# Are Emergent Abilities in Large Language Models just In-Context Learning?

**Sheng Lu**[1*], **Irina Bigoulaeva**[1*], **Rachneet Sachdeva**[1],
**Harish Tayyar Madabushi**[2], and **Iryna Gurevych**[1]

[1] Ubiquitous Knowledge Processing (UKP) Lab, Technische Universität Darmstadt
[2] Department of Computer Science, The University of Bath

`www.ukp.tu-darmstadt.de`
`htm43@bath.ac.uk`

## Abstract

Large language models have exhibited *emergent* abilities, demonstrating exceptional performance across diverse tasks for which they were not explicitly trained, including those that require complex reasoning abilities. The emergence of such abilities carries profound implications for the future direction of research in NLP, especially as the deployment of such models becomes more prevalent. However, one key challenge is that the evaluation of these abilities is often confounded by competencies that arise in models through alternative prompting techniques, such as in-context learning and instruction following, which also emerge as the models are scaled up. In this study, we provide the first comprehensive examination of these emergent abilities while accounting for various potentially biasing factors that can influence the evaluation of models. We conduct rigorous tests on a set of 18 models, encompassing a parameter range from 60 million to 175 billion parameters, across a comprehensive set of 22 tasks. Through an extensive series of over 1,000 experiments, we provide compelling evidence that emergent abilities can primarily be ascribed to in-context learning. We find no evidence for the emergence of reasoning abilities, thus providing valuable insights into the underlying mechanisms driving the observed abilities and thus alleviating safety concerns regarding their use.[1]

## Contents

---

# 1 Introduction, Motivation and Context

One of the most captivating aspects of pre-trained language models (PLMs) is their remarkable capacity to acquire a wide range of knowledge across different domains, when trained primarily through masked language modelling. Early PLMs like BERT (Devlin et al., 2019) have demonstrated their access to a significant amount of linguistic information, encompassing syntax (Lin et al., 2019; Tenney et al., 2019), semantic roles (Ettinger, 2020), and, to some extent, pragmatic inference, role-based event knowledge (Ettinger, 2020) and world knowledge (Petroni et al., 2019). The diverse abilities of PLMs can be categorised into two broad types: formal linguistic abilities and functional linguistic abilities (Mahowald et al., 2023). The former encompasses the understanding of language rules and patterns, while the latter includes a range of cognitive abilities necessary for real-world language comprehension and utilisation (this distinction is further discussed in Section 2). However, while PLMs appeared to excel in formal linguistic abilities, they have until recently faced challenges in developing functional linguistic abilities (Mahowald et al., 2023) (see also Section 2).

The introduction of Large Language Models (LLMs), which are typically generative PLMs scaled up to billions of parameters and trained on vast amounts (web-scale) of data, is changing this landscape (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023a,b). Recent works indicate that LLMs exhibit *emergent abilities*, as measured by their performance on specific tasks. The emergence of a significant number of abilities in LLMs, including those that explicitly involve some form of reasoning, was first described by Wei et al. (2022b). They defined an *emergent ability* as an ability to solve a task which is absent in smaller models, but present in LLMs. The works of Wei et al. (2022b) and Srivastava et al. (2023), published approximately concurrently, base their definition of emergent abilities on the more general definition of emergence in physics: "Emergence is when quantitative changes in a system result in qualitative changes in behaviour" (Anderson, 1972). As such, emergent abilities in LLMs are characterised by the fact that they typically do not present themselves in smaller models, which makes their emergence unpredictable. This inherent unpredictability associated with emergent abilities, and to a certain degree, the challenge in compre-

hensively explaining their origins, hold substantial implications for the discussion surrounding safety and security when utilising LLMs. The ability of LLMs to perform well above the random baseline on tasks that cannot be solved through memorisation and are indicative of certain "abilities", *without explicit training on those tasks*, is central to the idea of emergence.

## 1.1 Safety and Security Implications

While Wei et al. (2022b) do not explicitly make the distinction between formal and functional linguistic abilities, they observed, through a review of LLM literature, a significant number of functional linguistic abilities to be emergent in LLMs. The emergence of such functional linguistic abilities in LLMs, especially those which cannot be based on mere memorisation, has profound significance, with the potential to shape the future research direction of NLP (Bommasani et al., 2021). If there are indeed multiple abilities that emerge with scale, it suggests that further scaling has the potential to unlock a wide array of additional abilities, including those that have not yet been explored or tested (Wei et al., 2022b).

It has been argued that the unpredictability associated with such emergent abilities underscores the possibility that these models might harbour latent hazardous abilities that manifest unexpectedly. It's important to emphasise that the development of explicit linguistic proficiencies does *not* pose inherent risks of this nature. The same can be said for the capacity to efficiently handle information retrieval tasks. The real concern lies in potential harmful capabilities relating to functional linguistic abilities, especially reasoning and planning (Hoffmann, 2022), which we refer to in this work as latent hazardous abilities. However, it must be emphasised that this does not include other dangers posed through the misuse of these models, such as the use of LLMs to generate fake news.

## 1.2 Emergent Abilities vs Prompting Techniques

The scaling up of LLMs facilitates the acquisition of diverse competencies, which can be generally grouped into two categories: The first group encompasses *abilities*, already described. The second group encompasses various *techniques*, which LLMs can benefit from, but which prove ineffective in smaller models. Among these techniques are in-context learning, instruction tuning, and chain-

of-thought prompting. In-context learning is the technique in which LLMs are provided with a limited number of examples from which they learn how to perform a task (Brown et al., 2020; Liu et al., 2023). Recent investigations into the theoretical underpinnings of in-context learning and its specific manifestation in LLMs indicate that it might bear resemblance to the process of fine-tuning (Akyürek et al., 2022; Dai et al., 2023; Von Oswald et al., 2023; Wei et al., 2023). Another technique exclusive to LLMs is instruction tuning, alternatively known as instructional fine-tuning. This technique involves training LLMs on datasets of instructional content, which enables the models to follow explicit instructions in prompts (Chung et al., 2022; Wei et al., 2022a; Taori et al., 2023) (see also Section 2). Finally, chain-of-thought prompting is the technique in which models are provided with a sequence of intermediate reasoning steps to boost their 'reasoning skills' (Wei et al., 2022c).

