# Understanding Large Language Models

## Generative AI and Prompt Engineering
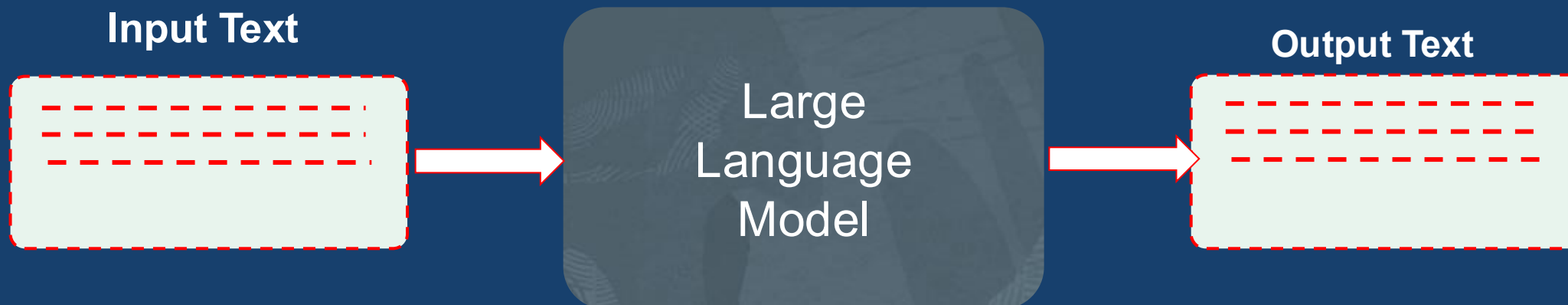
Ram N Sangwan

- Large Language Models
- Transformers, Sequence Models, RNN, Encoder - Decoder
- Embeddings, Tokenization

# Large Language Models (LLMs)

- Understand, generate and process human language at massive scale.
- Designed for sequence-to-sequence tasks such as machine translation

**Input Text**

**Output Text**

Large Language Model

# Large Language Models Example

**1** Translate "How are you" into French.

**2** What is the capital of France?

**3** Write an essay on French Revolution.

Large Language Model

**1** "Comment allez-vous?"

**2** The capital of France is Paris.

**3** **Title:** The French Revolution: A Turning Point in History
**Introduction:** The French Revolution, which took place between 1789 and 1799, stands as one of the most influential and transformative events in human history...

# LLMs – What you need to Know.

- LLM Architectures
  - What else can LLMs do?
- Prompting and Training
  - How do we affect the distribution over the vocabulary?
- Decoding
  - How do LLMs generate text using these distributions?

# Model Size and Parameters

**Model Size:** Memory required to store the model's parameters.

**Parameters:**

- Tuneable variables that are adjusted during training to minimize the difference between predicted outputs and actual labels (e.g., correct answers).
- Used to represent relationships between input tokens and output text.
- The goal is to find the optimal values such that the model can accurately predict new, unseen inputs.

**Two Types of Parameters**

1. **Learned Embeddings** are dense vectors that represent each token in a fixed-dimensional space. Learned embeddings capture semantic relationships between tokens and are used as input features for the model.
2. **Model Weights** - are the adjustable values within the neural network's layers, such as fully connected layers, convolutional layers, or recurrent layers.
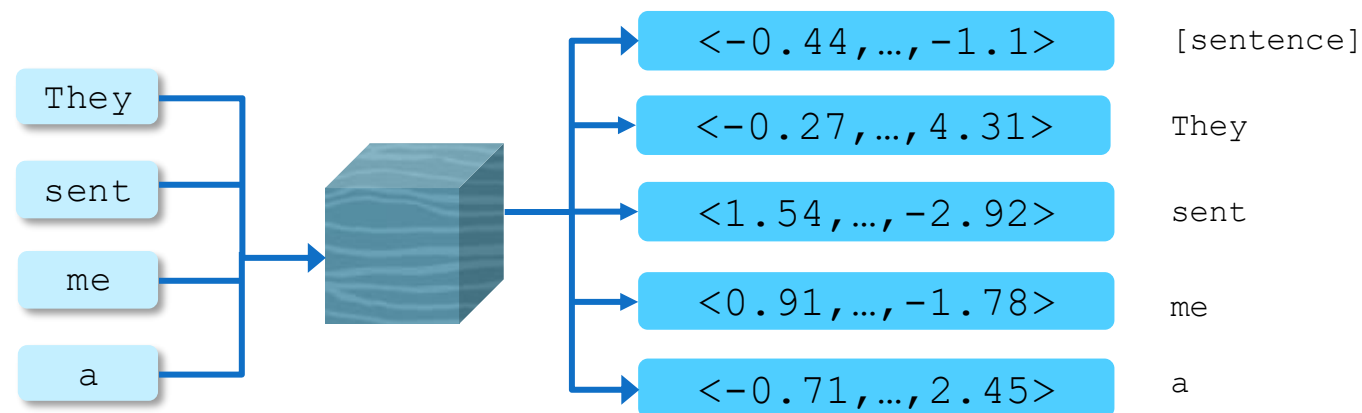


