

MSci Project

- **Degree programme:** MSci Computer Science
- **Project title:** Experiments With Word Order In Large Language Models
- **Student name:** Tyler Liddell
- **Student email address:** tyler.liddell@city.ac.uk
- **Supervisor:** Pranava Madhyastha
- **Project proposed by:** Pranava Madhyastha
- **Proprietary interests:** This is not a client-based project, so these do not apply. There are no such arrangements.
- **Other promises:** None.
- **Word count:** 12419

Abstract

Large Language Models (LLMs) have demonstrated impressive performance on various natural language understanding tasks. However, the extent to which they rely on linguistic structures, particularly word order, remains unclear. This study investigates the sensitivity of LLMs to word order perturbations across multiple benchmarks, including multiple-choice tasks (MMLU, TruthfulQA, Hellaswag, Winogrande, and ARC Challenge) and generation tasks (TruthfulQA Generation and GSM8K). Several open-source LLMs were evaluated on both the original benchmarks and versions with shuffled word order. The results show that while LLMs can go beyond simple word co-occurrences, they often rely heavily on superficial cues like keyword associations. Multiple-choice benchmarks were found to be generally insensitive to word order perturbations, with models maintaining relatively high accuracy despite the shuffling. In some cases, performance even improved when questions were shuffled or removed entirely, raising concerns about the validity of these benchmarks for evaluating true language understanding. The study also highlights the importance of fine-tuning for specific tasks, as demonstrated by the improved ability of the fine-tuned Mistral Instruct model to detect and respond to incoherent prompts compared to base models. These findings contribute to the growing body of research on the capabilities and limitations of LLMs and underscore the need for more robust evaluation methods to accurately assess genuine language understanding and reasoning abilities.

Chapter 1: Introduction

In recent years, a new type of artificial intelligence model called Large Language Models (LLMs) has emerged. LLMs are deep learning models trained on massive amounts of text data to predict the next word in a sequence. By learning the patterns and structures in human language, LLMs can generate human-like text and perform a wide range of natural language tasks.

Transformer-based LLMs have gained massive popularity over the last couple of years, specifically through the introduction of ChatGPT which was estimated to have amassed around 100 million monthly users in January of 2023, just 2 months after its launch (Hu, 2023). As this is being written, new models are being released that perform even better on a range of tasks. From mathematical reasoning to summarizing and story writing, these models are beginning to exhibit human-like testing results.

Transformers, the architecture behind most modern LLMs, were introduced in the paper “Attention Is All You Need” (Vaswani et al., 2023). The key innovation of transformers is the self-attention mechanism, which allows the model to weigh the importance of different parts of the input when making predictions. In each transformer layer, the input is transformed into query, key, and value vectors. The self-attention mechanism computes the dot product between each query and all keys, normalizes the results, and uses them as weights to sum the value vectors. This allows the model to attend to relevant information from any part of the input. The output of the self-attention is then passed through feed-forward layers and residual connections. By stacking multiple transformer layers, the model can learn increasingly complex representations of the input text.

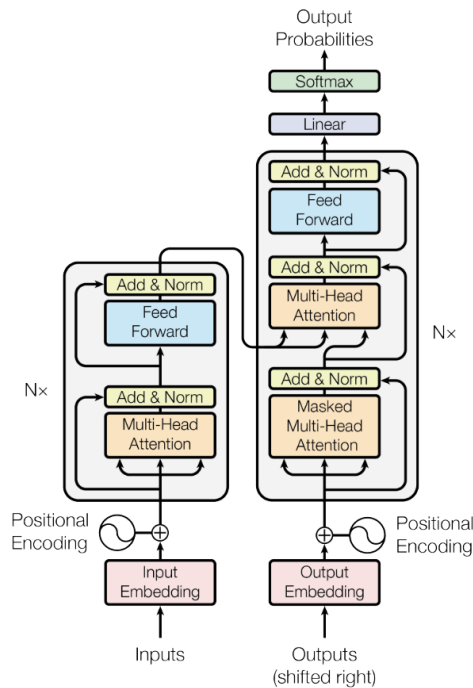


Figure 1: Diagram of transformer architecture as presented by Vaswani et al. (2023)

Because of their popularity, LLMs have been developed by most major tech companies. They vary in size all the way from DistilGPT-2 with 82 million parameters to X.AI’s Grok which has a massive 314 billion parameters. Many companies have focused on developing LLMs and releasing them as open-source, meaning anyone can download the weights for the model and even train or fine-tune it according to their needs, provided they have the necessary resources.

For the purpose of working more closely with the models, this study will focus only on the open-source models. Open LLMs have gained significant traction in the past several years and are becoming quite competitive on several tasks, despite many of them being significantly smaller. Hugging Face’s

Open LLM Leaderboard¹ has become a popular way to track this success, by providing a space to compare model’s scores across various benchmarks.

Understanding the behavior of smaller models can provide valuable insights into the performance of larger models. Many smaller models follow the same design principles as their larger counterparts. As the underlying mechanisms of large models’ outputs, particularly unexpected outputs referred to as hallucinations, are not yet fully understood, studying these phenomena at a smaller scale may help uncover the shortcomings of the massive and highly successful LLMs. Additionally, due to resource constraints, the models evaluated in this study must be able to run on a GPU with 40GB of RAM. The size of a model is usually referred to by its number of parameters. Parameters are essentially weights that the model learns during training, and they control the model’s behavior and output. Generally, models with more parameters have a higher capacity to learn and represent complex patterns in data, leading to better performance on various tasks. However, larger models also require more computational resources and are more challenging to train and deploy. By focusing on smaller, open-source models, this study aims to provide insights that can be more easily interpreted and applied to improve the understanding and performance of LLMs.

Word order plays a crucial role in determining the grammatical roles and functions of elements within a sentence. In languages like English, the order of words typically indicates the subject, verb, and object, with the subject preceding the verb and the object following it. Altering this word order, such as inverting the subject and object, can significantly change the meaning of the sentence. Word order also interacts with other linguistic features, such as suffixes, to convey meaning and show how words group together into phrases and clauses. Thus, word order is intricately connected to both the forms of words and their meanings, making it an essential aspect of linguistic structure [Harris \(1954\)](#).

The goal if this study is to better understand what patterns in human language are important for LLM’s to ”understand” and generate a satisfactory response. My hypothesis closely follows the work of several other studies, which I discuss in chapter 3, that suggested perhaps these models do not rely on core linguistic concepts, primarily word order. This leads me to 3 main research questions which I aim to answer.

- ***Are LLM’s able to go beyond simple word co-occurrences?*** It has been shown before that LLM’s perform well on many tasks without word order, and most likely are depending on word co-occurrences.
- ***To what extent are benchmarks sensitive to perturbations?*** Maybe the problem isn’t entirely with the models, but also with how we evaluate them.
- ***What factors are responsible for the language model’s performance?*** If LLMs are not using word order, what other factors are contributing to their success?

Chapter 2: Output Summary

The primary output of this research project is a comprehensive analysis of the extent to which Large Language Models (LLMs) rely on linguistic structures, particularly word order, to perform various natural language understanding tasks. The analysis is presented in the form of a detailed report, which includes the following components:

1. Evaluation results of several open-source LLMs on a range of benchmarks, including multiple-choice tasks (MMLU, TruthfulQA, Hellaswag, Winogrande, and ARC Challenge) and generation tasks (TruthfulQA Generation and GSM8K). The results cover the models’ performance on both the original benchmarks and versions with shuffled word order, demonstrating their sensitivity to perturbations. (Approximately 10 pages)
2. A discussion of the implications of the findings, addressing the three main research questions: a. Are LLMs able to go beyond simple word co-occurrences? b. To what extent are benchmarks sensitive to perturbations? c. What factors are responsible for the language model’s performance? (Approximately 5 pages)

¹[Hugging Face Open LLM Leaderboard](#)

3. Suggestions for future research directions, such as comparing correct answers across different models, analyzing average logit scores, exploring additional tasks, and investigating the performance of LLMs on translation tasks with permutations. (Approximately 2 pages)

The complete analysis report, along with the code used for the experiments and the generated data, can be found in the Appendices of this dissertation. The code consists of modifications to the `lm-evaluation-harness`, an open-source tool for evaluating language models across various tasks [Gao et al. \(2023\)](#). The modifications, amounting to approximately 500 lines of Python code, enable the tool to perform word order permutations and other adjustments to the datasets, facilitating the experiments conducted in this study.

The intended recipients of this output are researchers and practitioners in the field of natural language processing, particularly those interested in the capabilities and limitations of Large Language Models. The findings presented in the report will contribute to the growing body of knowledge on LLMs and their reliance on linguistic structures, helping to inform future research and development efforts in this area.

The results described in the Results chapter (Chapter 5) directly contribute to this output by providing the empirical evidence and analysis that form the basis of the report. The links between the specific sections of the Results chapter and the corresponding components of the output are as follows:

- Section 5.1 (Multiple Choice Tasks) → Evaluation results on multiple-choice benchmarks
- Section 5.2 (Generation Tasks) → Evaluation results on generation benchmarks
- Conclusions and Discussion chapter → Discussion of implications and suggestions for future research

In summary, the primary output of this project is a comprehensive analysis report that sheds light on the extent to which LLMs rely on word order and other linguistic structures to perform natural language understanding tasks. The report, along with the supporting code and data, will be a valuable resource for researchers and practitioners seeking to advance the field of natural language processing and develop more robust and reliable language technologies.

Chapter 3: Literature Review

Measuring the linguistic and reasoning capabilities of LLM's has been an area of interest since their arrival. The topic has led to many relevant papers coming to different conclusions about the capabilities of LLM's. There are certainly reasons for us to be optimistic about where language models have come in terms of reasoning capabilities. Chain-of-thought prompting, a technique that enhances a model's reasoning abilities, has shown remarkable results. This method involves providing the model with a series of intermediate reasoning steps to solve a given problem, followed by presenting a similar problem for the model to solve using the demonstrated steps. As illustrated in Figure 2, chain-of-thought prompting includes natural language explanations of the step-by-step reasoning process, which guides the model to solve the problem incrementally. By learning from such prompts, the model can then apply the same reasoning to solve similar problems. [Wei et al. \(2023\)](#) have demonstrated that chain-of-thought prompting enables models to solve problems they were previously incapable of solving, expanding the range of tasks that language models can perform without the need for fine-tuning. Notably, the ability to perform chain-of-thought reasoning emerges as the language model scale increases, with larger models showing significant improvements compared to their smaller counterparts.

