

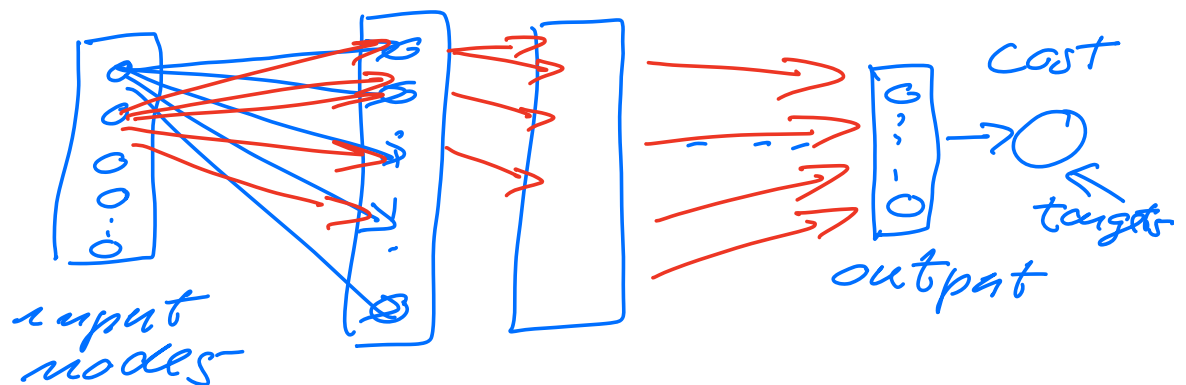
ML ERASMUS, JAN 15, 2023

---

Deep learning methods:

— FFNN = Feed Forward NN

Fully connected



- Dense matrix with weights + biases, Affine operations
- Data is well structured and homogeneous
- CNN reduce dimensionality by filtering (convolution) does not work well with

data of unknown

- RNN recurrent NN can be used on data set with unknown length.
- time series
- Natural language studies

Time series :

$$m \frac{d^2 x}{dt^2} + \eta \frac{dx}{dt} + x(t) = F(t)$$

initial conditions,  $x_0 = x(t_0)$   
 $v_0 = v(t_0)$

$$v(t) = \frac{dx}{dt}$$

$$m \frac{dv}{dt} + \eta v + x = F$$

Discretize using Euler's method

$$x_{i+1} = x(t_i + \Delta t) = x_i + \Delta t \cdot v_i$$

$$v_{i+1} = v_i + \Delta t \cdot \dot{v}?$$

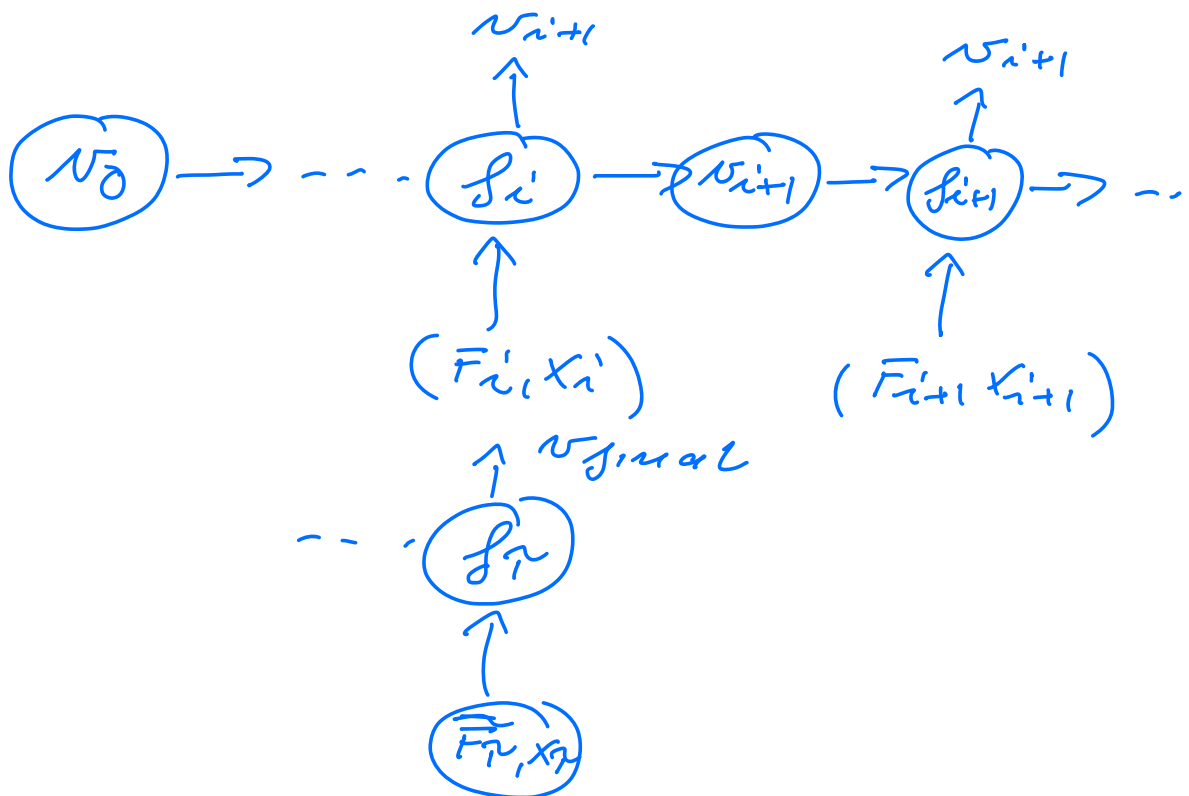
$$\begin{aligned} \left( \frac{dv}{dt} \right) &= - \underbrace{\left( \frac{v}{m} \right)}_{\alpha} v - \underbrace{\left( \frac{1}{m} \right)}_{\downarrow} x + \frac{F}{m} \\ &= \tilde{F} - \alpha v + \delta x \end{aligned}$$

