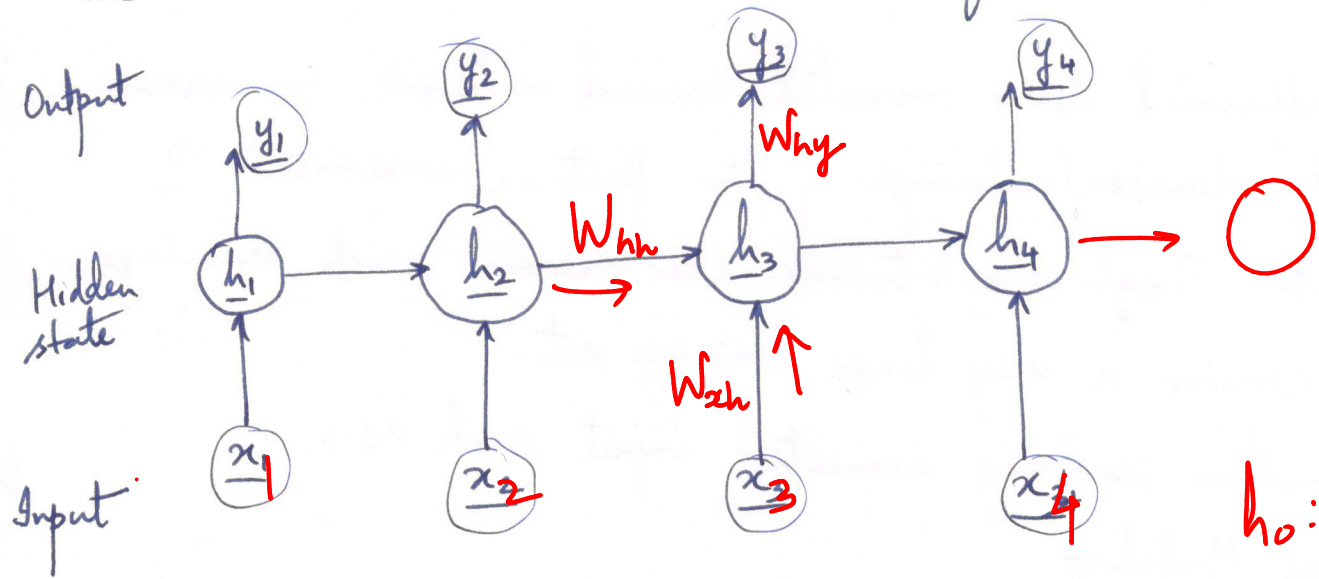


The recurrent network can be unfueled in time



$$\underline{h}_t = f(\underline{h}_{t-1}, \underline{x}_t)$$

$$\underline{y}_t = g(\underline{h}_t)$$

$$\underline{h}_1 = f(\underline{h}_0, \underline{x}_1)$$

$$\underline{h}_2 = f(f(\underline{h}_0, \underline{x}_1), \underline{x}_2)$$

$$\vdots$$

