

# CS725

So far, all neural network families worked with fixed-length inputs and fixed-length outputs

E.g. Image classification. Input  $x \in \mathbb{R}^{h \times n}$  or  $x' \in \mathbb{R}^d$

Need to support variable-length inputs and outputs

# Examples of variable length input/output tasks

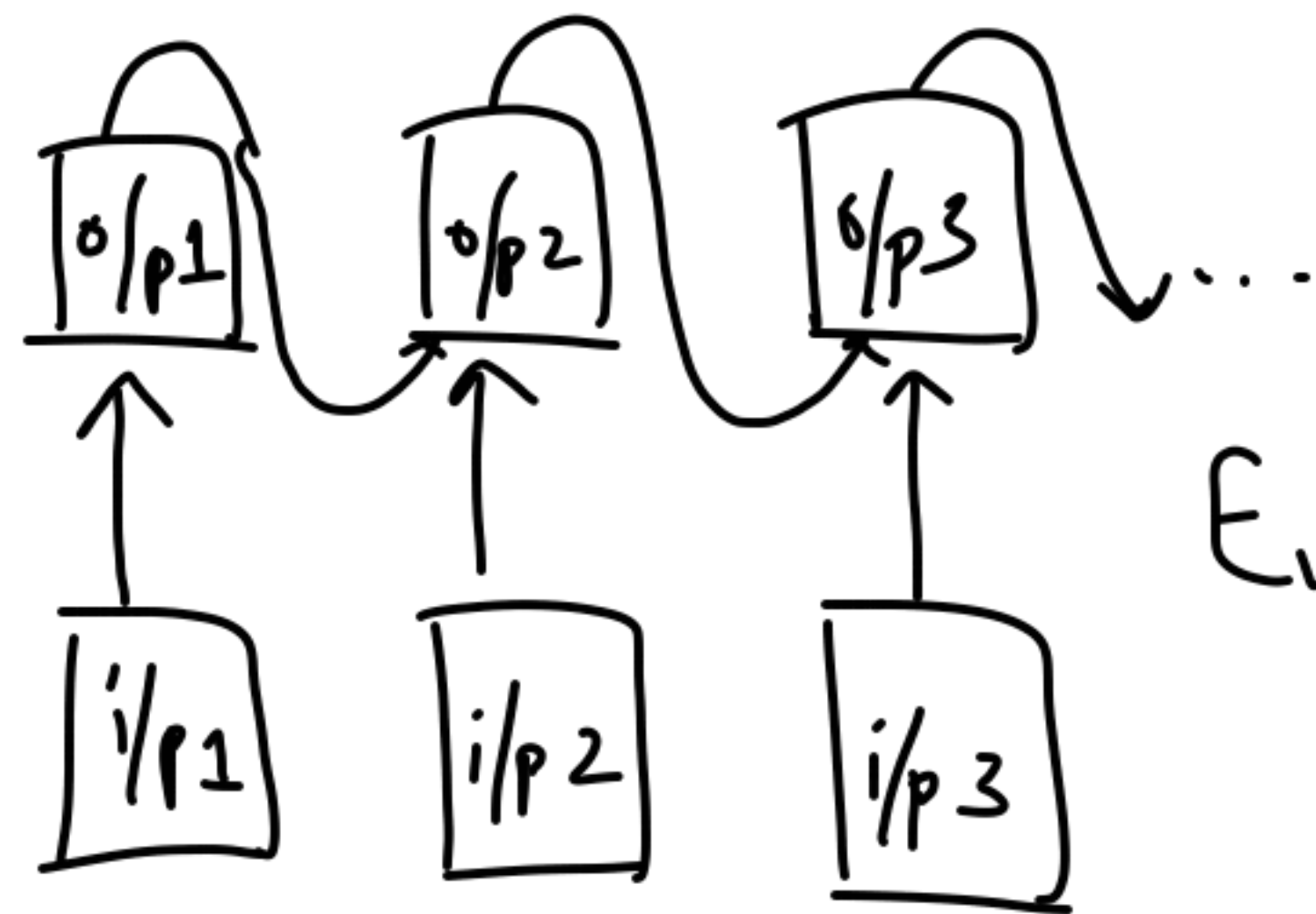
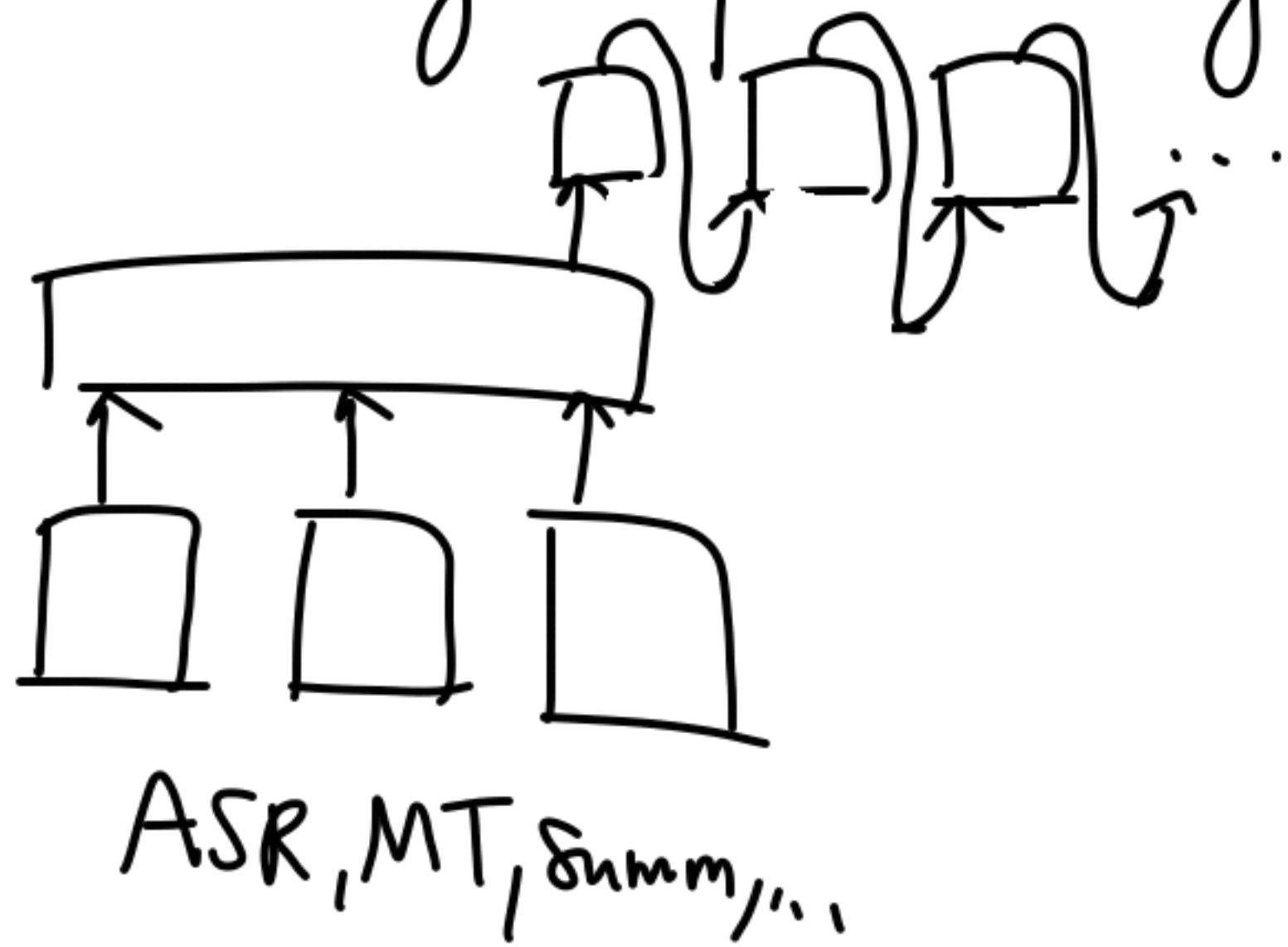
① Fixed-length i/p & Variable-length o/p  
E.g., Image captioning



② Variable-length i/p & fixed length o/p  
E.g., Image generation, Speech emotion classification, Sentiment analysis, etc  
or text