2018
BERT

2020
GPT-3

2022
PaLM

3.7B Tokens
240M Parameters

499B Tokens
175B Parameters

780B Tokens
540B Parameters

# Encoders

**Encoder –** models that convert a sequence of words to an embedding (vector representation)

**Examples**

*Embed-light, BERT, RoBERTA, DistillBERT, SBERT,…*



They
sent
me
a

<-0.44,…,-1.1>    [sentence]
<-0.27,…,4.31>    They
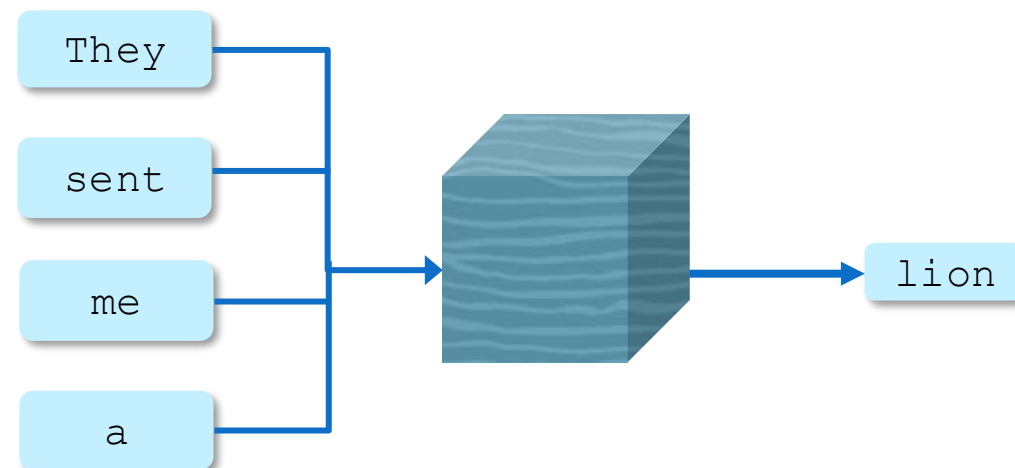<1.54,…,-2.92>    sent
<0.91,…,-1.78>    me
<-0.71,…,2.45>    a

