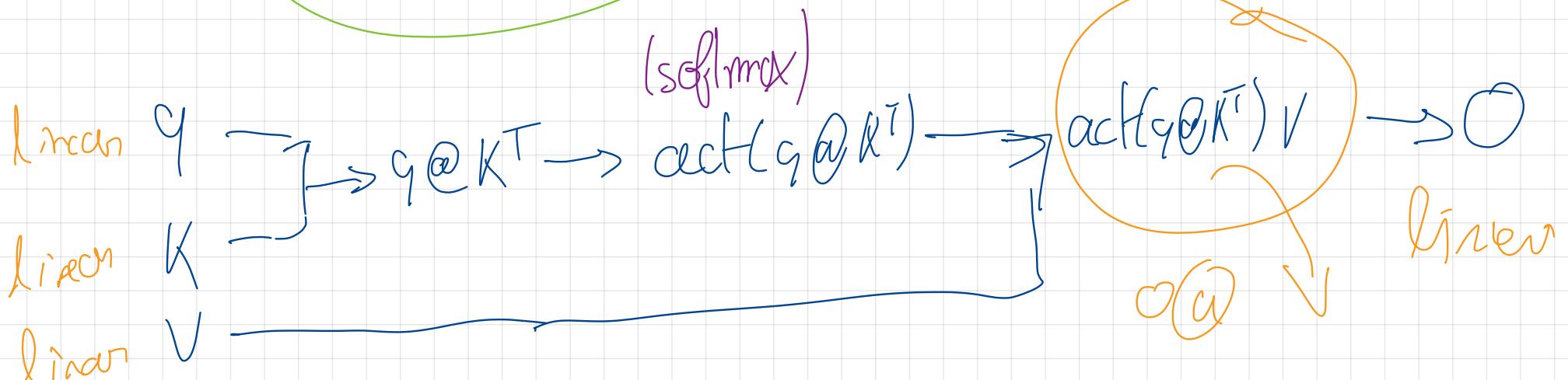
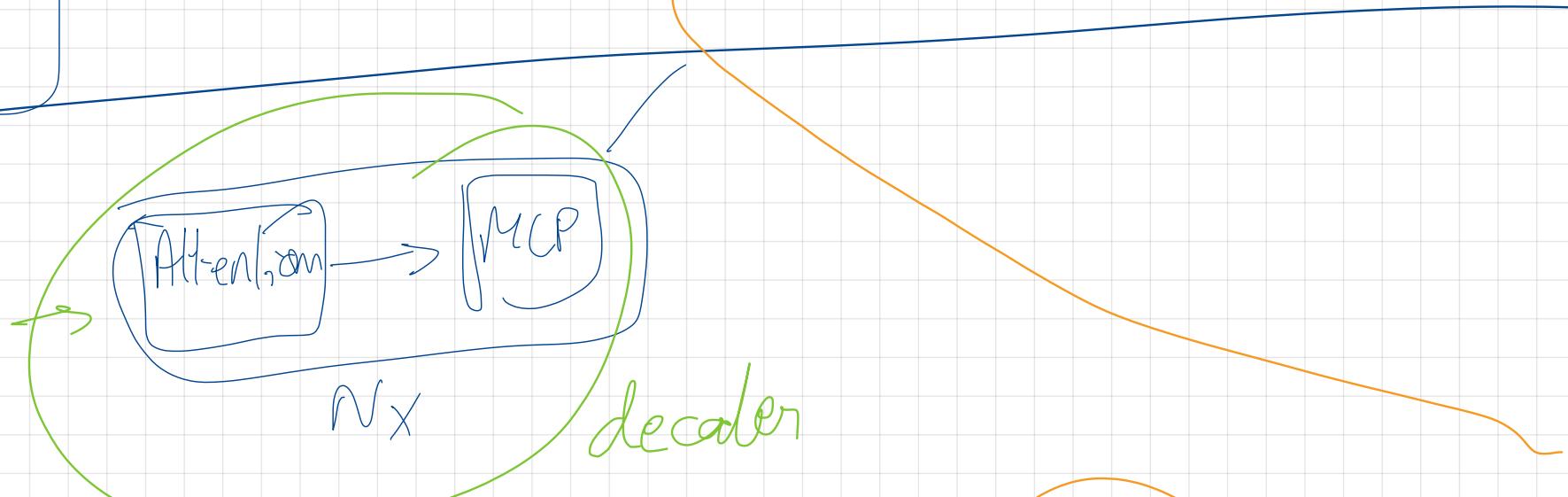
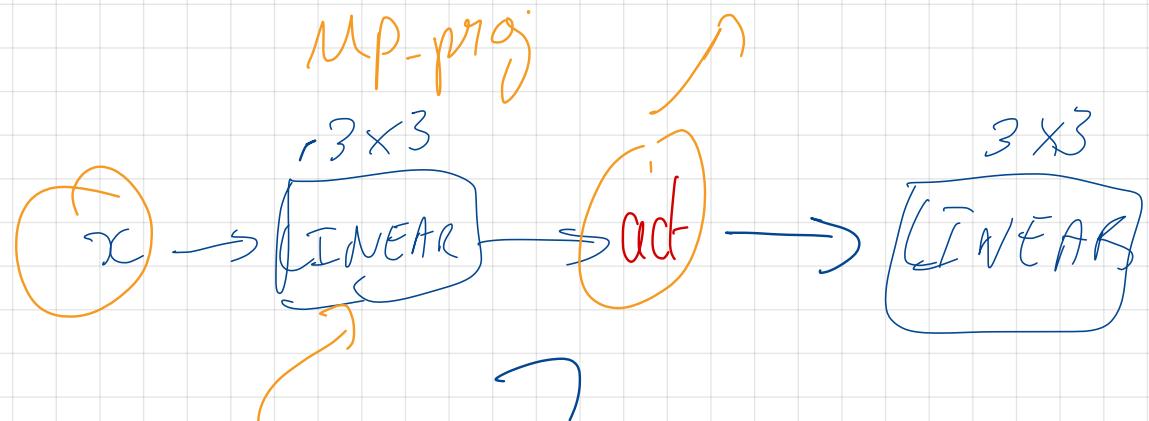
