

Encoder/Decoder architecture

Two RNN's

- Encoder: takes input sequence \mathbf{x}
- Decoder: creates output sequence $\hat{\mathbf{y}}$

RNN's process sequences using the "loop architecture"

Consider the task of

- constructing the *next* element \hat{y}_{tp} of sequence y
 - conditioned on some input sequence $x = x_{(1)} \dots x_{(t)}$
- $$\Pr(\hat{y}_{tp} | x_{(1)} \dots x_{tp})$$

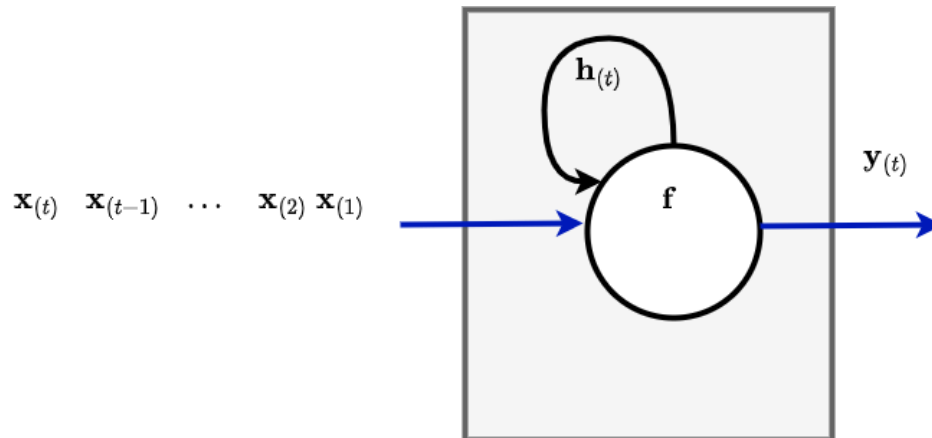
RNN Loop architecture

- Uses a "latent state" that is updated with each element of the sequence, then predict the output

$\text{pr } \mathbf{h}_{tp} | \mathbf{x}_{tp}, \mathbf{h}_{(-1)}$ latent variable \mathbf{h}_{tp} encodes $[\mathbf{x}_{(1)} \dots \mathbf{x}_{tp}]$

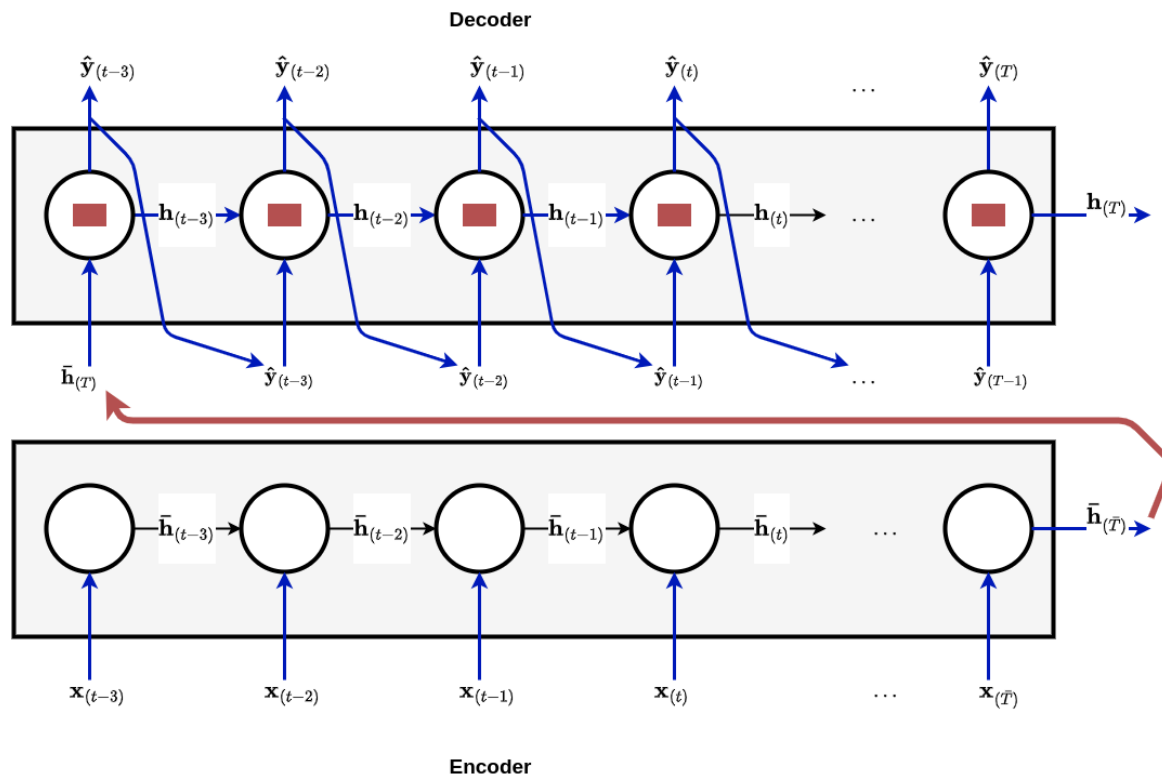
$\text{pr } \hat{\mathbf{y}}_{tp} | \mathbf{h}_{tp}$ prediction contingent on latent variable