$$v_{i+1} = v_i + \Delta t (\tilde{F}_i - \alpha v_i - \delta x_i)$$

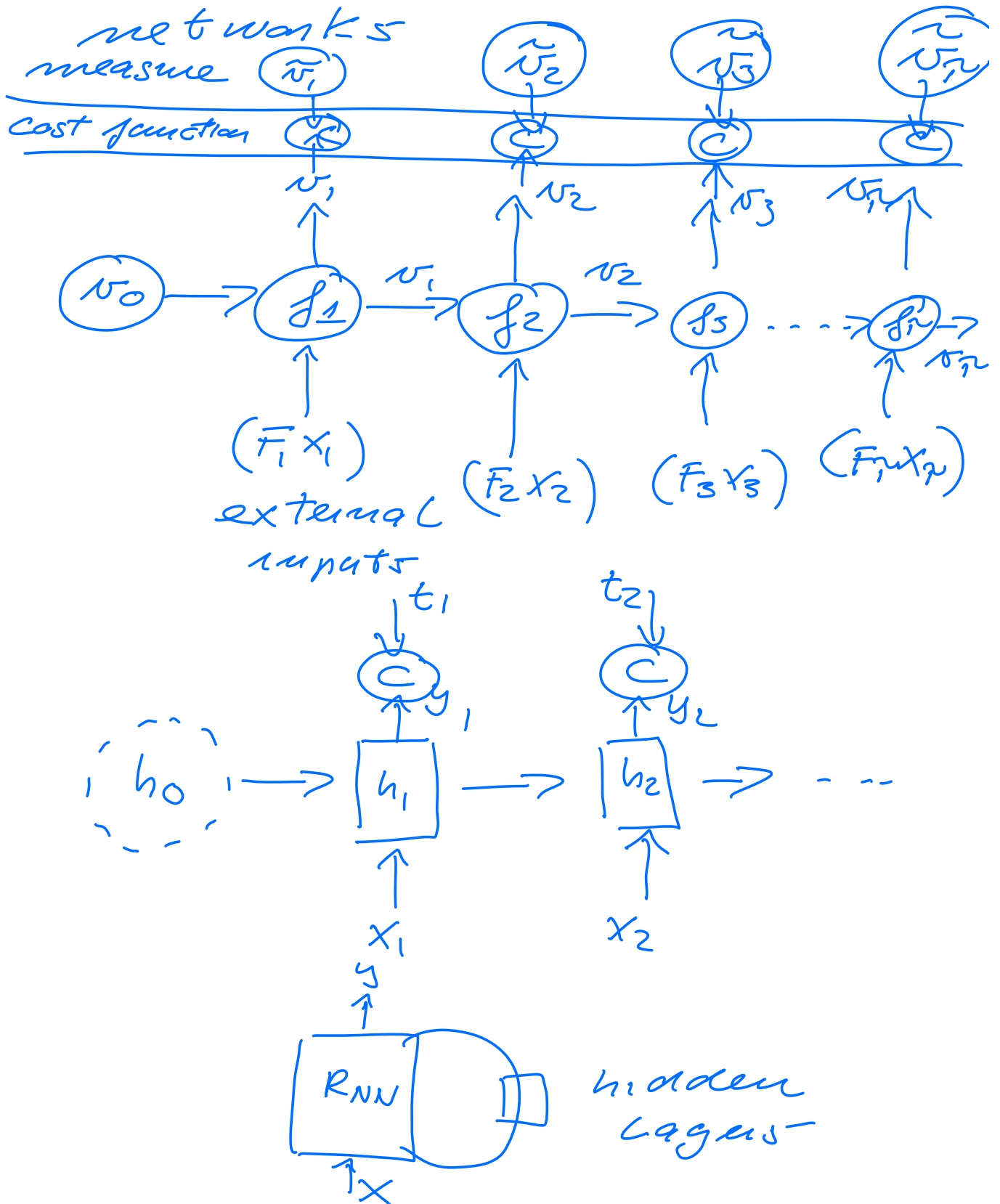
$$= v_i + f(v_i, \tilde{F}_i, x_i)$$