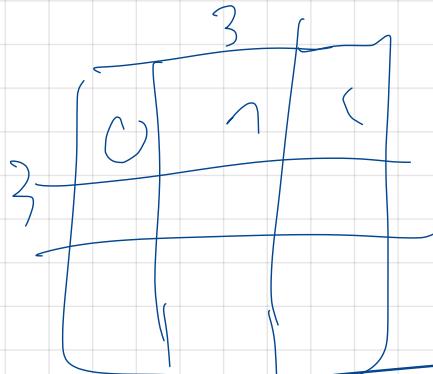
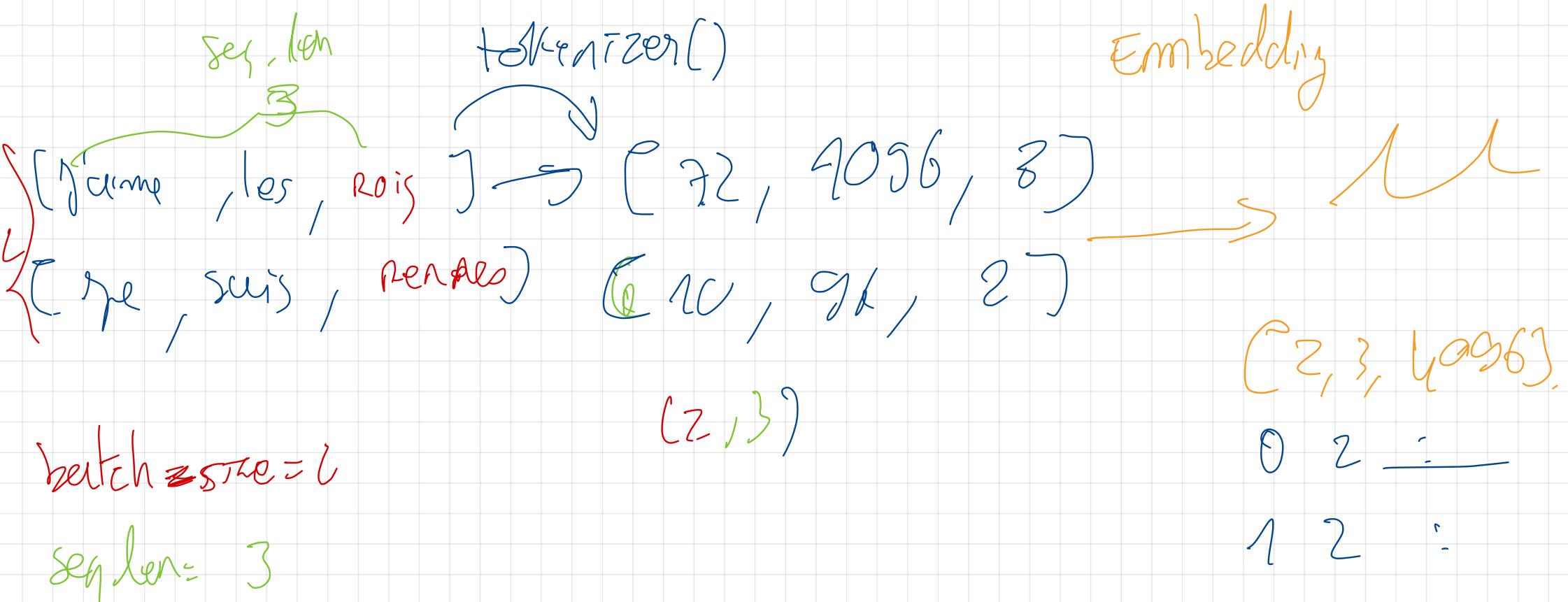
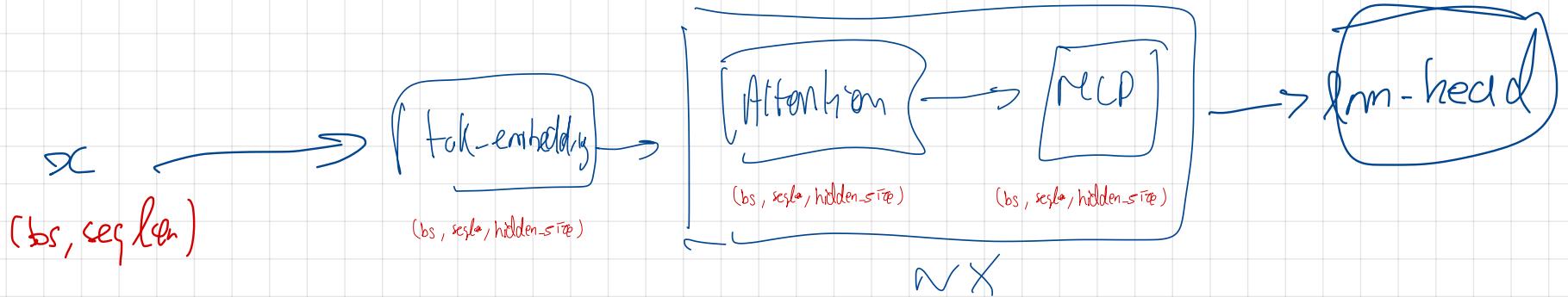