To avoid any potential confusion, we adapt the terminology proposed by Wei et al. (2022b) and shall refer to these *techniques* as "prompting techniques". It is noteworthy that a single *prompting technique* can be adaptable across multiple tasks. For example, in-context learning can be used in performing any task through the inclusion of a few illustrative examples within the prompt. We note that this contrasts with the notion of emergent abilities, which are implied to occur due to LLMs' capacity to perform above the random baseline on the corresponding tasks without explicit training on that task. As an example, the emergent ability to understand social situations in LLMs is inferred from LLMs' performing well above the random baseline on the Social IQa (Sap et al., 2019) task. This task serves to evaluate models' emotional and social intelligence and includes questions such as:

> **Question:** Carson was excited to wake up to attend school. Why did he do this?
> **Options:** "Take the big test", "Go to bed early", "Just say hello to friends"

See also Section 2 for more details on tasks.

Significant to our investigation is the observation that these prompting techniques and emergent abilities manifest within LLMs at a comparable scale (See Figure 1). It is crucial to emphasise that "prompting techniques" do not give rise to safety concerns with respect to latent hazardous abilities because prompting techniques provide a predictable outcome that is reliant on the user. These
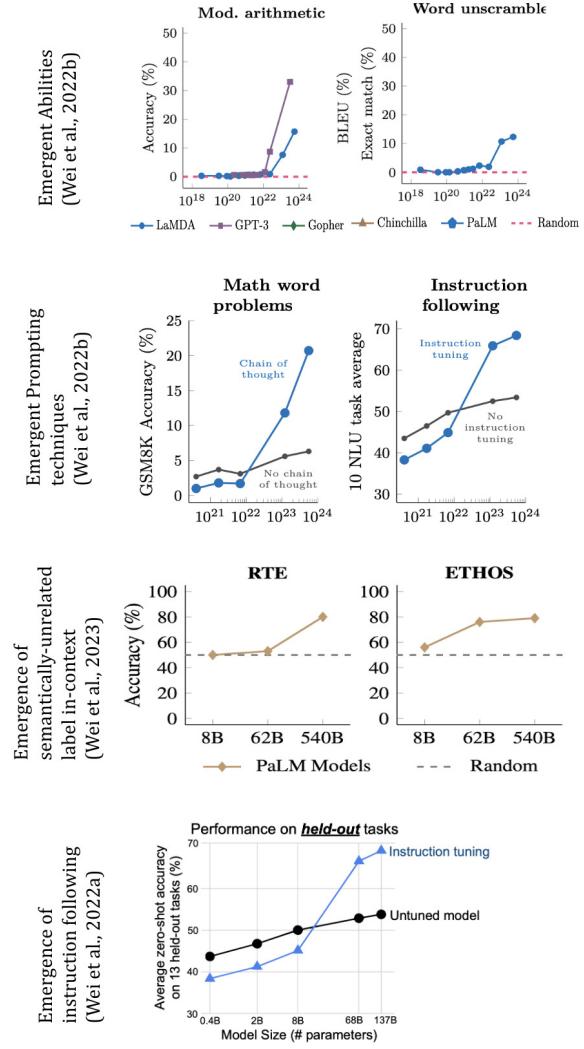


Figure 1: Row 1: The Emergence of various abilities as inferred through the ability of LLMs to 'solve' tasks. Row 2: the jump in effectiveness of prompting techniques such as chain-of-thought and instruction tuning (Wei et al., 2022b). Row 3: the effectiveness on flipped label in-context learning (Wei et al., 2023), and Row 4: the effectiveness on instruction following (Wei et al., 2022a). While the x-axis in Rows 1 and 2 corresponds to the number of training flops, and that in Rows 3 and 4 portray model parameters, it has been shown that there is a strong correlation between these two Wei et al. (2022b). In light of the fact that instruction following, in-context learning and, in particular, semantically-unrelated label in-context learning manifest themselves at a scale comparable to the emergence of emergent abilities, we contend that it is essential to evaluate emergent abilities in the absence of in-context learning and instruction tuning

prompting techniques offer novel avenues for leveraging models, but cannot lead to latent perilous abilities. Hence, this manner of risk is exclusively a consequence of emergent abilities. As previously

mentioned, this threat is not universal among all emergent abilities, but pertains exclusively to those involving reasoning or planning. To be explicit, formal linguistic abilities do not pose any threat, and nor does the ability of models to perform above the random baseline on tasks that can be solved through memory ('*memorisable tasks*') since such proficiency only indicates models' ability to memorise information. Given that prompting techniques manifest themselves at the same scale as the emergence of a significant number of abilities, coupled with the safety implications associated with emergent abilities but not with the prompting techniques, it becomes imperative to ascertain the extent of these emergent abilities in the absence of prompting techniques.

### 1.3 Fine-tuning, In-Context Learning and other Prompting Techniques

Artificial neural models have, for some time, exhibited the ability to achieve significant success on specific tasks when trained on those tasks (Devlin et al., 2019; Liu et al., 2019). PLMs in particular have demonstrated this even in the few-shot setting (Hofer et al., 2018; Radford et al., 2019; Brown et al., 2020; Gao et al., 2020). Such performance after training on a task is not considered "emergent", precisely because models are trained on that very task. Indeed, the fact that LLMs are *not* trained on the tasks used in evaluating their emergent abilities is central to identifying abilities which are truly *emergent*. The assertion that achieving satisfactory performance on a given task signifies the *emergence* of associated 'abilities' hinges on the condition that models are not explicitly trained for that specific task. Alternatively, it would suggest that a model possesses the expressive power required to undergo training for the said task. Only if an LLM has not been trained on a task that it performed well on can the claim be made that the model *inherently* possesses the ability necessary for that task. Otherwise, the ability must be learned, i.e. through explicit training or in-context learning, in which case it is no longer an ability of the model per se, and is no longer unpredictable. In other words, the ability is not emergent.

