$lights = 13600 \, g \, 0.01$

$f \times 0.17 = 13600 \times 0.01$

$$\frac{13600}{17} = f$$

$$\underline{800 = f}$$

Decoder-only architecture base are more scalable due to simplicity of training data, and engineering simplicity. They have such architecture which removes complexity of cross layer Attention layers and also removes balancing of encoder to decoder depth ratio.

Decoder only architecture uses Next Token prediction. This ensures that all of the training tokens contribute towards the learning of gradient.

Next Token Prediction is sufficient as, ~~on far so attra me~~ for accurate prediction ~~sta nut note it a repures~~, the model must learn thes underlying structure, meaning of the language/ sentence.

In next token prediction, the model has to learn meaning, structure and reasoning patterns which eventually help the model to translate, summarize and reason

**Q.1** Decoder only architecture leave are more scalable due to simplicity of training data, and engineering simplicity. They have such architecture which removes complexity of Cross coded Attention layers and also removes balancing of encoder to decoder depth ratio.

Decoder only architecture uses Next Token prediction. This ensures that all of the training tokens contribute towards the training of gradient.

**Q.2** Next Token Prediction is sufficient as, ~~so far to attack to~~ for accurate prediction ~~of the next word in a sequence~~, the model must learn thee underlying structure, meaning of the language/sentence.

In next token prediction, the model has to learn meaning, structure and reasoning patterns which eventually help the model in translation, summarize and reason