COURS 1

Linear / MCP



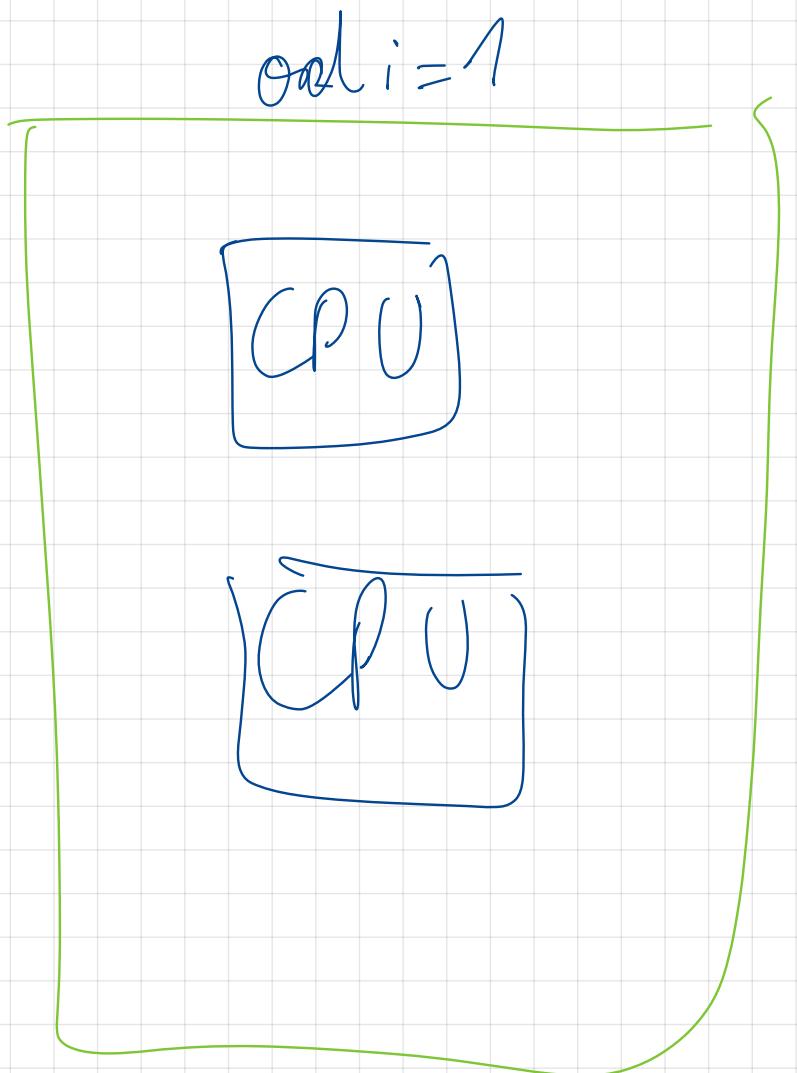


Elle aime les lois

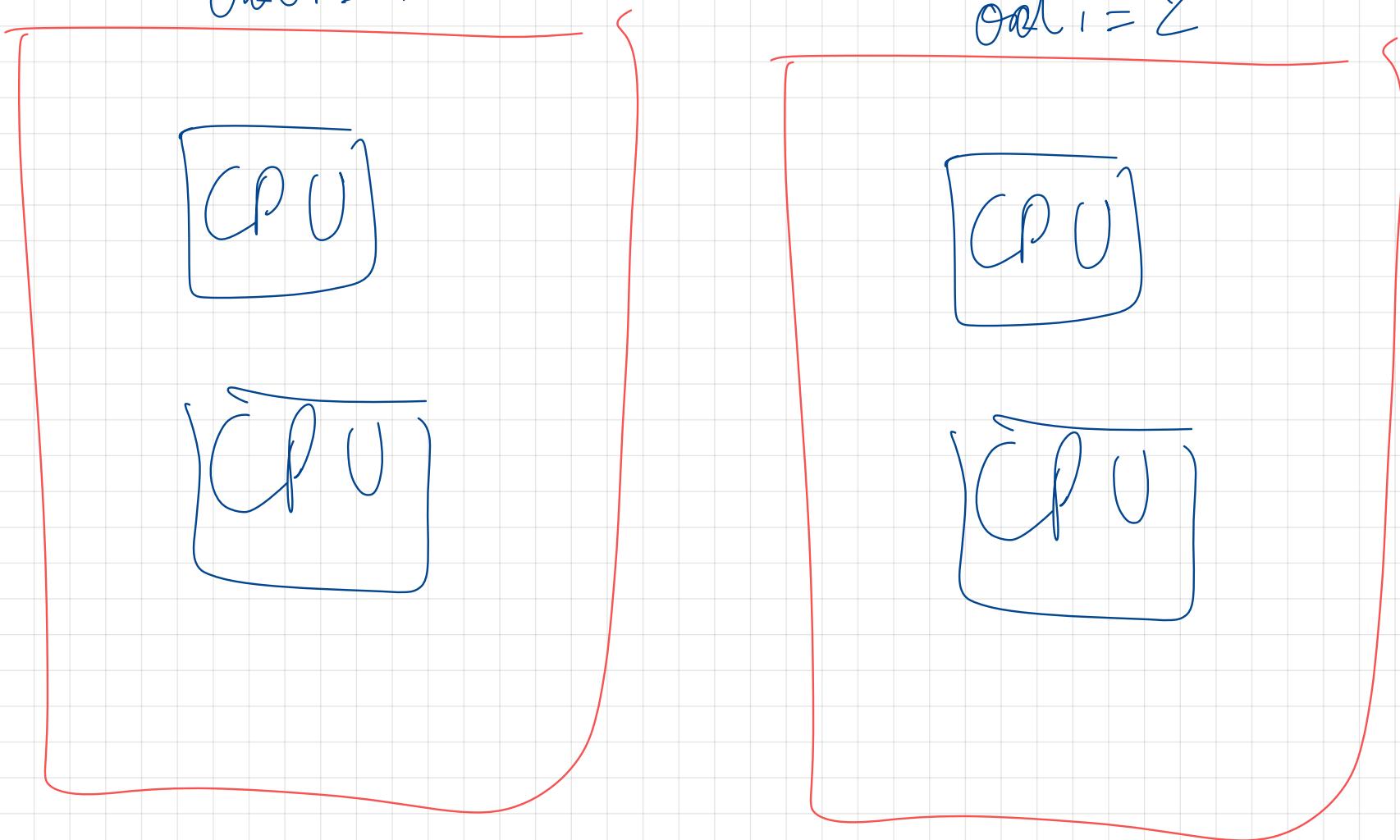
? /

COUR2:

torchrun --nproc_per_node=2 --nodes=1



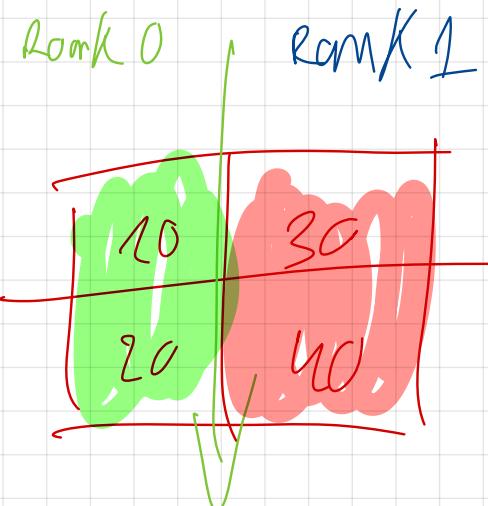
`torchrun` --
 |
 | $\text{--nproc_per_node} = 2$
 |
 | $\text{--local_rank} = 1$
 |
 | $\text{--local_rank} = 2$
 |
 | $\text{--world_size} = 2$