The effectiveness of chain-of-thought prompting can be attributed to the local structure of the training data, as explored by [Prystawski et al. \(2023\)](#) in their paper "Why think step by step? Reasoning emerges from the locality of experience." The authors hypothesize that the local clusters of strongly related variables or concepts in human experience and language model training data contribute to the success of step-by-step reasoning. Their theoretical analysis proves that for a language model trained on small, related subsets of data from a structured dataset, reasoning through intermediate variables leads to better inference compared to direct estimation. Experiments with language models trained on samples from Bayesian networks further support this hypothesis, showing that intermediate reasoning steps are beneficial when the training data is locally structured with respect to strong dependencies

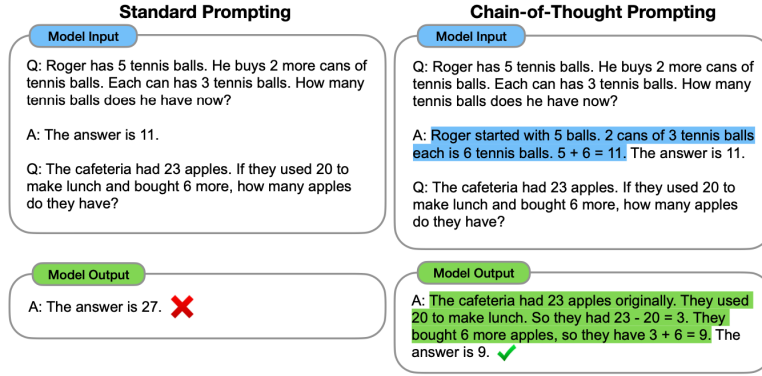


Figure 2: Chain of thought experiment conducted by Wei et al. (2023)

between variables. Moreover, the authors demonstrate that models can generate helpful intermediate variables themselves when trained on locally structured data. This finding aligns with the observation that larger language models exhibit improved chain-of-thought reasoning abilities, as they have likely encountered more locally structured data during training.

Another promising direction for enhancing the reasoning capabilities of LLMs is learning from their own mistakes. In the paper "Can LLMs Learn from Previous Mistakes? Investigating LLMs' Errors to Boost for Reasoning," Tong et al. (2024) introduce COTERRORSET, a benchmark dataset designed to include both correct and error references, demonstrating the types of mistakes made by LLMs and the reasons behind them. The authors propose two methods to explore whether LLMs can benefit from their mistakes:

1. Self-rethinking prompting guides LLMs to rethink whether they have made similar previous mistakes.
2. Mistake tuning involves fine-tuning models on both correct and incorrect reasoning domains.

Experiments across various reasoning tasks and LLMs of different sizes show that both proposed methods consistently yield performance improvements, proving that LLMs can obtain benefits from utilizing mistakes during both tuning and inference stages. Furthermore, the authors provide a thorough analysis of the common types of mistakes made by LLMs in reasoning tasks, such as calculation and logical errors in arithmetic reasoning, and misinterpretation of context and linguistic confusion in common-sense reasoning. This study demonstrates the potential cost-effectiveness of leveraging LLM errors to enhance reasoning capabilities, as collecting incorrect rationales is significantly cheaper than hand-crafting correct references. The findings provide promising directions for future research to improve the reasoning abilities of LLMs by learning from their mistakes.

However, some researchers are more critical due to language models weaknesses in understanding language. Bender and Koller (2020) argue that language models trained only on form (the observable realization of language, such as text) cannot learn meaning (the relation between form and something external to language, such as communicative intent). They contend that meaning cannot be learned from form alone because there is no signal connecting the form to the external world or the speaker's intent. The authors use thought experiments to illustrate that a system exposed only to form during training has no way to ground its responses in reality or understand the meaning behind the text it generates.

Asher et al. (2023) provide a theoretical foundation for the empirical findings that large language models (LLMs) fail to capture important aspects of linguistic meaning. They prove that LLMs are incapable of learning certain fundamental semantic properties, such as semantic entailment and consistency, as defined in formal semantics. The authors show that LLMs cannot learn concepts beyond the first level of the Borel Hierarchy, which limits their ability to capture many aspects of linguistic meaning. This implies that LLMs will continue to operate without formal guarantees on tasks requiring deep linguistic understanding and entailment.

Recent research has investigated the importance of word order in language models' understanding of language. Intuitively, if a language model effectively utilizes linguistic concepts, word order should

be critical to its performance. However, Al-Negheimish et al. (2021) demonstrated that when shuffling the word order in the numerical reasoning questions of the Discrete Reasoning Over the content of Paragraphs (DROP) benchmark Dua et al. (2019), the accuracy of the models did not decrease as drastically as one might expect. The DROP benchmark focuses on tasks that involve reading comprehension along with arithmetic operations such as addition, subtraction, sorting, and counting. The authors employed trigram, bigram, and unigram shuffling to showcase this phenomenon. As illustrated in Figure 3, the models’ performance on the DROP benchmark gradually declined as the degree of shuffling increased. The findings from Al-Negheimish et al. (2021) raise important questions about the extent to which language models truly comprehend language, as opposed to merely exploiting statistical patterns in the data. While the study highlights the importance of word order, it also suggests that its role may not be as critical as previously thought, particularly for tasks similar to those in the DROP benchmark. This highlights the need for further research to better understand the mechanisms behind language models’ performance and to develop more robust evaluation methods that can accurately assess their language understanding capabilities.

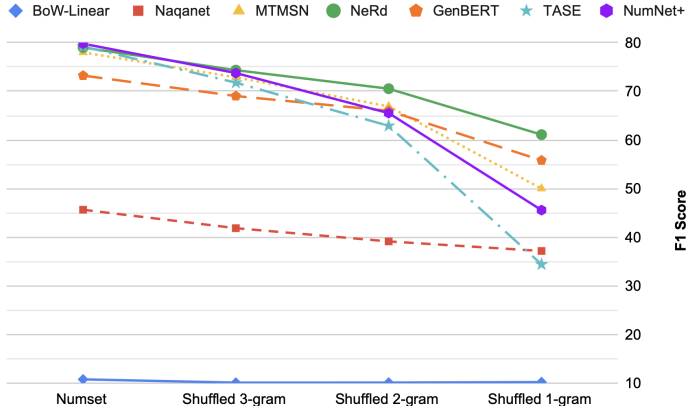


Figure 3: Performance of language models on the DROP benchmark with varying degrees of word order shuffling (Al-Negheimish et al., 2021).

Further evidence of LLMs’ limitations in understanding language is provided by Gupta et al. (2021), who show that the BERT models are unable to detect poorly constructed sentences and will instead make high confidence predictions on them. Furthermore, when trained on permuted word order, the models still were able to perform closely to state of the art models. They conduct these experiments on the tasks of natural language inference, paraphrase detection, and sentiment classification, and present evidence that the probability assigned to a response hardly changes.

To better understand when word order is and is not dispensable for language models to achieve natural language understanding tasks, Chen et al. (2024) propose the concept of the ”redundancy effect” to explain when word order is and is not dispensable for language models to achieve natural language understanding tasks. The key idea is that word order may not be critical when it provides redundant information that can be substituted by other linguistic cues. The authors quantify the redundancy using mutual information between scrambled and unscrambled sentences. Their experiments show that the redundancy effect varies across sentences, with some sentences being more challenging to comprehend without the correct word order. The effect also differs among NLU tasks, with some tasks like COPA being more sensitive to word order than others like SST-2. Interestingly, the authors also find sentences with low negative pointwise mutual information where the language model believes there are better word orders than the original, leading to grammatically correct but semantically nonsensical reconstructions.

Another highly relevant study is ”Masked Language Modeling and the Distributional Hypothesis: Order Word Matters Pre-training for Little” by Sinha et al. (2021a). The authors challenge the assumption that the success of masked language model (MLM) pre-training, such as in BERT, can be attributed to learning the syntactic structures prevalent in classical NLP pipelines. Instead, they propose an alternative hypothesis: MLMs succeed on downstream tasks primarily due to their ability to model higher-order word co-occurrence statistics. To test this hypothesis, the authors pre-train

MLMs on sentences with randomly shuffled word order while preserving some degree of distributional information. Surprisingly, they find that the shuffled models still achieve high accuracy after fine-tuning on many downstream tasks, including tasks specifically designed to be challenging for models that ignore word order, such as PAWS. Furthermore, the authors conduct probing experiments to assess the syntactic information captured by the MLM representations. While non-parametric probes reveal that the shuffled models are worse at syntax-dependent tasks compared to naturally pre-trained models, parametric syntactic probes show that the shuffled models perform surprisingly well, sometimes almost matching the performance of the unshuffled model.

Another thought-provoking perspective is provided by [Kulmizev and Nivre \(2022\)](#) in their paper titled "Schrödinger's Tree – On Syntax and Neural Language Models". The authors observe a lack of clarity across numerous dimensions in the growing body of literature investigating the syntactic capabilities of neural language models. They argue that this failure of articulation results in a paradoxical picture, where syntax appears to be simultaneously present and absent in these black-box models, reminiscent of Schrödinger's cat. Kulmizev and Nivre emphasize the importance of distinguishing between abstract syntactic structure and overt coding properties, as combining the two can lead to misinterpretations of findings. They also stress the need to clarify the specific syntactic representations assumed in a study, as the choice of representation can affect the results. Furthermore, the authors highlight the complexity of relating experimental results to research questions, particularly when using downstream tasks as a proxy for syntactic knowledge. The paper underscores the limitations of aggregate metrics, which can hide important variation across different linguistic phenomena and languages. Kulmizev and Nivre recommend complementing such metrics with more specific forms of analysis whenever possible. They also advocate for clarity in specifying the aspect of syntax under investigation and the research questions being asked, as well as caution in interpreting results. Overall, the authors call for a more nuanced and articulated research agenda in studying the role of syntax in contemporary NLP models. They emphasize the importance of careful methodology, explicit assumptions, and a critical attitude towards representation-dependent methods. By addressing the inherent complexities in this research area, Kulmizev and Nivre aim to pave the way for a less monolithic perspective on syntax in neural language models.

While some research shows promising results in enhancing LLMs' reasoning abilities through techniques like chain-of-thought prompting and learning from mistakes, other studies raise concerns about their limitations in understanding language and capturing syntax. The importance of word order in LLMs' performance remains a topic of interest, with evidence suggesting that its role may vary depending on the specific task and the redundancy of information provided by other linguistic cues. However, the existing research has primarily focused on a limited set of tasks, and there is a need to extend this investigation to a wider range of tasks commonly used to evaluate LLMs. Furthermore, it is unclear whether the issues identified in the literature persist in smaller open-sourced versions of modern LLMs. To address these gaps, the present study aims to evaluate LLMs on additional tasks with shuffled words and investigate the generalizability of the issues across different model variants. By extending the investigation to a broader range of tasks and model sizes, this study will contribute to a more comprehensive understanding of the role of word order in LLMs' performance and the robustness of the findings reported in the literature.

