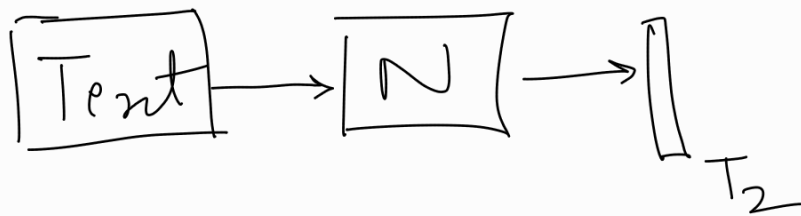
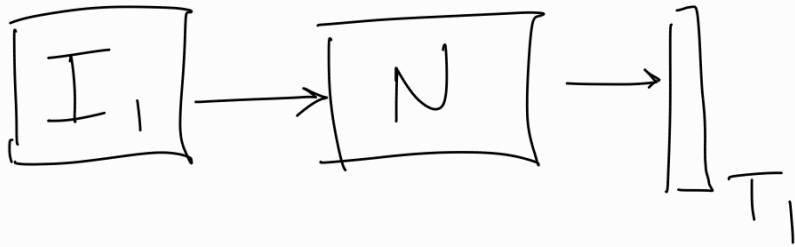


12/3/24

# Multi-Modal Self Supervised Learning

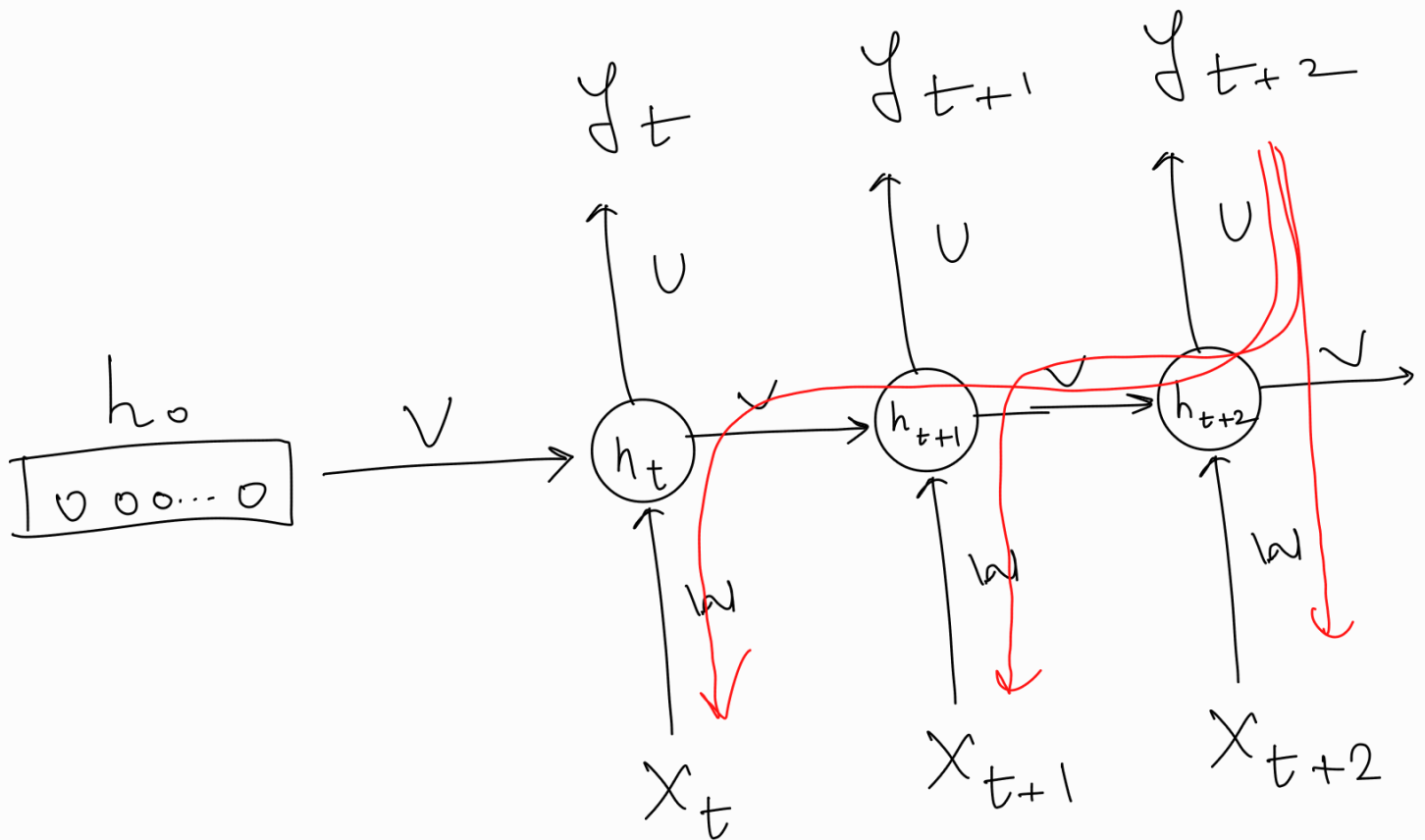


→ Uses: image captioning.

→ Search using keywords to narrow down search space. Now you can apply KNN, etc.

# RNNs

Keep content of previous sequence of characters/words.



BPTT: Back Prop  
Through Time