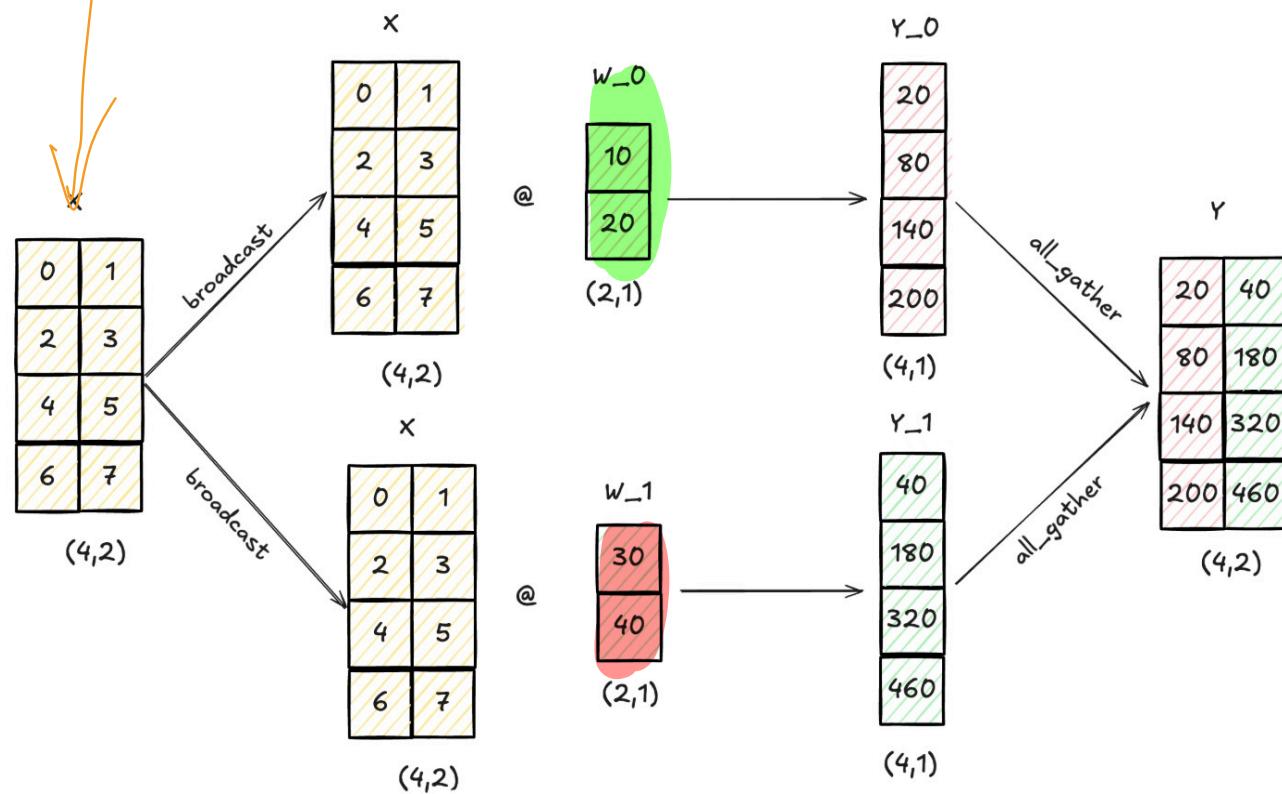
Cours 3:

Tensor parallel

$$f(x) = \underbrace{(xW)^T}_{\text{Column linear}} + b$$



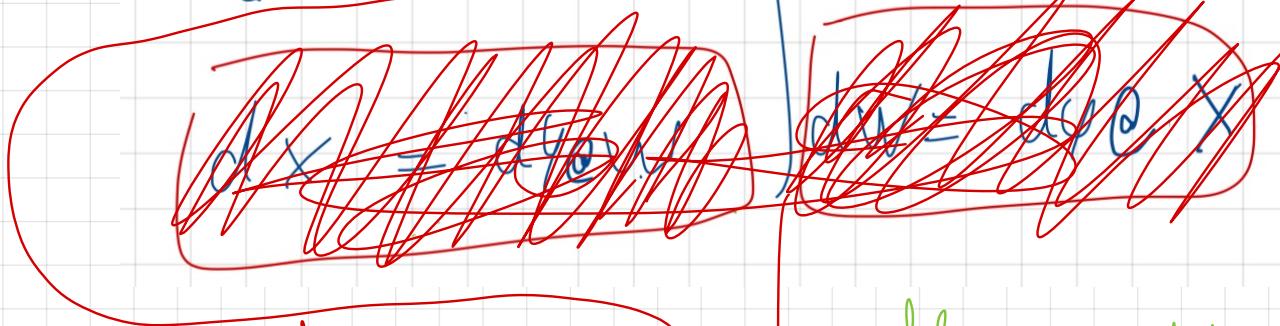
Column linear



$$\psi = x @ w \upharpoonright$$

$$\frac{dy}{dx} = w$$

$$\frac{dy}{dw} = X$$

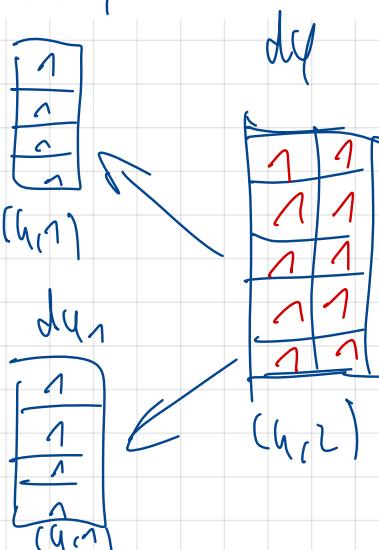
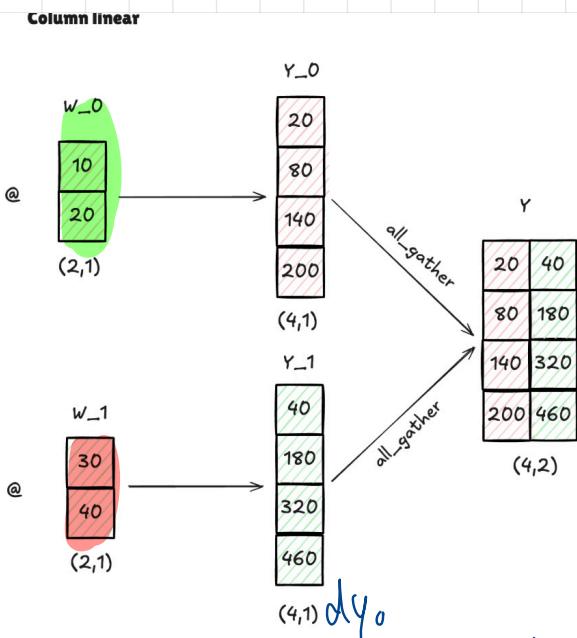
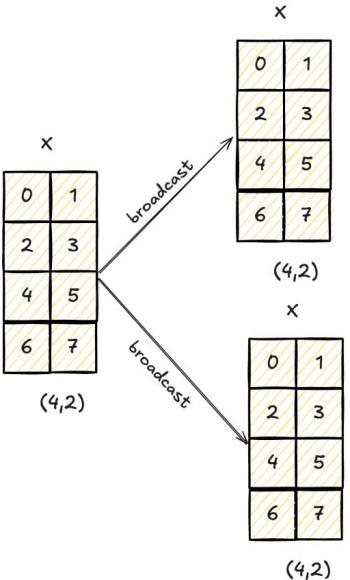


$$\frac{dL}{dx} = \frac{dL}{dy} \left(\frac{dy}{dx} \right) = \left(\frac{dL}{dy} \right) w$$

$$\frac{dL}{dw} = \frac{dt}{dy} \times \left(\frac{dy}{dw} \right) = \left(\frac{dt}{dy} \right) X$$

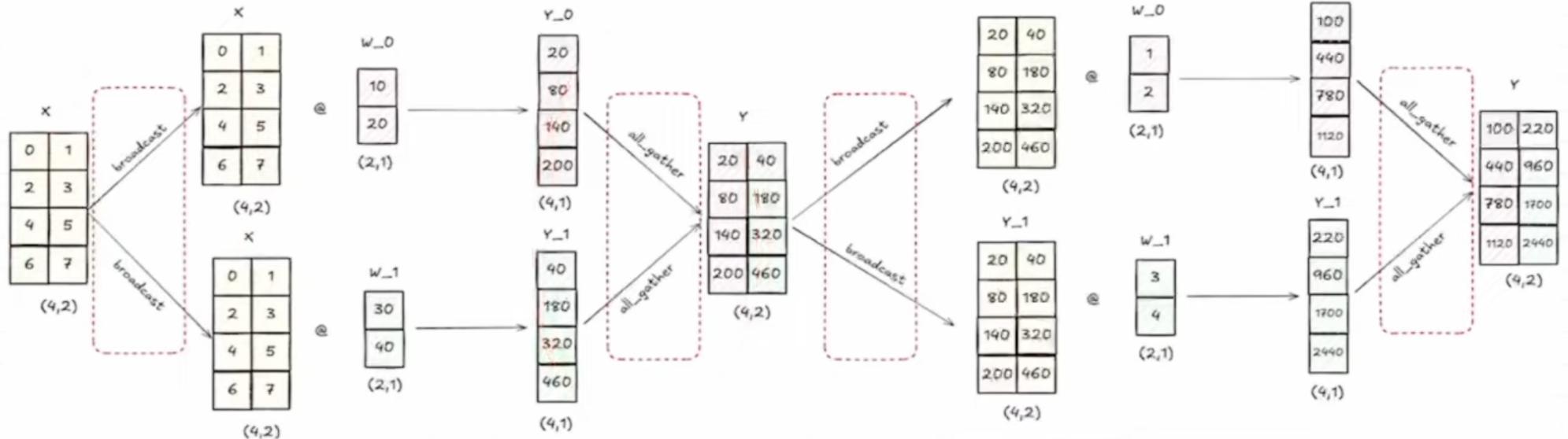
Backward Column Linear

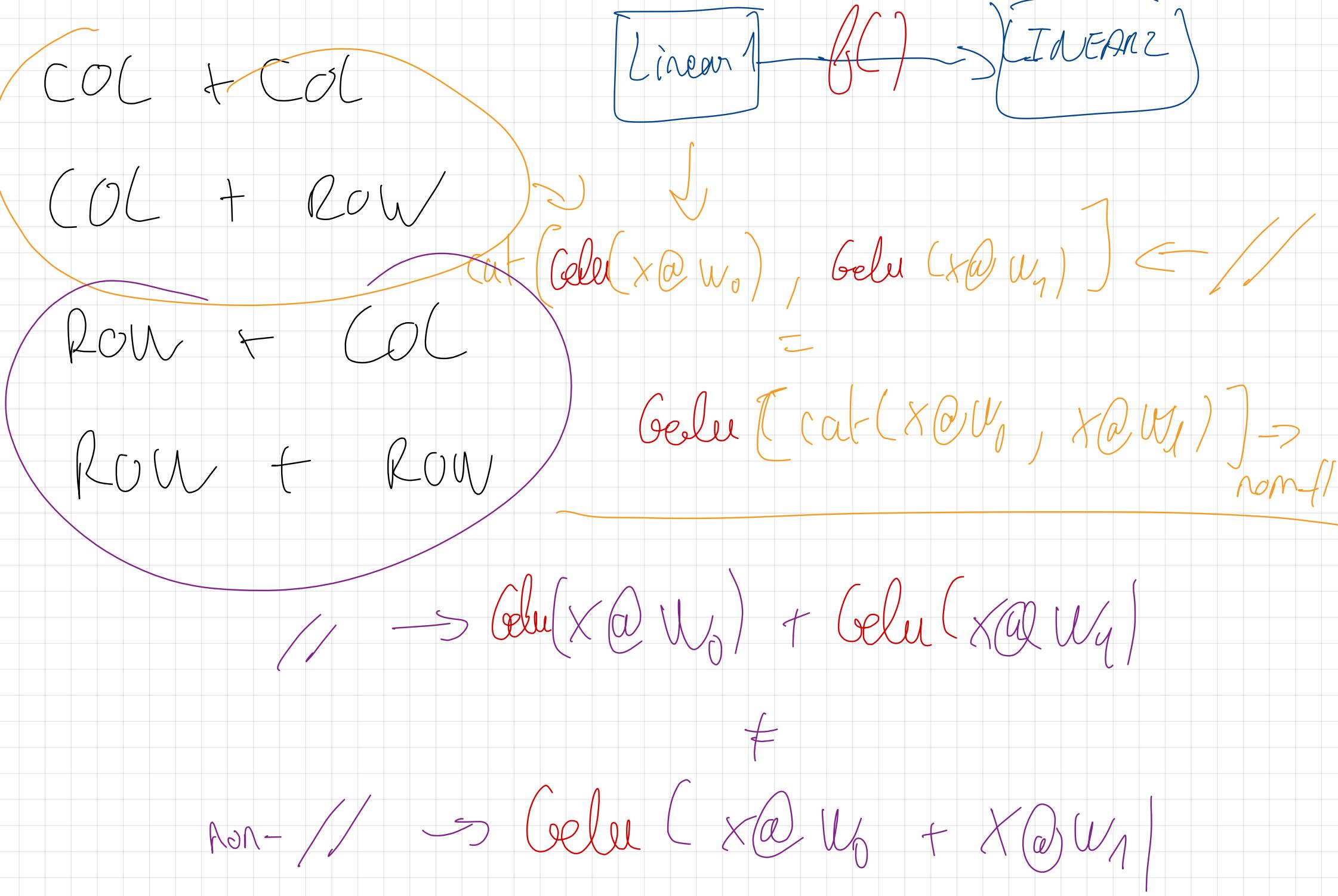
$$\frac{dL}{dX} = \frac{dL}{dy} \times w$$