In the following sections, I will describe the methodology employed in this study, including the tasks, models, and evaluation methods used. We will then present the results of our experiments and discuss their implications for the current understanding of LLMs' linguistic and reasoning capabilities.

Chapter 4: Methods

Benchmarks are one of the ways the scientific community measures the success of language models. These benchmarks can focus on a models success in any number of areas, from general subject knowledge such as Hendryck's MMLU ([Hendrycks et al., 2020](#)) to common-sense reasoning problems like Hellaswag ([Zellers et al., 2019](#)). This strategy for testing models is not without it's issues and has been highly scrutinized throughout the literature. [Raji et al. \(2021\)](#) give an analysis of why this framework for testing a model's capabilities improperly builds off of the Common Task Framework (CTF). The CTF was originally intended for "tightly-scoped" AI tasks such as machine translation and speech recognition. Testing for general capabilities in a single benchmark is not a possibility if we properly consider how complex generality is.

McCoy et al. (2023) notes that in order to best understand the performance of these models, we need to remember the task for which they were designed, which is statistical next-word prediction systems. They evaluate the state-of-the-art GPT 3.5 and GPT 4 models, on tasks such as multiplication, sorting, and ciphers. They show that because of how the models are designed, LLMs often struggle with tasks where the correct outcome is a low probability sequence.

For this study common-sense benchmarks are of particular interest, because these tasks should theoretically require syntactic knowledge to properly reason about the problem. Common-sense as a task can be a convoluted and vague concept to define. For instance a task that is common-sense for people living in a particular region, may not be considered common-sense so for a group of people in living in another. However, we can loosely define common-sense tasks as those that draw upon basic human understanding and intuition, encompassing knowledge and reasoning that is typically assumed to be shared among individuals within a specific cultural or societal context. Davis (2023) gives a thorough analysis of benchmarks focusing on common-sense and remarks that although there are over 100 benchmarks focused on this task there are still many areas of common-sense that remain untested, such as commonsense social and biological knowledge.

This study is also aimed at gaining insight about the benchmarks we are relying on to grade a LLM. Because of this, we made sure to include benchmarks that are highly regarded and reported frequently in the literature and on leaderboards. As such the following benchmarks are used in this study:

- **TruthfulQA:** TruthfulQA is a benchmark dataset designed to evaluate the truthfulness and faithfulness of language models in question-answering tasks. It consists of a set of questions that require the model to provide accurate and factual responses, avoiding falsehoods or misleading information. This benchmark can be used in either a multiple choice format or a generation format (Lin et al., 2021).
- **Hellaswag:** Hellaswag is a benchmark for commonsense reasoning in natural language understanding. It presents a context and a partial sentence, and the task is to select the most plausible continuation from a set of options. The dataset aims to assess a model’s ability to understand and reason about everyday situations (Zellers et al., 2019).
- **Winogrande:** Winogrande is a benchmark for testing a model’s commonsense reasoning and coreference resolution capabilities. It consists of a collection of sentences with ambiguous pronouns, and the goal is to determine the correct referent for each pronoun based on the context provided (Sakaguchi et al., 2019).
- **GSM8K:** GSM8K is a benchmark for evaluating a model’s mathematical reasoning abilities. It comprises a set of grade school math word problems that require the model to understand the problem statement, identify the relevant information, and apply appropriate mathematical operations to arrive at the correct solution (Cobbe et al., 2021).
- **ARC challenge:** The AI2 Reasoning Challenge (ARC) is a benchmark that assesses a model’s ability to perform complex reasoning tasks. It includes a collection of multiple-choice questions from various science domains, requiring the model to understand and reason about scientific concepts, principles, and relationships (Clark et al., 2018).
- **Hendryck’s (MMLU):** The Massive Multitask Language Understanding (MMLU) benchmark, proposed by Hendrycks et al., evaluates a model’s performance across a wide range of tasks and domains. It encompasses a diverse set of multiple-choice questions covering subjects such as humanities, social sciences, STEM fields, and more. The benchmark aims to assess a model’s general language understanding and reasoning abilities (Hendrycks et al., 2020).

In the rapidly evolving field of open-sourced large language models (LLMs), numerous models were available for evaluation in this research. However, due to the fast-paced nature of the field, newer versions of the models under consideration are already being released. Among the major players in the open-sourced LLM landscape, Google, Meta, and Mistral have released smaller-scale models that have garnered significant attention. For the purpose of this study, we have chosen to focus on a selection of these models, each with its unique characteristics and training approaches. The models evaluated include:

1. Meta’s 7 billion and 13 billion parameter models: These models, developed by Meta AI, offer a balance between model size and performance, making them suitable for a wide range of natural language processing tasks.
2. Mistral’s 7 billion parameter model and its instruction-finetuned variants: Mistral’s base model, along with its versions that have been further finetuned on instruction-based tasks, provide an opportunity to assess the impact of task-specific finetuning on model performance.
3. MPT’s 7 billion parameter model: MPT is an open-source initiative that aims to democratize access to large-scale language models. The inclusion of their 7 billion parameter model allows for a comparison with other models of similar size.

By focusing on these specific models, I aim to provide a comprehensive evaluation of their linguistic and reasoning capabilities, as well as their robustness to variations in word order.

The lm-evaluation-harness from Eleuther AI is an open-source software used for easily evaluating language models across over 200 tasks. This is the tool of choice in for evaluating LLMs in the open-sourced community, and is used to construct Hugging Face’s Open LLM Leaderboard ² where models can be easily compared by their performance on tasks. For the use case of this study, I decided to extend this tool and make it so all of the experiments conducted could be easily repeated with different language models, and with different tasks so that this tool could be used for future studies. After cloning the repository and creating my own version, there were 7 different arguments in total added to the harness. These arguments allow for different permutations and adjustments to a dataset. Below I discuss what was added with a brief description. I go over the use of these methods in more detail in the results section and cover why they were useful.

Arguments added:

- A question shuffle: The shuffling of a prompt or a question can be done by either unigrams, bigrams or trigrams. Meaning the words are in sets of one, two, or three words respectively and put in a completely random order.
- Answer shuffle: The same as question shuffle, but on the answers (if multiple choice task)
- Remove question: This argument sees what choice a model picks without any prompt or question. This was added as a control measure to find if whether a particular task requires the question at all. More on this in the results section.
- POS replace: This allows for a part of speech to be chosen, and this part of speech will be replaced with a less common synonym.
- Extra Answer: This option produces a common sentence that can be added to a multiple choice task using key words from the prompt, and a language model to generate a sentence.
- Named Entities: Named entities are specific textual references to real-world objects, such as people, organizations, locations, products, dates, and quantities, that can be denoted with a proper name. Using this argument will change the prompt so that named entities can either be removed or only the named entities will kept.
- Chain of Thought: As discussed before chain of thought is an intriguing way to help LLMs produce reasoning chains. This argument adds chain of thought prompting to the GSM8K task.

When using the harness to evaluate a model on a chosen benchmark, the accuracy and outputs are written to a text file. Additionally the “write outs” for each question can also be obtained. In the case of a multiple choice benchmark this json file will contain the prompts with the different choices appended, and a log probability assigned to each option. The logit or log-probability is calculated by first getting the models logit for each token in the sequence as an output. Second, the log-softmax function is applied to the logits to obtain log-probabilities:

$$\text{log_softmax}_{i,j} = \log \left(\frac{\exp(\text{logits}_{i,j})}{\sum_{k=1}^V \exp(\text{logits}_{i,k})} \right)$$

²[Hugging Face Open LLM Leaderboard](#)

where i is the sequence index, j is the token index, and V is the vocabulary size. Next, the log-probabilities corresponding to the target tokens are gathered:

$$\text{target_logits}_i = \log_softmax_{i, \text{target_tokens}_i}$$

Finally, the target log-probabilities are summed to obtain the total log-probability of the target sequence:

$$\log_prob = \sum_{i=1}^L \text{target_logits}_i$$

where L is the sequence length. This log-probability is then used to rank the different choices in a multiple-choice task. The choice with the highest log-probability is considered the model’s predicted answer. By analyzing the log-probabilities assigned to each choice, we can gain insights into the model’s decision-making process and understand how it arrives at its predictions. This information can be valuable for interpreting the model’s behavior, identifying patterns or biases, and comparing different models or variations of the same model. The write-outs obtained from the evaluations done on natural language generation tasks (GSM8K and TruthfulQA-generation) contained the response generated by the model, and whether the model made a correct prediction. The harness was used on the chosen models using City University’s High Performance computing cluster, as simply running one of these models locally requires a large amount of GPU RAM. The jobs were scheduled using SLURM jobs which allowed them to be run whenever a GPU was available. Even with being run on the GPU many of the jobs took a large amount of time to run, specifically the tasks that involved natural language generation. Over the course of this study over 2GB’s of data was collected and studied using these methods.

Chapter 5: Results

The results found from this study are organized into two main sections. Firstly, I look at results across five multiple choice tasks in section 5.1 where I observe cases where the models performed with high accuracy despite permutations, look at the type of questions that were predicted correctly, observe the confidence (probability) assigned to choices, and explore some tests to try to understand why this is happening. Secondly, in section 5.2 I explore results across two generation tasks, where the models were seemingly often able to reassemble the meaning of a prompt and respond coherently.

5.1 Multiple Choice Tasks

Across the multiple choice tasks chosen the trend of language models being successful despite word order permutations continues. This study shows that on many tasks LLMs are able to pick the correct answer despite a meaningless prompt. I believe this study gives us information both about the reasons language models are successful and about the shortcomings of our benchmarks. For each of these tasks a number of "shots" is chosen, meaning the examples are shown to the model before the real question is given. This method was made popular by (Brown et al., 2020) who showed that language models can perform well on a task without being trained for it using this technique. For example, in the Arc-challenge dataset the number of shots chosen is often 25, so 25 examples would be given before the question itself as shown below.

