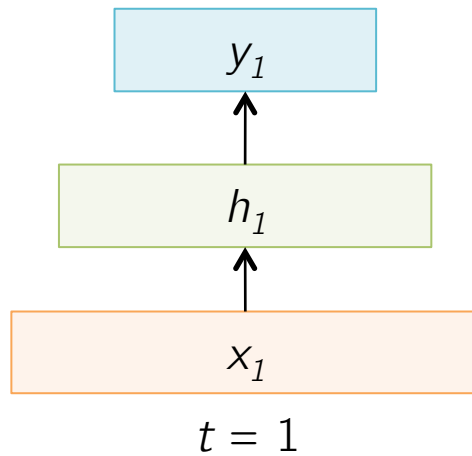# Motivation

- Not all problems can be converted into one with fixed-length inputs and outputs

- Problems such as Speech Recognition or Time-series Prediction require a system to store and use context information
  - Simple case: Output YES if the number of 1s is even, else NO
    1000010101 – YES, 100011 – NO, ...

- Hard/Impossible to choose a fixed context window
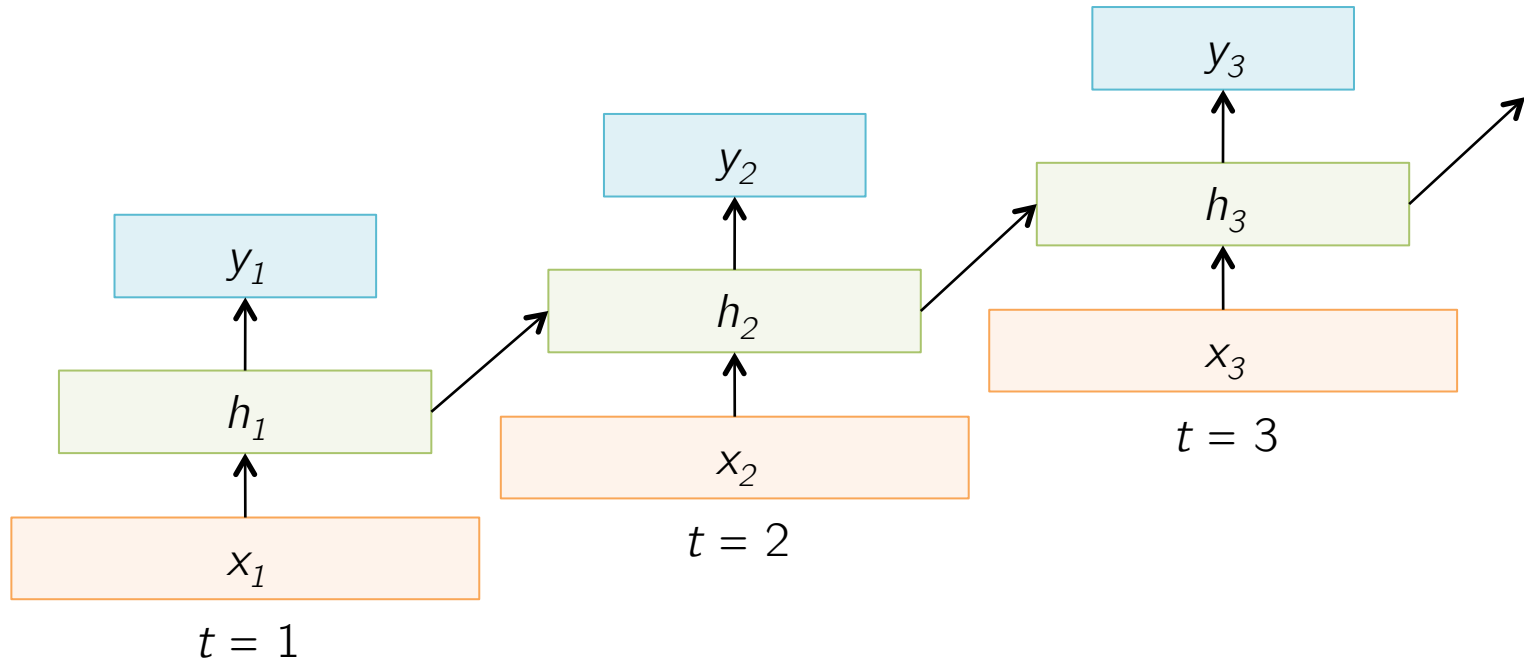  - There can always be a new sample longer than anything seen

# Recurrent Neural Networks (RNNs)

- Recurrent **N**eural **N**etwork**s** take the previous output or hidden states as inputs.
  The composite input at time $t$ has some historical information about the happenings at time $T < t$

