

Attention and Transformers

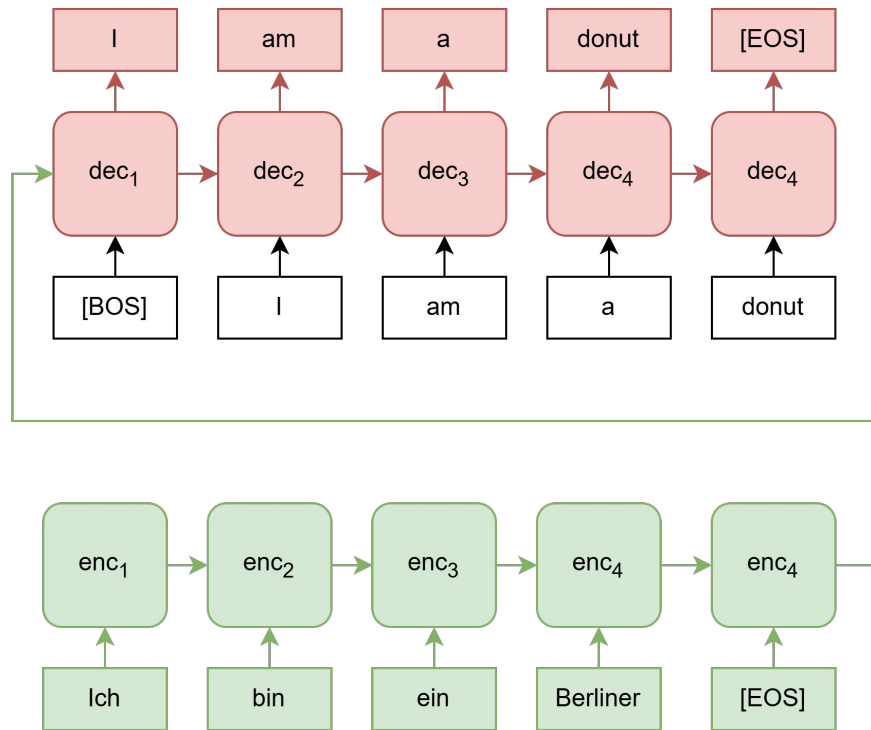
Lukas Edman

Layout

1. RNNs recap
2. RNNs with attention
3. Transformer

RNNs Recap

- Each token gets fed into the RNN one-by-one
- Hidden state is accumulated
- BOS token (or SEP token, whatever) initiates the decoding
- Given the hidden state and current token, output probabilities for the next token
- Take the argmax (if you're doing greedy search), or topk (if you're doing beam search)
- Repeat until reaching EOS



RNNs Recap

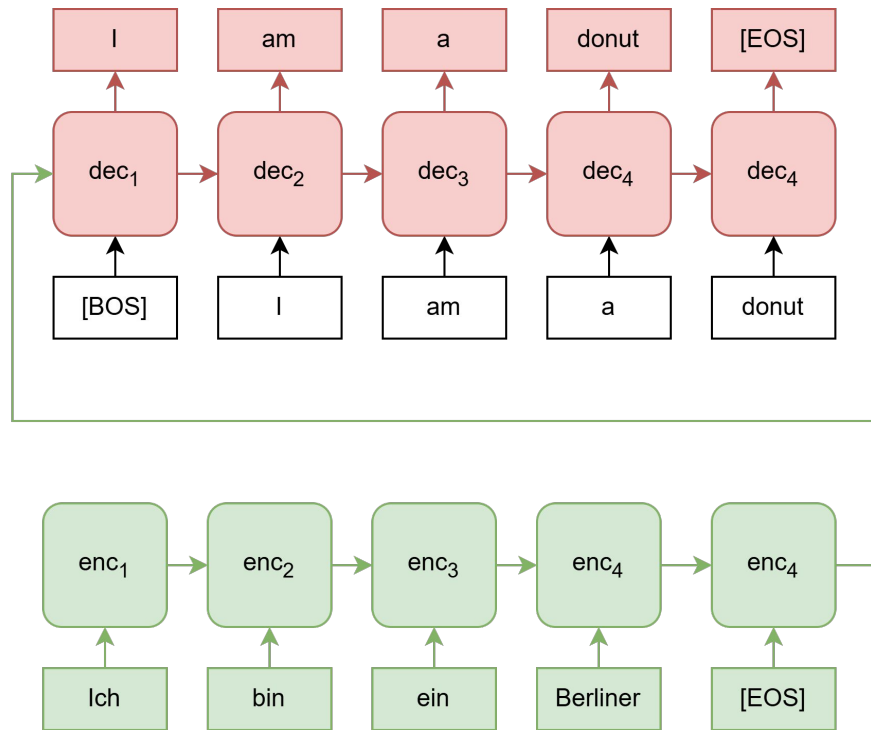
Training

- Get output
- Apply softmax
- Cross Entropy Loss:

$$\text{Loss} = - \sum_{i=1}^{\text{output size}} y_i \cdot \log \hat{y}_i$$

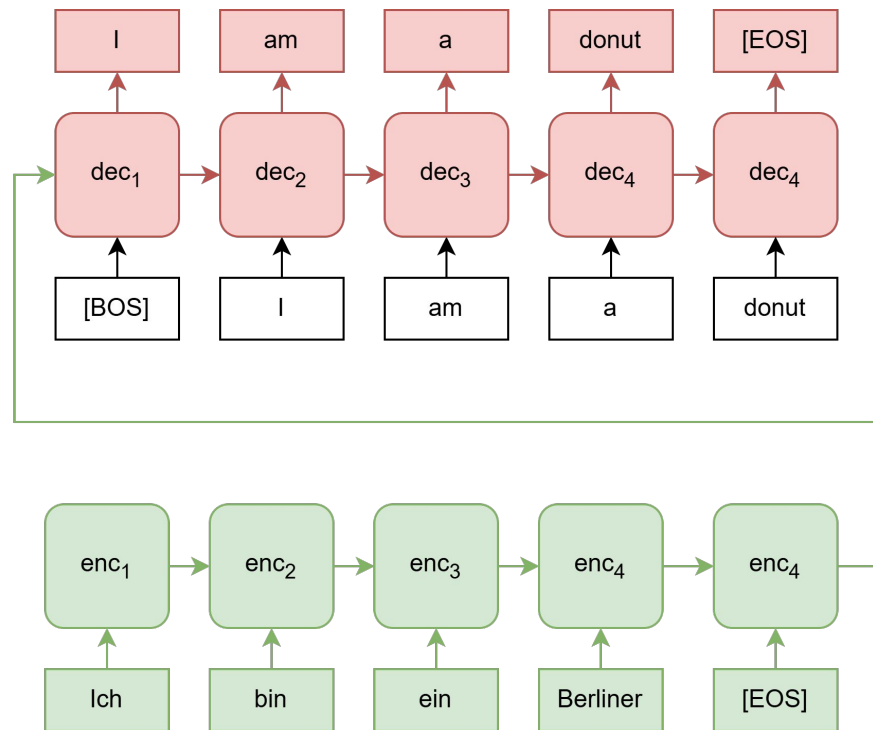
where

- \hat{y}_i = probability of token i
 - y_i = whether y_i is the correct token (1=yes, 0=no)
 - output size = vocabulary size
-
- Use teacher-forcing, e.g. if the model outputs "We" instead of "I", still train it as though it output "I"



Problems with Vanilla RNNs

- What's the problem with this structure?
- Why might this structure make it difficult for an RNN to perform well?



Problems with Vanilla RNNs

- It works well with short passages, but try feeding in:

Heute war ich in Zug (DB).

Es ist da was passiert, Geschichten die das Leben so schreibt. Sowas kann man sich nicht ausdenkeng.

1 alter Oma fuhr mit mir in ICE Zug DB von München nach Dresden, weiss selber nicht wie man darauf kommt durchzufahren, da ich Nürnberg raus wollte von noch in Bayern bleiben her. (Ich überquere die Freistaatsgrenzen nur selten, bin auch nicht geimpft etc). Vielleicht wollte sie dort Verwandte besuchen und bei Flucht helfen in Westen? KP aber es geht mir auch nicht an. Wie wir da so fahren, kommt 1 Schaffner im Sinne von Kontrolle. Jeder packt Smartphone aus, 1 connected Laptop mit Schaffnerkontrollgerät an Netzwerk her ich zeige mein Smartwatch (Appel) für QR Code zum scannen. Omer sieht man schon an dass sie normales Papierfahrkarte hat, alle in Abteil sind schon am sie ausliosen weil so rückschrittlich, passt einfach nicht mehr in die Zeit denken die wahrscheinlich. Ja und bissl recht haben ja auch.

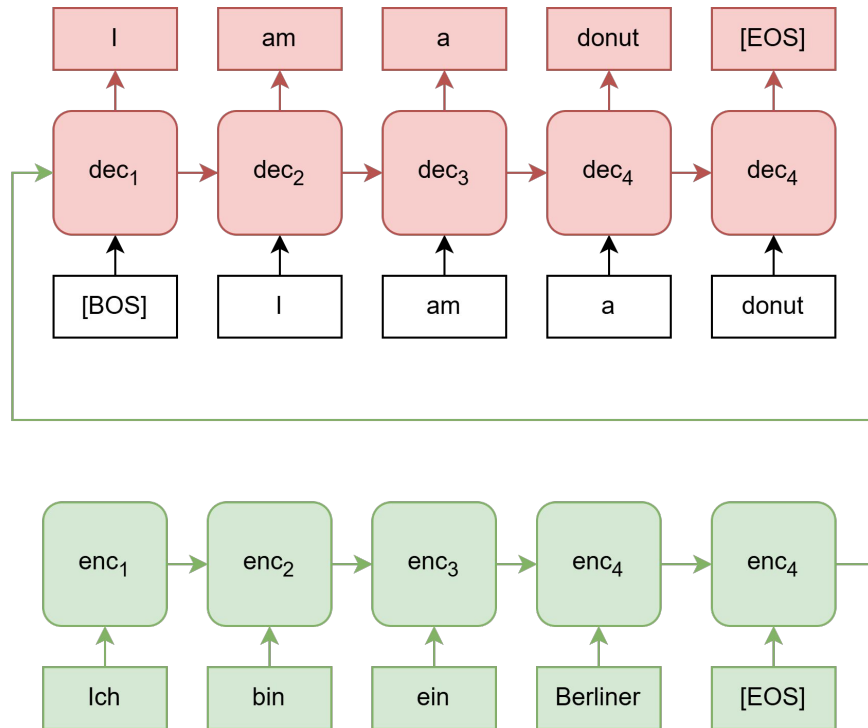
Aber langer Rede gar kein Sinn, es stellt sich raus die Omer hat auch noch 1 Ticket sich gekauft dass nicht im ICE Zug giltet. Dies hätte sie im Internet nachlesen können, aber das hat sie nicht weil es ist zu teuer und den VHS Kurs wo es erklärt wird hat nicht leisten können. Da sind die Enkel gefragt meiner Meinung nach aber hilft jetzt auch nix mehr. Der Schaffner besteht auf sein Geld, das die alte Frau zahlen muss in Sinne von Schwarzfahren her es sind 120 Euro. Sie hätte nur fahren dürfen in Regionalbahn wo es 5 Tage dauert bis man Augsburg ist lol (von München aus wohlgemerkt). Die Frau sagt hat kein Geld und sie muss vom Rente leben von ihrem verstirbten Mann und es ist wenig. Aber der Jockel von Schaffner lässt sie nicht erweichen. Er bleibt hart. Alle schauen verlegen auf ihrem Handy oder lesen Börsenkurse. Ich stehe auf und erhebe Wort. "Hör mal zu du Überjochen, die Frau fährt ja wohl jetzt unsonst, sonst haben wir hier rukizuki Rambazamber und es gibt 1 Bombe" sage ich den Schnauzbarträger im Gesicht. Es wird noch leiser in Abteil, was gar nicht möglich ist weil ja vorher schon so leise war. Er sagt kann er nicht machen wegen Privatisierung von Bahn früher schon da hätte es Steuerzahler gezahlt aber jetzt nicht und ihm sind auch die Finger gebunden. Ich erwiedere: "Schau mal her du Lauchkönig, ich zahl jetzt die Hälfte von dem Schwarzfahren aber dann ist auch gut, den Rest zahlt die DB der Knechtzirkus" Jetzt stehen Leute auf und wollen Schlegerei anfangen von sozialer Ungerechtigkeit her. Sie sind entzürmt weil die Bahn so 1 Geldverein ist. Einer rollt das DB Kundenmagazin gans fest zusammen, für dass er ordentlich zubatschen kann. Der Schaffner ist in der Unterzahl (logisch). Allgemeine Berlinstimmung jetzt in Abteil, Kaffeebecher und Kundenmagazine fliegen durch die Luft (Rigaerstrasse artig)

