



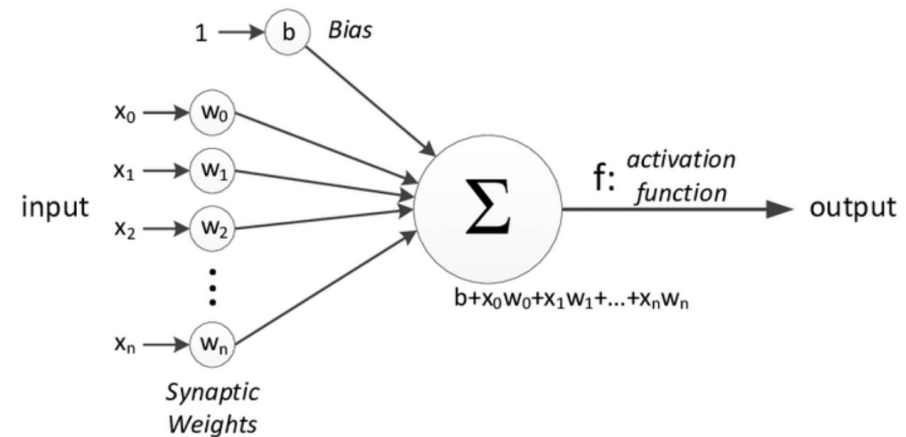
May 5th, 2026
Ay 119

Advanced Deep Learning: A Brief View to Transformers

**Matthew J. Graham
and Pavlos Protopapas**

Deep learning: recap from last seminar

- The basic neuron: $y = f(\mathbf{w} \cdot \mathbf{x} + b)$
- Choice of activation function
- MLP or feedforward NN trained by backpropagation
- Loss functions for different purposes
- Regularization to prevent overfitting
- Optimizers
- Different types of neural architecture:
 - Convolutional Neural Networks
 - Recurrent Neural Networks
 - Autoencoders
 - GANs



Sequential data is everywhere



[p.004] What clashes here of wills gen wonts, oystrygods gaggin fishygods! Brékkek Kékkek Kékkek Kékkek! Kóax Kóax Kóax! Ualu Ualu Ualu! Quauouauh! Where the Baddelaries partisans are still out to mathmaster Malachus Micgranes and the Verdons catapelting the camibalistics out of the Whoyteboyce of Hoodie Head. Assiegates and boomerangstroms. Sod's brood, be me fear! Sanglorians, save! Arms apeal with larms, appalling. Killykillkilly: a toll, a toll. What chance cuddleys, what cashels aired and ventilated! What bidimetoloves sinduced by what tegotetabsolvers! What true feeling for their's hayair with what strawng voice of false jiccup! O here here how



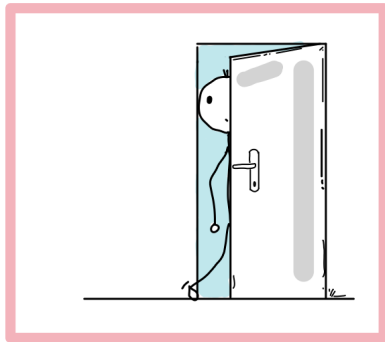
```

GAGCTCAGATACANCTTTGCGTGGNTTGACATACTGCITAGAGGNI GTAGSSGCACTGGGACCACACATCAGCTCATATTGACCCTGCTTAATCAATTTT 100
      |
      | EcoRI
      |
GCTTGGICAGTCAACAATGAATTCTTACTTAGATTATACGATTTATAACCGTGGTCTAACACTTATAGTTCGAAAGTGGGIGITTTCTCTGAGACCAA 200
      |
      | M N S Y L D Y T I Y N R G S N T I Y S S K V G C F P V E Q
      |
GAATACTTGCCAAGTGCITIGTSCAADTACAAATAGTTACATTCGGAGGGACGGCCAGTCGGAGGAAACACTTTCACATCGGCACCGCAGAAACTCATG 300
      |
      | E Y L P S A C A S T N S Y I P E G R P V G G N T F T S A P H E T H
      |
      | H18(72) →
    
```



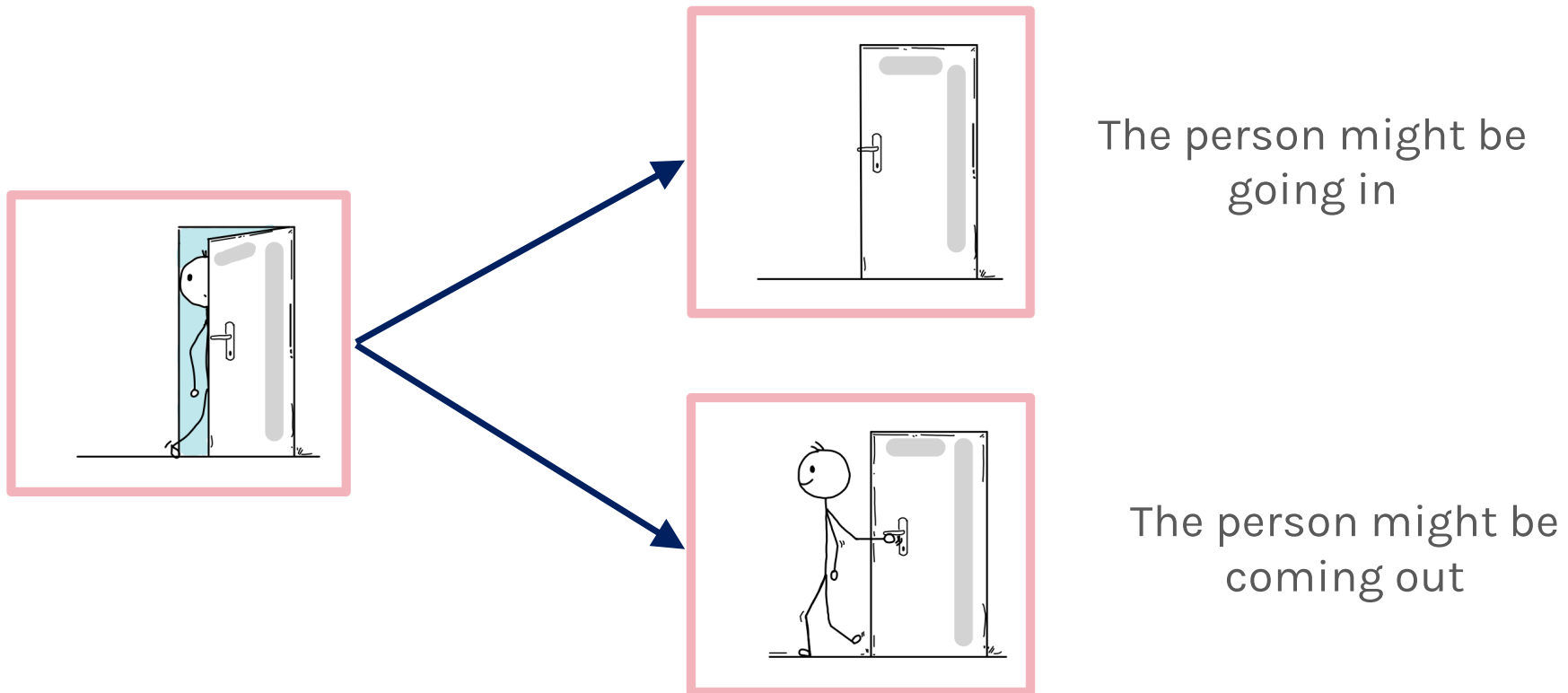
Let's think about sequences

Given this frame, what do you think is the next likely frame?



Let's think about sequences

Given this frame, what do you think is the next likely frame?

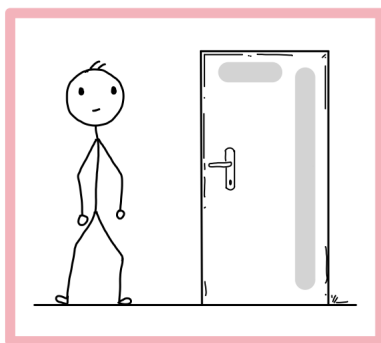


Based on only the first frame, it would be difficult to predict the next frame.

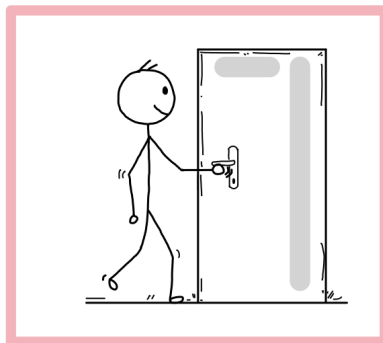
Let's think about sequences

However, if we give the model the previous frames, it is easy to predict the next:

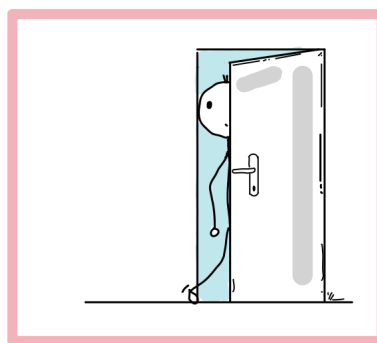
Frame 1



Frame 2



Frame 3



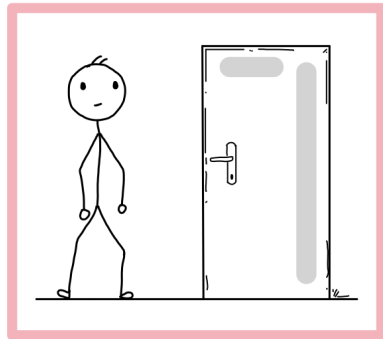
Frame 4



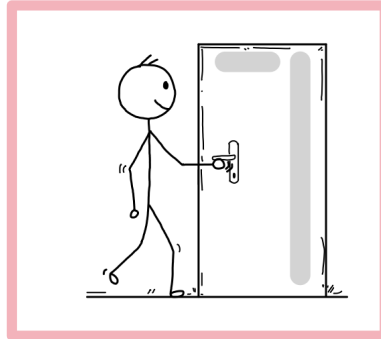
Let's think about sequences

However, if we give the model the previous frames, it is easy to predict the next:

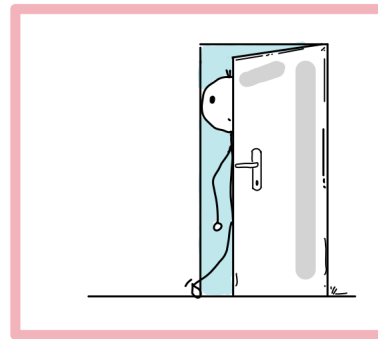
Frame 1



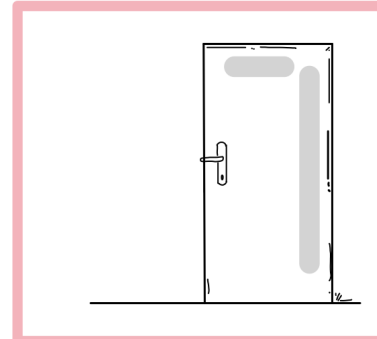
Frame 2



Frame 3



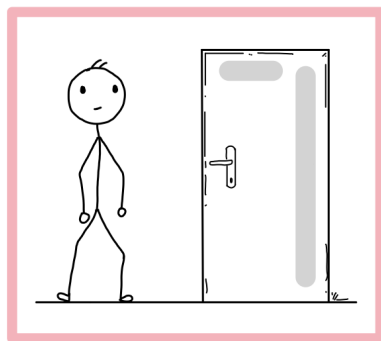
Frame 4



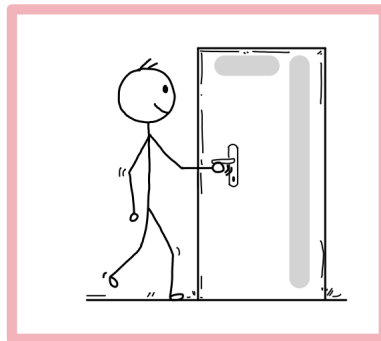
Let's think about sequences

However, if we give the model the previous frames, it is easy to predict the next:

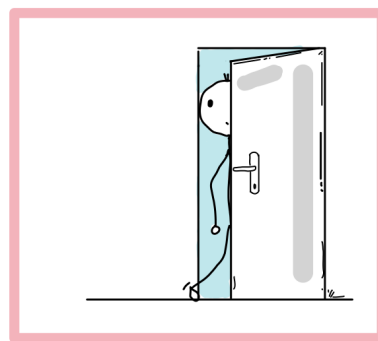
Frame 1



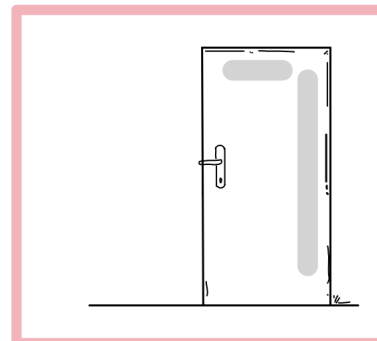
Frame 2



Frame 3



Frame 4

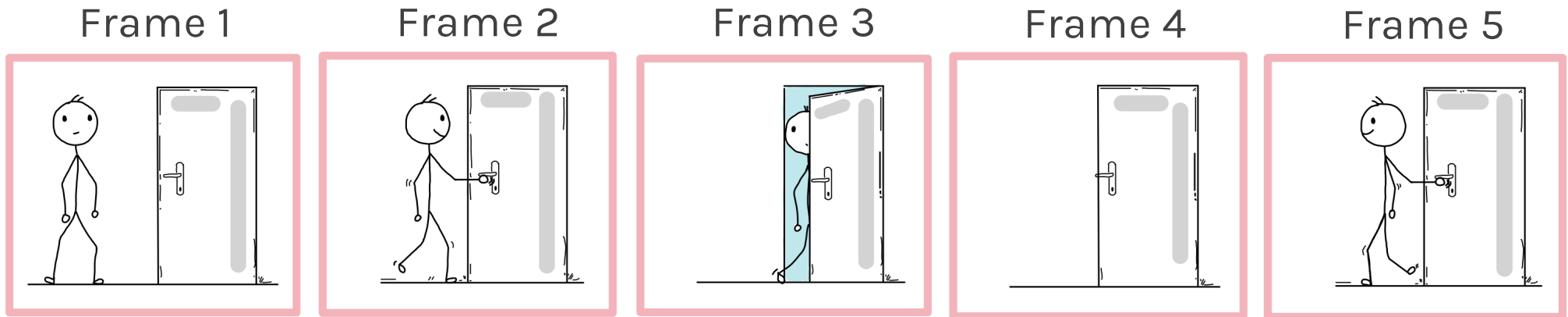


Frame 5



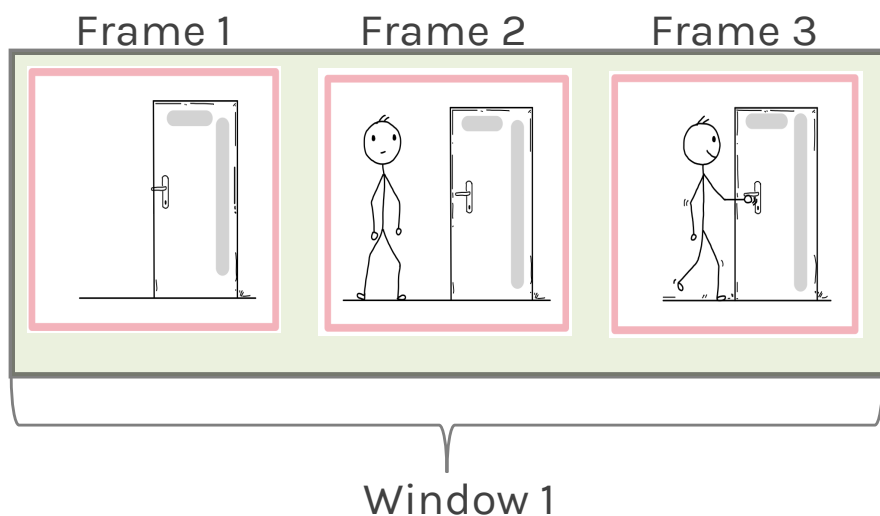
Let's think about sequences

However, if we give the model the previous frames, it is easy to predict the next:



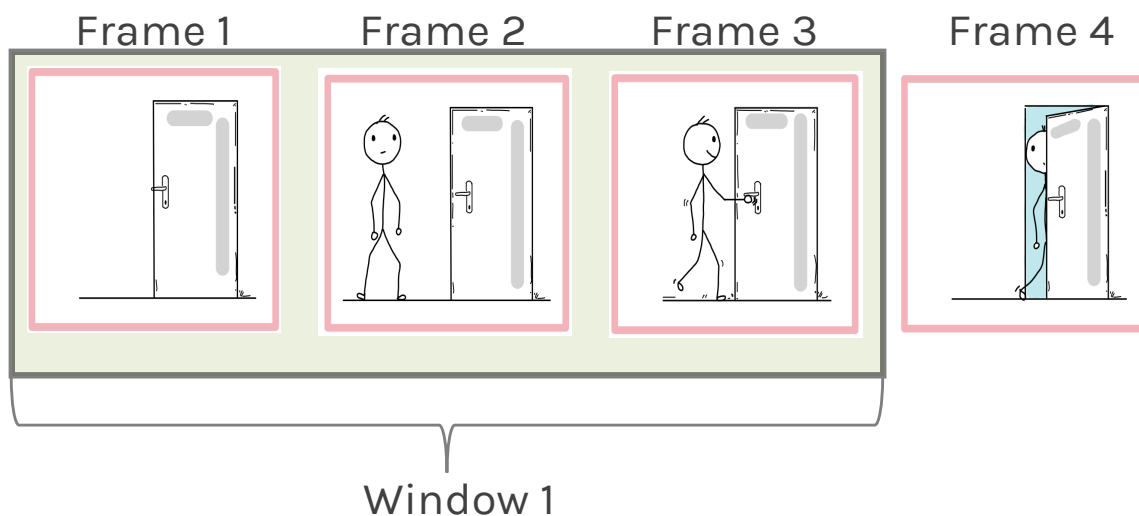
Sequences play an important role for forecasting and predictions.

Let's think about sequences



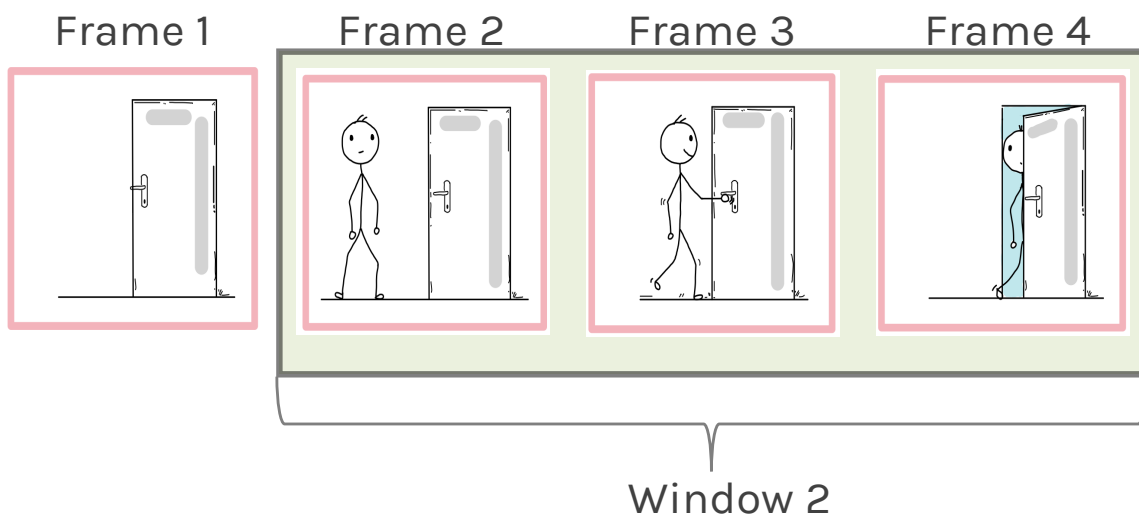
If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



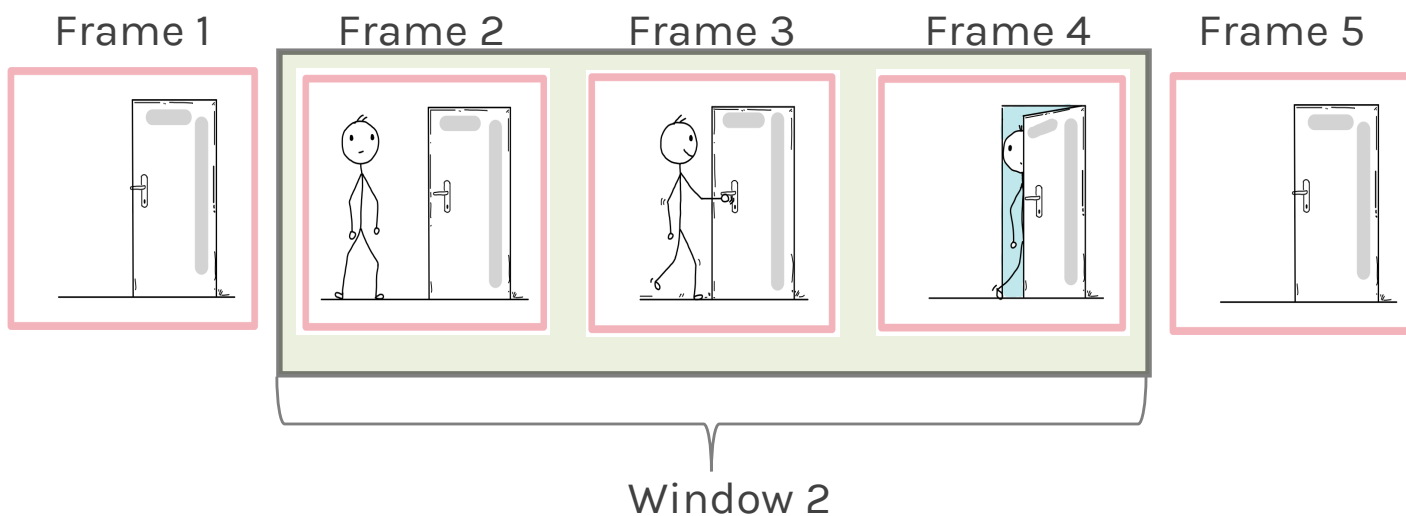
If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



