# Leveraging Benchmark Datasets for Holistic Model Evaluation

## Introduction

Large Language Models (LLMs) are intricate systems, and simple evaluation metrics like ROUGE and BLEU scores offer limited insights into their capabilities. To comprehensively measure and compare LLMs, researchers utilize pre-existing datasets and associated benchmarks. Selecting appropriate evaluation datasets is crucial for accurately assessing an LLM's performance across various tasks and scenarios.

## Evaluation benchmarks



## Considerations for Evaluation Datasets

- Isolation of Model Skills: Choose datasets that isolate specific model skills such as reasoning, common sense knowledge, or identifying potential risks like disinformation or copyright infringement.
- Novelty of Evaluation Data: Evaluate the model on data it hasn't encountered during training to gauge its true capabilities.
- Key Benchmarks for LLM Evaluation
- GLUE (General Language Understanding Evaluation)
- Introduced in 2018 to assess LLMs across multiple natural language tasks like sentiment analysis and question-answering.
- Aims to promote the development of models capable of generalizing across various tasks.

## SuperGLUE

Introduced in 2019 as an advancement over GLUE, featuring more challenging tasks including multi-sentence reasoning and reading comprehension.
Includes tasks not covered in GLUE and provides a benchmark for measuring LLM performance.

# Massive Multitask Language Understanding (MMLU)

Designed for modern LLMs, requiring extensive world knowledge and problem-solving abilities.
Tests models on a wide range of tasks including mathematics, US history, computer science, and law.

# BIG-bench

Consists of 204 tasks spanning diverse domains such as linguistics, childhood development, math, and common sense reasoning.
Available in three sizes to accommodate varying computational costs associated with running large benchmarks.

# Holistic Evaluation of Language Models (HELM)

Aims to enhance model transparency and offer guidance on model performance for specific tasks.
Utilizes a multimetric approach, assessing seven metrics across 16 core scenarios, including fairness, bias, and toxicity.

# Importance of HELM

Multimetric Assessment: Goes beyond basic accuracy measures to evaluate fairness, bias, and toxicity, crucial for assessing potentially harmful behavior exhibited by LLMs.
Continuous Evolution: HELM is a living benchmark designed to evolve with the addition of new scenarios, metrics, and models.

# Conclusion

By leveraging benchmark datasets like GLUE, SuperGLUE, MMLU, BIG-bench, and HELM, researchers can conduct comprehensive evaluations of LLMs, gaining insights into their performance across a wide range of tasks and scenarios. These benchmarks play a pivotal role in advancing the understanding of LLM capabilities and guiding the development of more robust and responsible language models.