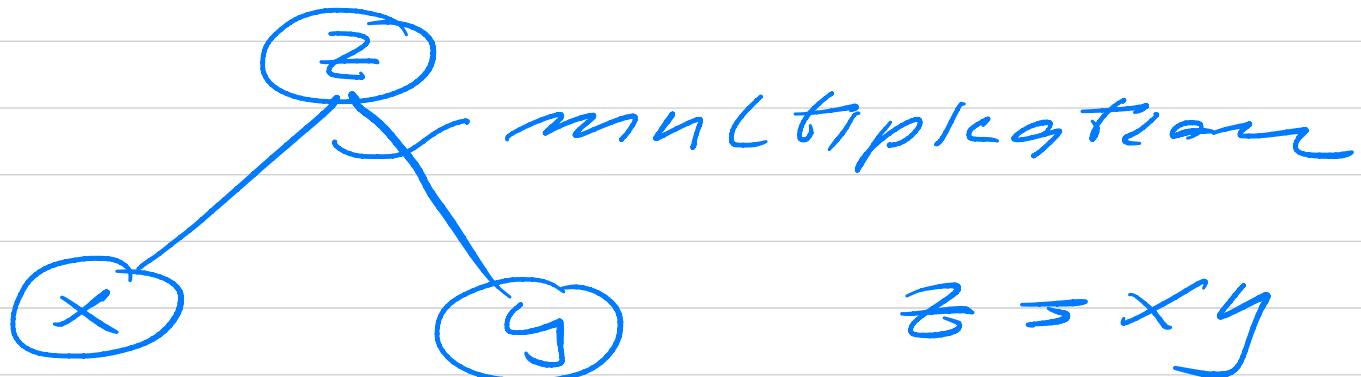


Lecture Fys- Stk3155/4155, November 9, 2023

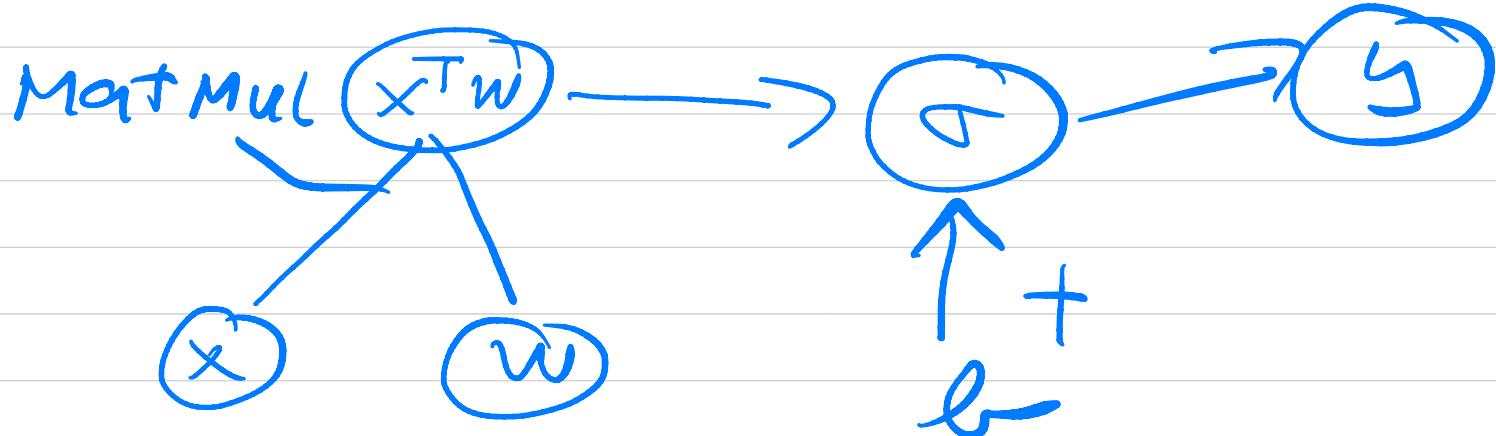
Recurrent NNs

Graphical representation



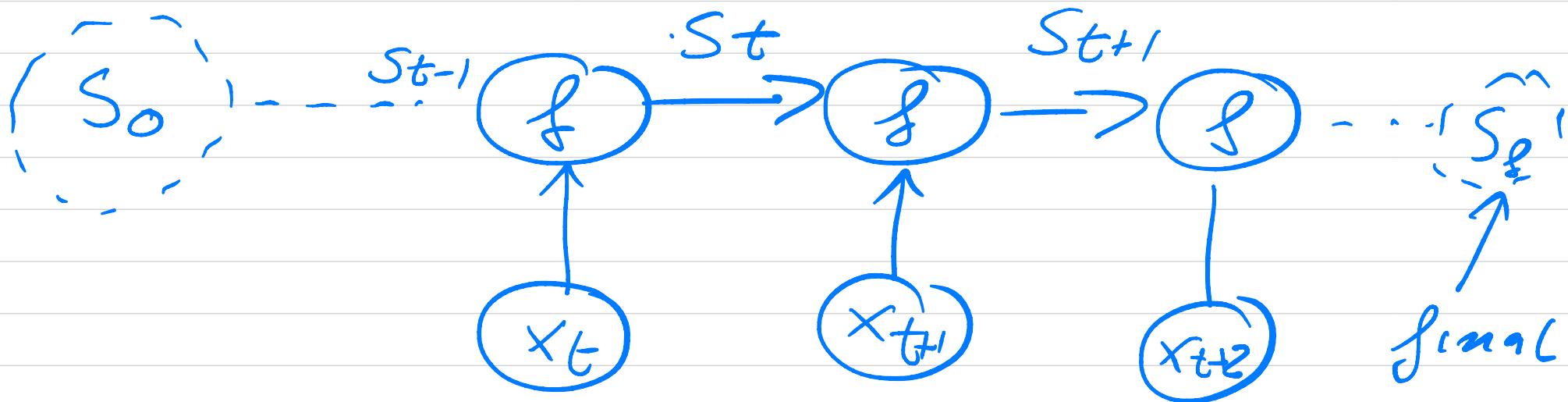
NN:

$$y = \sigma(x^T w + b) = \sigma(z)$$



Think of a dynamical system
driven by an external signal
at a time $-t-$

$$S_t = f(S_{t-1}, x_t; \theta)$$



Example

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

$$x_0 = x(t_0) \quad v_0 = v(t_0)$$

re write as two coupled
differential (Focus on v)

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

$$\begin{aligned}\frac{dv}{dt} &= -\frac{\gamma}{m} v - \frac{x}{m} + \frac{F(t)}{m} \\ &= g(F, x, v)\end{aligned}$$