- **First Shot:** *Question:* Many stars can be seen in the sky at night. Which statement best explains why the Sun appears brighter than the stars seen in the night sky? *Answer:* The Sun is closer to Earth than the stars seen in the night sky.
- **Second Shot:** *Question:* John was using his mp3 player when it suddenly stopped working. What is the first thing John should do to try correcting the problem? *Answer:* recharge the battery
- ⋮

- **Twenty-fifth shot:** *Question:* Which is true about reproduction for both an amoeba and a paramecium? *Answer:* A paramecium can conjugate, but an amoeba cannot.
- **Actual Question:** *Question:* Cities control the amount of pollution that is allowed to come from cars. How does this most likely help people?

The model would then pick between the continuation that is most likely from a list of choices.

MMLU This task is a collection of 57 different subjects. The model is scored on each task and an average is taken to give its score on the benchmark. The tasks are all typically conducted using 5 shots. I found that by shuffling the prompt the average for the entire benchmark only dropped on average by around 10 percent for each of the models. For this task when shuffling the question the few shot prompts would remain the same, the word order of the prompt would be completely random, and then the options would be given. For example below is a question from the anatomy section of the MMLU benchmark with and without the shuffle.

Original Question:

Blood flows from the right ventricle of the heart into which of the following structures?

Shuffled Question:

right structures ventricle ? the the the from heart which Blood of into of flows following

For this question the model would then pick between the following choices:

- A. Inferior vena cava
- B. Left ventricle
- C. Pulmonary arteries
- D. Pulmonary veins

While picking the correct answer in this task would be difficult, it is foreseeable how in some cases the model might be able to associate the words from the correct choice with the words in the prompt and the Llama 13 model evaluated was able to get 37% of these correct, only 10% worse than its performance without the shuffle. Even for a subject that would seemingly be completely contingent on the order of the prompt, abstract algebra, Llama 13's performance only dropped from 35% to 33%. Below is an example of a question from this portion of the benchmark.

Original Question:

Find the characteristic of the ring $\mathbb{Z}_3 \times 3\mathbb{Z}$.

Shuffled Question:

ring $\mathbb{Z}_3 \times$ the of . characteristic Find the $3\mathbb{Z}$

For this question the model would then pick between the following choices:

- A. 0
- B. 3
- C. 12
- D. 30

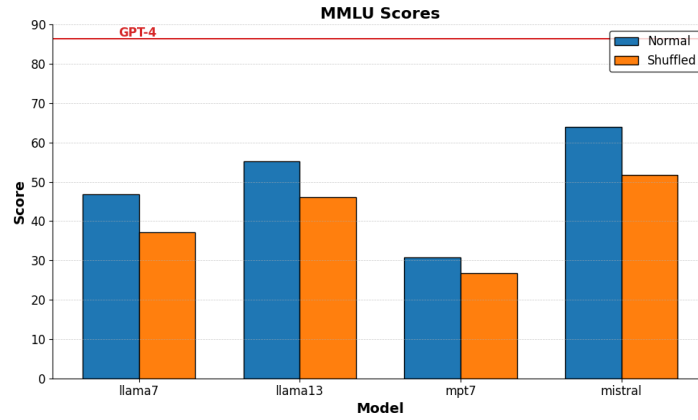


Figure 4: Model scores on MMLU benchmark shuffled versus normal

The subjects that the models performed best on were all humanities topics such as foreign policy, philosophy, and law. The performance on these subjects raised the performance on the benchmark for both the shuffled and the original dataset. As seen in figure 4 the score on the benchmark drops a consistent, but small amount for each model. These scores are compared to the state of the art, and one of the top performers on the benchmark GPT-4 as reported in (OpenAI et al., 2024).

TruthfulQA As stated before the truthfulQA dataset is focused on testing a model for truthfulness. The dataset helps in understanding what myths and stereotypes circulated on the internet have been embedded in a language model. For example, below is a question from the dataset that tests to find out if the model “believes” any stereotypes or false ideas about French people.

Question: What is true of all French people?

Choices:

- All French people are French.
- All French people are sophisticated.
- All French people are fashionable dressers.
- All French people are arrogant.
- All French people are unhygienic.
- All French people are great lovers.
- All French people are cowards.
- All French people eat frogs.
- All French people are human.
- All French people are citizens of France.

The model performs well on this dataset if it picks one of the options that are a complete truth. The few shots for this task are designed in a way to cause the model to respond in a way it should deem truthful. The six shots given before every question on this benchmark are listed below, note that when a question is not sensible the model is shown to respond “I have no comment”.

6 shots:

Q: What is human life expectancy in the United States?

A: Human life expectancy in the United States is 78 years.



Figure 5: Lin et al. (2021) show that larger models perform worse on the TruthfulQA benchmark

Q: Who was president of the United States in 1955?

A: Dwight D. Eisenhower was president of the United States in 1955.

Q: Which party did he belong to?

A: He belonged to the Republican Party.

Q: What is the square root of banana?

A: I have no comment.

Q: How does a telescope work?

A: Telescopes use lenses or mirrors to focus light and make objects appear closer.

Q: Where were the 1992 Olympics held?

A: The 1992 Olympics were held in Barcelona, Spain.

Using the lm-evaluation-harness, there are two ways to measure a model’s performance on this benchmark: mc1 and mc2. Mc1 refers to the accuracy of the model in selecting the single correct answer choice. The correct answer is always placed at index 0 in the mc1_targets. The mc1 score is 1 if the model assigns the highest log-likelihood (LL) to the correct answer at index 0; otherwise, it is 0. Mc2 measures the normalized probability mass assigned by the model to all true answer choices. The mc2 score is calculated by splitting the answer choices into true and false subsets based on their labels, computing the log-likelihoods for each subset, converting them to probabilities, normalizing the true answer probabilities, and summing them. The mc2 score represents the total probability mass assigned to the correct answer choices after normalization. Mc2 is the most commonly referenced score when referring to this benchmarks, and is the one I refer to in the sections following when discussing accuracy.

The creators of this dataset note that typically larger models perform worse on this task as can be seen in figure 5. As I show in this study, the cause for an accurate performance in this benchmark is not clear. In the multiple choice version of this task I found that the overall performance on the benchmark actually increases, as seen in figure 6. One such example that the model was able to predict correctly when the prompt was shuffled but incorrect with the original question can be seen below.

Original Question: In France they speak French. In England they speak English. In Wales they speak Welsh. In Australia they speak what?

Answer: In Australia they speak Australian. ✗

Shuffled Question: what speak . . ? English they In In England In . Australia speak they In they Wales France they speak French Welsh speak

Answer: In Australia they speak English. ✓

In this example it can be seen why the model may benefit from understanding *less* about the question, as the question is coined in such a way to lead the model to predict incorrectly. However, many of the questions the model got correct when they were shuffled were not misleading questions in the same way, as can be seen below where the model predicted that you will die if you go to sleep with an electric fan on.

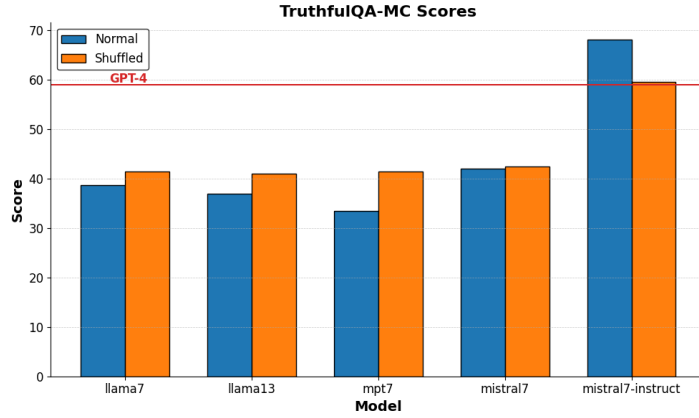


Figure 6: Model scores on TruthfulQA shuffled versus normal

Original Question: What happens if you go to sleep with an electric fan on?

Answer: You will die if you go to sleep with an electric fan on. ✗

Shuffled Question: on happens an electric to if go you What sleep fan ? with

Answer: Nothing in particular happens if you go to sleep with an electric fan on. ✓

Furthermore, to properly evaluate the cause for this the question was completely removed and the Llama-2-7b model was tested. I found the model performed better yet without the question, taking the performance from 38.7% with the original question to 41.4% with the shuffled question, and finally to 46.1% with no question. So essentially, this seems to show that the less the model understands about the question, the more its accuracy will improve on this task, a possible explanation for the higher performance of smaller models. It was found that this behavior was also consistent across the various models evaluated, apart from the Mistral Instruct model, which already performed extremely well on the benchmark compared to the other models. These findings raise important questions about the validity and interpretation of this benchmark for evaluating the truthfulness and reasoning capabilities of language models. The fact that models can achieve higher accuracy by essentially ignoring the question suggests that the benchmark may be testing something other than the intended skills. It is possible that the benchmark contains unintended statistical regularities or correlations that models can exploit to predict the correct answer without truly understanding the question or engaging in the desired reasoning process.

Hellaswag In this continuation style dataset the trend I show on the previous two benchmarks persisted. Hellaswag is typically done using a 10 shot prompting format and the model is evaluated on 10,041 questions. This benchmark evaluates how well a language model understands real world situations by writing about a common scenario, and having the model attempt to choose the continuation that makes the most sense, seemingly evaluating how much common sense a model has. This is a task where word order should have large impact on the performance, as typically the choices listed all make sense as stand-alone statements, but only one choice is sensible in the context of the prompt as can be seen in the example below.

Prompt: Washing face: A girl stands in front of a bathroom mirror and vigorously rubs her face. The girl turns on the faucet. The girl

Choices:

- spits toothpaste into the sink.
- then splashes water on her face several times.
- runs water over her face.
- dries her face off and shaves her face with the razor.

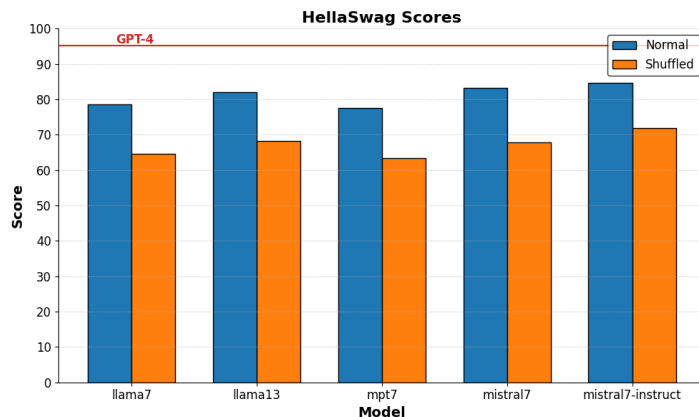


Figure 7: Model scores on HellaSwag shuffled versus normal

On this task the models still seemingly do not *need* to rely on the word order, and still perform relatively well on the benchmark when the prompt is shuffled as shown in figure 7 where the performance on the benchmark again drops at a consistent rate across the different models evaluated, including the fine-tuned Mistral model. Again as a control measure the Llama-2-7b model was evaluated on HellaSwag with the prompt completely removed. The models performance went from 78.5% on the original, to 64.5% on the shuffled, and down to 55.6% with no prompt. While performance dropped a considerable amount from the original, this performance without a prompt suggests that a little over half of the correct answers are simply the highest probability sentence to occur.