Loop with latent state



Original Encoder/Decoder architecture

RNN Encoder/Decoder without Attention
Bottleneck



Critique

- bottleneck
 - *all* information about input $\backslash x$ passes through out of Encoder (red line)
 - and must be carried over to every iteration of the Decoder loop (red box)
- loop architecture for Encoder and Decoder
 - dependency: horizontal line carrying latent state across time

Cross-Attention: removing the bottleneck

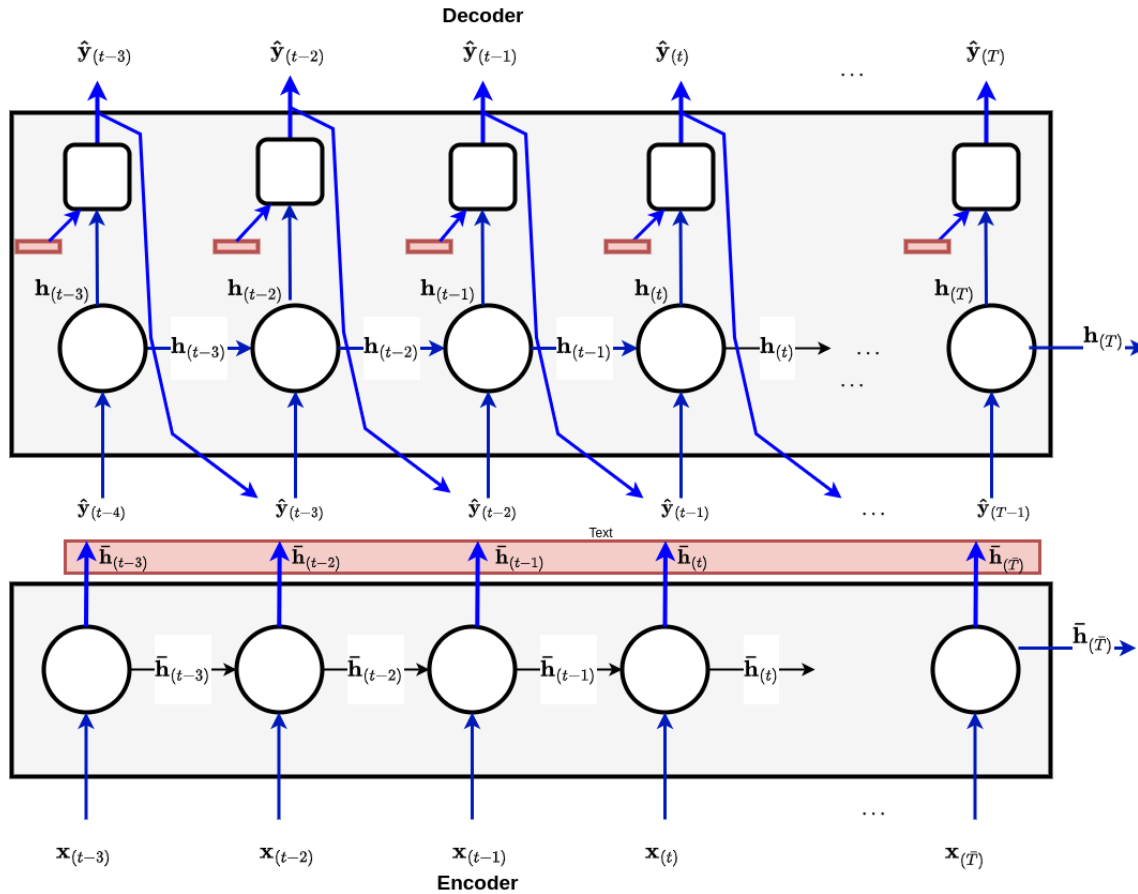
We removed the bottleneck via *Cross Attention*