Die Omer ist das peinlich, sie möchte das nicht das soviel Turbul um sie gemacht wird. Ich sage "Fresse jetzt Hexengesicht, es ist 1 Sache von Ehre jetzt. Dem Opferkönig sein Zahnbürste greift morgen ins Leere wenn es jetzt weitergeht und ich aufdrehe von Fausttanz her" Gegröhle in Wagen. 1 mitreissender Arzt hat sein Koffer aufgemacht und 1 Flasche Chloroform hervorgeholt für dass Betaubung da wäre wenn man bräuchte in Kampfgetümmel.

Jetzt gibt der Schaffner klein bei. Er sagt es passt so und er würde Omer nicht mehr belegen. Ich gebe ihm noch die 60€ und klatsche ihm auf die Stirn wie man das mit schlechter Schüler macht. Unter hemischen Rufen wie "Du Protojockel!" und "Hau ab du Vollgasotto!" verlässt er das Abteil. Die Omer dankt allem und sätzt sich verlegen hin. Es war viel für sie. Von Aufregung her.

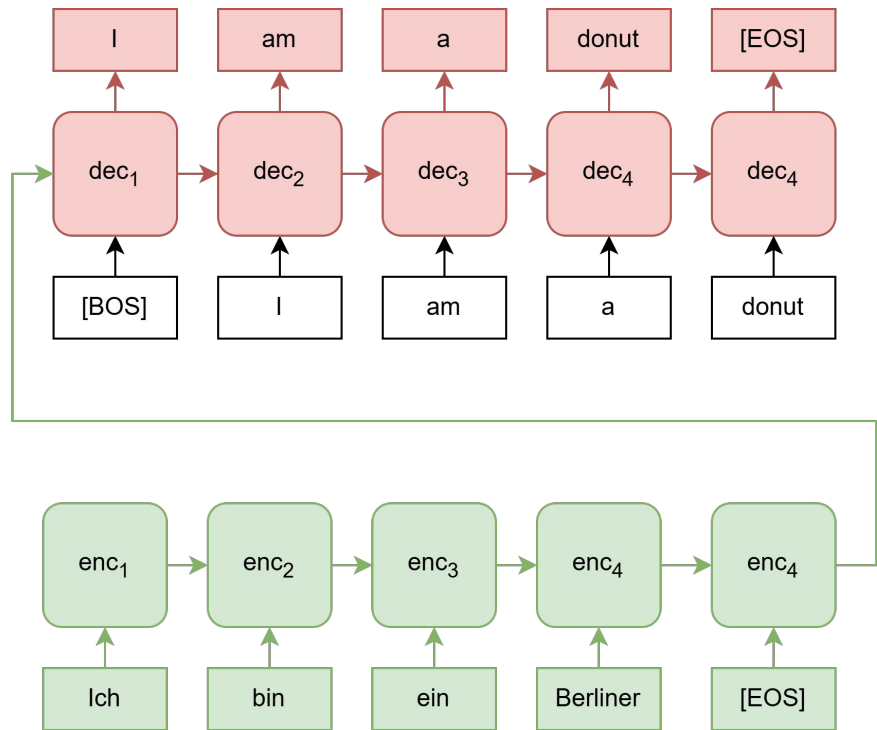
Die Leute wollen jetzt auch mir Geld zustecken weil sie Aktion gut fanden glaub ich. Ein modischer Geschäftsmann aus Bad Tölz steckt mir beileufig 1 Hunni zu und meint verschmitzt "Ich kann mir gönnen ich hab 1 Startup von vsganer Käse es boomt" An Ende habe ich 480 Euro plus. Da sieht man dass es sich auszahlt wenn man Solidarität und Zivilcourtag zeigt. Alleine sind wir schwach, gemeinsam sind wir mehr! #zivilcourtag #alleineindwirschwach #ottonhaftigkeitabschaffen

- Then it won't work so well.



