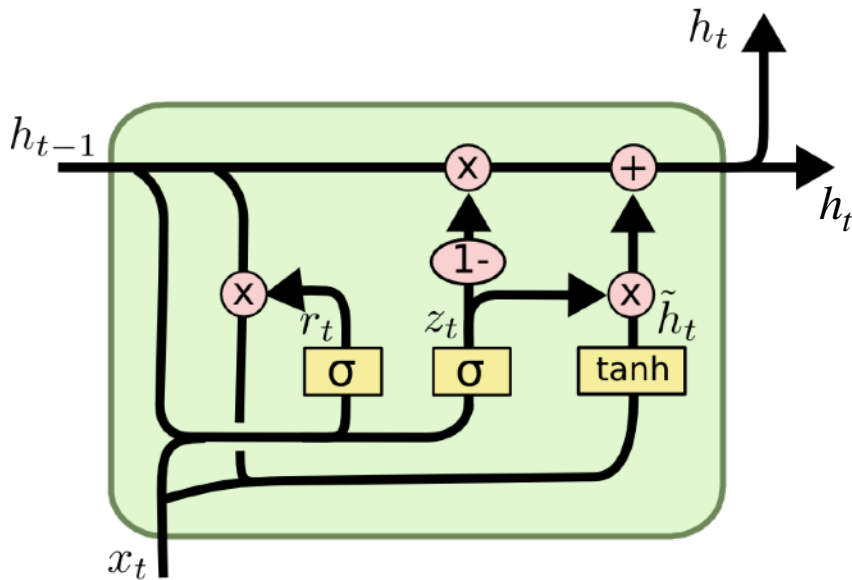
*element-wise multiplication*

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

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

(e.g., $\langle 1,2,3 \rangle \odot \langle 0.5,2,1 \rangle = \langle 0.5,4,3 \rangle$)

# Gated Recurrent Unit (GRU) architecture

*logistic/sigmoid activation function*

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1})$$
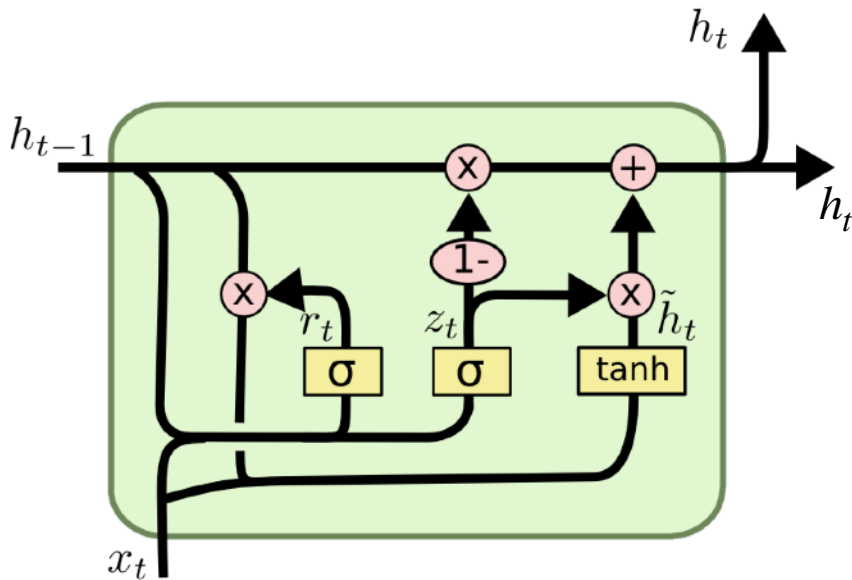
*element-wise multiplication*

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

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

(e.g., $\langle 1,2,3 \rangle \odot \langle 0.5,2,1 \rangle = \langle 0.5,4,3 \rangle$)

