# Code.Hub

The first Hub for Developers

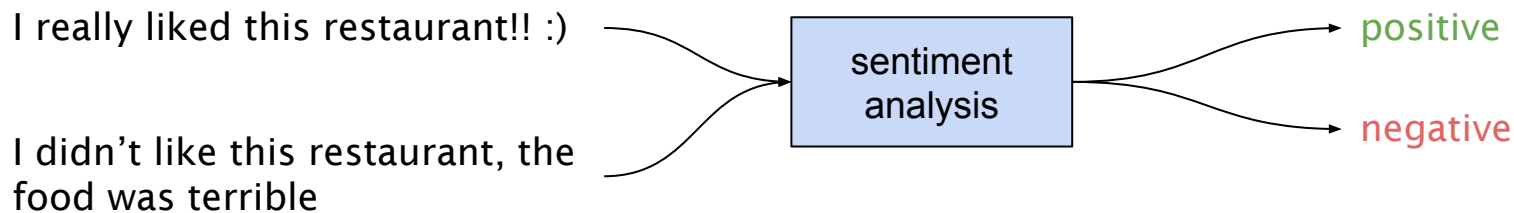# Sequential Models

Thanos Tagaris

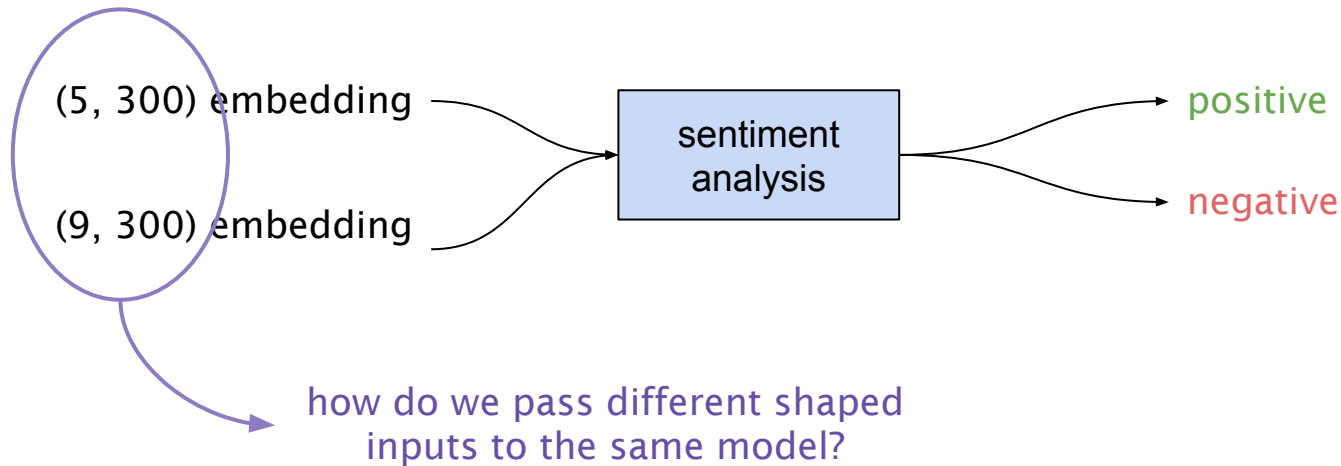# NLP timeline up till today…

1 - 2      3 - 4      5      6      7 - 8

BoW and TF-IDF as VSMs    |    Word Embeddings    |    **Sequential Models**    |    **Attention Mechanism**    |    Transformers    |    LLMs

Trainable word embeddings gave rise to extended Neural Network usage for NLP tasks.

Code.Hub

# How do we use word embeddings for our downstream tasks?



I really liked this restaurant!! :)

I didn't like this restaurant, the food was terrible

sentiment analysis

positive

negative

Code.Hub

# How do we use word embeddings for our downstream tasks?



(5, 300) embedding

(9, 300) embedding

sentiment analysis

positive

negative

how do we pass different shaped inputs to the same model?

