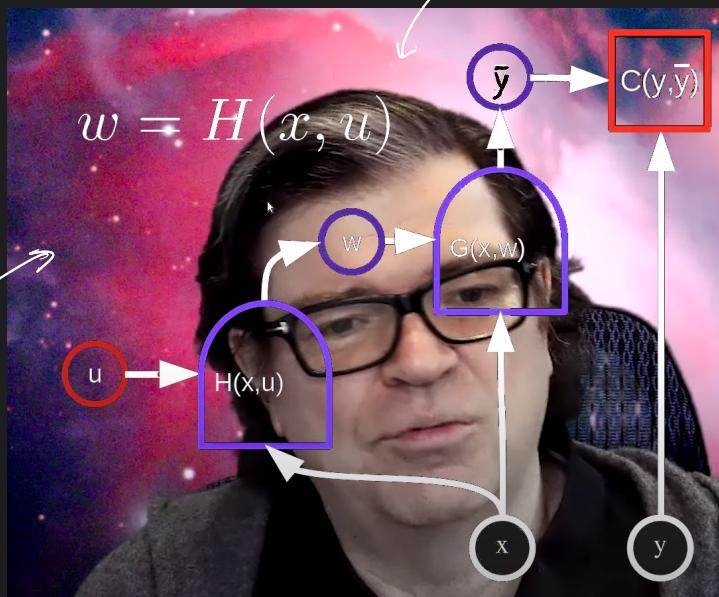


Hyper Networks

cool background

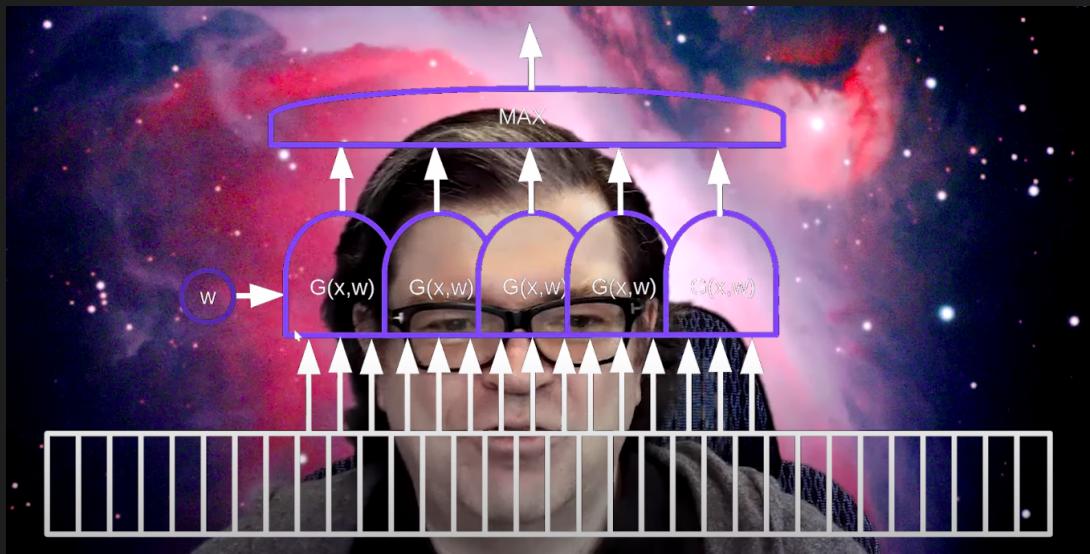
good Sir,



- Weights w are the output of another network $H(x,u)$

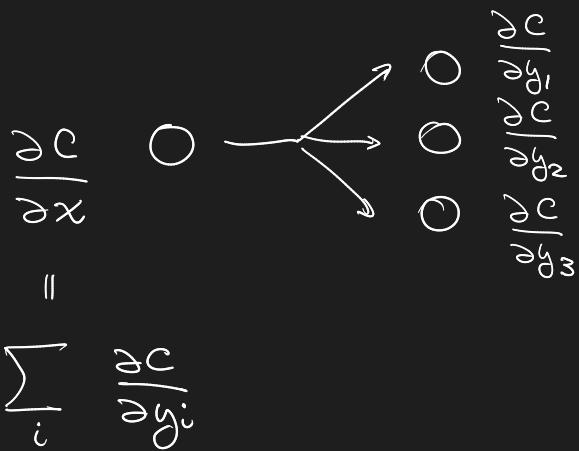
Another idea: Shared weights

- Searching for a motif / pattern.