Since  $\underline{y}_t = g(\underline{h}_t)$

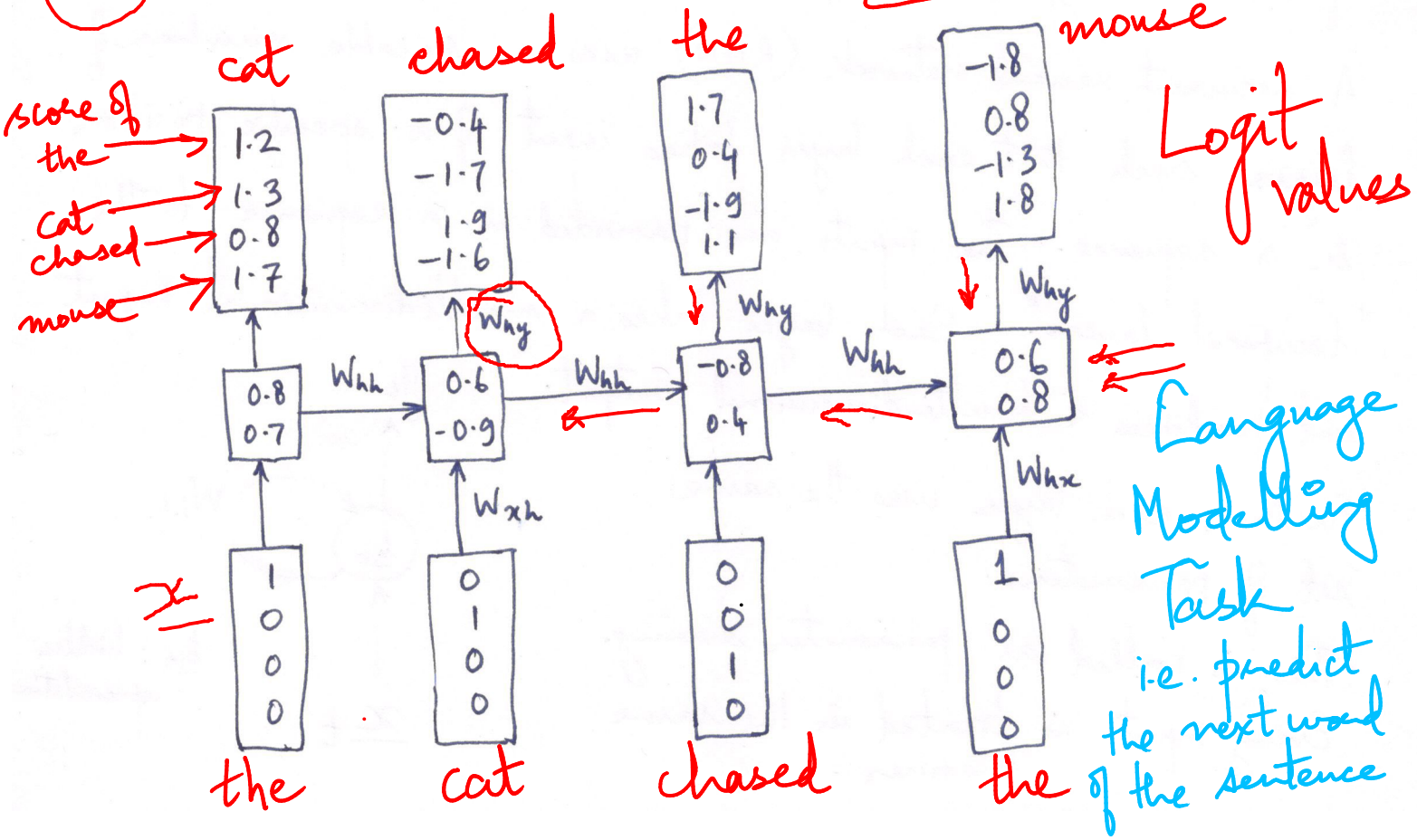
$$\underline{h}_t = \text{tanh}(W_{xh} \underline{x}_t + W_{hh} \underline{h}_{t-1})$$

nonlinearity

Affine

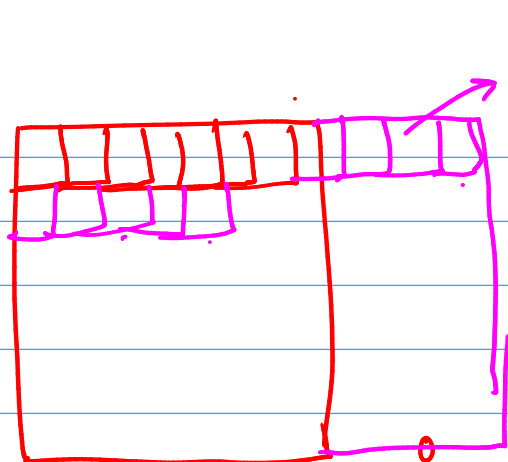
$$\underline{y}_t = W_{hy} \underline{h}_t$$

$$\underline{y}_t = F_t(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_t)$$