Code.Hub

# Average embedding

a →

quick →

brown →

fox →

jumped →

over →

the →

lazy →

dog →

average embedding


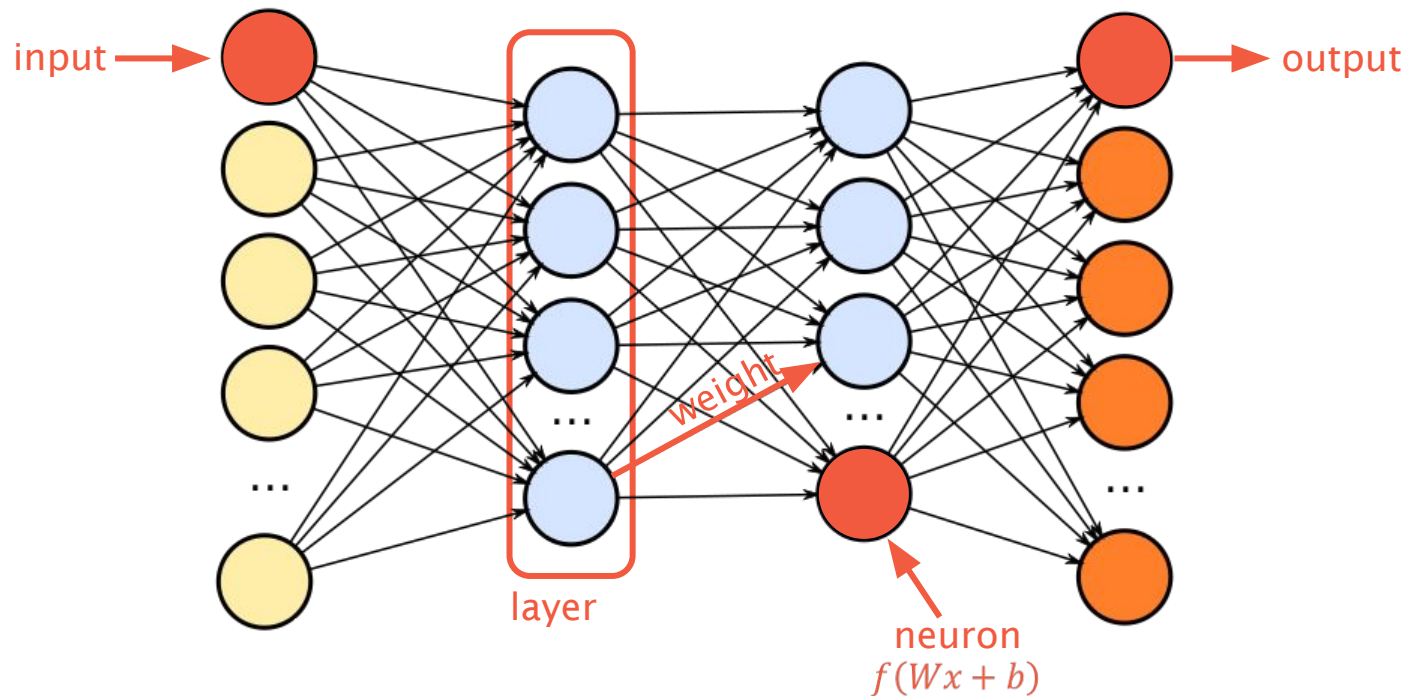
Code.Hub

# Average embedding

Pros:

- always the same dim, can be used to train ML models
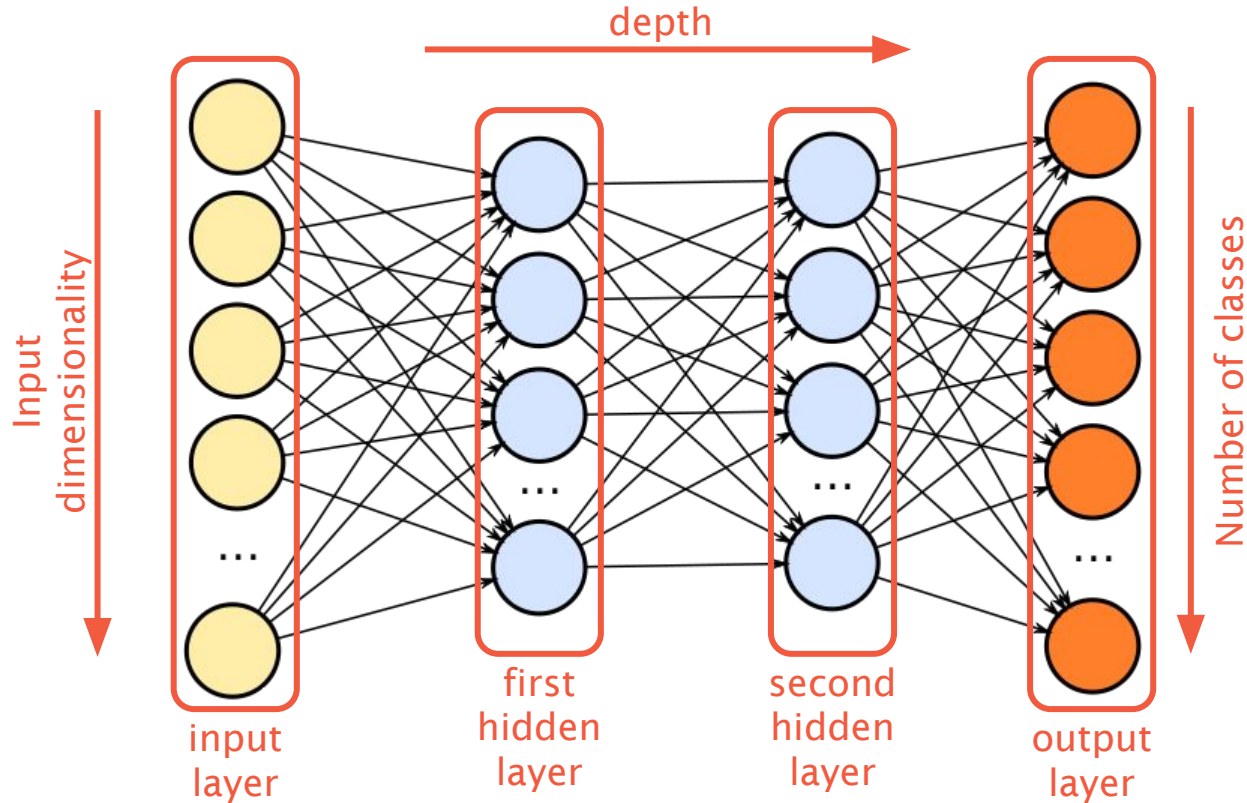- cheap to compute

Cons:

- cannot capture word order
- degenerates with long sentences
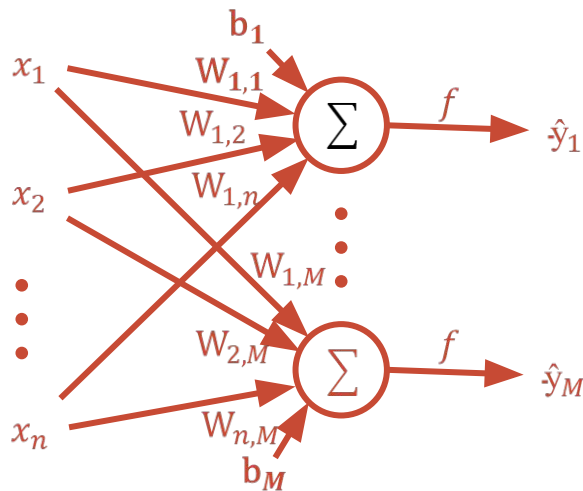
Is there a better way to encode a sentence?