The examination of the extent to which information pertaining to specific attributes – such as parts-of-speech (PoS) – is encoded within the representations of words and sentences in a model, commonly involves the use of classifiers. These classifiers, called *probes*, are trained to learn a mapping between said representations and the corresponding attribute information (Conneau et al., 2018; Hupkes and Zuidema, 2018), exemplified here by PoS details. Considerable efforts have been dedicated to investigating the attributes of classifier that are used as probes, including, for example, the complexity of the probing classifiers (the simpler the better), whether or not the model's parameters are updated during the training of the probing classifier (typically not), and the application of information-theoretic metrics that remain impartial to the specific classifier used (Voita and Titov, 2020). In the context of generative models, however, such classifiers can be done away with, and instead, prompts can be employed to probe models to discern their access to various forms of information. While significant research has gone into exploring the characteristics of classifier based probes, a similar analysis of the prompts used to probe LLMs is lacking.

The recent insights suggesting parallels between in-context learning and explicit training imply that the success on a task through in-context learning, much like models trained explicitly for task-solving, does not inherently imply a model possessing that *ability* (Dai et al., 2023) (see also Section 2). This bears resemblance to attempting to probe for specific inherent abilities using particular tasks following training on those same tasks. Therefore, it is crucial to conduct an independent evaluation of LLMs' abilities, detached from in-context learning. Crucially, a notable proportion of instruction tuning datasets encompass instances where instructions are linked to exemplar outputs (refer to Section 2). This characteristic underscores the possibility that instruction tuning potentially triggers in-context learning in LLMs, which, if true, would imply that the success of a model to solve a task in this scenario does not indicate the emergence of the corresponding ability. Conversely, although unlikely, due to the nature of instruction tuning described above, there remains the possibility that instruction tuning might inadvertently facilitate reasoning capabilities in LLMs.

### 1.4 Research Questions

Based on the observations above, our research seeks to address the following questions:

1. Given the potential influence of in-context learning on the purported 'emergence' of vari-

ous abilities in LLMs, which abilities are truly emergent in the absence of in-context learning, including instructional tuning?

2. Considering that instruction tuning datasets typically include mappings between instructions and exemplars, is there evidence for the emergence of 'reasoning' in instruction-tuned models? Alternatively, in line with Occam's Razor, can we assert that a simpler explanation is that instruction tuning enables these models to more efficiently and effectively leverage in-context learning?

It is also possibly the case that chain-of-thought prompting is a prompting strategy that provides an effective way of using in-context learning, especially in scenarios where tasks demand multi-step reasoning for inference. However, we do not perform experiments to test this possibility and leave this exploration for future work (see Section 7).

### 1.5 Contributions

We focus our efforts on providing a controlled, yet extensive evaluation of emergent abilities across models at various scales by controlling for various prompting techniques, specifically in-context learning and instruction tuning. By employing these controls, we aim to shed light on the degree to which the various abilities of LLMs are inherently emergent and to distinguish such emergence from the *appearance* of emergence that might arise due to the increased degree to which prompting techniques manifest themselves. Concretely, this work makes the following significant contributions to the understanding of what has thus far been discerned as emergent abilities in LLMs:

1. Recognising the parallels between in-context learning and training, we advocate for the probing of languages models using strategies that do not trigger in-context learning.

2. For the first time, we evaluate the abilities of LLMs independent of in-context learning and instruction tuning, providing a clear and precise measure of the abilities which are truly emergent.

3. We empirically investigate the hypothesis that the added capabilities of instruction-tuned models can be explained as their efficient use

of in-context capabilities. By applying Occam's Razor, we contend that our straightforward explanation obviates the need to attribute such capabilities to the emergence of reasoning.

4. We provide an explanation for the abilities exhibited by LLMs: the combination of formal linguistic abilities, the capacity to retain and recall large amounts of information and, significantly, in-context learning.

5. Our findings indicate that there are no emergent functional linguistic abilities in the absence of in-context learning, affirming the safety of utilising LLMs and negating any potential hazardous latent capabilities.

## 2 Related Work

This section explores literature that is of significance to our study. We begin by elaborating on the distinction between formal and functional linguistic abilities, previously introduced in Section 1. We provide further details of emergent abilities and then examine in-context learning, including recent findings that attempt to shed light on the theoretical foundations that make it possible. We then explore instruction tuning, including a description of the datasets employed for this purpose. Lastly, we highlight the singular study that has so far questioned the existence of emergent abilities in LLMs, with an emphasis on the aspects that set our work apart.

The existence of formal linguistic abilities in PLMs has been well established for some time now (Rogers et al., 2020). In fact, it has been shown that pre-training on as few as between 10 and 100 million tokens is sufficient for PLMs to capture a significant amount of formal linguistic information (Zhang et al., 2021). It has been argued that such linguistic information might essentially be theoretical constructs, a notion that gains traction especially given that their existence does not precede the use of language (Tayyar Madabushi et al., 2023). Given that it is possible that linguistic elements are statistical generalisations that arise out of language use, representing intricate patterns that are necessitated by various goals tied to efficient communication, it is not surprising that LLMs have access to such information when trained on large corpora. Even more unsurprising is the fact that LLMs, when trained more extensively, also capture

memorisable information, being able to, for example, perform on par with knowledge-bases (Petroni et al., 2019). These observations, coupled with the recognition that the formal linguistic competencies and the ability of models to perform above the random baseline on memorisable tasks do not pose the risk of latent hazardous abilities, thereby ensuring their safety, leads us to differentiate tasks that hinge on formal linguistic abilities and memory from functional linguistic abilities which demand more advanced processes, such as reasoning.

In this context, the work by Wei et al. (2022b) provided significant evidence for unpredictable emergence of several *functional linguistic abilities*. This evidence stems from their review of prior literature of LLMs including GPT-3, PaLM (Chowdhery et al., 2022), Chinchilla (Hoffmann et al., 2022), Gopher (Rae et al., 2021) and LaMDA (Thoppilan et al., 2022). Wei et al. (2022b) identified a total of 67 emergent abilities based on above-random performance of LLMs on tasks designed to test those abilities from the Beyond the Imitation Game benchmark, BIG-bench (Srivastava et al., 2023), a crowd-sourced evaluation benchmark consisting of over 200 tasks, and a further 51 of the 57 tasks from the Massive Multitask Language Understanding Benchmark (Hendrycks et al., 2020). However, Wei et al. (2022b), unlike this work, do not make the distinction between formal linguistic and functional linguistic abilities and also do not attempt to identify tasks that can be solved predominantly using memory and so do not require reasoning. Importantly, they do not consider the impact of in-context learning or other prompting techniques, instead listing them as alternative means of interacting with LLMs. Our work marks the first exploration of the emergent abilities of LLMs independent of various prompting techniques.

