# TOKEN EMBEDDINGS

# Representing Words Numerically

- Computers need numerical representation of words

- How can we represent words in numbers?

Can we assign random numbers to each word?

| | | |
|---|---|---|
| "cat" | $\longrightarrow$ | 34 |
| "book" | $\longrightarrow$ | 2.9 |
| "tablet" | $\longrightarrow$ | -20 |
| "kitten" | $\longrightarrow$ | -13 |

# The Problem With Using Random Numbers

"cat"     ⟶     34

"book"    ⟶     2.9

"tablet"  ⟶     -20

"kitten"  ⟶     -13

"cat" and "kitten" are semantically related.
However the associated numbers 34 and -13 cannot capture this relation.
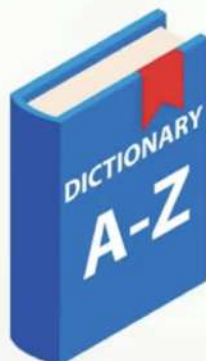
# What About One-Hot Encoding?

1) Create a dictionary of words
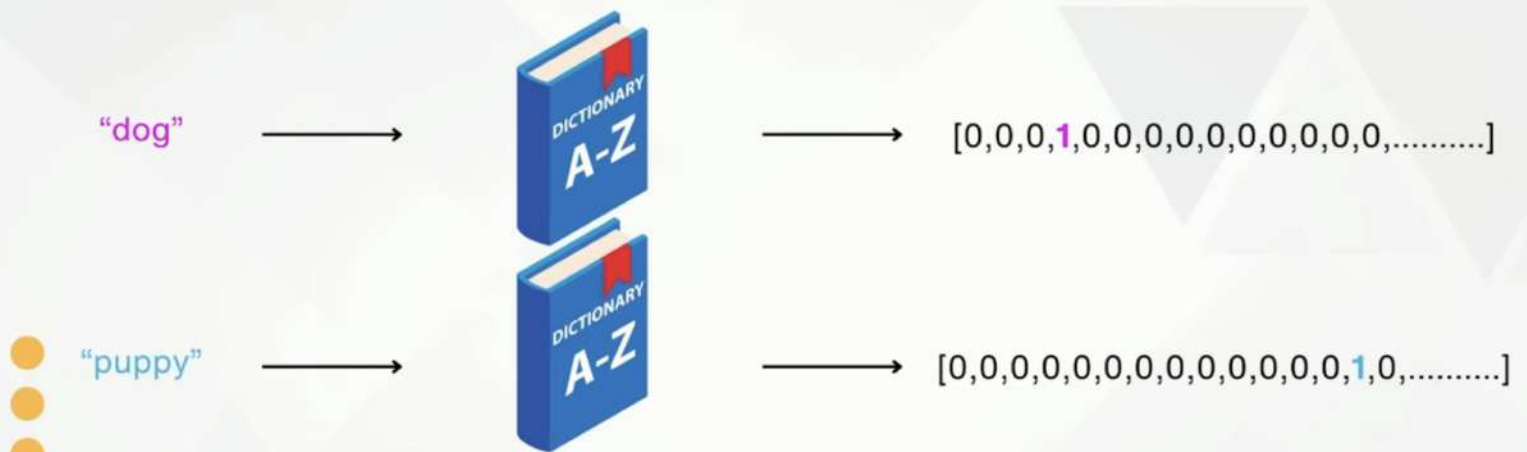2) Assign sequential one-hot encoding to each word

"dog" → DICTIONARY A-Z → [0,0,0,1,0,0,0,0,0,0,0,0,0,0,0,..........]

# The Problem With One-Hot Encoding

"dog" $\longrightarrow$  $\longrightarrow$ [0,0,0,**1**,0,0,0,0,0,0,0,0,0,0,0,..........]

"puppy" $\longrightarrow$  $\longrightarrow$ [0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,**1**,0,..........]

One-hot encoding also fails to capture semantic relationship

# Semantically Similar Words Should Have Similar Vectors

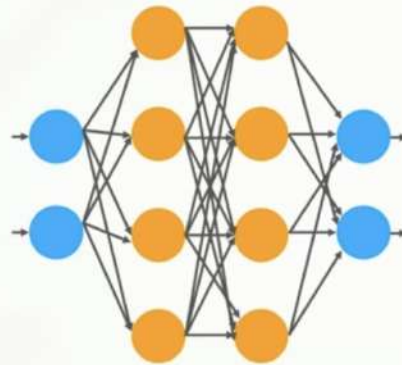|  | "dog" | "cat" | "apple" | "banana" |
|---|---|---|---|---|
| has_a_tail | 23 | 31 | 1 | 2 |
| is_eatable | 2 | 3 | 22 | 38 |
| has_4_legs | 19 | 21 | 0 | 0 |
| makes_sound | 12 | 18 | 0.5 | 0.2 |
| is_a_pet | 35 | 31 | 5 | 7 |

Vectors can capture semantic meaning

# We Can Train a Neural Network To Create Vector Embedding

"dog"

"cat"

"apple"

"banana"

[23, 2, 19, 12, 35]

[31, 3, 21, 18, 31]

[1, 22, 0, 0.5, 5]

[2, 38, 0, 0.2, 7]