③ Var-length i/p & var-length outputs

E.g. Speech recognition<sup>(ASR)</sup>, machine translation<sup>(MT)</sup>, Summarization, etc



E.g. POS-tagging

Sequence labeling tasks where length of the i/p & o/p seqs are the same

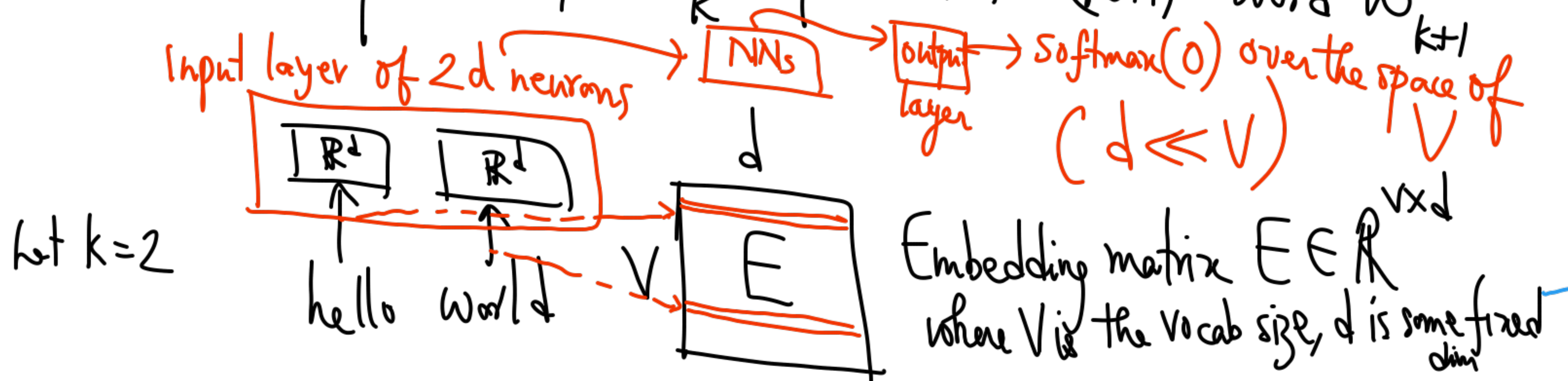
# LANGUAGE MODELING task

What is language modeling? Given a sequence of words  $w_1, \dots, w_{T-1}$ , what is the most likely next word?