Code.Hub

# Neural Networks – reminder



input → 

output

layer

weight

neuron
$f(Wx + b)$

Code.Hub

# Neural Networks – reminder

# Inside each layer…



In matrix format

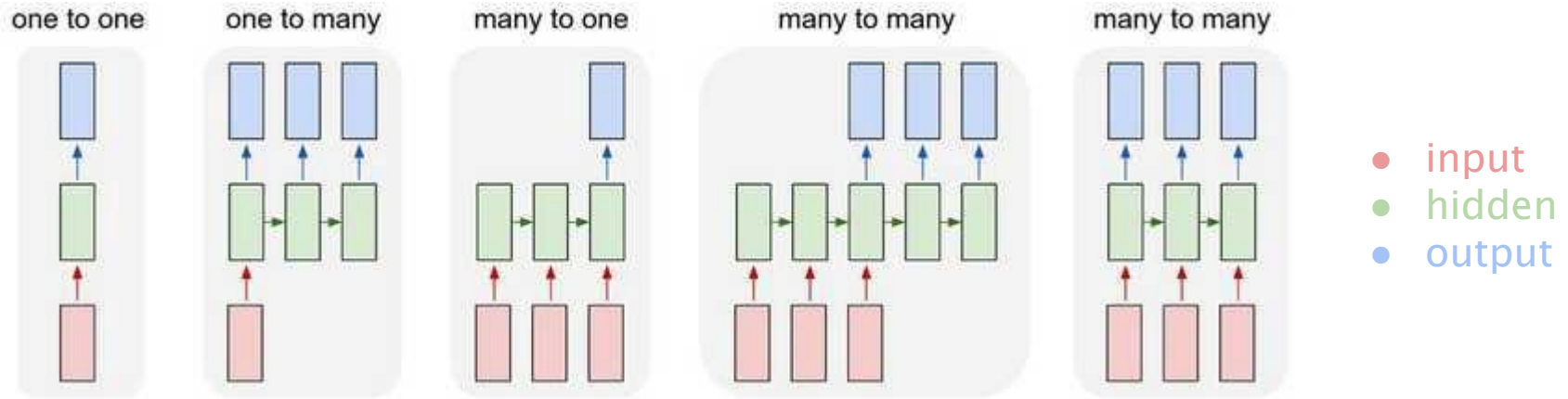$$\hat{y} = f(XW + b)$$

If we have n input features and M neurons:

- Weight matrix has $M \times n$ values.
- Bias matrix has $M \times 1$ values.

Code.Hub

# Inside each layer…



$$\hat{y} = f(xW + b)$$

# Recurrent Neural Networks



one to one

one to many

many to one

many to many

many to many

- input
- hidden
- output

Code.Hub

# Recurrent Neural Networks

- Case study: sentiment analysis
- Predict if a sentence is positive or negative.
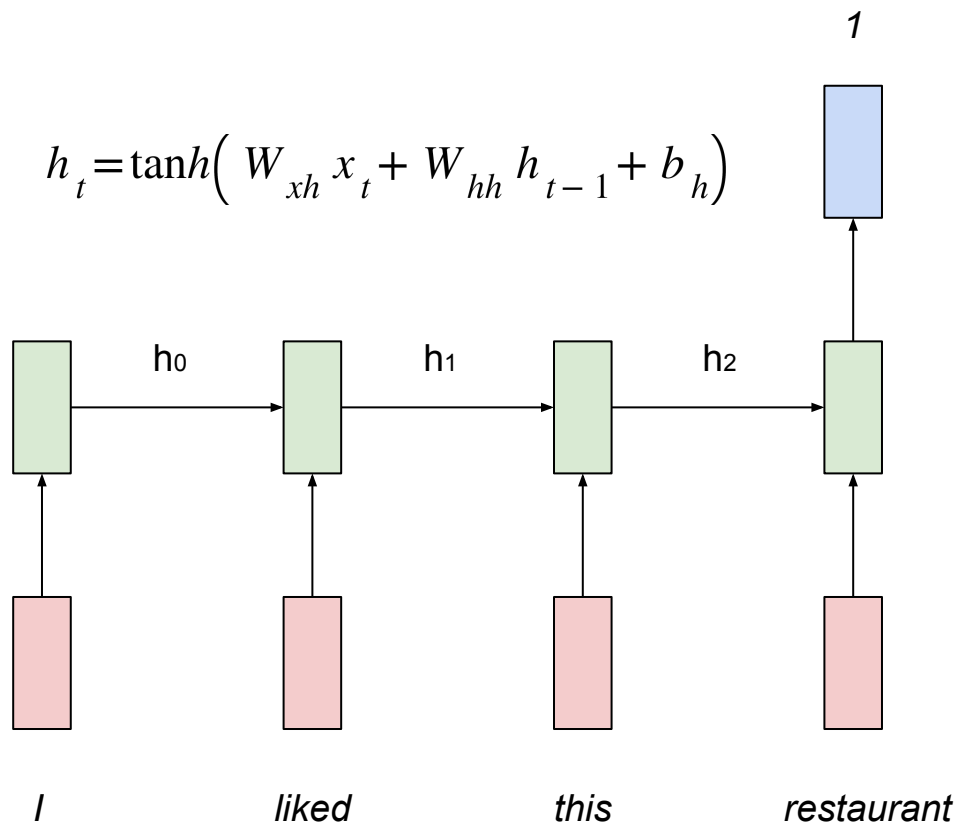- Example *"I liked this restaurant"*
- *Network's input:*

