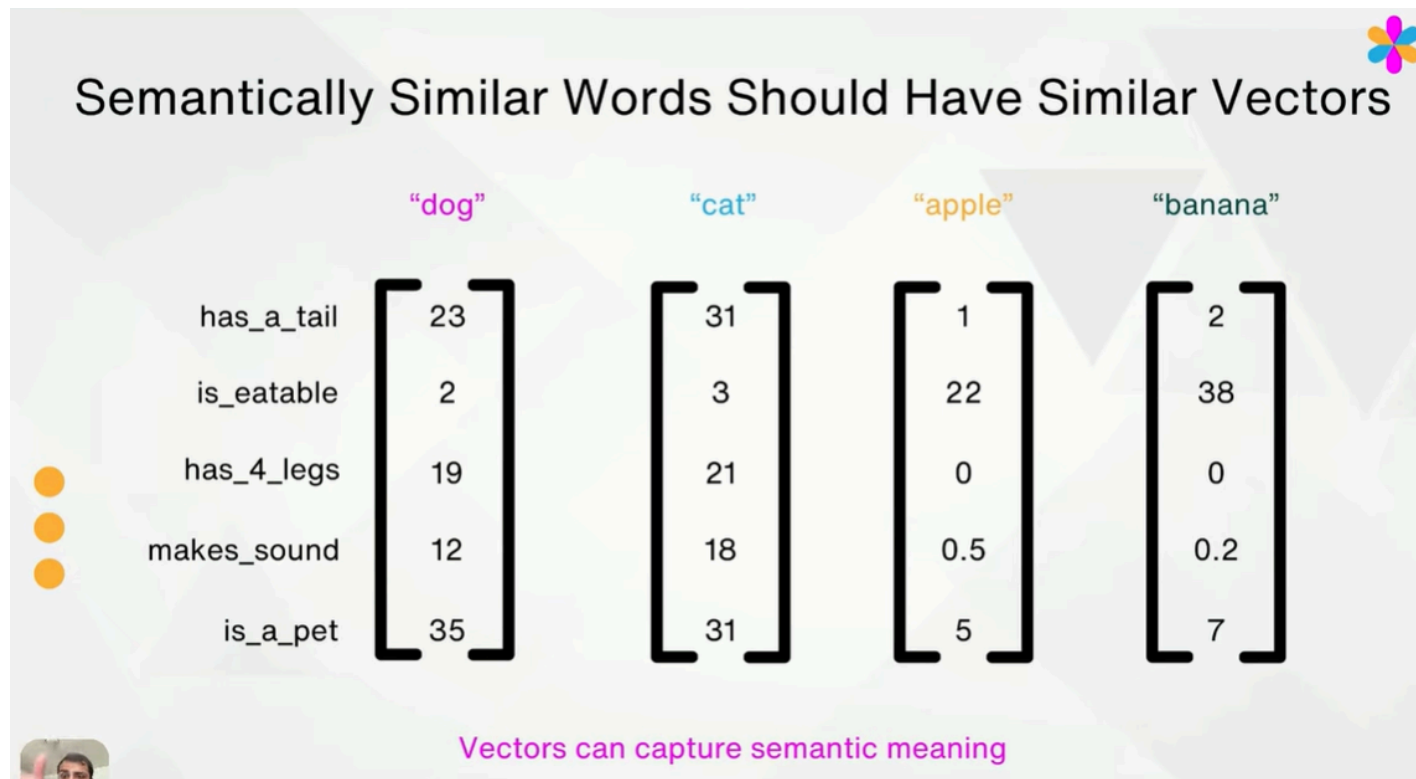# Lecture 10: Token Embeddings

Why not just use the token ids for each word?

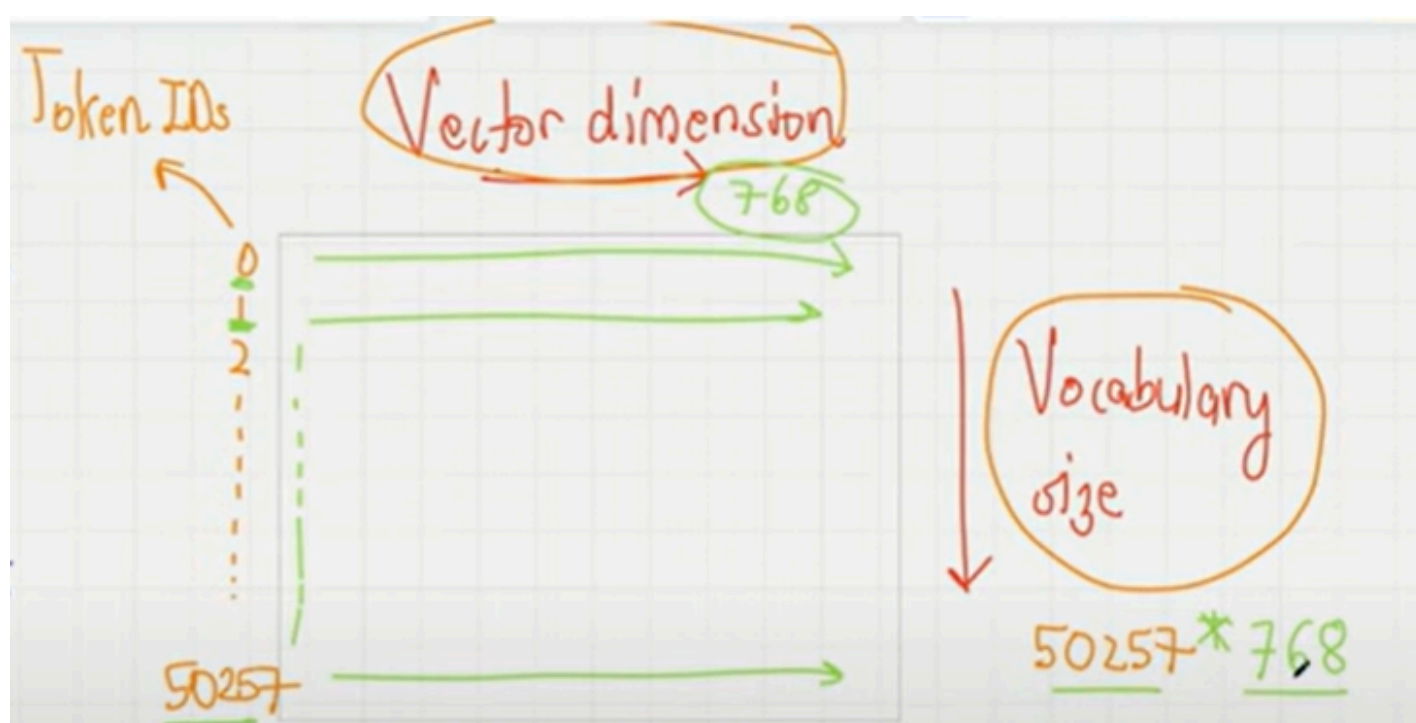→ Because, just randomly assigning token ids cannot capture the semantic meaning of words.



