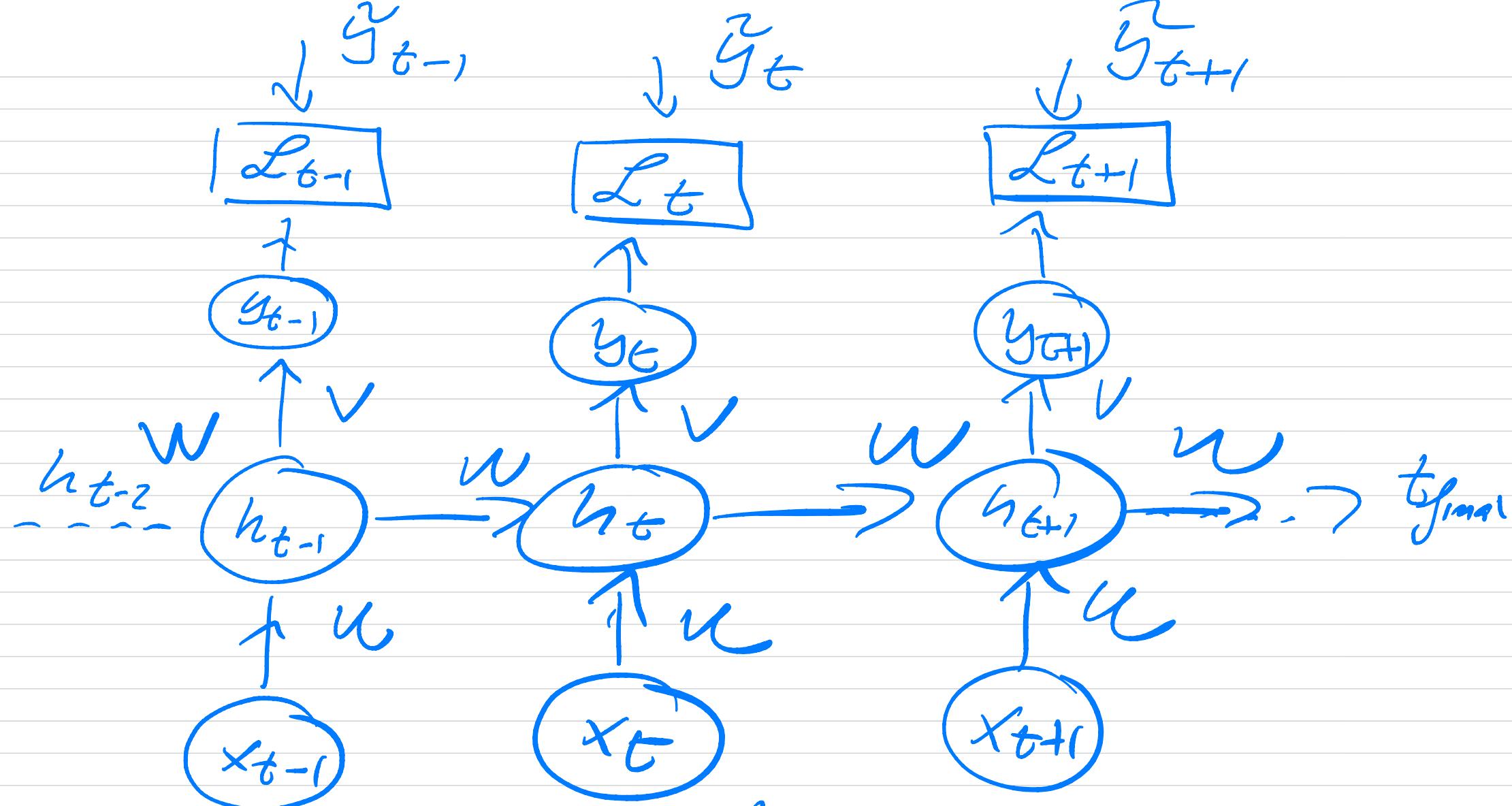


**Lecture FYS5429,
March 5, 2024**

FYS 5429, MARCH 5, 2029

Set of parameters

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



$$\mathcal{L}(\epsilon) = \sum_{t=0}^{t_{final}} \mathcal{L}_t$$

$$z_t = u x_t + w h_{t-1} + b_z$$
$$h_t = \sigma_h(z_t)$$

$$z_t = V h_t + c_t \leftarrow \text{bias}$$

$$y_t = \sigma_y(z_t)$$

$$L_t(y_t, \hat{y}_t)$$

Common strategies

- weight sharing
- truncation of number of steps where different weights are training
- Training is done iteratively
 - Feed forward pass
(randomly assigned weights & biases)
 - Back propagation in time (BPT)

BPT

$$L = \sum_t L_t$$

$$\nabla_c L$$

$$\nabla_b L, \nabla_w L, \nabla_v L$$

$$\nabla_u L$$

can give rise to exploding or vanishing gradients.

$$\frac{\partial L_t}{\partial w} = \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t}$$

$$\times \left\{ \sum_{k=1}^n \left[\frac{\partial h_t}{\partial h_k} \right] \frac{\partial h_k}{\partial w} \right\}$$



$$\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^n \frac{\partial h_i}{\partial h_{i-1}}$$

can lead to exploding or vanishing gradients

assume same w for all steps

$$h_1 = w^T h_0$$

$$h_i = w^T h_0 = w^T h_{i-1}$$

$$= w \cdot w^T h_{i-2}$$

assume we can diagonalize w

$$w = u D u^T \quad u u^T = u^T u = 1$$

eigenpairs are (λ_i, w_i)

$$h_0 = \sum_i \lambda_i w_i$$

$$w^{t_{ho}} = \sum_i \alpha_i w_i'$$

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

$$w^{t_{ho}} = \sum_i \lambda_i^t \alpha_i w_i'$$

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

$$\lambda_0 \geq 1$$
$$h^t = w^{t_{ho}} \succeq \lambda_0 \alpha_0 w_0$$

$\lambda_0 \geq 1 \Rightarrow$ exploding gradients

$\lambda_0 = 1$, ideal case

$\lambda_0 < 1 \Rightarrow$ vanishing gradient

exploding gradient,
gradient clipping

gradient \vec{g}

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

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

endif