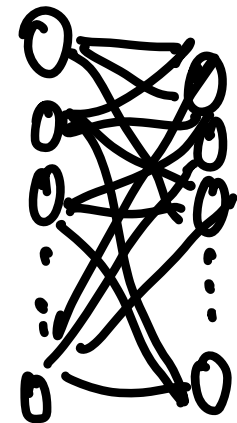# Plan

Recap
Reminders
Embedding
Motivation
Self-attention

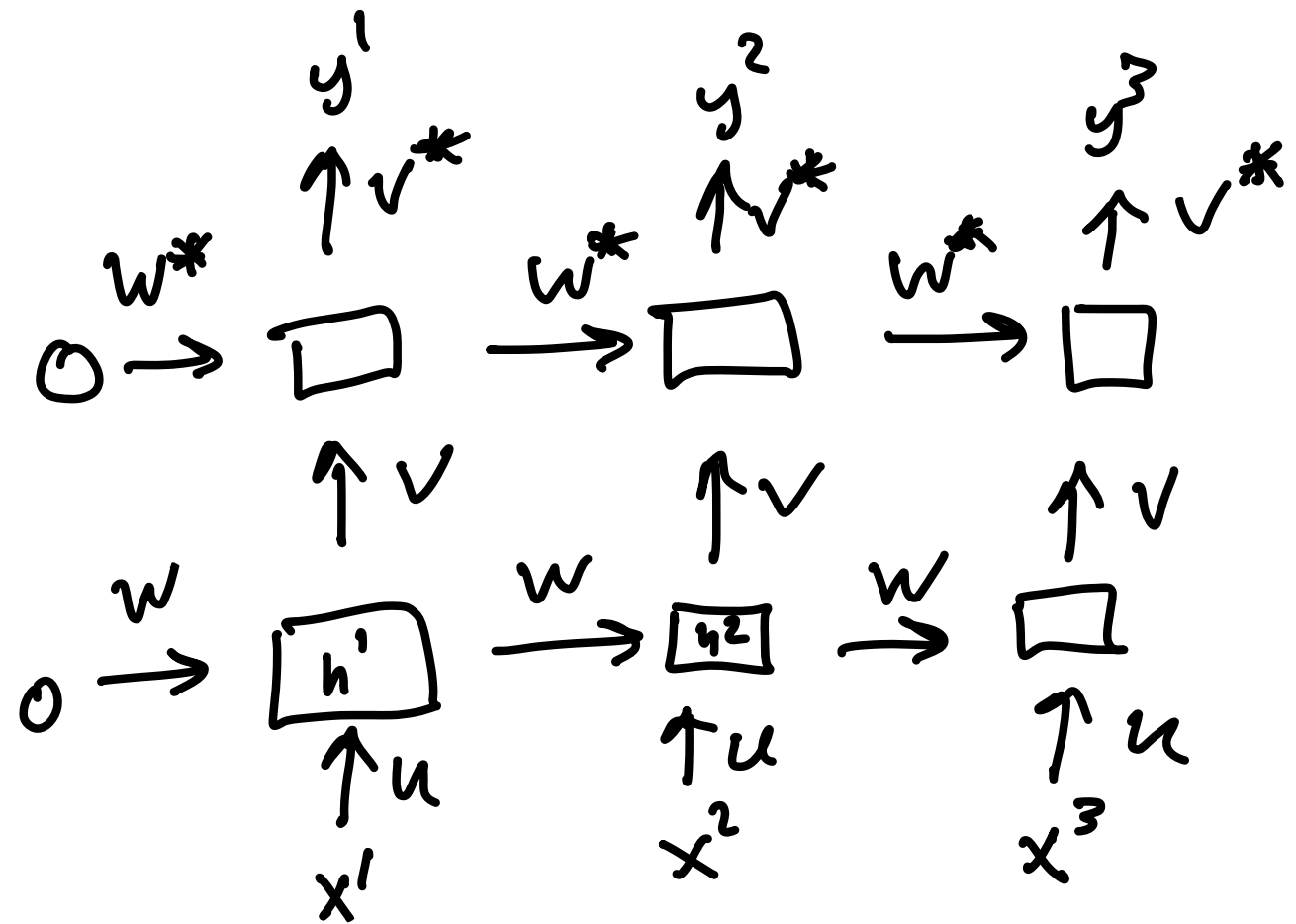$x$        $Wx$



# Recap

1) How to encode text?