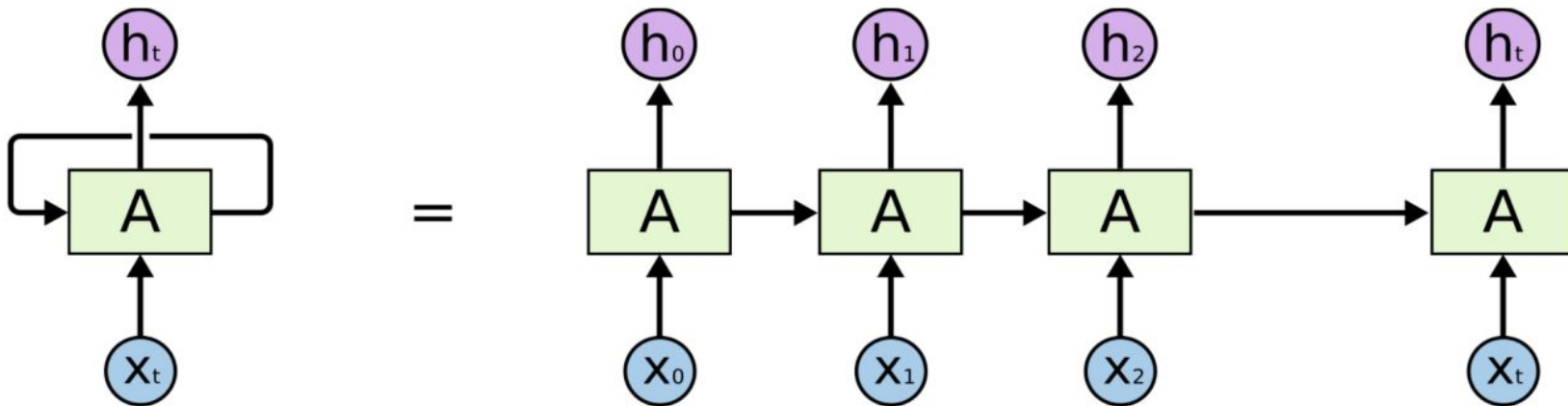|  |  |  |  |
|---|---|---|---|
| *I* | *liked* | *this* | *restaurant* |

Code.Hub

# Recurrent Neural Networks

$$h_t = \tanh\left( W_{xh}\, x_t + W_{hh}\, h_{t-1} + b_h \right)$$

1

h0     h1     h2

I     liked     this     restaurant

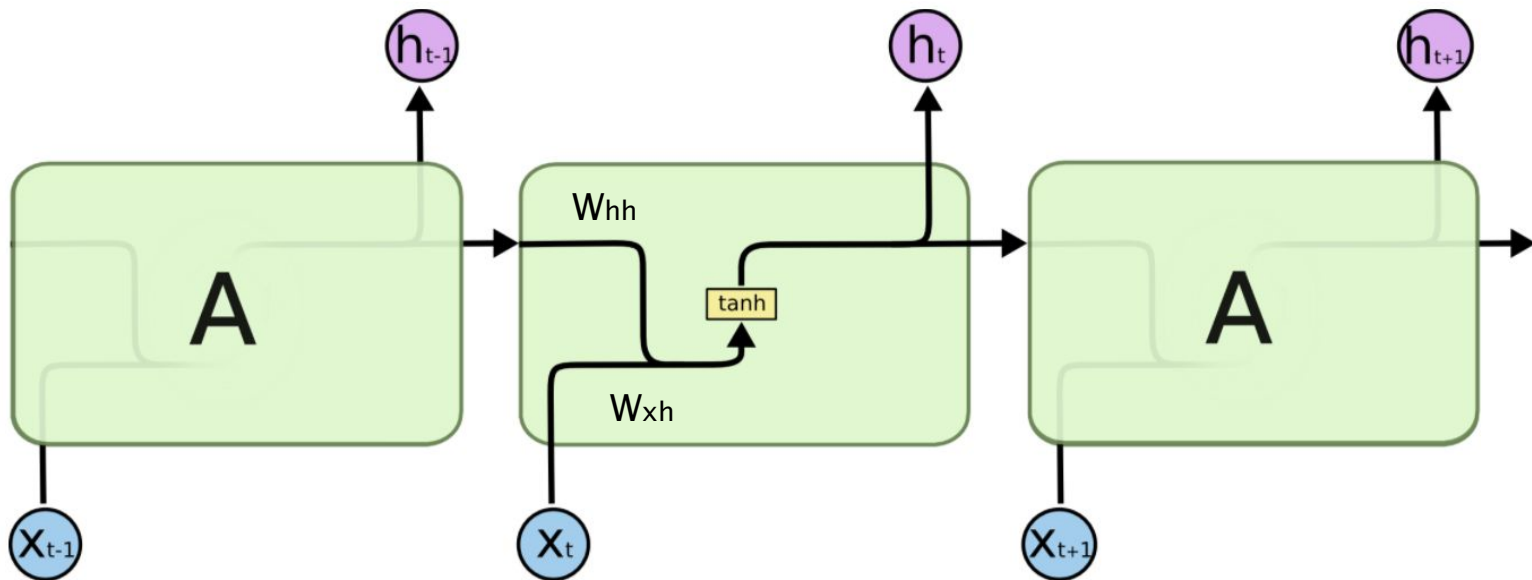- Trained with Backpropagation Through Time
- Issues with vanishing/exploding gradients
- Partially solved with different "neuron" types, e.g. LSTM

Code.Hub

# Recurrent Neural Networks



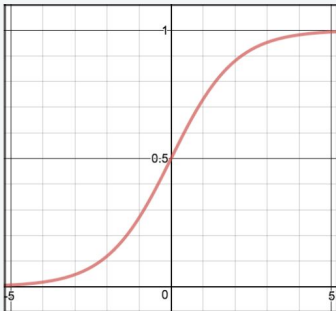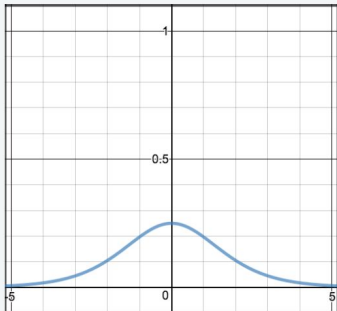Must read: *"The unreasonable effectiveness of RNNs"* - Andrej Karpathy