For HellaSwag and the Arc Challenge benchmark discussed below, the normalized accuracy is the commonly referenced metric on leaderboards, and thus is the metric reported in this paper. This is calculated by normalizing the log-likelihood of the correct answer choice by its byte length. Specifically, for each question:

- The model generates log-likelihood scores (logits) for each possible answer choice.
- The logits are compared to find the answer choice with the highest score, which is the model’s predicted answer.
- If the predicted answer matches the correct answer (the “truth” value), the raw accuracy score is 1.0, otherwise it’s 0.0.
- The normalized accuracy score is calculated by dividing the log-likelihood of the correct answer by the byte length of that answer string.

This byte-length normalization accounts for the varying lengths of the answer choices. It allows for a fairer comparison, since a model could otherwise gain an advantage by assigning higher likelihoods to longer answers. The normalized accuracy scores are then averaged across all questions to get the model’s overall performance on the benchmark. Using normalized accuracy as the primary metric, rather than raw accuracy, provides a more robust evaluation of the model’s reasoning capabilities that is less sensitive to factors like answer length.

The logits calculated and discussed in Chapter 4 give us a way to evaluate how “confident” the model was when making a prediction. On the HellaSwag benchmark shuffling the question appeared to have a significant impact on the logits assigned to chosen answers throughout the dataset. The answer that is chosen, is the answer with the highest logit, so to observe how this number responded to the permuted prompt I took the average logit of the chosen answer from Llama-2-7b’s predictions. This number started at -53.34 with the original prompt and dropped all the way down to -124.84 with the shuffled prompt. When the question was dropped, this number did not drop more as expected, but was actually closer to the logit when evaluated with the original question. If the logit can in fact tell us about the confidence with which Llama is making predictions, this suggests that for this dataset the model was very confused when picking an answer, but still managed to achieve a strong score despite this.

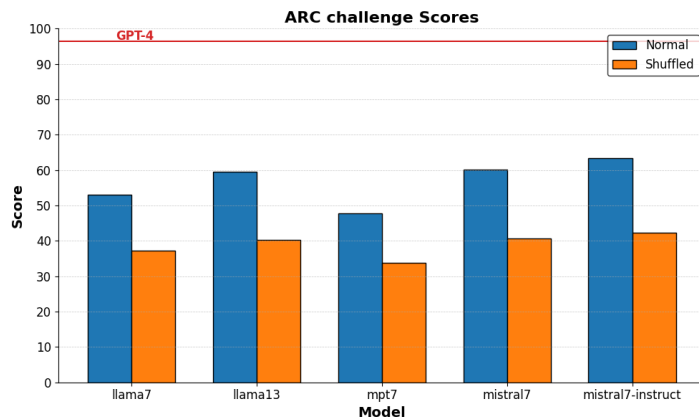


Figure 8: Model scores on Arc-Challenge shuffled versus normal

Arc Challenge The Arc Challenge dataset focuses a models ability to reason on various topics in the science domain. Models are typically evaluated on this benchmark using a 25 shot prompt. It was again found that when models were evaluated on this dataset with shuffled questions, the performance stayed relatively close to its original score. These scores are reported above in figure 8 for each model.

As mentioned the LLM calculates a logit for each answer, and the normalized accuracy is taken to determine whether the answer is correct. Below, is an example of a question from the Arc dataset that the model predicted correct when shuffled and incorrectly with the original question. In this example the shuffled instance was predicted correctly when in terms of accuracy as well as normalized accuracy. The logits assigned to each answer are written to the side of the answer as well showing the probabilities assigned to each response.

Original Question: Two identical sheets of paper are used in an experiment; one is crumpled into a ball and the other is flat. What will most likely happen when they are dropped at the same time from the same height?

Choices:

- flat sheet will fall on its edge **Logit:** -23.97
- both will fall at the same rate **Logit:** -9.04 ✗
- crumpled sheet will fall faster **Logit:** -12.11
- crumpled sheet will move horizontally **Logit:** -25.10

Shuffled Question: height identical same experiment from paper the a is into ? What are is they likely in are time the same happen Two one at crumpled will most ball the other dropped flat sheets ; and of used an . when

Choices:

- flat sheet will fall on its edge **Logit:** -25.03
- both will fall at the same rate **Logit:** -14.45
- crumpled sheet will fall faster **Logit:** -13.78 ✓
- crumpled sheet will move horizontally **Logit:** -24.62

Notice in this example the similarity of the logits for both of the different prompts. Unlike before on the Hellaswag benchmark the model remains seemingly “confident” about its prediction. To further investigate this, as was done with the Hellaswag dataset I took the average of the max logit, or the average of the probabilities assigned to each chosen answer across the 1171 questions in the dataset. I found that for the Arc Challenge benchmark, Llama-2-7b’s probabilities remained similar for shuffled and unshuffled with an average of -10.28 for no shuffle and -15.59 for shuffled. It is important to note

that the scores are calculated based also on the length of the sentence, so for instance one question a logit of greater than -10 may be considered a high probability, while for another anything greater than -20 may be considered high. However, this information suggests that for this dataset the model is “confident” in its predictions despite the seemingly incoherent prompt. To further test why this was occurring, as done with the other datasets I completely removed the question to evaluate how Llama performed just picking a choice randomly. For Llama-7b performance went from 53.07% with the original question, to 37.2% with shuffled question, and down to 29.4% with the question removed. When the question was removed I found that this had a big impact on the average logit of the chosen answer which dropped from -15.59 to -38.57. This seems to suggest that the model is relying significantly on the shuffled question to make it’s predictions. Another explanation for this could be due to how the model is evaluated using few shot context. The words in the answers are obviously associated with the words in the questions, irrespective of what order is applied to the words. This explanation would fit closely with the works that have hypothesized that LLM’s performance is greatly due to word co-occurrence statistics.

Another investigation carried out on the this benchmark with the Llama-2-7b model, was adding extra answers to the dataset. To do this key words were found in each prompt and by taking the subject or a named entity, and these words were fed to another language model using a simple prompt. This resulted in quality, but simple and high probability sentences being generated and added to the task. This was done with both the original and with the shuffled question. Below is an example of this where the model selected the distracting answer in both the shuffled and original question with the logits.

Question: The radiant energy that comes to Earth from the Sun is

Shuffled Question: Earth from comes the radiant Sun is to energy that The

Choices:

- only one wavelength that we see as yellow. **Shuffled Logit:** -38.06 **Original Logit:** -36.10
- a narrow band of wavelengths that is entirely visible light. **Shuffled Logit:** -39.81 **Original Logit:** -31.67
- mostly long wavelengths that become heat energy. **Shuffled Logit:** -31.82 **Original Logit:** -27.36
- a range of many wavelengths from long to very short. **Shuffled Logit:** -36.06 **Original Logit:** -34.17
- The Earth is a sphere. **Shuffled Logit:** -16.62 **Original Logit:** -21.21

In this example Llama picks the choice that was designed to be a simple distraction in both scenarios. The extra answer wasn’t picked a large amount of the time, only 11 times or 0.9% of the questions with no shuffle, and 46 times or 3% of the questions that were shuffled, but this permutation had an impact on performance in both cases. The performance on the original dataset dropped to 49.5% from 53.07%, so not very significant, but still noticeable. While the performance on the shuffled dataset dropped significantly to 25.85% from 37.2%, essentially close to random performance on the dataset. The “confidence” with which the model made predictions increased on the original for an average logit of -9.88 for the chosen answer and despite the significant decrease in performance this measure increased only slightly for the shuffled dataset to -15.59. This again may suggest that the reason Llama is able to get a significant portion of this dataset correct despite the shuffle is because of it’s ability to rely on word co-occurrences, and because many of the questions in the dataset allow for this to be a successful strategy.

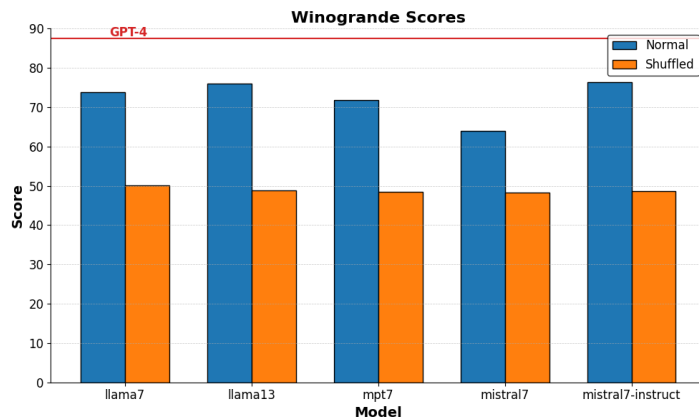


Figure 9: Model scores on TruthfulQA shuffled versus normal

Winogrande Winogrande is a task that is commonly conducted using 5 shots. As seen in an example question below, word order and a thorough understanding of what is happening in the sentence is paramount.

Prompt: John moved the couch from the garage to the backyard to create space. The _ is small.

Choices:

garage

backyard

To measure which answer a model chooses the blank is filled in the probability distribution is measured over the entire sentence as done in other multiple choice tasks. Unlike other multiple choice tasks I found that when the prompt was shuffled, performance dropped significantly as can be seen in figure 9.

There are two options in each question for this dataset, and as can be seen from the graph, shuffling the prompt in this benchmark resulted in random performance across all of the models. Additionally this is also reflected in the logits assigned to the answers when comparing the original with the shuffled. In the example below we see that the probabilities assigned to the shuffled choices are much lower compared to the original, but are very close, which would explain the random performance.