- RNNs are useful as their intermediate values (state) can store information about past inputs for a time that is not fixed a priori
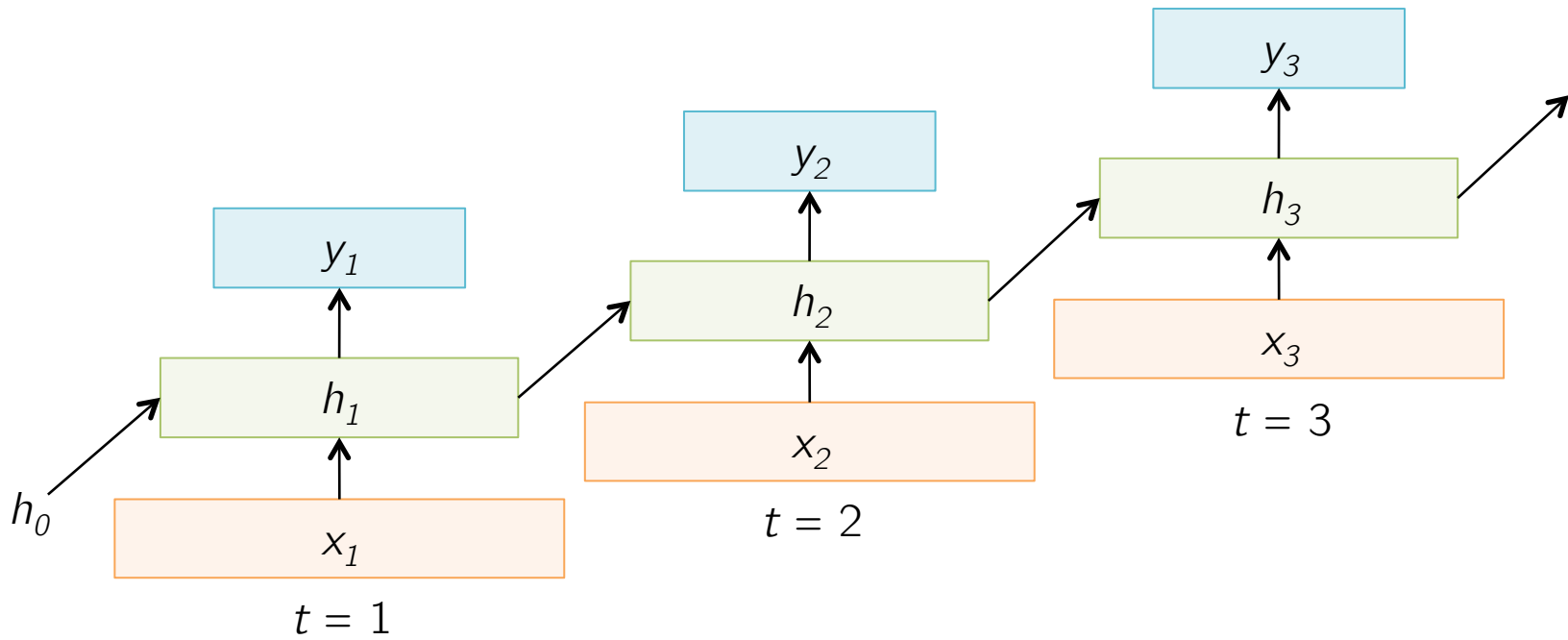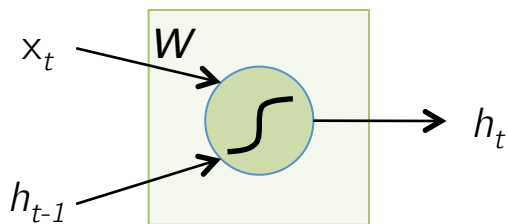
# Sample Feed-forward Network
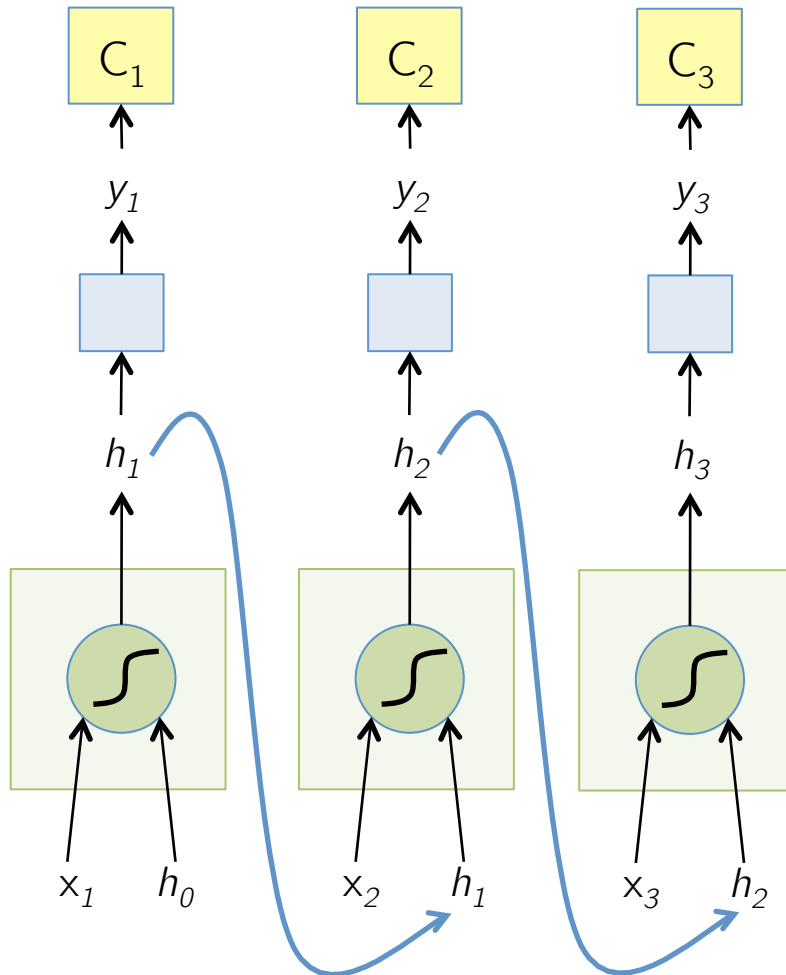
# Sample RNN