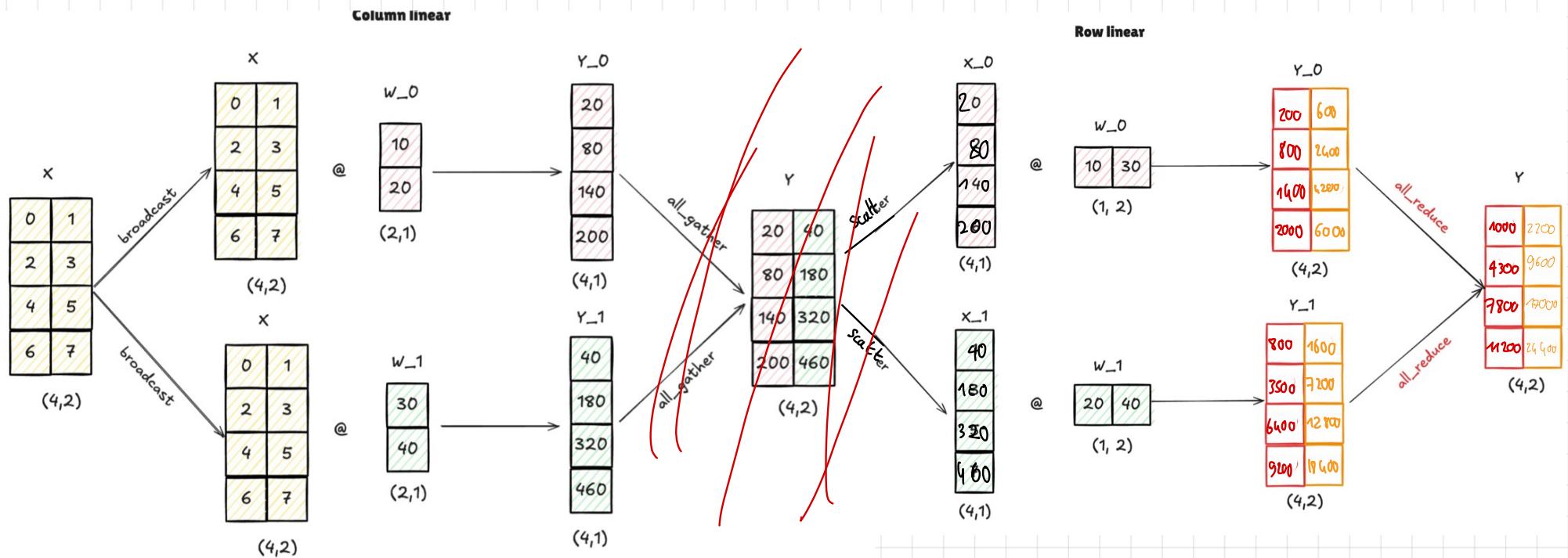
Week 3 + A

Column linear + Column Linear (forward)

