



Understanding Positional Encoding in Transformers

How Transformers Learn Word Order Without RNNs—A Simple, Intuitive Guide for Indian Engineering Students

Why Do We Need Positional Encoding?

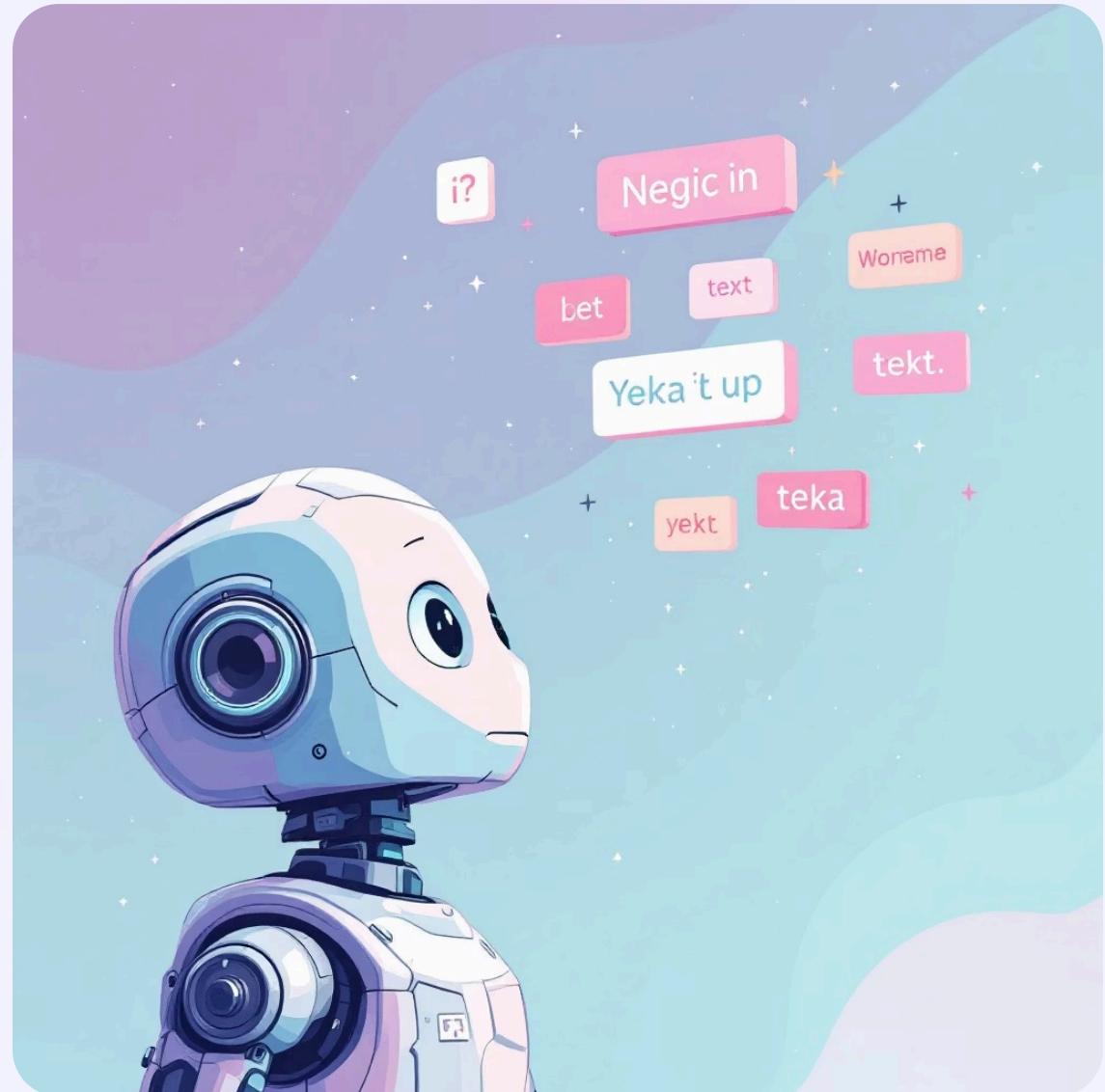
The Problem

Transformers process all words **simultaneously** (in parallel). Unlike RNNs, they don't naturally know the **order** of words.

Real-World Impact

"The cat sat" vs "Sat the cat"

→ Same words, completely different meanings!



- ❑ ✨ **The Solution:** We add Positional Encoding (PE) to tell the model *which word came first, second, third*, and so on –giving it a sense of sequence and structure.



What Exactly Is Positional Encoding?