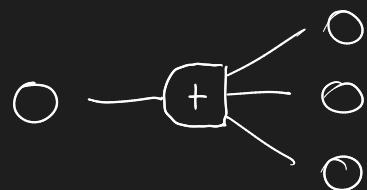
Because the weights are shared, while doing backprop, the contribution of each replica of the network (each gradient) is added, to update the only weights w

Backward pass of replicated variable



Branching up in the forward prop.

correspond to a sum in the backward prop.



And vice versa,

- Q : "Why we sum the gradients?"

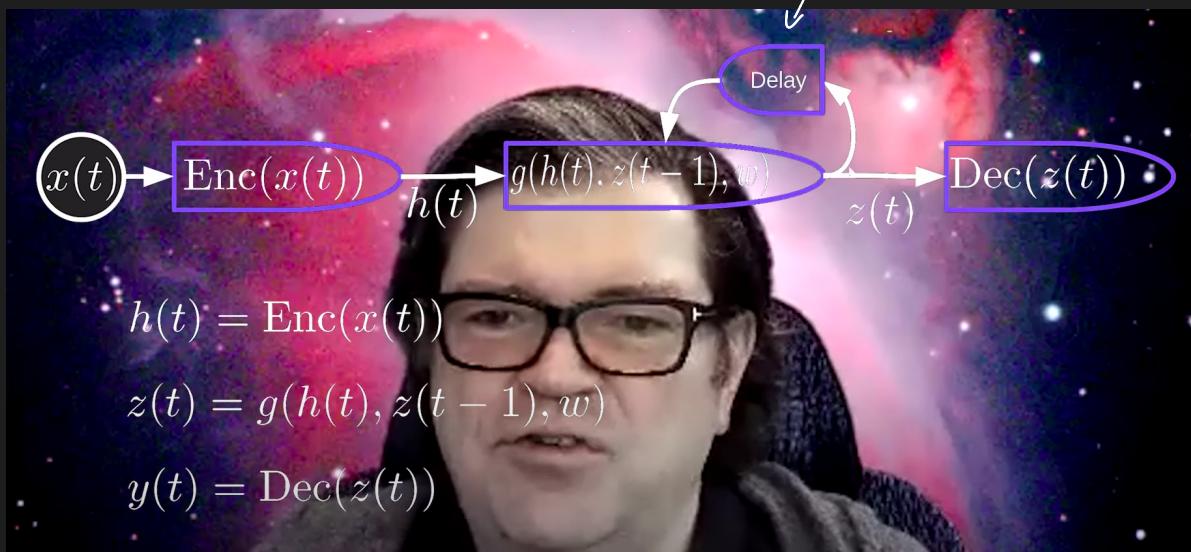
$$\delta C = \frac{\partial C}{\partial y_1} \delta y_1 + \frac{\partial C}{\partial y_2} \delta y_2 + \frac{\partial C}{\partial y_3} \delta y_3$$

$$\delta C = \delta x \left[\frac{\partial C}{\partial y_1} + \frac{\partial C}{\partial y_2} + \frac{\partial C}{\partial y_3} \right]$$

