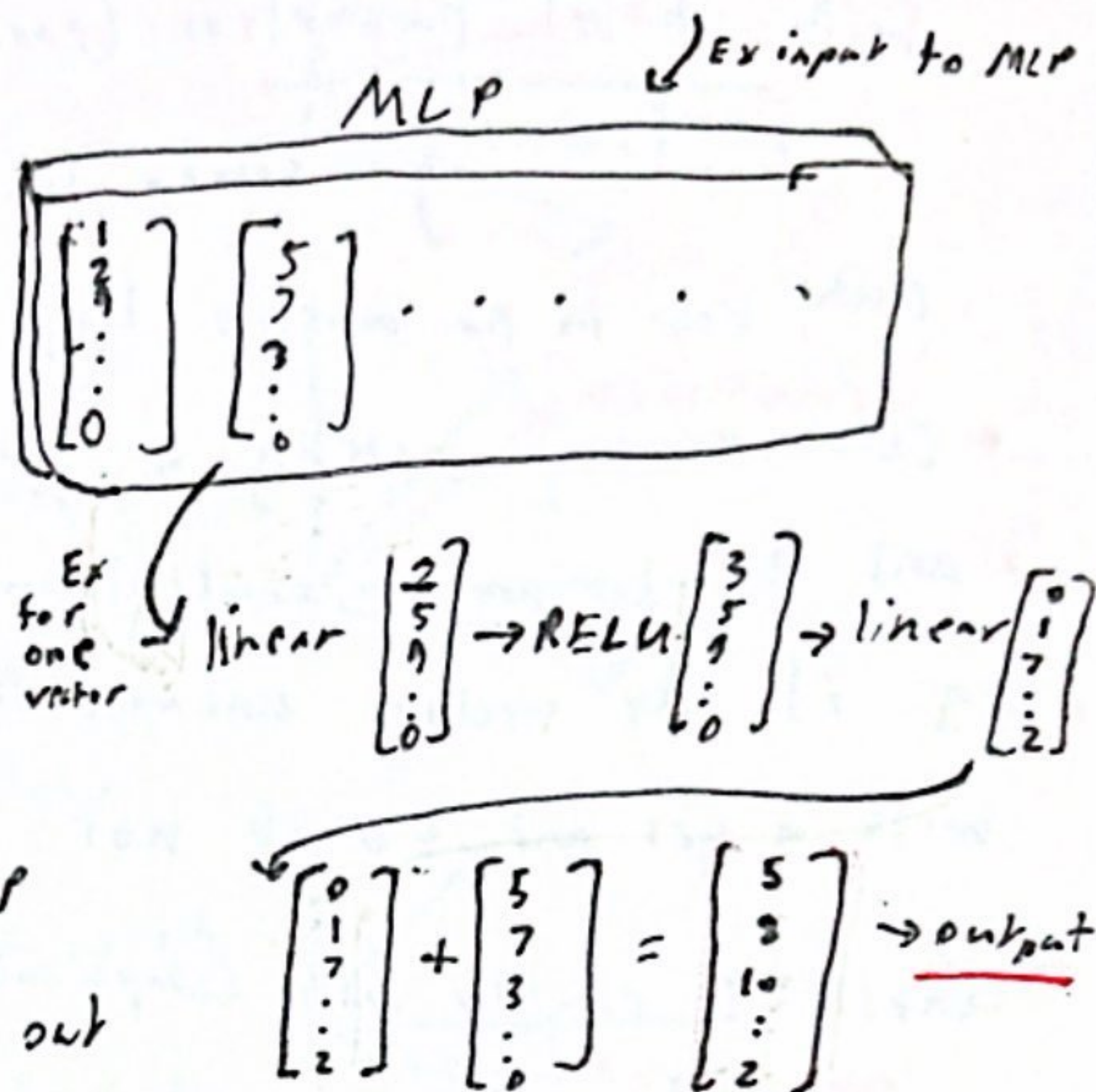# MLP (P-3) (LLM notes)

- inside a MLP: the sequence of vectors which where each vector is associated with a input token.

- Then each vector in that sequence inside the MLP goes through a sequence of operations

- at the end we get a new updated vector this gets added to the original input vector in the MLP and that sum is the output that flows out

- you then apply this operation to all vectors in the sequence ie for every token in the input, all this is in parallel and indipendent of one another

- What this means is that if the input "Jordan" includes context of first name michael last name jordan then the last operation before adding will produce a vector for the direction "Basket-Ball" which will be added to the vector for jordan ie the input, this produces the result MJ + basket Ball which gives the token jordan context of being associated with basket ball

MLP:

$$\left[ \begin{array}{c} E \\ \end{array} \right] \rightarrow linear \left[ \begin{array}{c} \\ \end{array} \right] \rightarrow Relu \left[ \begin{array}{c} \\ \end{array} \right] \rightarrow linear \left[ \begin{array}{c} \\ \end{array} \right] + \underline{\quad} \rightarrow output$$

64

Ex: Micheal Jordan plays...

$$\downarrow \qquad \downarrow \qquad \qquad \downarrow$$

$$\left[ \begin{array}{c} 1 \\ 2 \\ 9 \\ \vdots \\ 0 \end{array} \right] \quad \left[ \begin{array}{c} 5 \\ 7 \\ 3 \\ \vdots \\ 0 \end{array} \right] \quad \left[ \begin{array}{c} \vdots \\ \end{array} \right] \quad \ddots$$

MLP ← Ex input to MLP

$$\left[ \begin{array}{c} 1 \\ 2 \\ 9 \\ \vdots \\ 0 \end{array} \right] \left[ \begin{array}{c} 5 \\ 7 \\ 3 \\ \vdots \\ 0 \end{array} \right] \quad \cdots \quad \cdot \cdot$$

Ex for one vector

$$\rightarrow linear \left[ \begin{array}{c} 2 \\ 5 \\ 9 \\ \vdots \\ 0 \end{array} \right] \rightarrow RELU \left[ \begin{array}{c} 3 \\ 5 \\ 9 \\ \vdots \\ 0 \end{array} \right] \rightarrow linear \left[ \begin{array}{c} 0 \\ 1 \\ 7 \\ \vdots \\ 2 \end{array} \right]$$

$$\left[ \begin{array}{c} 0 \\ 1 \\ 7 \\ \vdots \\ 2 \end{array} \right] + \left[ \begin{array}{c} 5 \\ 7 \\ 3 \\ \vdots \\ 0 \end{array} \right] = \left[ \begin{array}{c} 5 \\ 3 \\ 10 \\ \vdots \\ 2 \end{array} \right] \rightarrow output$$

Note vector jordan has context of MJ from attention block (see last page)

Micheal Jordan ↓

$$\left[ \begin{array}{c} 5 \\ 7 \\ 3 \\ \vdots \\ 0 \end{array} \right] \xrightarrow{linear} \left[ \begin{array}{c} \vdots \\ \end{array} \right] \xrightarrow{Relu} \left[ \begin{array}{c} \vdots \\ \end{array} \right]$$

linear

$$\rightarrow \left[ \begin{array}{c} 0 \\ 1 \\ 7 \\ \vdots \\ 2 \end{array} \right] + \left[ \begin{array}{c} 5 \\ 7 \\ 3 \\ \vdots \\ 0 \end{array} \right] = output \left[ \begin{array}{c} 5 \\ 10 \\ \vdots \\ 2 \end{array} \right]$$

Basket Ball      MJ