$$w^* = \operatorname{argmax}_w P(w \mid w_1, \dots, w_{T-1})$$

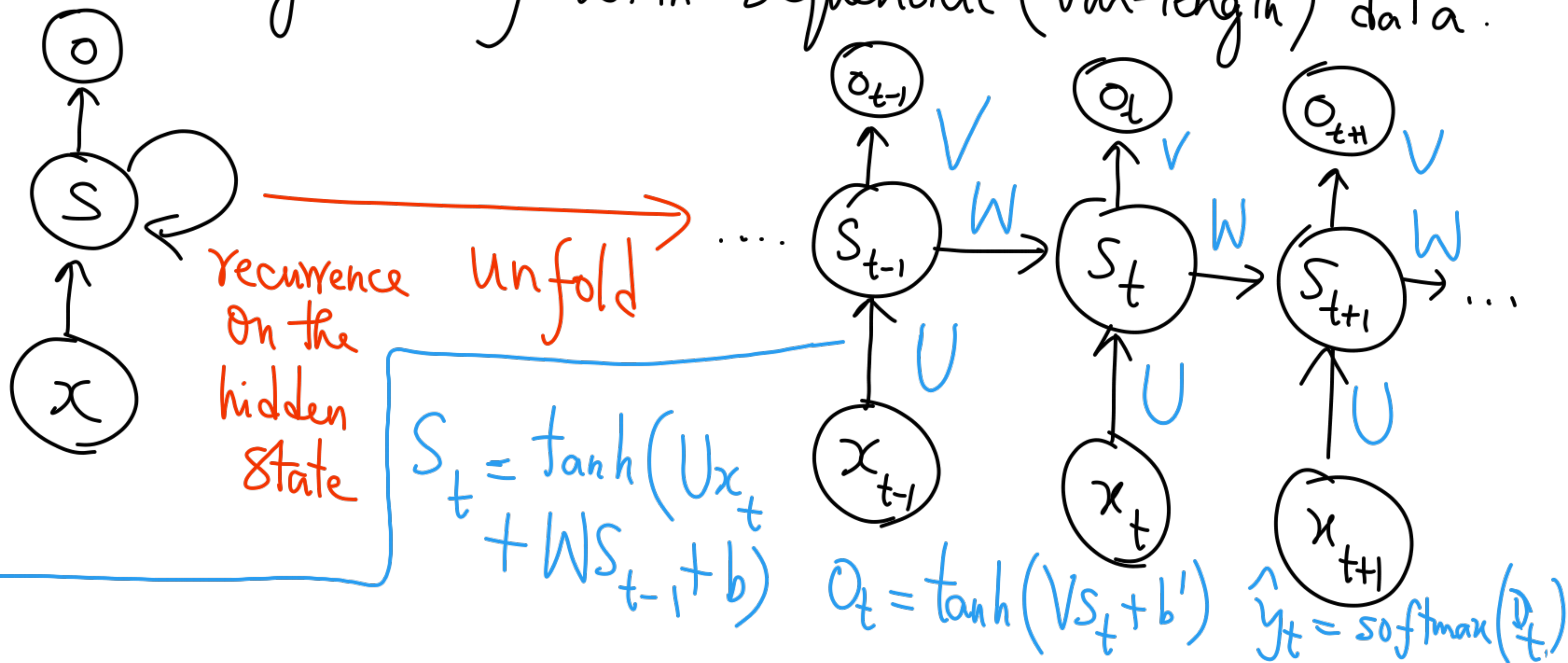
How can we use feedforward neural networks for language modeling?

fix a window size, say  $k$  and feed  $k$  words at a time as input  $w_1, \dots, w_k$  to predict the  $(k+1)^{\text{th}}$  word  $w_{k+1}$



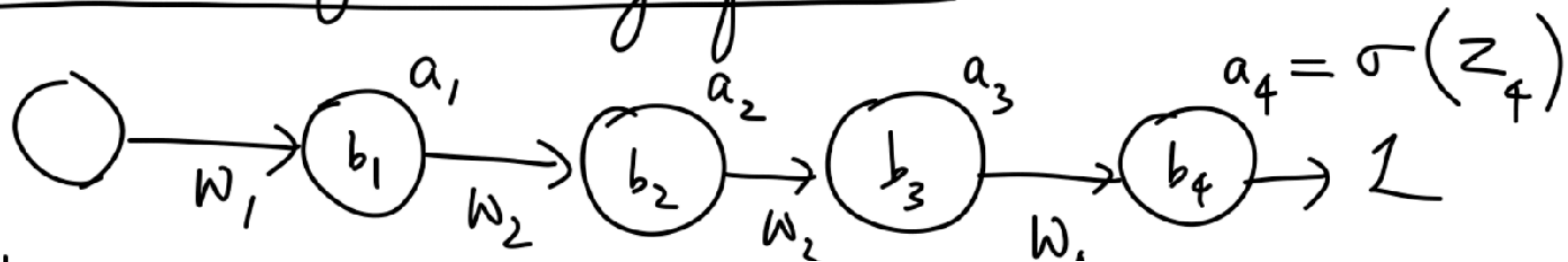


Unlike FFNs, recurrent neural networks (RNNs) deal organically with sequential (var-length) data.



A challenge of RNNs  $\Rightarrow$  Vanishing or exploding gradients

Illustration of vanishing gradient:



$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial a_4} \cdot \frac{\partial a_4}{\partial z_4} \cdot \frac{\partial z_4}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$= \left[ \frac{\partial L}{\partial a_4} \cdot \sigma'(z_4) \cdot w_4 \cdot \sigma'(z_3) \cdot w_3 \cdot \sigma'(z_2) \cdot w_2 \cdot \sigma'(z_1) \right]$$

Gradients become smaller in early layers (if wts are typically  $< 1$ )



Vanishing gradients when we consider gradients from early layers and the wts are initialized to values  $< 1$

Exploding gradients can happen when you start with large weights



Clipping threshold to keep exploding gradients in check!