# Recurrent Neural Networks

## Sequence models

Examples:

- Stock prices: $x_1, \ldots, x_n$
  Predict $\hat{x}_t \sim Pr(x_t | x_1, \ldots, x_{t-1})$
- Language models:
  Given phrase, what is its prob?
  $Pr[\text{The cat is black}] = Pr[\text{The}] \cdot Pr[\text{cat} | \text{The}] \cdot Pr[\text{is} | \text{The}, \text{cat}]$
  $\cdot Pr[\text{black} | \text{The}, \text{cat}, \text{is}]$

Next word predict:
$\hat{x}_t \sim Pr[x_t | x_1, \ldots, x_{t-1}]$

next word      prev words

Idea: Summarize past observations.
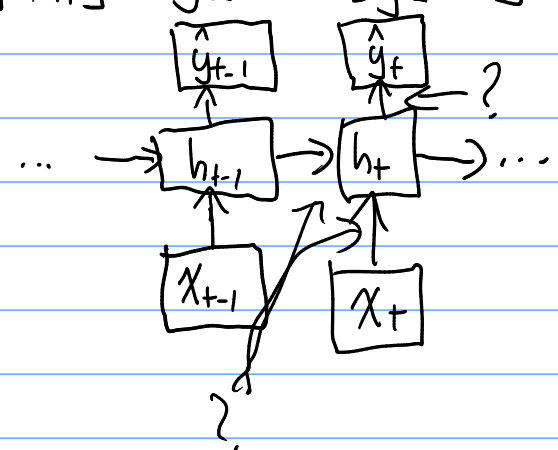Summary $h_t$ (vector). Then $\hat{x}_t \sim Pr[x_t | h_t]$

$h_t$: fn of $x_1, \ldots, x_t$. Update over time.
$h_1 = g(\vec{0}, x_1)$ . $h_2 = g(h_1, x_2)$ ...
$\hat{y}_1 \sim Pr[y_1 | h_1]$    $\hat{y}_2 \sim Pr[y_2 | h_2]$

(e.g. $y_1 = x_2$)
Next word prediction



How to compute hidden state?

$$h_t = f(h_{t-1}, x_t, \theta)$$

NN → (arrow to $h_t$)

prev hidden state (arrow to $h_{t-1}$)

current input (arrow to $x_t$)

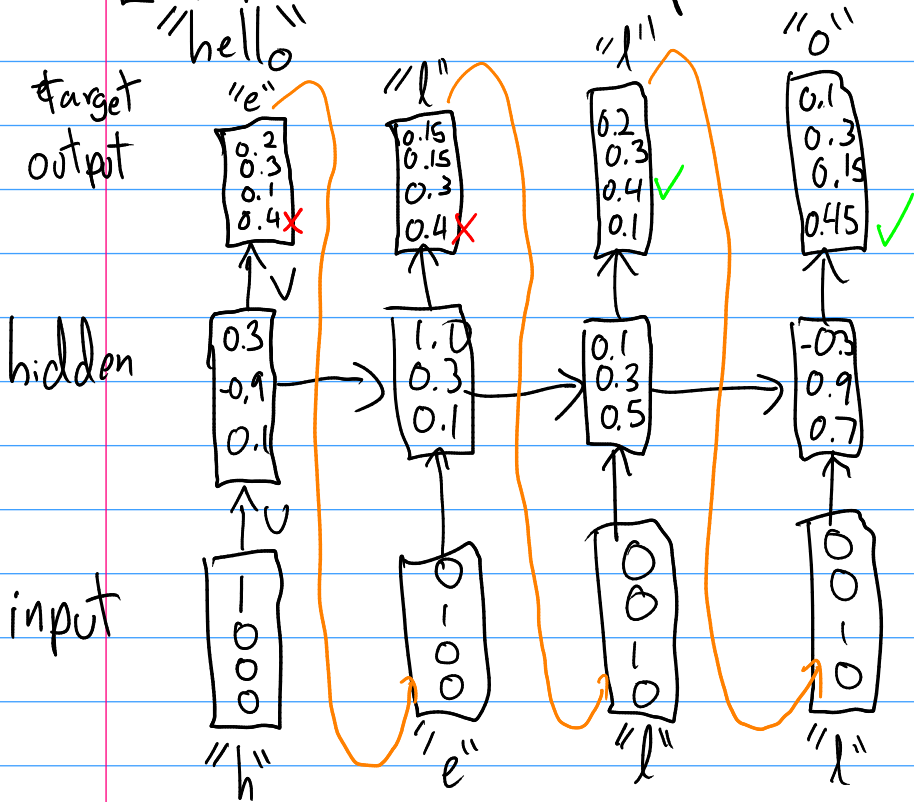parameter vector (arrow to $\theta$)