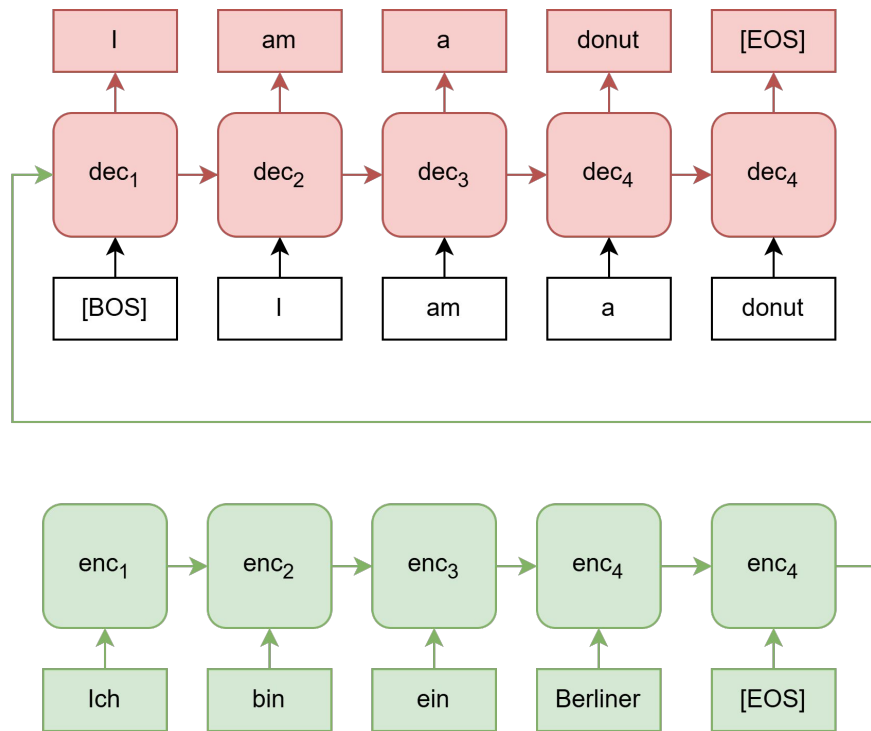
Problems with Vanilla RNNs

- There is a bottleneck
- The hidden state has to accumulate all of the input, so it becomes harder to accumulate everything perfectly (losslessly) the longer the text gets



Problems with Vanilla RNNs

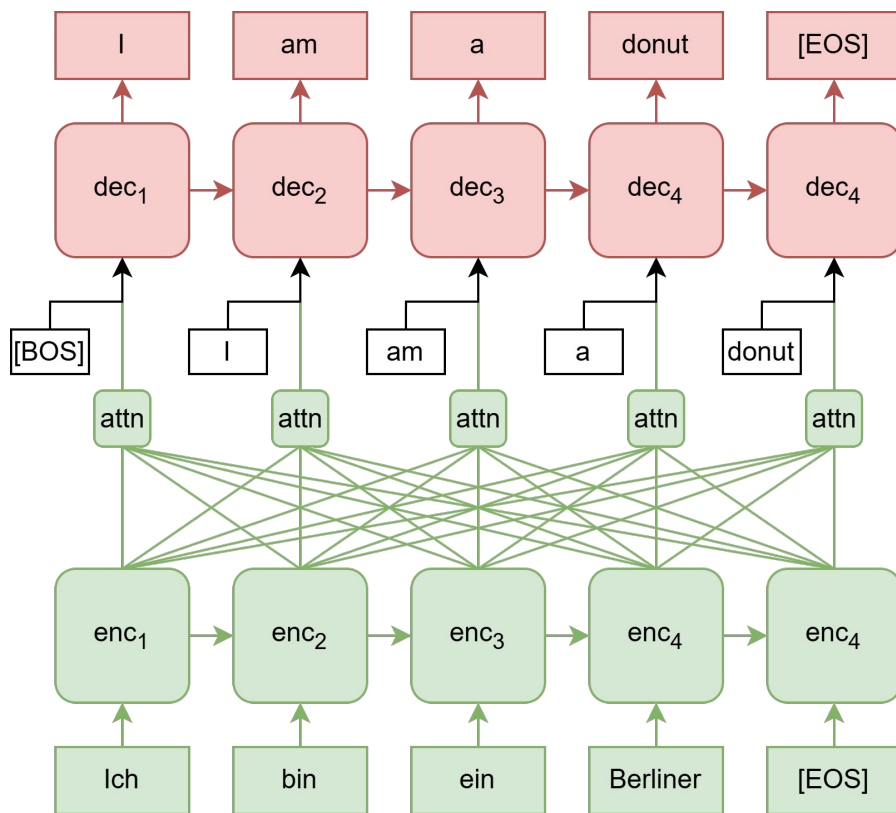
- There is a bottleneck
- The hidden state has to accumulate all of the input, so it becomes harder to accumulate everything perfectly (losslessly) the longer the text gets
- So what's the solution?



RNNs with Attention

- Attention lets the model see everything from the input altogether
- All of the input states get aggregated and are fed in to the RNN decoder at each step
- The additional connections remove the bottleneck

But what's the "attn" block doing?

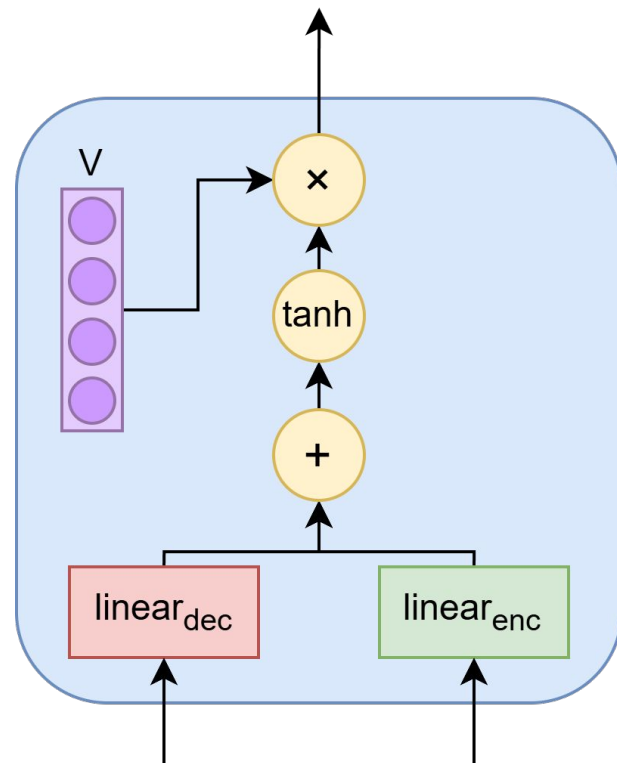


Attention Block

- First we calculate an attention score $e_{t,i}$ for each input i , at timestep t