# Sample RNN

# The Vanilla RNN Cell



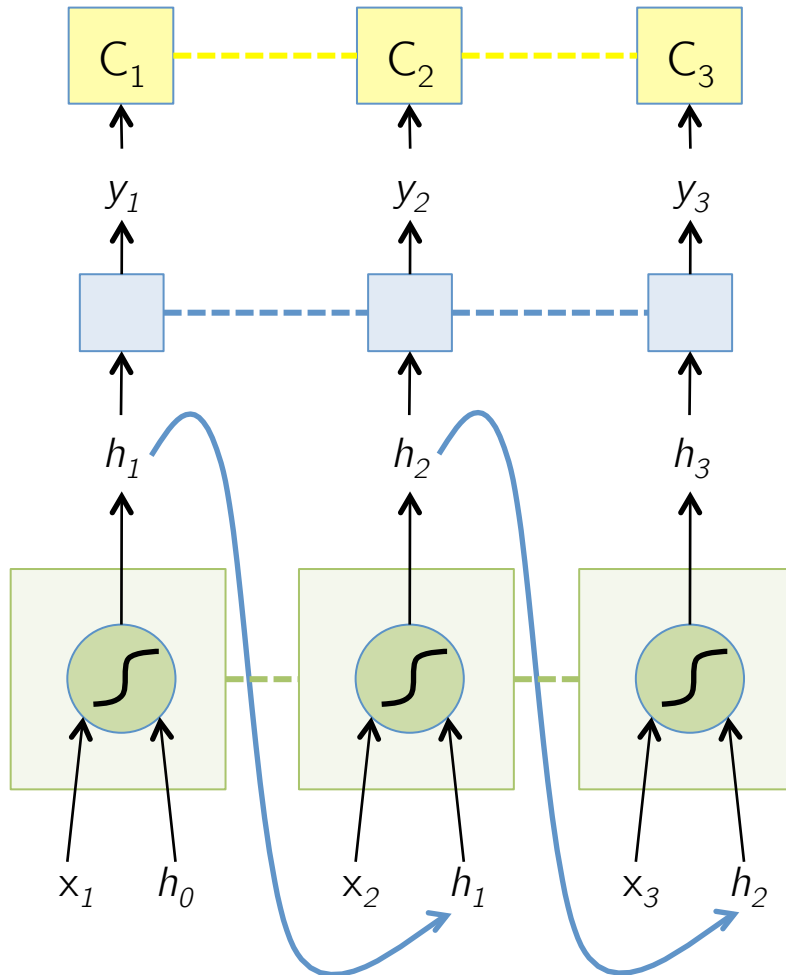$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

# The Vanilla RNN Forward



$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = \mathrm{F}(h_t)$$

$$C_t = \mathrm{Loss}(y_t, \mathrm{GT}_t)$$

# The Vanilla RNN Forward



$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = \mathrm{F}(h_t)$$

$$C_t = \mathrm{Loss}(y_t, \mathrm{GT}_t)$$

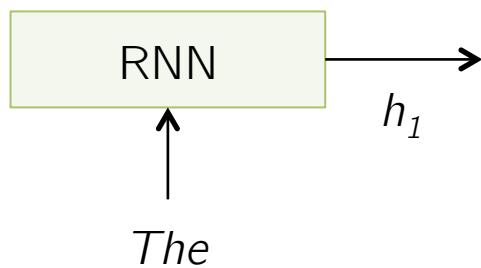------ indicates shared weights

10

# Recurrent Neural Networks (RNNs)

- Note that the weights are shared over time

- Essentially, copies of the RNN cell are made over time (unrolling/unfolding), with different inputs at different time steps
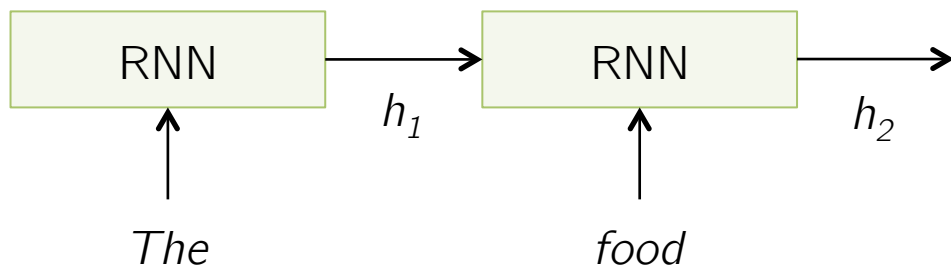
# Sentiment Classification

- Classify a
  restaurant review from Yelp! OR
  movie review from IMDB OR
  …
  as positive or negative

- **Inputs:** Multiple words, one or more sentences
- **Outputs:** Positive / Negative classification