2D  
Grid

CNN



Convolution kernel  
remains the same

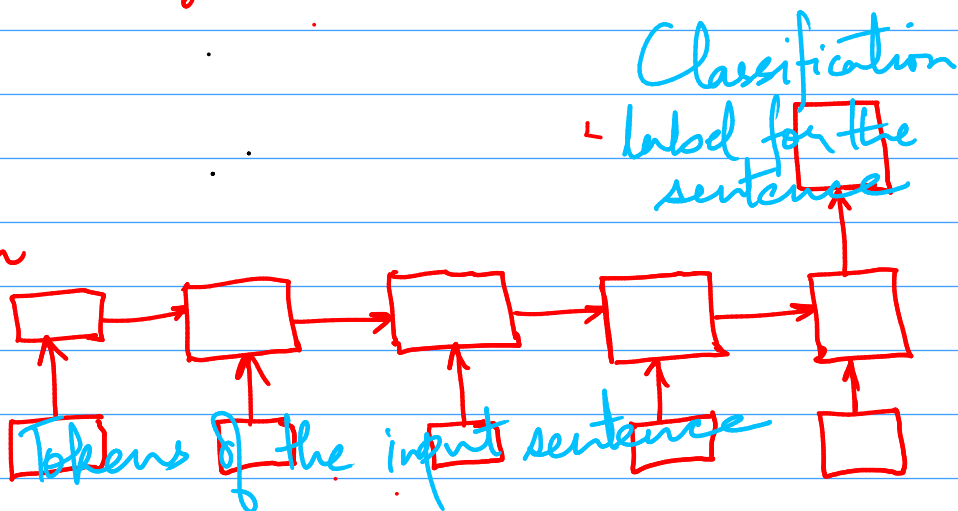
Image

Classification

0  
0  
0  
0  
0  
0

## Architectures

Sentence Classification



Machine  
Translation

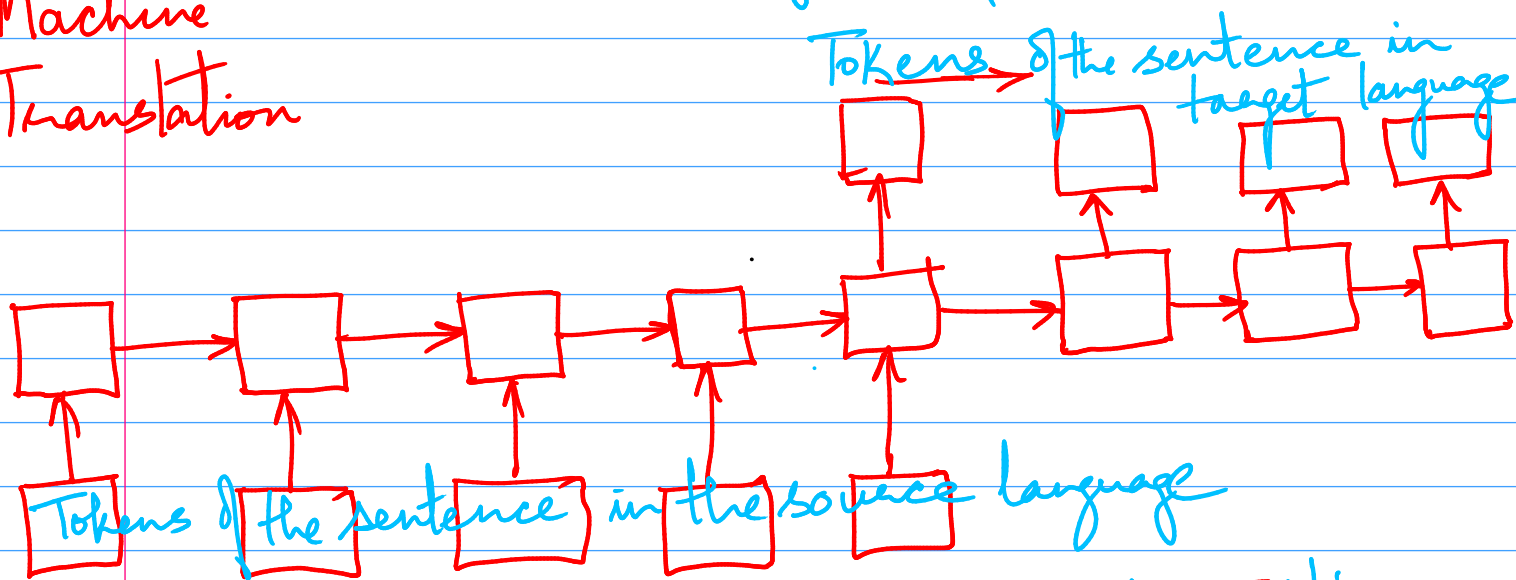
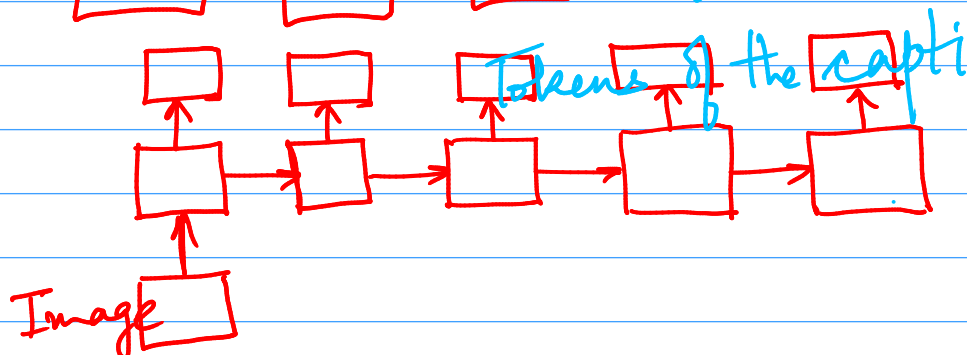