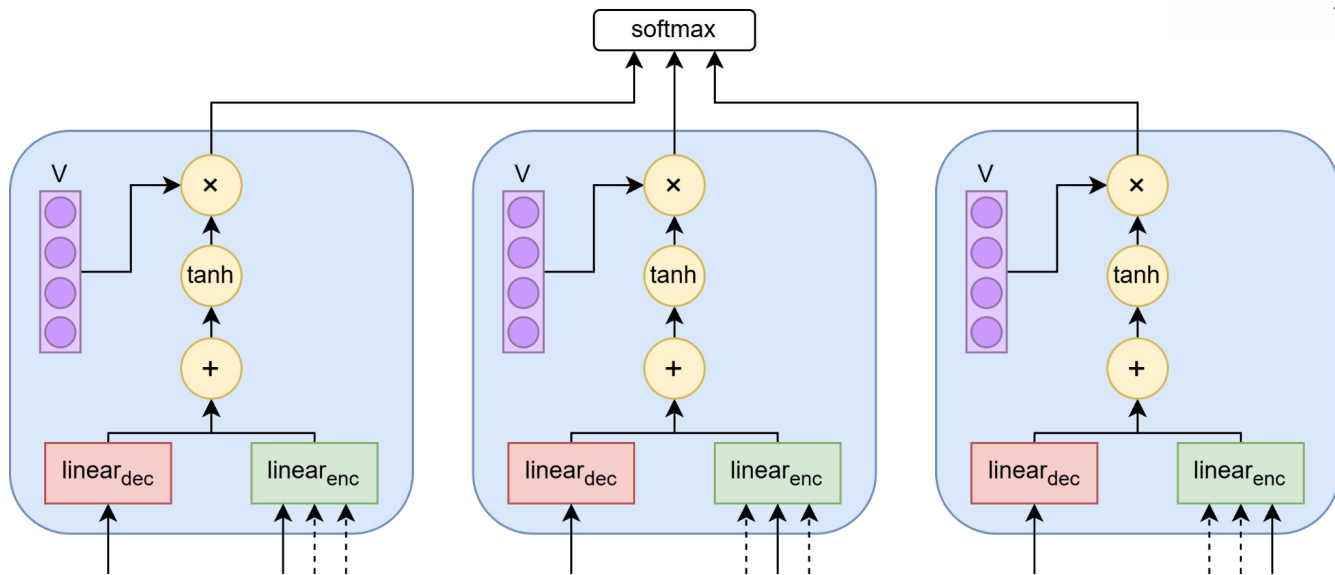
$$\begin{aligned} e_{t,i} &= \text{score}(h_t^{dec}, h_i^{enc}) \\ &= v^T \tanh(W_1 h_t^{dec} + W_2 h_i^{enc}) \end{aligned}$$

- Repeat this for every token in the input (i in $(0, N)$)



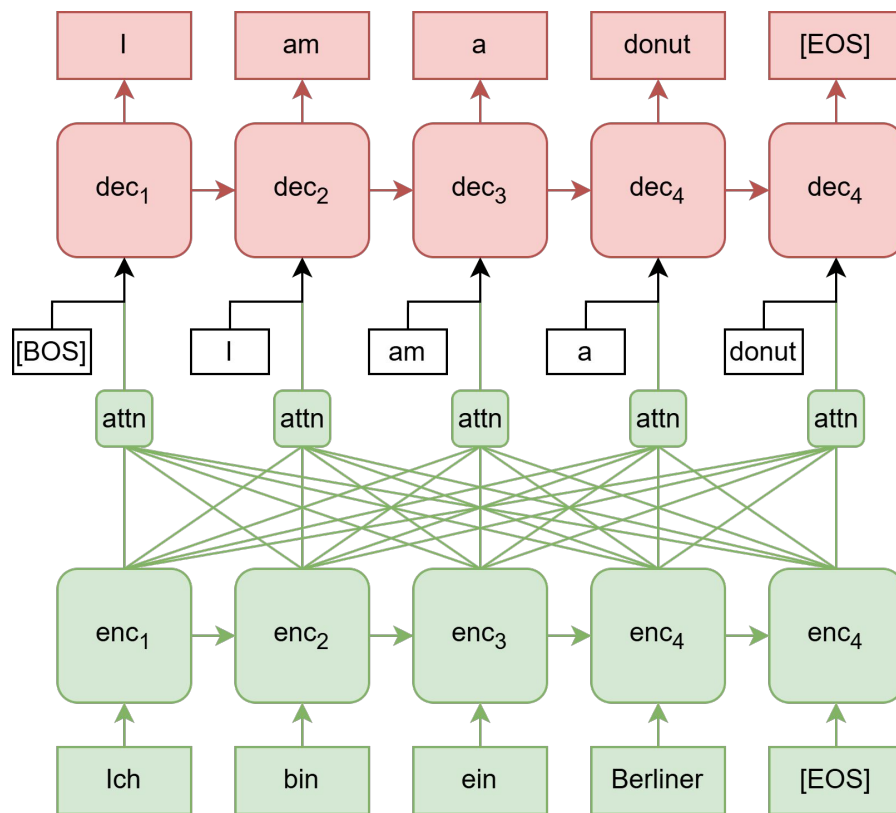
Attention Block

- Take the softmax of each encoding: $\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_j \exp(e_{t,j})}$
- Use the attention scores as weights for a weighted average: $c_t = \sum_i \alpha_{t,i} h_i^{enc}$



RNNs with Attention

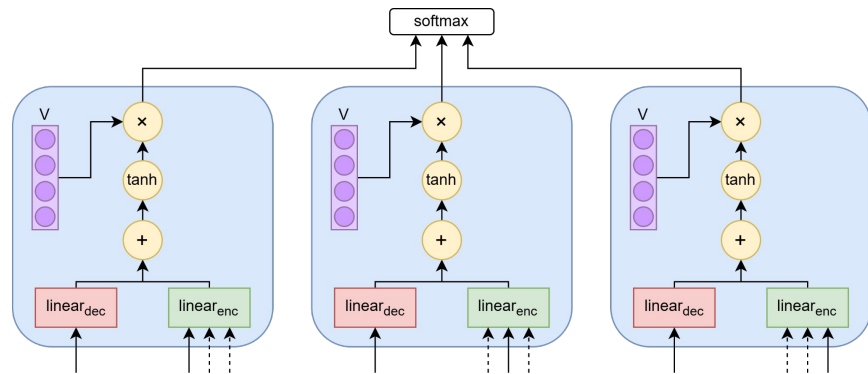
- RNNs with attention work a lot better, especially for longer sequences
- But we can do even better...



Transformers

Is Attention All We Need?

