

Personal Note

https://www.tensorflow.org/text/tutorials/transformer?authuser=0	Material
https://poloclub.github.io/transformer-explainer/	LLM Transformer Model Visually Explained
https://www.youtube.com/watch?v=wjZofJX0v4M	How large language models work, a visual intro to transformers

Review lại một số note từ week 6:

Transf ormer	Encoder	nhiều lớp, mỗi lớp bao gồm một lớp self-attention và một lớp feedforward (học các features phức tạp từ input)
	Decoder	tương tự như Encoder nhưng có thêm một lớp Masked Multi-Head Attention để kết nối với output của encoder
Thành phần chính	Self-Att ention	giúp mô hình tập trung vào các từ quan trọng trong câu, cho phép mô hình xác định những từ nào cần "chú ý" trong quá trình dịch.
	Multi-He ad Attentio n	tập trung vào các phần khác nhau của đầu vào; nắm bắt các mối quan hệ dài hạn
	Position al Encodin g	=> Transformer không có cấu trúc tuần tự như RNN, thêm positional encoding vào để cung cấp tt về vị trí của các token trong chuỗi
	Feedfor ward Network s	

Encoder-Decoder Model (e.g. T5, MBART)

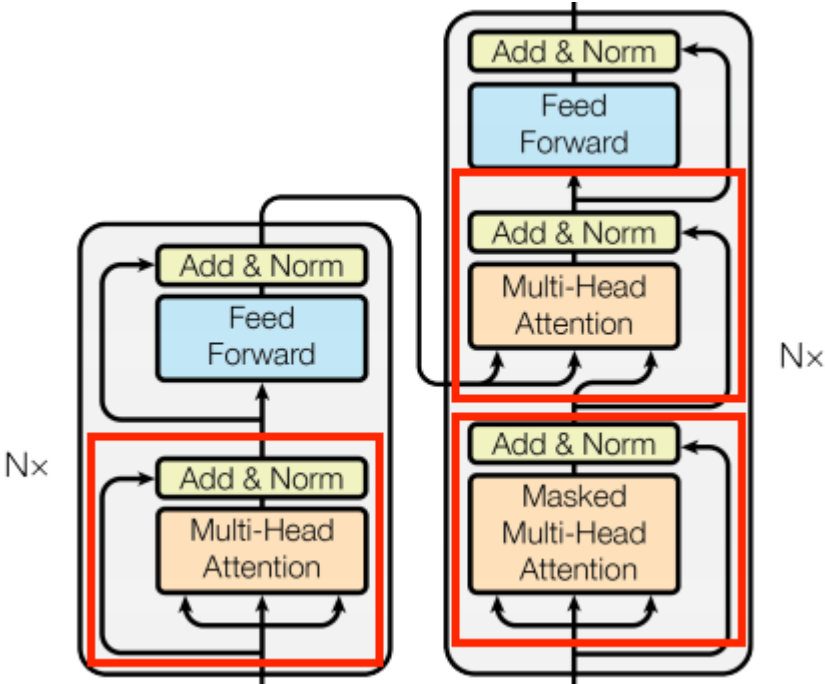
Dataset trong bài Material: TED Talks Open Translation Project, ngôn ngữ Portuguese-English từ tensorflow_datasets (50000 train, 1100 validate và 2000 test)
Các step:

Test the Dataset	<p>Input: (pt,en) Label: chuỗi (en) được dịch chuyển 1 vị trí (=> mỗi token đầu vào tiếng Anh sẽ khớp với nhãn là token tiếp theo)</p> <p>⇒ “Teacher forcing”: bắt kể đầu ra của mô hình là gì, chuỗi tiếp theo sẽ luôn dùng giá trị đúng (ground truth) làm đầu vào cho bước tiếp theo</p> <p>⇒ mô hình huấn luyện nhanh hơn vì không cần xử lý từng bước tuần tự, mà các đầu ra ở các vị trí khác nhau trong chuỗi có thể được tính toán song song</p>	<div><pre>print(en[0][:10]) print(en_labels[0][:10])</pre><p>tf.Tensor([[1369 43 3 201 245 5 22 6 1219</p></div> <p>(trong bài material không có ví dụ, nên em lấy tạm trong code phần assignment)</p>
Define the components	<ul style="list-style-type: none">- seq-2-seq- encoder, decoder- self-attention	
Embedding và Positional Encoding	<ul style="list-style-type: none">- Dùng lớp <code>tf.keras.layers.Embedding</code> để chuyển đổi từ chuỗi token => vector- Sau khi vector, thêm positional encoding vào vector embedding. Positional encoding dùng hàm sin và cos ở các tần số khác nhau để mã hóa vị trí từ trong chuỗi.⇒ đảm bảo các vị trí gần nhau có encoding tương đồng, giúp mô hình hiểu thứ tự từ trong câu- <code>en_emb._keras_mask</code>: <div><pre>> en_emb._keras_mask 133] ✓ 0.0s ... <tf.Tensor: shape=(64, 161), dtype=bool, nu array([[True, True, True, ..., False, Fa [True, True, True, ..., False, Fa [True, True, True, ..., False, Fa ..., [True, True, True, ..., False, Fa [True, True, True, ..., False, Fa [True, True, True, ..., False, Fa</pre></div>	<p>PositionalEmbedding:</p> <ul style="list-style-type: none">- <code>vocab_size</code>: kích thước từ vựng- <code>d_model</code>: chiều sâu của vector embedding.- Lớp <code>tf.keras.layers.Embedding</code> dùng để tìm vector embedding của các token.- <code>pos_encoding</code> là giá trị positional encoding được tạo ra trước đó.- Cộng vector embedding và positional encoding lại để có đầu ra hoàn chỉnh.- Dùng <code>en_emb._keras_mask</code> để mô hình bỏ qua các giá trị padding khi huấn luyện. <p>Padding là việc thêm các giá trị đặc biệt (thường là số 0) vào cuối (hoặc đầu) chuỗi để các chuỗi đều có độ dài như nhau. Ví dụ, nếu câu "how are you" được thêm padding để đạt chiều dài 6, ta có thể biểu diễn nó thành "how are you [0] [0] [0]".</p> <p>⇒ bỏ qua (mask) các vị trí chứa padding khi tính toán, đảm bảo mô hình chỉ học từ các giá trị thật trong chuỗi.</p>

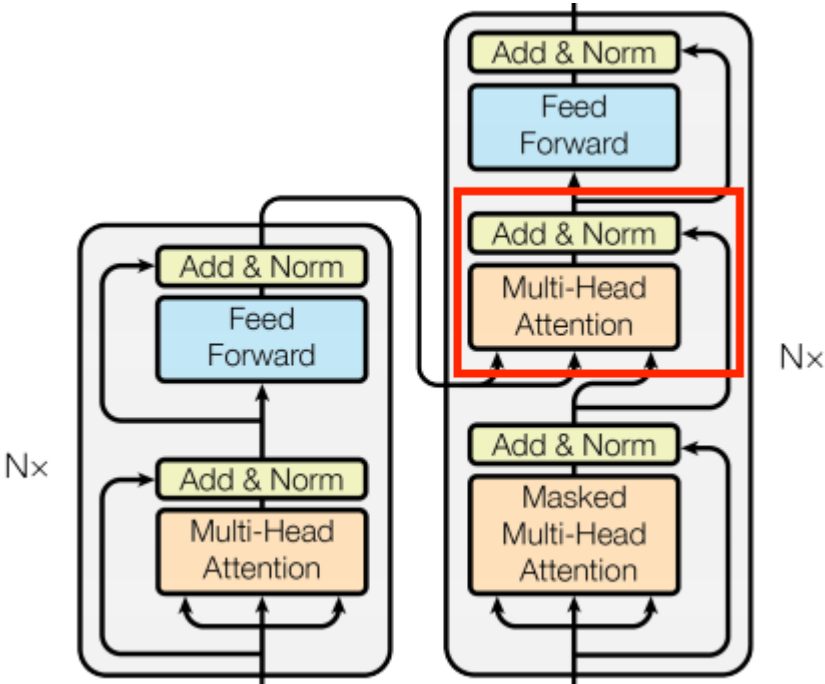
		<p>en_emb._keras_mask là một mask được tạo ra tự động để đánh dấu các vị trí có giá trị padding. Giá trị True trong mask cho biết token thật, còn False cho biết vị trí padding.</p>
--	--	--

