# BART Model Overview

1910.13461v1.pdf
PDF

The document is a research paper titled **"BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension"** by **Facebook AI researchers** Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer.

## Key Takeaways:

1. **BART Model Overview:**

   - BART (Bidirectional and Auto-Regressive Transformers) is a **denoising autoencoder** designed for **sequence-to-sequence (seq2seq) pretraining**.

   - It is a **generalization of BERT and GPT**, combining a **bidirectional encoder** (like BERT) and an **autoregressive decoder** (like GPT).

   - The model **corrupts text** with noising functions and learns to **reconstruct the original text**.

2. **Noising Strategies for Pretraining:**

   - **Token Masking**: Some tokens are replaced with [MASK] (similar to BERT).

   - **Token Deletion**: Some tokens are removed entirely.

   - **Text Infilling**: Random spans of text are replaced with a single [MASK] token.

   - **Sentence Permutation**: The order of sentences is randomly shuffled.

   - **Document Rotation**: The document is rotated to start from a random token.

3. **Fine-tuning and Applications:**

   - **Text Generation:** BART excels in **summarization**, **dialogue generation**, and **abstractive question answering**.

   - **Text Comprehension:** It achieves **state-of-the-art performance on SQuAD (QA) and GLUE (NLU benchmarks)**.

   - **Machine Translation:** It improves **back-translation** and translation tasks by using BART as a **pretrained target-side language model**.

4. **Performance Comparisons:**

   - Matches **RoBERTa** on **GLUE and SQuAD benchmarks**.

   - **Outperforms previous models** in **text generation** tasks (e.g., achieving a **6 ROUGE point improvement** in summarization).

   - Improves **machine translation BLEU scores** over **back-translation baselines**.

5. **Architectural Details:**

   - Uses **Transformer-based encoder-decoder architecture**.

   - Employs **GeLU activation functions**.

   - Has **6 encoder-decoder layers (base model)** and **12 layers each (large model)**.

6. **Ablation Studies:**

   - **Text infilling and sentence shuffling** yield the best performance.

   - **Token deletion is better for generation tasks than token masking**.

   - **Left-to-right pretraining improves generation quality**.

## Conclusion:

- BART is a **flexible and powerful pretraining method** for both **discriminative** (classification, QA) and **generative** (summarization, translation) NLP tasks.

- Future work could explore **task-specific corruption strategies** to further enhance pretraining.

Would you like a **simplified summary**, or do you need **specific insights on an aspect** of the paper? 😊

Architecture of BART

## Architecture of BART (Bidirectional and Auto-Regressive Transformers)

BART is built using a **Transformer-based sequence-to-sequence (seq2seq) model** that consists of two key components:

1. **Bidirectional Encoder** (like BERT)

2. **Autoregressive Decoder** (like GPT)

**1. Encoder:**

- Similar to **BERT**, it is a **bidirectional Transformer encoder** that processes the **entire input sequence simultaneously**.
- It takes a **corrupted (noised) text input** and learns to **understand the context**.
- Unlike GPT, which processes words **left-to-right**, BART's encoder **looks at all words at once**.

**2. Decoder:**

- Similar to **GPT**, it is an **autoregressive (left-to-right) Transformer decoder**.
- It generates the **reconstructed (denoised) output text** step by step.
- The decoder **attends to the encoder's hidden states** using **cross-attention**.

## Key Architectural Features:

- **Transformer-based:** Uses **self-attention and feed-forward networks** like other Transformer models (e.g., BERT, GPT).
- **Pretraining with Noising Functions:** BART **corrupts input text** (e.g., deleting words, shuffling sentences) and learns to **reconstruct** it.
- **Flexible Input/Output Alignment:** The input to the encoder does **not have to be aligned** with the output of the decoder.
- **Cross-Attention:** The decoder has **cross-attention layers** that allow it to reference encoded representations when generating output.
- **Activation Function:** Uses **GeLU (Gaussian Error Linear Unit)** instead of ReLU for better performance.
- **Layer Configurations:**
  - **Base Model:** 6 encoder layers + 6 decoder layers.
  - **Large Model:** 12 encoder layers + 12 decoder layers.