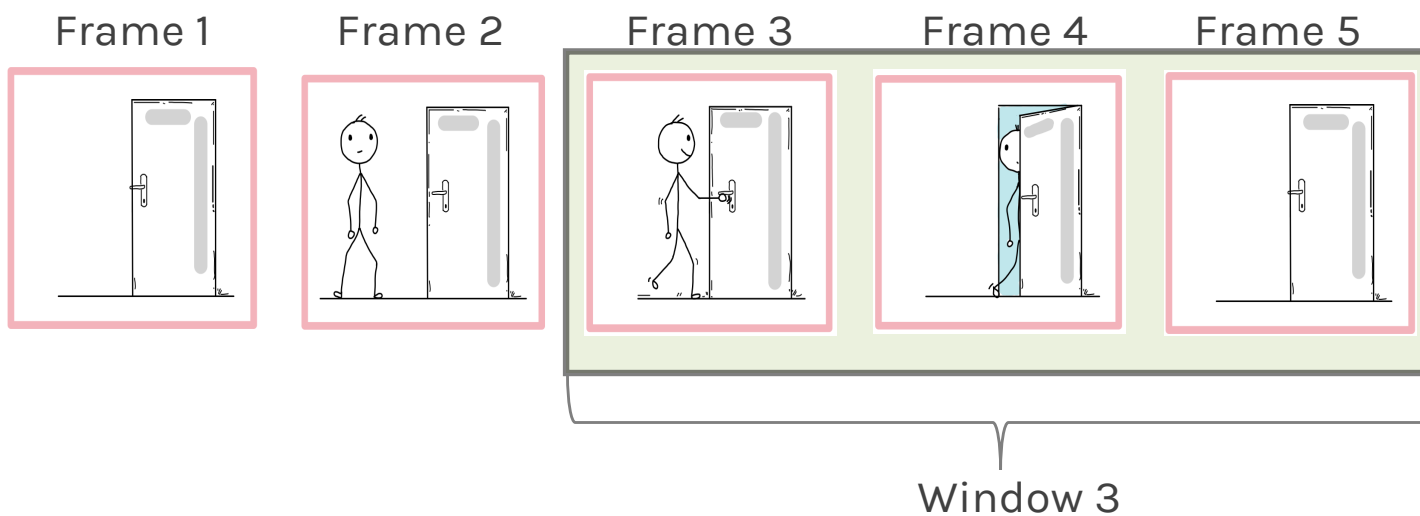
If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



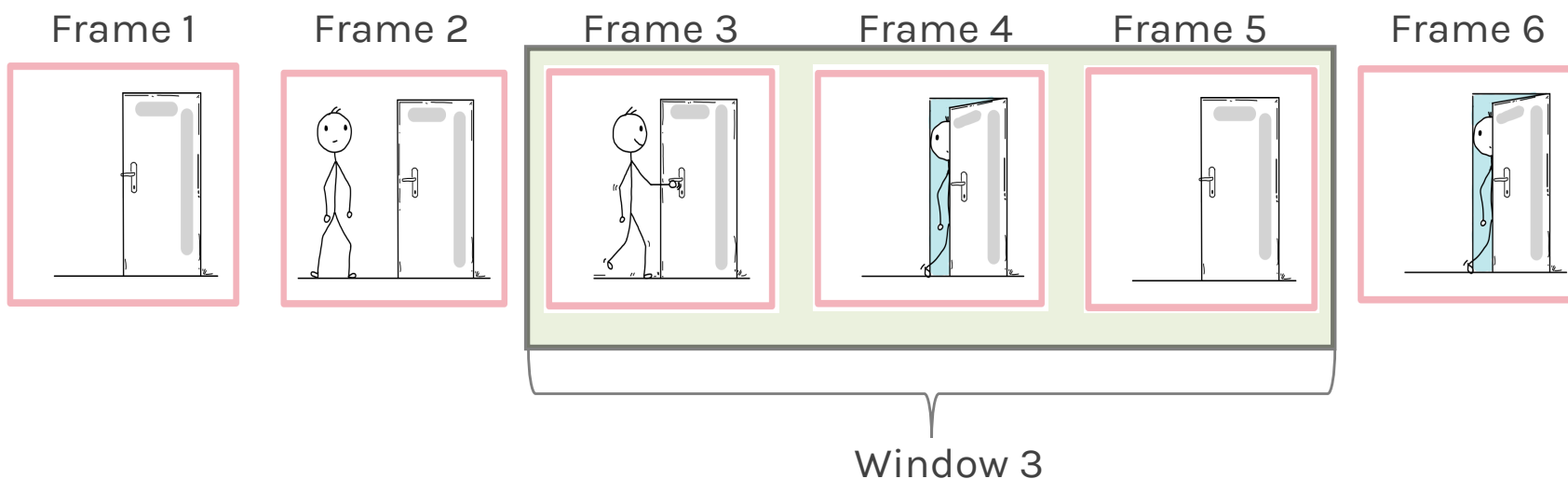
If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



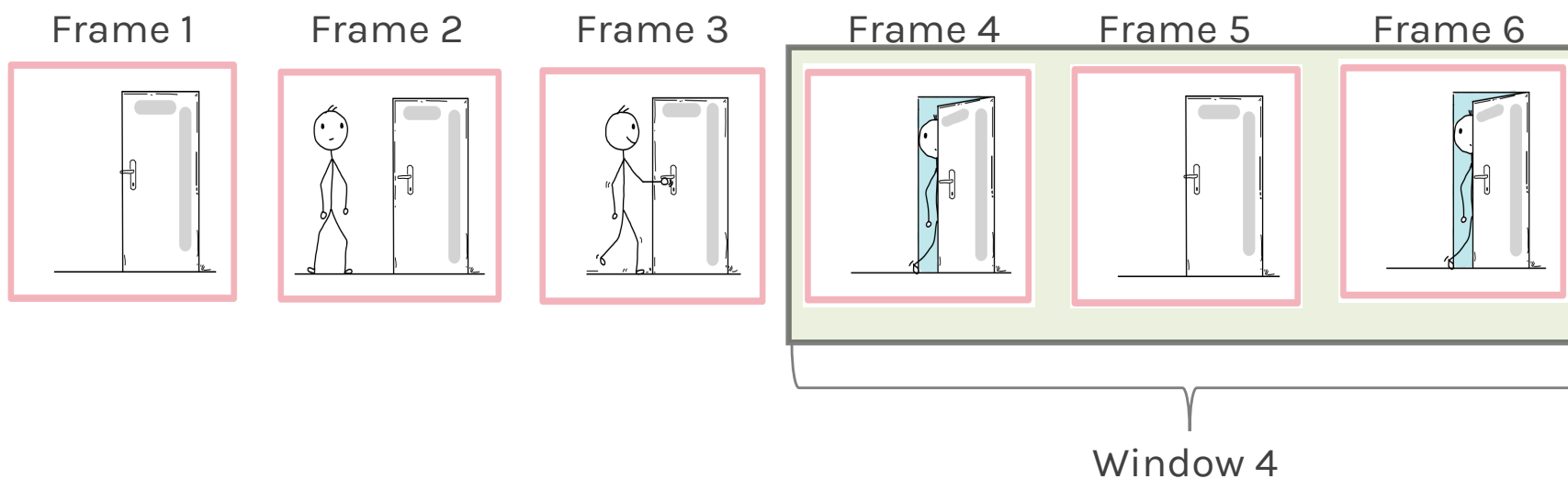
If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



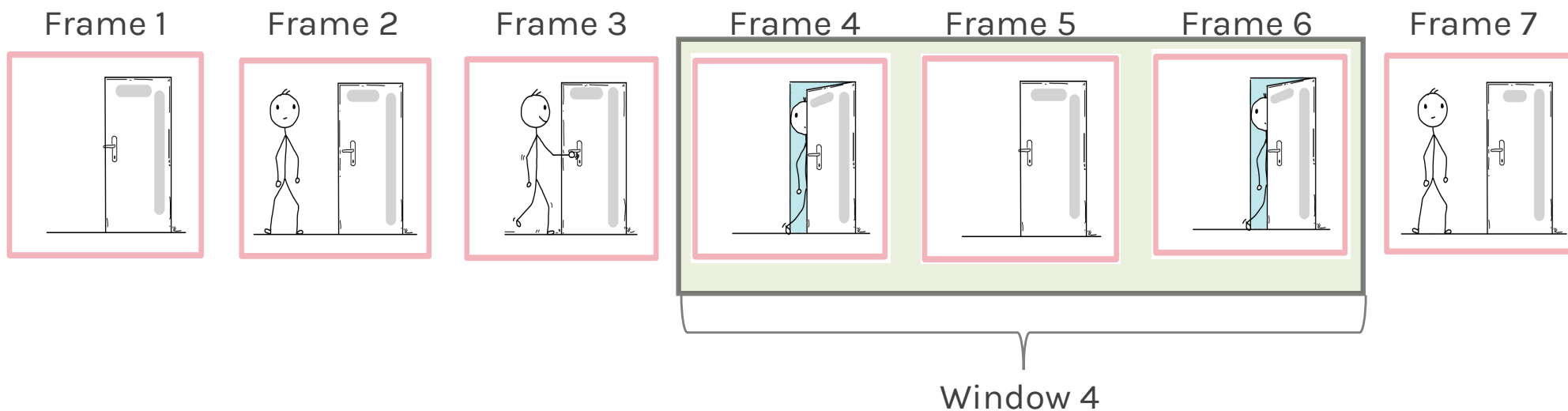
If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences



If we window a fixed number of frames as input to a NN like MLP or CNN then prediction will work

Let's think about sequences

However, consider the following sequence of frames:



Let's think about sequences

What is the next frame?



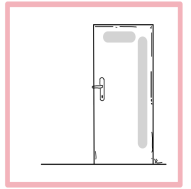
Let's think about sequences



Accelerated AI
Algorithms for
Data-Driven
Discovery



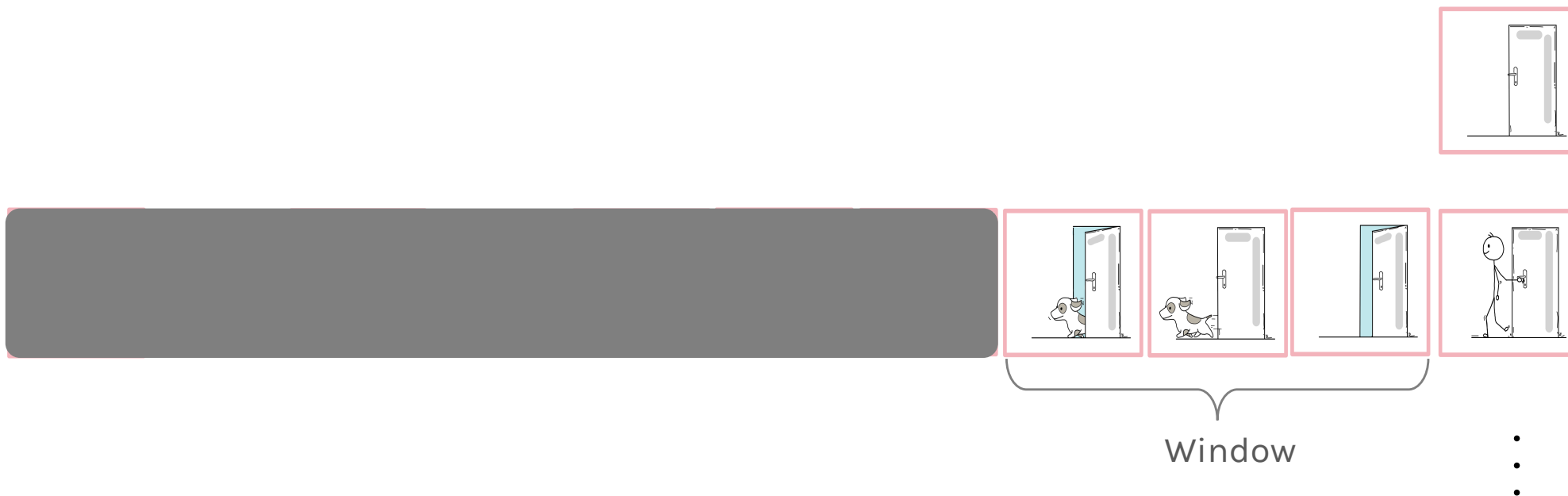
What is the next frame?



