# Training recurrent networks

# Recurrent Networks

- Processes a sequence

- Feedback connections

$$\mathbf{h}_{t-1} \qquad \mathbf{x}_t$$

$$f_{\mathbf{h}}$$

# How do we train RNNs?

- No longer a simple forward and backward pass

  - Cycles

$$\mathbf{h}_{t-1} \quad \mathbf{x}_t$$
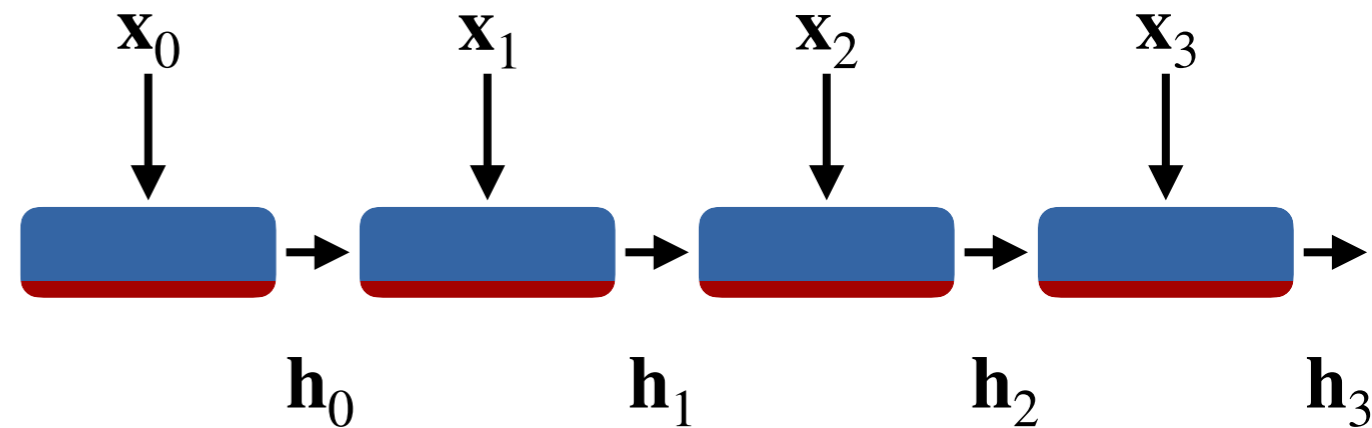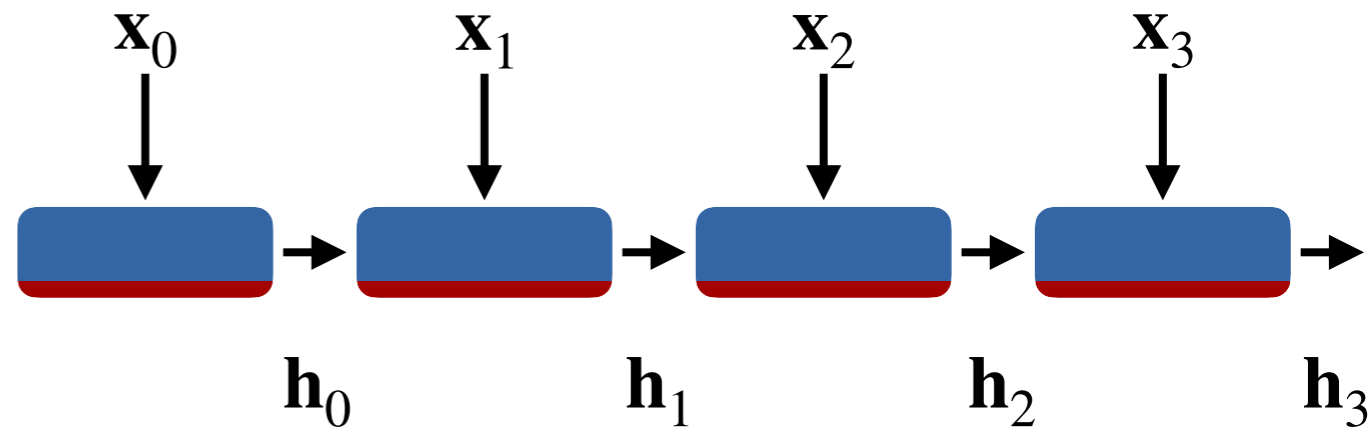
$$f_{\mathbf{h}}$$

# Solution: unrolling through time

# Unrolling through time

- Unrolled RNN (T steps)