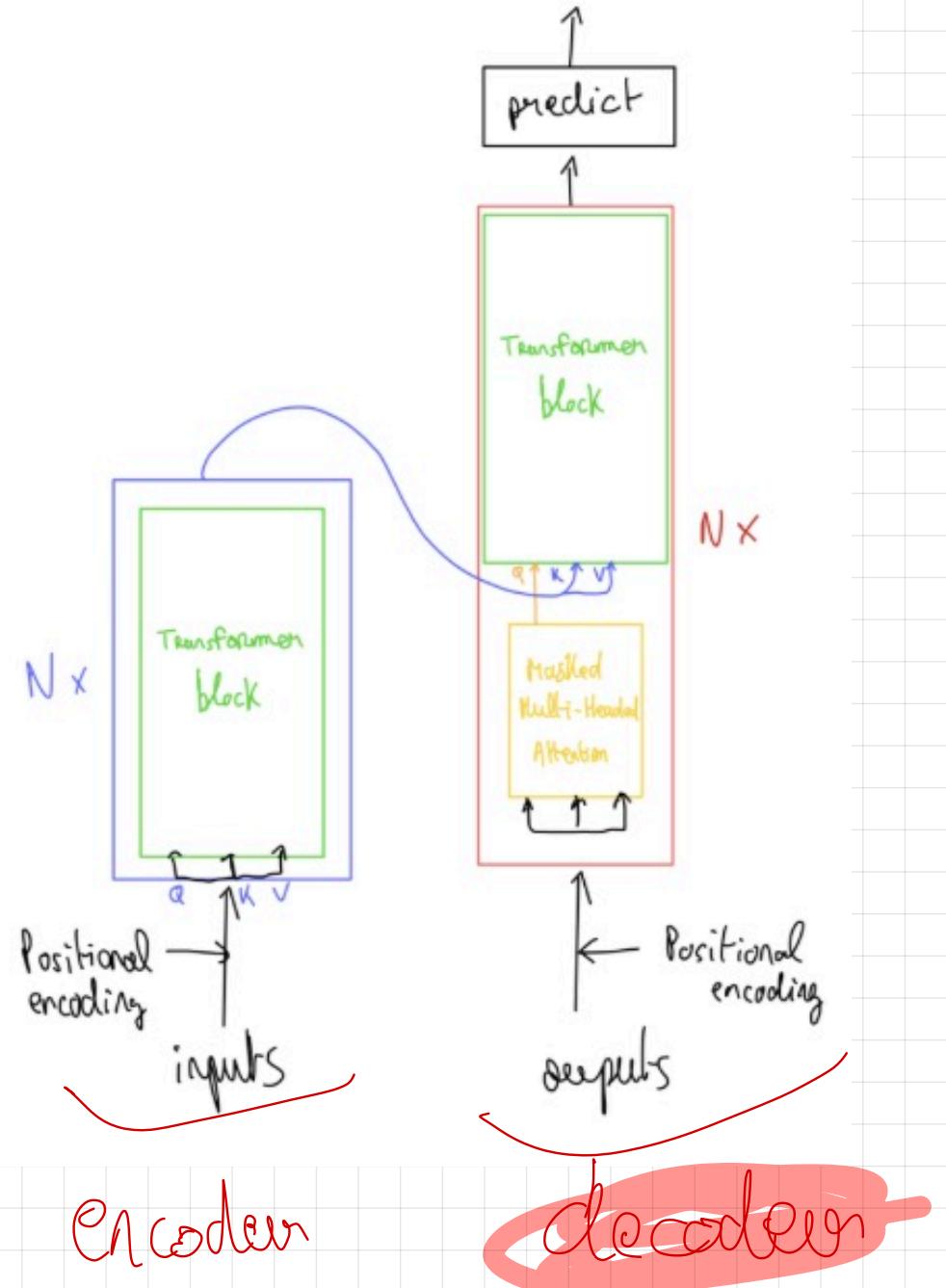
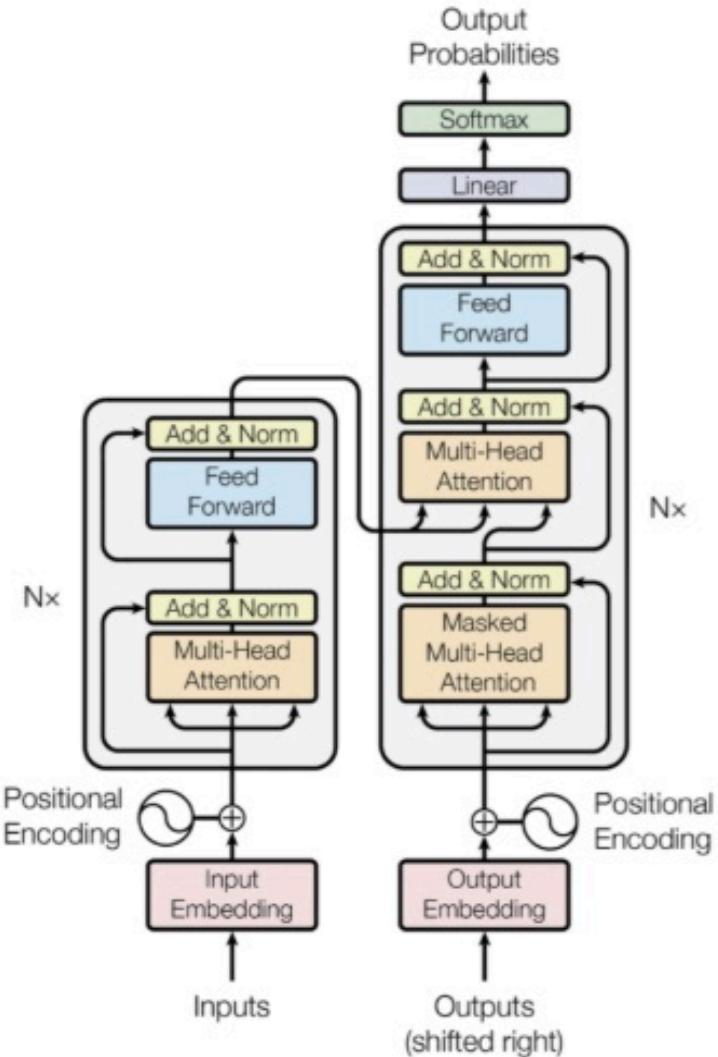




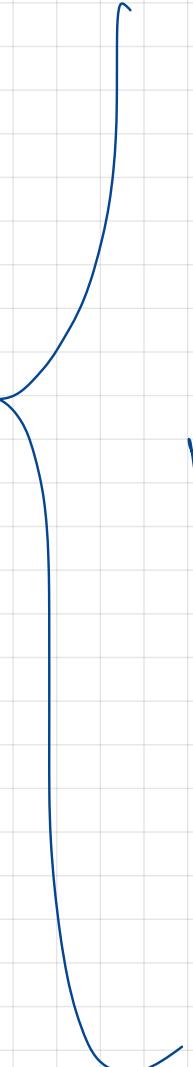
COLUMN + ROW linear



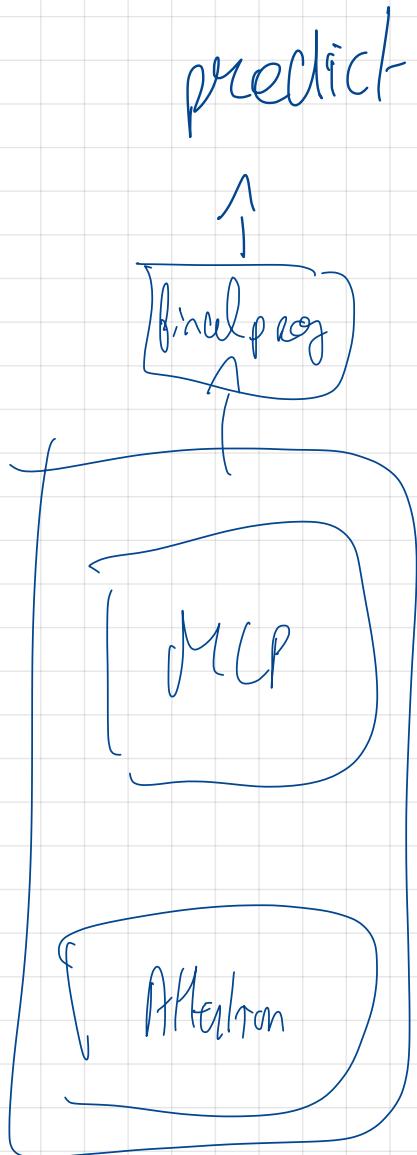
Saves 2 communications
(compared to Col + Col)