There are many options: the door shutting, the person coming out, etc.:

Let's think about sequences

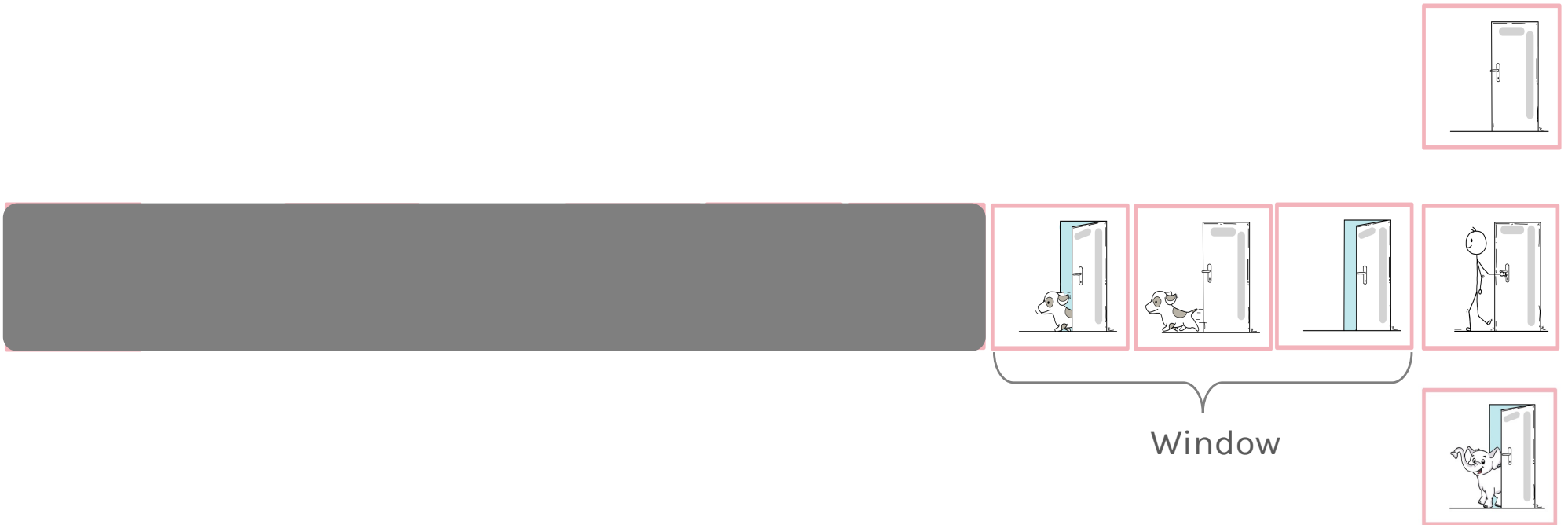
What is the next frame?



There are many options: the door shutting, the person coming out, etc.⋮

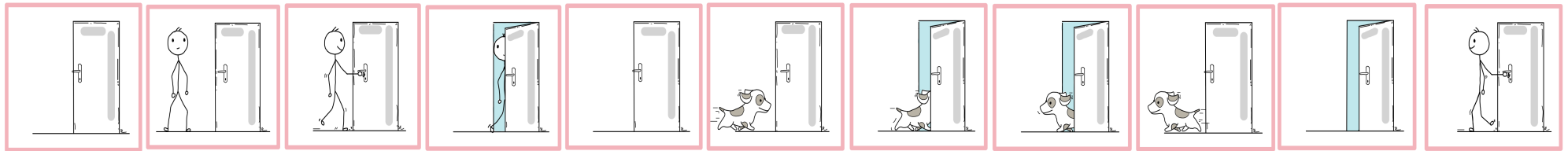
Let's think about sequences

What is the next frame?



There are many options: the door shutting, the person coming out, etc.:
So a longer memory is needed.

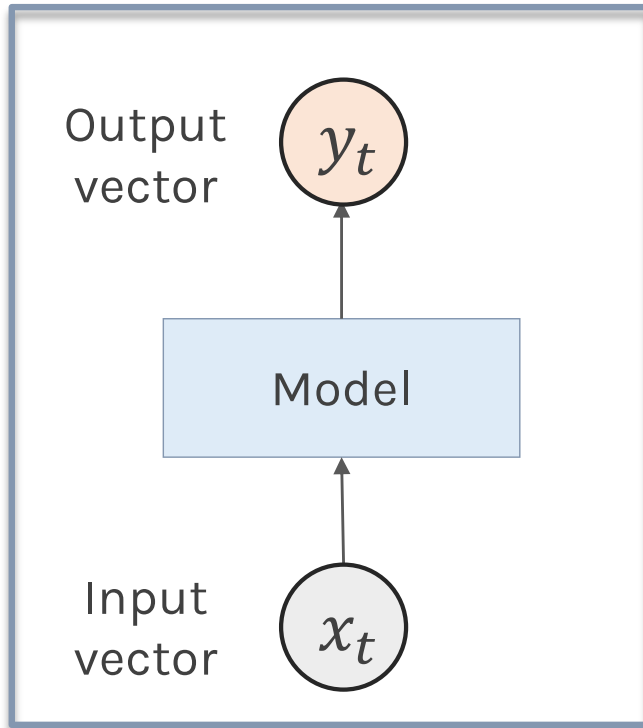
Let's think about sequences



With longer memory, it is easier to predict

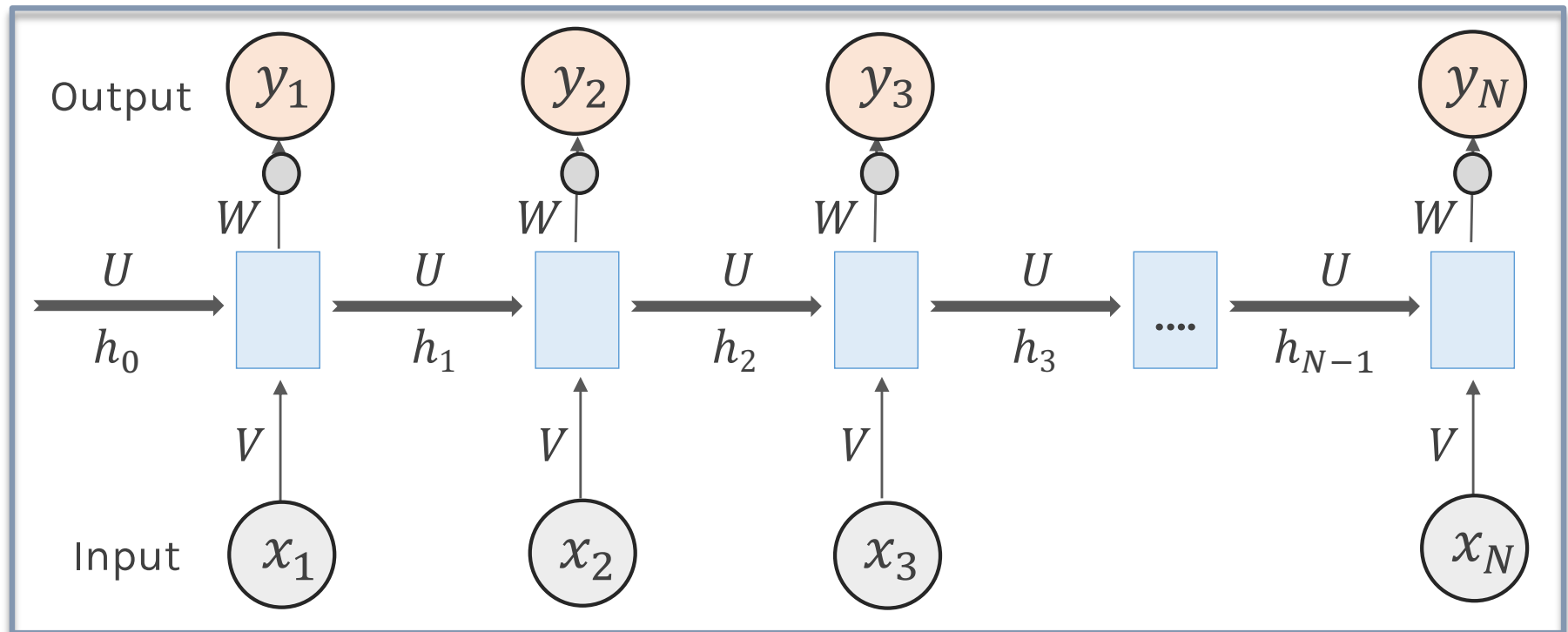
Intro to RNNs

FEED FORWARD NEURAL NETWORK



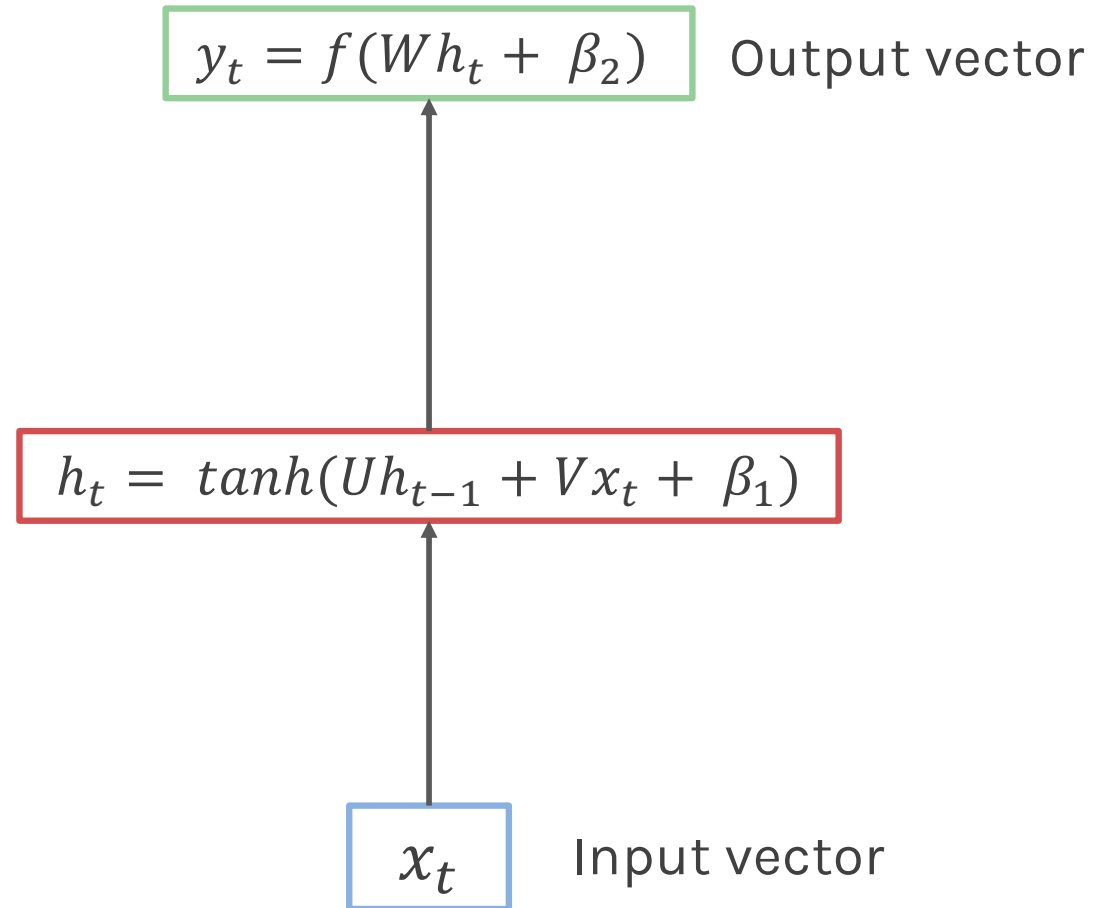
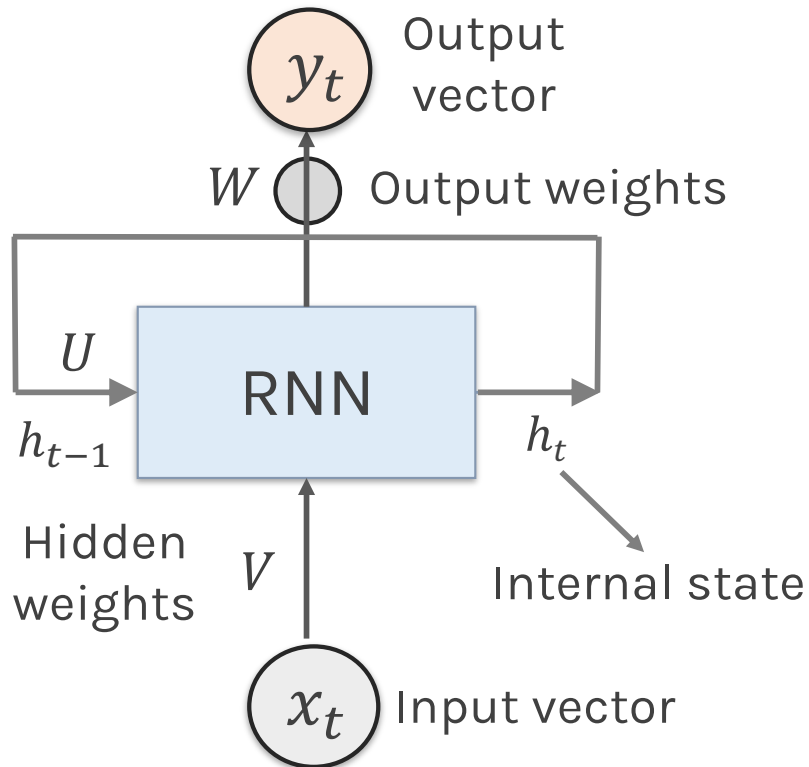
Cannot maintain
previous information