$$x \rightarrow x_i \quad i=0, 1, 2, \dots, n$$

$$v \rightarrow v_i \quad \Delta t = \frac{t_f - t_0}{n}$$

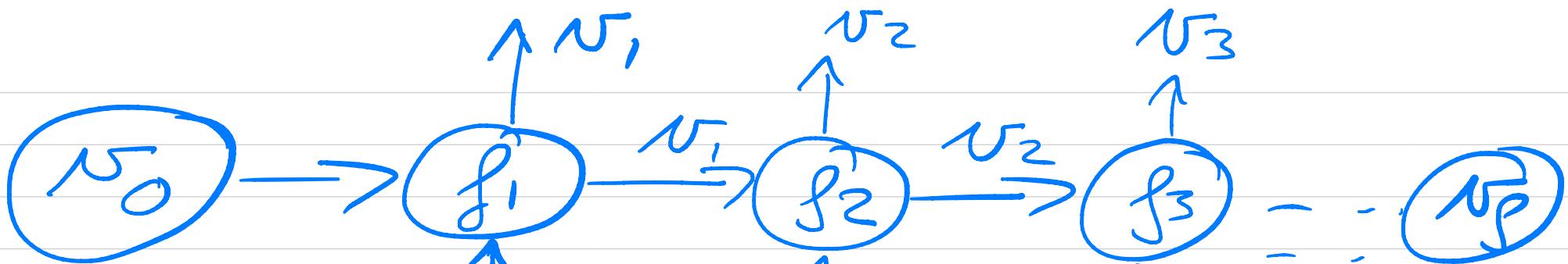
$$t_i = t_0 + i \cdot \Delta t$$