- "The food was really good"
- "The chicken crossed the road because it was uncooked"
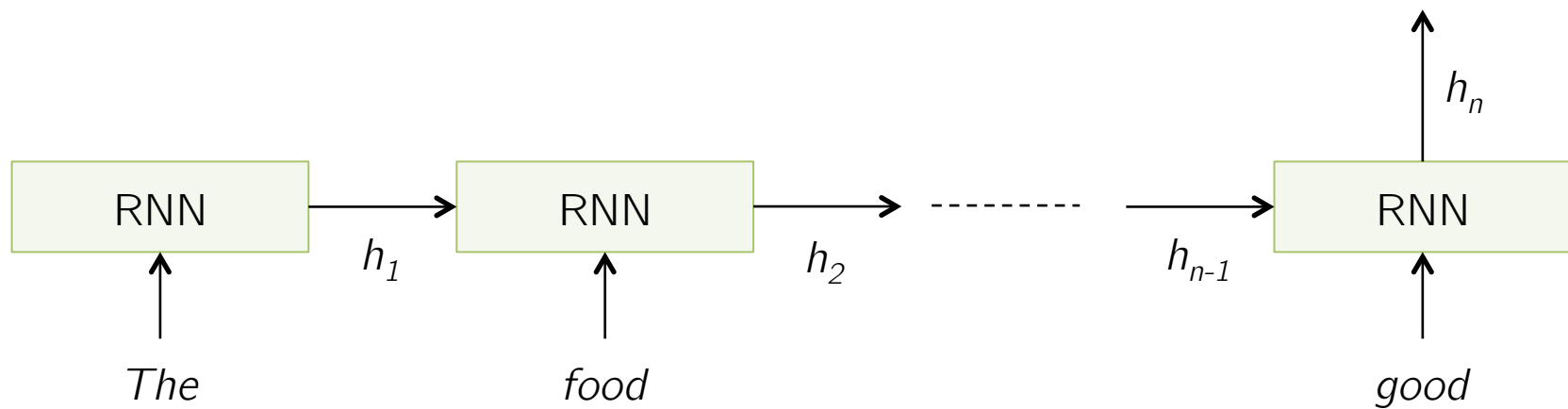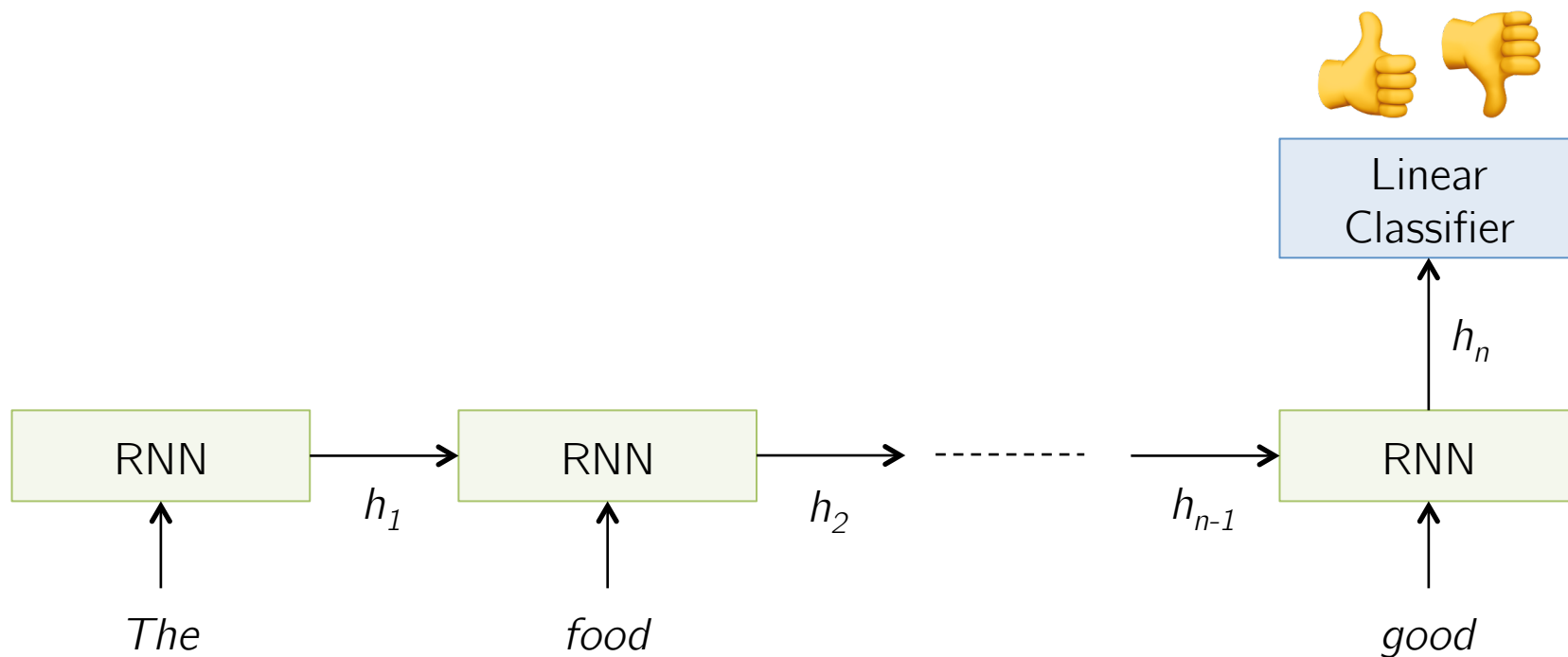
# Sentiment Classification