$$\frac{\partial C}{\partial x} = \sum_i \frac{\partial C}{\partial y_i}$$

Loops : RNNs

input at next timestep



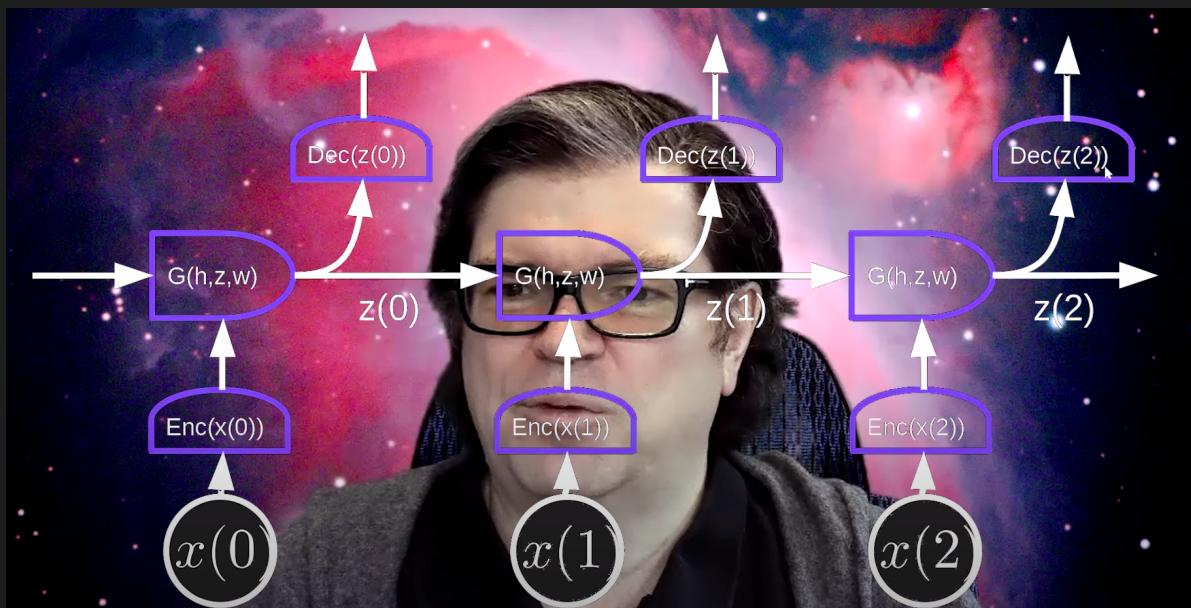
$$h(t) = \text{Enc}(x(t))$$

$$z(t) = g(h(t), z(t-1), w)$$

$$y(t) = \text{Dec}(z(t))$$

Allows the network to have \hookrightarrow memory.

Unrolling the network,



Difficult to train.

\hookrightarrow Long sequences : Vanishing/exploding gradients.

- Exploding : easy to solve (ie: add \Rightarrow sigmoid)
- Vanishing : not so easy.

Example

$$z_t \xrightarrow{W} z_{t+1} = \omega z_t$$

ω

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

eigen values
of 1

$$\begin{bmatrix} \omega & & \\ \cos \theta & \sin \theta & \\ -\sin \theta & \cos \theta & \end{bmatrix}$$

Rotation matrix

They do not change the norm after doing the prod
 $\|z_t\| = \|\omega z_t\| = \|z_{t+1}\|$

Now, if $\omega = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

$$\|z_{t+1}\| = 2 \cdot \|z_t\|$$

explosion

$$\omega = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$$

$$\|z_{t+1}\| = \frac{1}{2} \|z_t\|$$

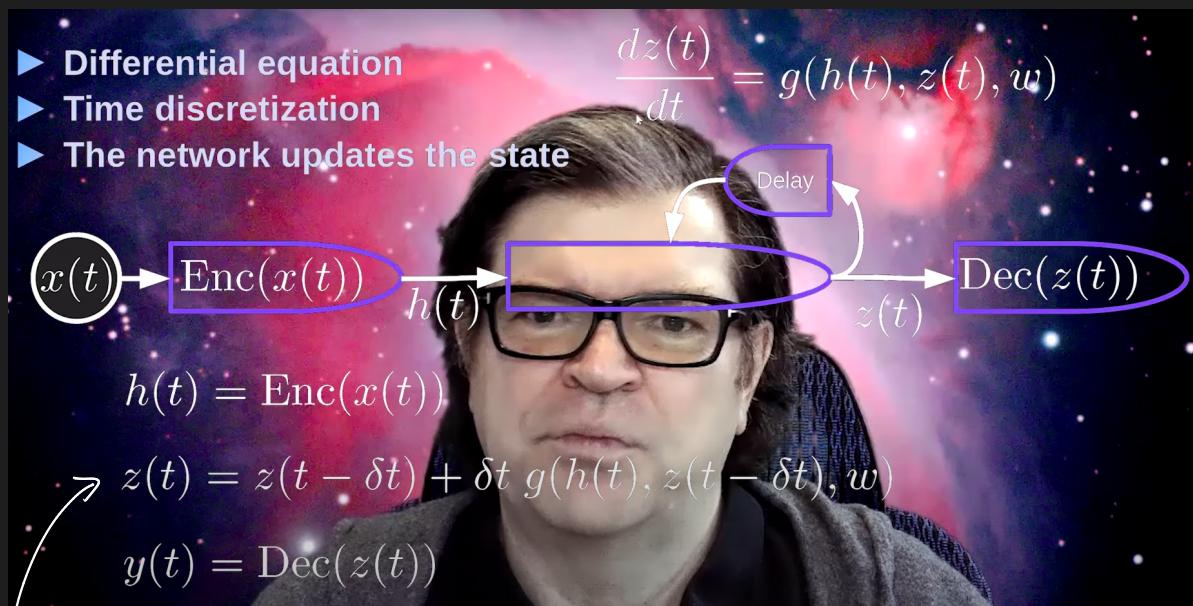
Vanishing
(shrinkes)