It's a Vector

A mathematical representation of each word's position in a sentence

Added to Embeddings

Combined with word embeddings to include both order and meaning

Enables Understanding

Gives each word vector knowledge of its meaning *and* its position



The Core Formula

Input to Transformer = Word Embedding + Positional Encoding



Example: Adding PE to Words

Let's see how positional encoding works with a simple sentence: "The Cat Sat"

Position	Word	Word Embedding	Positional Encoding	Final Input
0	The	E_the (512)	P ₀ (512)	E_the + P ₀
1	Cat	E_cat (512)	P ₁ (512)	E_cat + P ₁
2	Sat	E_sat (512)	P ₂ (512)	E_sat + P ₂

💡 Key Insight

Adding (not stacking) keeps the same vector size—typically 512 dimensions. Each dimension now elegantly mixes **word meaning + position information**.





The Mathematical Formula

Transformers don't just assign "Position 1, 2, 3..."—they use **smooth wave patterns** based on sine and cosine functions 🎵

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$



pos

Word position in the sentence (0, 1, 2, 3, ...)



i

Index of dimension inside the word vector (0 to $d_{model}/2$)



d_{model}

Embedding dimension size (typically 512)

- ❑ Together, these sine and cosine functions create a **unique wave pattern** for each word position—like a mathematical fingerprint!

Breaking It Down Simply

01

Vector Dimensions

Each word has a 512-dimensional vector—think of each dimension as a "slot" holding a number

02

Wave Assignment

Each slot gets a sine or cosine value based on its position (pos) and index (i) in the vector

03

Slow vs Fast Waves

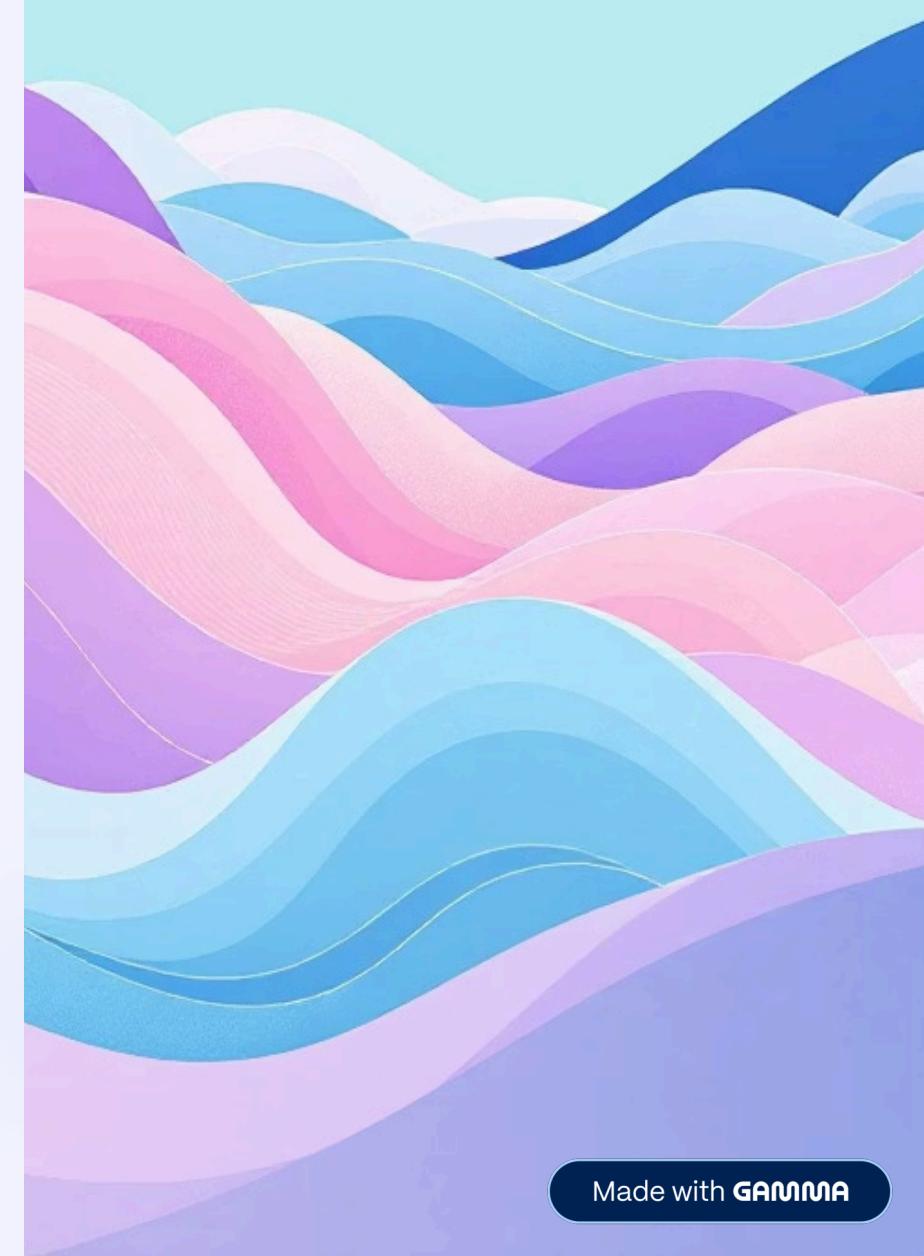
Small i values create slow waves (capturing big position differences); large i values create fast waves (capturing fine position details)