RECURRENT NEURAL NETWORK



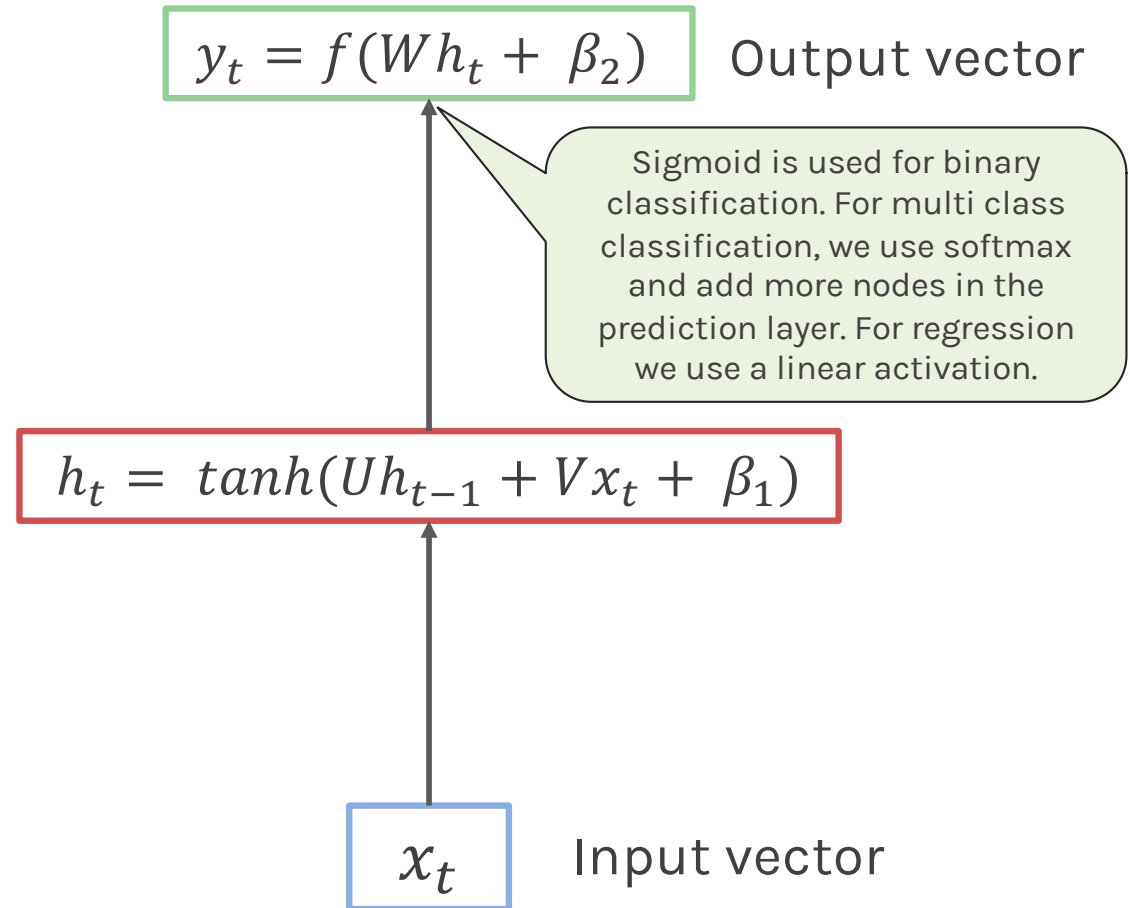
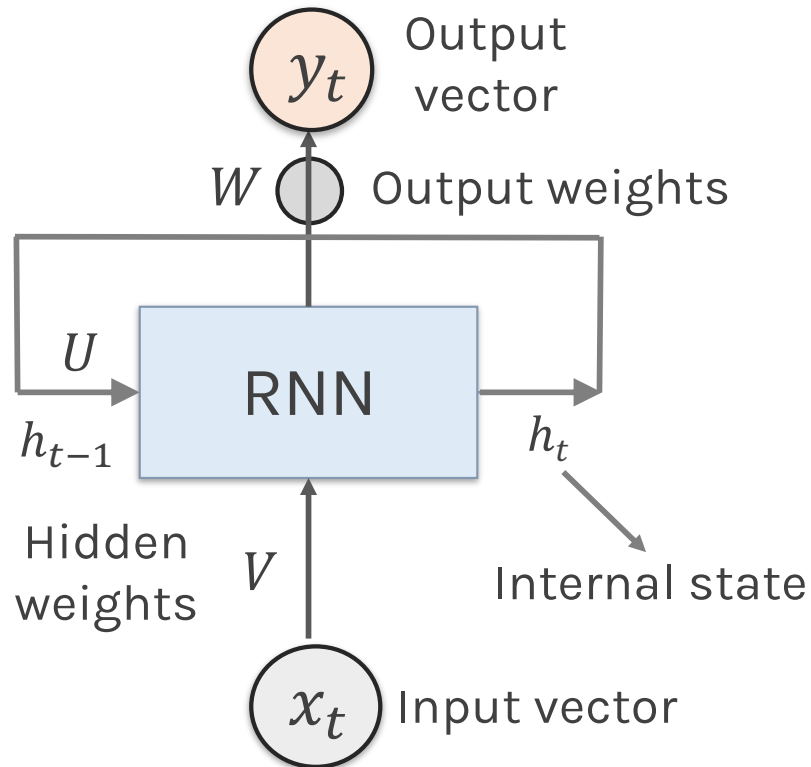
- The term recurrent comes from the fact that information is being passed from one time step to the next internally within the network.
- Network has loops for information to persist over time.

Intro to RNNs



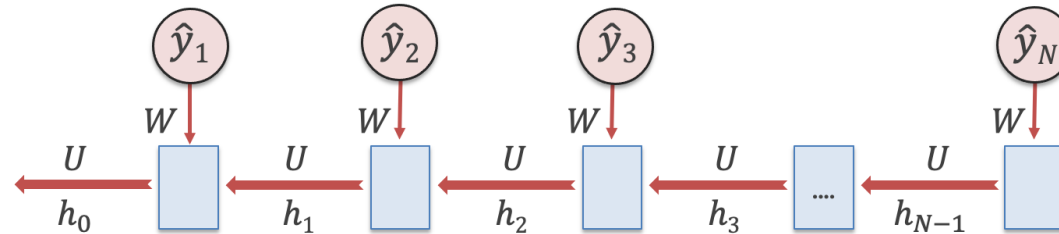
U , V and W are three different weight matrices learned during training

Intro to RNNs



U , V and W are three different weight matrices learned during training

Intro to RNNs: training issues



For longer sequences, we must backpropagate through more time steps.

This requires the gradient to be multiplied many times which causes the following issues:

If many values < 1 , then the product, i.e., the gradient, will be close to zero. This is called the **vanishing gradient problem**.

This causes the parameters to update very slowly.

If many values > 1 , then the product, i.e., the gradient, will explode. This is called the **exploding gradient problem**.

This causes an overflow problem.

Let's think about tokens

How do we process a sequence? **Tokenization** converts a sequence into a symbolic or structured representation (token) that can modeled

For **language**, this can be:

- a word: "I am here" -> ["I", "am", "here"]
- rule-based (regex): "can't" -> ["ca", "n't"]
- subword: "unhappiness" -> ["un", "happiness"]
- character: "cat" -> ["c", "a", "t"]

GPT models use Byte-Pair Encoding (BPE): subwords are formed iteratively from mergers of the most frequently adjacent character pairs

For **time series**:

- discretization
- event-based
- chunking (fixed window)
- autoencoder to a discrete latent space (codebook token)

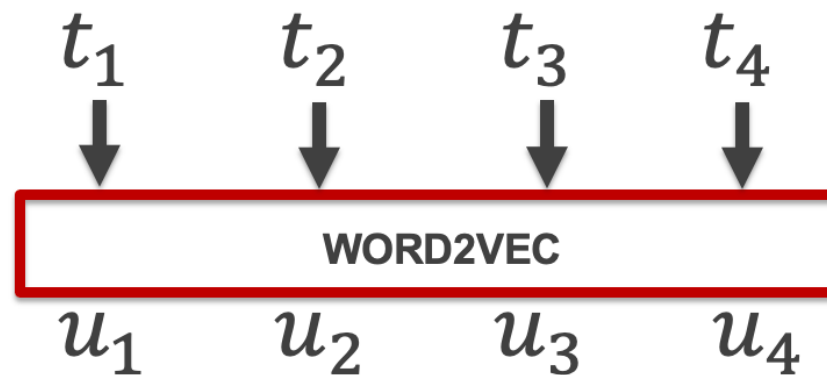
Let's think about embedding

How do we represent high dimensional discrete data (like words, categories, or nodes in graph) into a dense, continuous vector representation in a lower dimensional space?

An **embedding** is a learned mapping from a discrete input space to a continuous vector space. Similar items have closer vectors.

Word2Vec is a neural network model developed by Mikolov et al. (2013) at Google that learns word embeddings.

they crossed the street



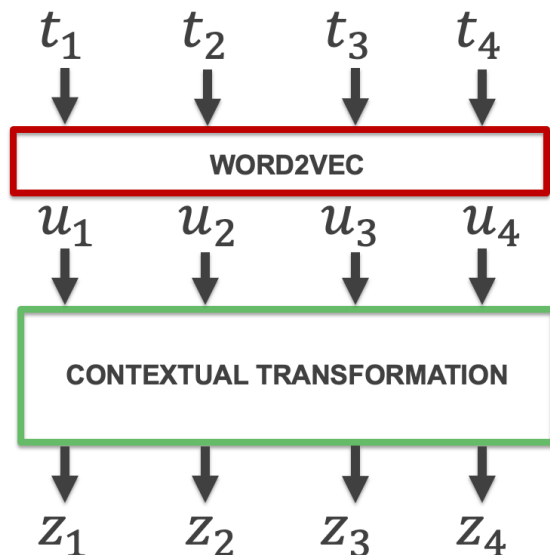
Let's think about context

Donald was hit by a bus because **they** crossed the street

How do we find what “they” is referring to?