**Primary uses:** embedding tokens, sentences, & documents

# Decoders

- **Decoder –** models take a sequence of words and output next word


- **Examples**

- *GPT-4, Llama, BLOOM, Falcon, …*

They

sent

me

a

lion

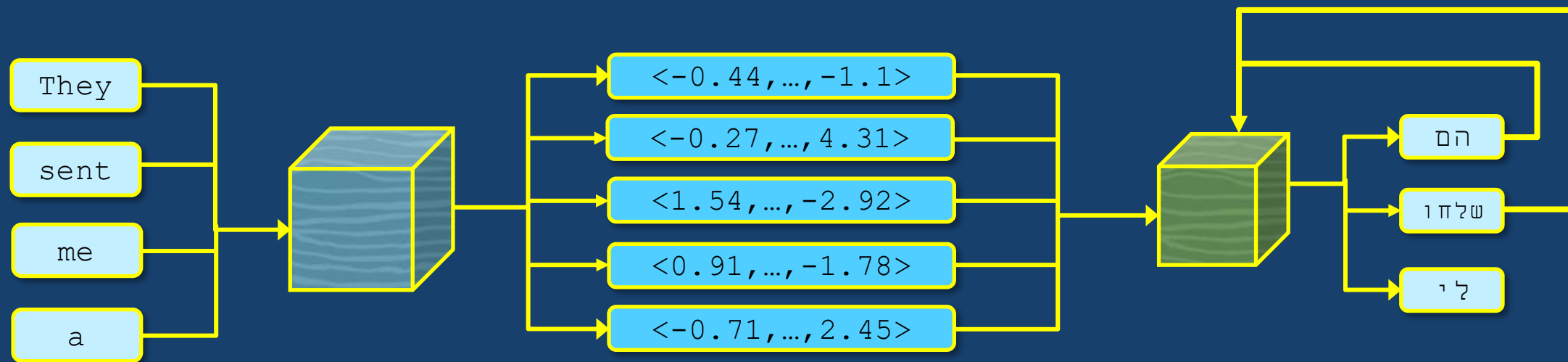**Primary uses:** text generation, chat-style models, (including QA, etc…)