2) How to include long range dependency?
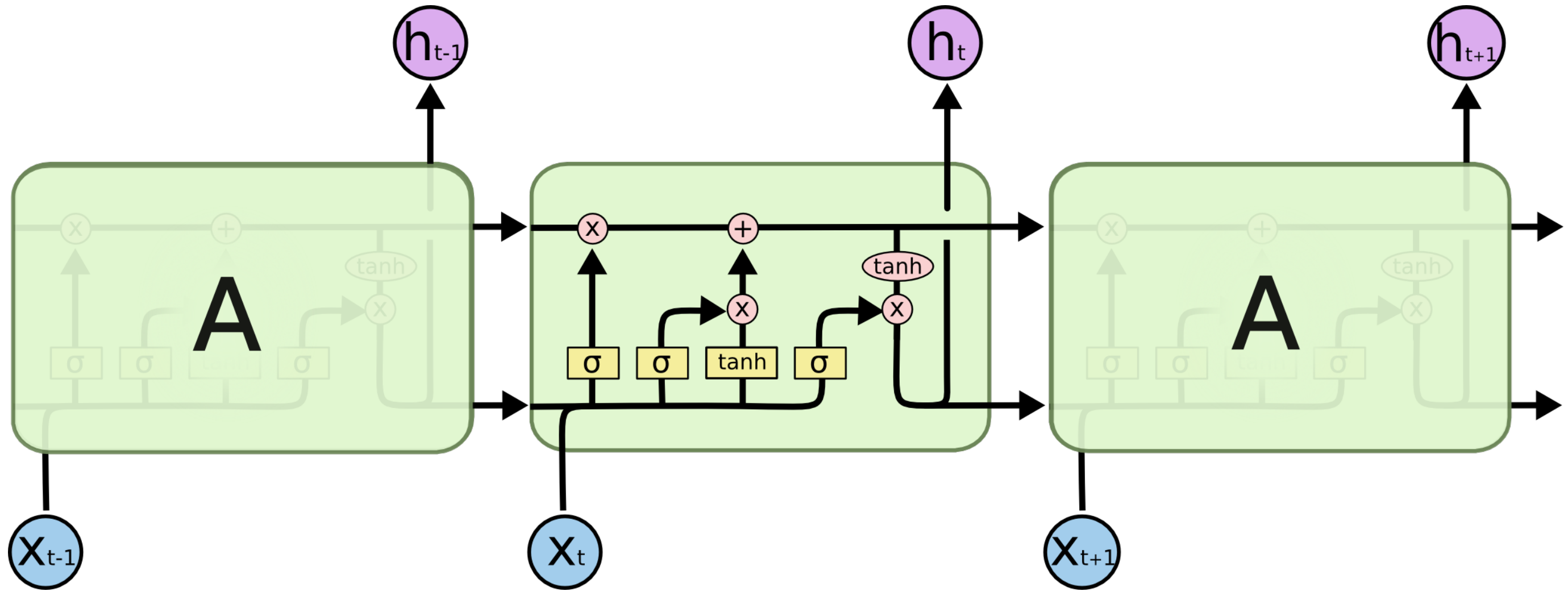


$$h^2 = \sigma(Wh^1 + Ux^2)$$

RNN
- variable length
- memory (but overwritten)

LSTM

cell state

$$\begin{bmatrix} \\ \\ \\ \end{bmatrix} * \begin{bmatrix} 1 \\ 0 \\ .5 \end{bmatrix}$$