Add & Normalize (Residual Connections và Layer Normalization)

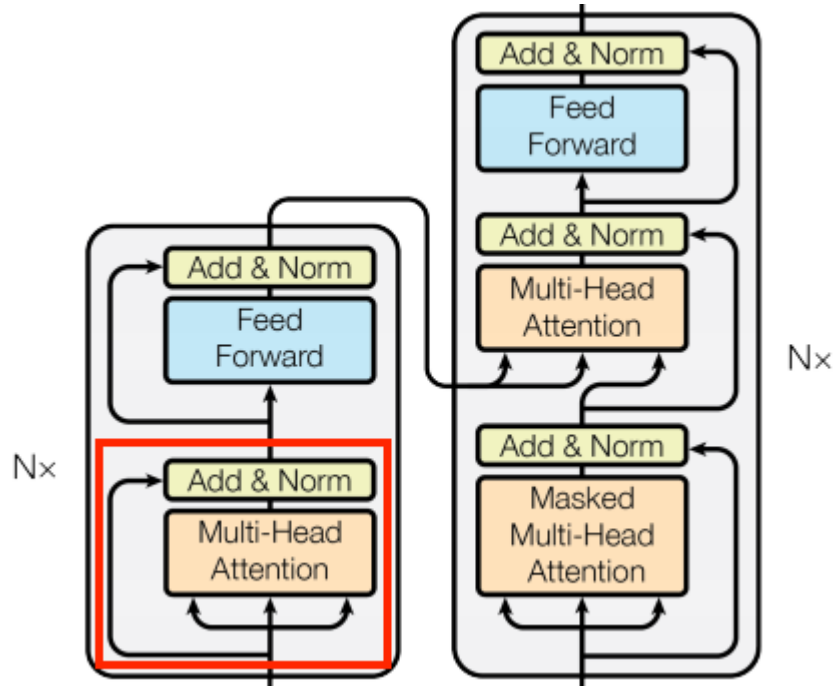
Base Attention Layer



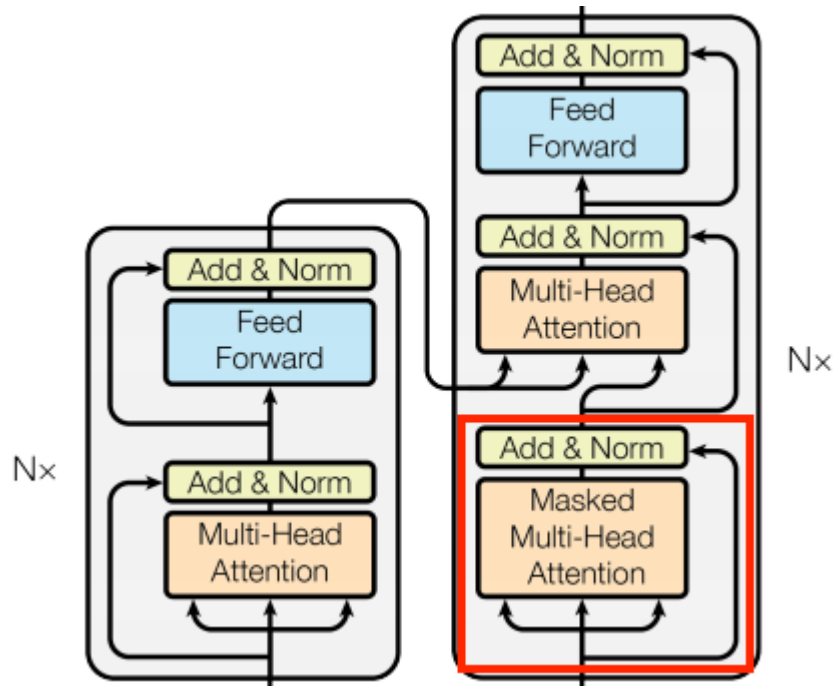
The cross attention layer



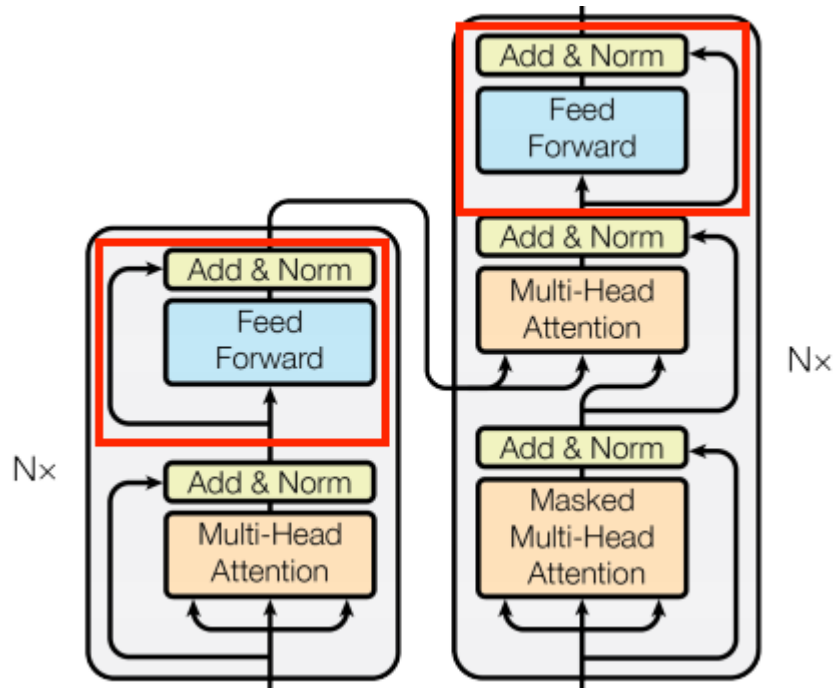
The global self-attention layer



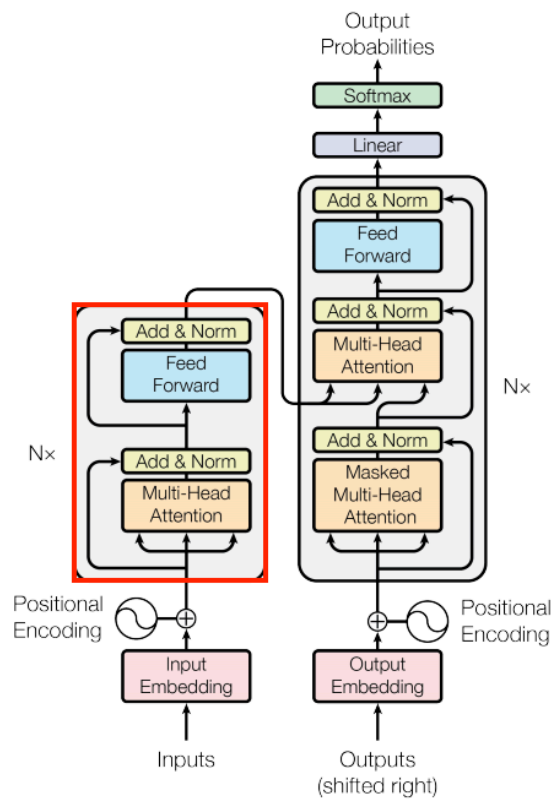
The causal self-attention layer



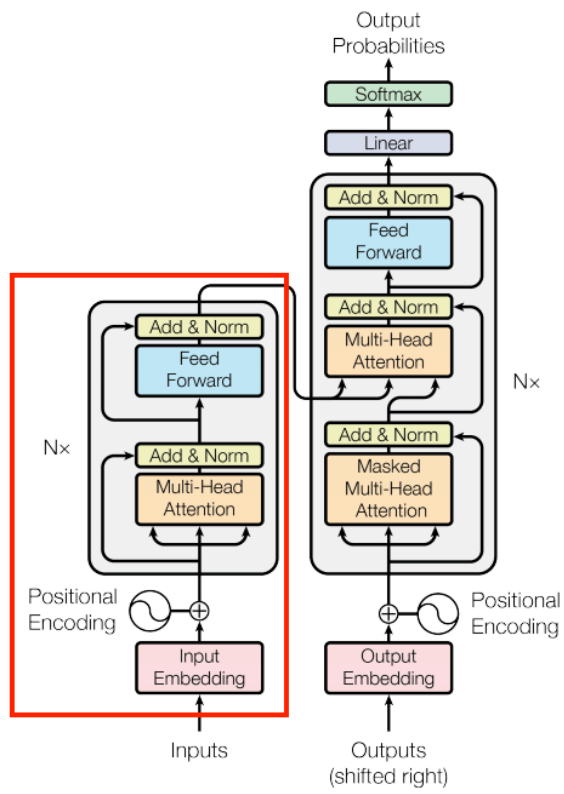
The feed forward network



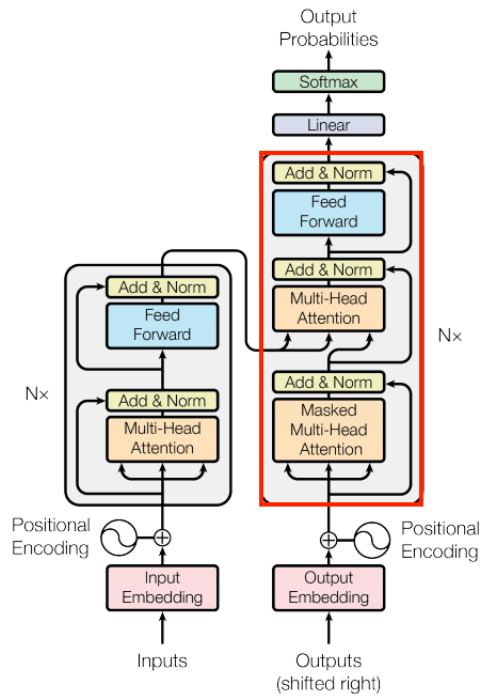
The encoder layer

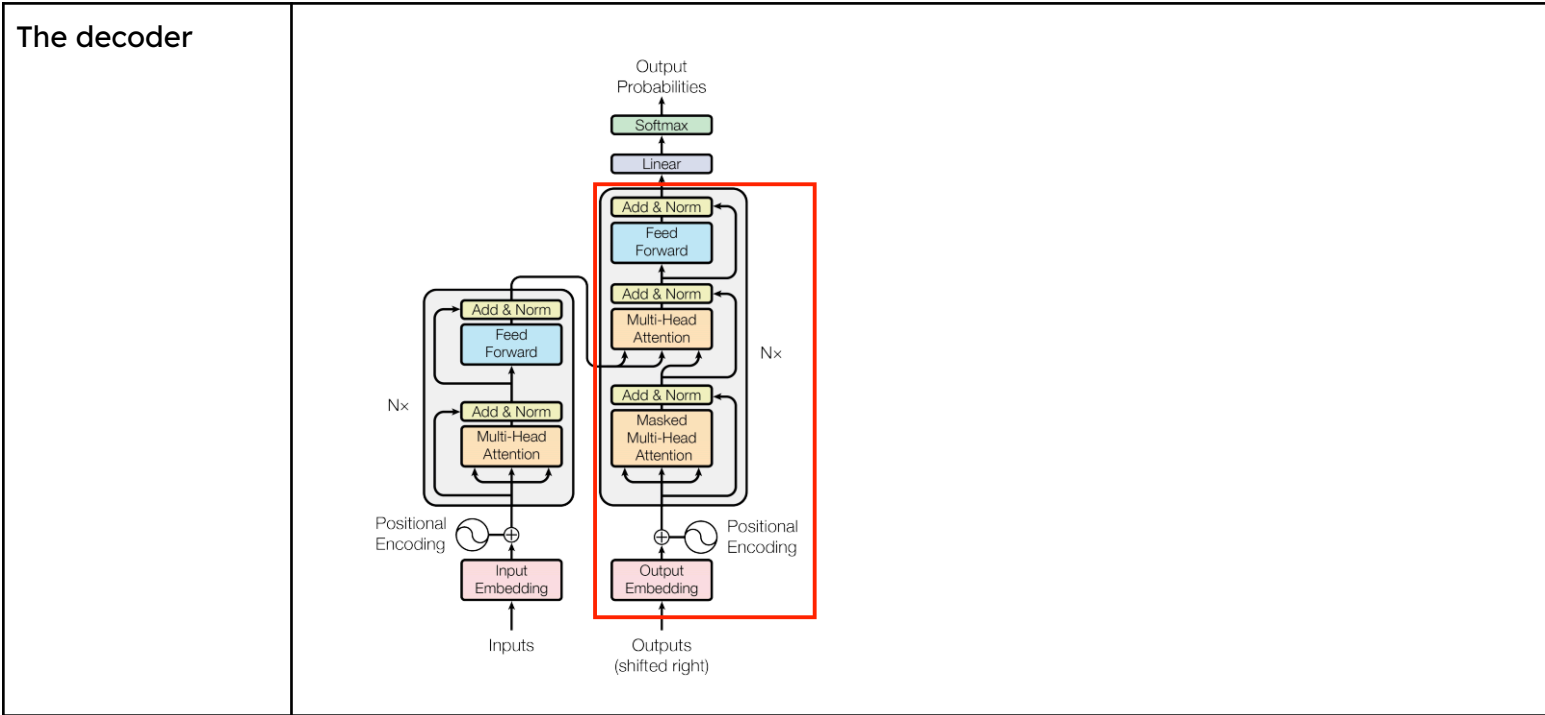


The encoder



The decoder layer





The transformer		
	Create the Transformer by extending <code>tf.keras.Model</code>	