# Encoders - Decoders

**Encoder-decoder** - encodes a sequence of words and use the encoding +  to output a next word

**Examples**
*T5, UL2, BART,…*

# Understanding Language for Machines Can be Tricky
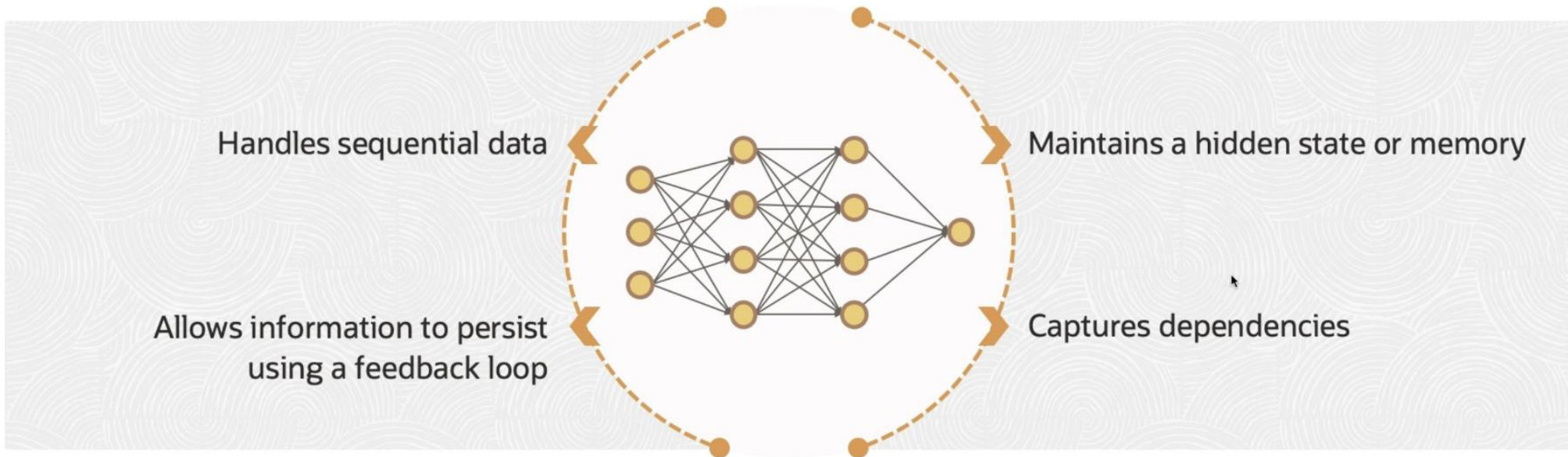
Jane is doing the throwing

Dog is doing the fetching

It refers to the frisbee

*Jane threw the frisbee and her dog fetched it.*

As human, we understand easily that "it" refers to the frisbee. But for a machine, it can be tricky.
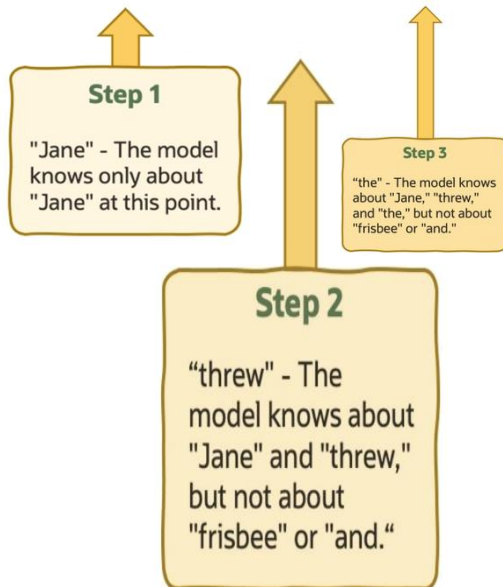
# Recurrent Neural Network (RNN)



Recurrent Neural Network works on the principle of saving the output layer and feeding this back to the
input in order to predict the output of the layer.

# RNNs Struggle with Long-Range Dependencies.

Jane threw the frisbee and her dog fetched it.

**Step 1**

"Jane" - The model knows only about "Jane" at this point.

**Step 3**

"the" - The model knows about "Jane," "threw," and "the," but not about "frisbee" or "and."

**Step 2**

"threw" - The model knows about "Jane" and "threw," but not about "frisbee" or "and."

- RNNs process input data one element at a time, maintaining a hidden state

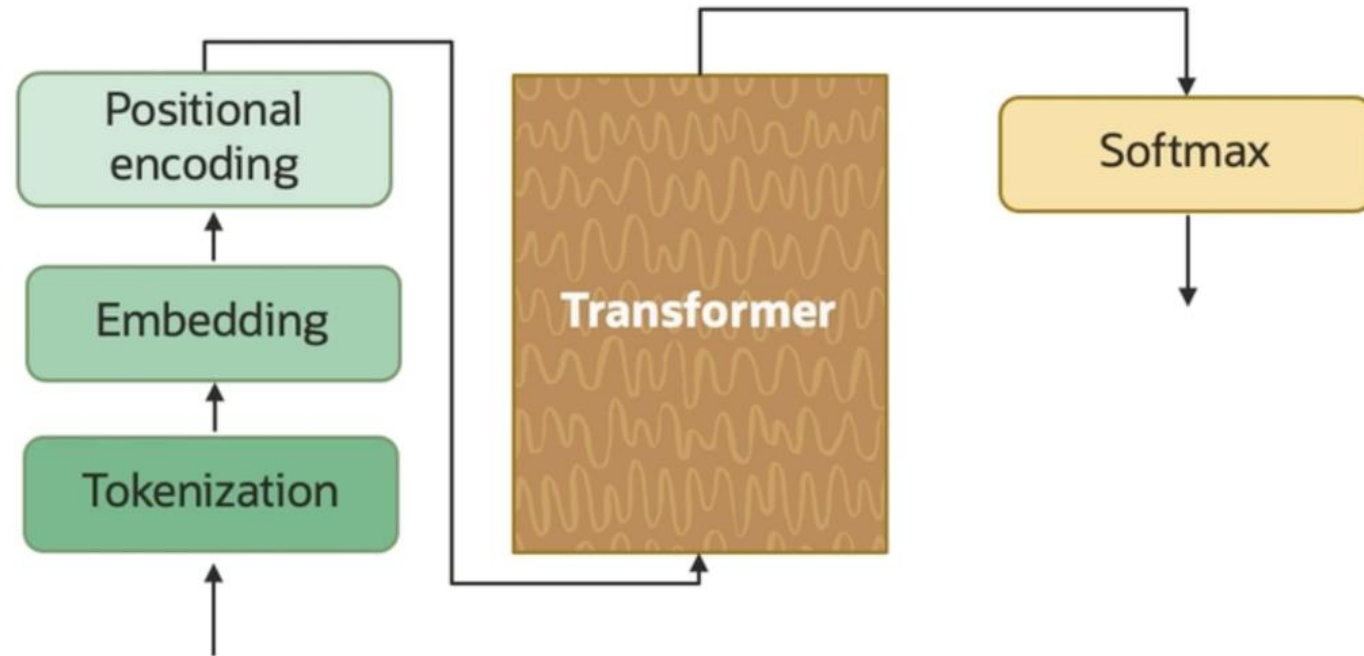# Transformers Understand Sentences as a Whole