N decoder
layer



W X



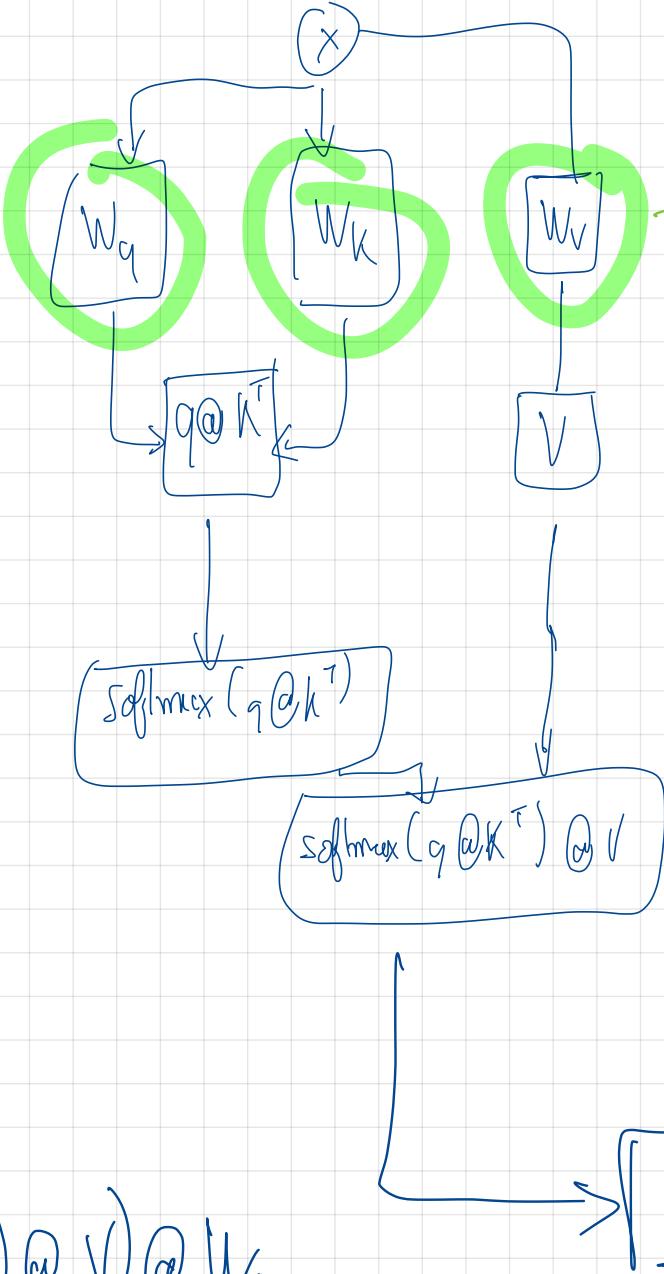
↑ positional encodig

Attention

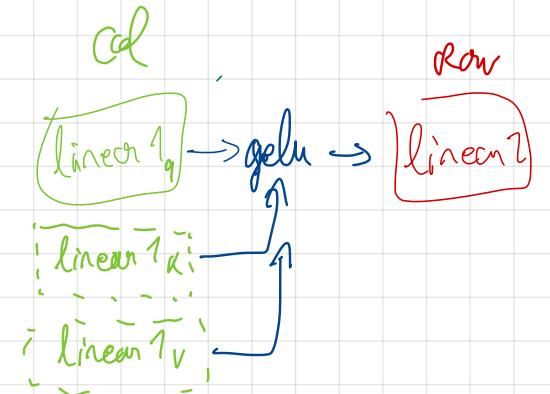
$$q: x @ W_q$$

$$k: x @ W_k$$

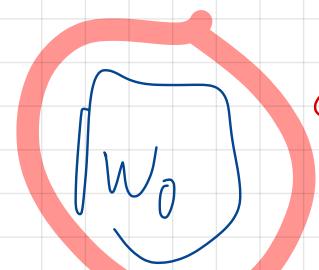
$$v: x @ W_v$$



Column Parallel

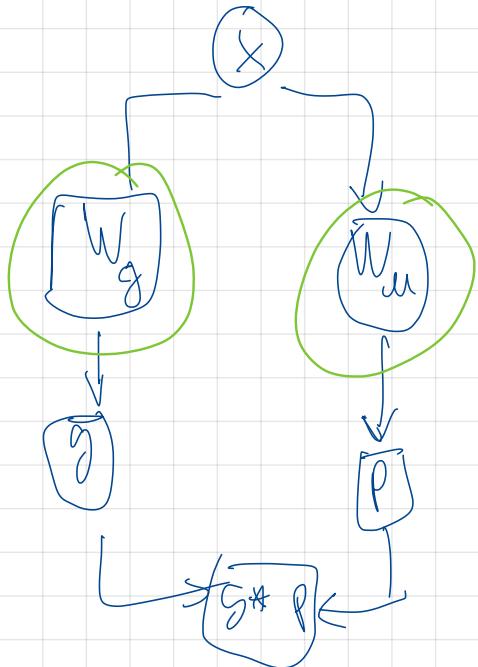


Row Parallel



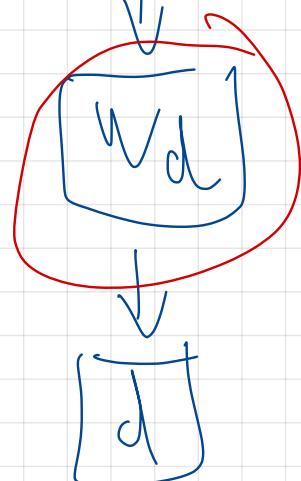
$$O = (\text{softmax}(q @ K^T) @ v) @ W_O$$

MUL



Column

$\text{Silu}(g * p)$



Row

$d = \text{Silu}(g * p) @ W_d$