# Reminders

- Form 24/26

- 24 hours <u>total</u> lateness

- work in 202 during office hours

- Project proposal due Monday

# Embedding

input    encoder    embedding    decoder    output

1

100
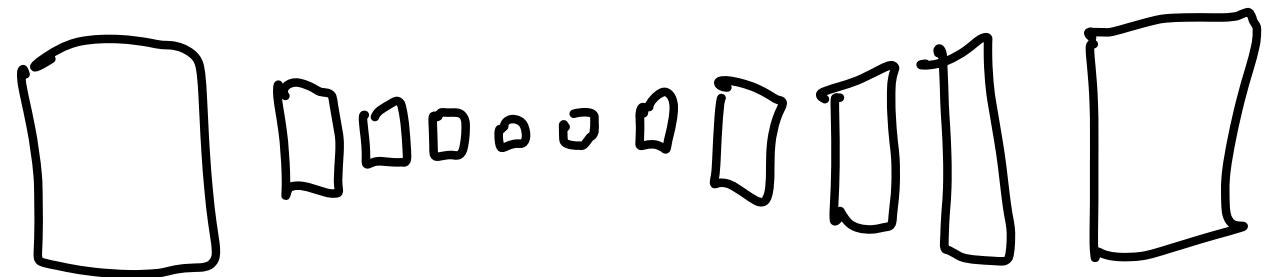
10,000

loss =. cross entropy

- $\ell_2$ norm = $\|input - output\|_2^2$

RNN work well on

- next word prediction

- sequence classification

but less well on

- translation

- sentence generation

---

The dog ran really fast

RNN :⋅

El perro corrio muy rapido

I love you a lot
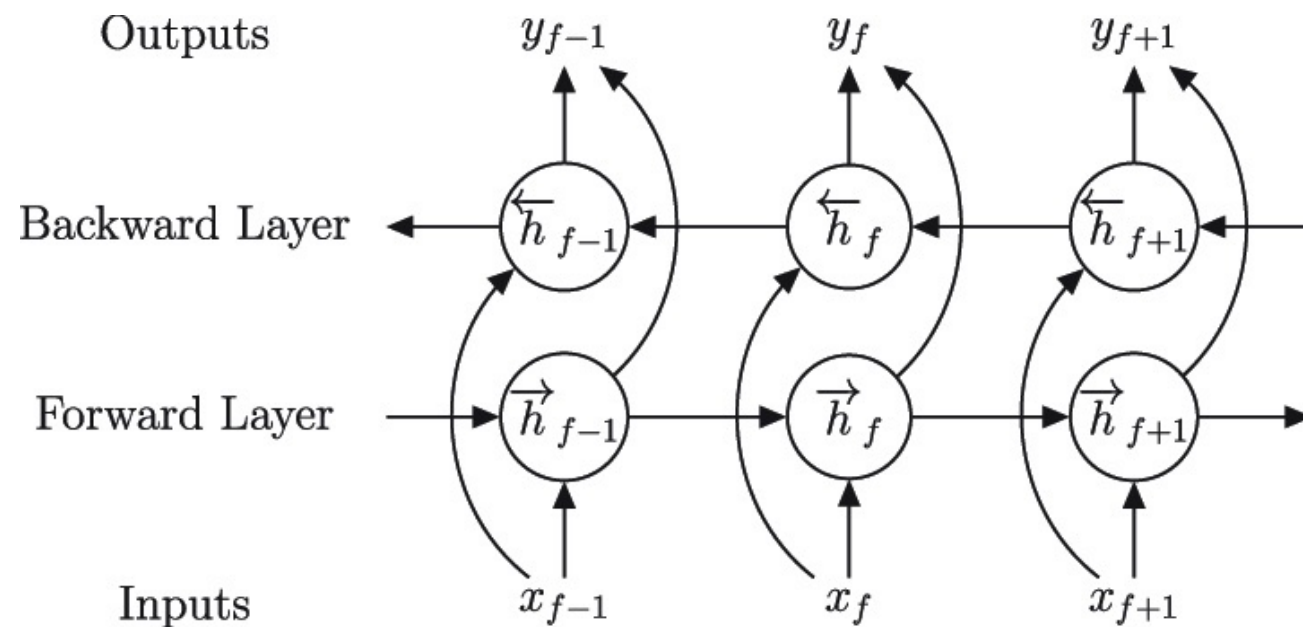
RNN :⋅

Te ~~amo~~ mucho

1. different # words 5 vs 3

2. mixed order

3. need whole context

**Approach 1:** sentences instead of words

\# sentences of length 5

$$5^{\#\ words}$$

**Approach 2:** bidirectional RNN



| | | | | |
|---|---|---|---|---|
| Outputs | $y_{f-1}$ | $y_f$ | $y_{f+1}$ | |
| Backward Layer | $\overleftarrow{h}_{f-1}$ | $\overleftarrow{h}_f$ | $\overleftarrow{h}_{f+1}$ | |
| Forward Layer | $\overrightarrow{h}_{f-1}$ | $\overrightarrow{h}_f$ | $\overrightarrow{h}_{f+1}$ | |
| Inputs | $x_{f-1}$ | $x_f$ | $x_{f+1}$ | |

no long range dependency

Wishlist:

- long range dependency
- variable length
- association between output and input

input: $x_1, \ldots, x_n \in \mathbb{R}^d$

output: $y_1, \ldots, y_n \in \mathbb{R}^d$

$$y_i = \sum_{j=1}^{n} w_{ij} \underbrace{x_j}_{\text{value}}$$

row normalized

$$\sum_{j=1}^{n} w_{ij} = 1$$

$W$

$$\begin{bmatrix} \phantom{xxxxxxx} \end{bmatrix}$$

$n \times n$