Code.Hub

# Recurrent Neural Networks



$$h_t = \tanh\left( W_{xh}\, x_t + W_{hh}\, h_{t-1} + b_h \right)$$
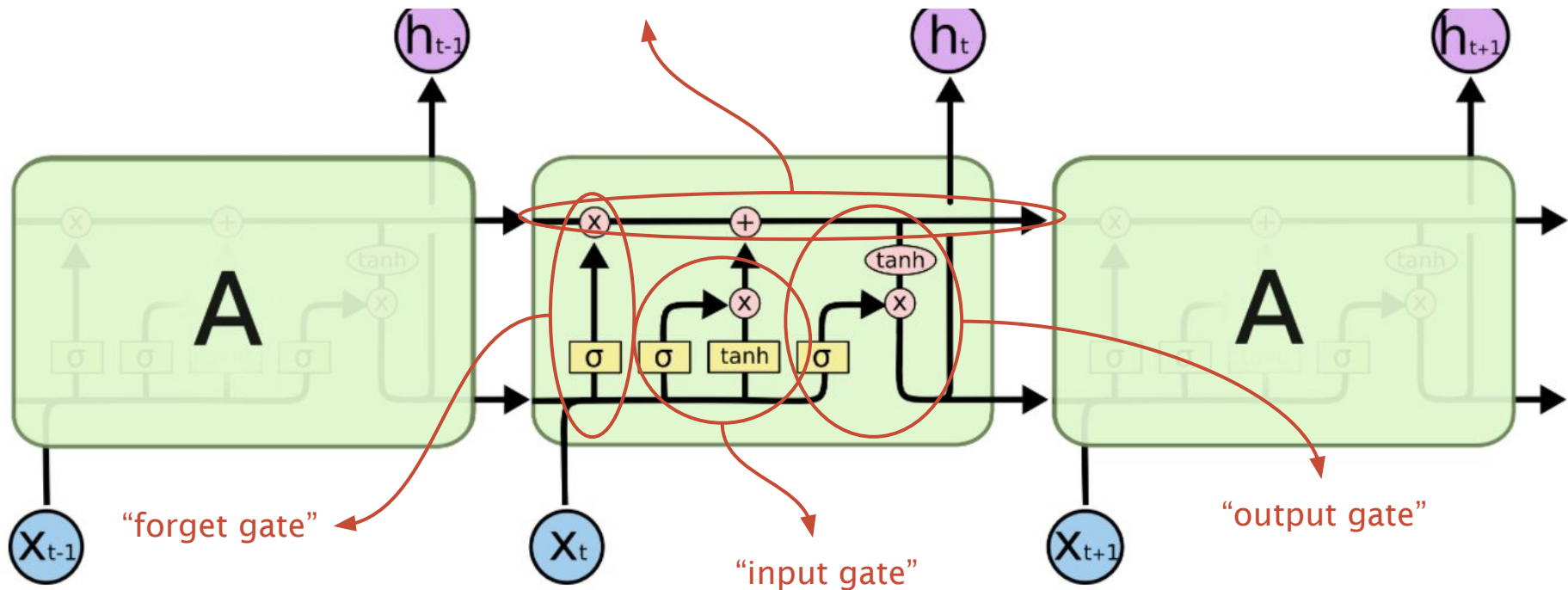
Code.Hub

# Vanishing/Exploding Gradients

- Historically, we couldn't model very long sequences due to numerical instability.
- Example: derivative of sigmoid
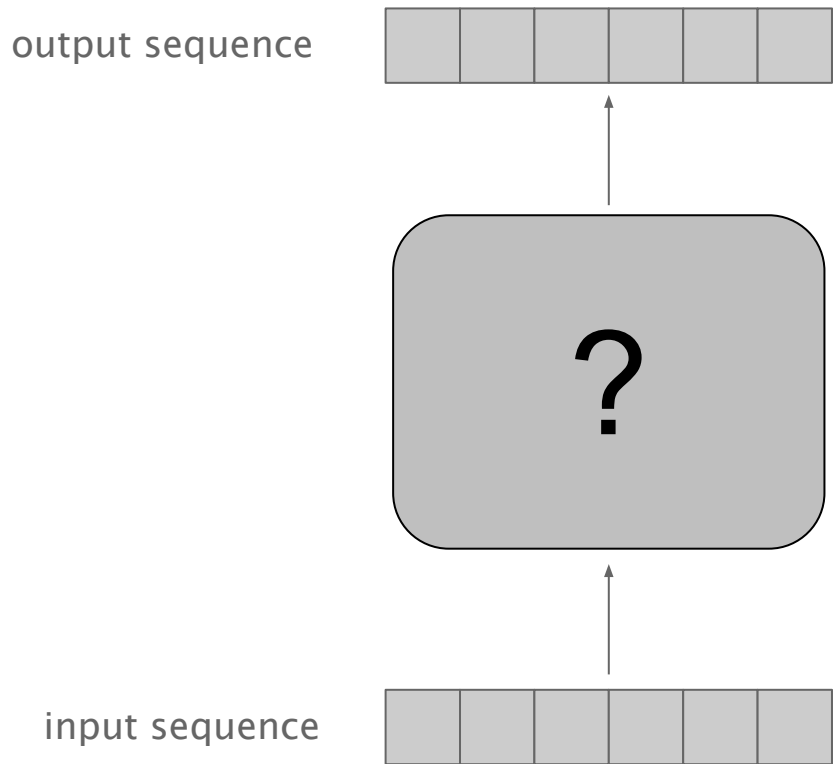
| Function | Derivative |
| --- | --- |
| $S(z) = \dfrac{1}{1 + e^{-z}}$ | $S'(z) = S(z) \cdot (1 - S(z))$ |

```
def sigmoid(z):
  return 1.0 / (1 + np.exp(-z))
```

```
def sigmoid_prime(z):
  return sigmoid(z) * (1-sigmoid(z))
```

Code.Hub

# Long Short-Term Memory (LSTM)



allows for information to pass straight through

"forget gate"

"input gate"