Jane
threw
the
frisbee
and
her
dog
fetched
It

Jane
threw
the
frisbee
and
her
dog
fetched
It

- Can look at all the words in the sentence at the same time and understand how they relate to each other

- Can simultaneously understand the connections between "Jane" and "dog," even though they are far apart in the sentence

# Understanding Transformer Architecture



*Jane threw the frisbee and her dog fetched it.*

# Tokenization

"Jane", "threw", "the", "frisbee", "and",
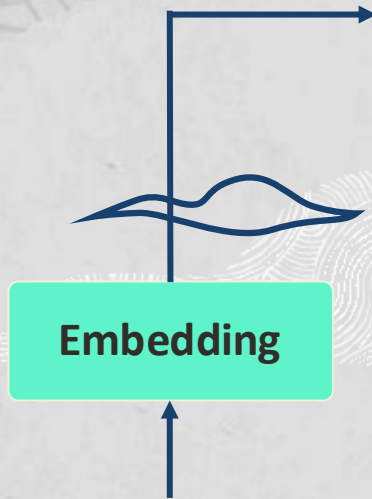"her","dog", "fetched", "it","."

**Tokenization**

**Jane threw the frisbee and her dog fetched it.**

- The Sentence is broken down into smaller pieces called tokens.

- The choice of how to break down the sentence depends on the specific tokenizer used.
- Learn to Create a Tokenizer from Scratch *here*

# Embedding

| | | |
|------|------|-----|
| Jane | 2.1 | ... |
| threw | 4.2 | ... |
| the | 6.8 | ... |
| frisbee | 9.2 | ... |
| and | -3.2 | ... |
| her | 2.3 | ... |
| dog | 2.1 | ... |
| fetched | 1.2 | ... |
| it | 1.5 | ... |
| . | 2.3 | ... |

**Embedding**

*"Jane", "threw", "the", "frisbee", "and", "her","dog", "fetched", "it","."*

- Each token is then converted into a numerical form (a vector) that the model can understand.
- These vectors are created in a way that captures the meaning of the word.
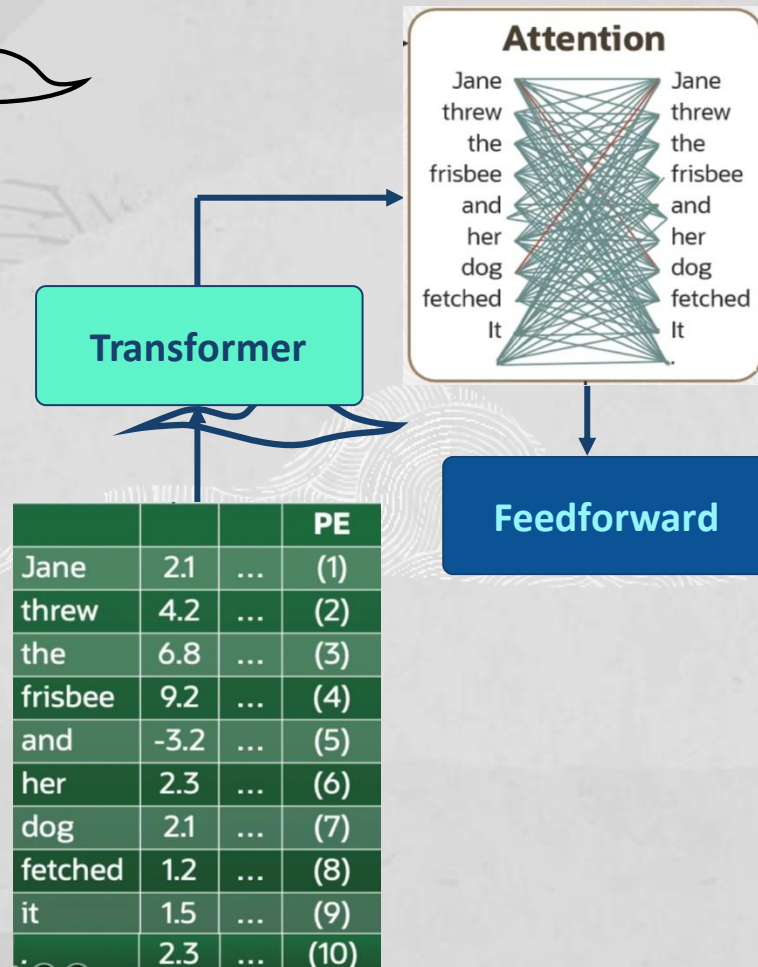- More on Vectors and Embeddings [here](here)