16

# Gated Recurrent Unit (GRU) architecture

*logistic/sigmoid activation function*

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1})$$
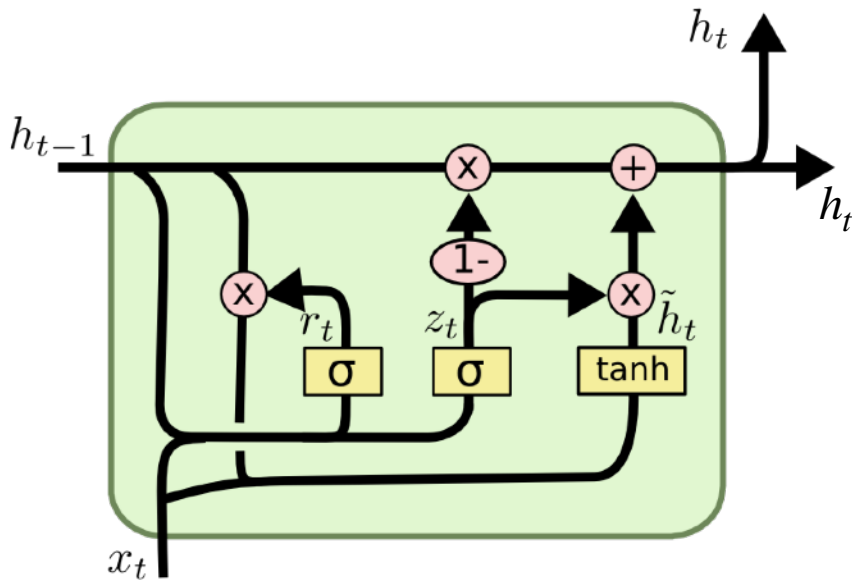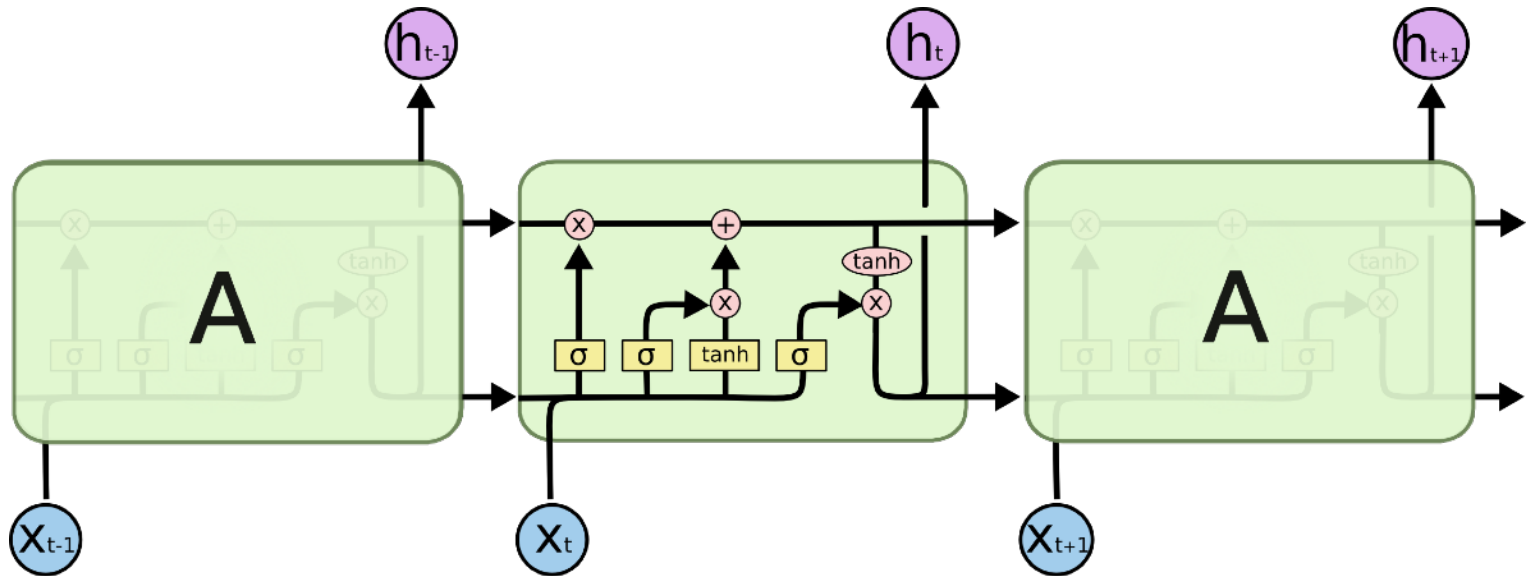
*element-wise multiplication*