  - Feed forward network

  - Shared parameters

- Trained with backprop

$$\mathbf{x}_0 \quad\quad \mathbf{x}_1 \quad\quad \mathbf{x}_2 \quad\quad \mathbf{x}_3$$

$$\mathbf{h}_0 \quad\quad \mathbf{h}_1 \quad\quad \mathbf{h}_2 \quad\quad \mathbf{h}_3$$
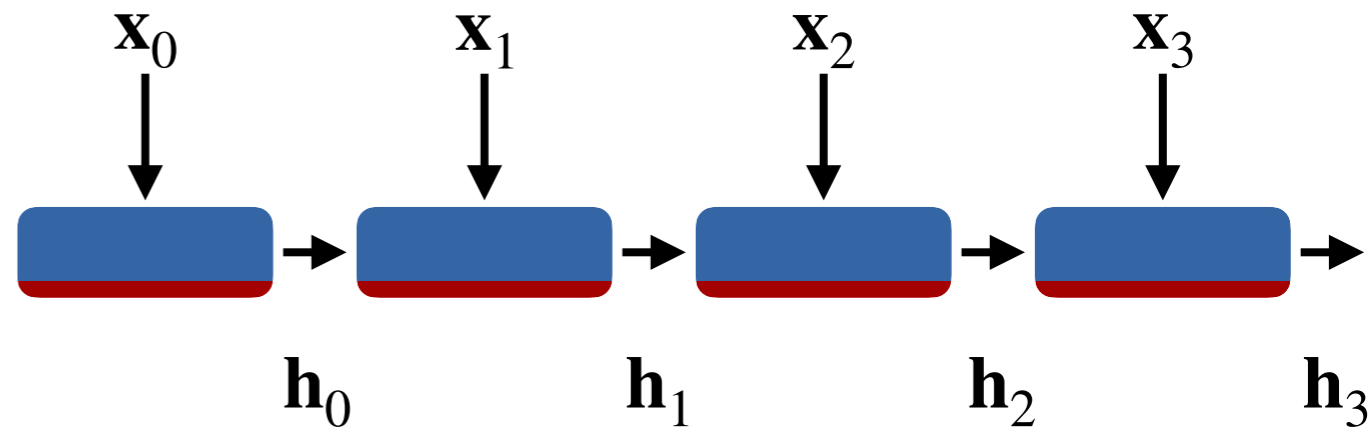
# Unrolling through time - Issues

- Long unrolling

  - Computationally expensive

  - Vanishing or exploding gradients

# Computation

- Solution (hack)

  - Cut RNN after n steps

    - Set h=0

- Works well in practice

$$\mathbf{x}_0 \qquad \mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3$$

$$\mathbf{h}_0 \qquad \mathbf{h}_1 \qquad \mathbf{h}_2 \qquad \mathbf{h}_3$$

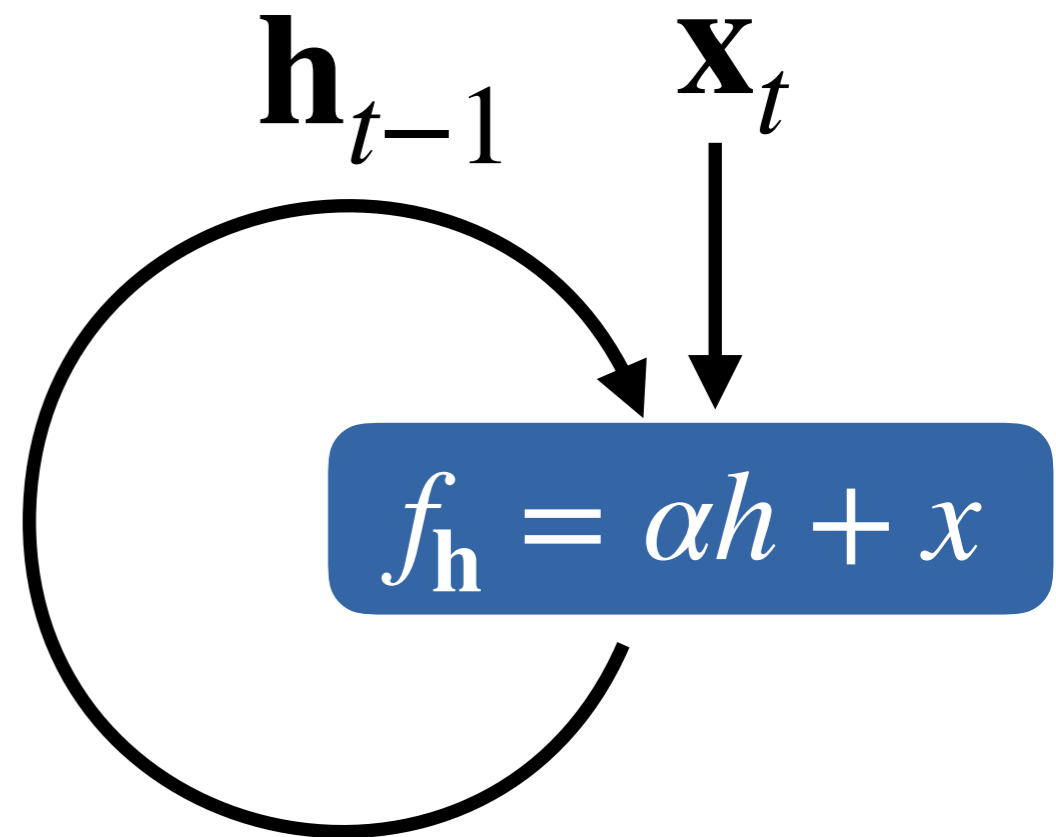# Vanishing and exploding gradients - Simple example

