

## BERT কী?

**BERT (Bidirectional Encoder Representations from Transformers)** একটি গভীর শিক্ষণ (Deep Learning) মডেল যা প্রাকৃতিক ভাষা প্রক্রিয়াকরণ (Natural Language Processing, NLP) এর ক্ষেত্রে বিপ্লব ঘটিয়েছে। এটি মূলত ভাষার মধ্যে শব্দের প্রসঙ্গ (context) বুঝতে সক্ষম, যা এটিকে অনেক NLP কাজের জন্য অত্যন্ত কার্যকরী করে তোলে, যেমন: টেক্সট শ্রেণীবিভাগ, প্রশ্ন-উত্তর, নামকরণী সত্ত্বা চিনহিতকরণ (Named Entity Recognition), ইত্যাদি।

BERT মূলত **Transformer** আর্কিটেকচারের ওপর ভিত্তি করে কাজ করে, যা ২০১৭ সালে ভাসওয়ানি (Vaswani) et al. দ্বারা প্রবর্তিত হয়েছিল। পূর্ববর্তী মডেলগুলির তুলনায় (যেমন RNN, LSTM), BERT শব্দগুলোকে একসাথে এবং সাযুজ্যপূর্ণভাবে বিশ্লেষণ করতে সক্ষম, যেহেতু এটি **Bidirectional** অর্থাৎ দুই দিক থেকে প্রসঙ্গ বুঝতে পারে।

## BERT কিভাবে কাজ করে?

BERT-এর কিছু গুরুত্বপূর্ণ বৈশিষ্ট্য:

### 1. Bidirectional Attention:

- পুরানো মডেলগুলি (যেমন RNN, LSTM) সাধারণত এক দিকে (বাম থেকে ডান অথবা ডান থেকে বাম) টেক্সট প্রক্রিয়া করে। কিন্তু BERT **Bidirectional** বা দুই দিক থেকে প্রসঙ্গ দেখার জন্য ডিজাইন করা।
- উদাহরণস্বরূপ, "bank" শব্দটি যদি "river bank" এবং "financial bank" এ ব্যবহৃত হয়, তবে BERT দুই দিক থেকে প্রসঙ্গ দেখে শব্দটির সঠিক অর্থ বুঝতে পারে।

### 2. Transformer Architecture:

- BERT মূলত **Transformer Encoder** আর্কিটেকচার ব্যবহার করে, যা **Self-Attention** মেকানিজমের মাধ্যমে কাজ করে। এটি শব্দের সম্পর্ক এবং প্রসঙ্গের উপর গুরুত্ব দেয়, ফলে BERT দীর্ঘ-পারস্পরিক সম্পর্ক (**long-range dependencies**) বুঝতে পারে।

### 3. Pre-training Tasks: BERT দুটি প্রধান প্রশিক্ষণ কাজ (tasks) নিয়ে প্রাক-প্রশিক্ষণ (pre-training) করা হয়:

- Masked Language Model (MLM):** একটি বাক্যে কিছু শব্দ এলোমেলোভাবে মুছে ফেলা হয় এবং BERT সেগুলো পুনরুদ্ধারের চেষ্টা করে। উদাহরণস্বরূপ:
  - ইনপুট: "The cat sat on the [MASK]."
  - BERT "mat" শব্দটি আগের ও পরের প্রসঙ্গ দেখে অনুমান করতে শিখে।
- Next Sentence Prediction (NSP):** দুটি বাক্য দেওয়া হলে, BERT এই দুটি বাক্যের মধ্যে সম্পর্ক বুঝতে শেখে। এটি বিশেষ করে প্রশ্ন-উত্তর বা টেক্সট এনটেইলমেন্ট কাজের জন্য গুরুত্বপূর্ণ।

### 4. Fine-tuning:

- BERT যখন প্রাক-প্রশিক্ষিত হয়, তখন এটি নির্দিষ্ট কাজে **fine-tune** করা হয়, যেমন: sentiment analysis, text classification, ইত্যাদি। Fine-tuning-এ খুব কম

পরিমাণে লেবেলযুক্ত ডেটা প্রয়োজন, যা মডেলকে আরও স্পেসিফিক কাজের জন্য প্রস্তুত করে।

## BERT কেন এত শক্তিশালী?

- **Contextual Word Representation:** BERT প্রতিটি শব্দের জন্য একটি আলাদা কনটেক্সটভিত্তিক উপস্থাপনা তৈরি করে। উদাহরণস্বরূপ, "bat" শব্দটি "baseball bat"-এ এবং "flying bat"-এ আলাদা হবে।
- **Transfer Learning:** BERT-এর মত প্রাক-প্রশিক্ষিত মডেলগুলি বড় পরিমাণের ডেটার ওপর প্রশিক্ষিত হয়ে থাকে, তাই এগুলি খুব দ্রুত এবং কার্যকরীভাবে নির্দিষ্ট কাজে ব্যবহার করা যায়।