RNN $\longrightarrow$

$h_1$

$\uparrow$

*The*

# Sentiment Classification

# Sentiment Classification

# Sentiment Classification

�👍 👎

| Linear Classifier |

$h_n$

| RNN | | RNN | | RNN |

$h_1$     $h_2$     - - - - - - - -     $h_{n-1}$

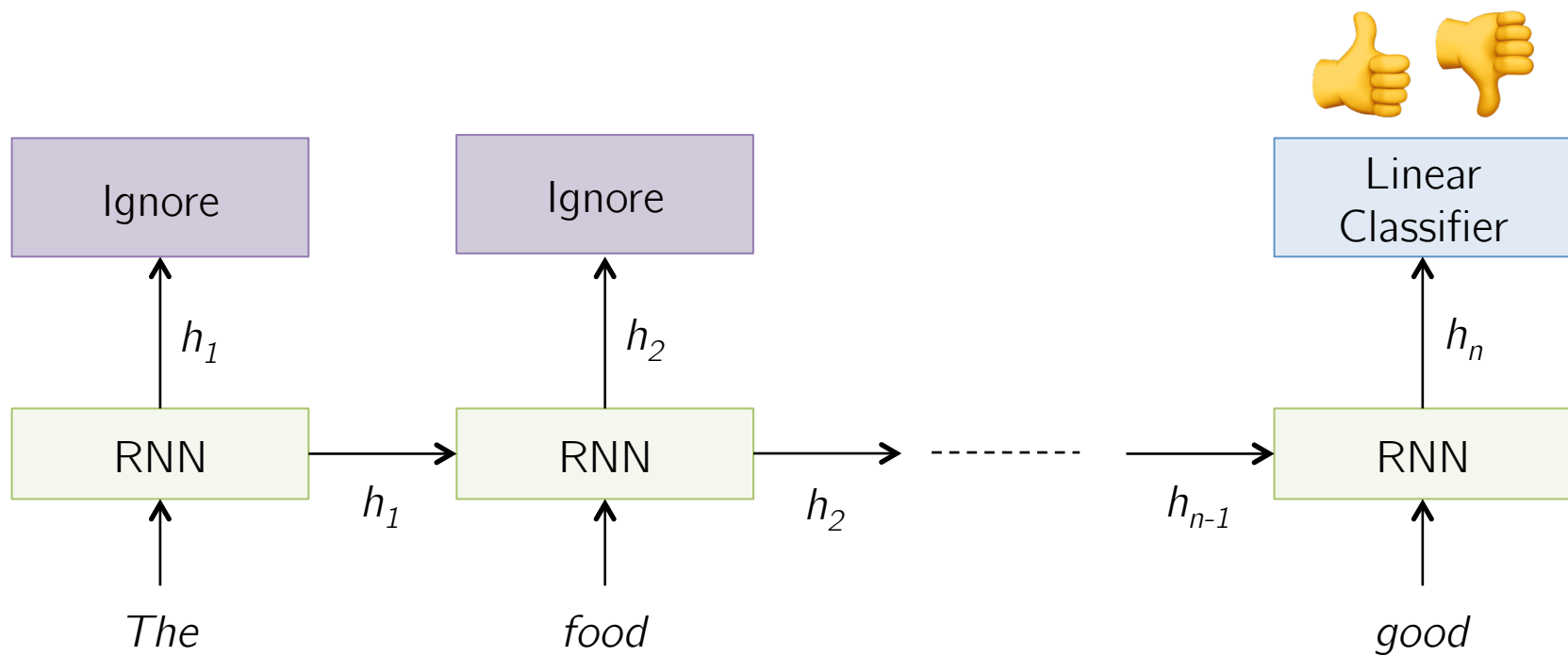*The*        *food*        *good*

