

Steps in Unsupervised Pretraining

1. Dataset:

- The training dataset used for GPT-1 was **BooksCorpus**, a collection of over 7,000 books, totaling around 800 million words. The reason for choosing this corpus was that books often have rich and long-range dependencies between sentences and paragraphs, making it ideal for the model to learn broader contextual relationships.

2. Model Architecture:

- GPT-1 is based on the **Transformer decoder architecture**, introduced by Vaswani et al. in 2017. Unlike the original Transformer, which has both an encoder and a decoder, GPT-1 only uses the decoder side.
- The Transformer architecture uses **self-attention mechanisms** to allow the model to weigh the importance of different words in the input sequence when predicting the next word.

3. Language Modeling Objective:

- The model is trained using the standard **maximum likelihood estimation (MLE)** objective for language modeling. This means it tries to maximize the probability of the correct next word, given all the previous words.

- Mathematically, for a sequence of words $w_1, w_2, w_3, \dots, w_{T-1}, w_T, w_{T+1}, \dots$, the model tries to maximize:

$$P(w_1, w_2, \dots, w_T) = \prod_{t=1}^T P(w_t | w_1, w_2, \dots, w_{t-1})$$

$$P(w_1, w_2, \dots, w_T) = \prod_{t=1}^T P(w_t | w_1, w_2, \dots, w_{t-1})$$

$$P(w_1, w_2, \dots, w_T) = \prod_{t=1}^T P(w_t | w_1, w_2, \dots, w_{t-1})$$

- Here, w_t is the word at position t , and the model learns to predict w_t based on the preceding words w_1, w_2, \dots, w_{t-1} .

4. Training Process:

- During training, the input is a sequence of tokens (words or subwords) from the text corpus. The model generates a

prediction for each token based on the previous tokens in the sequence.

- **Self-attention mechanism:** The model uses self-attention to focus on different parts of the input sequence. When predicting the next word, it weighs the importance of all previous words in the sequence.

- **Masked Self-Attention:** Since GPT-1 is autoregressive, it only attends to the previous tokens in the sequence and not to future ones. This ensures that it predicts the next word solely based on the words before it, not after.

5. Positional Encoding:

- Unlike traditional sequence models (like RNNs), Transformers don't inherently understand the order of tokens in a sequence. To handle this, **positional encodings** are added to the input embeddings to give the model a sense of word order.

- These positional encodings are fixed vectors that are added to the word embeddings, allowing the model to differentiate between, say, the first word in a sentence and the last word, even though both might be the same word.

6. Unsupervised Learning:

- The learning here is unsupervised because the model is not given any explicit information about the tasks it will be used for later (such as translation, summarization, or question answering). It's simply trying to learn to predict the next word in the sequence as accurately as possible.
- By doing so, the model learns a vast array of language patterns: syntactic rules (grammar), semantic relationships (word meanings), and even world knowledge (common phrases, facts, etc.), all without any explicit task-specific supervision.

Benefits of Unsupervised Pretraining

1. **General Language Understanding:**
 - By training on a large corpus of text data, GPT-1 learns to develop a general understanding of language that is not tied to any specific task. This general language knowledge can later be adapted to various tasks through fine-tuning.
2. **Contextual Awareness:**
 - The model can handle long-range dependencies between words, allowing it to capture context that spans entire paragraphs or chapters (especially when trained on data like the BooksCorpus).
3. **Transfer Learning:**
 - The key advantage of this stage is the ability to transfer the knowledge gained from unsupervised pretraining to downstream tasks (via supervised fine-tuning). This means the model doesn't have to start from scratch for every new task but can instead leverage the patterns and knowledge learned during pretraining.

Challenges and Considerations

- **Data Efficiency:** GPT-1 showed that large-scale unsupervised pretraining could be incredibly effective in improving task performance, even when fine-tuned with limited labeled data.
- **Compute Requirements:** Pretraining such a large model requires significant computational resources and time. This was an early demonstration of how scaling up model size and data could lead to performance gains.