"output gate"
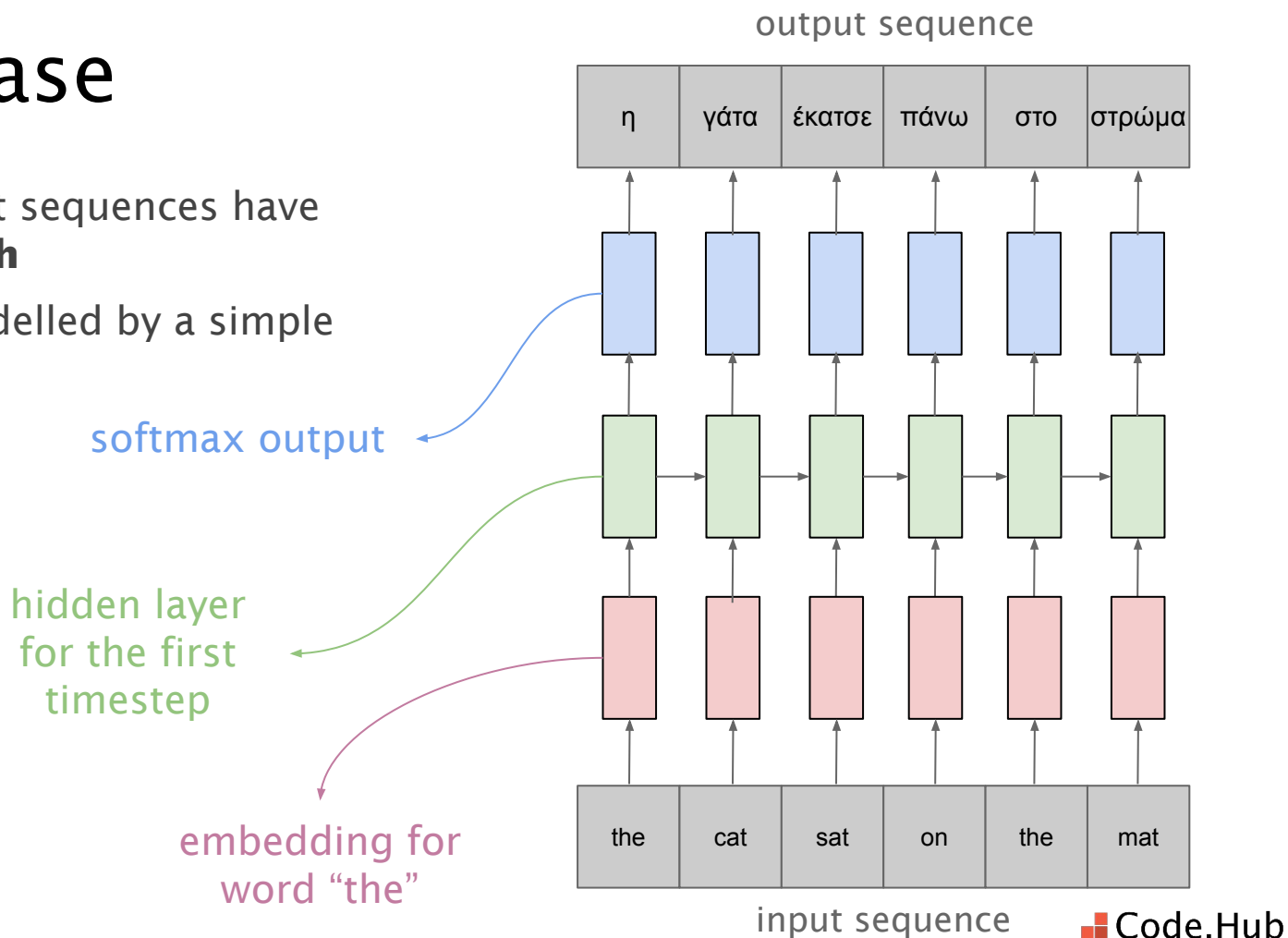
Code.Hub

# Sequence-to-sequence problems

output sequence

? 

input sequence

- Input is a sequence, target is also a sequence
- Most language problems are like this
- E.g.
  - machine translation
  - question answering
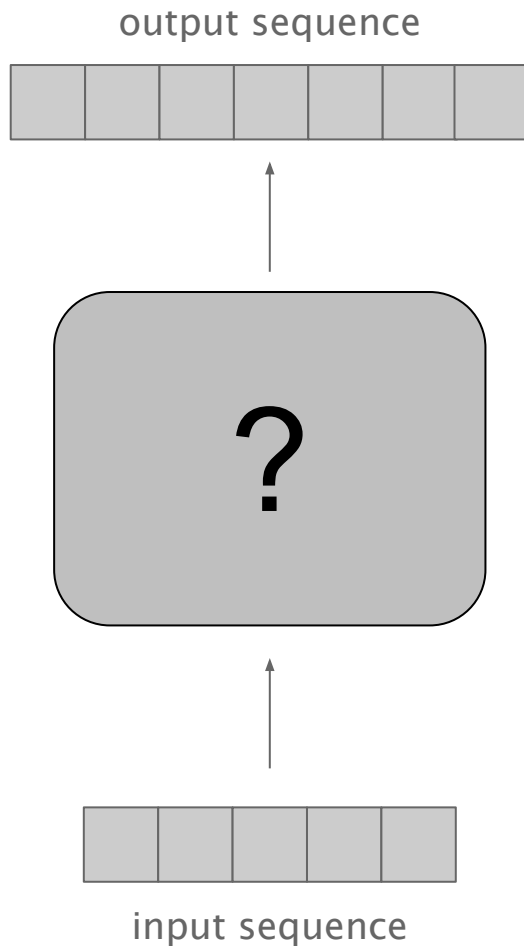  - summarization
- How do we model such tasks?

Code.Hub

# Trivial case

- Input and output sequences have the **same length**
- They can be modelled by a simple LSTM network

| η | γάτα | έκατσε | πάνω | στο | στρώμα |
|---|------|--------|------|-----|--------|

softmax output

hidden layer for the first timestep

embedding for word "the"

