

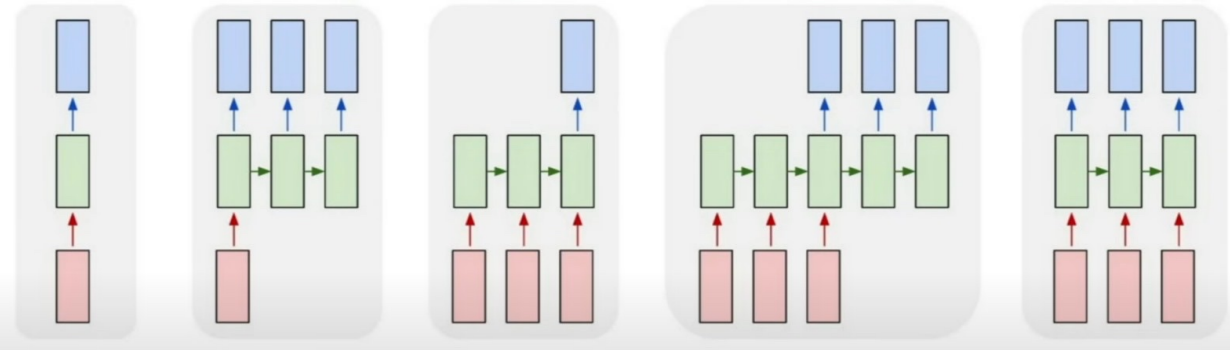
one to one

one to many

many to one

many to many

many to many



↳ image
↓
sequence of words

↳ sequence of words
↓
sentiment
↳ + / -

↳ seq of words
↓
Translated seq of words

↳ Video classification on frame level

Process sequence:

- ↳ Take a fixed input
- ↳ Network looks at glimpses of image
- ↳ After a series of glimpses, it makes a decision

RNN:

- ↳ Read an input
- ↳ update its internal hidden state
- ↳ Produce an output

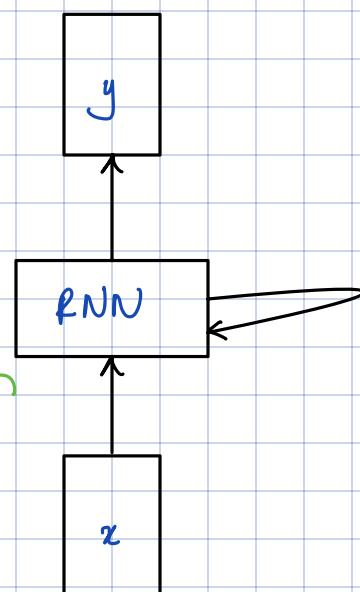
old state
↑

↓
new state

↓
some func with params W

↳ input vector at current time step

} functional form of recurrence relation



Vanilla RNN:

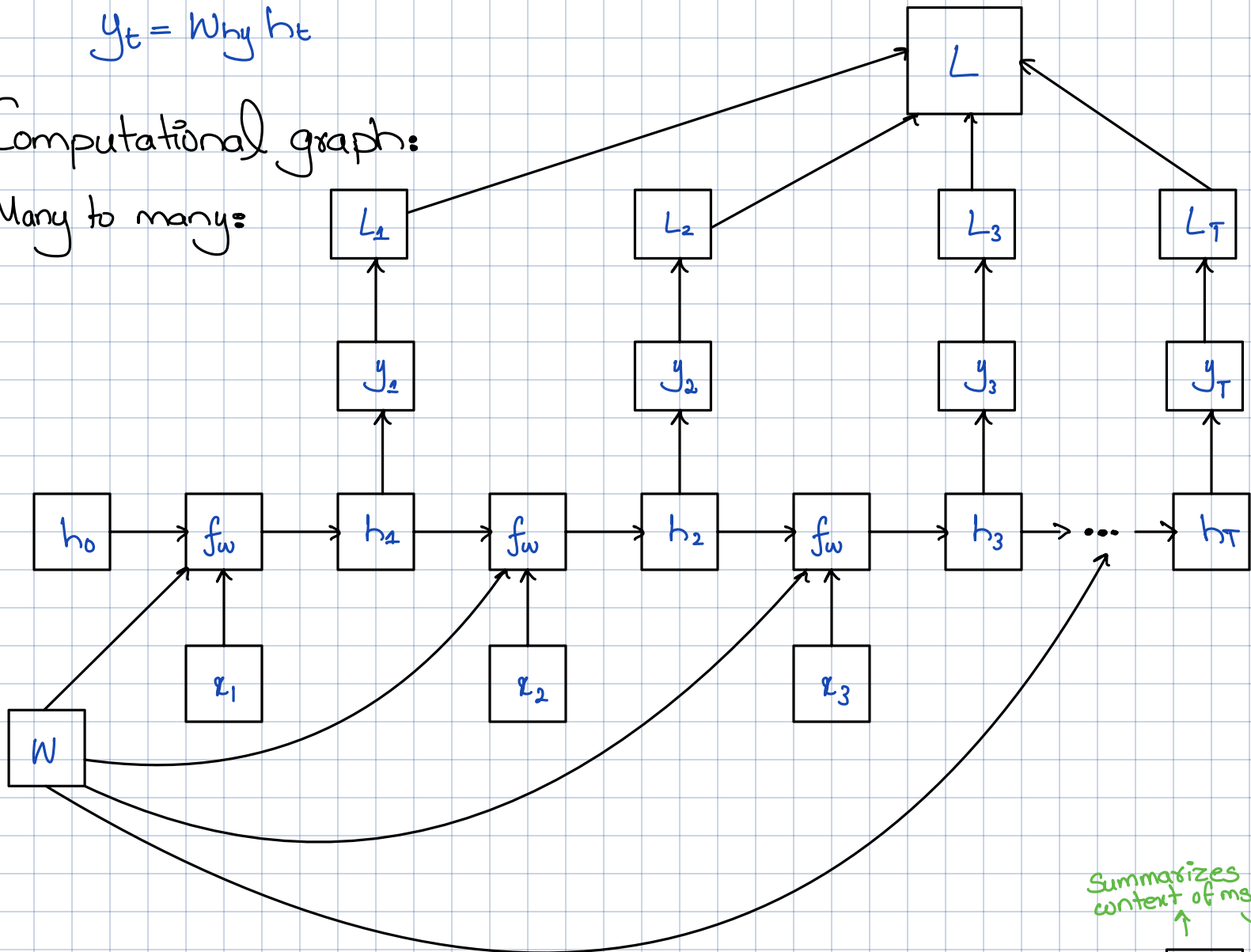
$$h_t = f_w(h_{t-1}, x_t)$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t) \quad \text{Non linearity}$$

