

# BERT and its Applications

This file is meant for personal use by diegorosenberg@gmail.com only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

# Agenda

- Key Questions
- What is BERT?
- What are the variants of BERT?
- How is BERT trained?
- What are the downstream applications of BERT?

# Key Questions

**What is BERT?**

**What are the variants of BERT?**

**How is BERT trained?**

**What are the downstream applications of BERT?**

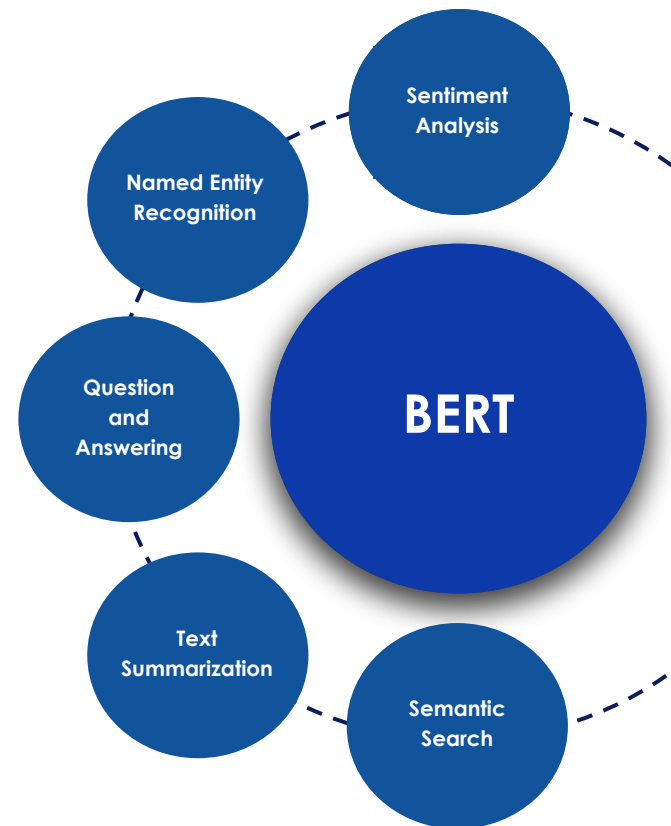
# What is BERT?

BERT stands for **Bidirectional Encoder Representations Transformer**

It was introduced in a **research paper** by **Jacob Devlin** and his **colleagues** at Google Research in **2018**.

BERT is an **encoder only transformer** model that leverages a **stack of encoders** to get an deeper **understanding of language context**

BERT **revolutionized** Natural Language Processing (NLP) by exhibiting **state-of-the-art performance** in a variety of NLP tasks



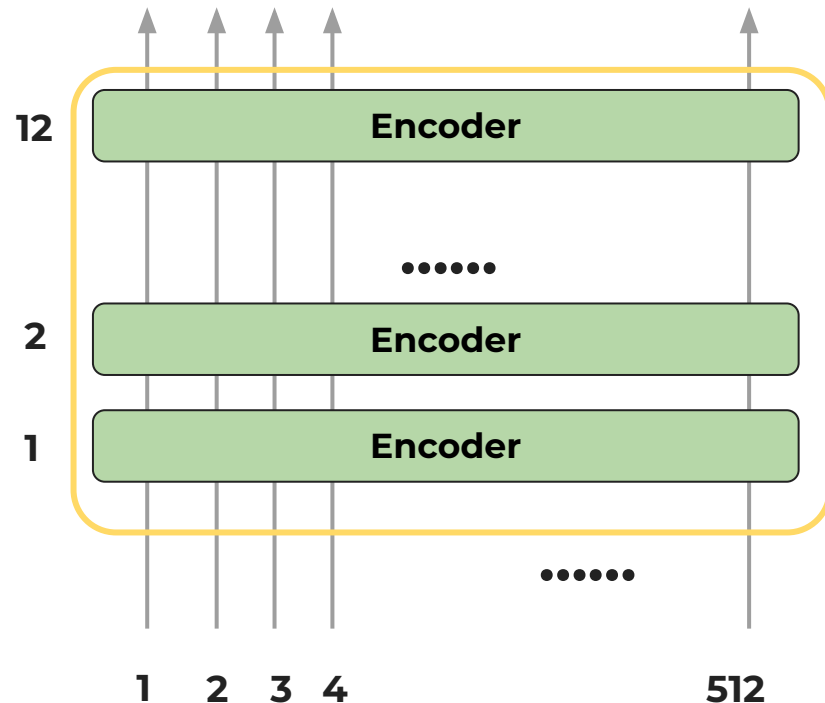
This file is meant for personal use by [diegorosenberg@gmail.com](mailto:diegorosenberg@gmail.com) only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

# What are the variants of BERT?

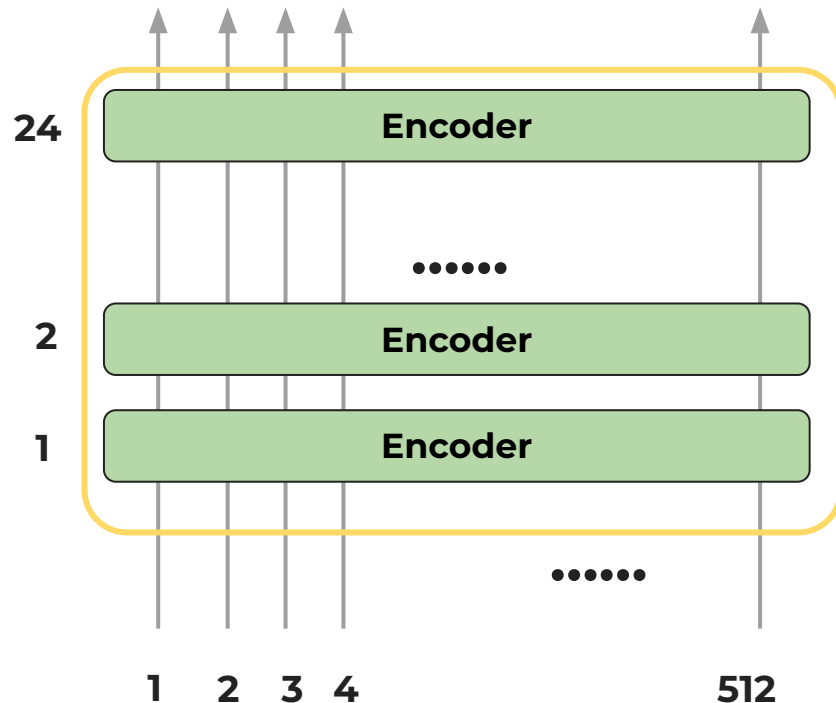
	<b>BERT<sub>BASE</sub></b>
# parameters	110M
# encoders	12
# attention heads	12
word embeddings dimension	768



BERT<sub>BASE</sub> has a **limitation** on the **number of input tokens** - it can take a maximum of **512** input tokens

# What are the variants of BERT?

	<b>BERT<sub>LARGE</sub></b>
# parameters	340M
# encoders	24
# attention heads	16
word embeddings dimension	1024



BERT<sub>LARGE</sub> has a similar limitation on the number of input tokens as BERT<sub>BASE</sub>

This file is meant for personal use by [diegorosenberg@gmail.com](mailto:diegorosenberg@gmail.com) only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

# How is BERT trained?

BERT's training essentially consists of two stages - a **Pre-training** Stage and a **Fine-tuning** Stage

## Pre-training Stage

The model builds a foundational understanding of language

Involves exposing the model to a vast amount of text data, where it learns about word relationships and gains contextual knowledge



## Fine-tuning Stage

The model adapts its foundational understanding of language to perform better on downstream tasks

Involves adjustment of model parameters (weights) to specialize in a particular task.

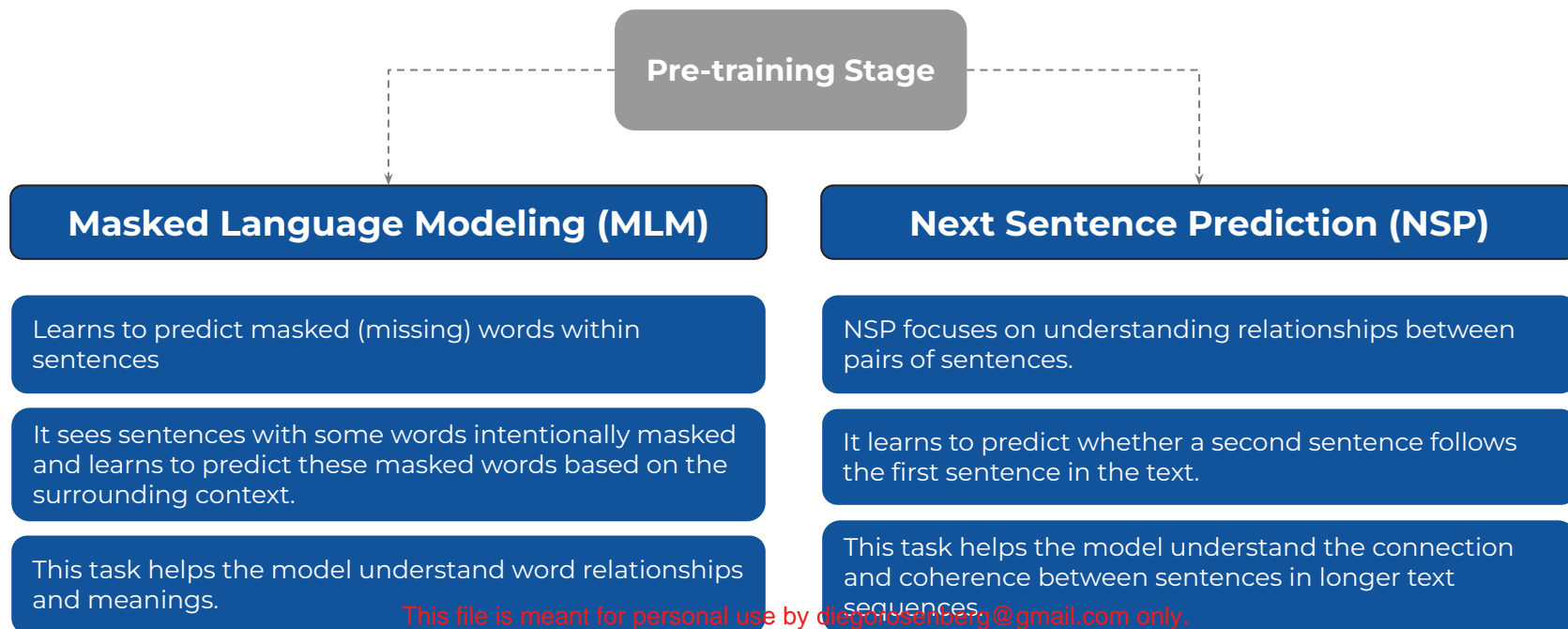
This file is meant for personal use by [diegorosenberg@gmail.com](mailto:diegorosenberg@gmail.com) only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

# BERT - Pre-training Stage

BERT's **pre-training** stage is done in **two parts**



This file is meant for personal use by [diego.senaberg@gmail.com](mailto:diego.senaberg@gmail.com) only.

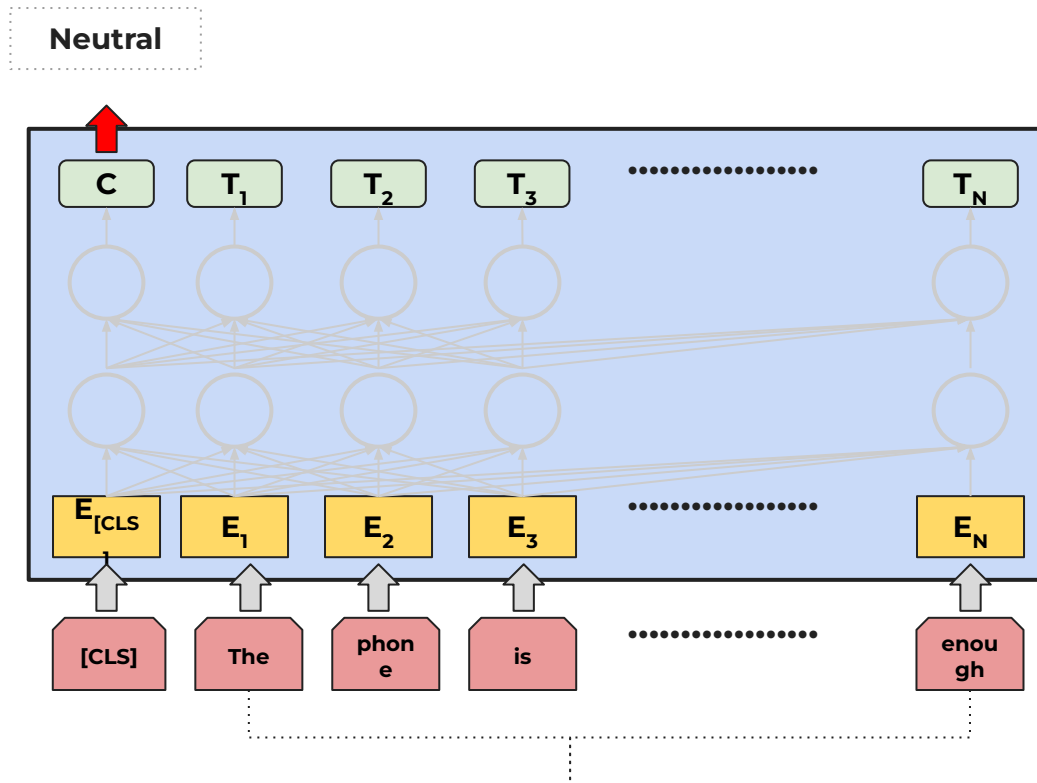
Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.



# What are the downstream applications of BERT?

## Sentiment Classification



Sentence

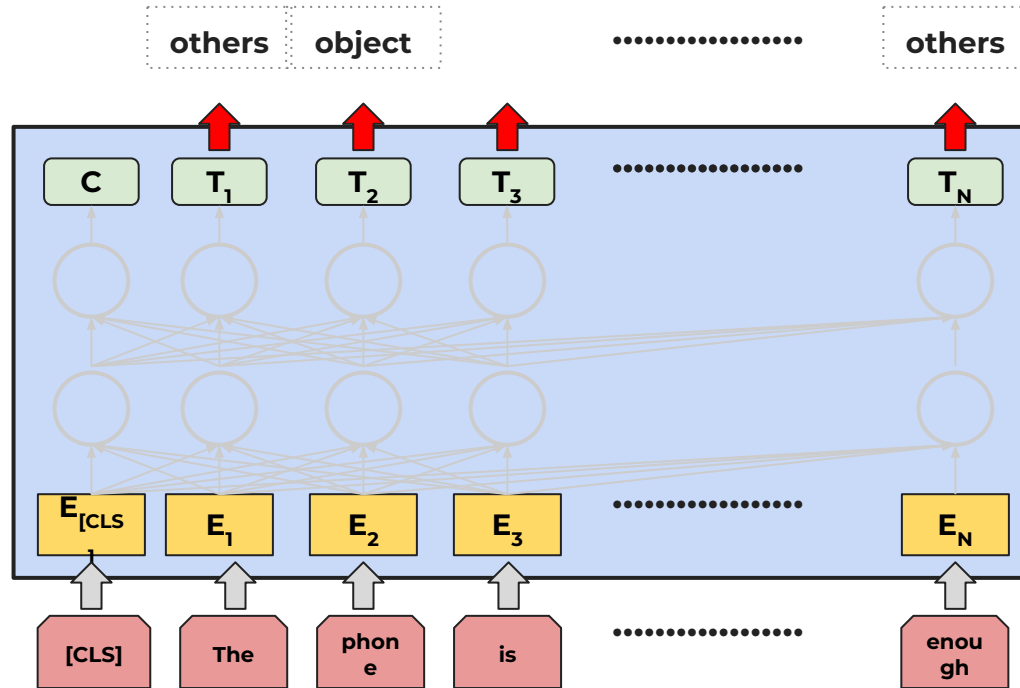
This file is meant for personal use by [diegorosenberg@gmail.com](mailto:diegorosenberg@gmail.com) only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

# What are the downstream applications of BERT?

## Named Entity Recognition



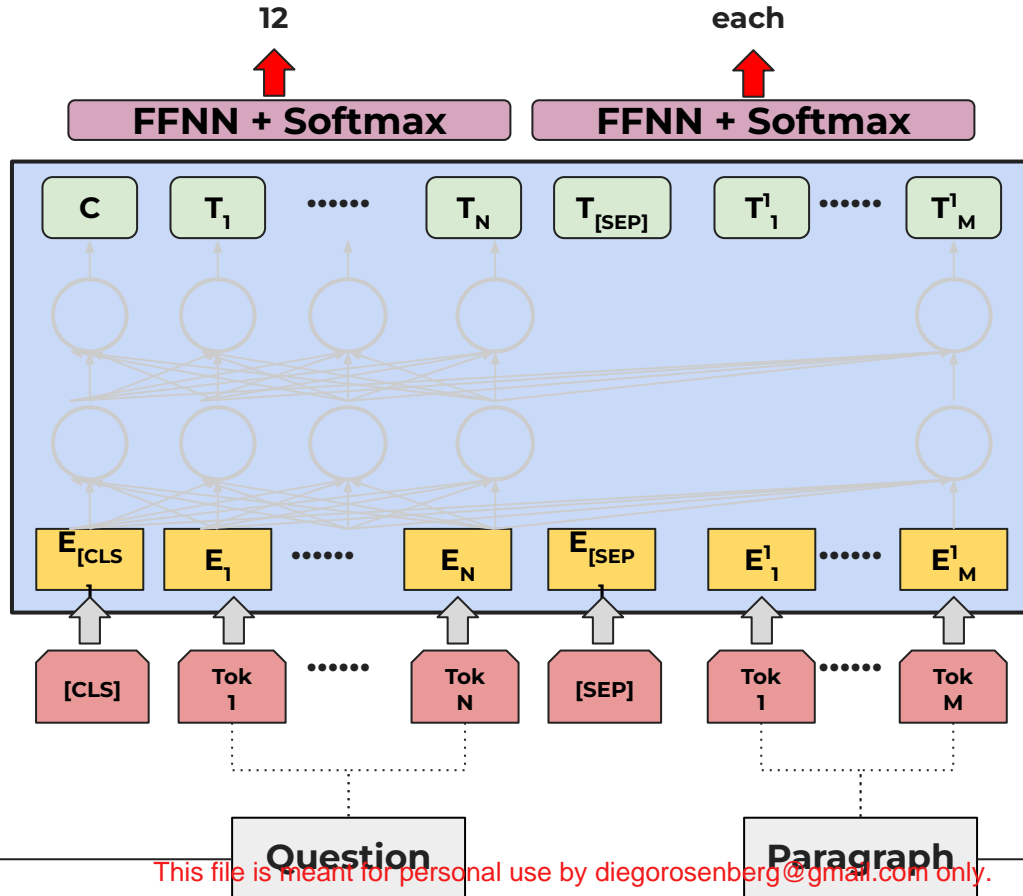
This file is meant for personal use by [ragan.sen@gmail.com](mailto:ragan.sen@gmail.com) only.

Sharing or publishing the contents in part or full is liable for legal action.

Proprietary content. © Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

# What are the downstream applications of BERT?

## Question Answering



How many attention heads does each transformer encoder in BERT have?

BERT consists of 12 transformer encoders with **12 attention heads in each**, totaling approximately 110 million parameters. This design allows BERT to deeply understand language nuances and relationships, making it adept at various language tasks.



# Happy Learning !