- Decoder has *direct access* to **all** outputs (i.e., Latent states) of the Encoder
 - each Encoder output is proxy for a prefix of the input

The pink box is the sequent of Encoder outputs

$\bar{\mathbf{h}}_{(1:\bar{T})}$

RNN Encoder/Decoder with Cross Attention



Encoder Self-Attention: removing the Encoder loop

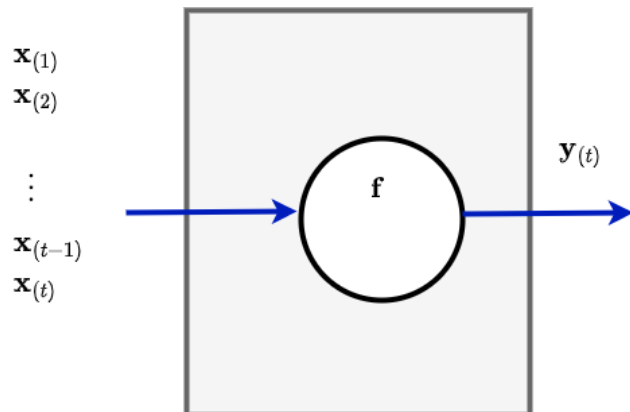
There is an alternative to the loop architecture for processing sequences

- the direct function approach

The alternative to the loop was to create a "direct function"

- Taking a **sequence** $\mathbf{x}_{(1\dots)}$ as input
- Outputting $\hat{\mathbf{y}}_{tp}$

Direct function



Can output *all* elements of sequence $\hat{\mathbf{y}}$ *simultaneously*

- each output position is independent of previous output
- only dependent on input

We removed the "loop" architecture of the Encoder by using the direct function approach

- the mechanism enabling each position of the Encoder output to *attend* to the entire sequence x is called *Self-Attention*
 - Notice: no dependency arrow between circles in the Encoder
- Encoder output is a direct function of **all** positions in the input
 - all Encoder output positions can be computed *in parallel*

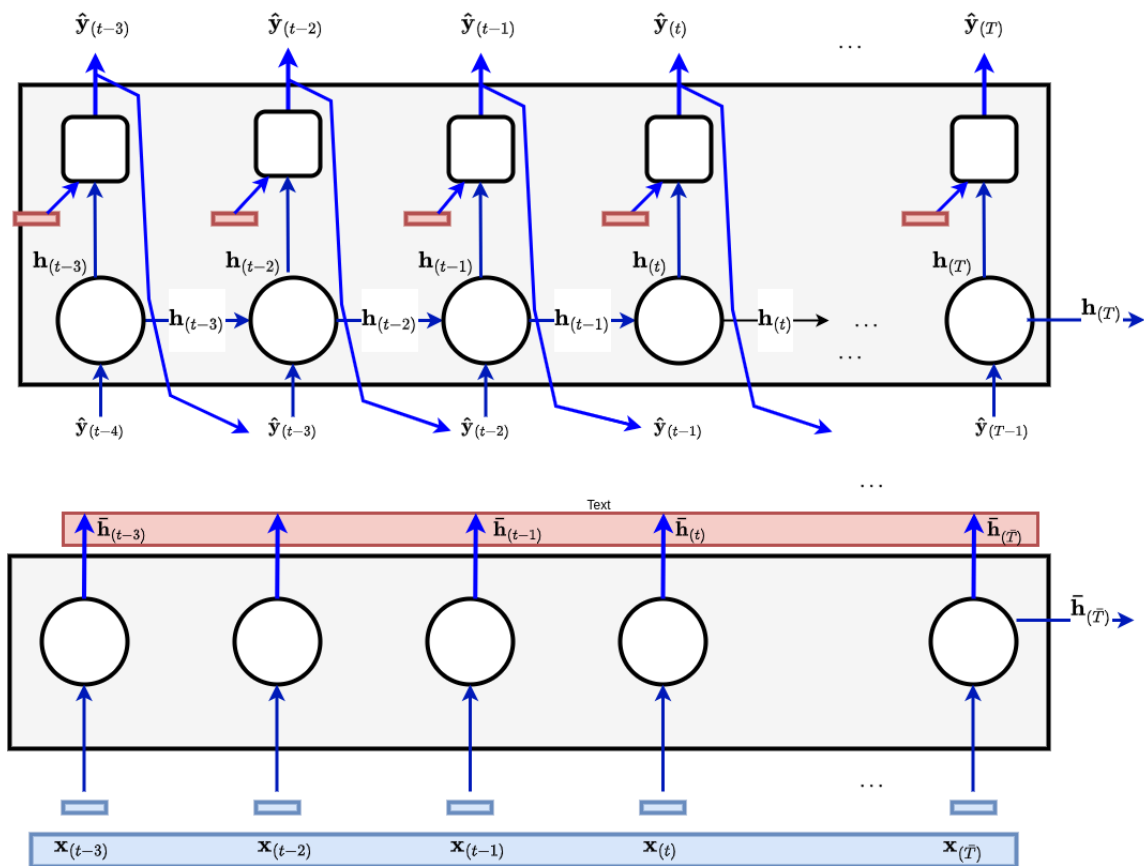
The blue box represents the *entire* input sequence

$\mathbf{x}_{(1:\bar{T})}$

We no longer refer to the Encoder output as a Latent state

- no more loop !

RNN Encoder/Decoder with Cross Attention/Decoder Self Attention



Observe that

- by removing the looping architecture from the Encoder
- the Encoder is no longer called an RNN

Masked Self Attention

With *unmasked* Self Attention

- Encoder output $\bar{\mathbf{h}}_{tp}$ at position
- is a function of **all** inputs $\mathbf{x}_{(1:\bar{T})}$
 - including positions after

This is useful, for example, when the meaning of a word depends on its *entire* context.

- as in our motivating example

For certain tasks (not so for our motivating example), full visibility of all inputs is not permissible

- "looking into the future"
 - e.g., predict stock return based only on **past** information

In this case, we use *masked* Self Attention

- we use a mask to hide inputs from position onwards so that
- output $\bar{\mathbf{h}}_{tp}$ at position
- is a function only of **preceding** inputs $\mathbf{x}_{(1:-1)}$

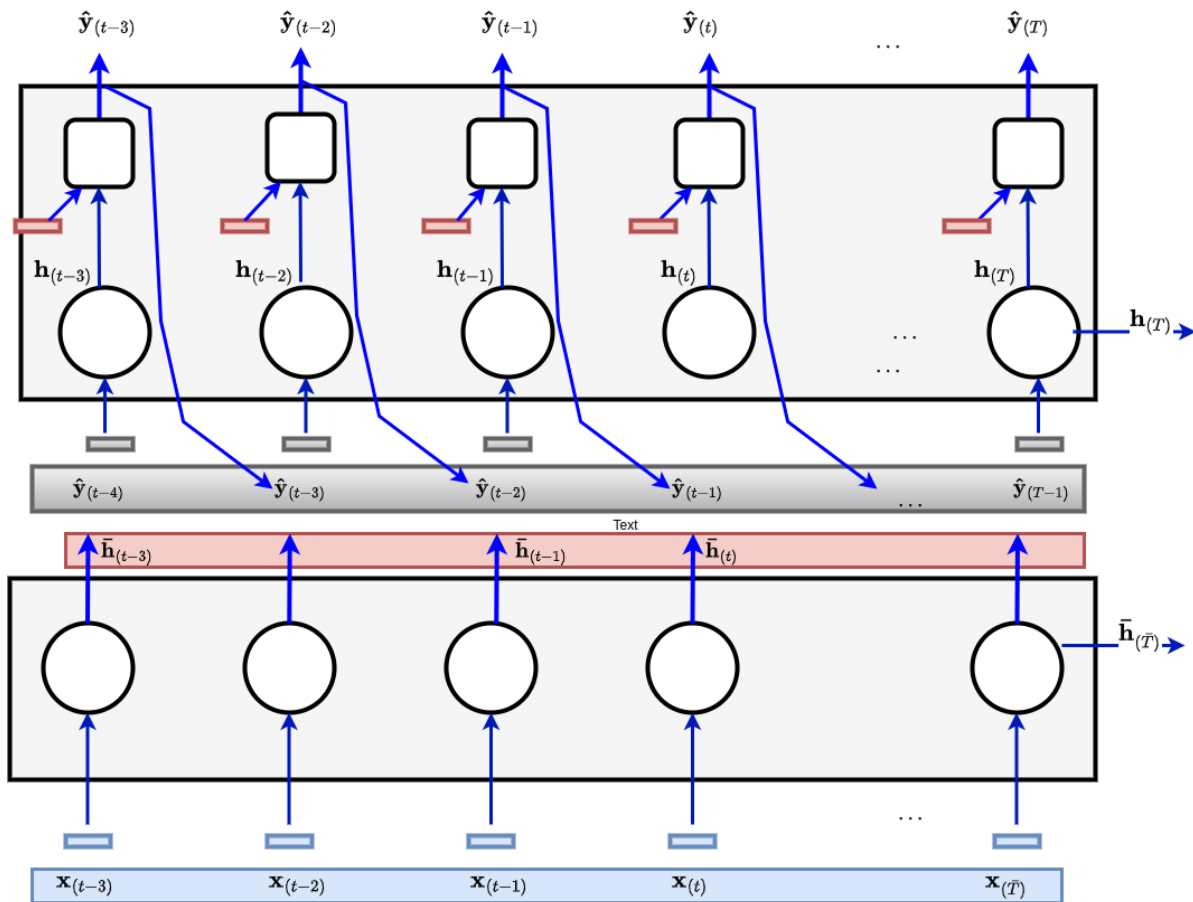
We will see the use of masking in the next section.

Causal Masked Self Attention: removing the Decoder loop

Finally we remove the loop architecture for the Decoder as well using a different "flavor" of Self-Attention

- Masked Self-Attention.

Encoder/Decoder with Cross Attention and Self Attention (Encoder/Decoder)



The grey box represents the *entire* output sequence

$$\hat{\mathbf{y}}_{(1:T)}$$

From this diagram: it appears that

- the Encoder/Decoder can produce output $\hat{\mathbf{y}}_{tp}$
- while attending to outputs *that have not yet been generated* at the start of step
- "looking into the future"

$$\hat{\mathbf{y}}_{(:T)}$$

That is, it is computing

$$\text{prc} \hat{\mathbf{y}}_{tp} \hat{\mathbf{y}}_{(1:T)}$$

What is going on ?

Teacher forcing at training time

An explanation of this strange behavior is that the behavior of the model is *different*

- at training time
- versus at test/inference time

Teacher Forcing alters the training behavior in order to improve the ability of a model to learn.

Let's examine [Teacher Forcing \(Teacher Forcing.ipynb\)](#) in depth.

Masked attention

Hopefully it is clear that, regardless of whether we are computing \hat{y}_{tp}

- at training time
- at inference time

the computation should depend only on positions $1 : -1$ of the output.

- can't peek into the future

To enforce this

- we *mask* the outputs
- so that only positions $1 : -1$ are visible when generating output position

The general mechanism of hiding some inputs is called

- **Masked Self-Attention**

The specific masking of only future positions is called

- **Causal Masked Self-Attention or Causal Self-Attention**

In [2]: `print("Done")`

Done