$$h_t = \tanh(Wh_{t-1} + Ux_t + b)$$

$$\hat{y}_t = \text{softmax}(Vh_t + c)$$

$$l(\{x_1, \ldots, x_n\}, \{y_1, \ldots, y_n\}) = \sum_{t=1}^{n} l(y_t, \hat{y}_t)$$

# Ex: Next char prediction

"hello"

target
output

| | | | |
|---|---|---|---|
| "e" | "l" | "l" | "o" |
| 0.2 | 0.15 | 0.2 | 0.1 |
| 0.3 | 0.15 | 0.3 | 0.3 |
| 0.1 | 0.3 | 0.4 ✓ | 0.15 |
| 0.4 ✗ | 0.4 ✗ | 0.1 | 0.45 ✓ |

V

hidden

| | | | |
|---|---|---|---|
| 0.3 | 1.0 | 0.1 | -0.3 |
| -0.9 | 0.3 | 0.3 | 0.9 |
| 0.1 | 0.1 | 0.5 | 0.7 |

U

input

| | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 |

"h"    "e"    "l"    "l"

# Ex. of oth seq. problems:

## Many to one



Sentiment classification
"this movie is great" → +1
"I hated this movie" → -1

## One to many
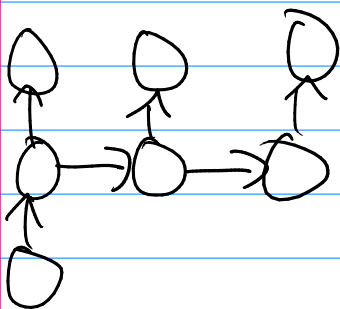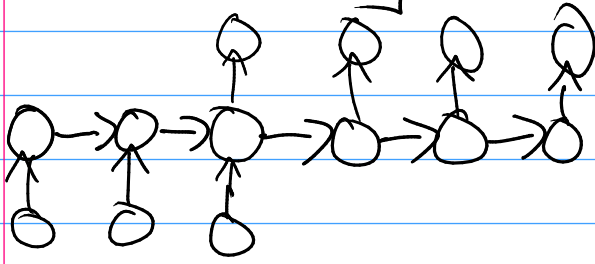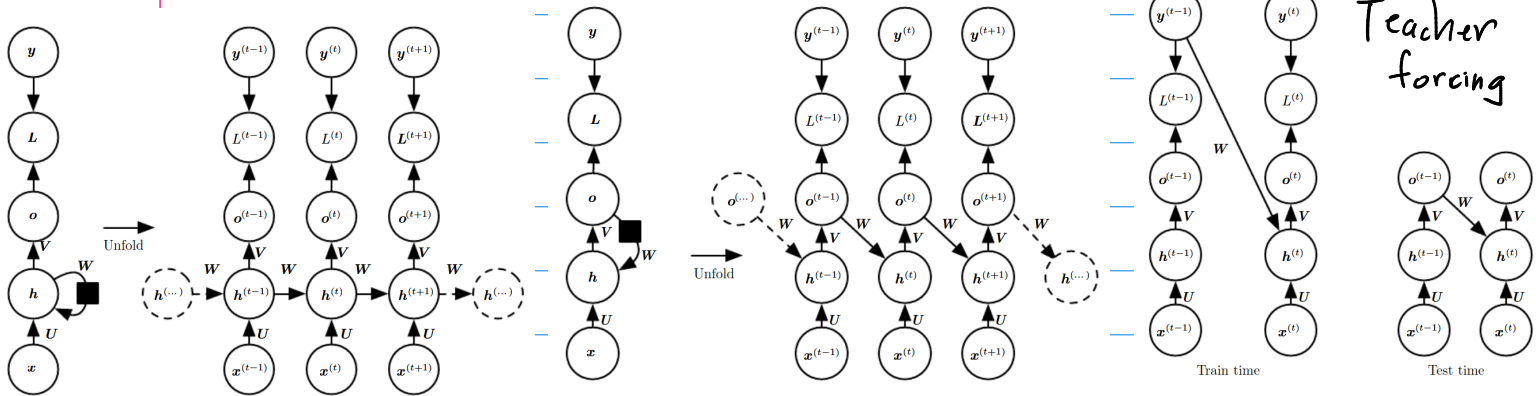