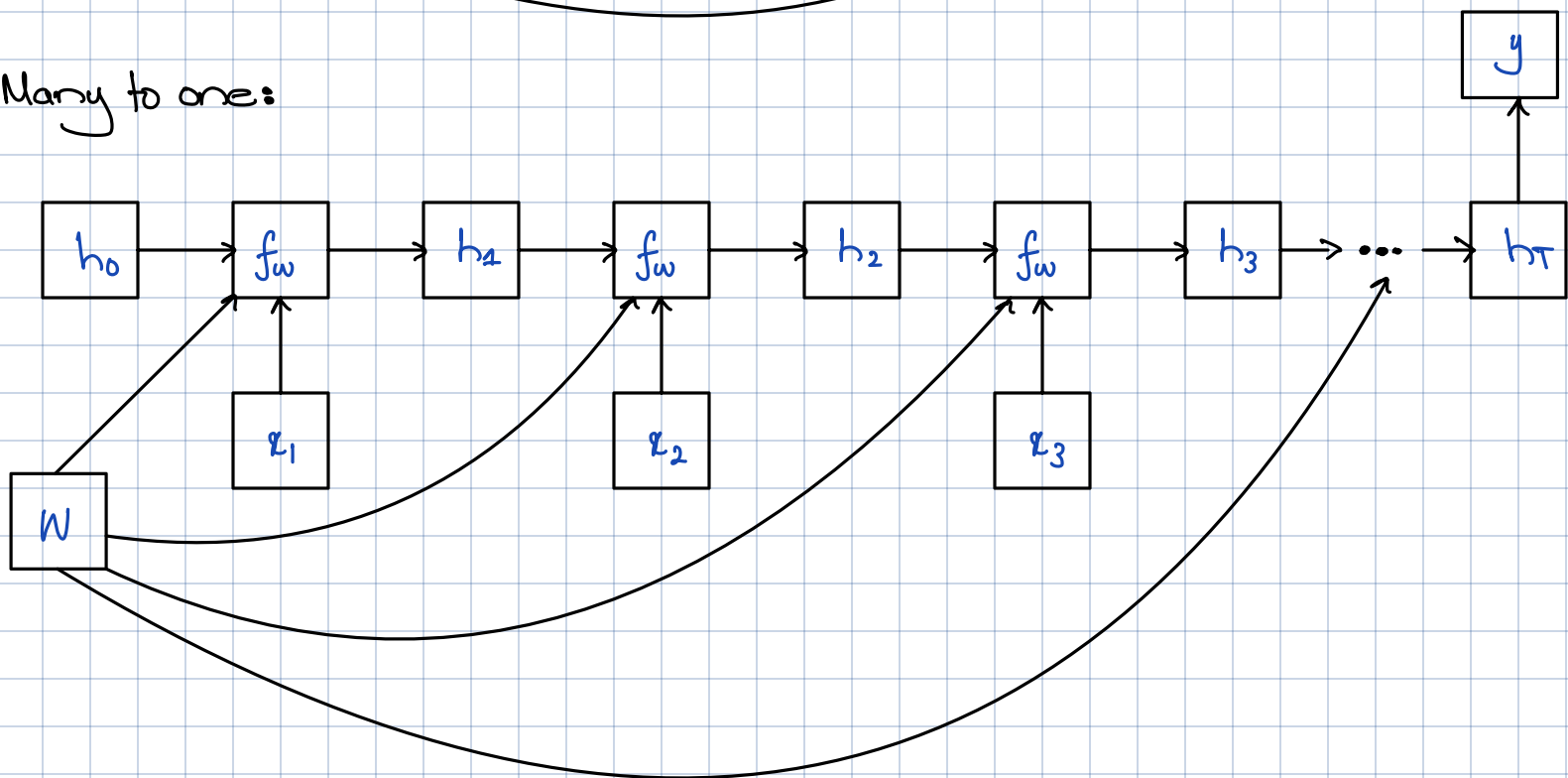
$$y_t = W h_t$$

Computational graph:

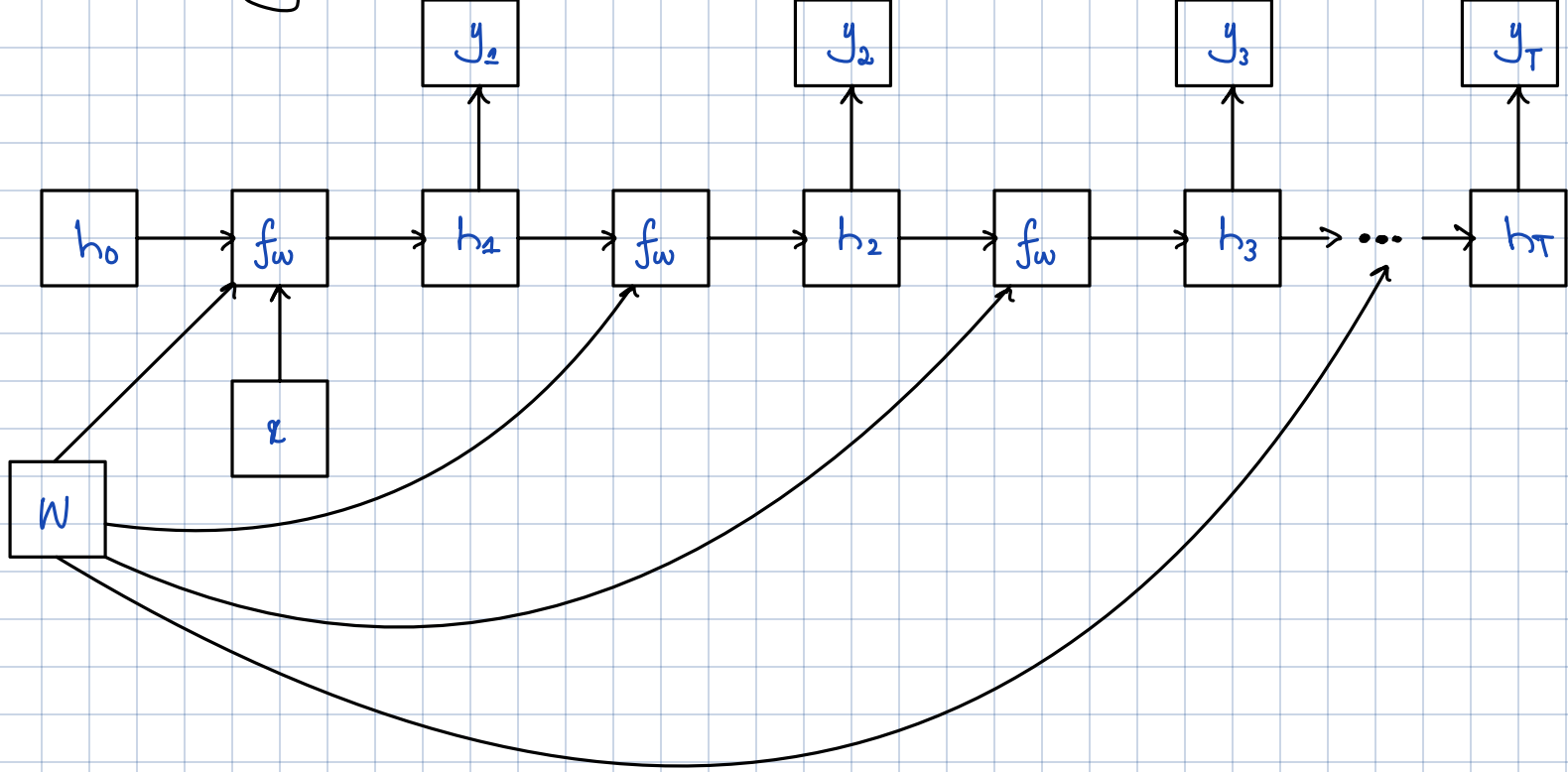
Many to many:



Many to one:

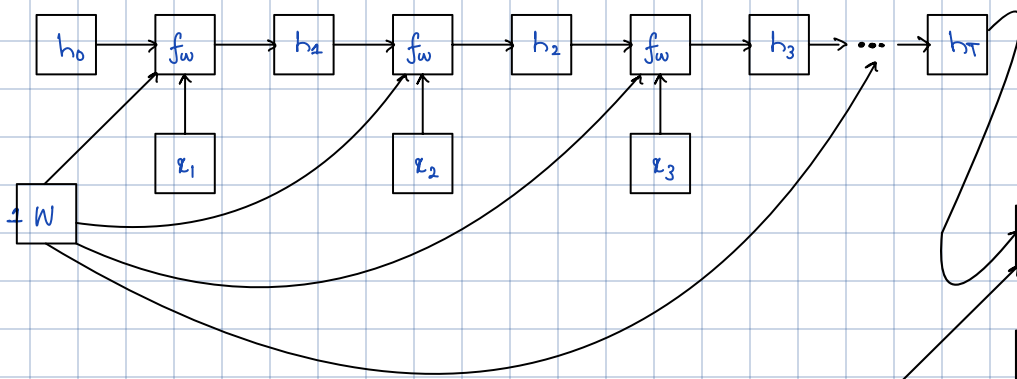


One to many:

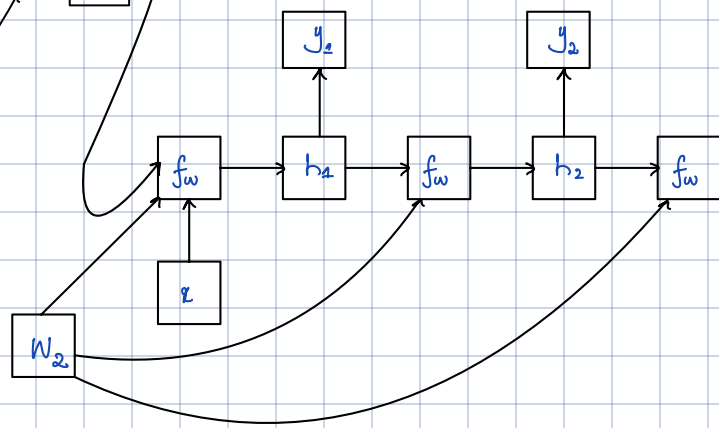


Sequence to Sequence (Many to one + One to many)

"Encoder"



"Decoder"