$$w_{ij} = \frac{\exp(\omega_{ij})}{\sum_{j'} \exp(\omega_{ij'})} \quad \xleftarrow{\text{softmax}}$$

$w_{ij}$ = similarity between $y_i, x_j$

$$= \underbrace{x_i^T}_{\text{query}} \underbrace{x_j}_{\text{key}}$$

I love you a lot

Te amo mucho

$$amo = \underset{1/4}{I} + \underset{3/4}{love} + \underset{0}{you} + \underset{0}{a} + \underset{0}{lot}$$

- set to set ☺
- see all input
  ↳ no distance
- parameters

$x_j$ appeared as

- part of output
- weight for own output
- weight for other output

$q_i = W_q x_i$ # query

$k_i = W_k x_i$ # key

$v_i = W_r x_i$ # value

$W_{ij} = q_i^T k_j$ $\quad w_{ij} = softmax(w_{ij})$

$y_i = \sum_{j=1}^{n} W_{ij} v_j$



| | I | love | you | a | lot |
|---|---|---|---|---|---|
| I | 3/4 | 1/4 | | | |
| love | 1/2 | 1/2 | | | |
| you | | | 1 | | |
| a | | | | 3/4 | 1/4 |
| lot | | | | 0 | 1 |

$W_{ij}$ where $i \to love$
$\quad j \to you$

$= q_i^T k_j$

$W_{ii}$
$= q_i^T k_i$

# Transformer

| input | transformer block | output |

self attention  $\oplus$  layer norm  MLP MLP MLP MLP  $\oplus$  layer norm

## Positional Encoding

①    ②    ③ ← t  ④  ⑤

Jack  gave  water  to  Jill

Jill  gave  water  to  Jack

---

$x_i$  ← embedding "Jack"

$\begin{bmatrix} \\ \\ \end{bmatrix}$

← encodes position
  ↳ put in index
  ↳ one hot encoding

Jack
$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

gave
$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

water
$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

← omega
$\begin{bmatrix} \sin(\omega_1 t) \\ \sin(\omega_2 t) \\ \sin(\omega_3 t) \\ \vdots \end{bmatrix}$

Jack gave water