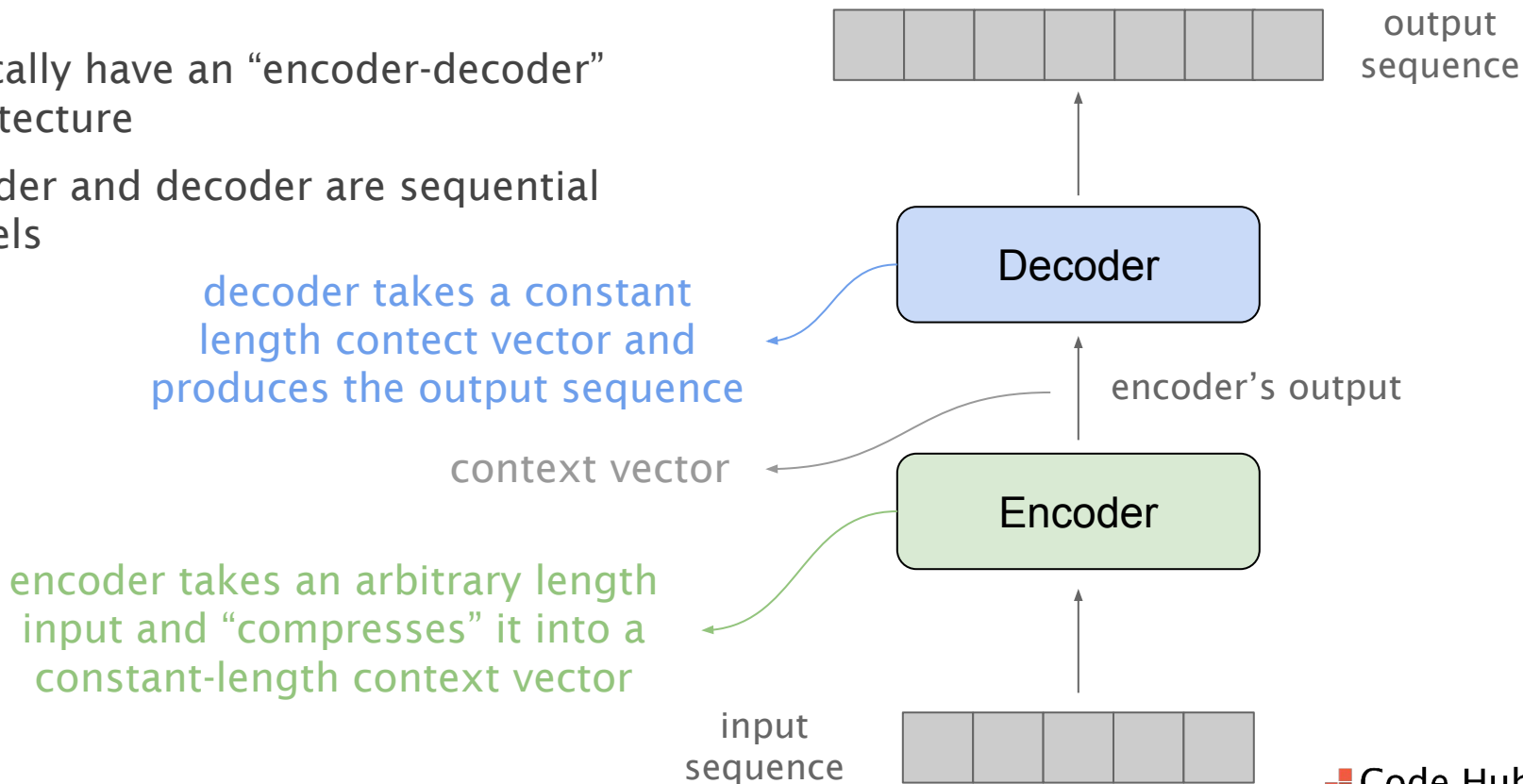| the | cat | sat | on | the | mat |
|-----|-----|-----|----|-----|-----|

input sequence

Code.Hub

# General case

output sequence

- What if we have a different length of input and output sequences?

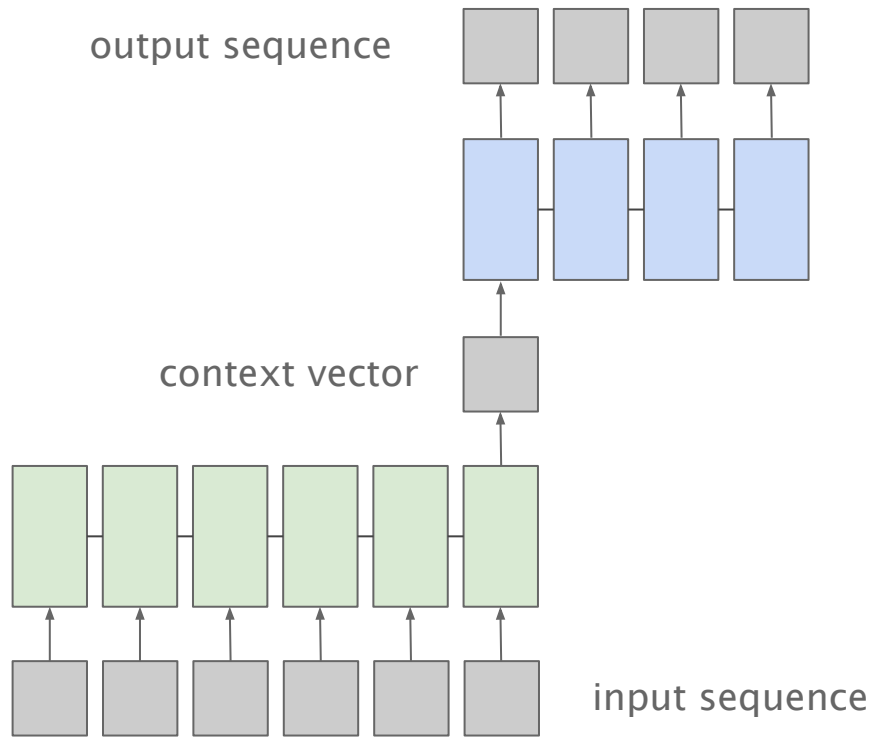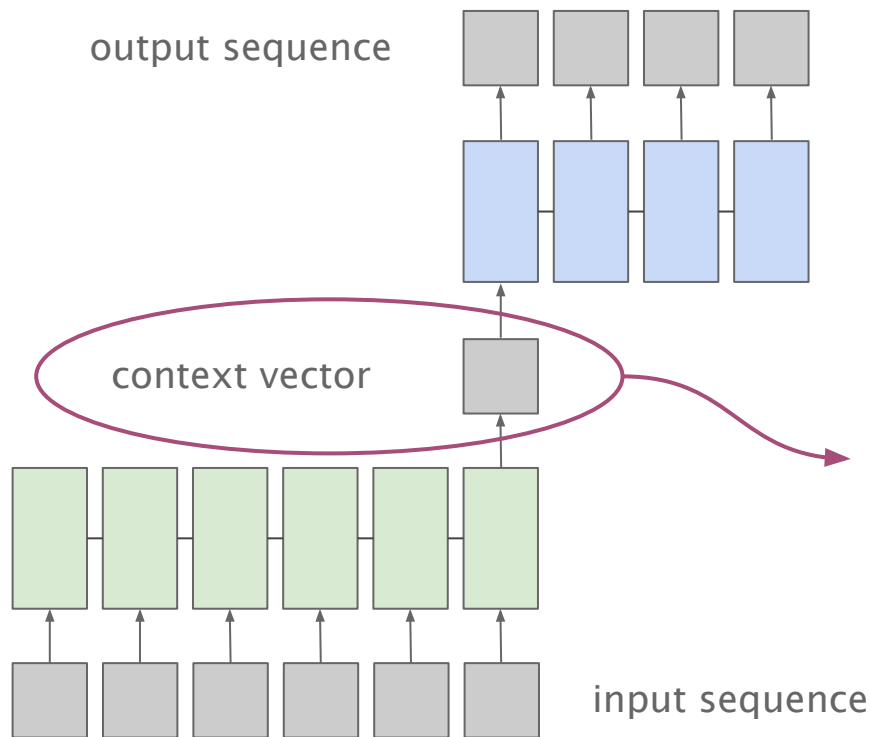- What architecture would we need to model these cases?