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

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

(e.g., $\langle 1,2,3 \rangle \odot \langle 0.5,2,1 \rangle = \langle 0.5,4,3 \rangle$)

# Gated Recurrent Unit (GRU) architecture

*logistic/sigmoid activation function*

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1}))$$

*element-wise multiplication*

$h_{t-1}$

$h_t$

$h_t$

$r_t$

$z_t$

$\tilde{h}_t$

$\sigma$    $\sigma$    tanh

$x_t$

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1})$$

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

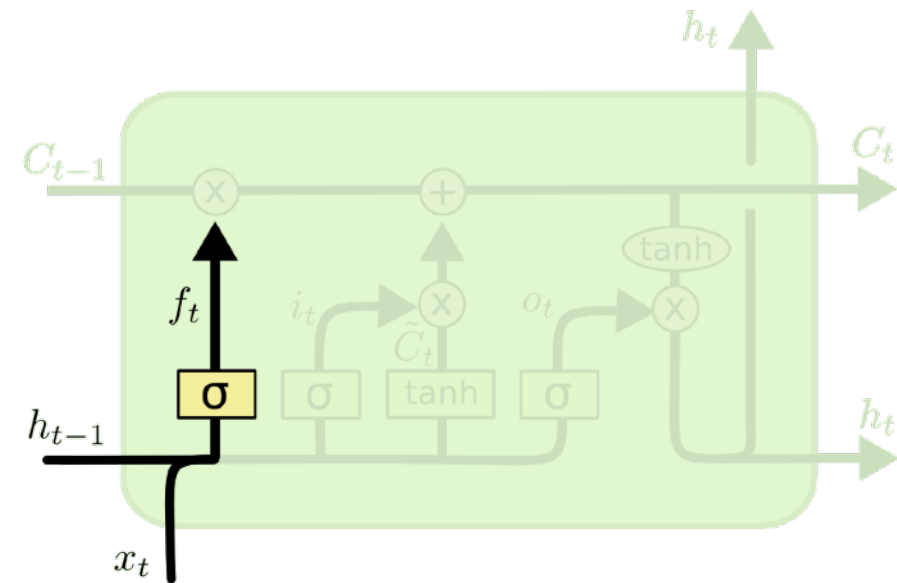(e.g., $\langle 1,2,3 \rangle \odot \langle 0.5,2,1 \rangle = \langle 0.5,4,3 \rangle$)

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

16

# Long short-term memory (LSTM) units



*(Hochreiter & Schmidhuber, 1997)*

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

# Inside the LSTM unit

- The "hidden layer" $\mathbf{h}_{t-1}$ was used to predict element $t$ of the sequence
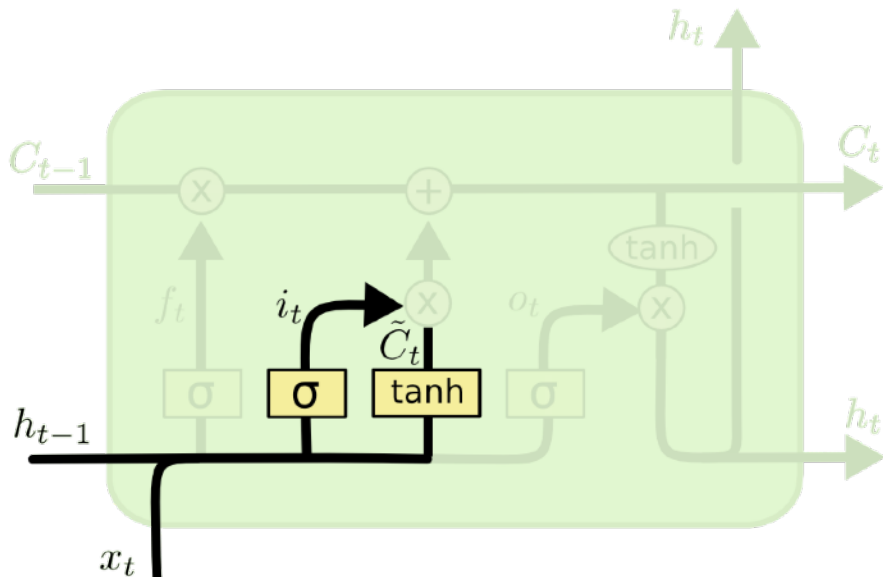
- It now gets passed through a "forget gate"