## Bangla ভাষায় BERT ব্যবহার

বাংলা ভাষার জন্যও BERT ব্যবহার করা সম্ভব। তবে কিছু চ্যালেঞ্জ আছে, যেমন:

- **Pre-training:** বাংলা ভাষায় BERT প্রাক-প্রশিক্ষণ করতে হলে, বিশাল পরিমাণ বাংলা টেক্সট ডেটার প্রয়োজন, যা পুরোপুরি English-এর মতো উপলব্ধ না হলেও কিছু বাংলা ডেটাসেট (যেমন, বাংলা উইকিপিডিয়া) ব্যবহার করে এটি করা যেতে পারে।
- **Fine-tuning:** বাংলা BERT মডেলকে নির্দিষ্ট কাজের জন্য fine-tune করতে হয়, যেমন: sentiment analysis, text classification, প্রশ্নোত্তর ইত্যাদি।

## বাংলা BERT ব্যবহার করার পদ্ধতি

1. **Pre-trained Bangla BERT Models:** বাংলার জন্য কিছু প্রাক-প্রশিক্ষিত BERT মডেল রয়েছে। যেমন: Bangla BERT যেটি Hugging Face থেকে পাওয়া যায়। এই মডেলটি বাংলা ডেটার ওপর প্রশিক্ষিত এবং বাংলা টেক্সট বুঝতে সক্ষম।
2. **Fine-tuning:** বাংলা BERT মডেলটি টাস্ক অনুযায়ী fine-tune করা যেতে পারে। যেমন: একটি বাংলা সেন্টিমেন্ট এনালাইসিস ডেটাসেট নিয়ে BERT মডেলকে প্রশিক্ষিত করা।
3. **টেক্সট শ্রেণীবিভাগে BERT ব্যবহার (বাংলা):** ধরুন, আমাদের একটি বাংলা বাক্য আছে এবং আমরা এর সেন্টিমেন্ট (পজিটিভ অথবা নেগেটিভ) নির্ধারণ করতে চাই। BERT এর মাধ্যমে এটি করা যেতে পারে।

## বাংলা BERT দিয়ে সেন্টিমেন্ট এনালাইসিসের কোড

```
from transformers import BertTokenizer, BertForSequenceClassification
from torch.utils.data import DataLoader
import torch

# Load pre-trained Bangla BERT model and tokenizer
tokenizer = BertTokenizer.from_pretrained('ai4bharat/indic-bert')
model = BertForSequenceClassification.from_pretrained('ai4bharat/indic-bert', num_labels=2)

# Example Bangla text (e.g., sentiment analysis task)
sentence = "এটি একটি খুব ভালো সিনেমা।"

# Tokenize the sentence
inputs = tokenizer(sentence, return_tensors='pt', truncation=True, padding=True)

# Predict sentiment (0 = Negative, 1 = Positive)
with torch.no_grad():
    outputs = model(**inputs)
    logits = outputs.logits
    predicted_class = torch.argmax(logits, dim=-1).item()

# Display the result
if predicted_class == 1:
    print("Sentiment: Positive")
else:
    print("Sentiment: Negative")
```

এই কোডে, আমরা একটি বাংলা বাক্যকে টোকেনাইজ করেছি এবং প্রাক-প্রশিক্ষিত বাংলা BERT মডেল ব্যবহার করে সেন্টিমেন্ট বিশ্লেষণ করেছি। মডেলটি "এটি একটি খুব ভালো সিনেমা" বাক্যটির সেন্টিমেন্ট পজিটিভ হিসেবে চিহ্নিত করবে।

## চ্যালেঞ্জসমূহ

- **Tokenization:** বাংলা ভাষায় টোকেনাইজেশন কখনো কখনো ইংরেজির চেয়ে একটু জটিল হতে পারে, বিশেষত যেখানে জটিল বা সংযুক্ত শব্দ রয়েছে।
- **ডেটাসেটের অভাব:** বাংলা ভাষার জন্য পর্যাপ্ত পরিমাণে লেবেলযুক্ত ডেটাসেট (যেমন সেন্টিমেন্ট এনালাইসিস) ইংরেজির তুলনায় কম পাওয়া যায়।

## উপসংহার

BERT এর দুটি প্রধান বৈশিষ্ট্য হল **Bidirectional Attention** এবং **Transformer Architecture**, যা এটি NLP কাজে বিশেষভাবে কার্যকরী করে তোলে। বাংলা ভাষায় BERT ব্যবহারের জন্য প্রাক-প্রশিক্ষিত মডেলগুলি ব্যবহার করে এবং সেই অনুযায়ী ফাইন-টিউনিং করে বিভিন্ন NLP কাজ করা যেতে পারে।