RNN tricks

- ▶ [Pascanu, Mikolov, Bengio, ICML 2013; Bengio, Boulanger & Pascanu, ICASSP 2013]
- ▶ Clipping gradients (avoid exploding gradients)
- ▶ Leaky integration (propagate long-term dependencies)
- ▶ Momentum (cheap 2nd order)
- ▶ Initialization (start in right ballpark avoids exploding/vanishing)
- ▶ Sparse Gradients (symmetry breaking)
- ▶ Gradient propagation regularizer (avoid vanishing gradient)
- ▶ LSTM self-loops (avoid vanishing gradient)

Other type of RNNs

- Neural ODEs



discretization of the differential equation.

GRU (Gated Recurrent Units)

► Recurrent nets quickly “forget” their state

► Solution: explicit memory cells

► GRU [Cho arXiv:1406.1078]

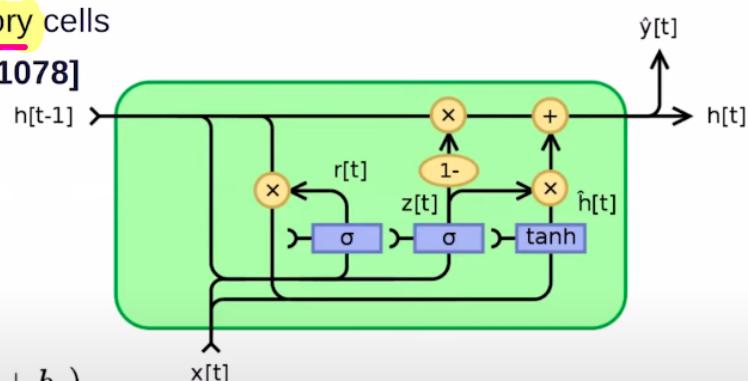
x_t : input vector

h_t : output vector

z_t : update gate vector

r_t : reset gate vector

W, U and b : parameter mat

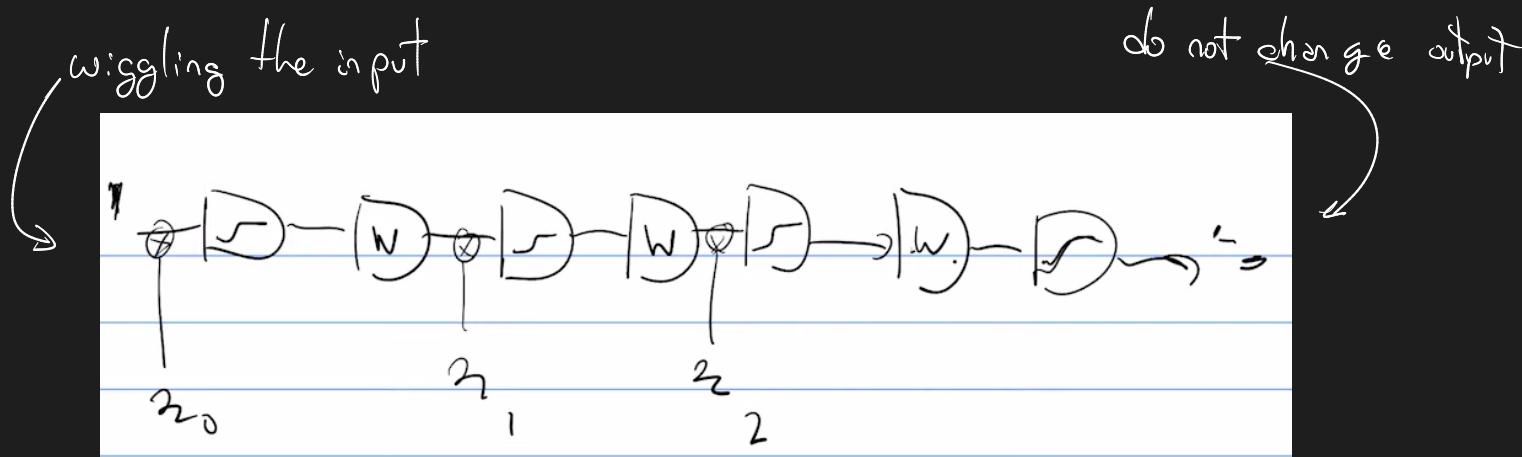


$$z_t = \sigma_g(W_z x_t + U_z h_{t-1} + b_z)$$

$$r_t = \sigma_g(W_r x_t + U_r h_{t-1} + b_r)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \phi_h(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

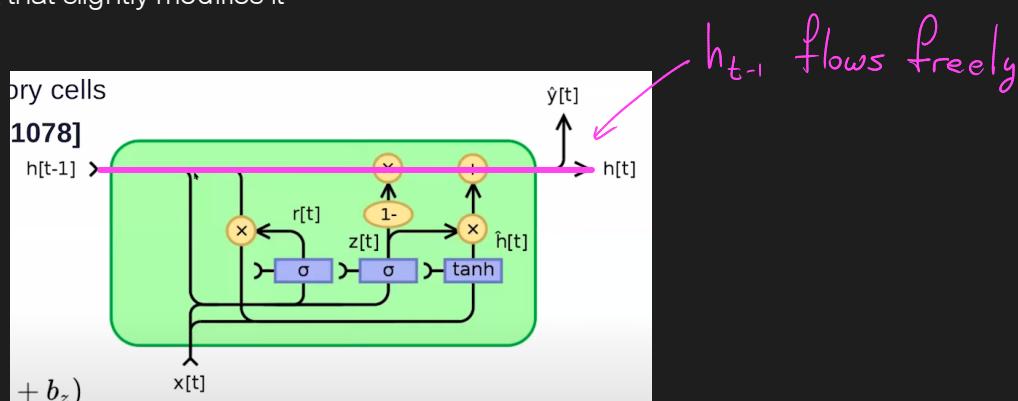
On RNNs

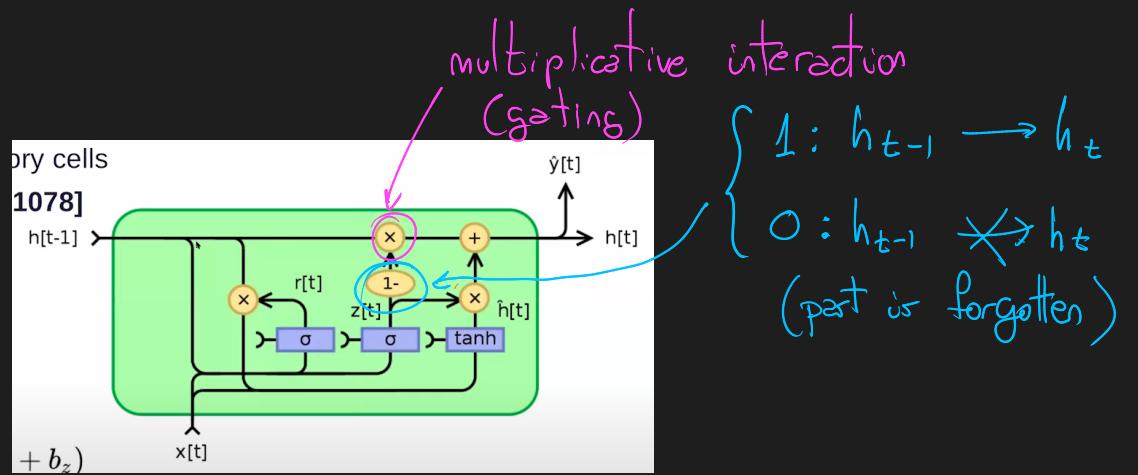


if the network state is saturated,

With GRUs

we can make the system remember by default (system=identity function by default: new state is equal to old state) and have a neural network that slightly modifies it