$$\mathbf{f}_t = \sigma(\mathbf{U_f}\mathbf{h}_{t-1} + \mathbf{W}_f\mathbf{x}_t)$$

*(Hochreiter & Schmidhuber, 1997)*

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

# Inside the LSTM unit

- Other information from $h_{t-1}$ gets put into the memory store

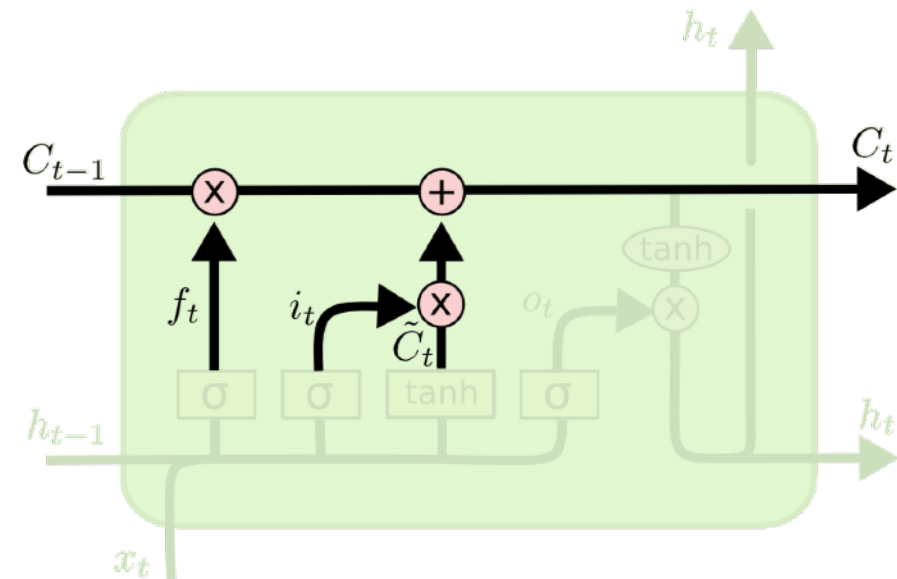$$\mathbf{i}_t = \sigma(\mathbf{U_i}\mathbf{h}_{t-1} + \mathbf{W_i}\mathbf{x}_t)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{U_C}\mathbf{h}_{t-1} + \mathbf{W_C}\mathbf{x}_t)$$

# Inside the LSTM unit

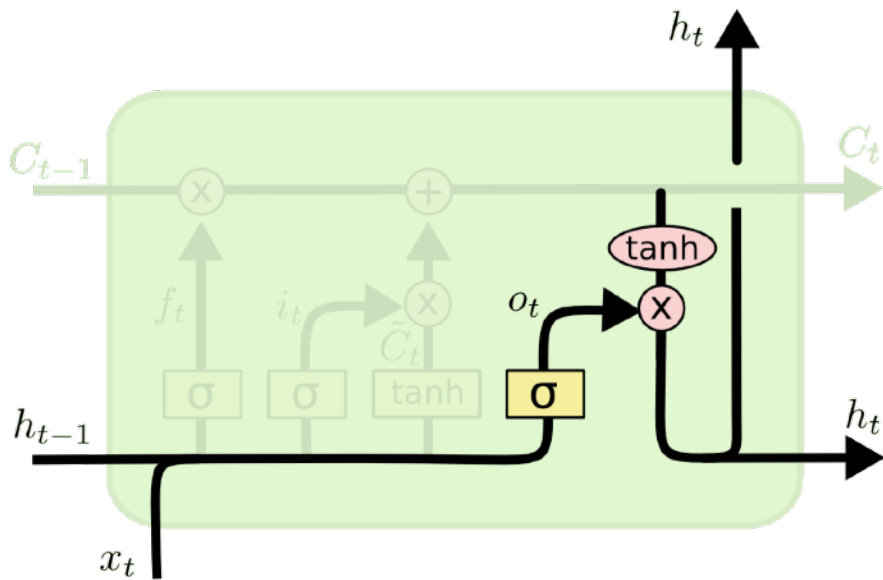- That information gets integrated into the memory store (which also gets passed on to the future