In-context learning is the ability of LLMs to perform a task with only a minimal set of exemplars presented within the context of the input prompt (Brown et al., 2020; Dong et al., 2022; Liu et al., 2023). While this ability of LLMs has been known for some time (Kojima et al., 2022; Srivastava et al., 2022), recent work has shown that LLMs are capable of in-context learning even in cases where labels are flipped or semantically unrelated, as in the case of flipped labels for sentiment analysis (Wei et al., 2023). Crucially, it should be noted that the capacity to excel in a flipped labelling task, such as attributing a negative sentiment to

sentences labelled as positive and vice versa, inherently relies on in-context learning. Without this, the model would be unable to adapt its sentence labelling accordingly.

One plausible theoretical rationale for this phenomenon is furnished by Dai et al. (2023), which indicates that in-context learning in LLMs might share similarities with fine-tuning, in that it might allow models to 'learn' from the examples presented in their prompt. Similarly, it has been shown that in-context learning implements gradient descent implicitly and constructs a function at inference time on regression problems (Akyürek et al., 2022; Li et al., 2023; Zhang et al., 2023a), which may be related to gradient-based meta-learning (Von Oswald et al., 2023). Another line of work shows that in-context learning is driven by the distributions of the pre-training data (Chan et al., 2022; Hahn and Goyal, 2023). Other theoretical explorations of why in-context learning works include work by Xie et al. (2022) and Zhang et al. (2023b), who attempt to explain this capability in terms of Bayesian inference, that by Li et al. (2023) who propose a PAC based framework for explaining this capability and, Jiang (2023)[2], Wang et al. (2023), and Wies et al. (2023) who provide an explanation based on latent space theory. Pertinent to this study, though, is the overarching observation in these theoretical explorations that in-context learning becomes more powerful with increased scale, and so is predictable and thus not emergent or indicative of the possibility of the emergence of other latent hazardous abilities. As mentioned earlier, however, the specific mechanisms governing in-context learning do not impact our argument: The fact of its functionality suffices to underscore the necessity of assessing emergent abilities in the absence of in-context learning. This becomes particularly relevant considering that the capacity of LLMs to execute in-context learning with flipped labels (Wei et al., 2023) becomes evident at a similar scale to the emergence of abilities (Wei et al., 2022b) (see also Section 1.2 and Figure 1).

In exploring the effectiveness of training LLMs on datasets of instructions, Wei et al. (2022a) cluster tasks based on the type of problem being addressed (e.g., natural language inference, sentiment . . . ) and train models to learn to follow instructions on clusters of tasks. They then eval-

---

[2]It should be noted that Jiang (2023) refer to in-context learning and instruction tuning as emergent abilities. This is in stark contrast to our definition of emergence.

**Template 1**

<premise>

Based on the paragraph above, can we conclude that <hypothesis>?

<options>

**Template 2**

<premise>

Can we infer the following?

<hypothesis>

<options>

**Template 3**

Read the following and determine if the hypothesis can be inferred from the premise:

Premise: <premise>

Hypothesis: <hypothesis>
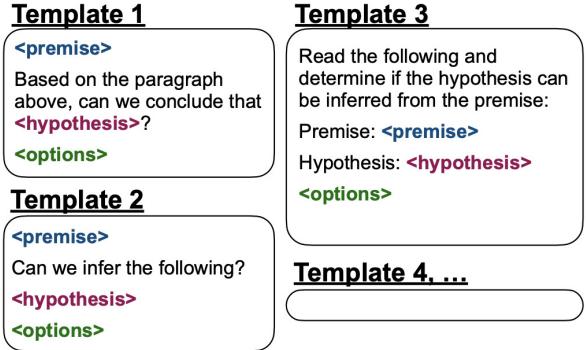
<options>

**Template 4, …**

Figure 2: The process of instruction tuning. Source Wei et al. (2022a).

uate the performance of modules thus fine-tuned (or instruction-tuned) on previously-unseen clusters of tasks and find that the ability to follow such previously-unseen tasks is emergent in LLMs. Figure 2 illustrates the process of instruction tuning, specifically highlighting how this approach trains models to translate instructions into exemplars as required by in-context learning.

Prior comparison of instruction-tuned and non-instruction-tuned models, for example by Chung et al. (2022), who introduce the Flan instruction tuning dataset, have shown that instruction tuning is powerful in improving downstream performance. However, to the best of our knowledge, no previous work has indicated the need to consider the possibility of instruction tuning triggering in-context learning in evaluating inherent abilities of LLMs, emergent or otherwise.

Prior studies have also explored other abilities emergent in LLMs, such as the examination of Theory of Mind (Kosinski, 2023). Similarly, Biderman et al. (2023), Tefnik and Kadlčík (2023), Wu et al. (2023), and Zheng et al. (2023) have explored the extent to which memory plays a role in LLMs' abilities. By using counterfactual tasks, Zheng et al. (2023) find that performance degrades in adversarial settings, thus indicating that, although LLMs exhibit some level of proficiency in abstract task-solving, frequently lean towards employing context-specific approaches. However, to the best of our knowledge, none of the previous evaluations of emergent abilities have been conducted in a manner that explicitly distinguished between the in-context learning and instruction-tuning settings and prompting in the setting wherein these abilities are not triggered (i.e., zero-shot evaluation of non-instruction-tuned models).