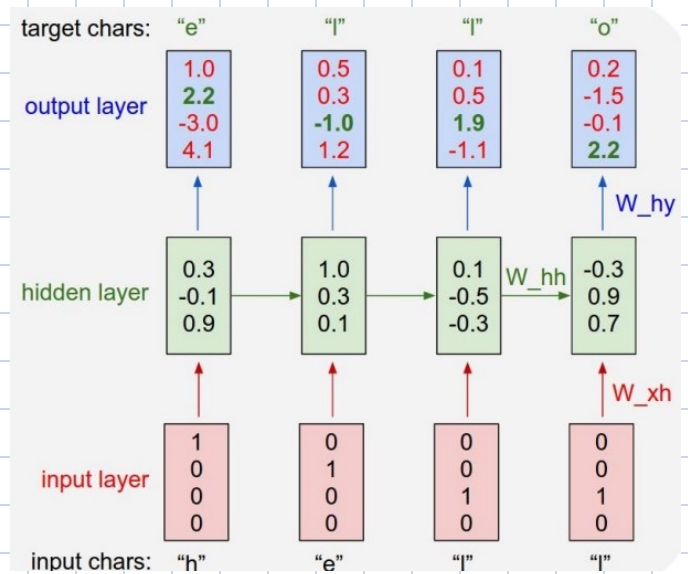


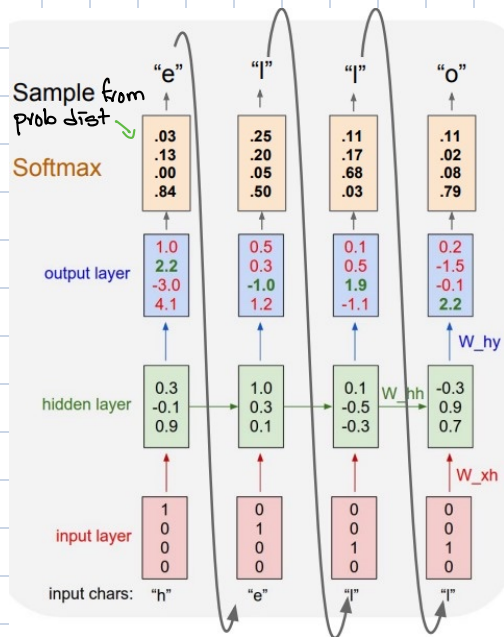
Character level language model:

Vocabulary: $[h, e, l, o]$

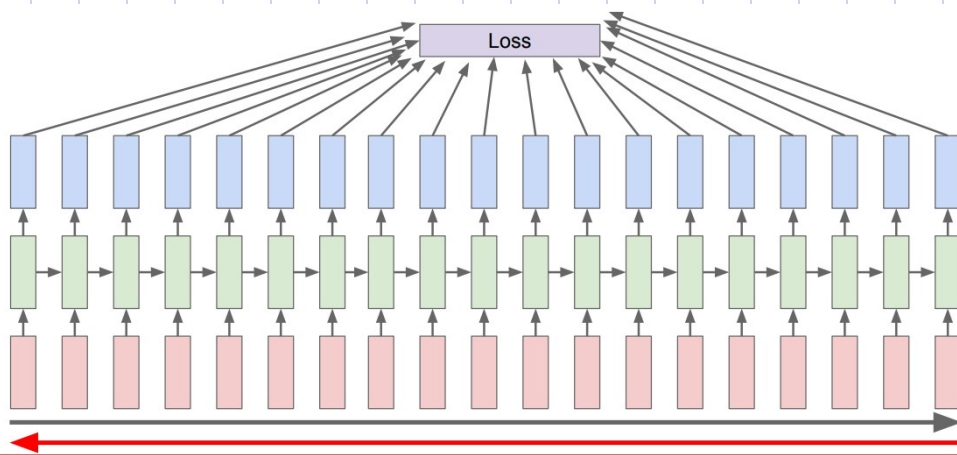
Example training sequence: "hello"

At test time, sample chars one at a time & feed back to model





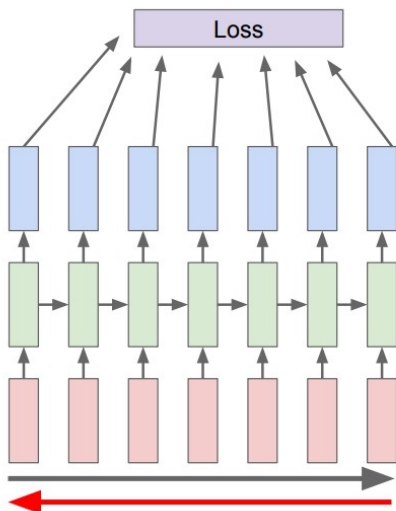
Backpropagation through time:



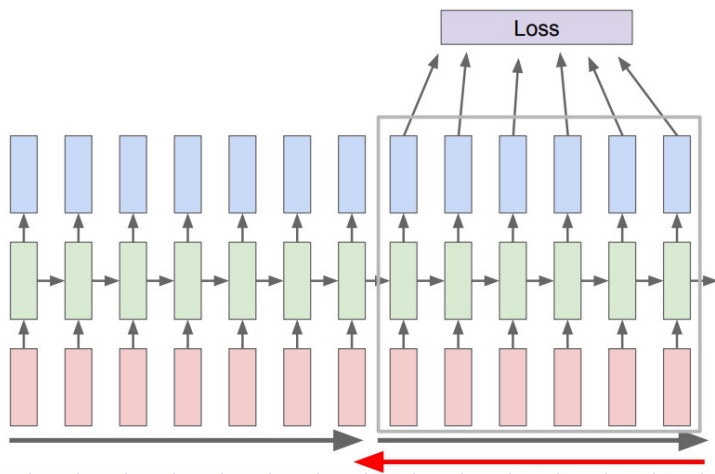
↳ forward through a entire sequence to compute loss, then backward through entire sequence to compute gradient

↳ Very bad. Very large sequence can slow model. So,

Truncated Backpropagation through time:

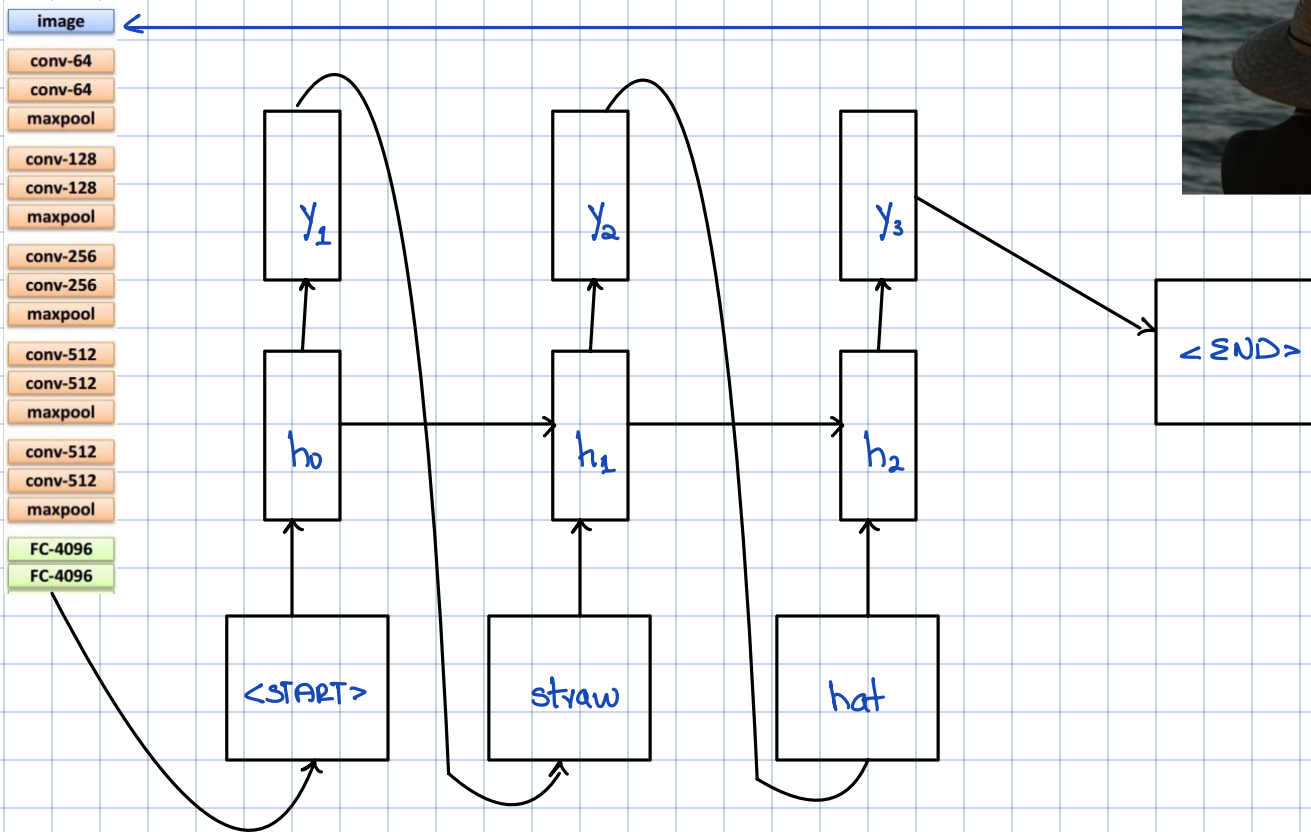


↳ run forward & backward through chunks of the sequence instead of whole sequence



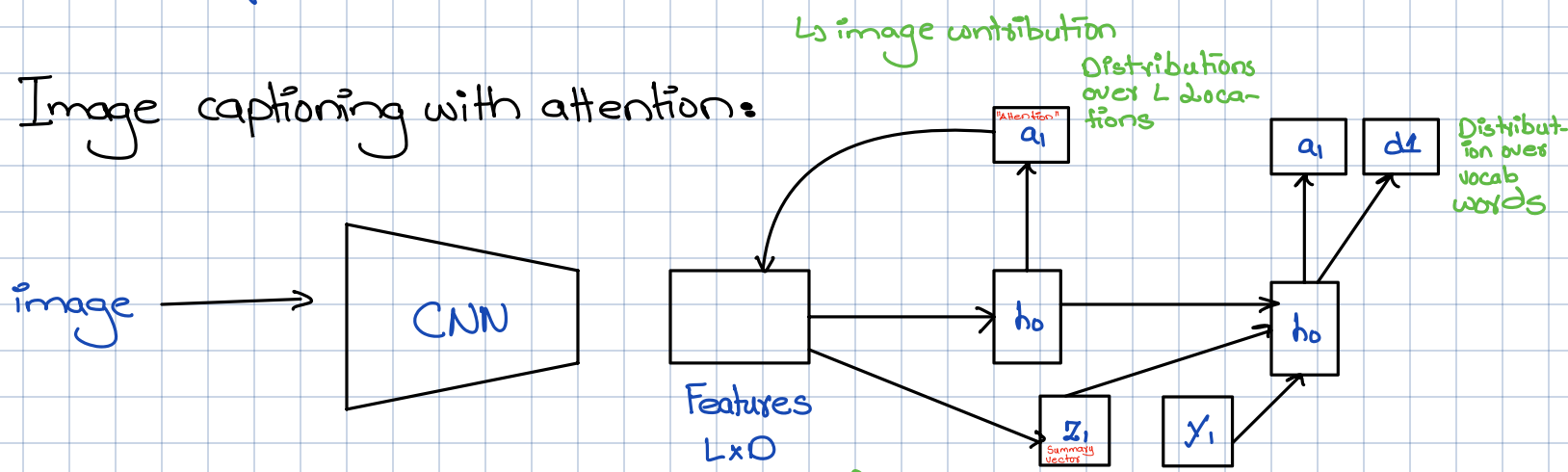
↳ Carry hidden states forward in time forever, but only backprop for some smaller no. of steps

Image Captioning:



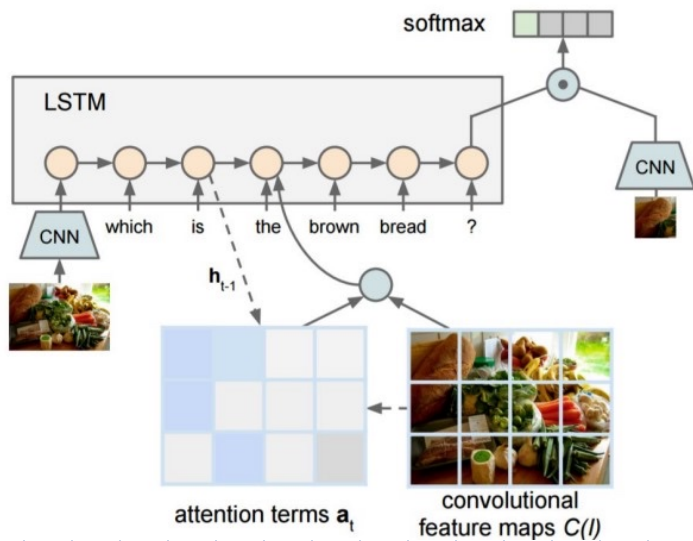
$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{vh} * v)$$

Image captioning with attention:



↳ One Vector for each spatial loc in image $z = \sum_{i=1}^L p_i \cdot v_i$

Visual Question Answering



Multi-layer RNN:

$$h_t^l = \tan W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

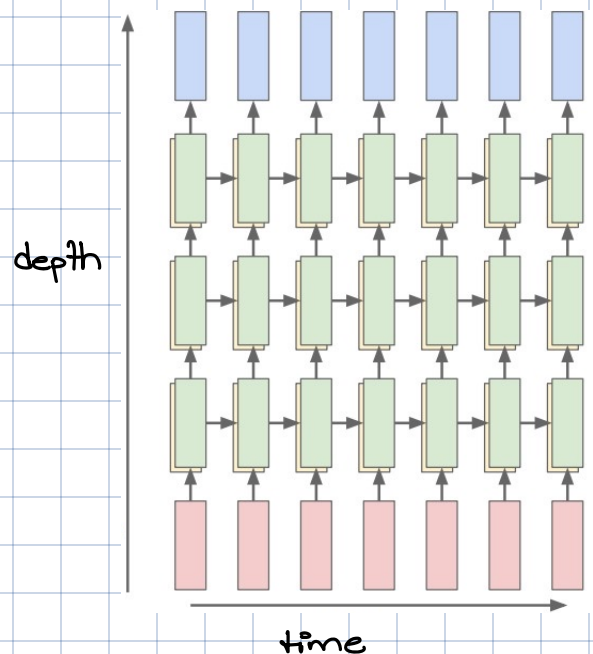
$$h \in \mathbb{R}^n \quad W^l [n \times 2n]$$

2STM:

$$\begin{pmatrix} i \\ f \\ o \\ q \end{pmatrix} = \begin{pmatrix} G \\ G \\ G \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

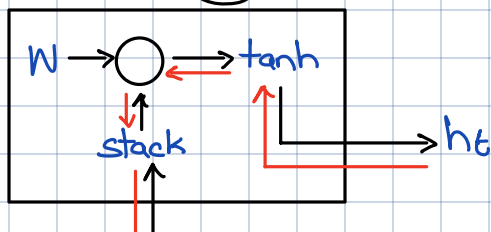
$$c_t^l = f \cdot c_{t-1}^l + i \cdot g$$

$$h_t^l = o \cdot \tanh(c_t^l)$$

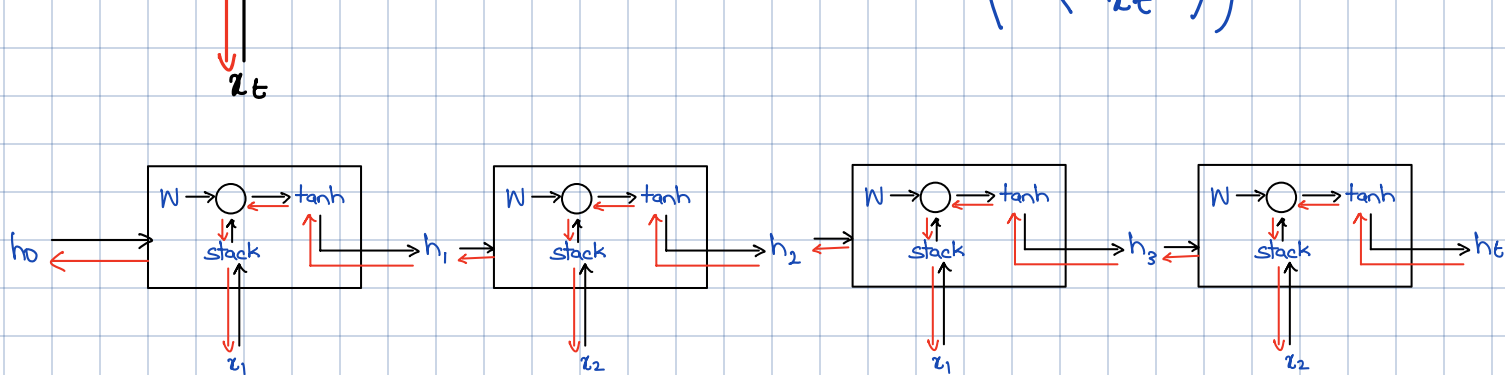


↳ input goes through RNN
↳ produces sequence of hidden states
↳ Use hidden states sequence as input

Vanilla RNN gradient flow:



$$\begin{aligned} h_t &= \tanh(W_{hh} h_{t-1} + W_{hx} x_t) \\ &= \tanh((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}) \\ &= \tanh(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}) \end{aligned}$$