Hyperparameters	<code>num_layers</code> : số lượng các lớp encoder và decoder trong mô hình Transformer.	Trong bài báo gốc dung 6 <code>num_layers</code> , có thể dung 4 lớp để nhẹ hơn
	<code>d_model</code> : kích thước của các vector embedding	512 trong mô hình gốc, có thể giảm xuống 128
	<code>dff</code> : kích thước của lớp FeedForward trong Transformer.	mô hình gốc là 2048 (trong bài material dung 512)
	<code>num_heads</code>	số lượng "heads" trong lớp attention
	<code>dropout_rate</code>	tỷ lệ dropout để ngăn overfitting
<code>attn_scores</code>	<code>attn_scores = transformer.decoder.dec_layers[-1].last_attn_scores</code>	các trọng số attention từ lớp cuối cùng trong decoder, dùng để theo dõi mức độ mà mô hình tập trung vào từng từ trong câu khi tạo bản dịch
Training		
Set up the optimizer	Use the Adam optimizer with a custom learning rate scheduler <ul style="list-style-type: none"> - Learning Rate Schedule - Warmup Steps: tăng dần tốc độ học qua một số bước nhất định 	Optimizer (Adam) điều chỉnh trọng số của mô hình = cách giảm dần loss dựa trên Learning Rate Schedule tùy chỉnh

	=> mô hình dần quen với quá trình điều chỉnh trọng số trước khi học với tốc độ tối đa, giảm thiểu nguy cơ gradient vanishing/exploding.	Sau một số bước, learning rate tăng dần, sau đó giảm khi mô hình đã qua giai đoạn "khởi động" (warmup).
Set up the loss and metrics	<p>Trong quá trình tính toán loss, cần có mask để loại bỏ các giá trị padding nhằm tránh ảnh hưởng đến kết quả</p> <ul style="list-style-type: none">- Masking: loại bỏ các giá trị padding trong dữ liệu đầu vào. (Trong neural machine translation, các câu đầu vào có độ dài khác nhau => padding => khi tính loss và accuracy, cần bỏ qua giá trị này => tập trung vào những từ thực tế trong câu)- Sparse Categorical Crossentropy: hàm mất mát thích hợp cho các bài toán phân loại nhiều lớp với nhãn đầu ra được mã hóa dưới dạng số nguyên- Attention Weights: mức độ liên quan giữa từ ở ngôn ngữ nguồn và từ đang được dịch	masked_loss & masked_accuracy
Train the model	transformer.compile	
	transformer.fit	