Finally, prior to this work, the only study to question the existence of emergent abilities has been the work by Schaeffer et al. (2023), who argued that emergence is likely to be a consequence of the discrete evaluation metrics commonly employed for assessing LLMs. Typically, the performance of LLMs is evaluated using F1 scores based on exact match or, in the case of multiple-choice questions, that of the most probable option (as defined by, for example, the minimal perplexity of the model). Schaeffer et al. (2023) argue that this approach may lead to the appearance of unpredictable emergence of capabilities, when none exist. This implies that the performance of LLMs does improve progressively with scale, but, when using discrete evaluation metrics, such improvements are only detectable when they tip over a threshold, thus giving the illusion of 'emergence'. They demonstrate that the use of continuous metrics, such as string edit distance, makes emergent abilities more predictable, thereby challenging their emergent status. Recall that the predictability of what LLMs are capable of makes them safe, as this negates the possibility that latent hazardous abilities might also be emergent. Importantly, however, some (Wei et al., 2022b) argue against this by pointing out that there are tasks on which LLMs are able to perform well above the random baseline where smaller models can only perform below the random baseline, suggesting that these abilities are still emergent and not just a consequence of discrete evaluation metrics. In this work, we investigate emergent abilities from a markedly distinct standpoint – that of in-context learning. Moreover, we account for a range of potential biases that could potentially impact our research, including those associated with evaluation metrics (see Section 3.3).

## 3 Methods

To answer the research questions presented in Section 1.4, we examine each of the following: 1) The impact of removing the factor of *in-context learning* by working in the zero-shot setting and removing the factor of *instruction tuning* by working with non-instruction-tuned models. 2) The interplay between in-context learning and instruction tuning to determine whether 'reasoning' or similar explanations are necessary to explain the additional capabilities of instruction-tuned models. We achieved this by comparing the capabilities of models in the absence of instruction-tuning and

models at different scales instruction-tuned to different extents. 3) If each of the tasks typically used to evaluate models assesses functional linguistic abilities, formal linguistic abilities, and whether the task can be memorised, which we evaluated through a manually examining each task.

In order to assess the true extent of emergent abilities in LLMs, we meticulously structure our experiments to minimise any factors that might create the appearance of emergence. Our experimental design specifically focuses on assessing models in a manner that does not trigger in-context learning. It is important to note that instruction tuning, which entails training models to translate instructions into exemplars (refer to Section 2), could potentially trigger in-context learning. Therefore, we take precautions to account for the possibility that instruction tuning triggers in-context learning by using non-instruction-tuned models.

### 3.1 Models

We perform our evaluation on a range of models of different sizes (parameter counts) from four model families: GPT, T5, Falcon and LLaMA. We choose these model families as GPT and LLaMA have previously been found to have emergent abilities, and Falcon is at the top of LLM leader-boards at the time of writing. Finally, we select T5 as it is an encoder-decoder model and its instruction-tuned version (Flan) is trained using an extensive instruction tuning dataset. In the GPT family, we use instruction-tuned and non-instruction-tuned versions of both GPT-2 and GPT-3 (Brown et al., 2020), and in the T5 family, we use T5 (Raffel et al., 2020) and its instruction-tuned counterpart FLAN-T5 (Chung et al., 2022). Similarly, we use both the instruction-tuned and non-instruction-tuned versions of Falcon[3]. In the case of LLaMA (Touvron et al., 2023a), which is not instruction-tuned, we were unable to access the instruction-tuned version of the model. In addition, we also evaluate GPT-3 text-davinci-003, an InstructGPT model. InstructGPT models are initially fine-tuned on annotator authored prompts and corresponding desired behaviours. This model is then used to collect a dataset of ranked model outputs and is further fine-tuned using reinforcement learning from human feedback (RLHF). This training regime has been shown to improve model performance (Ouyang et al., 2022). Table 1 enumerates the models that

---

[3]https://huggingface.co/tiiuae/falcon-40b

| Family | Model | Instruction-tuned | Size |
|---|---|---|---|
| GPT | GPT-2 | GPT-2-IT* | 117M |
| | GPT-2-XL | GPT-2-XL-IT* | 1.6B |
| | GPT-J | GPT-JT | 6.7B |
| | davinci | text-davinci-001 *text-davinci-003* | 175B |
| T5 | T5-small | FLAN-T5-small | 60M |
| | T5-large | FLAN-T5-large | 770M |
| Falcon | Falcon-7B | Falcon-7B-IT | 7B |
| | Falcon-40B | Falcon-40B-IT | 40B |
| LLaMA | LLaMA-7B | – | 7B |
| | LLaMA-13B | – | 13B |
| | LLaMA-30B | – | 35B |

Table 1: Details of the models we use in our experiments. InstrucTuned refers to the instruction fine-tuned version of the model. *text-davinci-003* is the only model additionally trained using RLHF.

we use in our experiments.

Where instruction-tuned versions of a model were not available (as in the case of GPT-2 and GPT-2-XL and identified in Table 1 with a star) we train the corresponding models using the FLAN (Wei et al., 2022a) instructional dataset. The T5 models we choose are intentionally below the 1B parameter count as emergent abilities have *not* been observed in models as small as this (Wei et al., 2022b) and serve as crucial control in our experiments. Of the models that we select, GPT-3 'davinci' (non-instruction-tuned), GPT-3 text-davinci-001 (instruction-tuned) and GPT-3 text-davinci-003 (InstructGPT) are the models of the scale at which emergent abilities have been previously observed. This choice is primarily a practical one due to the availability of models. Other model families that have demonstrated emergent abilities include PaLM (Chowdhery et al., 2022), Chinchilla (Hoffmann et al., 2022), Gopher (Rae et al., 2021) and LaMDA (Thoppilan et al., 2022) which we do not evaluate due to the unavailability of corresponding APIs.