Image caption



A very poorly drawn cat

# Many to many

Machine translation

The cat is black
$\Rightarrow$ Le chat est noir

对不起 $\rightarrow$ Sorry

## Sequential vs Parallel



More powerful                    Parallelizable
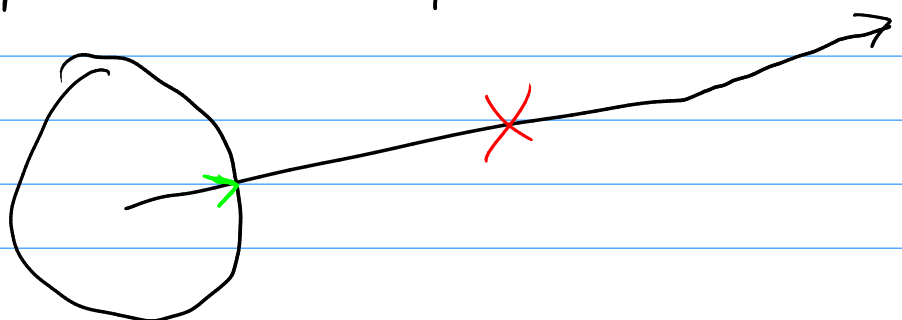
Teacher forcing

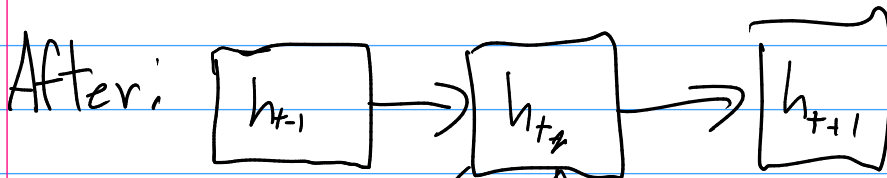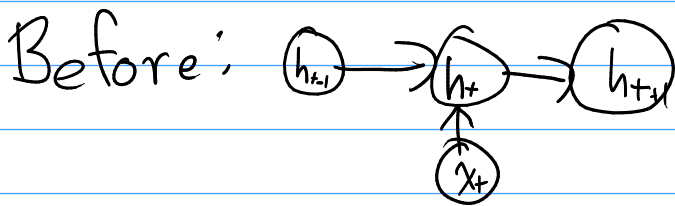# Optimization for RNNs
## "Backpropagation through time"

Long chains cause issues
Vanishing or exploding gradients.

- Truncate gradient chains
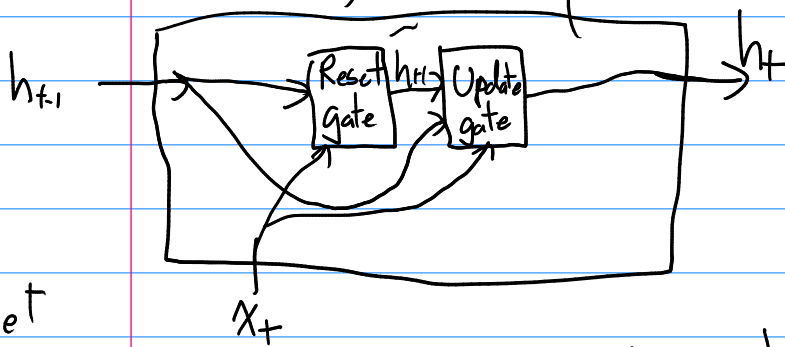  $\hookrightarrow$ stop grad comps after $\tau$ steps

- Gradien clipping
  If $\|g\|_2 > v \Leftarrow$ threshold
  then $g \Leftarrow \frac{g \, v}{\|g\|}$