# Sentiment Classification

# Sentiment Classification
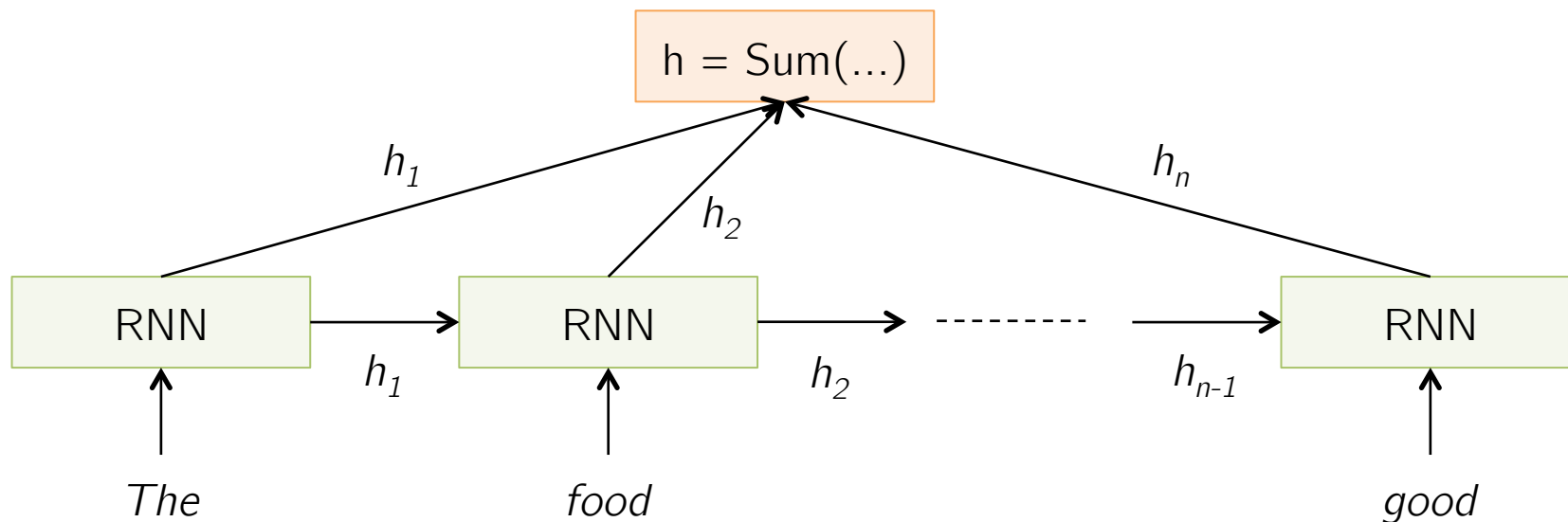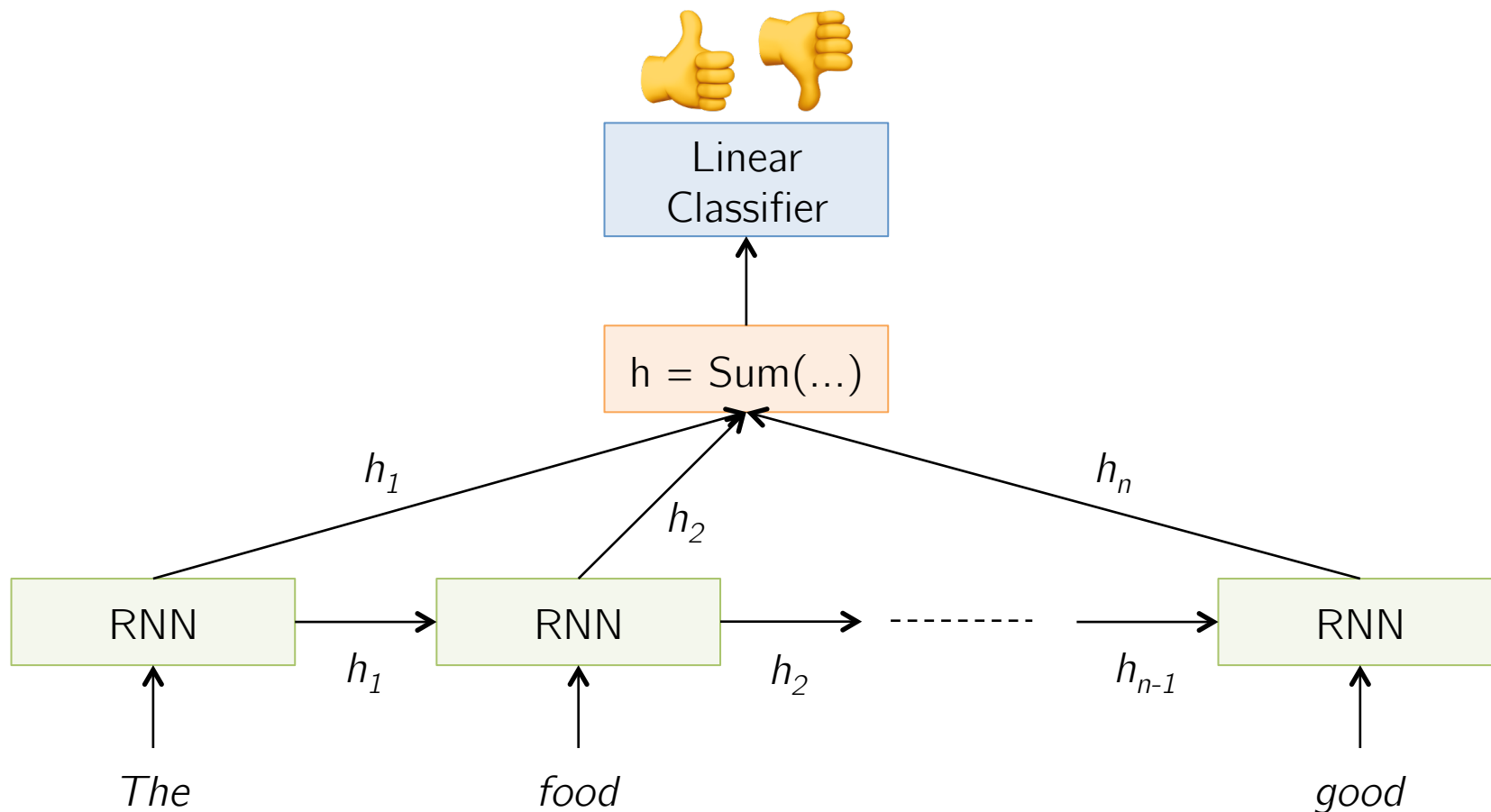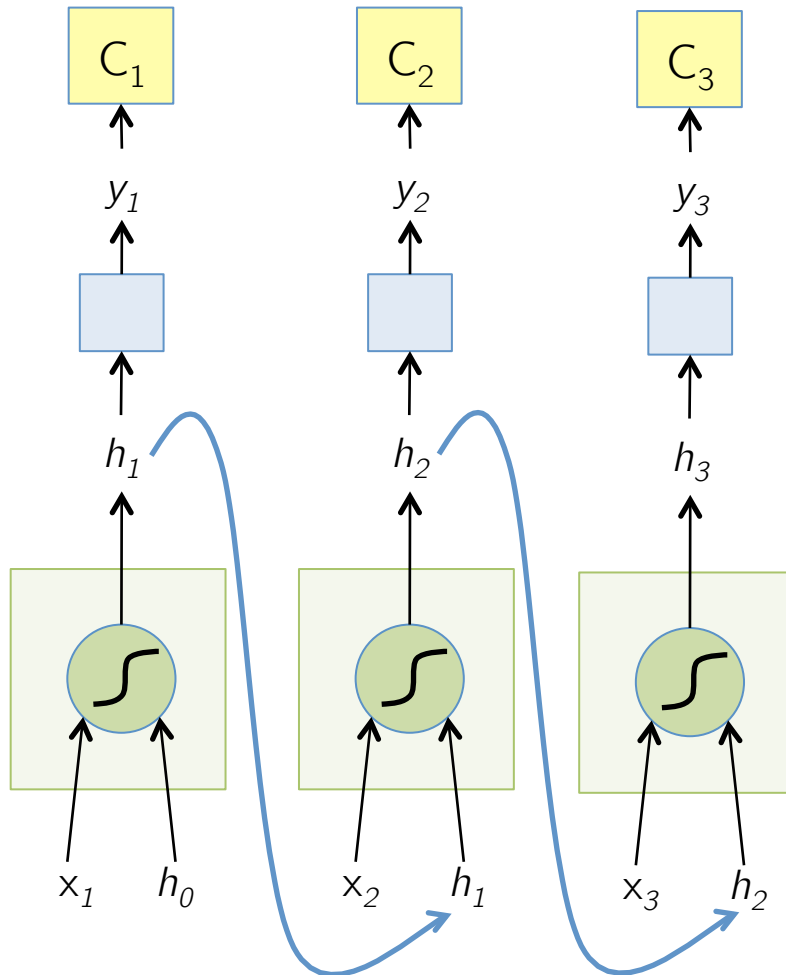
# Sentiment Classification
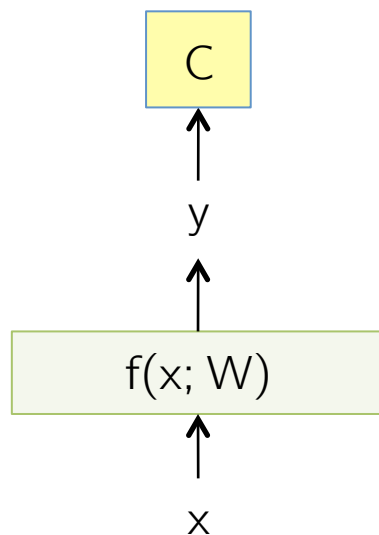
# The Vanilla RNN Forward