### 3.2 Tasks

A complete list of the tasks we use in our experiments is presented in Table 2. Since GPT-3 is the model that possesses the ability to demonstrate emergent abilities among the models we select for experimentation, our choice of tasks includes those tasks which were found to be emergent in GPT-3. Of the 17 BIG-Bench tasks recognised as emergent in GPT-3 (Wei et al., 2022b), we select 14. Three tasks previously identified as emergent are excluded from our analysis. These tasks are: Question-Answer Generation from COPA, which

involves generating questions, thereby complicating automated evaluation; Self Evaluation of Tutoring, also challenging to assess automatically due to its generative nature; and IPA Transliteration, omitted due to its reliance on BLEU score for output similarity evaluation, lacking an equivalent exact match metric, thus making our experiments inconsistent. Additionally, we randomly select a limited number of tasks (7) from the same dataset that were *not* found to be emergent in GPT-3, serving as a control group for comparison. Finally, we also include GSM8K (Cobbe et al., 2021), which comprises a set of grade-school mathematics word problems and is noteworthy because GPT-4 (OpenAI, 2023) attains an impressive accuracy of 94% on this task, while GPT-3.5 achieves only 57.2% accuracy. The 22 tasks thus selected provide a range that enables us to evaluate the emergent abilities of language models across various categories: a) tasks that are considered emergent for GPT-3, b) those considered not to be emergent despite LLMs' ability to solve them as this is predictable based on the success of smaller models, and c) tasks on which LLMs completely fail. GSM8K, and one of the emergent tasks from BIG-Bench, Modified Arithmetic, as originally structured, do not adopt a multiple-choice format. However, for alignment with our other experiments, we adapt them into a multiple-choice setup by including numeric distractors as follows: To ensure adversarial nature, we automatically select numbers in proximity to and randomly distant from the correct answer.

Table 3 provides a description and an illustrative example of three of the twenty two tasks utilised in our experiments: one task demonstrating emergence, one task that has not exhibited emergence, and one task previously identified as emergent solely in GPT-4 (GSM8K). A similar set of descriptions and corresponding examples of all 22 of the tasks that we use in our experiments is presented in Appendix A, Table 8.

## 3.3 Controls for Possible Bias

So as to ensure that our evaluation is fair, we identify potential biases that could influence our findings and design our experiments to mitigate such biases. In cases where this is not possible, we shape our experiments to maximise our chances of identifying such emergent abilities, if they do indeed exist.

### 3.3.1 Prompt Formats

Table 4 shows an example of each of our prompt formats. We make two significant changes to the prompting strategies used: First we refine all prompts to ensure their solvability even in the absence of instruction comprehension. We call this adjusted prompt format the *completion-style prompt*, and use it for all models (See Table 4. We experiment with minor variations to these prompts so as to find the most optimal format.

This change is necessary, since in order to assess the true abilities of non-instruction-tuned models in the zero-shot setting, it is imperative to evaluate their ability to accurately perform tasks without relying on explicit instructions. As outlined in Section 2, many tasks involve prompts that inherently require an understanding of explicit instructions. Since LLMs in their base form are trained to perform next-word prediction, it is unreasonable to expect that without instruction tuning, they will respond adequately to multiple choice question prompts requiring them to pick the correct answer from a set of options. We hypothesised that using such a prompt style would give an unequal advantage to the instruction-tuned models. Indeed, our initial prompt experiments demonstrated that non-instruction-tuned models merely try to "complete" the text of the prompt by generating additional answer choices, sometimes even additional new questions. However, once the prompt itself was adjusted to take the form of a sentence to be completed, non-instruction-tuned models were likelier to output one of the answer choices. We confirm that these changes do not skew our results by replicating prior results using instruction-tuned models, which we use as a baseline.

The second change we make to our prompting strategy involves the exploration of two types of completion-style prompts: *closed* and *open*. In the closed prompt format, we provide answer choices alongside the prompt, while in the open prompt format, the answer choices are withheld. We find that when models are prompted using the open prompt strategy, their generated results often exhibit little or no resemblance to the provided answer choices. Consequently, evaluating the correctness of the generated answers becomes challenging. As a result, experiments utilising the open prompt setting are completely excluded from our analysis. However, we provide access to these responses in the data accompanying this study, allowing other researchers

| Task | Prev. Emergent | Competence type | Memorisable (of 50) |
|---|---|---|---|
| Causal judgement | No | Functional | 0/50 |
| English Proverbs | No | Functional | 0/50 |
| Implicatures | No | Functional | 0/50 |
| Nonsense words grammar | No | Formal | 38/50 |
| Rhyming | No | Formal | 50/50 |
| Tracking shuffled objects | No | Functional | 0/50 |
| Commonsense QA (Talmor et al., 2019) | No | Functional | 3/50 |
| GSM8K | No | Functional | 0/50 |
| Analytic entailment | Yes | Functional | 4/50 |
| Codenames | Yes | Functional | 0/50 |
| Common morpheme | Yes | Formal | 0/50 |
| Fact checker | Yes | Functional | 50/50 |
| Figure of speech detection | Yes | Functional | 0/50 |
| Hindu knowledge | Yes | Functional | 50/50 |
| Logical deduction | Yes | Functional | 0/50 |
| Misconceptions | Yes | Functional | 50/50 |
| Modified arithmetic | Yes | Functional | 0/50 |
| Phrase relatedness | Yes | Functional | 50/50 |
| Physical intuition | Yes | Functional | 50/50 |
| Social IQa (Sap et al., 2019) | Yes | Functional | 0/50 |
| Strange stories | Yes | Functional | 0/50 |
| Strategy QA (Geva et al., 2021) | Yes | Functional | 27/50 |

Table 2: A list of the tasks used in our experiments, along with their previous identification as emergent or otherwise, accompanied by a categorisation of the nature of the requisite ability for solving the task. This classification is determined through a manual inspection of the data, employing the categorisation framework provided by (Mahowald et al., 2023). We evaluate memorisability of 50 examples from each task and we assume no data leakage of the task data. For an example-based justification of our classification of 'memorisable' examples, see Table 9 in the Appendix.

| Task name | Description | Example |
|---|---|---|
| Tracking shuffled objects | This task tests a model's ability to work out the final state of a system given its initial state and a sequence of modifications. | **Input:** Alice, Bob, and Claire are playing a game. At the start of the game, they are each holding a ball: Alice has a orange ball, Bob has a white ball, and Claire has a blue ball. As the game progresses, pairs of players trade balls. First, Alice and Bob swap balls. Then, Bob and Claire swap balls. Finally, Alice and Bob swap balls. At the end of the game, Alice has the **Options:** "orange ball", "white ball", "blue ball" **Target:** blue ball. |
| Logical deduction | This task requires deducing the order of a sequence of objects from a minimal set of conditions. | **Input:** On a shelf, there are three books: a black book, an orange book, and a blue book. The blue book is to the right of the orange book. The orange book is to the right of the black book. **Options:** "The black book is the leftmost", "The orange book is the leftmost", "The blue book is the leftmost" **Target:** The black book is the leftmost. |
| GSM8K | The dataset supports the task of question answering on basic mathematical problems that require multi-step reasoning. | **Input:** Weng earns \$12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How much did she earn? **Options:** 13, 17, 10, 8, 25 **Target:** 10 |

Table 3: Descriptions and examples from one task not found to be emergent (Tracking Shuffeled Objects), one task previously found to be emergent (Logical Deductions), and one task found to be emergent only in GPT-4 (GSM8K). A similar list of all 22 of the tasks that we use in our experiments is presented in Appendix A, Table 8.

to experiment with it.