LSTM (Long Short-Term Memory)

Y. LeCun

- ▶ Recurrent nets quickly “forget” their state

▶ Solution: explicit memory cells

- ▶ LSTM [Hochreiter & Schmidhuber 97]

$x_t \in \mathbb{R}^d$: input vector to the LSTM unit

$f_t \in \mathbb{R}^h$: forget gate's activation vector

$i_t \in \mathbb{R}^h$: input/update gate's activation vector

$o_t \in \mathbb{R}^h$: output gate's activation vector

$h_t \in \mathbb{R}^h$: hidden state vector also known as output

$c_t \in \mathbb{R}^h$: cell state vector

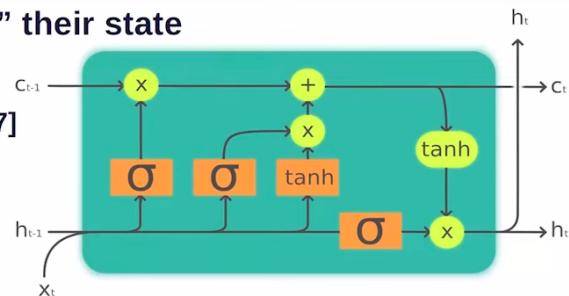
$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \sigma_c(W_c x_t + U_c h_{t-1} + b_c)$$

$$h_t = o_t \circ \sigma_h(c_t)$$



Legend:



Guillaume Chevalier <https://commons.wikimedia.org/w/index.php?curid=71836793>

- ▶ Sequence encoder → sequence decoder

Applications:

- ▶ Translation [Sutskever NIPS 2014]

- ▶ Multi-layer LSTM



- ▶ Sequence encoder → sequence decoder with attention
- ▶ Applications:
- ▶ Translation [Bahdanau, Cho, Bengio ArXiv:1409.0473]