$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = F(h_t)$$

$$C_t = \text{Loss}(y_t, \text{GT}_t)$$

"Unfold" network through time by making copies at each time-step

# BackPropagation Refresher

C

y

f(x; W)

x

$$y = f(x; W)$$

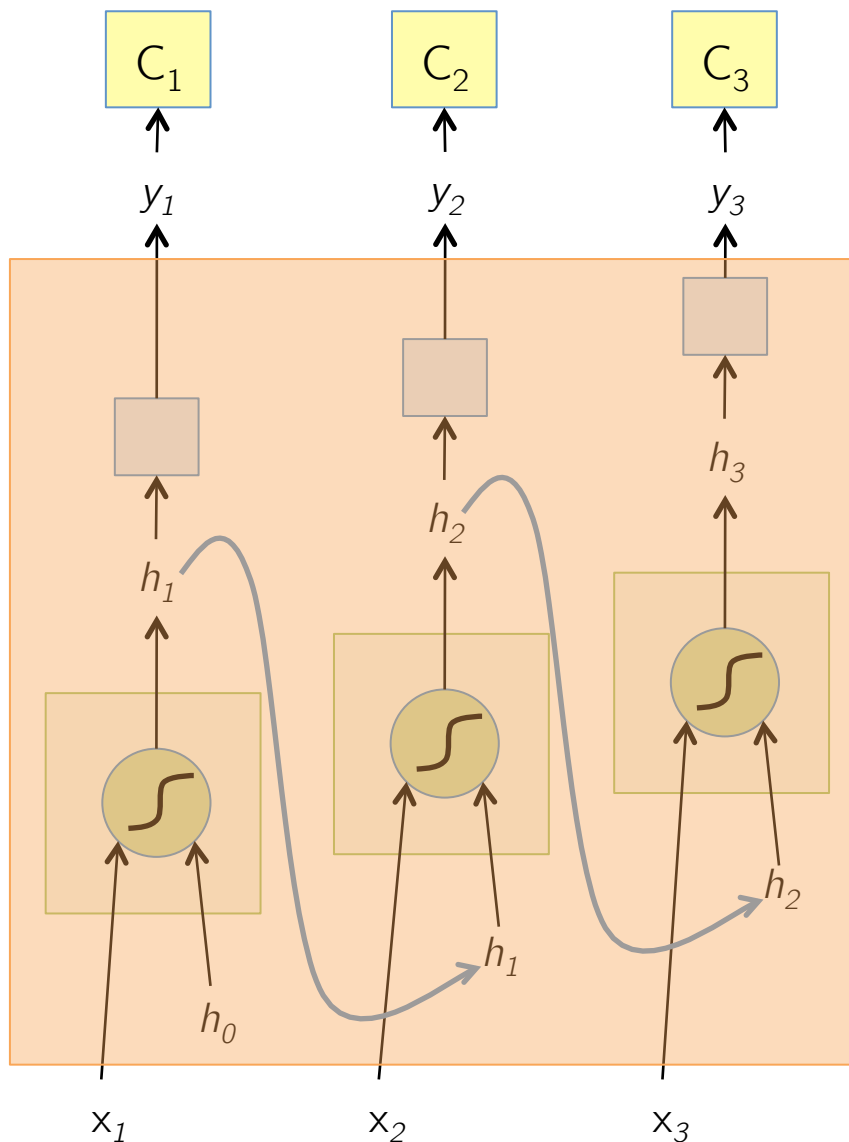$$C = \text{Loss}(y, y_{GT})$$

SGD Update

$$W \leftarrow W - \eta \frac{\partial C}{\partial W}$$

$$\frac{\partial C}{\partial W} = \left( \frac{\partial C}{\partial y} \right) \left( \frac{\partial y}{\partial W} \right)$$
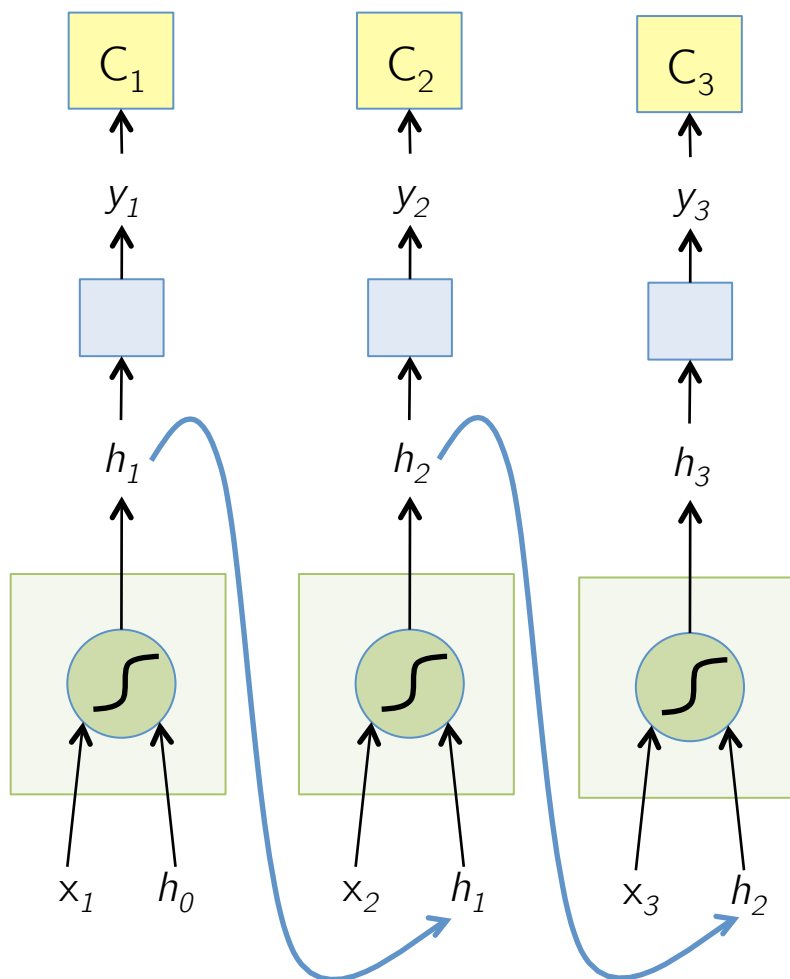
# BackPropagation Through Time (BPTT)

- One of the methods used to train RNNs
- The unfolded network (used during forward pass) is treated as one big feed-forward network
- This unfolded network accepts the whole time series as input