Image  
Captioning



# Backpropagation Through Time

The output vector  $y_t \equiv [y_t^1, y_t^2, \dots, y_t^d]$  is converted to a vector of probabilities using a softmax function

$$\begin{bmatrix} p_t^1 & \dots & p_t^d \end{bmatrix} = \underline{\text{Softmax}}([y_t^1, y_t^2, \dots, y_t^d])$$

The loss function for all  $T$  time stamps is

$$L = - \sum_{t=1}^T \log(p_t^{j_t})$$

$j_t$ : index of the correct label as given in GT.

The derivative of the loss w.r.t. the raw outputs is computed

as

$$\frac{\partial L}{\partial y_t^k} = \underline{p_t^k} - \underline{I(k, j_t)}$$

(refer to the derivation of the gradient of the CE loss when  $k$  and  $j_t$  are

The indicator function  $I$  gives an output 1 when  $k$  and  $j_t$  are the same.

To handle the shared weights we introduce temporal variables  $W_{xh}^{(t)}$ ,  $W_{hh}^{(t)}$  and  $W_{hy}^{(t)}$  for time-stamp  $t$ .

We first perform conventional backpropagation pretending that the variables are distinct from one another.

Finally, a unified update for each weight parameter is computed by adding the contributions of the temporal versions of the variables.



## Steps

1. Run the input sequentially in the forward direction through time and compute the errors and the negative log loss of softmax layer at each time-stamp.
2. Use conventional backpropagation to compute

$$\frac{\partial L}{\partial W_{xh}^{(t)}} \quad \frac{\partial L}{\partial W_{hh}^{(t)}} \quad \frac{\partial L}{\partial W_{hy}^{(t)}}$$

3. Compute  $\frac{\partial L}{\partial W_{xh}} = \sum_{t=1}^T \frac{\partial L}{\partial W_{xh}^{(t)}} \frac{\partial W_{xh}^{(t)}}{\partial W_{xh}}$  } Temporal aggregation wrapped around conventional backpropagation

$$\frac{\partial L}{\partial W_{hh}} = \sum_{t=1}^T \frac{\partial L}{\partial W_{hh}^{(t)}} \frac{\partial W_{hh}^{(t)}}{\partial W_{hh}}$$
$$\frac{\partial L}{\partial W_{hy}} = \sum_{t=1}^T \frac{\partial L}{\partial W_{hy}^{(t)}} \frac{\partial W_{hy}^{(t)}}{\partial W_{hy}}$$

Implicitly we have set  $\frac{\partial W_{xy}^{(t)}}{\partial W_{xy}} = 1. = \frac{\partial W_{xh}^{(t)}}{\partial W_{xh}} = \frac{\partial W_{hh}^{(t)}}{\partial W_{hh}} = \frac{\partial W_{hy}^{(t)}}{\partial W_{hy}}$

## Truncated Backpropagation through time:

Backpropagation updates are done only over segments of the sequence over fixed (modest) length.

Only the portion of the loss over the relevant segment is used to compute the gradients and update the weights.