# Navigating the Text Summarization Landscape:
# LLM Models vs. Leading Transformers

In an era saturated with information, professionals often find themselves immersed in vast amounts of text, ranging from reports to scientific papers. While these written pieces delve into specific themes, the need for efficient summarization tools becomes evident. Artificial intelligence and machine learning, particularly in Natural Language Processing (NLP) tasks, play a pivotal role in addressing this need. This blog post aims to shed light on the performance of LLM models and leading transformers, unveiling insights into their summarization capabilities.

## Introduction: Evolution of Automatic Summarization

The world of automatic summarization is ever evolving, relying on traditional metrics like ROUGE and cutting-edge language models as benchmarks. In this exploration, we focus on the latest ROUGE-L metrics for LLM models on the SAMSUM dataset, comparing them with the performance of leading transformers like FLAN-T5, Bart-large, Pegasus, Bart-finetuned, and Instruct DS with the same dataset SAMSUM.

## Types of Text Summarization

Text summarization training typically employs supervised learning, aligning text passages with golden annotated summaries. Summaries fall into three categories:

### 1. Abstractive Summarization

**Definition:** Crafting a concise summary not bound by the original text's exact words, interpreting content to express main ideas in a novel way.

**Example:**

Original Text:

The new restaurant in town serves a variety of cuisines, from Italian pasta to Japanese sushi. The ambiance is vibrant, and the service is exceptional.

Abstractive Summary:

Experience a diverse culinary journey at the latest restaurant in town, offering everything from delicious Italian pasta to authentic Japanese sushi, all within a vibrant ambiance and with exceptional service.

### 2. Extractive Summarization

**Definition:** Selecting and presenting existing sentences or phrases directly from the source text, retaining original wording for accuracy.

**Example:**

Original Text:

The benefits of regular exercise are numerous, including improved cardiovascular health, increased energy levels, and better mental well-being.

Extractive Summary:

Regular exercise brings numerous benefits, such as improved cardiovascular health, increased energy levels, and enhanced mental well-being.

## 3. Hybrid Summarization

**Definition:** Integrating elements of both abstractive and extractive methods to produce coherent and contextually relevant summaries.

**Example:**

Original Text:

The latest smartphone boasts an impressive camera with advanced image stabilization technology. Additionally, it has a sleek design and a long-lasting battery.

Hybrid Summary:

Highlighting the latest smartphone's key features, including its impressive camera with advanced image stabilization, sleek design, and long-lasting battery, makes it a top choice for tech enthusiasts.

## Comparative Performance Overview

Our journey begins with a comprehensive comparison of LLM models and leading transformers, each rigorously assessed using the latest ROUGE-L metrics on the SAMSUM dataset. The tables below provide a detailed breakdown of their performance, allowing us to draw meaningful insights.

**Leading Transformers on SAMSUM dataset**

| Models | ROUGE-1 | ROUGE-2 | ROUGE-L | Overview |
|---|---|---|---|---|
| FLAN-T5 | 47.24 | 23.51 | 39.63 | FLAN-T5 excels in ROUGE metrics. |
| bart-large | 53.43 | 28.74 | 44.18 | BART-Large stands out in summarization tasks. |
| pegasus | 44 | 21 | 35 | Pegasus achieves high ROUGE scores. |
| bart-finetuned | 52.82 | 28.13 | 43.72 | BART-Finetuned demonstrates robust summarization. |
| InstructDS | 55.3 | 31.3 | 46.7 | InstructDS showcases versatility and robust performance. |

## LLAMA Models on SAMSUM Dataset

| Models | ROUGE-L (Precision) | ROUGE-L (Recall) | ROUGE-L (F1 Score) | Overview |
|---|---|---|---|---|
| LLAMA-7B Finetuned Model | 52 | 46 | 48.81 | LLAMA-7B Finetuned Model achieves a harmonious balance of precision and recall. |
| LLAMA-7B Base Model | 15 | 28 | 19.53 | LLAMA-7B Base Model indicates areas for improvement in summarization quality. |
| LLAMA-70B Base Model | 25 | 53 | 33.9 | LLAMA-70B Base Model achieves a balanced performance in generating accurate content. |
| **Mistral 7B** | 16 | 57 | 37 | Mistral 7B demonstrates a distinctive recall while maintaining a reasonable balance between precision and recall. |

## Context Length Considerations:

Context length plays a crucial role in the effectiveness of summarization models. For Mistral models & LLAMA models, the context length is set at 4,000 tokens. It's essential to note that the context length for leading transformers varies:

- **FLAN-T5:** Common tokenization limits range from 512 to 2048 tokens.
- **Bart-large and Bart-finetuned:** Token limits typically range from 1024 to 4096 tokens, depending on the version and configuration.
- **Pegasus:** Token limits are often in the range of 1024 to 2048 tokens.

- **Instruct DS (trained on FLAN-T5):** Context length aligns with T5 models, ranging from 512 to 2048 tokens.

## Strategic Recommendations

As we navigate through the results, it becomes evident that the Mistral-7B finetuned & LLAMA-7B Finetuned Model stands out, particularly for precision and recall on the SAMSUM dataset. This revelation becomes the cornerstone of our strategic recommendation.

## Conclusion

As we carefully examine the performance of the latest Mistral, LLAMA models and leading transformers, Instruct DS emerges as the standout performer on the SAMSUM dataset, closely trailed by BART-Large and BART-Finetuned. Mistral & LLAMA models showcase varying degrees of performance, with the Finetuned Model leading in precision and recall.

The recent metrics underscore the importance of continuous evaluation and fine-tuning. For those aiming for a balance between precision and recall, the Mistral-7B & LLAMA-7B Finetuned Model stands out as an appealing choice. However, for applications requiring cutting-edge performance with a specific focus on abstractive summarization, Instruct DS or BART-Finetuned emerge as robust contenders. It's worth noting that Instruct DS involves instructive dialogue summarization with query aggregations (LLM), where the task is to condense dialogue information into concise text, following specific instructions. The output is tailored to the input, disregarding user preferences, and the challenge lies in the scarcity of relevant training data. To address this, the MIstral-7B Finetuned & LLAMA-7B Finetuned models becomes instrumental.

## Fine-Tune with AIQ for Summarization Excellence

In our ever-evolving landscape, the call to action is clear: fine-tune and evaluate LLM models, with a special emphasis on the LLAMA-7B Finetuned Model & Mistral-7B Finetuned Model, for your specific data. For those seeking optimal results in summarization tasks, consider leveraging the power of AIQ fine-tuning. Detailed guidance is available in our blog post, ["How to Fine-Tune LLM's for Summarization in AIQ"](), providing a step-by-step walkthrough using AIQ, enriched with screenshots.

These insights, combined with a nuanced understanding of context length considerations, empower practitioners and researchers to navigate the complexities of automatic summarization effectively, ensuring optimal performance for diverse summarization tasks.