Euler's method

$$x_{i+1} = x_i + \Delta t \cdot v_i$$

$$v_{i+1} = v_i + \Delta t \cdot f_i$$

$$f_i = f(F_i, x_i, v_i)$$

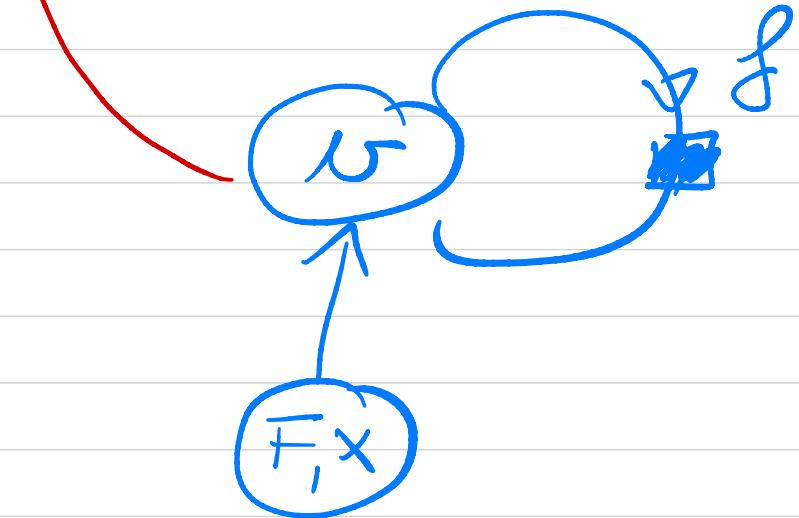


t_1

