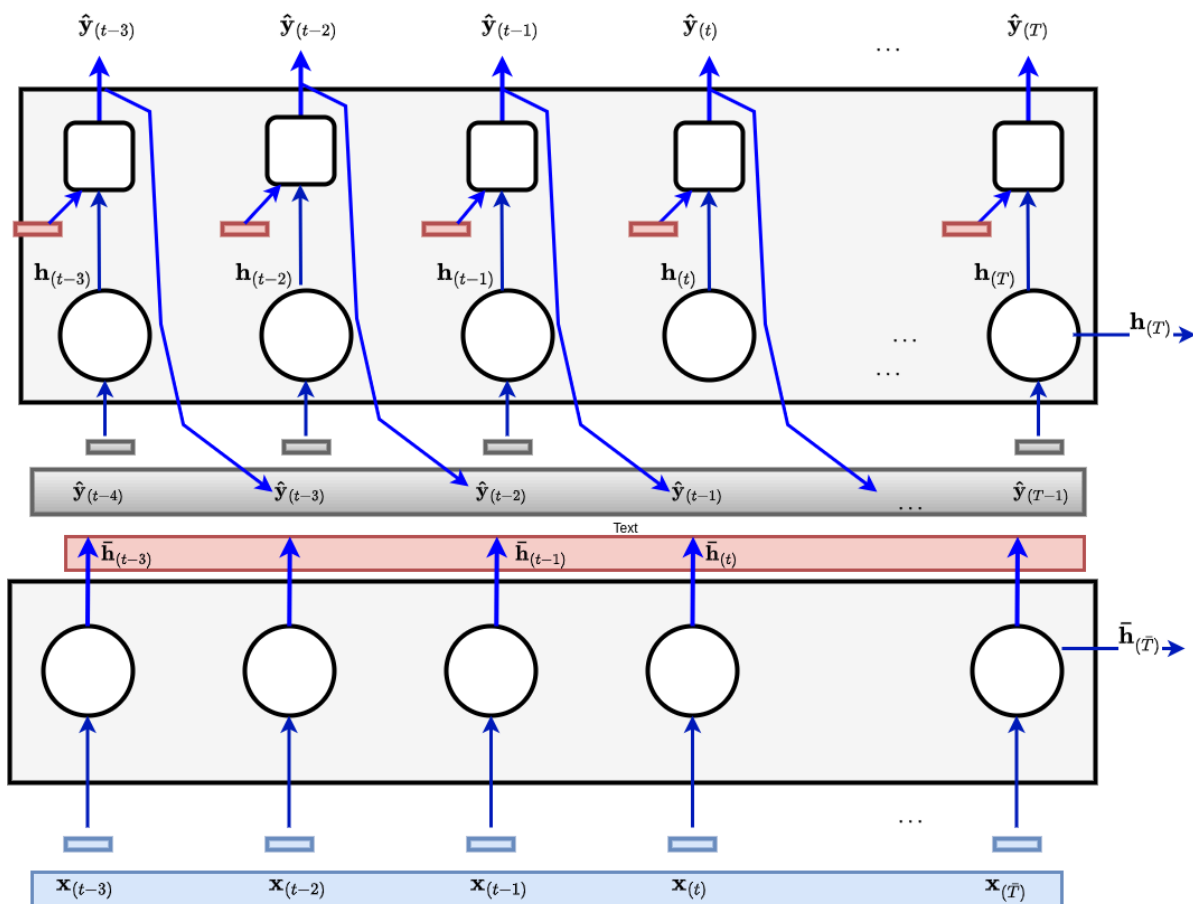# Teacher Forcing

Here is a diagram of the Encoder/Decoder architecture.

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

The grey box represents the *entire* output sequence

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

From this diagram: it appears that

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

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

- "looking into the future"

That is, it is computing

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

What is going on ?

# Teacher forcing at training time

The *Auto-Regressive* Decoder behavior

- constructs output sequence $\hat{\mathbf{y}}$ sequentially

  - output $\hat{\mathbf{y}}_{\backslash \mathbf{tp}}$ of step
  - becomes part of the input $\hat{\mathbf{y}}_{(1:)}$ of step +1

  So how can the entire output sequence $\hat{\mathbf{y}}_{(1:T)}$ be available before final step $T$ ?

We need to distinguish between behavior

- during *training*
- versus during *inference*

During inference, clearly we can only compute
$$\prc\hat{y}_{\tp}\hat{y}_{(1:-1)}$$

because we haven't generated the future outputs yet.

But at training time, example is

$$\langle \backslash \mathbf{y}_{(1:-1)}, \backslash \mathbf{y}_{\backslash tp} \rangle$$

We predict *only the immediate next* target $\backslash \mathbf{y}_{\backslash tp}$

- *not* the full suffix $\backslash \hat{\mathbf{y}}_{(:T)}$

The prediction is conditioned on the *true* target prefix $\backslash \mathbf{y}_{(1:-1)}$

- *not* the prefix generated during training $\backslash \hat{\mathbf{y}}_{(1:-1)}$

The Auto-Regressive behavior is eliminated during training !

Training in this manner has a big advantage.

In a a perfect world, when predicting $\backslash \hat{\mathbf{y}}_{\backslash \mathrm{tp}}$

- the Auto-Regressive behavior during training would result in
- the prefix $\backslash \hat{\mathbf{y}}_{(1:-1)}$ of the generated output
- matching the true target

$$\backslash \hat{\mathbf{y}}_{(1:-1)} = \backslash \mathbf{y}_{(1:-1)}$$

But

- if any element of the *generated* prefix is wrong
$$\hat{\mathbf{y}}_{(t')} \neq \mathbf{y}_{(t')} \text{ for } t' < t$$
- it is likely that all *subsequent predicted outputs* $\hat{\mathbf{y}}_{(t'+1:T)}$ will be wrong

- because each subsequent output **is conditioned on incorrect** $\hat{\mathbf{y}}_{(t')}$

That is

- a single mis-predicted element of the sequence
- causes a catastrophic chain of errors

Training a model under such conditions would be difficult

- the incorrect sequences don't even come from the true distribution of inputs !
- violating the Fundamental Theorem of Machine Learning

To avoid this, our training examples compute

$$\prc{\hat{y}_{\tp}}{y_{(1:-1)}}$$

rather than

$$\prc{\hat{y}_{\tp}}{\hat{y}_{(1:-1)}}$$

That is: we train on *target* prefixes rather than train-time generated prefixes.

This is called *Teacher Forcing*.

Teacher forcing can be implemented

- by making the entire target output sequence
$$\textbf{y}_{(1:T)}$$
i.e., the grey box in the diagram

- available at training time via setting example to
$$\left\langle \textbf{y}_{(1:T)}, \textbf{y}_{\backslash tp} \right\rangle$$

- and using *Causal masking*
- to make only prefix $\textbf{y}_{(1:-1)}$ visible during the prediction of $\textbf{y}_{\backslash tp}$

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

Done