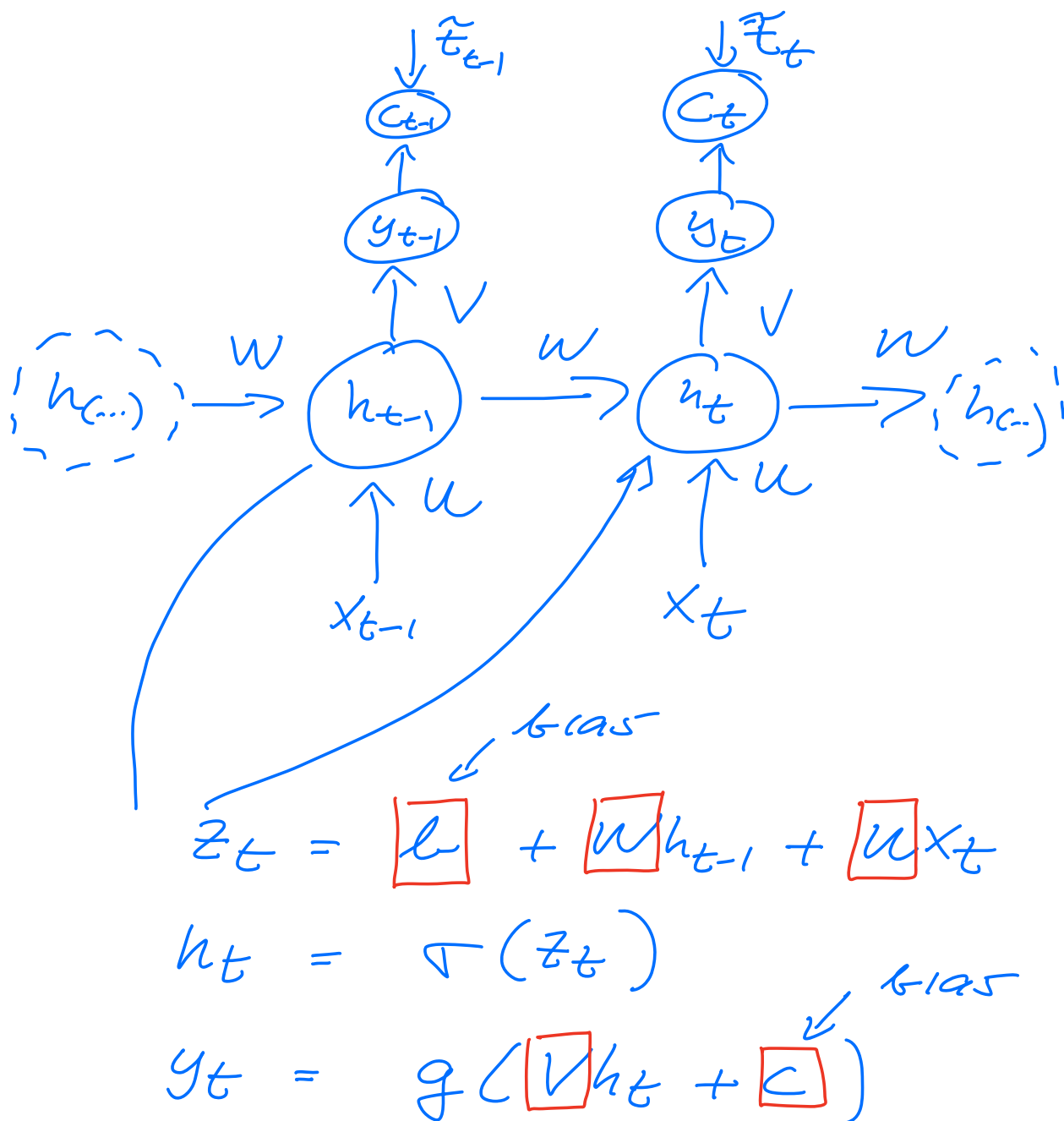
$$v_{i+1} = f(f(v_{i-1}, \tilde{F}_{i-1}, x_{i-1}), \tilde{F}_i, x_i)$$