- RNNs are pretty slow, given that they have to process each token sequentially
- Hold on, this attention thing just let us aggregate all of the input side in parallel steps...
- Why don't we make the whole architecture parallel and save massive amounts of time?



We can do these 3 all at the same time!

Attention /S All We Need!

- Model architecture based on attention
- Multiple layers of attention for lots of contextualization
- Every token can be processed in parallel in training
- Generation must be done sequentially, but only for the target side

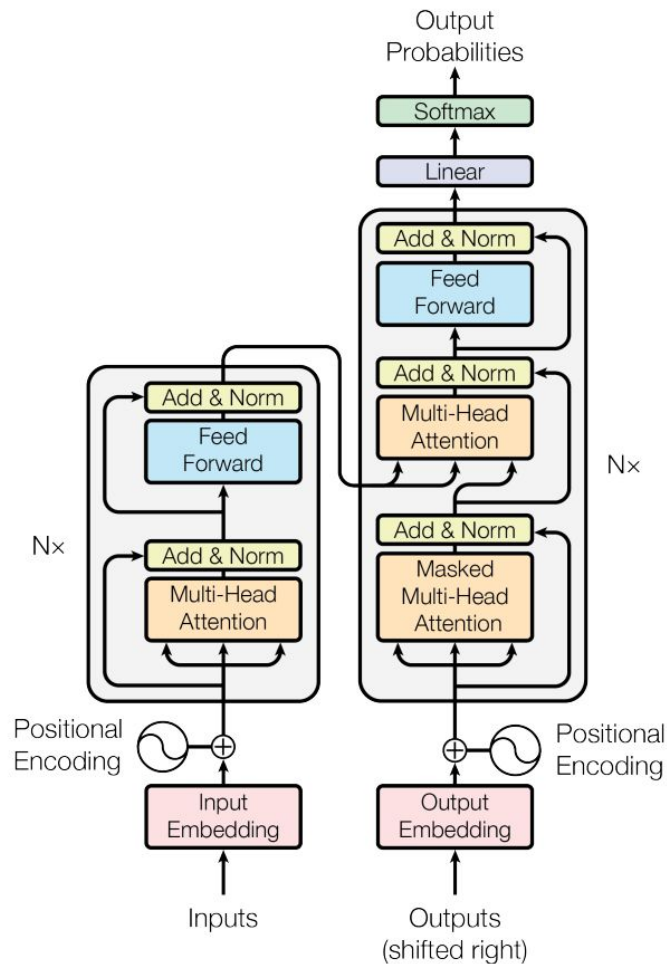


Figure 1: The Transformer - model architecture.

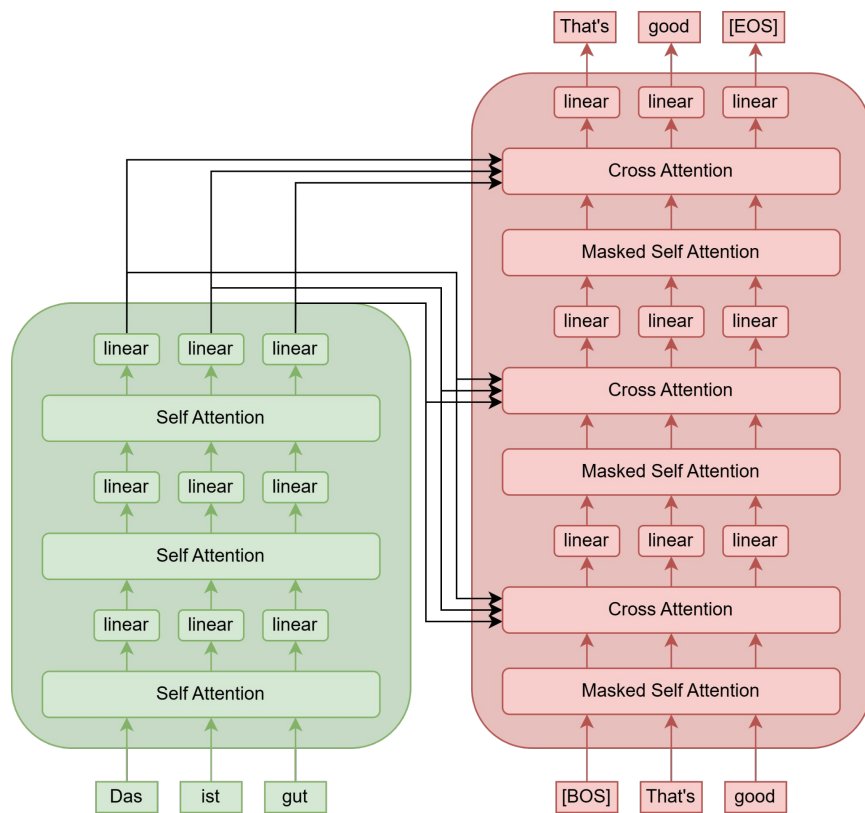
Transformer Architecture

The architecture has 3 main components:

1. Self attention
2. Cross attention
3. Linear layers

The encoder is bi-directional (or more accurately, omni-directional)

The decoder is still uni-directional (left to right)



Self attention

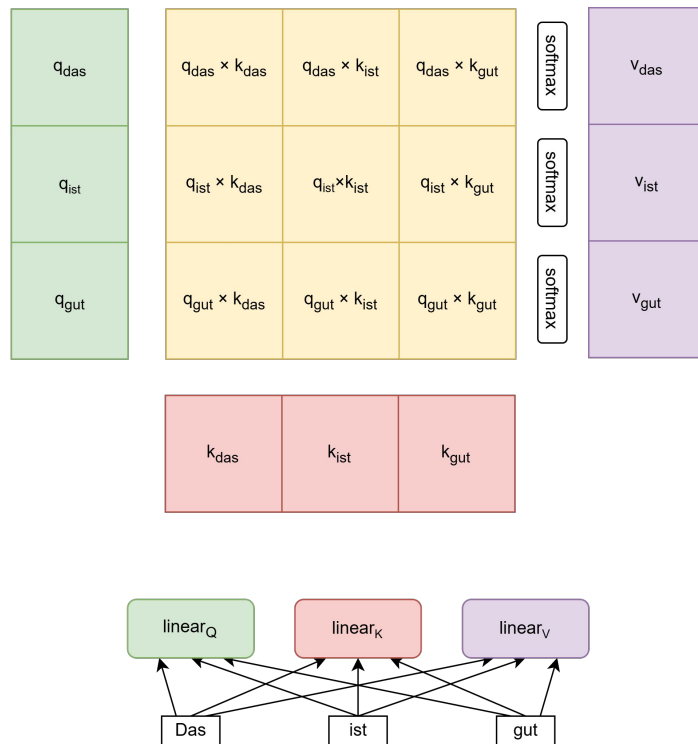
Very similar to the attention we saw earlier