Whatever is higher in dog, is higher in cat. Similarly for the low numbers. So it means dog and cat are similar things.

Whatever is higher in apple, is higher in banana. Similarly for the low numbers. So it means apple and banana are similar things.

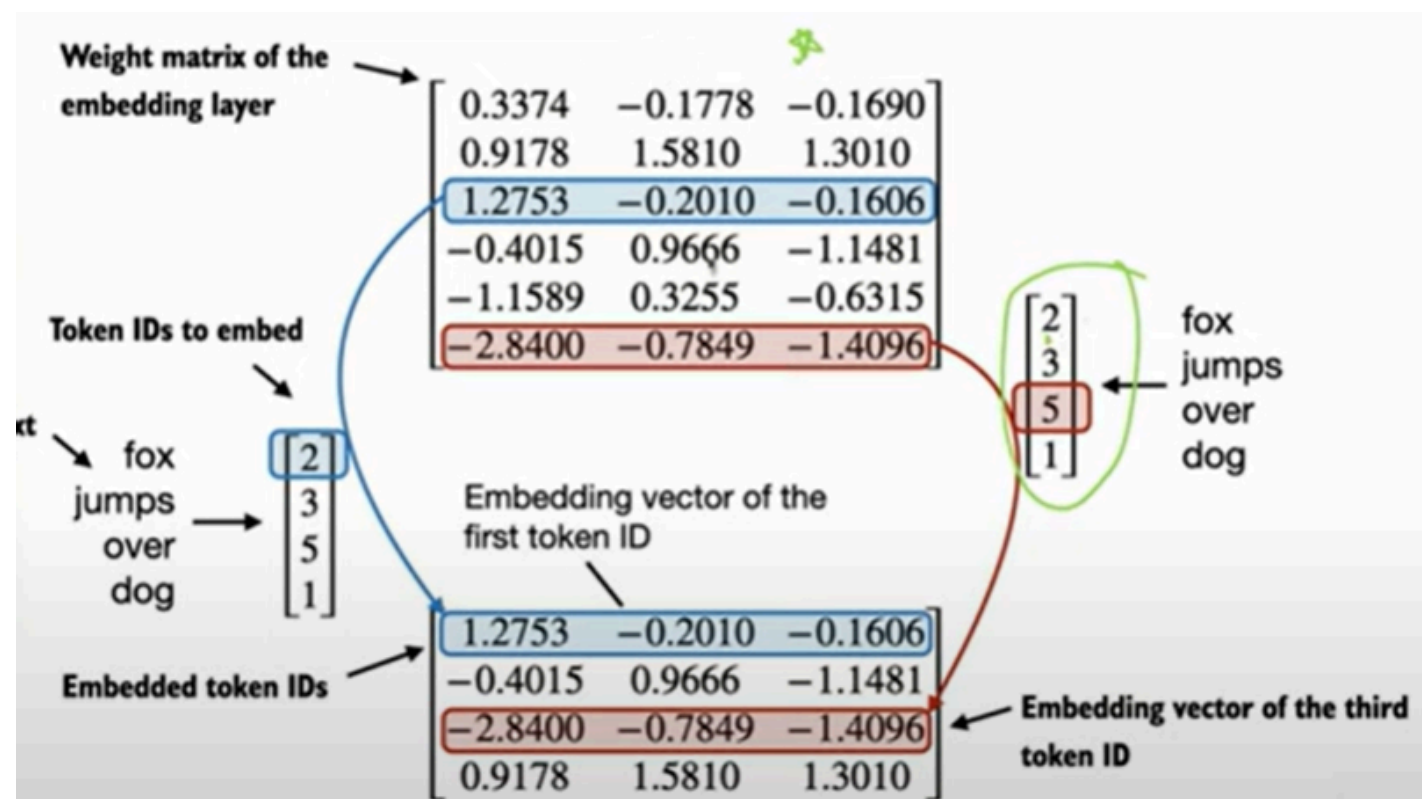These are called as vector embeddings or token embeddings.

We have to train a neural network to create vector embeddings.

In GPT-2 model:



This is why embedding vector is also called a **lookup table** that stores embeddings of a fixed dictionary and size.

The embedding layer is a **lookup operation** that retrieves rows from the embedding layer weight matrix using a token ID.

Both embedding layer and Neural network layer (linear layer) lead to the same output.

But embedding layer is used instead because it is much more computationally efficient, neural network has so many unnecessary multiplication with zeros.