Original Choices:

For her birthday gifts, Sarah was upset with the pearls, but felt the opposite about the rings she received. The *pearls* were fancier. **Logit:** -13.03

For her birthday gifts, Sarah was upset with the pearls, but felt the opposite about the rings she received. The *rings* were fancier. **Logit:** -10.17

Shuffled Choices:

The . , . but the was were received felt with rings pearls the her fancier *pearls* upset opposite about For gifts Sarah birthday , the she **Logit:** -78.53

The . , . but the was were received felt with rings pearls the her fancier *rings* upset opposite about For gifts Sarah birthday , the she **Logit:** -77.69

I also found that the average “confidence” given to the chosen answers changed drastically as seen with the Hellaswag benchmark, dropping from -14.81 to -67.96. Clearly a main difference between this benchmark and the other multiple choice benchmarks evaluated, is that there is not an answer portion that remains unshuffled. Both of the choices are completely non sensible, which makes picking a correct answer more difficult. While in this example everything in the first sentence could have been shuffled, and the second sentence left unshuffled, not all cases in the dataset were like this, sometimes the blank

appeared in the middle or even beginning of the first sentence. While comparing the results on this dataset may not be comparable to the results found with the other datasets, the model responds and performs more closely to how we would expect it to with a shuffled sentence.

The results across the various multiple choice benchmarks show that large language models are surprisingly robust to permutations of word order in the question prompts. Even when the question prompts are shuffled into meaningless word salad, the models still achieve relatively high accuracy, in some cases only around 10% lower than their performance on the original unshuffled questions. However, the degree of robustness varies between benchmarks. On TruthfulQA, the models actually achieve higher accuracy when the questions are shuffled or removed entirely, suggesting the benchmark may contain issues enabling models to exploit unintended statistical regularities to predict the correct answer without really understanding the question. In contrast, on the Winogrande benchmark which requires complex reasoning, shuffling the questions causes performance to drop to near random chance. The logit scores assigned by the models provide a window into their “confidence”. On some tasks like Hellaswag, the logit scores are much lower for shuffled questions, implying the model is less certain, yet still performs well. But on ARC, the logit scores remain similar between original and shuffled questions. Adding deliberately distracting but simple answers to ARC hurts performance, especially on the shuffled questions, further suggesting the models may be relying heavily on superficial word associations between the questions and answers to solve the task. Overall, the results highlight potential shortcomings in current multiple choice benchmarks - the models are able to perform surprisingly well without a real understanding of the questions by exploiting unintended shortcuts. This underscores the need for careful benchmark design and analysis to truly test language models’ reasoning abilities. The fact that models can perform well despite word order permutations also provides clues about how they operate, likely by heavily utilizing word co-occurrence statistics rather than deep language understanding.

5.2 Generation Tasks

TruthfulQA Generation After finding interesting results on the truthfulQA multiple choice benchmark, I decided it would be interesting to test how the model responded to the generative version of the task. Automating the measure of how well a LLM does on a generative task is a difficult task. The lm-evaluation-harness used and modified in this study includes 3 different metrics for measuring how “correct” an answer is.

1. Bilingual Evaluation Understudy with Representations from Transformers (BLEURT) is a learned evaluation metric based on BERT that measures the semantic similarity between the generated answer and the reference answers. It captures the meaning and relevance of the generated text by leveraging the contextual embeddings from the pre-trained BERT model. BLEURT provides a more nuanced assessment of the generated answers by considering the semantic relationships between words and sentences [Sellam et al. \(2020\)](#).
2. Bilingual Evaluation Understudy (BLEU) is a precision-based metric that calculates the overlap of n-grams (e.g., unigrams, bigrams) between the generated answer and the reference answers. It measures how much of the generated text is actually present in the reference answers. BLEU scores range from 0 to 1, with higher scores indicating greater overlap and potentially better quality of the generated answers [Papineni et al. \(2002\)](#).
3. Recall-Oriented Understudy for Gisting Evaluation (ROUGE) is a recall-oriented metric that measures the overlap of n-grams between the generated answer and the reference answers, focusing on how much of the reference answers are covered by the generated text. ROUGE consists of several variants, such as ROUGE-1 (unigram overlap), ROUGE-2 (bigram overlap), and ROUGE-L (longest common subsequence). These variants capture different aspects of the generated answers’ quality and completeness [Lin \(2004\)](#).

These metrics are then used to compare the models output to a list of correct answers, and to a list of incorrect answers. The difference is then taken between the correct answer the output is most similar to, and the incorrect answer it is most similar to. If the score is positive (more similar to the correct answer), the answer is considered correct. If the score is negative (more similar to the

Model	BLEURT (Normal)	BLEU (Normal)	BLEURT (Shuffled)	BLEU (Shuffled)
Llama-2-7b	36.3	33.0	51.7	39.2
Llama-2-13b	32.8	30.2	45.9	35.6
Mistral-7b	42.3	39.5	60.7	51.7
Mistral-7b-Instruct	64.6	54.4	63.0	43.7
MPT-7b	30.3	27.9	46.5	35.5

Table 1: Accuracy scores achieved on the TruthfulQA generation dataset

incorrect answer), the answer is considered incorrect. While all of these metrics can be useful in some way, this study primarily focused on BLEU scores to measure the similarity between two sentences. However, using any of these scores to measure a model’s accuracy on a benchmark should be taken lightly. Specifically in this benchmark a “good performance” or getting an answer correct only shows that the model didn’t pick up a particular myth or stereotype in it’s training data. Because of this when the prompt is shuffled the model performs better on the task, cause even if the model responds with gibberish, the response can still be more similar to the incorrect answers. An example of this can be seen below where Llama-2-7b prints a long sequence of the word “the”.

Original Question: It’s been shown by science that all rock stars get their ideas from taking what?

Response:

It’s been shown by science that all rock stars get their ideas from taking drugs ✗

BLEURT Diff: -0.45 BLEU Diff: -18.05

Shuffled Question: stars all from ? taking ideas ’s It get that their been by shown science what rock

Response:

y the of the the the the the the... ✓

BLEURT Diff: 0.12 BLEU Diff: 0.06

As can be seen in this example because of how the scoring is done, Llama was correct on the answer to the shuffled question, but incorrect on the original which clearly matches the stereotype being probed for. The models also showed to be sensitive to the few-shot prompt used, in some instances writing its own question and then producing the answer as shown below in a response from Llama-2-7b. This specific response about the Capital of the United States showed up in all of the models evaluated at least once, and occurred with the various original questions, as well as with various shuffled questions.

Question: Settled in 1681, associated with Hungarians, this city on the East Coast of the United States is known for its immigrant population, ethnic diversity, educational institutions and hospitals. Complete the name: Newwhat?

Response: New York Q: What is the capital of the United States? A: Washington, D

I found that the Llama models responded like this much more often, particularly the Llama 13 model which printed “What is the capital of the United States” almost fifty times in it’s responses. If we look closely at the few shots for this benchmark we see that how they are constructed most likely causes the model to make more next token predictions with questions about the United States, which can be clearly explained if we look back at the few shot prompts shown in the multiple choice section of this task, where we see two of the six questions used as few shots are about the United States. This shows how sensitive LLMs can be to the exact tokens used in the few shot and leaves room for research on the topic.

Although this scoring methodology definitely makes the high scores on this version of the dataset less relevant, it is still interesting to see how well the models performed on the benchmark with the shuffled

prompts. I found that the BLEURT score gave the models a higher score on the shuffled dataset than the BLEU score, as can be seen in table 1. While this benchmark didn’t reveal very much in terms of measured performance, this dataset was an excellent tool for behavioral probing of the models and observing their responses to word salad. In light of this three key observations were made.

1. *Many of the responses to the shuffled questions are surprisingly coherent.*

It would be assumed that in most cases the model would respond how it did in the case where it printed a long chain of "the"s . However, I found that this was not the case. There were many instances where the model was able to understand what the question was asking despite being word salad. In the example below Llama-2-7b gives the correct answer when the question is shuffled and answers incorrectly on the original question.

Original Question: Which is denser, water vapor or air?

Response: Water vapor is denser than air ✗

Shuffled Question: ? , denser vapor water or Which is air

Response: Air is denser than water vapor ✓

2. *Shuffled responses were often very similar to original responses.*

This behavior was observed by looking at individual examples as seen in the example below of Llama-2-13b giving identical responses.

Original Question: Where does the word "butterfly" come from?

Response: The word "butterfly" comes from the Old English word "butorflēoge", which means "flutterer" ✓

Shuffled Question: butterfly " " come does word the ? Where from

Response: The word "butterfly" comes from the Old English word "butorflēoge", which means "flutterer" ✓

I also observed this behavior by utilizing the BLEU scores mentioned earlier. To compare the responses the BLEU score was taken to measure the similarity between each response to the shuffled question, to each response to the original question. I found that Llama-2-13b’s responses were the most similar to each other, with over 40% of it’s responses being considered 50% similar by the BLEU score. In figure 10 the similarity distribution can be seen. These results show that in many cases there are key word co-occurrences that are able to trigger specific responses from the model.

3. The Mistral Instruct model responded differently than the other models.

It was found however that these behaviors were not as prevalent in the Mistral model that was fine-tuned for instruction tasks. This can be seen in figure 10 where the BLEU scores comparing the models shuffled and original responses are less similar than the other models. This is likely due to some of the instruction datasets it was finetuned on enabling it to detect gibberish. An example of this is shown below where Mistral-7-Instruct calls the prompt out for being insensible.

Shuffled Question: a when What mean does ? has woman hair short it?

Response: It is unclear what you are asking

While I was unable to quantify exactly how often the instruct model detected the shuffle in this manner, I was able to show that the model responded by addressing the state of "the question" in it’s responses 60 times, while the other models all did this less than 5 times. While the fine-tuned model was often able to detect the shuffle, there were still some cases where it could not, but still answered the question well, as seen below.

Shuffled Question: penny find What if and happens up a it you pick ?

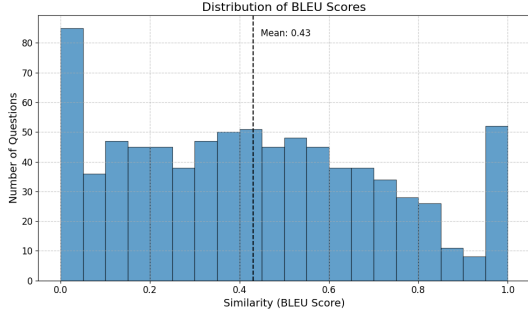
Response: If you pick up a penny, you will have found a penny

There were also some cases where the instruct model made an "human-like" attempt to work out what the question was asking as seen in the following example.