$$\mathbf{C}_t = \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{C}}_t$$
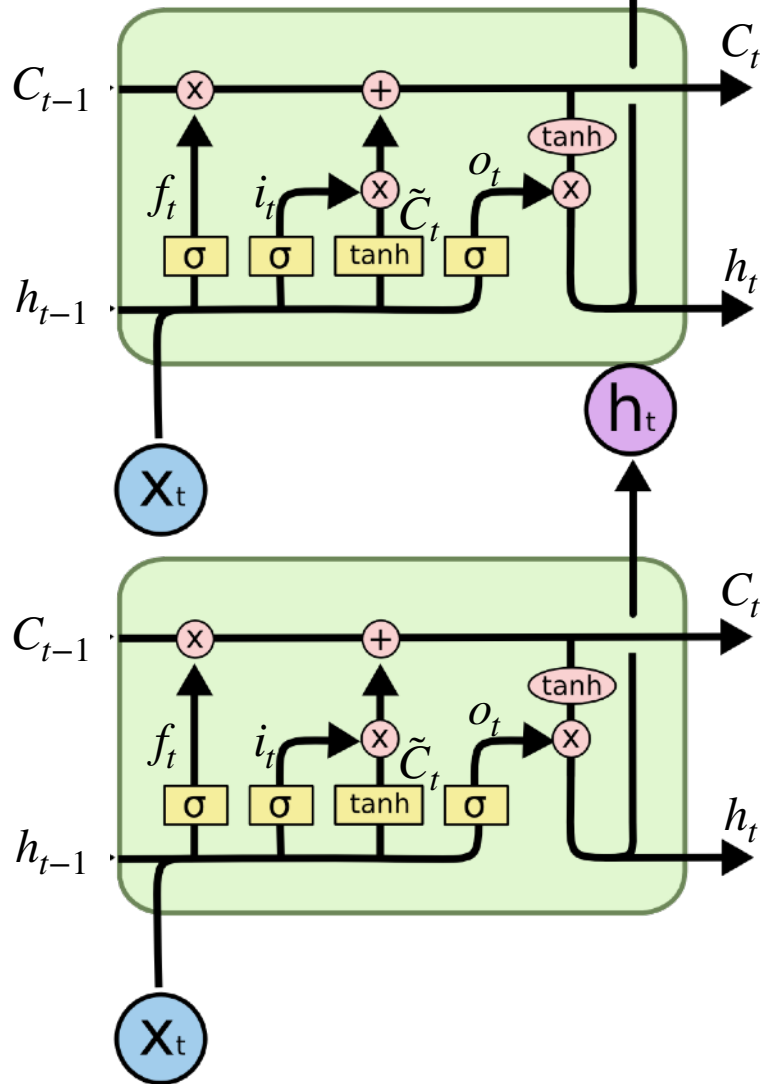
# Inside the LSTM unit

- Finally, we determine the new hidden layer to predict input *t+1*



$$\mathbf{o}_t = \sigma(\mathbf{U_o}\mathbf{h}_{t-1} + \mathbf{W_o}\mathbf{x}_t)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{C}_t)$$

# The LSTM unit, complete



$$\mathbf{f}_t = \sigma(\mathbf{U_f}\mathbf{h}_{t-1} + \mathbf{W}_f\mathbf{x}_t)$$

$$\mathbf{i}_t = \sigma(\mathbf{U_i}\mathbf{h}_{t-1} + \mathbf{W_i}\mathbf{x}_t)$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{U_C}\mathbf{h}_{t-1} + \mathbf{W_C}\mathbf{x}_t)$$

$$\mathbf{C}_t = \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{C}}_t$$

$$\mathbf{o}_t = \sigma(\mathbf{U_o}\mathbf{h}_{t-1} + \mathbf{W_o}\mathbf{x}_t)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{C}_t)$$