# Gated Recurrent Units (GRUs)

Before: $h_{t-1} \rightarrow h_t \rightarrow h_{t+1}$

$x_t$

After: $h_{t-1} \rightarrow h_t \rightarrow h_{t+1}$

$x_t$

**Simplification:**
Reset, Update are bools
$= 0$ or $1$

$h_{t-1} \rightarrow$ [Reset gate] $\tilde{h}_{t+1} \rightarrow$ [Update gate] $\rightarrow h_t$

$x_t$

If "update" $= 1$:
$h_t = h_{t-1}$ (use $\underset{\text{state}}{\text{old}}$)

Else:
   If "reset" $= 1$:
   $h_t = \tanh(U x_t + W h_{t-1} + b)$
   (Standard RNN update)

Reset
$\downarrow$
$R_t = \text{sigmoid}(U^{(r)} x_t + W^{(r)} h_{t-1} + b^{(r)})$

$Z_t = \text{sigmoid}(U^{(z)} x_t + W^{(z)} h_{t-1} + b^{(z)})$

Vectors $\in [0,1]$
coordinate wise
same dim as $h_t$

   Else:
   $h_t = \tanh(U x_t + b)$
   (Drop old state)
      MLP

Update
$\tilde{h}_t = \tanh(U x_t + W(R_t \odot h_{t-1}) + b)$
$\tilde{\phantom{h}}$ candidate

coordinate wise product

$h_t = Z_t \odot h_{t-1} + (1 - Z_t) \odot \tilde{h}_t$
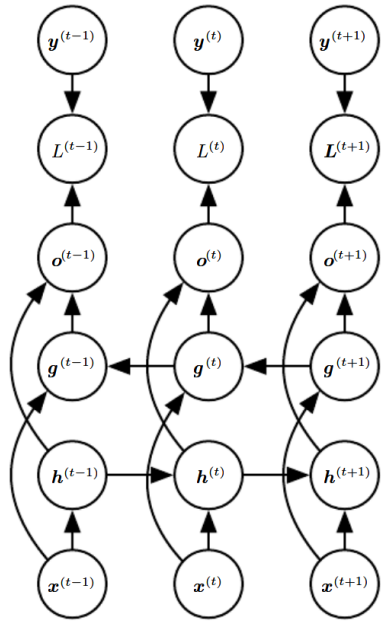
# Long Short Term Memory (LSTM)

# Bidirectional RNNs

Eg. : I went to the bank
   a) of the river
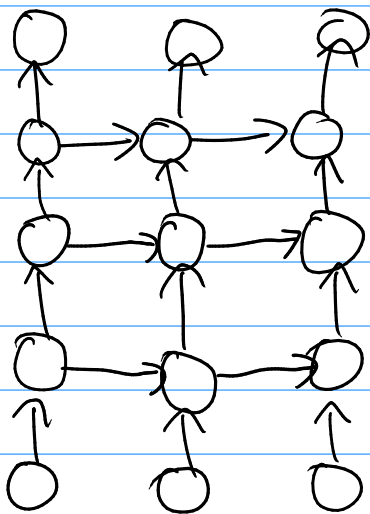   b) to withdraw cash

# Deep RNNs

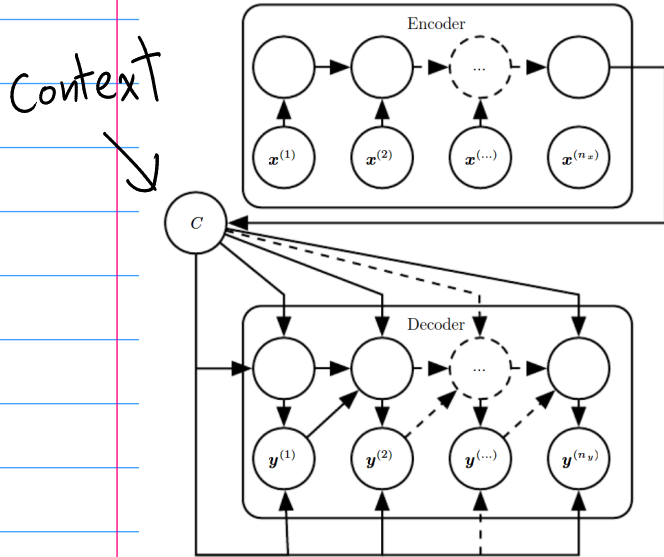$\hat{y}$

$f$

$g$

$h$

$x$

# Encoder-Decoder

Context



Machine translation