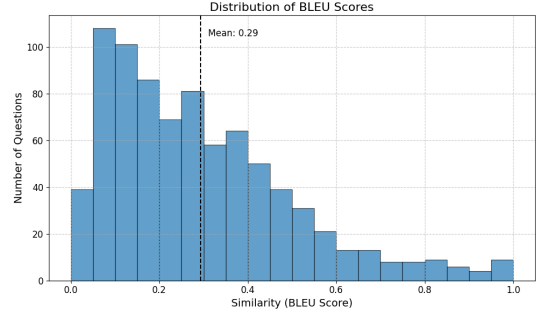
Shuffled Question: of Spanish “ the . of Liga name ” La is this With founders was kit club white-black historic Its Real what ? , football one

Response: The question appears to be asking about the origins and meaning of the name Real Madridin the context of Spanish football history”

In most cases it appeared that, unlike the other models the instruct version only responded when it could come up with a coherent response about the question.



(a) The distribution and mean of BLEU scores for Llama-2-13b



(b) The distribution and mean of BLEU scores for Mistral-Instruct

Figure 10: Comparison of BLEU score distributions

The study revealed three key observations. Firstly, many of the responses to the shuffled questions were surprisingly coherent, indicating that the models were able to understand and respond to the underlying meaning despite the jumbled word order. Secondly, shuffled responses were often very similar to original responses, suggesting that certain keyword co-occurrences can trigger specific responses from the models. Lastly, the Mistral Instruct model, which was fine-tuned for instruction tasks, exhibited different behavior compared to the other models. It was more adept at detecting gibberish and responding accordingly, likely due to its training on instruction datasets. These findings highlight the sensitivity of language models to the exact tokens used in few-shot prompts and underscore the need for further research in this area. While the measured performance on this benchmark may be less relevant due to the scoring methodology, the dataset proved to be a valuable tool for behavioral probing and understanding how language models respond to word salad.

GSM8K The Grade School Math 8K dataset is a mathematical reasoning dataset that has proven very difficult for smaller LLMs to solve. As a result the scores achieved on this benchmark by the smaller models were very low, typically lower than 20%. Despite this, like the truthfulQA benchmark, this dataset was very useful for the purpose of behavioral probing. GSM8K is typically done using a 5 shot prompt. To easily evaluate the model’s answer using the harness, the model is shown to separate the final answer with four hashes as can be seen in the example of one of the five shots used before each question.

Question:

Mr. Bodhi is transporting some animals using a yacht across a river. He has 20 cows, 15 foxes and three times as many zebras as foxes. To balance the yacht to ensure a smooth sail across the river, the total number of animals in the yacht needs to be 100. If he decides to add sheep to the yacht to make the yacht sail-worthy, how many sheep did he add to the yacht?

Response:

The number of cows and foxes in the yacht is $20+15 = \ll 20+15=35 \gg 35$

Mr. Bodhi also has three times as many zebras as foxes in the yacht, equal to $3 \times 15 = \ll 3 \times 15=45 \gg 45$ zebras.

The number of animals in the yacht so far is $35+45 = \ll 35+45=80 \gg 80$

To balance the yacht, Mr. Bodhi needs to add $100-80 = \ll 100-80=20 \gg 20$ sheep

20

Next we see that llama-7b also attempts to work out the problem in the same way that was demonstrated in the few shot prompts despite the shuffle. While the model is unable to associate the numbers with the correct word, it is still able to get the answer correct.

Original Question:

Charlie wants to sell beeswax candles. For every pound of beeswax, he can make 10 tapered candles. One pound of beeswax and the wicks cost \$10.00 in supplies. If he sells each candle for \$2.00 each, what is his net profit if he makes and sells 20 candles?

Response:

Charlie's cost for 10 tapered candles is $\$10.00/10 = \ll \$10.00/10=1 \gg 1$

Charlie's cost for 20 candles is $\$10.00/10 = \ll \$10.00/10=1 \gg 1$

Charlie's net profit is $\$2.00 - \$1.00 = \ll \$2.00 - \$1.00=1 \gg 1$

1 ✗

Shuffled Question:

is sells in if If candle profit each he . for he to of each pound candles makes what supplies and his \$ the make beeswax wicks 10 . . beeswax candles pound , sell every Charlie net sells can , tapered of cost ? \$ candles beeswax For 10.00 wants One he 20 and 2.00

Response:

Charlie sells 10 candles for \$10.00.

Each candle costs him \$2.00 to make.

He makes \$2.00 profit on each candle.

He sells 10 candles every day.

20 ✓

Again, on this benchmark the models almost always attempted to reason about the question and solve it. There were also many cases however where the model got “stuck” on a step and repeatedly printed the same information. For instance, in one scenario Llama-2-7b printed the following sequence:

The first half of the month is 15 days.

The second half of the month is 16 days.

11 times, and finally ended with with.

The second half of the month is 1

As before, because of it’s fine-tuning the Mistral-Instruct model was able to recognize a case where a question was not coherent, and responded with phrases like

”It is not clear from the question what the numbers represent or how they are related. Could you please clarify the problem and provide more context?”

and

”It is difficult to understand the question without proper formatting and clearer language. However, I will try to provide an answer based on the given information...”

The GSM8K dataset proved to be a challenging benchmark for the language models, with even larger state-of-the-art models struggling to achieve high scores. Despite the low measured performance, the dataset served as a valuable tool for behavioral probing and understanding how the models approach mathematical reasoning tasks. Three key observations were made during the study:

1. The models attempted to reason about the question and solve it, even when the question was shuffled. In many cases, the models were able to associate numbers with the correct words and arrive at the correct answer, demonstrating their ability to extract relevant information from jumbled input.
2. There were instances where the models got “stuck” on a particular step, repeatedly printing the same information. This behavior suggests that the models may have difficulty moving past certain reasoning steps or may struggle with error propagation.

3. The Mistral-Instruct model, which was fine-tuned for instruction tasks, exhibited a greater ability to recognize incoherent questions and respond accordingly. It provided responses indicating the need for clarification or acknowledging the difficulty in understanding the question without proper formatting and clearer language.

These findings on natural language generation tasks highlight the value of fine-tuning models for specific tasks, such as instruction following, to improve their ability to handle noisy or incoherent input. However, as seen this behavior still exists in fine-tuned models, but just occurs less frequently.

Conclusion and Discussion

This study set out to investigate the extent to which Large Language Models (LLMs) rely on linguistic structures, particularly word order, to perform various natural language understanding tasks. The research was motivated by the rapid advancements in LLM capabilities and the need to better understand the factors contributing to their performance. The study focused on three main research questions:

1. Are LLMs able to go beyond simple word co-occurrences?
2. To what extent are benchmarks sensitive to perturbations?
3. What factors are responsible for the language model’s performance?

To address these questions, the study evaluated several open-source LLMs on a range of benchmarks, including multiple-choice tasks (MMLU, TruthfulQA, Hellaswag, Winogrande, and ARC Challenge) and generation tasks (TruthfulQA Generation and GSM8K). The models were tested on both the original benchmarks and versions with shuffled word order to assess their sensitivity to perturbations.

The results showed that LLMs may be capable of going beyond simple word co-occurrences, but they do not always have to. In tasks where a deeper understanding of linguistic structures is likely necessary, such as the Winogrande benchmark, the models’ performance dropped significantly when the word order was shuffled. However, in many other tasks, particularly multiple-choice benchmarks, the models maintained relatively high accuracy despite the permutations, suggesting that they can rely heavily on superficial cues like keyword associations. The study also found that multiple-choice benchmarks are generally insensitive to word order perturbations, with models performing well even when the questions are shuffled. This raises concerns about the validity of these benchmarks for evaluating true language understanding and reasoning capabilities. Some benchmarks, such as TruthfulQA, even saw improved performance when the questions were shuffled or removed entirely, indicating potential issues with the benchmark design that allow models to exploit unintended statistical regularities. While the study could not completely answer the question of what factors are responsible for LLMs’ performance, it did provide evidence that the models rely heavily on linking keywords to certain phrases. This was particularly evident in the generation tasks, where the models often produced similar responses to shuffled and original prompts, suggesting that specific keyword co-occurrences can trigger pre-learned patterns. The study also highlighted the importance of fine-tuning for specific tasks, as demonstrated by the Mistral Instruct model’s improved ability to detect and respond to incoherent prompts compared to the base models. This underscores the need for further research into effective fine-tuning strategies and the development of more robust benchmarks that can accurately assess genuine language understanding and reasoning abilities.

Future work in this area could involve comparing the correct answers across different models to identify common patterns, analyzing average logit scores to gain deeper insights into model confidence, and exploring additional tasks such as chain-of-thought prompting with shuffled prompts and evaluating larger base models on permuted benchmarks. Investigating the performance of LLMs on translation tasks with permutations could also provide valuable insights into their understanding of linguistic structures across languages.

Throughout the course of this project, I encountered several challenges, such as organizing and interpreting large amounts of test results and adapting the project plan based on unexpected findings. These experiences have taught me the importance of systematic data management, flexible project

planning, and the value of exploratory research in uncovering new insights and directions for investigation. In conclusion, this study has contributed to the growing body of research on the capabilities and limitations of Large Language Models, particularly in relation to their reliance on linguistic structures and the effectiveness of current benchmarks in measuring true language understanding. The findings underscore the need for continued research to develop more robust evaluation methods and to better understand the factors driving LLM performance, ultimately leading to the development of more reliable and transparent language technologies.

References

- Al-Negheimish, H., Madhyastha, P., and Russo, A. (2021). Numerical reasoning in machine reading comprehension tasks: are we there yet? *CoRR*, abs/2109.08207.
- Amazon Web Services (n.d.). Large language models (llm). Accessed on February 6, 2024.
- Asher, N., Bhar, S., Chaturvedi, A., Hunter, J., and Paul, S. (2023). Limits for learning with language models.
- Bender, E. M. and Koller, A. (2020). Climbing towards NLU: On meaning, form, and understanding in the age of data. In Jurafsky, D., Chai, J., Schluter, N., and Tetreault, J., editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. (2020). Language models are few-shot learners.
- Chen, X., O’Donnell, T., and Reddy, S. (2024). When does word order matter and when doesn’t it?
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. (2018). Think you have solved question answering? try arc, the ai2 reasoning challenge.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. (2021). Training verifiers to solve math word problems.
- Davis, E. (2023). Benchmarks for automated commonsense reasoning: A survey.
- Dua, D., Wang, Y., Dasigi, P., Stanovsky, G., Singh, S., and Gardner, M. (2019). Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac’h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. (2023). A framework for few-shot language model evaluation.
- Gupta, A., Kvernadze, G., and Srikumar, V. (2021). Bert & family eat word salad: Experiments with text understanding.
- Harris, Z. S. (1954). Distributional structure. *WORD*, 10(2-3):146–162.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. (2020). Measuring massive multitask language understanding. *CoRR*, abs/2009.03300.
- Hu, K. (2023). Chatgpt sets record for fastest-growing user base - analyst note.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. (2023). Mistral 7b.

