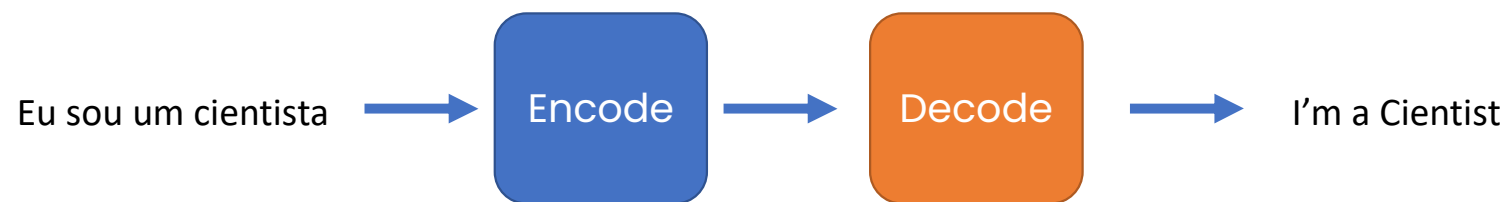


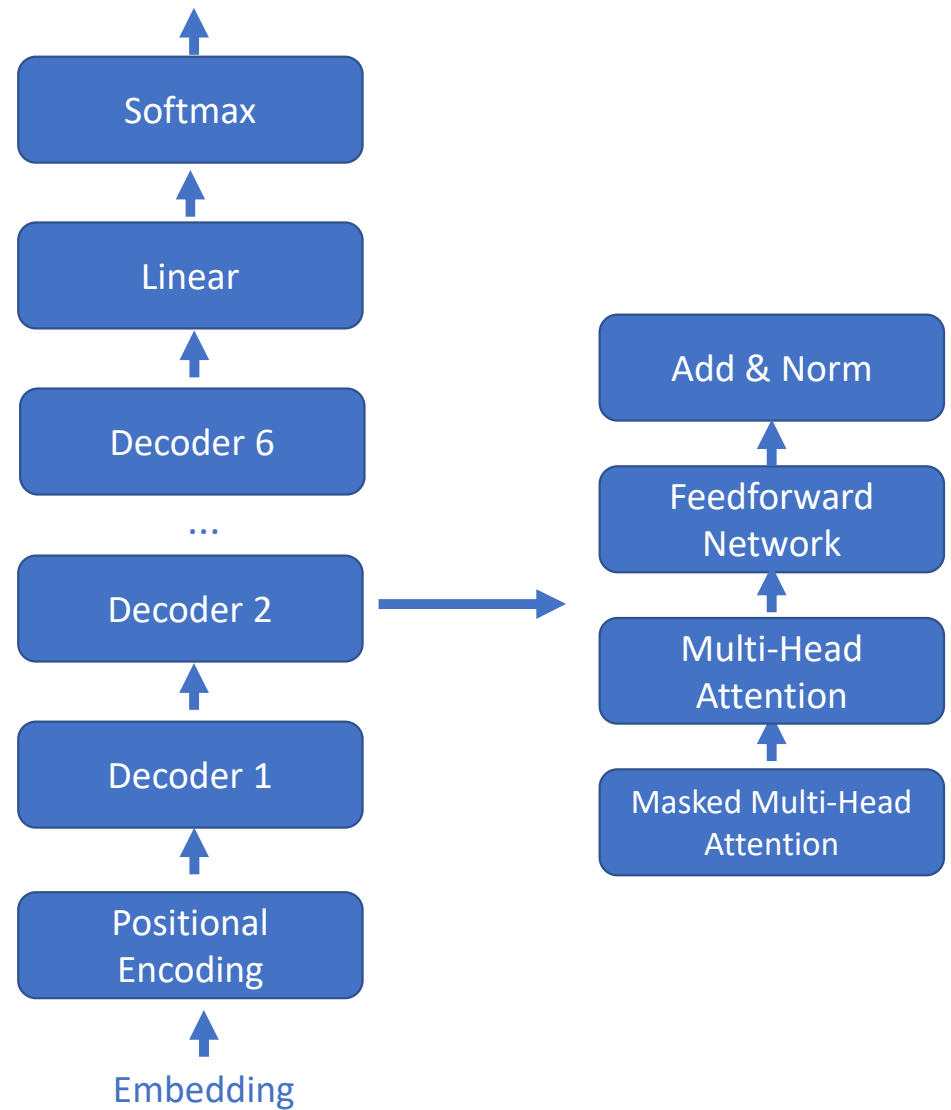
# Decoder



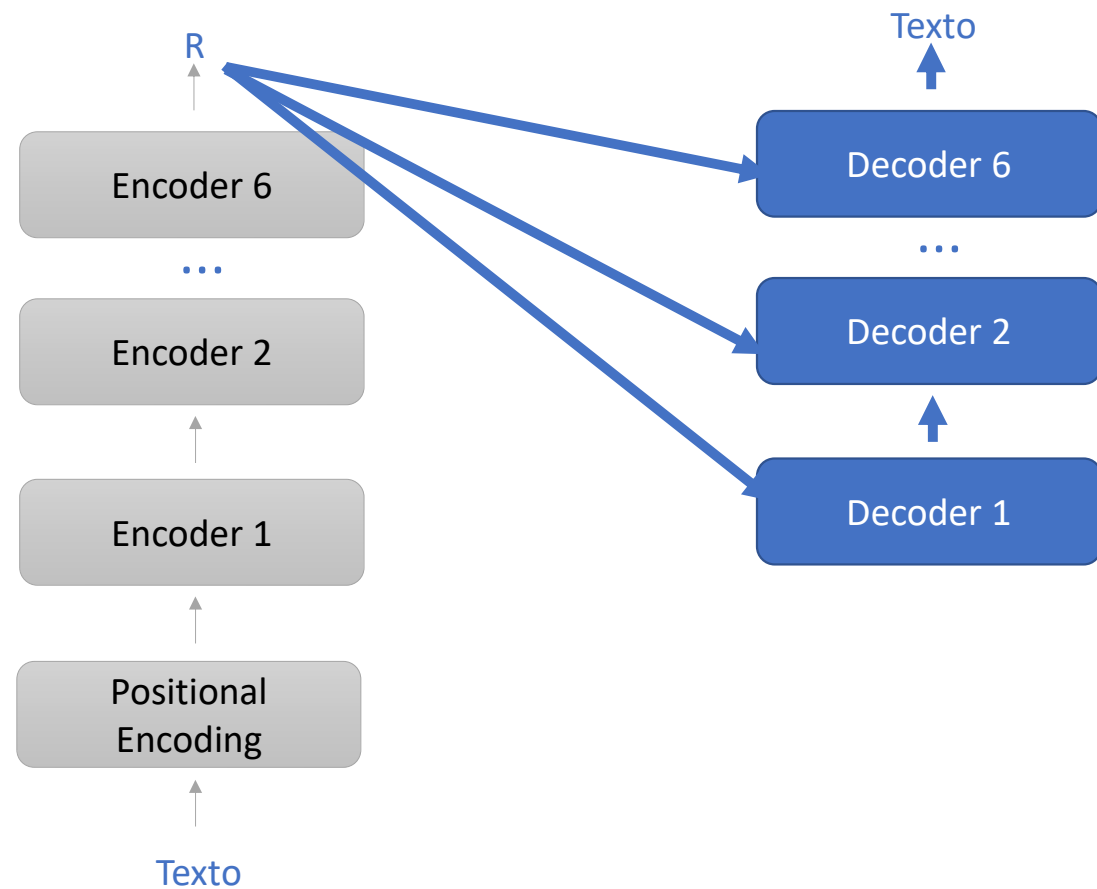