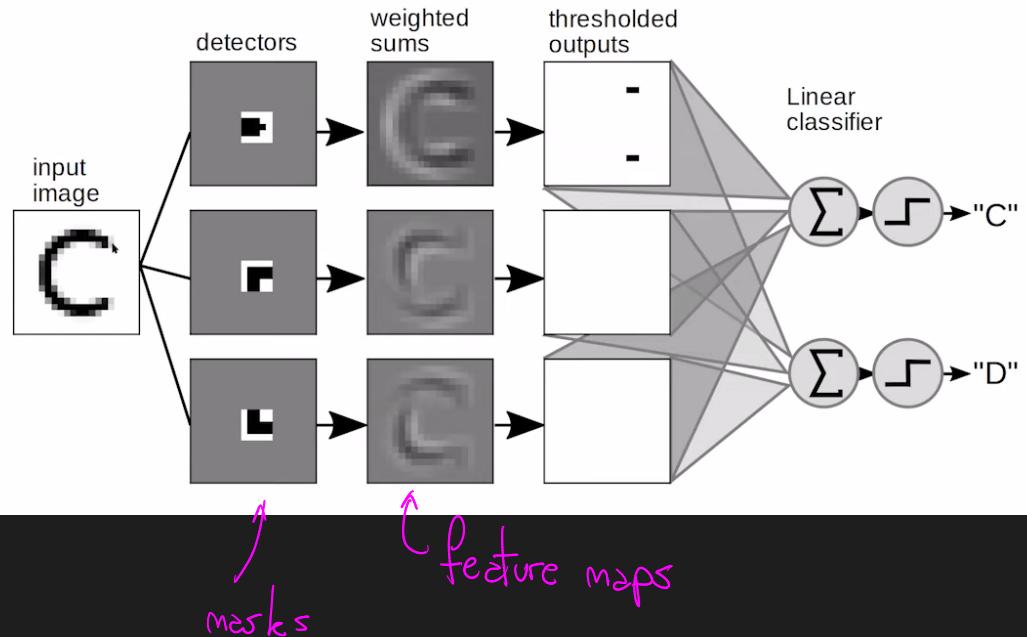


Convolutional Networks

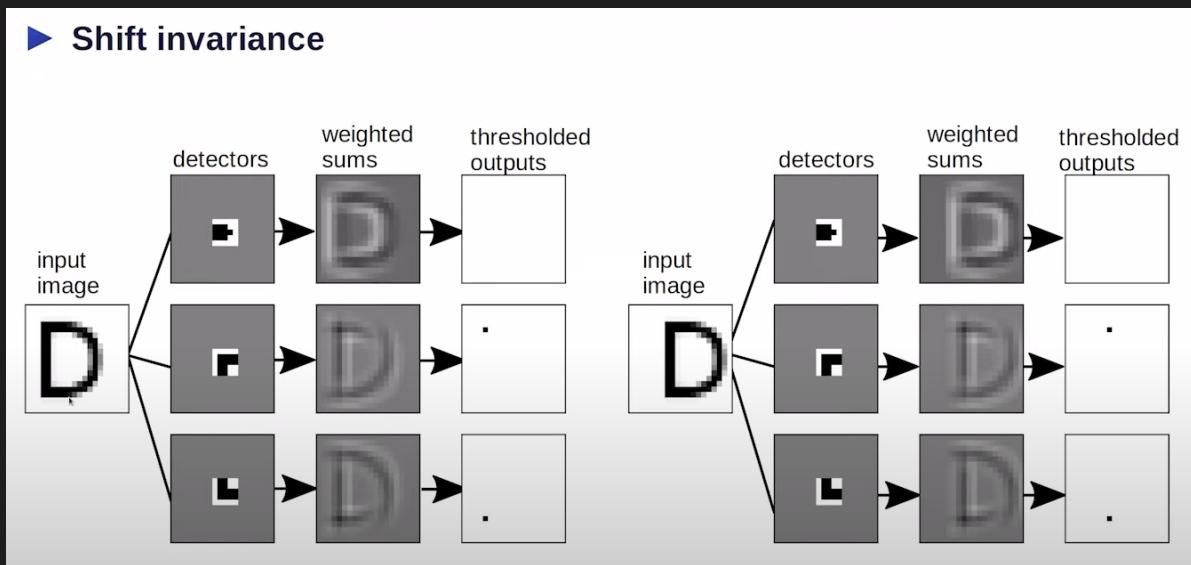
Hierarchical representation of the world

Detecting Motifs in Images

- ▶ Swipe “templates” over the image to detect motifs



► Shift invariance



Discrete Convolution: Swipe a small pattern of coefficient over an image

Def: Convolution (1D)

$$y_i = \sum_j w_j x_{i-j}$$

In practice: Cross-correlation

(PyTorch)

$$y_i = \sum_j w_j x_{i+j}$$

In 2D

$$y_{ij} = \sum_{kl} w_{kl} x_{i+k, j+l}$$

Backpropagating through convolutions

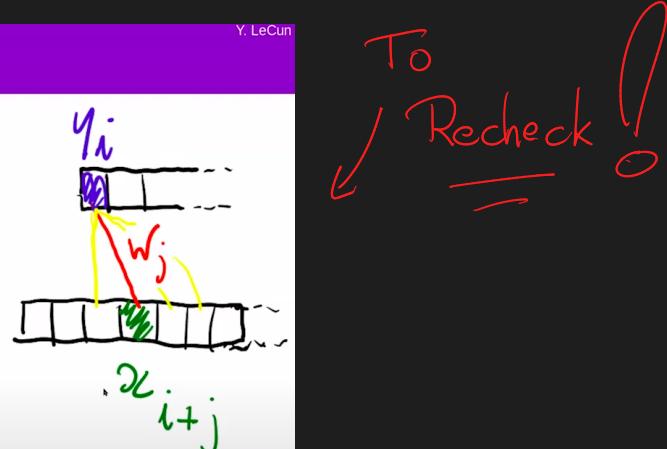
Y. LeCun

- Convolution $y_i = \sum_j w_j x_{i+j}$
- (really: cross-correlation)

