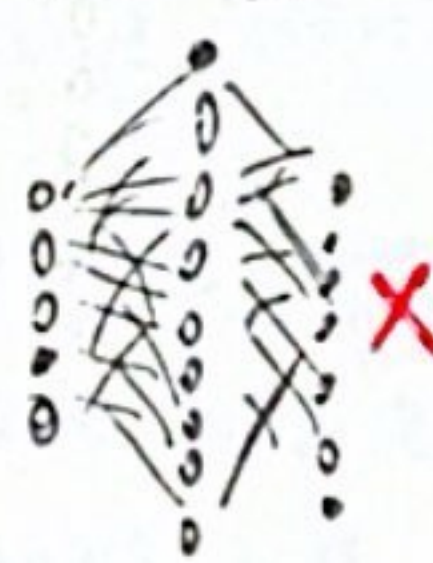


# MLP (P-7) (Klm notes)

## - Superposition + Recap

- in real LLM facts are not stored exactly as described in the past few pages, while the rows of the Matrix can be thought as directions or directions in embedding space, that helps get facts for context. and a neuron tells you how much a given vector aligns with some direction. and the columns of the second matrix tell you what will be added to the result if that neuron is active. But individual neurons rarely represent a single clean feature or a single row of  $M_2$  does not always ask a clean single direction. (like it does in our ex. rather multiple neurons come together to store a fact or a question. If it might take many rows to answer one question or store some fact (Fig 1))

Michael Jordan



Michael Jordan



- This happens theoretically by superposition

in short it allows to store ideas > dimensions if we had a basketball direction but that 1 dimension can store other closely related ideas like a sport will balls for ex. hence this can explain why LLM scale so well but are still hard to interpret, and the features in the NN would have to come together to form facts not just one lit neuron (fig 2) that specific combination of 68 neurons is a superposition.