Largest singular value $> 1 \rightarrow$ Grad clipping

Exploding grads

$$\text{grad-norm} = \text{np.sum}(\text{grad} * \text{grad})$$

if $\text{grad-norm} > \text{threshold}$:

$$\text{grad}^* = \left(\frac{\text{threshold}}{\text{grad-norm}} \right)$$

Largest singular value $< 1 \rightarrow$ Change architecture (LSTM)

Vanishing grads

LSTM:

$$\begin{pmatrix} i \\ f \\ 0 \\ o \end{pmatrix} = \begin{pmatrix} G \\ G \\ G \\ \tanh \end{pmatrix} W \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^{l-1} \end{pmatrix}$$

cell state

$$c_t = f \cdot c_{t-1} + i \cdot g$$

$$h_t^l = o \cdot \tanh(c_t)$$

\hookrightarrow hidden state

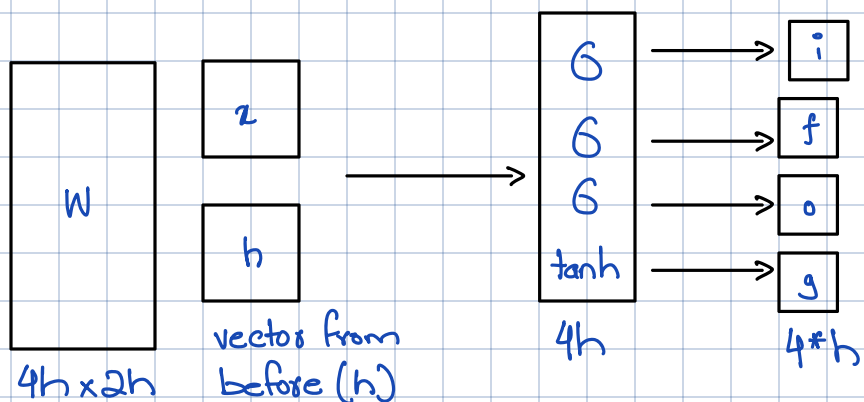
How much of prev cell state we need to forget + write if $i=1$ vice versa

f : forget gate, whether to erase cell

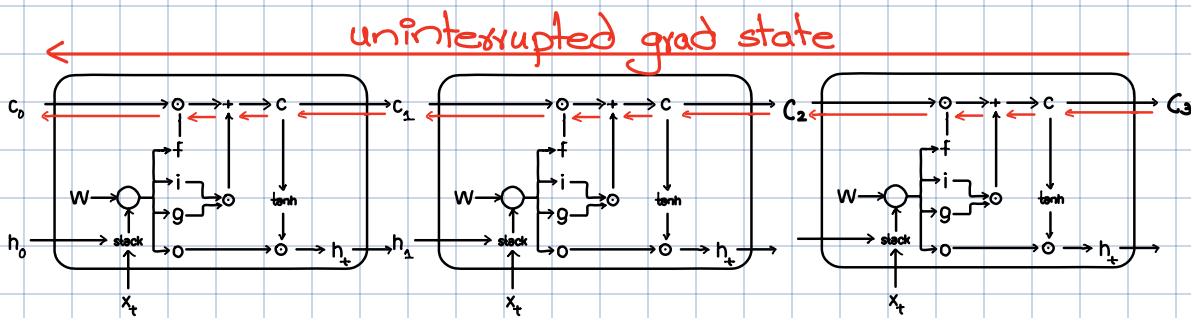
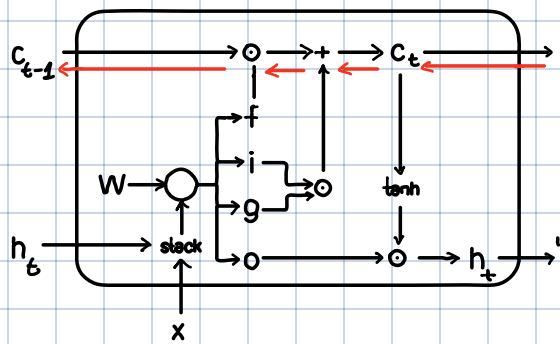
i : input gate, how much do we input into cell

g : gate gate, how much to write to cell

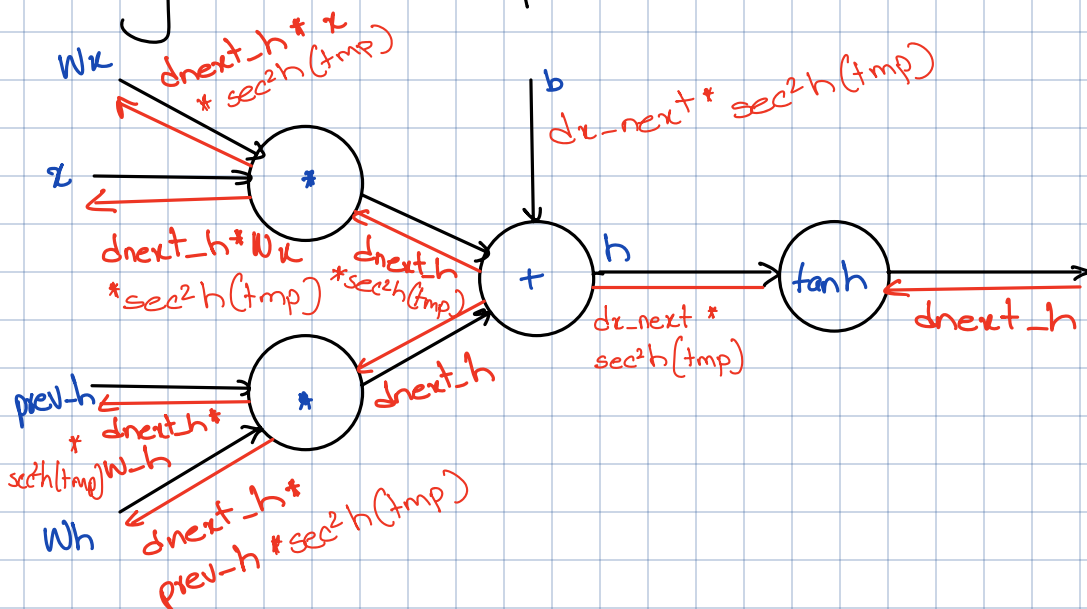
o : Output gate, how much to reveal cell