Make a contextual embedding that is a linear combination of the other word embeddings:

they crossed the street



$$z_1 = \alpha_{11}u_1 + \alpha_{12}u_2 + \alpha_{13}u_3 + \alpha_{14}u_4$$

$$z_2 = \alpha_{21}u_1 + \alpha_{22}u_2 + \alpha_{23}u_3 + \alpha_{24}u_4$$

$$z_3 = \alpha_{31}u_1 + \alpha_{32}u_2 + \alpha_{33}u_3 + \alpha_{34}u_4$$

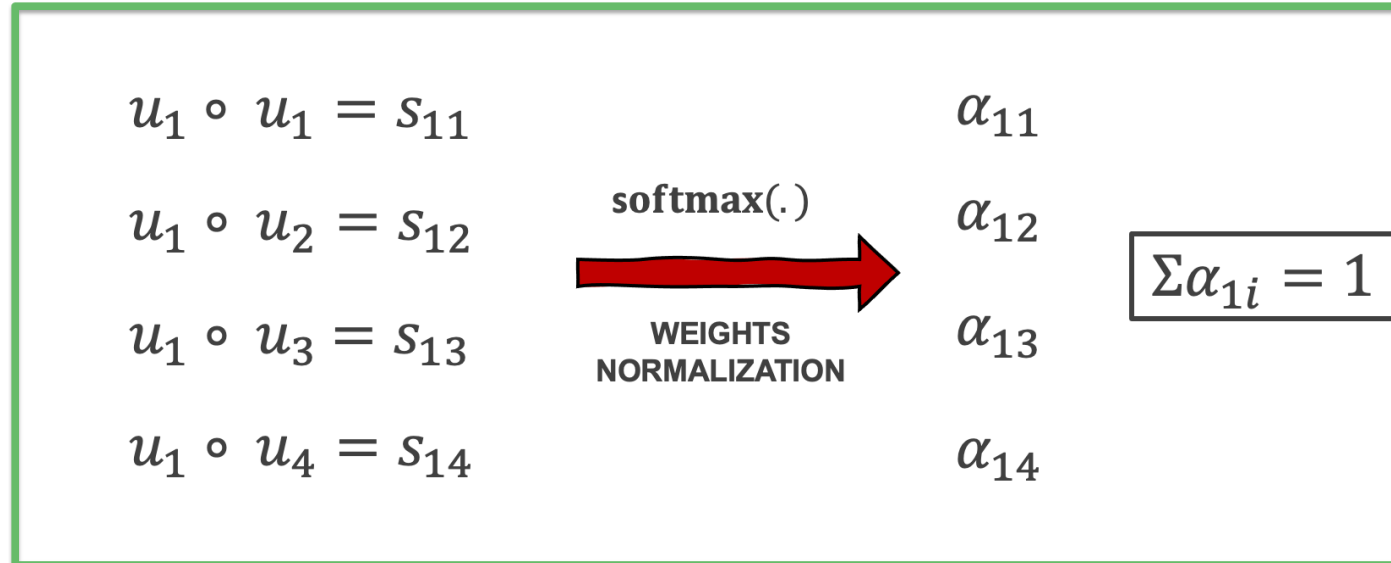
$$z_4 = \alpha_{41}u_1 + \alpha_{42}u_2 + \alpha_{43}u_3 + \alpha_{44}u_4$$

Let's think about context

The coefficients are weights that encapsulate the degree to which new contextual embeddings should incorporate the influence of other tokens

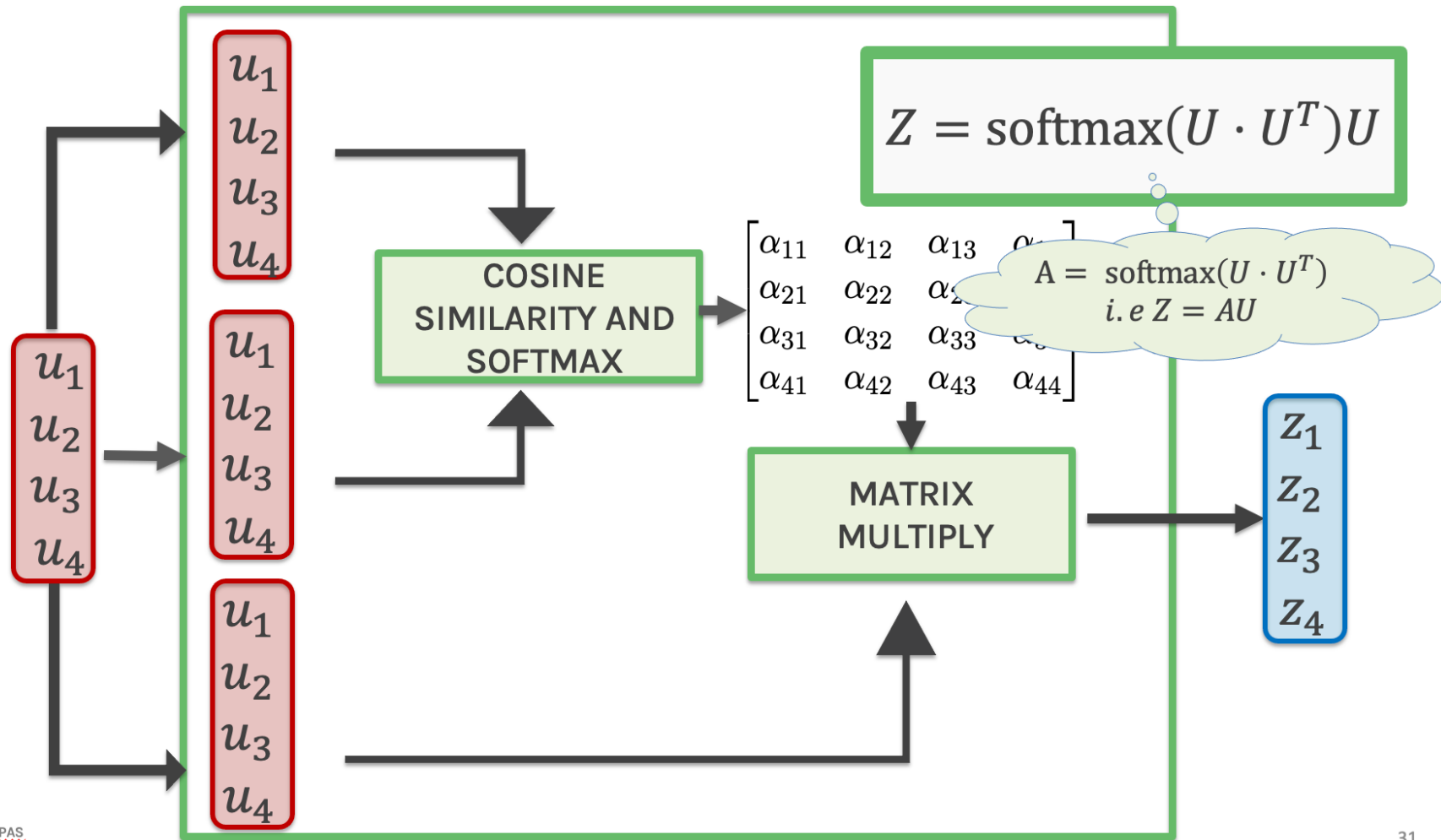
For words, we use the cosine similarity:

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



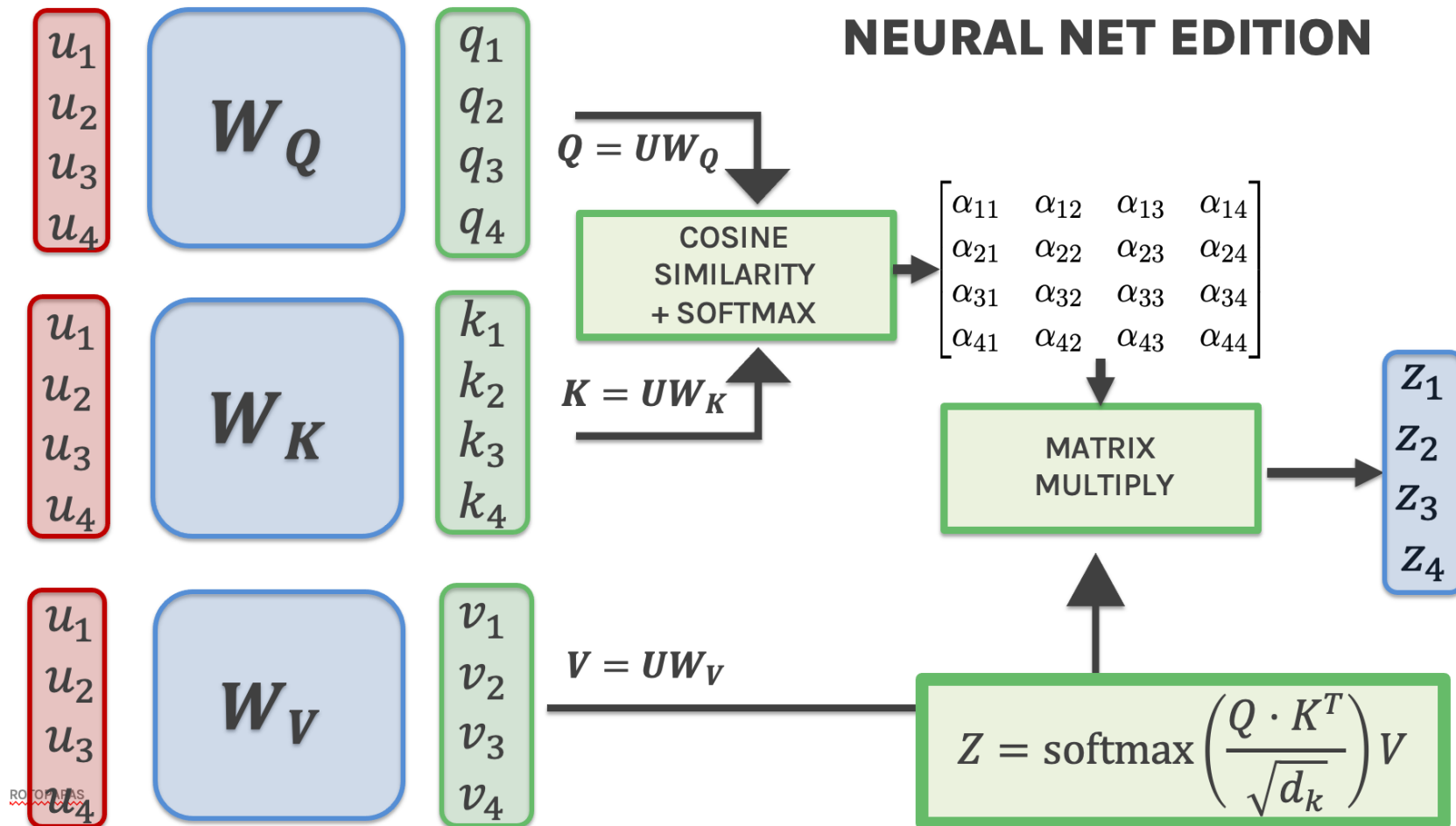
$$z_1 = \alpha_{11}u_1 + \alpha_{12}u_2 + \alpha_{13}u_3 + \alpha_{14}u_4$$