Get scores (e): $e_i = \frac{q \cdot k_i}{\sqrt{d_k}}$

Softmax: $\alpha_i = \frac{\exp(e_i)}{\sum_j \exp(e_j)}$

Multiply by value: $c = \sum_i \alpha_i v_i$

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$



Masked Self attention

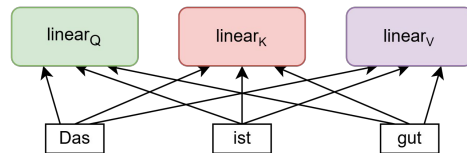
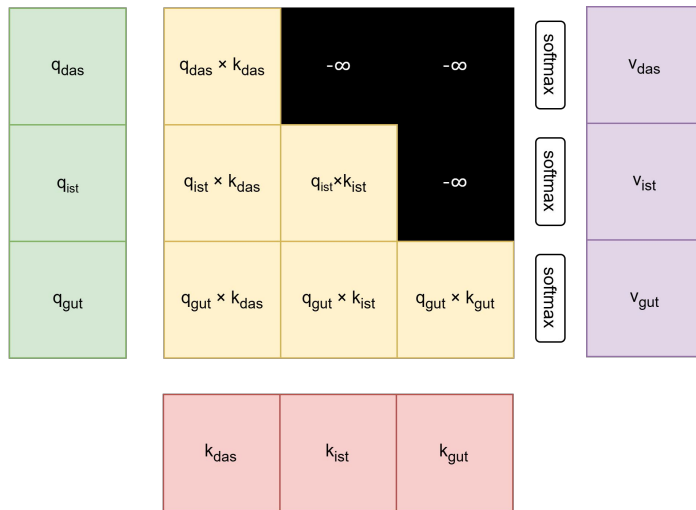
In the decoder, we can't look ahead, that would be cheating!

Mask QK^T such that $q_i k_j$ is $-\infty$

whenever $j > i$

Because this is the only time tokens are mixed, the model cannot look ahead!

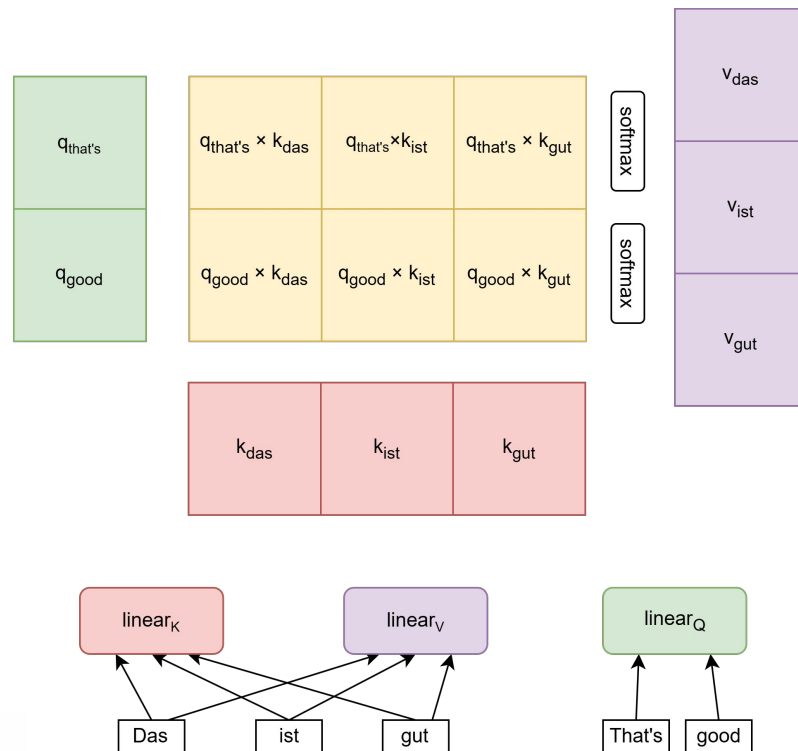
$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$



Cross attention

Same as self attention between source and target, except that

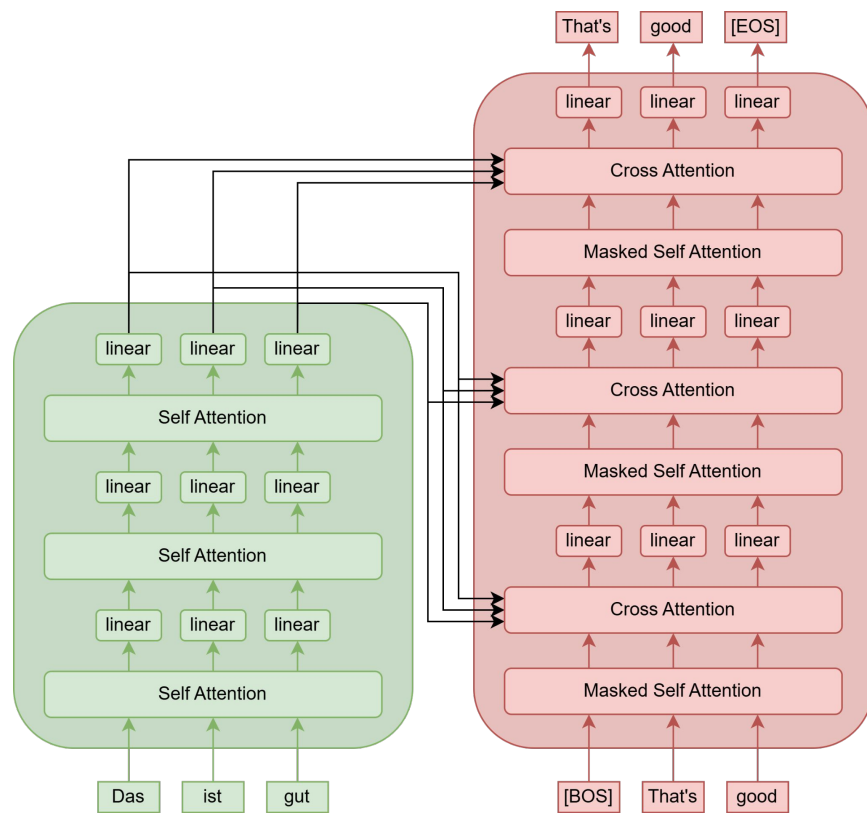
- K and V come from the **source** side
- Q comes from the **target** side