Computational graphs:

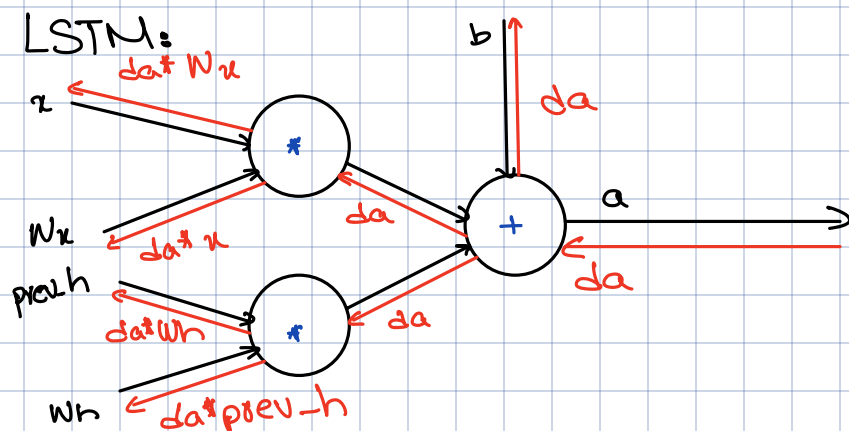


Assignment helpers:

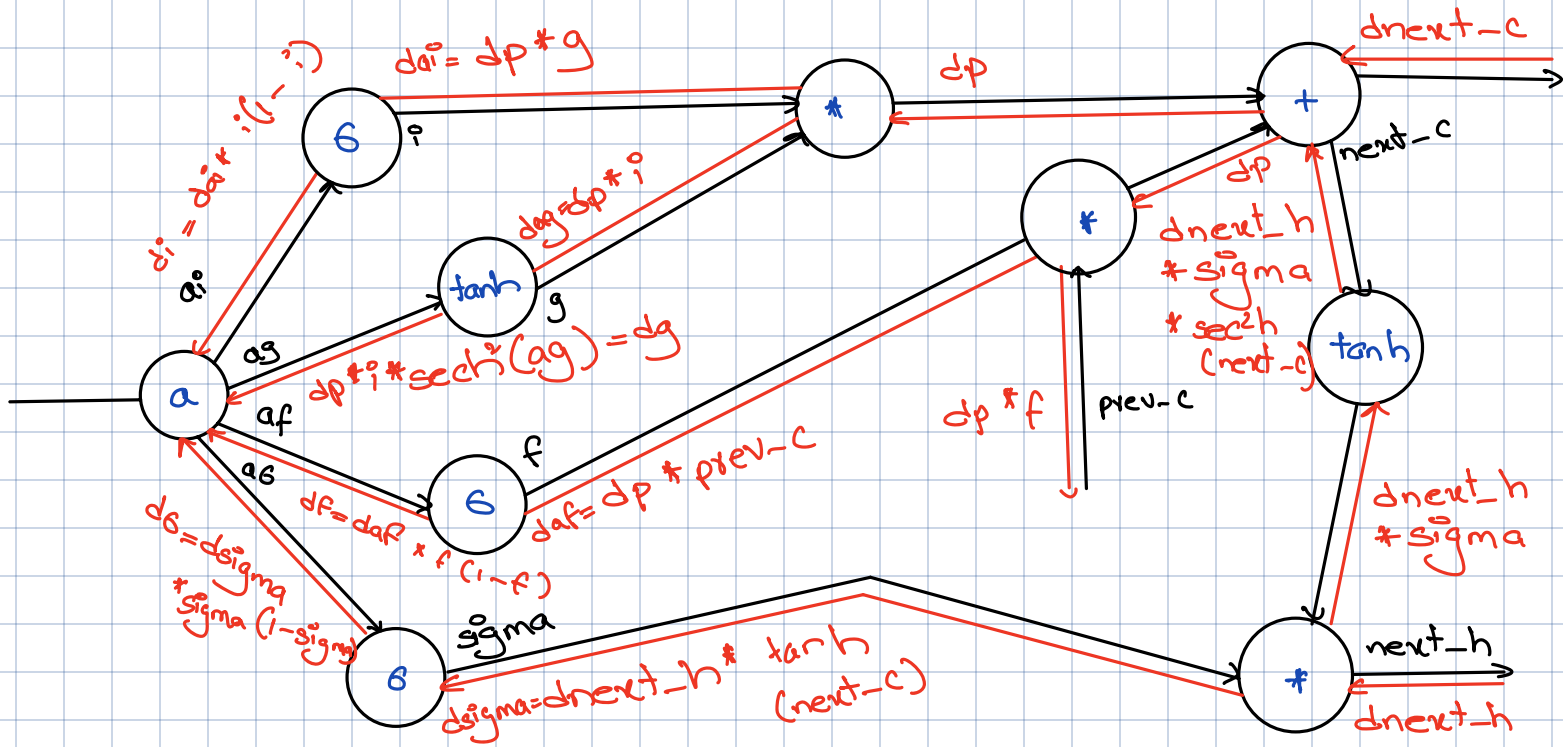


$$\frac{\partial \tanh x}{\partial x} = \text{sech}^2 x$$

$$tmp = Wh * prev_h + Wx * x + b$$



$$dp = dnext_h * \text{sigma} * \text{sech}^2(h(tmp)) + dnext_c$$



$$da = di + df + ds + dg$$