- The weight updates are computed for each copy in the unfolded network, then summed (or averaged) and then applied to the RNN weights
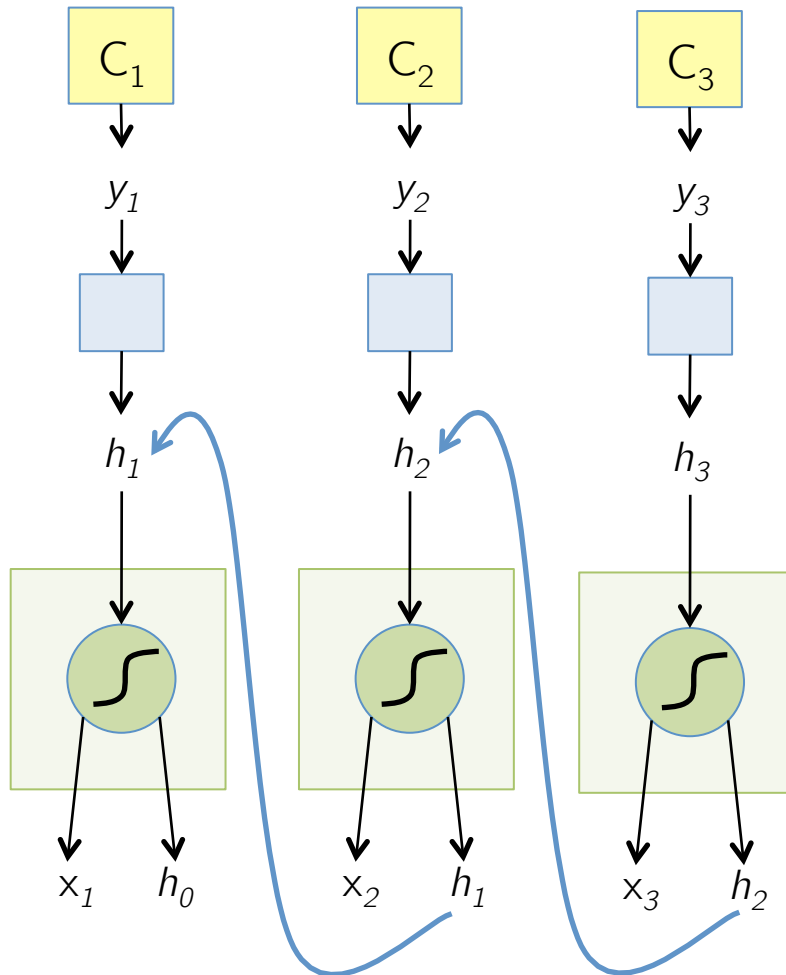
# The Unfolded Vanilla RNN



- Treat the unfolded network as one big feed-forward network!

- This big network takes in entire sequence as an input

- Compute gradients through the usual backpropagation
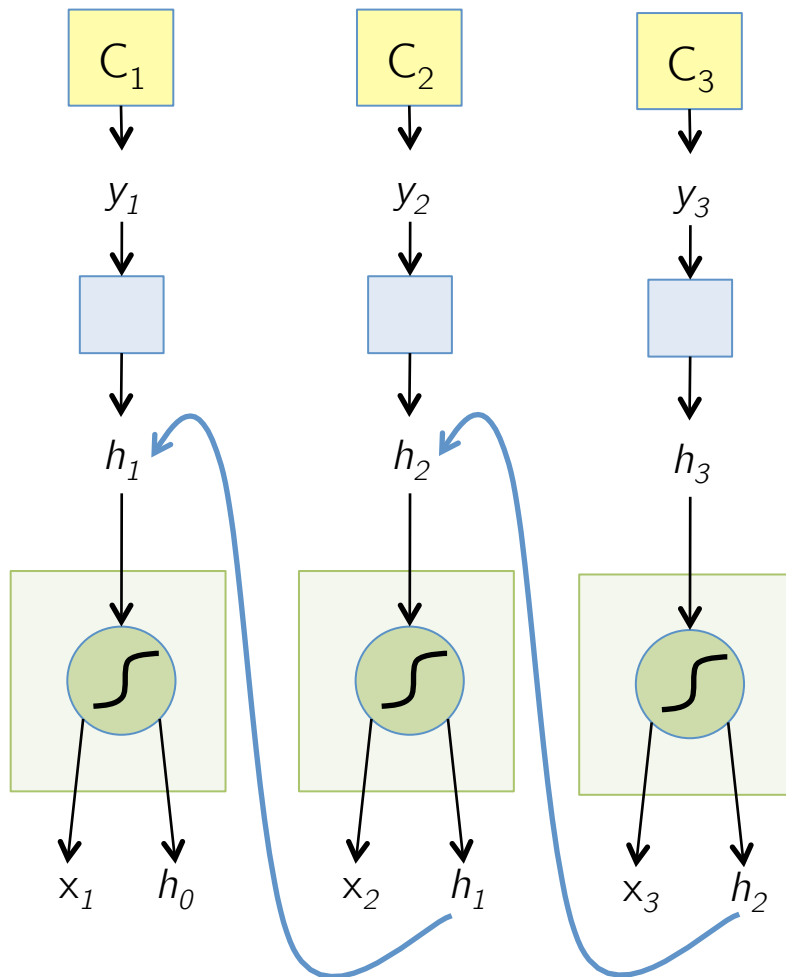
- Update shared weights

38

# The Unfolded Vanilla RNN Forward

# The Unfolded Vanilla RNN Backward

# The Vanilla RNN Backward



$$h_t = \tanh W \begin{pmatrix} x_t \\ h_{t-1} \end{pmatrix}$$

$$y_t = \mathrm{F}(h_t)$$

$$C_t = \mathrm{Loss}(y_t, \mathrm{GT}_t)$$

$$\frac{\partial C_t}{\partial h_1} = \left( \frac{\partial C_t}{\partial y_t} \right) \left( \frac{\partial y_t}{\partial h_1} \right)$$

$$= \left( \frac{\partial C_t}{\partial y_t} \right) \left( \frac{\partial y_t}{\partial h_t} \right) \left( \frac{\partial h_t}{\partial h_{t-1}} \right) \cdots \left( \frac{\partial h_2}{\partial h_1} \right)$$