**Q1** A single perceptron cannot solve XOR because as we can see and conclude from the truth table that XOR is not linearly separable function.

Whereas a single perceptron can only solve linearly separable functions, because it uses linear activation function.

For eg: $z = \sum w_i x_i + b$ , $y = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$

XOR

| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |
| 1 | 0 | 1 |

With introducing multi-layer perceptrons we bring in hidden layers which allow us to even represent functions which are not linear, by using activation functions like ReLU, sigmoid, tanh, etc. that are non-linear.

$$z = \sum w_i x_i + b, \quad h = \sigma(z) \qquad \sigma \to \text{non-linear activation function}$$

**Q2** Stacking linear layers : $y_1 = w_1 x + b_1$

$$y_2 = w_2 y_1 + b_2 = w_2(w_1 x + b_1) + b_2$$
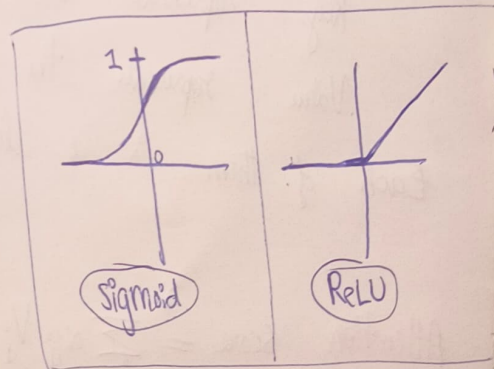$$= \underline{w_1 w_2} x + \underline{w_2 b_1 + b_2}$$
$$y_2 = w_3 x + b_3 \longleftarrow \text{Linear}$$

∴ Each layer just does some linear function on the input resulting in the stack still being linear as shown for arbitrary linear function above.

Since in backpropagation we use chain rule, which involves multiplication and if the derivative is <1, repeated multiplication results in the shrinking of gradients (deep networks)

sigmoid $\quad \sigma_1(x) = \dfrac{1}{1+e^{-x}} \qquad \sigma_1'(x) = \sigma(x)\big[1 - \sigma(x)\big]$

ReLU $\quad \sigma_2(x) = \max(0, x) \qquad \sigma_2'(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$



Sigmoid          ReLU

as we can see max derivative in case of sigmoid $\to 0.25$

ReLU $\to 1$

So ReLU solves vanishing gradient better than sigmoid since in case of active
newrons ⊗ derivative = 1 hence allowing gradients to go through many layers.

(Q3) We need positional encoding since transformers do not know the order of words
by default. If two sentences have same words but different order the
transformer will not be able to differentiate without positional encodings.

Sinusoidal PE ⇒ cannot be learned. It is a fixed mathematical function. It has
no max sentence length and is better in understanding relative
awareness. Extrapolation = good.

Absolute PE ⇒ It is learned from training examples. It has a max training length
and learns vectors for encoding. Poor at extrapolation.

RoPE ⇒ Instead of directly adding positional encoding with token embeddings, it
rotates query and key vectors based on position.

$$\begin{pmatrix} q'_1 \\ q'_2 \end{pmatrix} = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \begin{pmatrix} q_1 \\ q_2 \end{pmatrix}$$
← rotates like
this

$\theta \rightarrow$ relation with
position $i$

$\theta \rightarrow f(i)$

In this way, relative position is naturally encoded and learned. Since rotation
is periodic and continuous (∴ not fixed encodings) so better for long contexts
and extrapolation since attention score now depends on relative positions.

(Q4) Query represents what information the token needs
Key represents what info the token has/offers to other tokens
Value represents the actual info that is being passed forward.

Each of them is a learned linear projection $Q = XW_Q$     $X \rightarrow$ Input token
$K = XW_K$       embeddings
$V = XW_V$

Attention Score $= \sum_j \alpha_{ij} V_j$     $\alpha_{ij} = \text{softmax}\left(\dfrac{Q_i \cdot K_j}{\sqrt{d_k}}\right)$

Assuming they are normally distributed so mean $= 0$
variance $= 1$

but dot product $Q \cdot K = \sum_{i=1}^{d_k} q_i k_i$ ← variance = $d_k$

if not scaling then as $d_k$ increases, dot product increases, gradients vanish.

by scaling    Var $\left( \dfrac{Q \cdot K}{\sqrt{d}} \right) \simeq 1$    resulting in better gradient flow by avoiding softmax saturation ✪

Diagonal entry ⇒ $Q_i \cdot K_i$ ⇒ a token attention to itself ∴ highest attention score (self)

(Q5) Splitting attention into multiple heads allows us to learn many patterns at the same time by running in parallel. So it is fast because parallel and learns complex relations better.

$head_i = $ Attention $(Q_i, K_i, V_i)$    Multihead = concat $(head_1, head_2, \ldots head_h) W_0$

✪ $d_{model}$ ⇒ dimension of token embedding

$h$ ⇒ no. of heads

$d_{head}$ ⇒ dimensions per head $= \dfrac{d_{model}}{h}$

(Q6) Greedy decoding is suboptimal ~~because~~ because it selects the token with max probability at each time step t. But what we actually want is to maximise the joint probability of all tokens. Basically the problem with greedy, although it feels the best intuitively ~~local maximum~~

└───────→ Local maximum ≠ Global maximum

Beam search is better because instead of selecting the best one, it keeps the 'top k' sequences at each step then expand each, score all and keep the best k. So it goes through multiple options in parallel and chooses the globally optimum one.

(Q1) (a) $d\_head = \dfrac{d\_model}{h} = \dfrac{768}{12} = \boxed{64}$

(b) Parameters $\Rightarrow$ $W_Q, W_K, W_v \in \mathbb{R}^{d_{model} \times d_{model}} = 768 \times 768$

(One layer)
$$= \boxed{589824} \text{ each}$$

(Q2) $softmax = \dfrac{e^{z_i}}{\sum\limits_{i} e^{z_i}}$

$e^2 = 7.389$

$e^1 = 2.718$     $\left.\begin{array}{r} \end{array}\right\}$ add $= 11.107$

$e^0 = 1$

$softmax\ (2.0) = \dfrac{7.389}{11.107} = \boxed{0.665}$

$softmax\ (1.0) = \dfrac{2.718}{11.107} = \boxed{0.245}$

$softmax\ (0.0) = \dfrac{1}{11.107} = \boxed{0.090}$