- Decoder

Também são camadas empilhadas

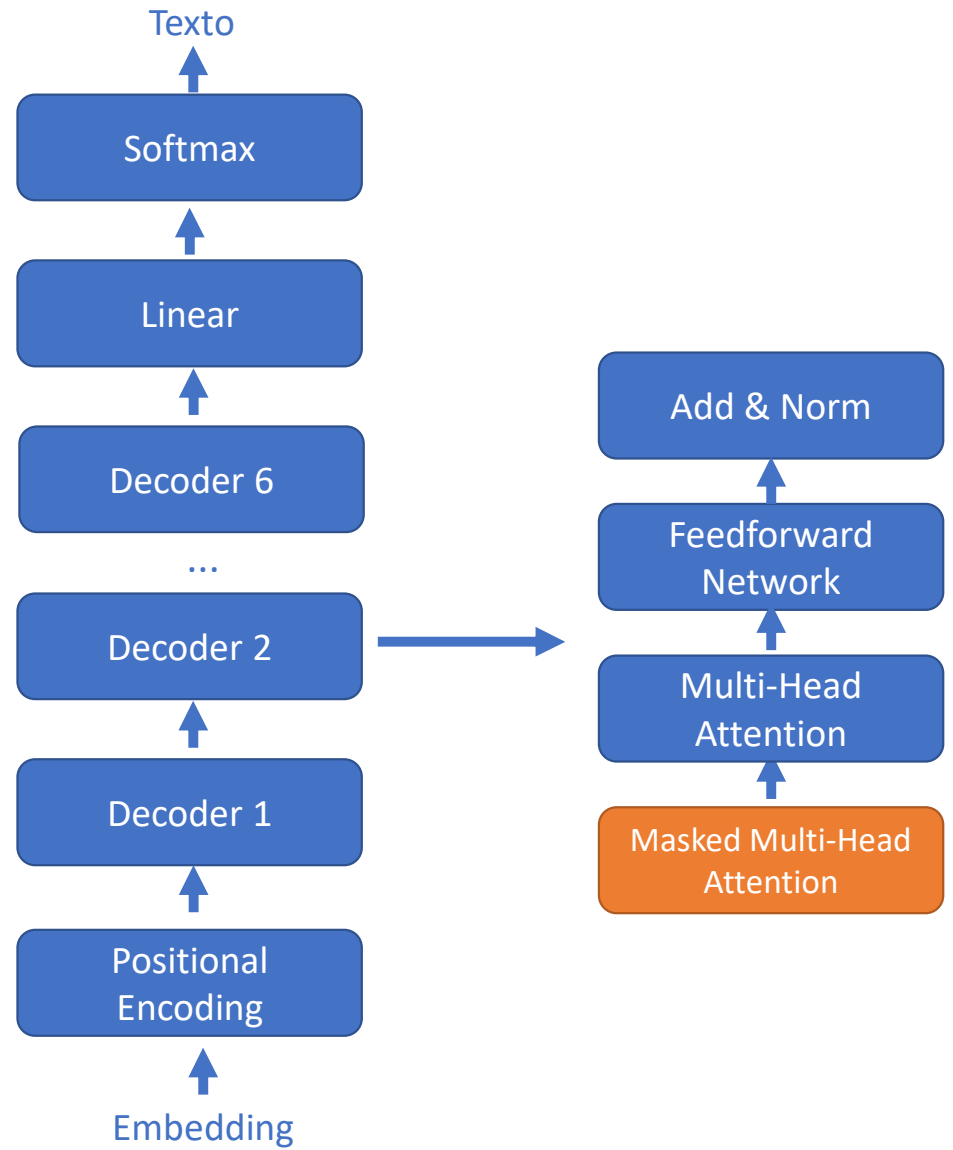
Também ocorre a conversão para embedding e positional encoding



- Cada Decoder recebe 2 entradas:
  - decoder anterior
  - a saída do encoder
- A conexão é com o elemento Multi-head attention do decoder