Furthermore, we introduce an additional adversarial prompt in the closed format, where the correct answer choice is indicated by a corresponding letter rather than the text itself. This adversarial style of prompting provides us with two advantages:

a) it requires the model to rely on its instructional understanding, as it needs to associate the letter with the appropriate choice, and b) it allows models to output answers with a single letter, simplifying the process of generating an answer, allowing for the possibility that they possess the knowledge to

| Prompt format | Example |
| --- | --- |
| default, closed | **Question:** Austin's family was celebrating their parents 50th anniversary during dinner at a new restaurant. What would Austin's family do next? From the following choices, choose the correct answer: "Refuse to eat dinner with the family", "Eat dinner at the restaurant", "Happy"<br>**Answer:** |
| completion, closed | Austin's family was celebrating their parents 50th anniversary during dinner at a new restaurant. What would Austin's family do next? The possible answers are "Refuse to eat dinner with the family", "Happy", "Eat dinner at the restaurant", but the correct answer is |
| adversarial, closed | **Question:** Austin's family was celebrating their parents 50th anniversary during dinner at a new restaurant. What would Austin's family do next?<br>**Options:** (a) "Refuse to eat dinner with the family", (b) "Eat dinner at the restaurant", (c) "Happy"<br>**Answer:** |

Table 4: Sample prompts of the three formats we use. The samples are from the Social IQa task of BIG-Bench.

answer the question, but encounter difficulties in its generation.

### 3.3.2 Evaluation Metrics

The typical method of evaluating model performance on tasks which are multiple choice is to match the model output to the correct answer option. However, when evaluating non-instruction-tuned models using the completion-prompts, we must account for the possibility that the outputs generated do not match the provided answer choices exactly. Therefore, we additionally evaluate models based on the semantic similarity between the output text and the provided answer choices using BERTScore (Zhang et al., 2020) to estimate the model's answer choice. In this setting, the answer is considered correct if the generated answer is most similar (semantic text similarity) to the correct answer choice, and incorrect if it is closer to any of the others. We call this metric *BERTScore accuracy*. It's worth noting that this process is akin to selecting the answer where the model has exhibited lowest perplexity. Since calculating this perplexity for models that are exclusively

accessible through APIs is not practical, we adopt this alternative metric. We opt for BERTScore over alternatives like BLEURT (Sellam et al., 2020) because the latter are additionally trained to assess the fluency of the output text, a factor which is not our focus, and one that renders them computationally resource-intensive.

Some tasks, however, require models to generate numerical outputs (refer to Section 3.2). In these instances, we limit our evaluation to exact matching, as measuring semantic similarity between numbers does not accurately reflect their numerical proximity. Although it is possible to employ a metric that considers the proximity between numbers (e.g., the numeric difference between the prediction and the correct output), we refrain from doing so due to the significant variation in magnitude of the numbers across different questions in the same task, which would complicate the evaluation process. Similarly, the task *'codenames'*, is a task which requires the model to output a word from a given list of words and consists of questions such as: "input: Try to identify the 1 word best associated with the word FOURTH from the following list: pad, wedding, quarter, brain, wish, halloween, einstein, helmet, sun, tip, laundry, judge. Give your answer in alphabetical order. target: quarter". Given the nature of this task, we also evaluated it only using exact match.

Notice that BERTScore accuracy is a discrete evaluation metric and does *not* constitute the use of a continuous metric. Given that recent work has indicated that emergence might be a result of the discrete evaluation metrics (See Section 2), we also include string edit distance and continuous-BERTScore as additional evaluation metrics. It's important to note that the application of the continuous evaluation metric serves the purpose of establishing whether the performance of larger models is genuinely unpredictable. However, our investigation reveals that even when employing a discrete metric, especially in the absence of in-context learning, there is no sudden performance increase in large models. As such, we do not include results using continuous metrics in our analysis. However, we do include them in the dataset accompanying this study, offering the potential for further exploration.

### 3.3.3 Manual Evaluation of Responses

To further mitigate the possibility that models which are not instruction-tuned are not unfairly

penalised, we perform a manual analysis of their responses. This is required, due to the possibility that these models generate the correct answer, albeit in a format that is not semantically most similar to the option provided. To account for this possibility, we preform a manual post-hoc analysis of a subset of the outputs generated by models which are not instruction-tuned.

## 3.4 Experimental Setup

| Family | Model | Tasks |
|--------|-------|-------|
| GPT | GPT-2<br>GPT-2-IT<br>GPT-2-XL<br>GPT-2-XL-IT<br>GPT-J<br>GPT-JT<br>davinci<br>text-davinci-001<br>*text-davinci-003* | All 22 Tasks |
| T5 | T5-small<br>FLAN-T5-small<br>T5-large<br>FLAN-T5-large | |
| Falcon | Falcon-7B<br>Falcon-7B-IT<br>Falcon-40B<br>Falcon-40B-IT | Logical Deductions, Social IQA, GSM8K, Tracking Shuffled Objects |
| LLaMA | LLaMA-7B<br>LLaMA-13B<br>LLaMA-30B | |

Table 5: An overview of the experimental setup. Models in the GPT and T5 families are evaluated on all tasks and those in the Falcon and LLaMA families on a subset of representative tasks. In addition, each evaluation is performed in the closed and closed adversarial prompting strategies.