# Positional Encoding

—

| | | | PE |
|---|---|---|---|
| Jane | 2.1 | … | (1) |
| threw | 4.2 | … | (2) |
| the | 6.8 | … | (3) |
| frisbee | 9.2 | … | (4) |
| and | -3.2 | … | (5) |
| her | 2.3 | … | (6) |
| dog | 2.1 | … | (7) |
| fetched | 1.2 | … | (8) |
| it | 1.5 | … | (9) |
| . | 2.3 | … | (10) |

**Positional encoding**

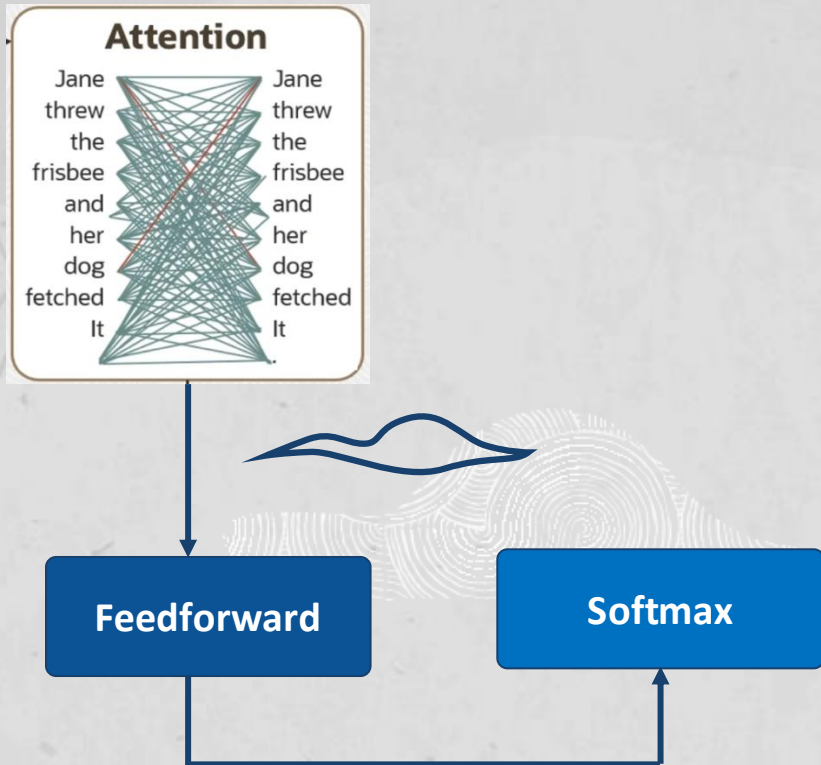| | | |
|---|---|---|
| Jane | 2.1 | … |
| threw | 4.2 | … |
| the | 6.8 | … |
| frisbee | 9.2 | … |
| and | -3.2 | … |
| her | 2.3 | … |
| dog | 2.1 | … |
| fetched | 1.2 | … |
| it | 1.5 | … |
| . | 2.3 | … |

- The Position Encoding layer represents the position of the word.

- The Model needs to know the order of the words in the sentence.