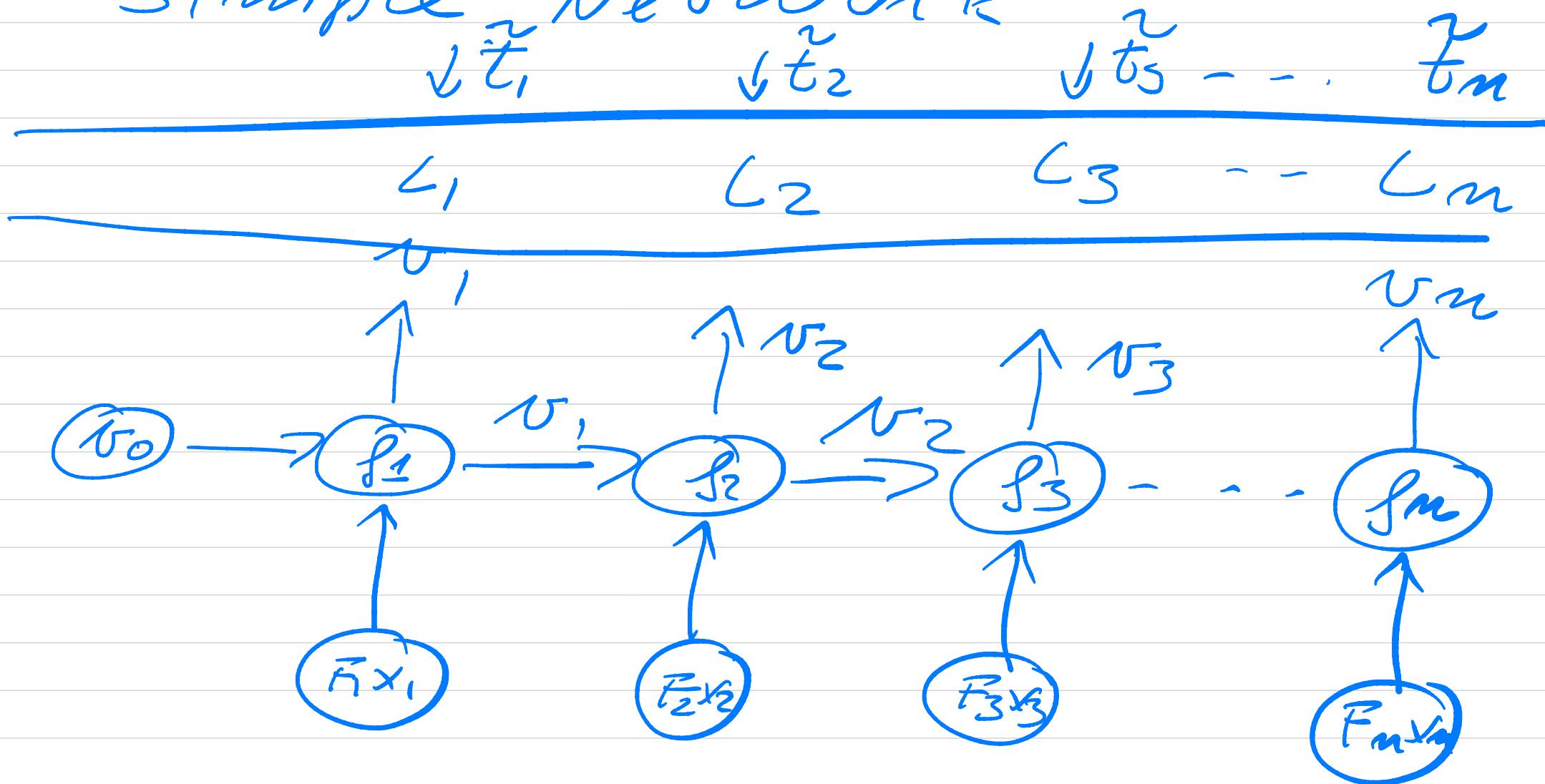
t_2

t_3

compact rewrite



Simple Network



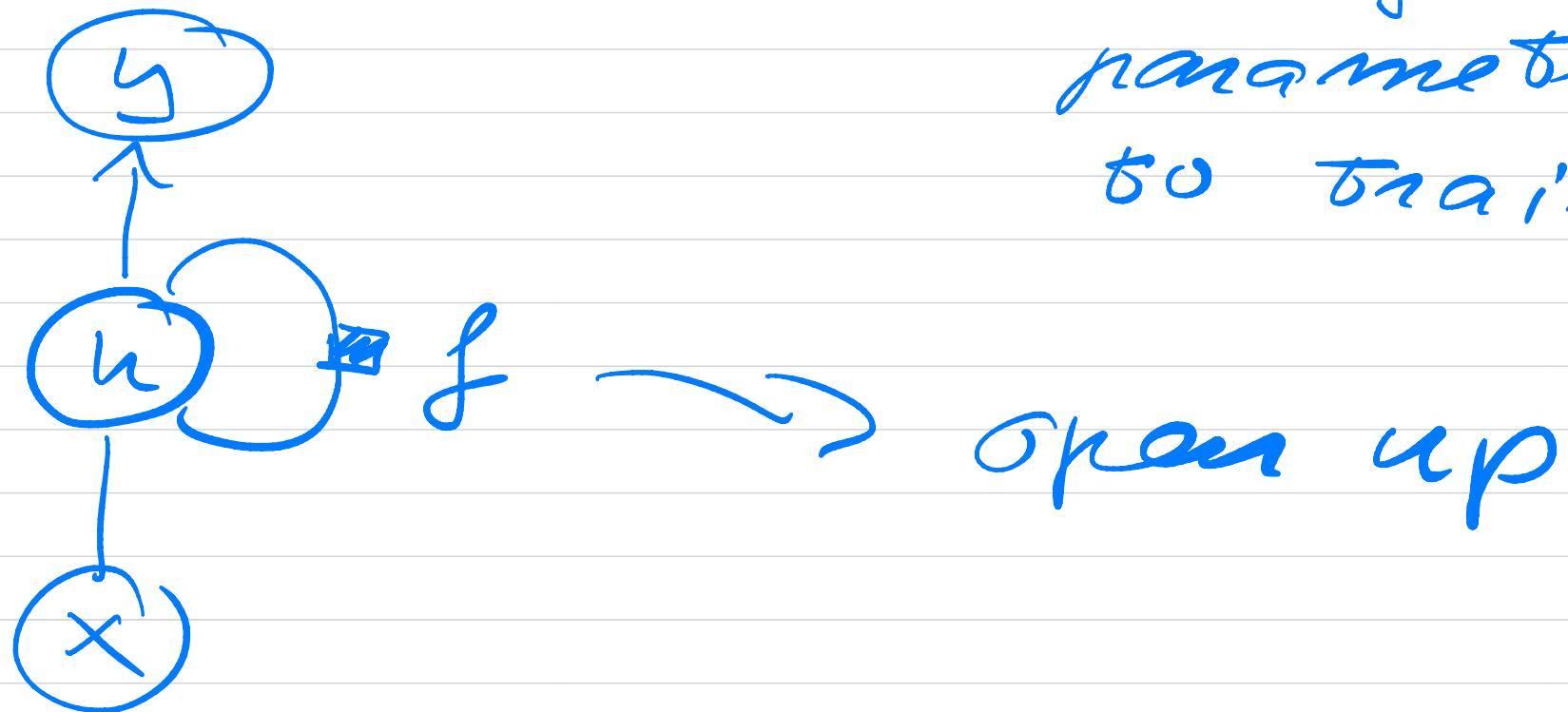
\tilde{t}_i = observation at time
 t_i \ (target, output)