## Comparison with BERT & GPT

| Feature | BERT (Encoder) | GPT (Decoder) | BART (Seq2Seq) |
|---|---|---|---|
| **Encoder** | ✅ Yes (Bidirectional) | ❌ No | ✅ Yes (Bidirectional) |
| **Decoder** | ❌ No | ✅ Yes (Left-to-right) | ✅ Yes (Left-to-right) |
| **Cross-Attention** | ❌ No | ❌ No | ✅ Yes |

| Feature | BERT (Encoder) | GPT (Decoder) | BART (Seq2Seq) |
|---|---|---|---|
| Pretraining | Masked Language Model | Left-to-right Language Model | Denoising Autoencoder |
| Output Task | Classification, QA | Text Generation | Both Classification & Generation |

## Visual Representation:

BART combines both bidirectional and autoregressive modeling, making it **powerful for both understanding and generating text**.

📌 **Would you like a diagram or a more in-depth explanation of any component?** 🚀

yes

Here's a visual breakdown of **BART's architecture** and its components:

---

## 🖼️ BART Architecture Diagram

```pgsql
        ┌─────────────────────┐
        │   Noised Input Text     │  <-- Corrupted with masking, deletion,
shuffling, etc.
        └─────────────────────┘
                   ↓
        ┌─────────────────────┐
        │  **Bidirectional Encoder** │   (Similar to BERT)
        │   (Processes entire input)  │
        └─────────────────────┘
                   ↓
        ┌─────────────────────┐
        │    **Cross-Attention**     │   (Encoder-Decoder Attention)
        └─────────────────────┘
                   ↓
        ┌─────────────────────┐
        │  **Autoregressive Decoder** │   (Similar to GPT)
        │    (Generates text step-by-step) │
        └─────────────────────┘
```

```
                        ↓
        ┌──────────────────────────────┐
        │   **Reconstructed Output**    │   <-- Model predicts original text
        └──────────────────────────────┘
```

---

## ✳️ Detailed Explanation of Components

### 1️⃣ Input Corruption (Pretraining Step)

- The input text is **corrupted** using various **noising functions**:

  - **Token Masking** (Replace words with `[MASK]` like BERT)

  - **Token Deletion** (Remove words completely)

  - **Text Infilling** (Replace spans of text with a single mask)

  - **Sentence Shuffling** (Rearrange sentence order)

  - **Document Rotation** (Change where the text starts)

### 2️⃣ Encoder (Bidirectional like BERT)

- The **encoder processes the entire input sequence simultaneously**.

- It learns contextual representations of words from **both left and right directions**.

- The **noised input text** is passed through the encoder.

### 3️⃣ Cross-Attention (Bridges Encoder & Decoder)

- The decoder receives **contextual information** from the encoder using **cross-attention layers**.

- This helps the decoder generate the correct output based on the **original (uncorrupted) input**.

### 4️⃣ Decoder (Autoregressive like GPT)

- The decoder is **left-to-right autoregressive**, meaning it **predicts words one by one**.

- It generates the **original (denoised) text step-by-step**.

- Unlike BERT (which only predicts missing words independently), BART's decoder **generates full sequences**.

## 5️⃣ Final Output

- The decoder reconstructs the **original, uncorrupted text**.

- This allows BART to be useful for tasks like **summarization, text generation, and translation**.

---

## 🚀 Why is BART Powerful?

- **Combines BERT and GPT Strengths**:

  - Encoder **understands text deeply** like BERT.

  - Decoder **generates text sequentially** like GPT.

- **Flexible Pretraining**:

  - Works well for **text comprehension (classification, QA)** and **generation (summarization, translation)**.

- **Handles Large Text Inputs**:

  - Can model long-range dependencies **better than BERT**.

---

## 🔍 Want me to generate a more detailed image for you? 😊

Is this conversation helpful so far?