- The model adds information about the position of each word.
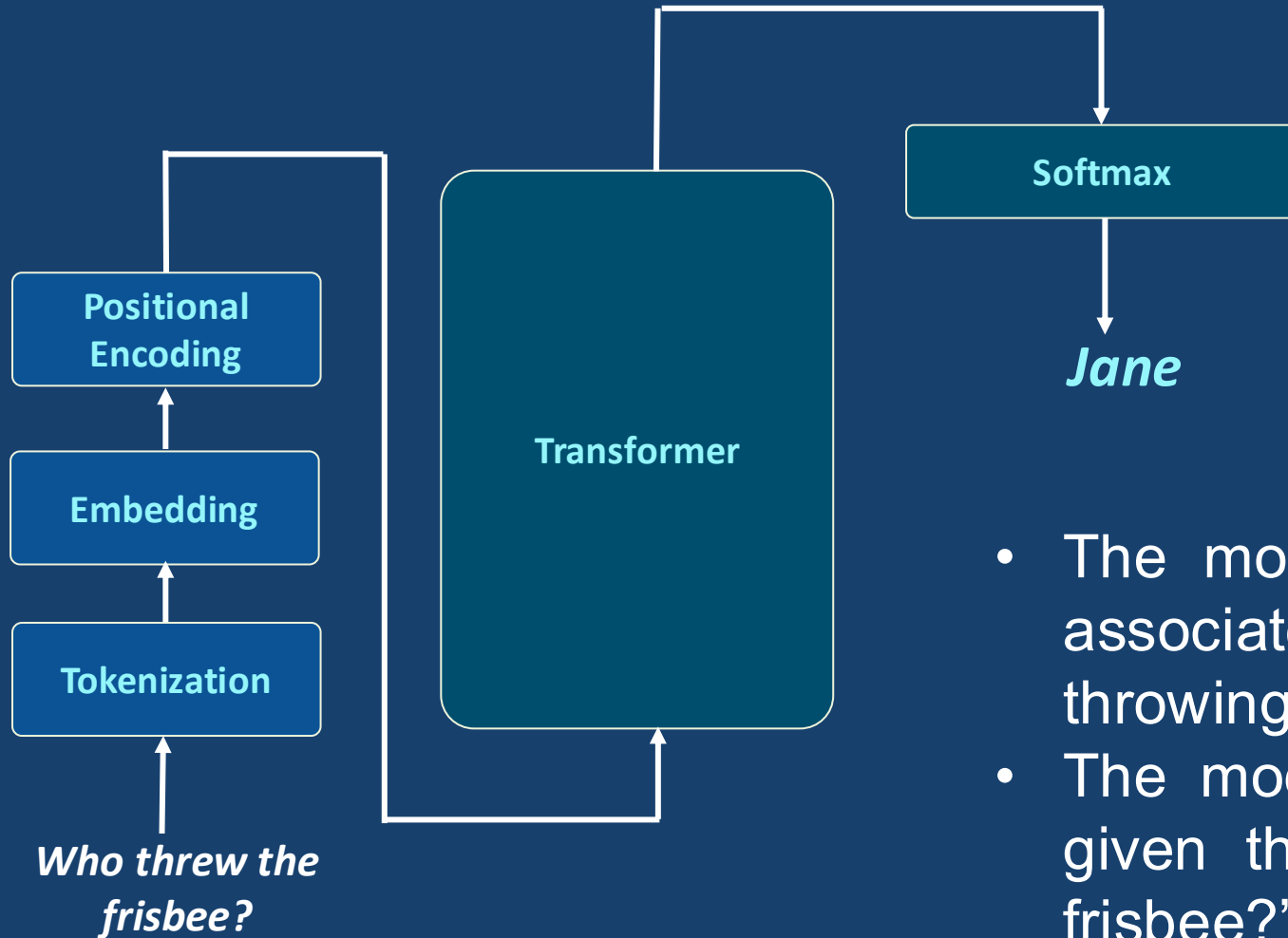
# Transformer



- Attention: Helps understand the context of each word.

- Feedforward: Applies a specific function to each word individually.

# Transformer



- The model generates a list of scores for each word in the vocabulary.

- The model uses the Softmax function to turn these scores into probabilities.

# Transformer

```
Positional
Encoding
    ↑
Embedding
    ↑
Tokenization
    ↑
```
*Who threw the frisbee?*

**Transformer**

**Softmax**
↓
*Jane*

Output from the transformer's Softmax layer is a <span style="color:red">probability distribution</span> across the entire dictionary of words.

- The model realizes that "Jane" is associated with the action of throwing the frisbee.
- The model predicts the next word given the context "who threw the frisbee?"

# Thank You

— Watch a Video on LLS Explained [here](#)