Let's think about attention

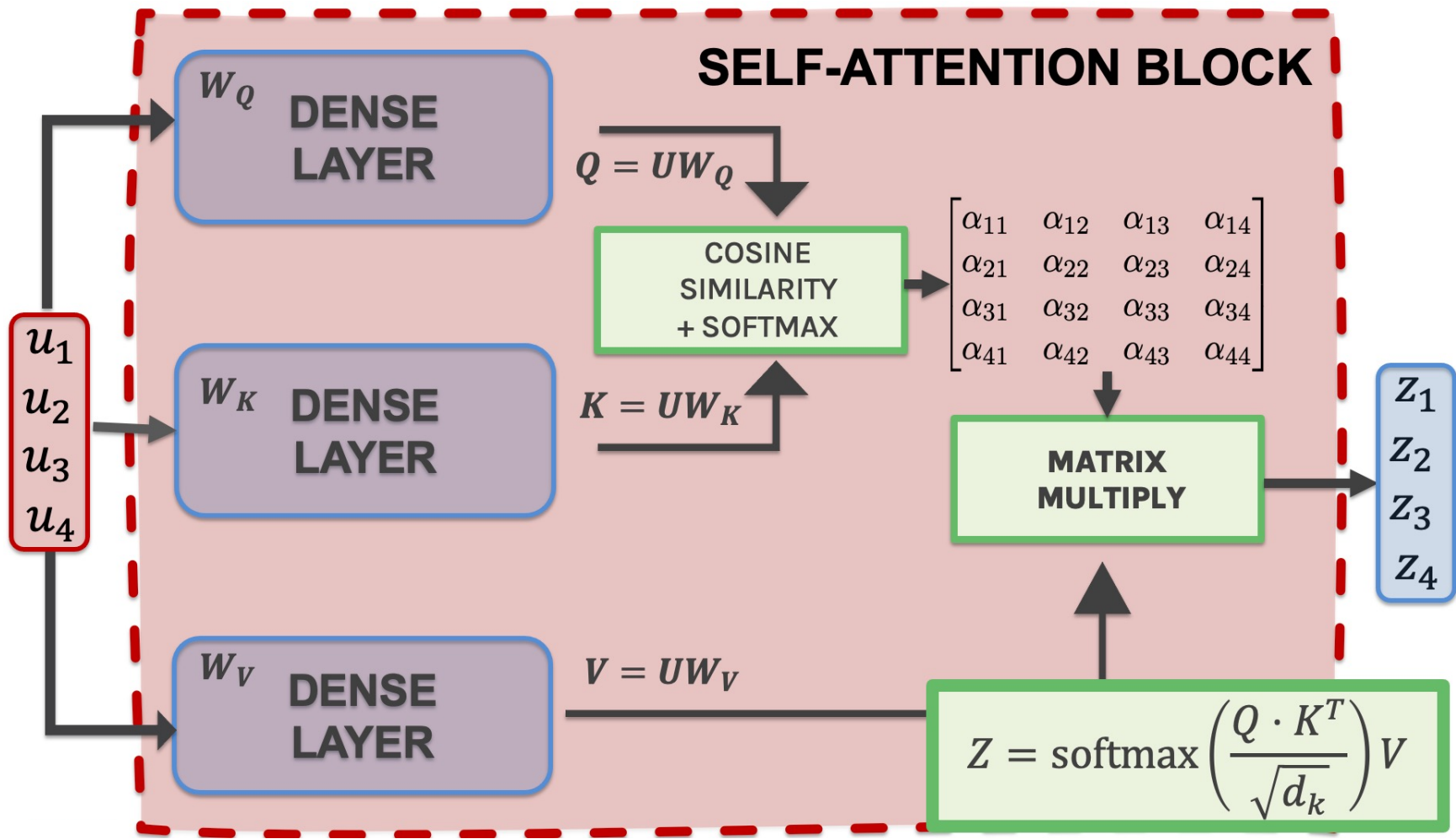


Let's think about attention

- As currently defined, two words will also have the similarity: no weights have been trained from a corpus

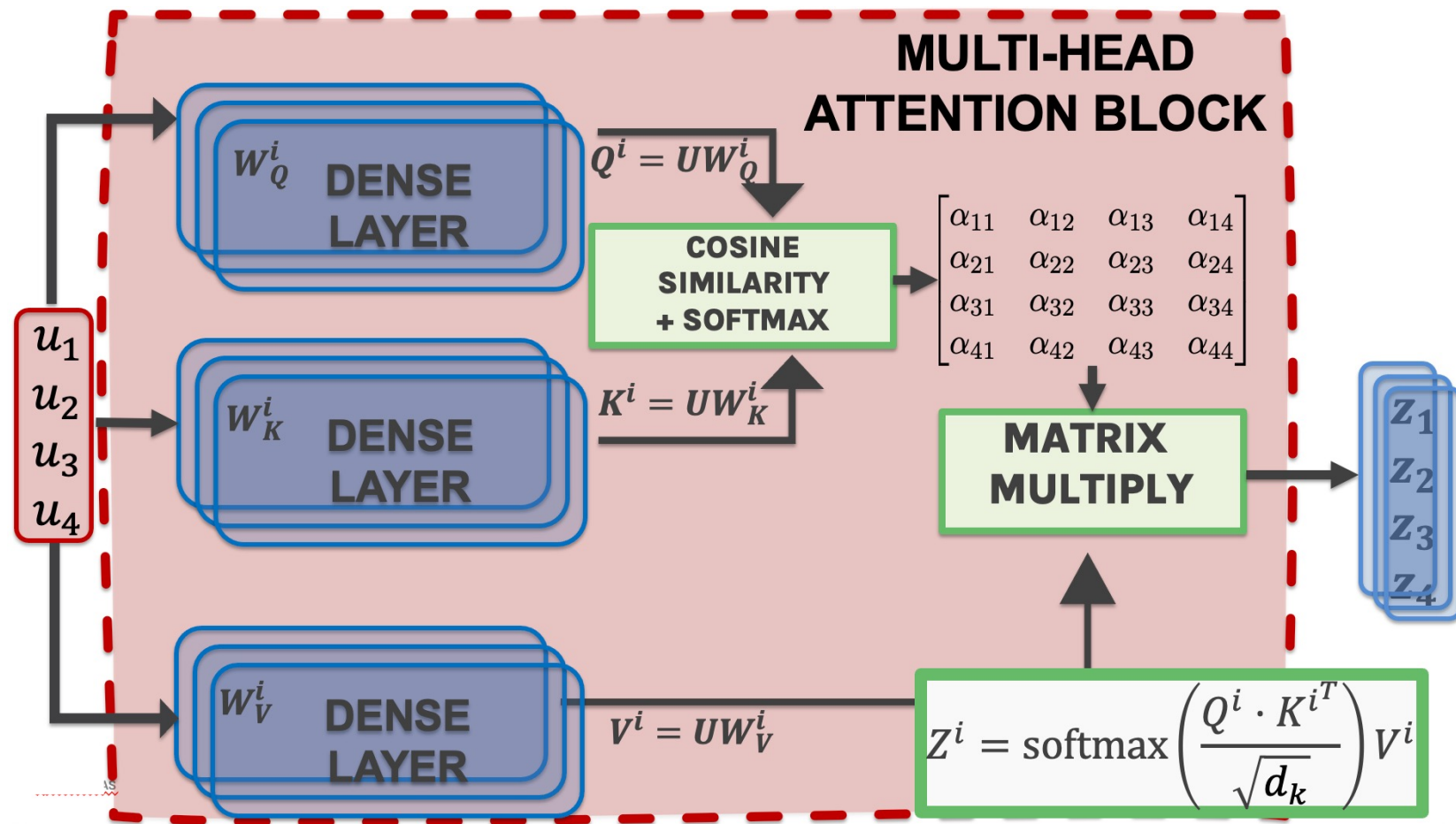


Self-attention block



Multi-head attention block

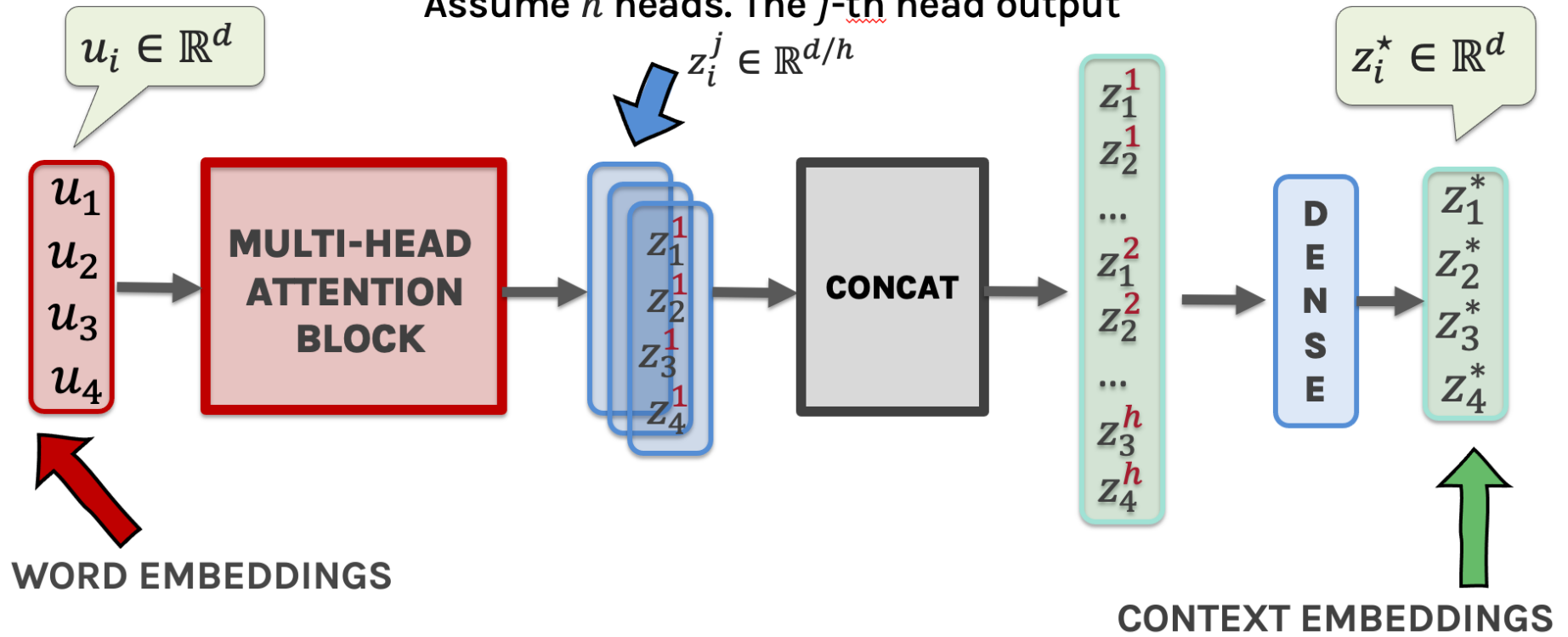
- Attention leads to limited contextual mapping
- We can have multiple attention blocks to look for different relations (cf. CNNs with multiple filters to learn different features)