04

Unique Patterns

Every position gets a unique combination of these waves—like a mathematical signature for its location

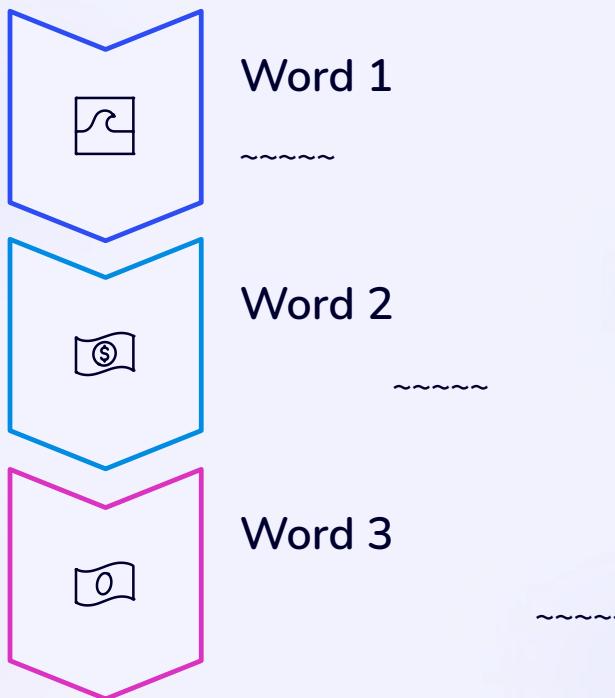
🎵 **Musical Analogy:** Think of it like different instruments playing together—some play low, slow notes (bass), others play high, quick notes (flute). Together, they create a unique "song" for each word position!





Visual Intuition: The Wave Patterns

Imagine sine and cosine waves moving smoothly across positions, creating unique patterns:



Similar Positions

Words close together have similar wave patterns—helping the model understand nearby relationships

Distant Positions

Words far apart have different patterns—clearly distinguishing their relative positions



This elegant design helps the model **"see" word order** naturally, even without sequential loops like RNNs use!

Why Sine and Cosine Functions?



Smooth & Continuous

Nearby words have similar encodings, preserving relative position relationships naturally



Works for Any Length

Wave patterns repeat naturally, handling sentences of any length without retraining



Math Friendly

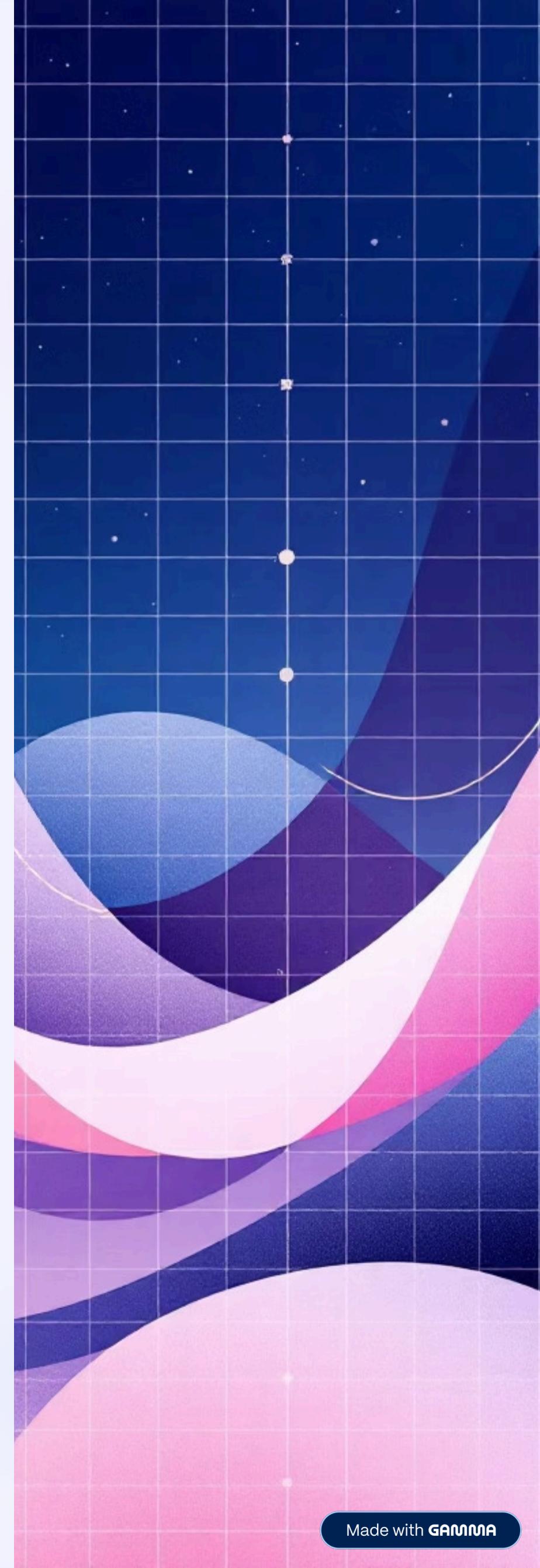
Values always stay between -1 and +1, making training stable and computations efficient



Relative Distance

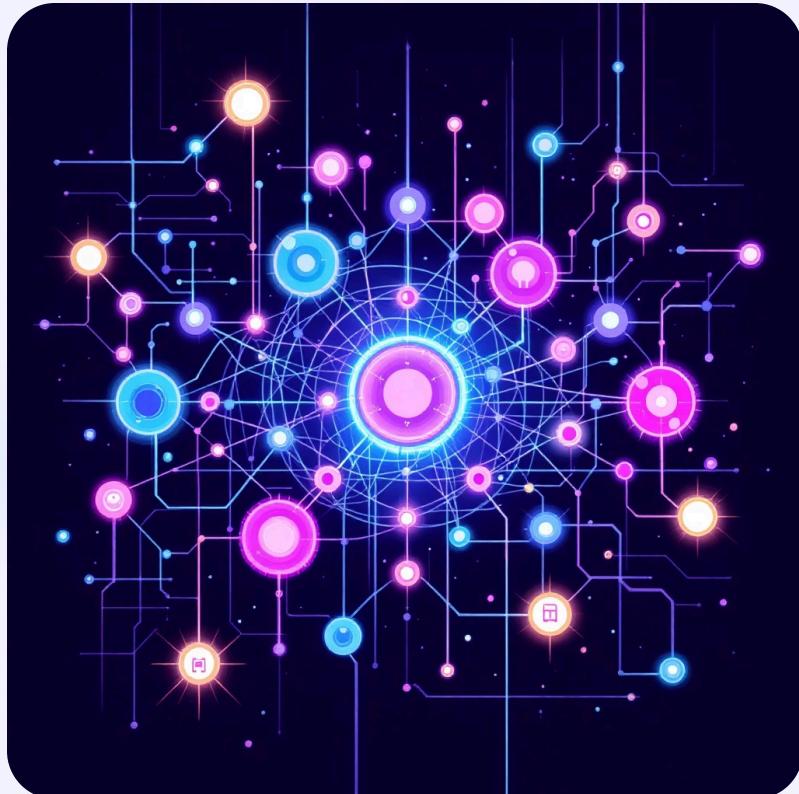
The model learns how far apart words are, not just their absolute positions

They provide a mathematical rhythm to encode order effectively—like giving each word a unique musical note that harmonises with its neighbours!





How It Empowers the Transformer



What the Model Now Knows



Sequential Order

Which word came before or after another in the sentence



Relative Distance

How far apart words are from each other



Long-Range Dependencies

Relative order across long, complex sentences

- ❑ **The Result:** Transformers can now understand **grammar, structure, and context** without using sequential RNNs—enabling parallel processing and much faster training!

Summary: Positional Encoding at a Glance

1

Core Goal

Give each word position awareness in parallel processing

2

Method Used

Add sine–cosine wave patterns to word embeddings

3

Key Formula

Uses pos, i, and d_{model} to create unique patterns

Why It's Important

Adds crucial order information for language understanding—without it, "cat sat the" would mean the same as "the cat sat"!

Final Result

The model **"sees" sentence structure** naturally, enabling state-of-the-art performance in NLP tasks



In Short:

Positional Encoding = The GPS of a Transformer Model