- Linear RNN

  - $\mathbf{h}_t = \alpha \mathbf{h}_{t-1} + \mathbf{x}_t$

- For large $t$

  - $a > 1$ $\quad \dfrac{\partial}{\partial \mathbf{x}_0}\mathbf{h}_t = \alpha^t \rightarrow \infty$
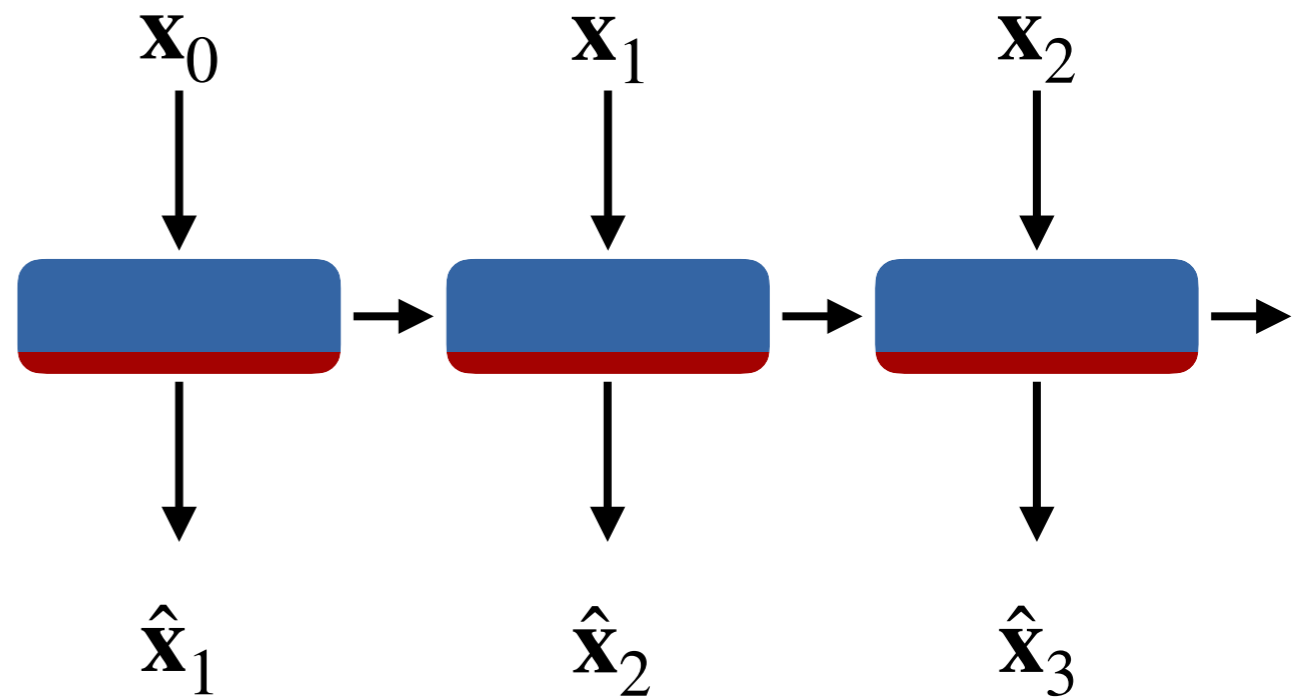
  - $a < 1$ $\quad \dfrac{\partial}{\partial \mathbf{x}_0}\mathbf{h}_t = \alpha^t \rightarrow 0$

$$\mathbf{h}_{t-1} \qquad \mathbf{x}_t$$

$$f_{\mathbf{h}} = \alpha h + x$$

# Preventing vanishing and exploding gradients

- Generative models

- Use ground truth inputs

# Preventing vanishing and exploding gradients

- Exploding gradients

  - Gradient clipping
  $$\nabla \ell' = \nabla \ell \min \left( 1, \frac{\epsilon}{|\nabla \ell|} \right)$$

- Vanishing gradients

  - Different RNN structure