Multi-head attention block

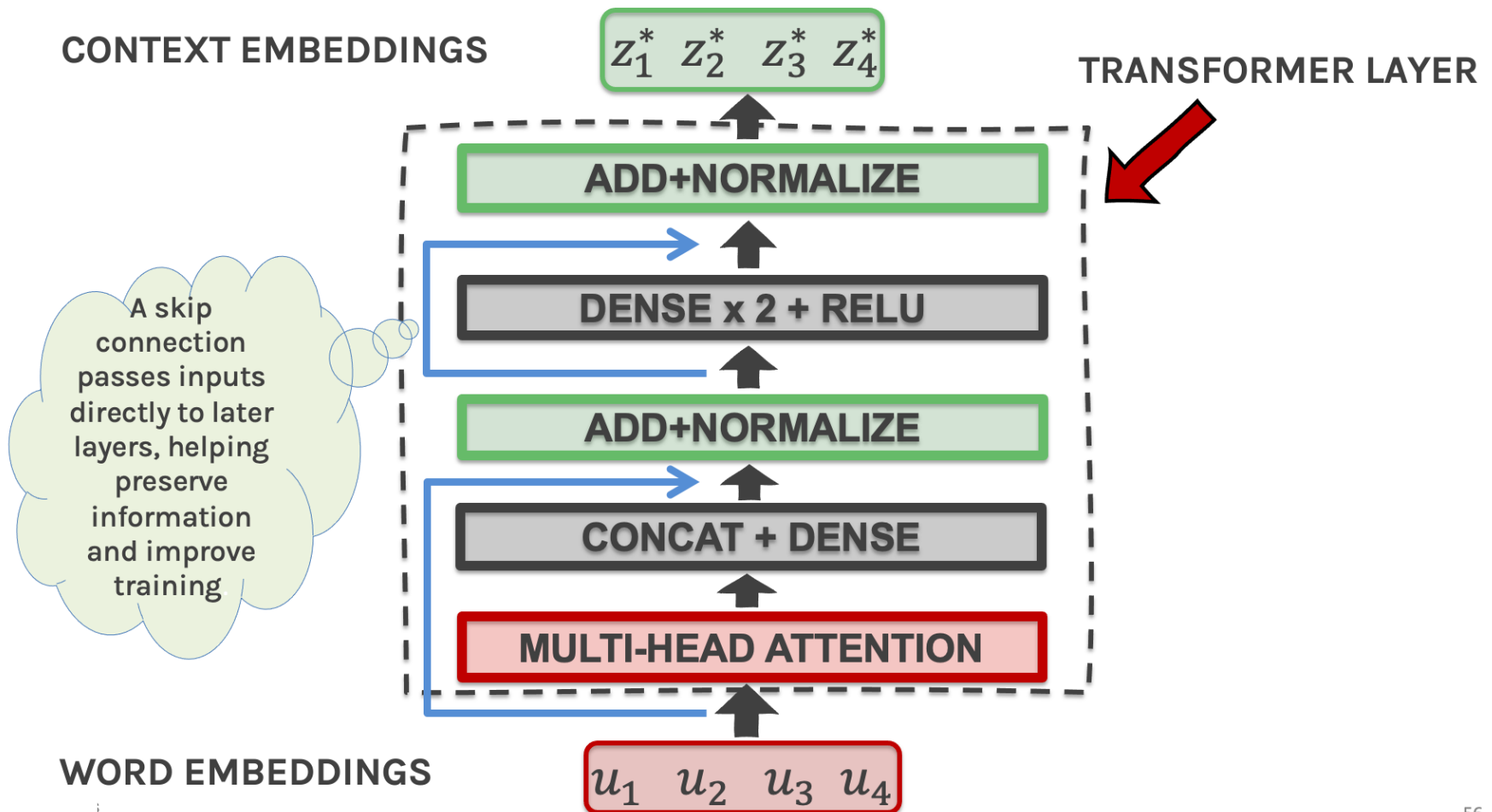
Multiple outputs. Each head produce a lower dimensional vector.

Assume h heads. The j -th head output



Transformer layer

- Add skip connections to avoid the issue of vanishing gradients
- Add normalization layers to help gradients flow better during backpropagation



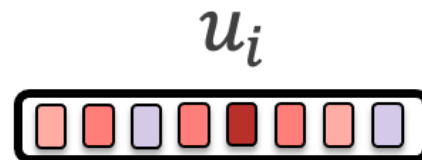
Let's talk about positional information

George spoke to Matthew about the Ay 119 class projects
Matthew spoke to George about the Ay 119 class projects

RNNs inherently take the **order** of words into account

Multi-head attention blocks do not take order into account

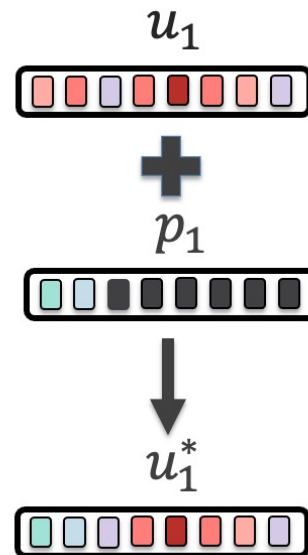
Consider the input embedding u_i for a word at position i :



This will be the same for any position in the sentence

Let's talk about positional information

Modify the input embedding with a positional vector:



The same input embedding will have a different value depending on its position in the sentence

Let's talk about positional information

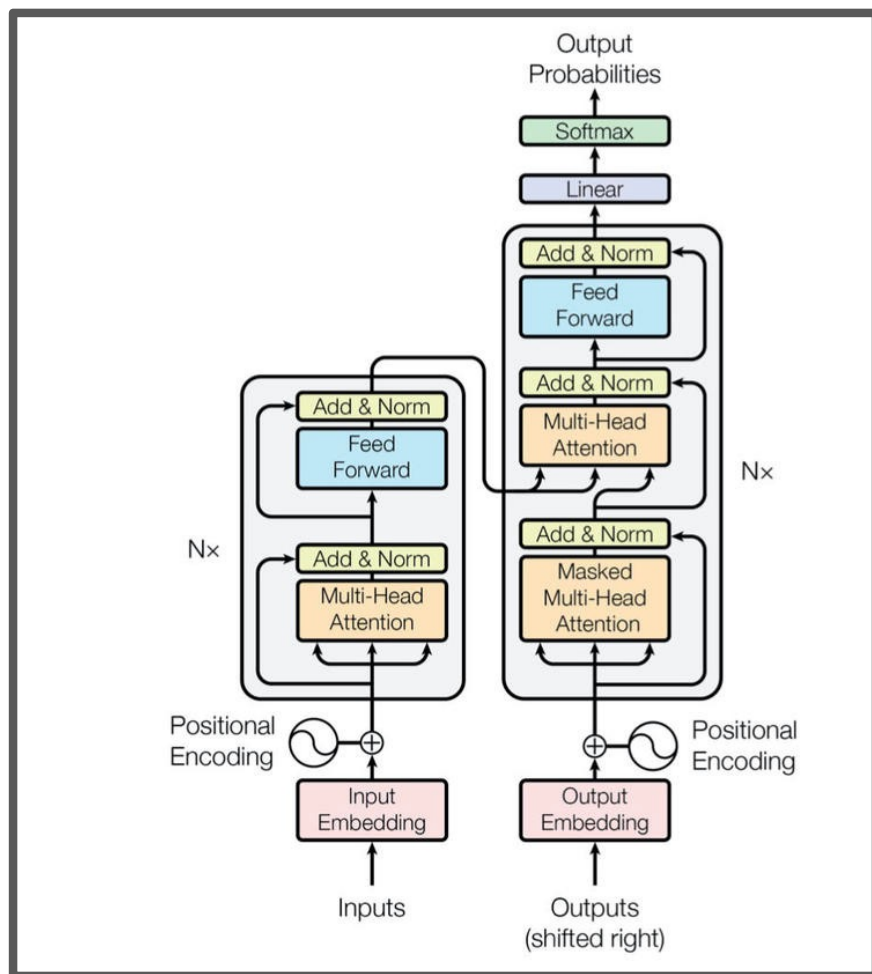


Each position in a sentence is assigned a vector that encodes its position (from Attention is All You Need):

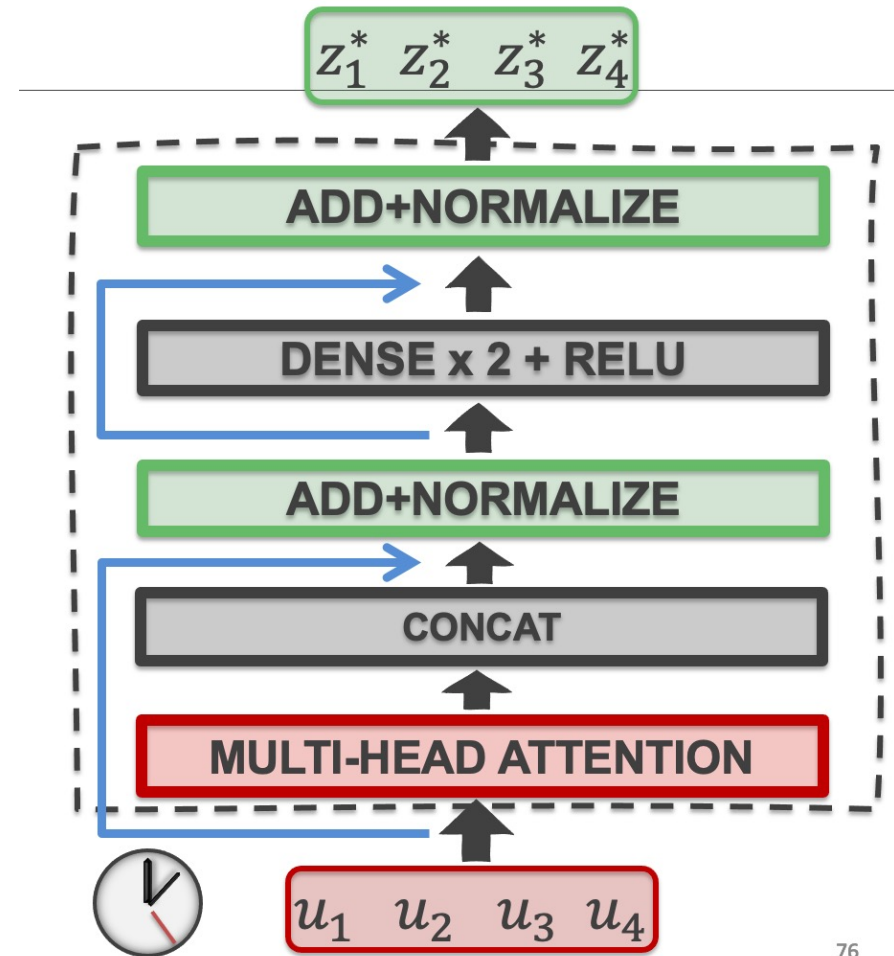
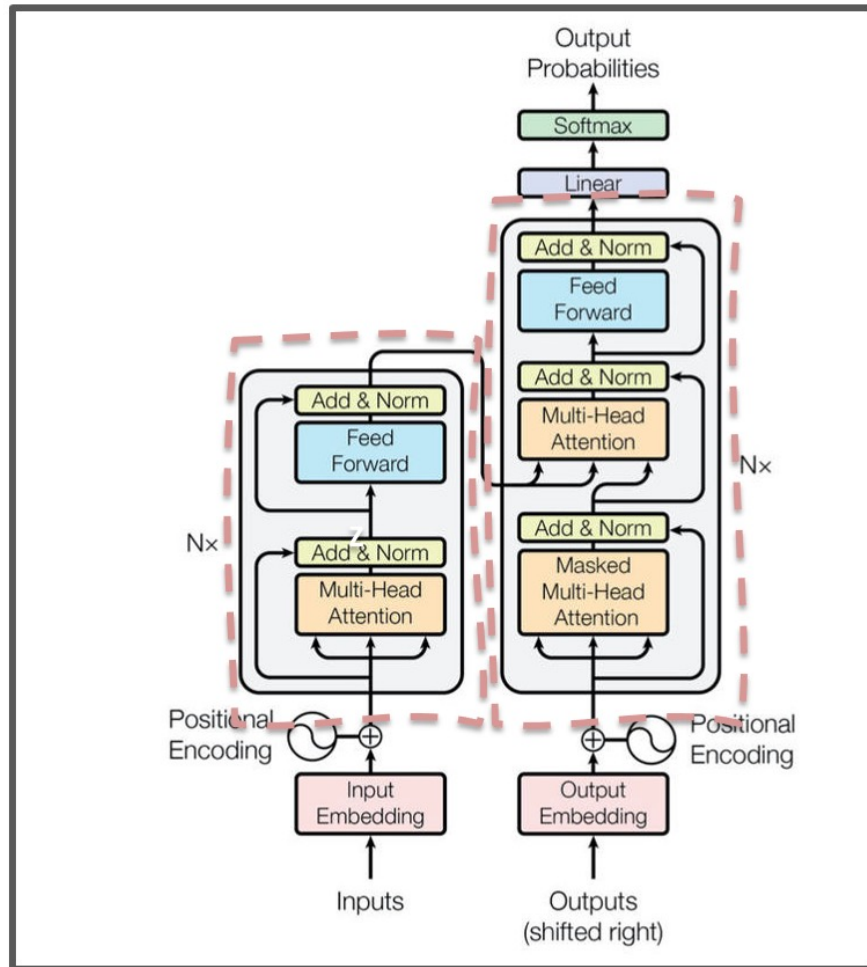
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

The encoding uses a mix of trigonometric functions to ensure each position has a unique yet repeatable pattern

Putting it all together: transformers



Putting it all together: transformers



76