?

input sequence

Code.Hub

# Sequence-to-sequence models

- Typically have an "encoder-decoder" architecture

- Encoder and decoder are sequential models

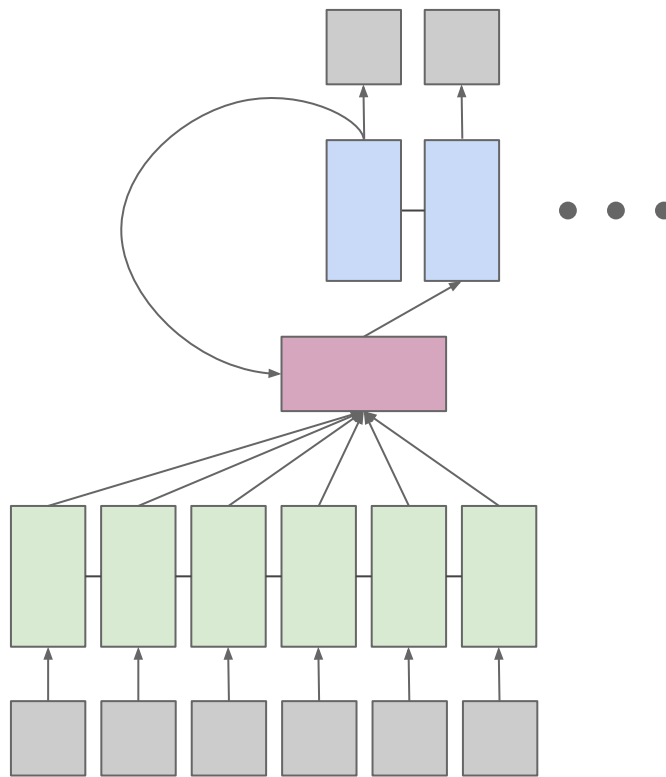decoder takes a constant length contect vector and produces the output sequence

encoder takes an arbitrary length input and "compresses" it into a constant-length context vector

output sequence

Decoder

encoder's output

context vector

Encoder

input sequence

Code.Hub

# A deeper look...



output sequence

context vector

input sequence
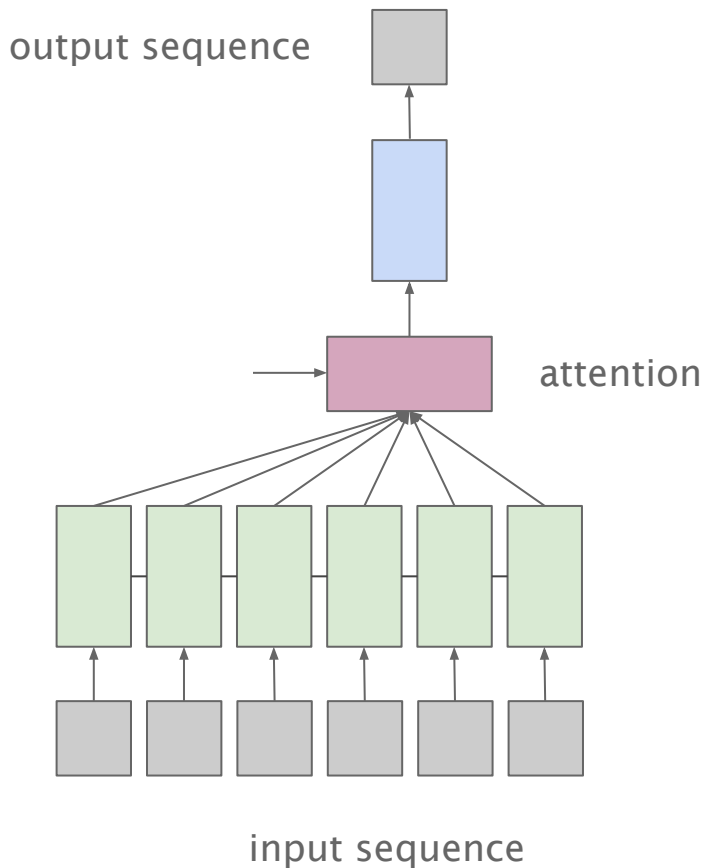
Code.Hub

# Limitations



output sequence

context vector

input sequence

the "context vector" needs to encode all relevant information about the input sequence

Code.Hub

# Attention Mechanism

output sequence

attention

input sequence

Code.Hub

# Attention Mechanism



Code.Hub