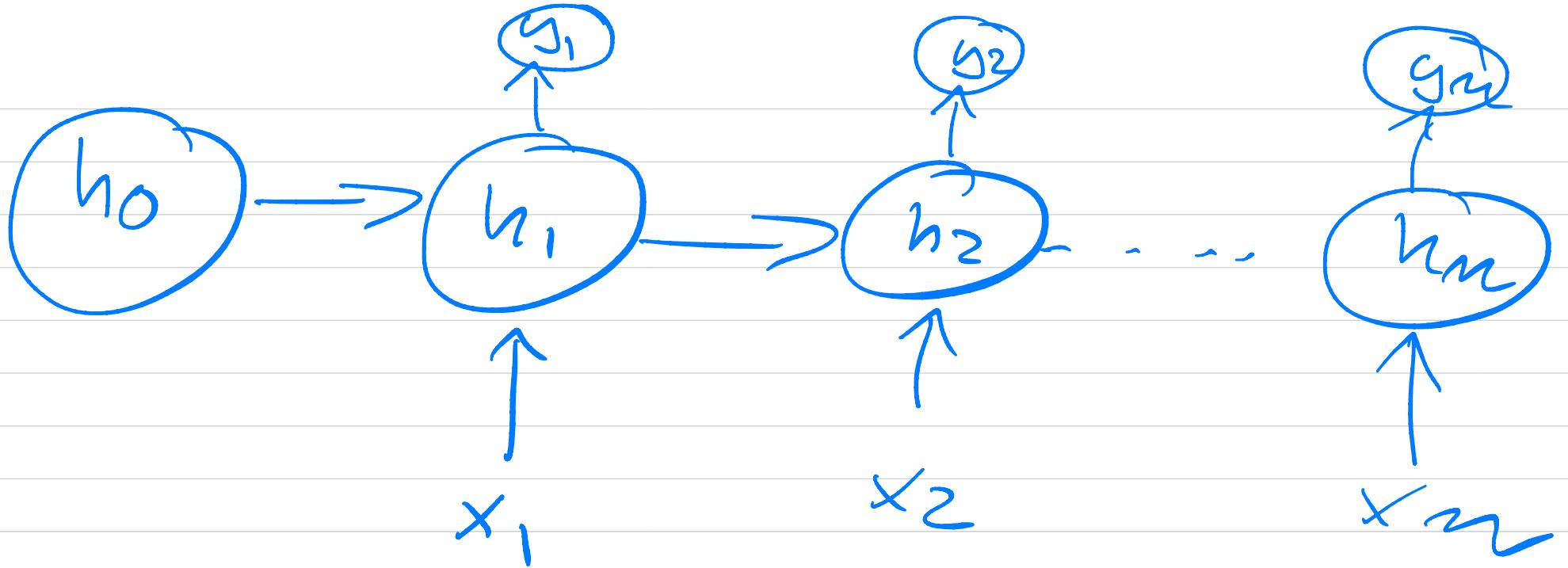
Now we replace the function

- f - with a neural net

$$h_t = f(h_{t-1}, x_t; \Theta)$$

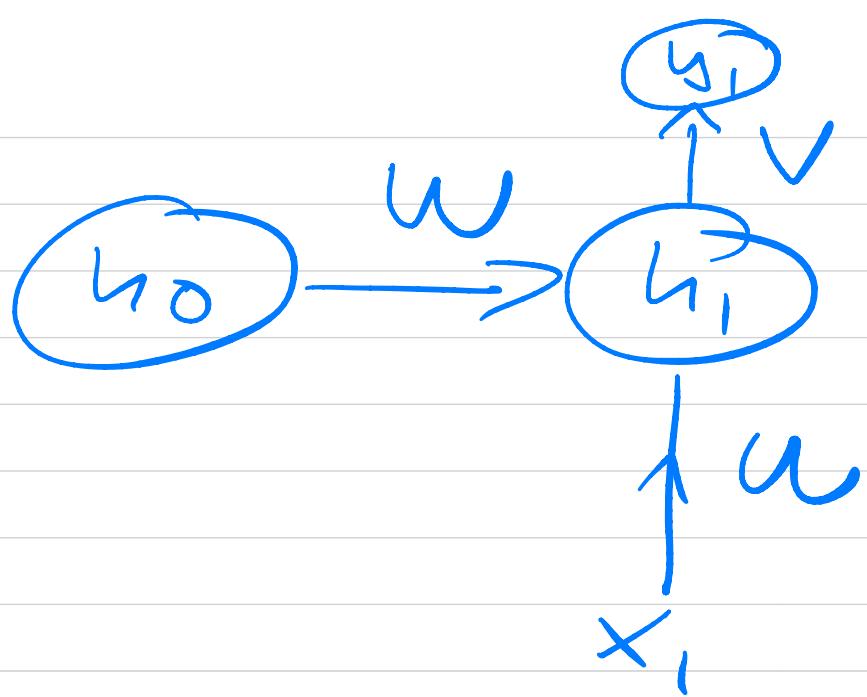
set of
parameters
to train





h_i refers to hidden part

$$\Theta = \{ u, w, v \}_{b, c}$$

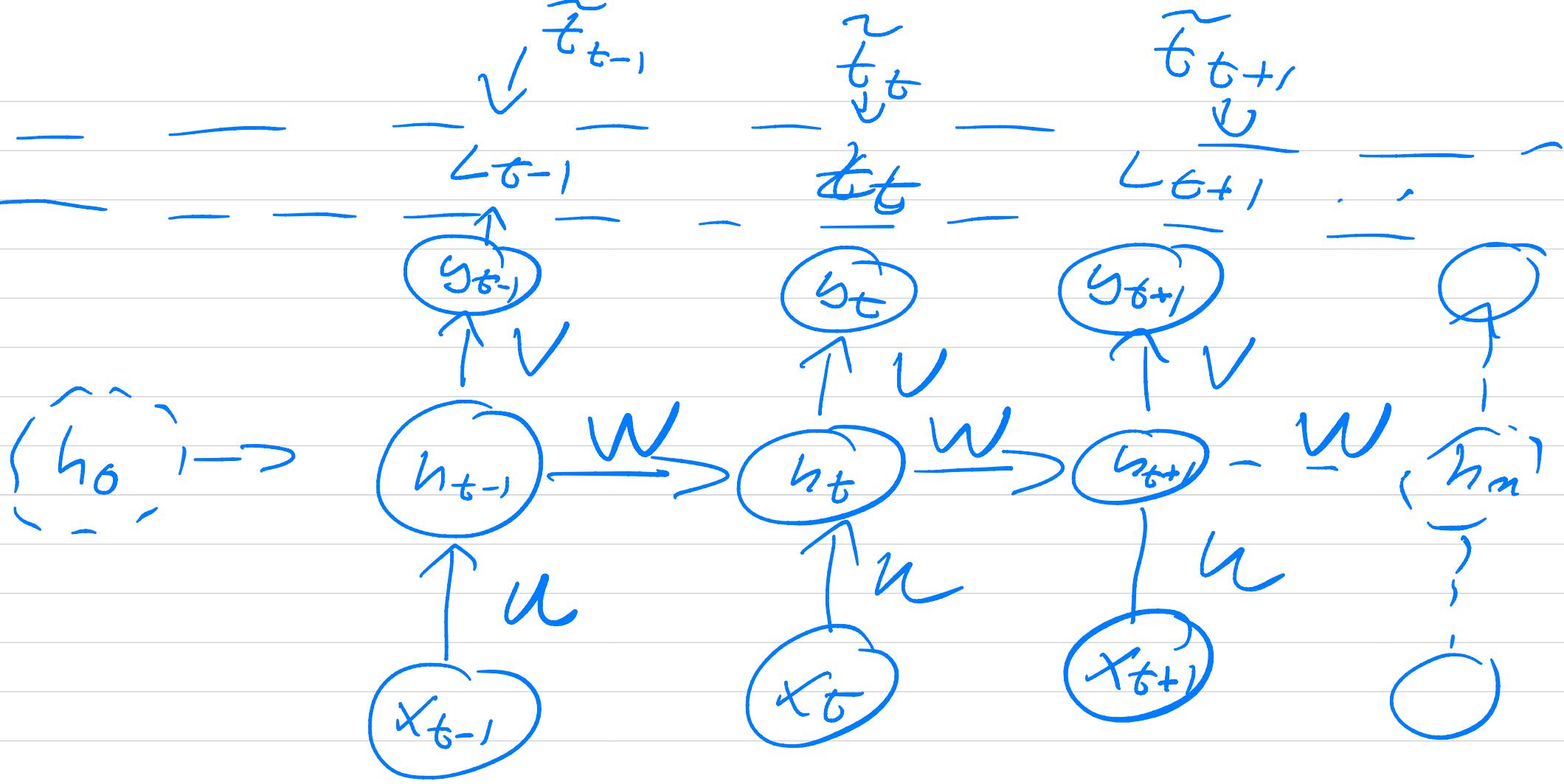


$$z_t = u x_t + w h_{t-1} + b \quad \leftarrow \text{bias}$$

$$h_t = f(z_t) \quad (\text{softmax tank}) \quad \text{or ReLU}$$

$$r_t = v h_t + c \quad \leftarrow \text{bias}$$