Focus on  $v_i$



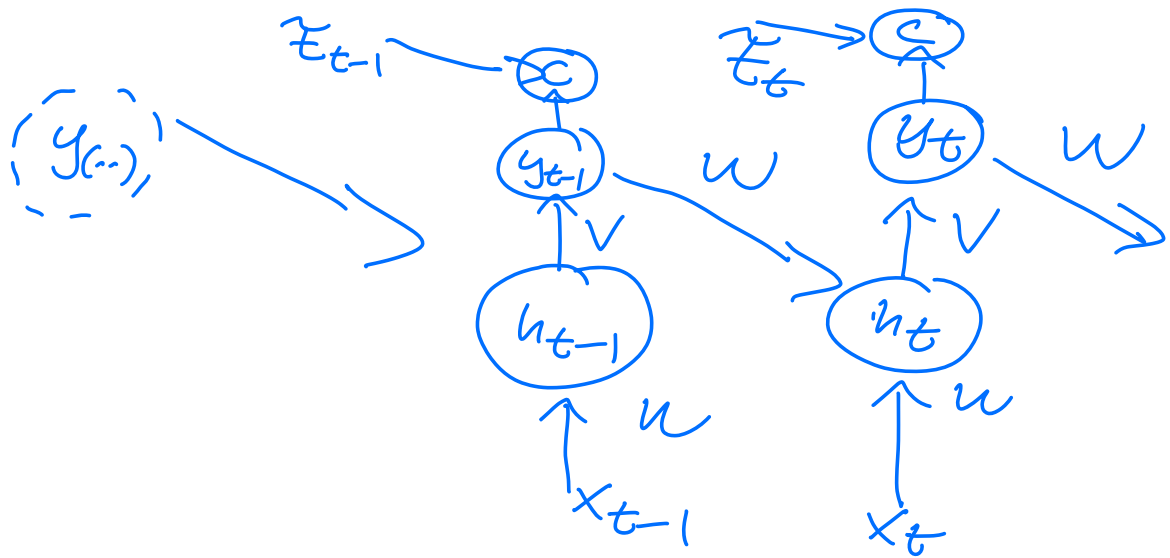
# Connection with a neural





BPTT = Back propagation through time,

Simplify the computational complexity



Problem with RNNs are often due to exploding gradients

$$h_t = W h_{t-1}$$

This operation is repeated for  $h_t$   $t$ -times

$$h_t = (W)^t h_0$$

$$W = S D S^T \quad S^T S = S S^T = \mathbb{I}$$

$$D = \begin{bmatrix} \lambda_0 & & 0 \\ & \ddots & \\ 0 & & \lambda_{d-1} \end{bmatrix} \quad Wx = \lambda x$$

eigenvalues  $\lambda_i$  and eigenvectors  $w_i$

$$h_0 = \sum_i \alpha_i w_i \quad \begin{array}{l} w_i w_i' \\ = \lambda_i w_i' \end{array}$$

$$W h_0 = h_1 = \sum_i \alpha_i \lambda_i w_i$$

repeat -  $t$  - times

$$(W)^t h_0 = h_t = \sum_i \alpha_i \lambda_i^t w_i$$

$$\lambda_0 > \lambda_1 > \lambda_2 \dots \lambda_{d-1}$$

$$h_t \approx \lambda_0^t w_0 \cdot \alpha_0$$

if  $\lambda_0 \geq 1$  we may get contributions to  $h_t$  which can be large  $\Rightarrow$  can give rise to exploding gradients

To avoid this, it is common  
to use gradient clipping

gradient  $\vec{g}$

if  $\|\vec{g}\|_2 \geq \text{specific value}$

$$\vec{g} \leftarrow \frac{\epsilon}{\|\vec{g}\|_2} \vec{g}$$

end if