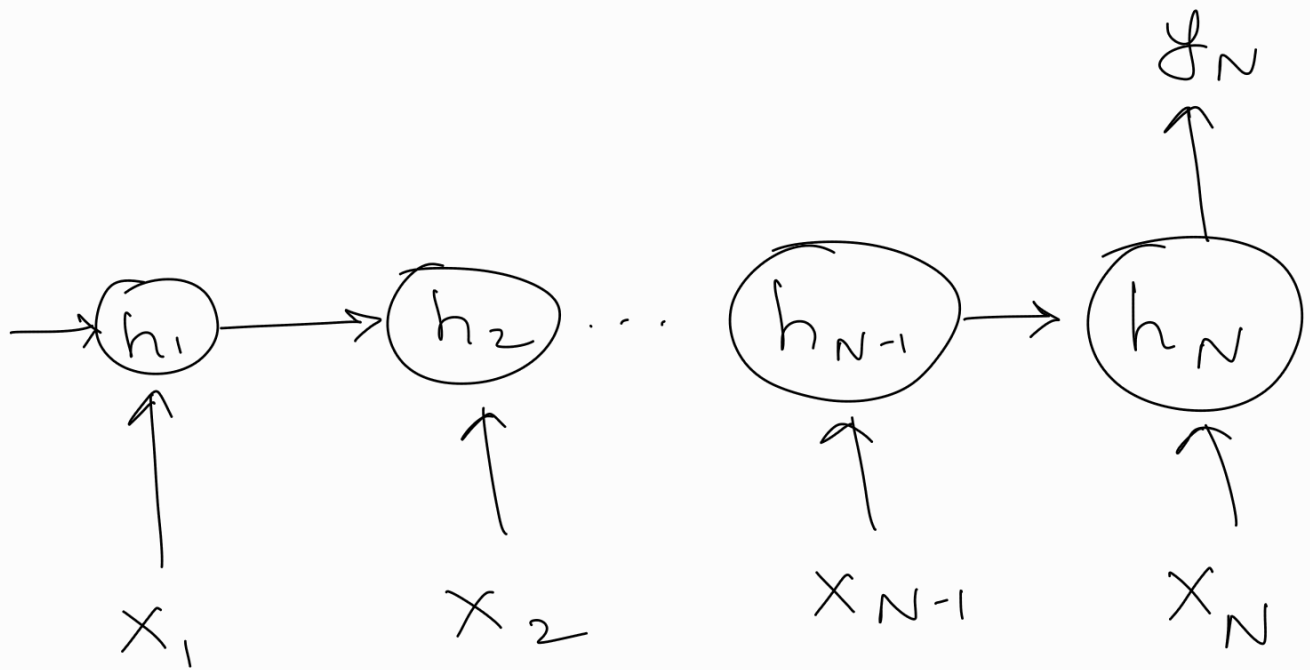
$$h_t = \phi(w^T x_t + v^T h_{t-1})$$

$$y_t = \phi(u^T h_t)$$

Vanishing  
Gradient problem.

The deeper you  
go, the lower  
the effect of  $\partial L$   
on earlier inputs.

# Sentiment Analysis Network



Single output

---

How to embed image?

Pass through CNN, take second-last layer. Use that as 0<sup>th</sup> context to the RNN (used for image captioning)