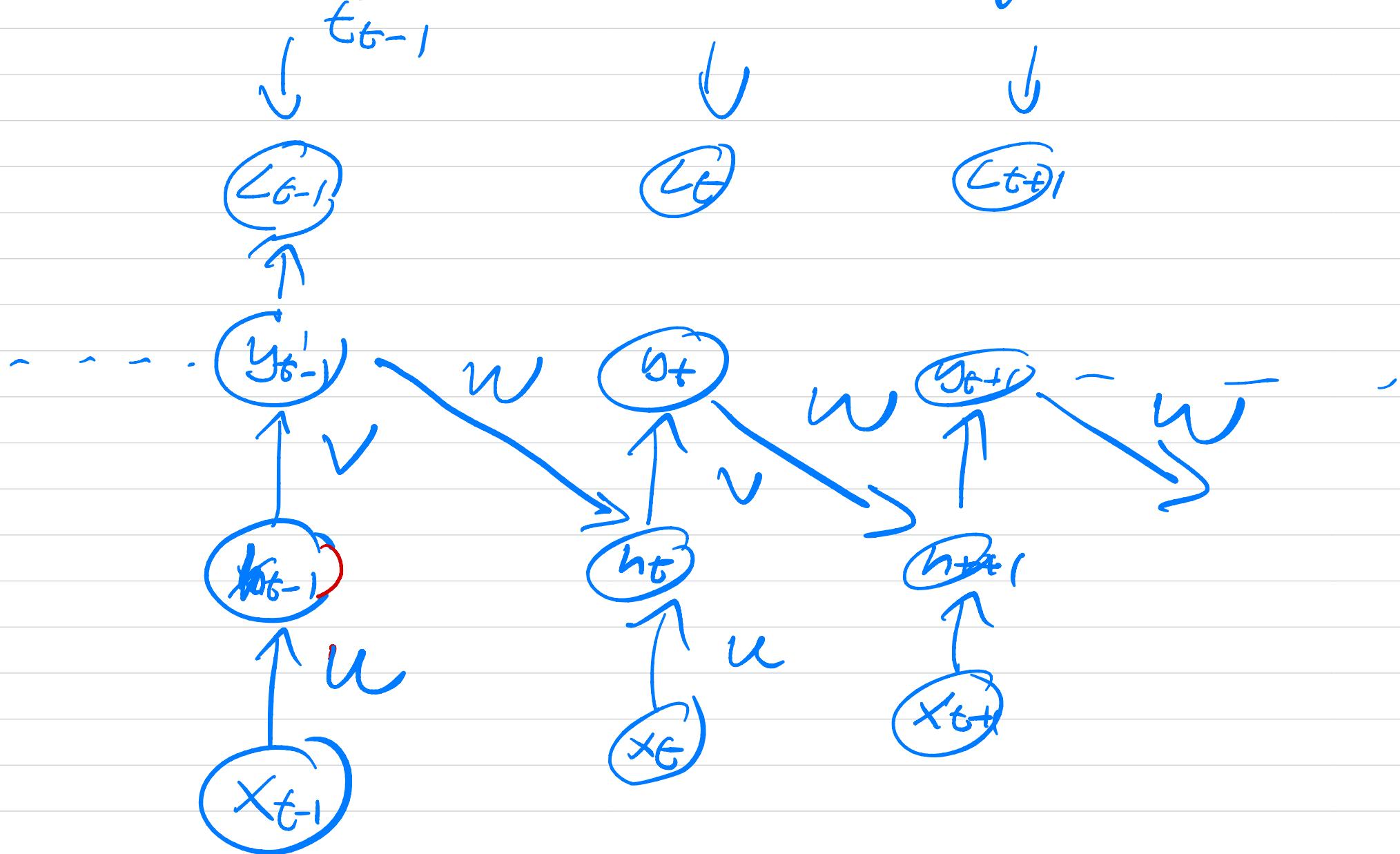
$$y_t = g(r_t)$$



$$L = \sum_{i=0}^m L_i$$

Back propagation
in time
(BPT) leads
to a sequential

One possible simplification



For each node (at time- t)

$\nabla_{\text{u}} \text{L}_t$, $\nabla_{\text{w}} \text{L}_t$, $\nabla_{\text{v}} \text{L}_t$

$\nabla_{\text{b}} \text{L}_t$, $\nabla_{\text{c}} \text{L}_t$

challenger isn't exploding
or vanishing gradients

$$\frac{\partial \text{L}_t}{\partial w} = \frac{\partial \text{L}_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \left[\sum_{k=1}^m \frac{\partial h_k}{\partial h_t} \frac{\partial \text{L}_t}{\partial w} \right]$$

Example

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

assume that w is the same for all previous $t -$

$$h_t = \underbrace{w \cdot w \cdot w \cdots}_{t\text{-times}} w h_0$$

$$= w^t h_0$$

suppose w is diagonalizable,
 $u^\top w u = D$

Expand h_0 in terms of eigenvectors w_i , eigenvalues λ_i

$$h_0 = \sum_i \alpha_i w_i \quad Ww_i = \lambda_i w_i$$

$$h_1 = Wh_0 = \sum_i \alpha_i w_i w_i'$$

$$= \sum_i \lambda_i \alpha_i w_i'$$

Repeat t-times

$$h_t = W^t h_0 = \sum_i \alpha_i \lambda_i^t w_i'$$

$$\lambda_0 > \lambda_1 > \lambda_2 \dots > \lambda_n$$

$$w_t \simeq \lambda_0 w_0$$

if $\lambda_0 > 1$, then gradients can explode,

if $\lambda_0 < 1$, can lead to vanishing gradients.

Gradient clipping can be used to avoid exploding gradients.

gradient \vec{g}

if $\|\vec{g}\|_2 \geq \epsilon$

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

end if.