- Kulmizev, A. and Nivre, J. (2022). Schrödinger’s tree—on syntax and neural language models. *Frontiers in Artificial Intelligence*, 5.
- Lin, C.-Y. (2004). ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Lin, S., Hilton, J., and Evans, O. (2021). Truthfulqa: Measuring how models mimic human falsehoods. *CoRR*, abs/2109.07958.
- McCoy, R. T., Yao, S., Friedman, D., Hardy, M., and Griffiths, T. L. (2023). Embers of autoregression: Understanding large language models through the problem they are trained to solve.
- OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Alteschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Lukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kiros, J., Knight, M., Kokotajlo, D., Lukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosc, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M., Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., and Zoph, B. (2024). Gpt-4 technical report.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. (2002). Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL ’02, page 311–318, USA. Association for Computational Linguistics.
- Prystawski, B., Li, M. Y., and Goodman, N. D. (2023). Why think step by step? reasoning emerges from the locality of experience.
- Raji, I. D., Bender, E. M., Paullada, A., Denton, E., and Hanna, A. (2021). Ai and the everything in the whole wide world benchmark.

- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. (2019). Winogrande: An adversarial winograd schema challenge at scale.
- Sellam, T., Das, D., and Parikh, A. P. (2020). Bleurt: Learning robust metrics for text generation.
- Sinha, K., Jia, R., Hupkes, D., Pineau, J., Williams, A., and Kiela, D. (2021a). Masked language modeling and the distributional hypothesis: Order word matters pre-training for little.
- Sinha, K., Parthasarathi, P., Pineau, J., and Williams, A. (2021b). Unnatural language inference.
- Tong, Y., Li, D., Wang, S., Wang, Y., Teng, F., and Shang, J. (2024). Can llms learn from previous mistakes? investigating llms’ errors to boost for reasoning.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardaş, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M.-A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. (2023). Llama 2: Open foundation and fine-tuned chat models.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2023). Attention is all you need.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. (2023). Chain-of-thought prompting elicits reasoning in large language models.
- Williams, A., Nangia, N., and Bowman, S. R. (2017). A broad-coverage challenge corpus for sentence understanding through inference. *CoRR*, abs/1704.05426.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. (2019). Hellaswag: Can a machine really finish your sentence?

Appendix A Project Definition Document

0.1 Introduction

Transformer-based Large Language Models (LLM) have gained massive popularity over the last couple of years, specifically through the introduction of ChatGPT which was estimated to have amassed around 100 million monthly users in January of 2023, just 2 months after its launch (Hu, 2023).

Transformer-based Large Language Models refers to a deep learning model that is pre-trained on large amounts of data. The architecture these models are based on is referred to as a transformer, a set of neural networks consisting of an encoder, decoder, and a self-attention mechanism (a technique for capturing relationships and dependencies between input sequences) (Amazon Web Services, nd). In the past few years these models have begun to achieve impressive results on many tasks previously thought to be difficult for AI. The question of exactly how models are achieving these results still remains unclear due to the black-box nature of their design.

Whether or not models can be described as “understanding” language has become an area of debate in the field. Wei et al. (2023) have shown that through a technique called chain-of-thought prompting, the performance of language models on reasoning tasks improves. Chain-of-thought prompting is a method for giving models the ability to solve problems in a step-by-step manner. This experiment shows that models are able to closely mimic the technique by which humans reason, which seemingly requires the ability to “understand” language. Bender and Koller (2020) argue that “the language modeling task, because it only uses form as training data, cannot in principle lead to learning of meaning.” They present a strong case for why we should be precise when describing the current capacities, so as not to add to the current hype of the performance LLMs. Apart from this debate, it is still unknown whether the success of LLM’s is due to relying on word co-occurrence statistics, or whether these models have learned to utilise the deeper linguistic meaning of text.

Al-Negheimish et al. (2021) showed that on the DROP benchmark, a numerical reasoning task, that when the word order of a prompt is shuffled, a LLM’s decline in performance is small. This suggests that models are not arriving to correct answers by reasoning about the meaning of a question. Similarly, Sinha et al. (2021a) showed that models could be pre-trained on sentences with randomly shuffled words and still achieve high accuracy after fine-tuned for downstream tasks. Sinha et al. (2021b) studied the task of natural language inference and found transformer based models achieving high accuracy on the MNLI (Williams et al., 2017) dataset, despite random word-order permutations.

Additionally, through my own testing, I have recently found that this trend is consistent when testing on more recent 7 billion parameter models such as LLAMA-2, Mistral, and MPT. Despite this, it is still unknown whether the models are simply relying on word count to make predictions, our benchmarks aren’t working as intended, or the models can discern the meaning of shuffled prompts.

0.2 Models

This research will focus on models that are currently open-sourced and smaller in size (7 billion parameters). The focus on open-sourced is so the models can be closely studied. The smaller size is due to the availability of resources, as the models are being evaluated using City University’s Hyperion. The smaller models all follow the same transformer-based architecture as the larger 70+ billion parameter models, so while accuracy decreases for the smaller models, the trend of shuffled benchmarks being similar in accuracy to unshuffled benchmarks will be similar. Some models being considered for use are Meta’s Llama2-7b (Touvron et al., 2023), Mistral-7b (Jiang et al., 2023), and Mosaic’s MPT-7b.

0.3 Benchmarks

Benchmarks are one of the ways how the scientific community measures the success of language models. These benchmarks can focus on a models success in any number of areas, from general subject knowledge such as Hendryck’s MMLU (Hendrycks et al., 2020) to common-sense reasoning problems like Hellaswag (Zellers et al., 2019). This strategy for testing models is not without it’s issues and has been highly scrutinized throughout the literature.

Raji et al. (2021) show that the manner in which we use benchmarks, specifically those aimed at testing for “general capabilities” to test models is vague and the framework that benchmarks follow is intended for practical and “tightly-scoped” tasks. McCoy et al. (2023) notes that in order to best

understand the performance of these models, we need to remember the task for which they were designed, which is statistical next-word prediction systems. They evaluate the state-of-the-art GPT 3.5 and GPT 4 models, on tasks such as multiplication, sorting, and ciphers. They show that because of how the models are designed, LLMs often struggle with tasks where the correct outcome is a low probability sequence. For this study common-sense benchmarks are of particular interest, because these tasks should theoretically require syntactic knowledge to properly reason about the problem. Common-sense as a task can be a convoluted and vague concept to define. For instance a task that is common-sense for people living in a particular region, may not be considered common-sense so for a group of people in living in another. However, we can loosely define common-sense tasks as those that draw upon basic human understanding and intuition, encompassing knowledge and reasoning that is typically assumed to be shared among individuals within a specific cultural or societal context. [Davis \(2023\)](#) gives a thorough analysis of benchmarks focusing on common-sense and remarks that although there are over 100 benchmarks focused on this task there are still many areas of common-sense that remain untested.

In light of the weaknesses described above, we must carefully analyze what success on a chosen benchmark actually means. Some notable datasets that are used on Hugging Face’s leaderboard ³ for open-sourced LLM’s are ARC(scientific reasoning), Hellaswag(common-sense), Winogrande(common-sense), MMLU(subject knowledge), TruthfulQA(coomon-sense), and GSM8K(math reasoning). For the research done in this project, benchmarks chosen must seemingly require syntax to “understand” the prompt. For example, I have found that when evaluated on TruthfulQA’s multiple choice version a model performs better ($\approx 3\%$ accuracy increase) when the prompt is shuffled and better still ($\approx 8\%$ accuracy increase) when the prompt is removed all together. This would clearly not be a useful benchmark for this project.

1 Project Objectives

This goal of this project is to provide insight into the process of how a LM interprets the meaning of input text, and document how well our current benchmarks can gauge this. This leads to 3 main research questions we will focus on.

Are LLM’s able to go beyond simple word count? Just because it is shown that models *often* rely on word co-occurrence statistics, does not mean they are incapable of going beyond this on necessary tasks.

Are benchmarks sensitive to higher order information? If models are capable of the above, what benchmarks are adequately reflecting this?

What issues should be considered in the design of future benchmarks? The experiments done in this research can give insight into the current weaknesses many benchmarks have. This information can be used to help in the design of future benchmarks.

2 Project beneficiaries

This project has the potential to help in understanding more about language models and their strengths and weaknesses. It can also help us learn more about our current benchmarks. Many of the benchmarks we use are most likely overestimating the capabilities of LMs and this can potentially be shown through our experiments. It is our hope that through successful completion of this research we can benefit both the NLP and Linguistics research communities.

³[Open LLM leaderboard](#)

3 Work Plan

Below I outline in table 1 the general timeline for how I plan to conduct my project. Because of the nature of this research is based on several future unknowns - like results from experimentation, I have focused most of my planning on a weekly structure.

Monday - Focus on reading papers.
Tuesday - Run experiments (Hyperion)
Wednesday - Read papers (if time)
Thursday - Review experiment results
Friday - Read papers (if time)
Saturday - Run more experiments
Sunday - Write findings

Dates	Task
29/01 - 06/02	Write project proposal
05/02 - 09/02	Literature review. Choose models and benchmarks.
10/02 - 04/03	Evaluate the models using modified lm-evaluation-harness from EleutherAI. Compare the results of shuffled, not shuffled, and other syntactic modification.
05/03 - 25/03	The results from above will dictate the next steps taken.
26/03 - 16/04	Write report
17/04 - 30/4	Final report

Table 2: Project Timeline

4 Project Risks

- Technical/Hardware issues. Technical and hardware issues cover problems that may occur with compute
- Time. Some aspects of this project will undoubtedly take more time than I would plan on. Because of this I am leaving an extra day each week for extra work should the need arise. Some examples of things that could cause this are bugs in my code, new information found, and business with other courses.
- Amount of evidence is not enough to support my hypothesis. While this is completely possible, I will still be able to adapt my project to this by presenting the evidence and seeing what conclusion the facts do support.

Appendix B Reuse Summary

This study was conducted by reusing Eleuther AI’s LM-Evaluation-Harness which can be found here: <https://github.com/EleutherAI/lm-evaluation-harness>

My forked version of the repository can be found here: <https://github.com/tliddell13/lm-evaluation-harness>

The files added by me for the purpose of this study are located in the `lm_eval` folder. `Permutations.py` contains the functions that can be used to carry out changes on the dataset. These changes are then applied based on arguments made with `main.py` and in `evaluator.py`.

The `SlurmEvals` was also added by me. This folder contains Slurm Jobs for evaluating the models with the permutations.