# Decoder



# Masked Multi-Head Attention

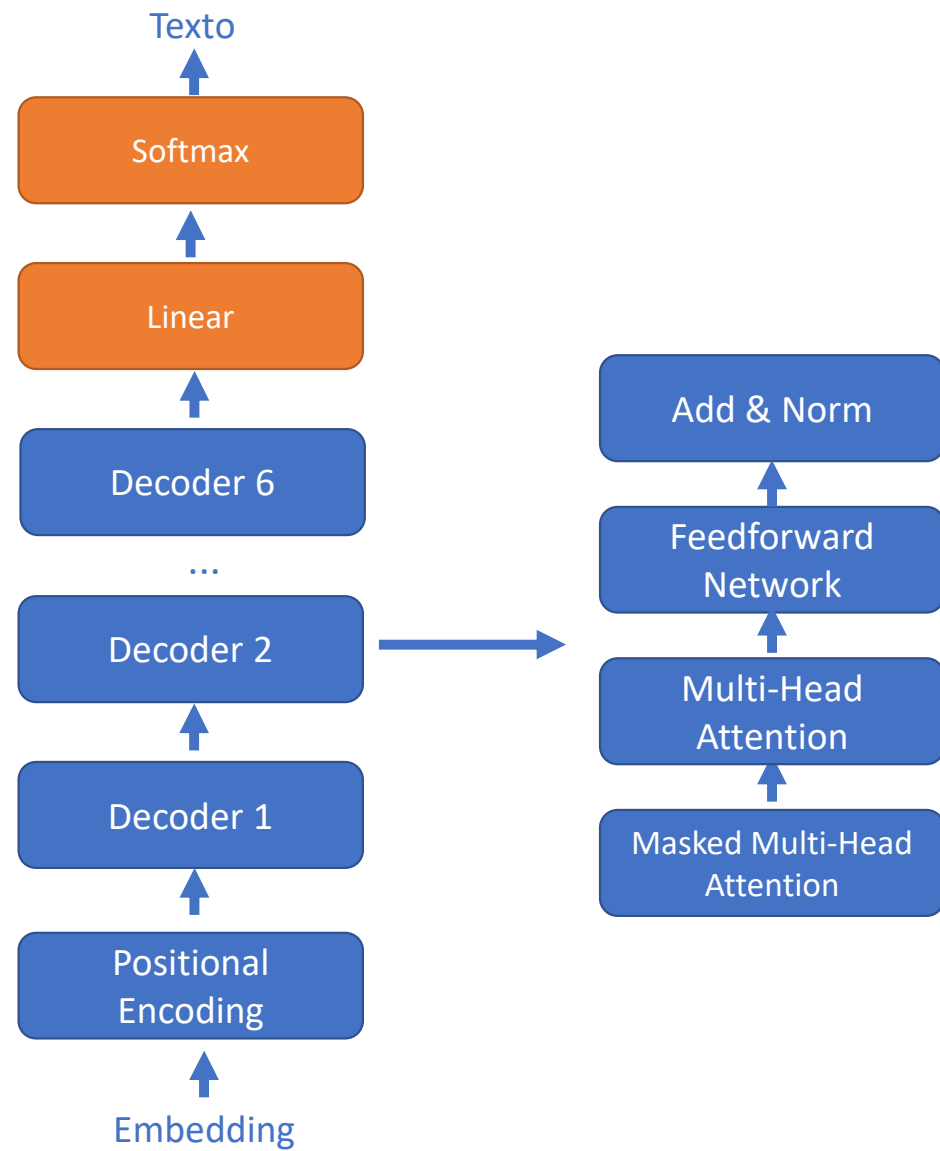
- `<sos>` Start of sentence
- `<eos>` End of sentence
- `<sos>` + “prever primeira sentença:” The
- `<sos>` + The “prever próxima:” animal
- `<sos>` + The + animal + ... + `<eos>`

# Masked Multi-Head Attention

Palavras a direita são mascaradas, de forma incremental...

	<b>The</b>	<b>Animal</b>	<b>Didn't</b>	<b>cross</b>	<b>the</b>
<b>The</b>	<b>0,9</b>	---	---	---	---
<b>Animal</b>	<b>0,57</b>	<b>0,79</b>	---	---	---
<b>Didn't</b>	<b>0,61</b>	<b>0,08</b>	<b>0,79</b>	---	---
<b>cross</b>	<b>0,79</b>	<b>0,39</b>	<b>0,60</b>	<b>0,39</b>	---
<b>the</b>	<b>0,60</b>	<b>0,57</b>	<b>0,39</b>	<b>0,57</b>	<b>0,39</b>

# Decoder



## Linear e Softmax

- Linear produz um vetor com o tamanho do vocabulário
- Softmax produz as probabilidades
- Decoder gera as palavras com maior probabilidade