$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Transformer Architecture

1. Self attention – Mixes tokens in encoder or decoder
2. Cross attention – Mixes tokens from encoder into the decoder
3. Linear layers – Transform the representation of each token *independently*



Transformer Architecture

More minor things:

- Multi-head attention – Splits the hidden dim into smaller groups, computes attention on them independently
- Positional encoding – adds another embedding that encodes the position of each element.
- Layer Normalization – helps stabilize training.

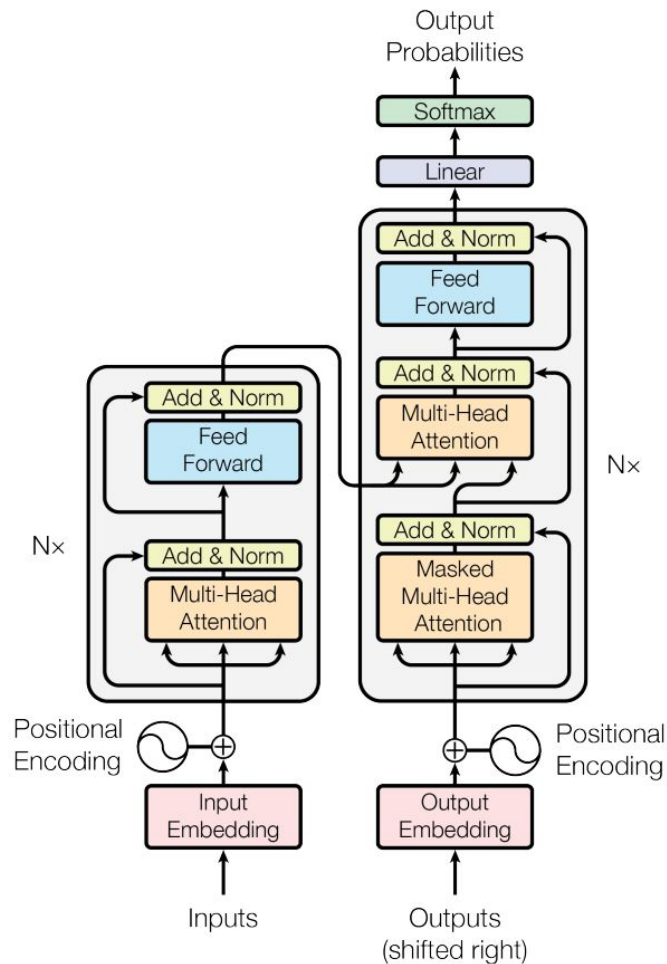
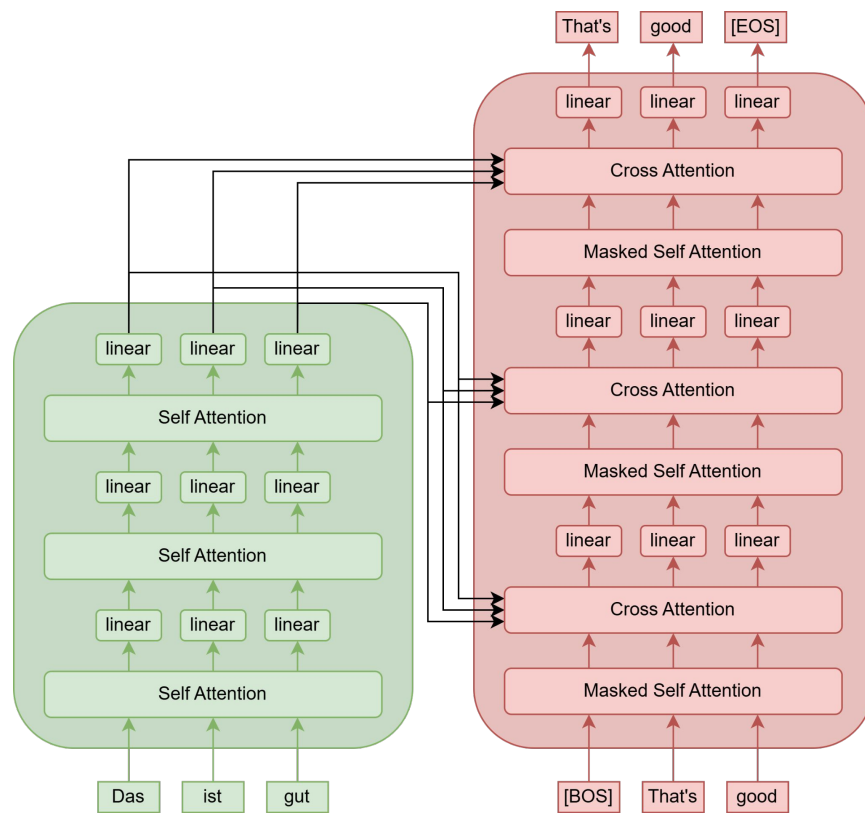


Figure 1: The Transformer - model architecture.

Transformer Training

- Training in parallel → much faster
- General setup the same as RNNs:
 - Outputs probabilities
 - Uses cross entropy loss



Attention in Transformers

You can probe attention in a Transformer and see what it's looking at

For a good translation model, you'll see something like this:

