



Great! You're referring to the **BERT (Bidirectional Encoder Representations from Transformers)** model, and the diagram you've shared is a classic representation of how BERT works, especially during **pretraining** with the **Masked Language Modeling (MLM)** objective.

Let me walk you through this diagram and the BERT process in detail:

🔧 1. Input Embedding

At the bottom of the diagram:

- You have a sequence of input tokens: W_1, W_2, W_3, W_4, W_5 .
- One of them (W_4) is replaced with [MASK]. This is the **Masked Language Model (MLM)** part.
- These tokens go through:
 - **Token embeddings**
 - **Segment embeddings** (if applicable, for tasks with multiple sentences)
 - **Position embeddings** (to retain order information)
- Resulting in embeddings: $W_1, W_2, W_3, [MASK], W_5$

🧠 2. Transformer Encoder

Middle section labeled **Transformer encoder**:

- This is the core of BERT.
- It consists of multiple layers of **multi-head self-attention** and **feed-forward networks**.
- All input embeddings are passed **simultaneously** (in parallel) through these layers.
- Because it's **bidirectional**, the model looks at both the **left and right** context of each token (unlike traditional left-to-right models).
- The output is a sequence of **contextualized representations**: O_1, O_2, O_3, O_4, O_5 .

For [MASK], the transformer produces O_4 , a representation that incorporates context from both sides — W_1, W_2, W_3 (left) and W_5 (right).

3. Classification Layer

The next layer is a **fully-connected (dense) layer**, followed by:

- **GELU activation** (Gaussian Error Linear Unit)
- **Layer normalization**

This processes the contextual vector O_4 (for the masked word) into a prediction.

4. Output Prediction (Softmax over Vocabulary)

At the top:

- The final representation is passed through a **softmax layer** over the entire vocabulary.
 - It tries to predict the **original masked word** (W_4).
 - The model is trained to minimize the **cross-entropy loss** between the predicted token and the true token (W_4).
 - In your diagram, the model outputs W'_4 , its guess for the masked token.
-

5. Training Objective

During **pretraining**, BERT is trained using two objectives:

1. **Masked Language Modeling (MLM)** – as shown in your diagram.
 2. **Next Sentence Prediction (NSP)** – where BERT is given two sentences and must predict whether the second follows the first.
-

Summary of Steps in Diagram

Step	Component	Description
1	Input Tokens	Sentence with [MASK] token
2	Embedding Layer	Adds position, segment, and word embeddings
3	Transformer Encoder	Learns contextual relationships
4	Classification Layer	Maps output to vocabulary space
5	Softmax Prediction	Predicts the original word for [MASK]

After Pretraining: Fine-Tuning

Once BERT is pretrained, it can be fine-tuned for tasks like:

- Text classification
- Question answering
- Named Entity Recognition (NER)
- Sentiment analysis

You just add a task-specific layer on top and train it on labeled data.

Important BERT Features

- **Bidirectional Attention:** Looks both forward and backward in context.
 - **Transfer Learning:** Pretrain on large corpora → fine-tune on small task-specific datasets.
 - **Masked Language Model:** Encourages deeper context understanding.
-

If you want, I can explain:

- A real example with actual tokens.
- The math behind the self-attention mechanism.
- How fine-tuning differs from pretraining in practice.