*visualization due to Christopher Olah, http://colah.github.io/posts/2015-08-Understanding-LSTMs/*

22

# Learning the classic counting language

$a^n b^n$

(Weiss et al. 2018, ACL)

# Learning the classic counting language

$a^n b^n$

Easily generable with a context-free grammar:

$$S \rightarrow a \ b$$
$$S \rightarrow a \ S \ b$$

(Weiss et al. 2018, ACL)

# Learning the classic counting language

$a^n b^n$      Easily generable with a context-free grammar:

```
^ab$

^aabb$

^aaabbb$

^aaaabbbb$

^aaaaabbbbb$

^aaaaaabbbbbb$

^aaaaaaabbbbbbb$
    ⋮

^aaaaaaaaaaaaaaaaaaaabbbbbbbbbbbbbbbbbbbb$
```

```
S  →  a  b
S  →  a  S  b
```

# Learning the classic counting language

$a^n b^n$      Easily generable with a context-free grammar:

$D_{train}$

```
^ab$

^aabb$

^aaabbb$

^aaaabbbb$

^aaaaabbbbb$

^aaaaaabbbbbb$

^aaaaaaabbbbbbb$

    ⋮

^aaaaaaaaaaaaaaaaaaaaabbbbbbbbbbbbbbbbbbbbbb$
```

```
S → a b
S → a S b
```

(Weiss et al. 2018, ACL)

# Learning the classic counting language

$a^n b^n$

Easily generable with a context-free grammar:

```
^ab$

^aabb$

^aaabbb$

^aaaabbbb$

^aaaaabbbbb$

^aaaaaabbbbbb$

^aaaaaaabbbbbbb$
        ⋮

^aaaaaaaaaaaaaaaaaaaabbbbbbbbbbbbbbbbbbbb$
```

```
S → a b
S → a S b
```

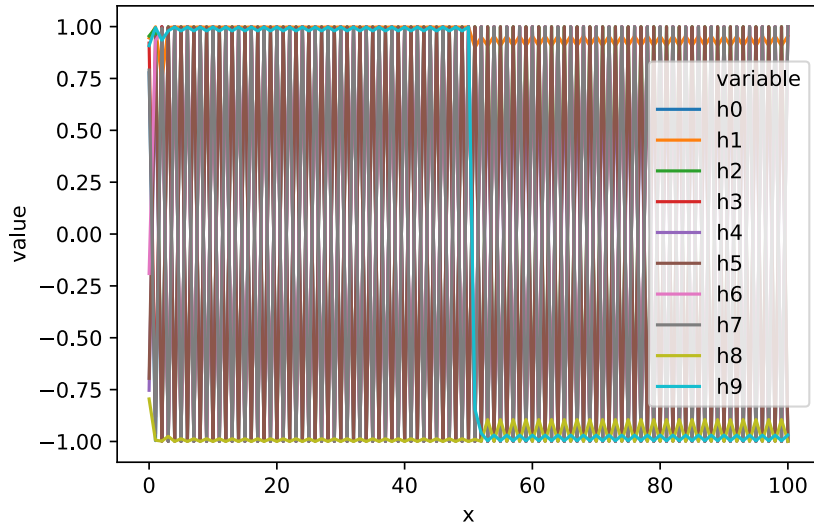*N=20*            *N=20*

(Weiss et al. 2018, ACL)                    23

# Training recurrent architectures on $a^n b^n$

# Hidden & cell state contents



**SRN hidden state**

**LSTM hidden state**

**GRU hidden state**

**LSTM cell state**

(Weiss et al. 2018, ACL)

# Summary

- Mechanisms for neural networks at the sentence level:
  - Learned word embeddings
  - Recurrent state representation
- Different units used for recurrent state representation:
  - Simple recurrent network (SRN)
  - Gated recurrent unit (GRU)
  - Long short-term memory (LSTM)
- For classic counting language, LSTM works the best