- Backprop to input
- Sometimes called "back-convolution"

$$\frac{\partial C}{\partial x_j} = \sum_k w_k \frac{\partial C}{\partial y_{j-k}}$$

- Backprop to weights $\frac{\partial C}{\partial w_j} = \sum_i \frac{\partial C}{\partial y_i} x_{i+j}$



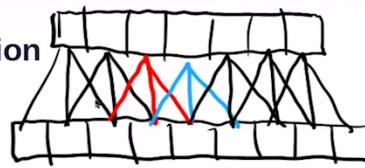
Stride and Skip: subsampling and convolution "à trous"

T. LeCun

- ▶ Regular convolution

- ▶ "dense"

$\text{kernel} = 3$
 $\text{stride} = 1$

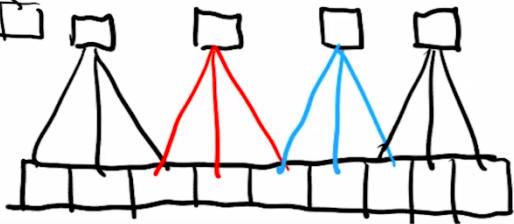


- ▶ Stride

- ▶ subsampling convolution

- ▶ Reduces spatial resolution

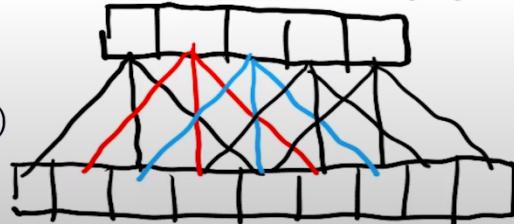
$\text{kernel} = 3, \text{ stride} = 2$



- ▶ Skip, convolution "à trous"

- ▶ pronounced "ah troo" (means "with holes")

- ▶ Dimensionality reduction without loss of resolution

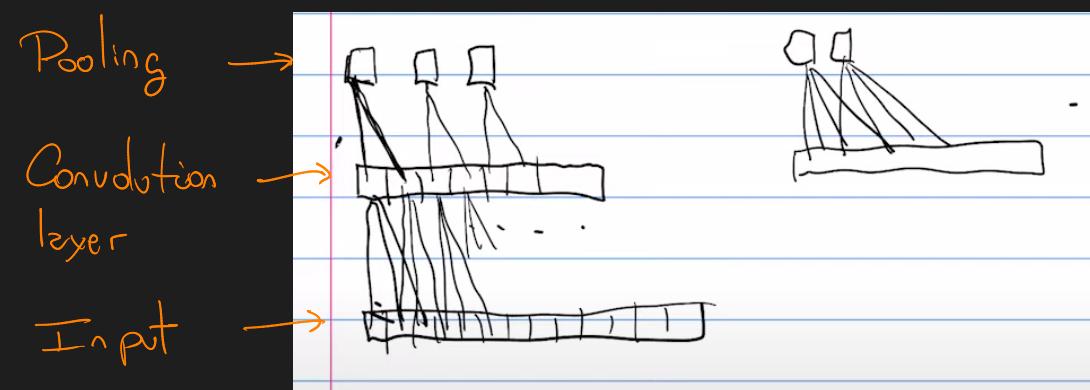


window of 5, but with holes

$\text{kernel} = 3$

$\text{skip} = 1$

$\text{stride} = 1$



Pooling

↳ Aggregation

↳ Permutation invariant ($f = \text{sum}, \text{avg}, \text{max}, L_p,$

most popular
↓ $P=2$

$$y = \left(\sum_n z_n^p \right)^{1/p} \quad L_p \text{ pooling}$$

$$y = \max_n (z_n) \quad \text{max pooling}$$

$$y = \frac{1}{P} \log \sum_n e^{P z_n} \quad \text{log sum exp pooling}$$

log sum exp

f invariant to $T : f(T(x)) = f(x)$

f equivariant to $T : f(T(x)) = T'(f(x))$

Convolution op. is equiv. to translation

$$\tau(x)_i = x_{i+\tau}$$

Pooling is locally invariant to translation