The overall experimental setup, including the different models tested, the different tasks used in our experiments and the evaluation settings employed, is detailed in Table 5. Given our objective of evaluating the emergent abilities of LLMs independent of other factors, we perform the following experiments. We evaluate each of the 12 models selected from the T5 and GPT families (Section 3.1) on all of the 22 selected tasks as described in Section 3.2. For each case, we employ the prompting strategies: closed, and closed adversarial, as discussed in Section 3.3.1. To consider the variability in responses, we conduct each experiment three times and calculate the average result. All experiments that we run locally are run on NVIDIA

A100 GPUs using a temperature of 0.01 and a batch size of 16. In the case of GPT-3 175B parameter models (davinci, text-davinci-001, and text-davinci-003), we make use of the official API for evaluation which is done once using a temperature of 0. While we restrict our evaluation to a single run due to cost constraints, it's improbable that this will impact the results of our experiments. This is because we also set the temperature to 0, which guarantees result reproducibility and minimises the possibility of hallucinations.

In addition, we evaluate six selected models from the LLaMA and Falcon families (See Section 3.1), on four of the 22 tasks chosen earlier. We pick these four tasks ensuring that two have been previously identified as emergent (Logical Deductions and Social IQA) and the other two have been determined to be non-emergent (GSM8K and Tracking Shuffled Objects). Once again we test these using the closed and adversarial prompting strategies and run each experiment thrice to account for variance. Given the variable number of options associated with some of the tasks under evaluation, we construct the baseline for each task by randomly selecting options for questions in that task multiple times and finding an average score.

## 4 Results: Emergent Abilities in the absence of In-context Learning

The detailed results of all the experiments conducted are provided in the Appendices. In this and the next section, we highlight a subset of the results that are particularly noteworthy based on our observations and analysis. These selected results aim to highlight the key findings and trends from our experiments. When assessing models which are not instruction-tuned, we only consider the BERTScore accuracy metric, which is flexible, and requires only that the generated answer has the highest semantic similarity to the correct option. Notice that this is essential as models that are not instruction-tuned might generate the correct answer, but in ways which do not align with the answer choices (See Section 3.3.2). This heightened flexibility, combined with the alterations we apply to the prompts (Section 3.3.1), and the manual evaluation we perform, ensures that our evaluation of non-instruction-tuned models is unbiased. The three exceptions to the use of BERTScore accuracy are the two numeric tasks and 'codenames', for which we employ an exact match metric, given the

nature of the tasks (Section 3.3.2).

## 4.1 Overview of Results on Emergent Abilities

The first set of results and corresponding analysis presented in this section are aimed at answering the first of our research questions: Given the potential influence of in-context learning on the purported 'emergence' of various abilities in LLMs, which abilities are truly emergent in the absence of in-context learning, including instructional tuning?

We begin by presenting the performance of non-instruction-tuned 175B parameter GPT-3 models in the zero-shot setting, which, as we have previously argued, is a more accurate approach of evaluating models' inherent abilities (Section 4.2). We then list the potential emergent abilities and examine each within the framework of functional and formal linguistic requirements of the corresponding tasks, as well as the tasks' memorisability (Section 4.3). To verify that these results generalise, we repeat the analysis on results obtained on the other model families that we experiment with: T5, LLaMA and Falcon (Section 4.4). Finally, to mitigate any potential biases that could influence our conclusions, we perform a manual assessment of a subset of model responses. Additionally, we explore the potential for emergence in the closed adversarial setting, wherein the model is tasked with producing a single output character, allowing for evaluation through conventional metrics (Section 4.5. These findings are summarised in Section 4.6.

## 4.2 Performance in the absence In-Context Learning

Figure 3 illustrates the performance of models from the GPT family on the tasks chosen for evaluation. This figure represents results obtained using the closed prompting strategy. Tasks listed in the first two rows, against a grey background, are tasks which have not been found to be emergent by Wei et al. (2022b), while the rest of the tasks are those which have been found to be emergent in prior literature (Wei et al., 2022b). It should be noted that the instruction-tuned GPT models exhibit performance that could lead to the appearance of emergence in the exact match, few-shot setting, on all of those tasks that have previously been found to be emergent, as illustrated in Figure 3. This outcome enables us to validate our experimental configuration and establish a baseline for comparison. These models also exhibit the *lack* of emergence on tasks which have previously been found not to

be emergent: while there is a performance over the random baseline on a few of the non-emergent tasks (e.g., english proverbs), these are not considered truly emergent, as this increased performance is predictable based on the performance of smaller models. This outcome, which aligns with previous results, serves as a baseline. It indicates that the modifications made to the prompts to ensure that non-instruction-tuned models are not disadvantaged – specially their conversion to 'completion style prompts' does not hinder the potential for detecting emergent abilities, and does not disadvantage instruction-tuned models.

Furthermore, the exact match accuracy (depicted in blue in Figure 3) tends to be consistently lower compared to the BERTScore accuracy (depicted in yellow). This is once again in line with expectations, since BERTScore accuracy considers the semantic similarly between the model's output and the answer options, and selects the option most similar to the output generated by the model. As a result, BERTScore accuracy is more forgiving than the exact match metric, allowing for partial matches or closely related responses to be considered as correct. The striking similarity in accuracies between the two metrics, exact match and BERTScore, particularly in the few-shot prompting strategy is worth noting. This indicates that in this setting, models tend to generate output that is exactly one of the options provided.

Additionally, we note that instruction-tuned models, even those trained using program code and reinforcement learning with human feedback (text-davinci-003), consistently perform worse in the zero-shot setting compared to text-davinci-001 in the few-shot setting (See Complete Results presented in Appendix C). This indicates that the instruction-tuned models, which are explicitly trained to translate instructions to exemplars, benefit from the use of in-context learning both indirectly (through instruction tuning) and also directly (in the few-shot setting).

## 4.3 Emergent Abilities

The most intriguing findings lie within the non-instruction-tuned GPT models in the zero-shot setting. Recall that in this setting we evaluate models using BERTScore accuracy. These results, depicted by the green line in Figure 3, stand out as they shed light on the inherent abilities of GPT models. Possible emergence in this setting is summarised in
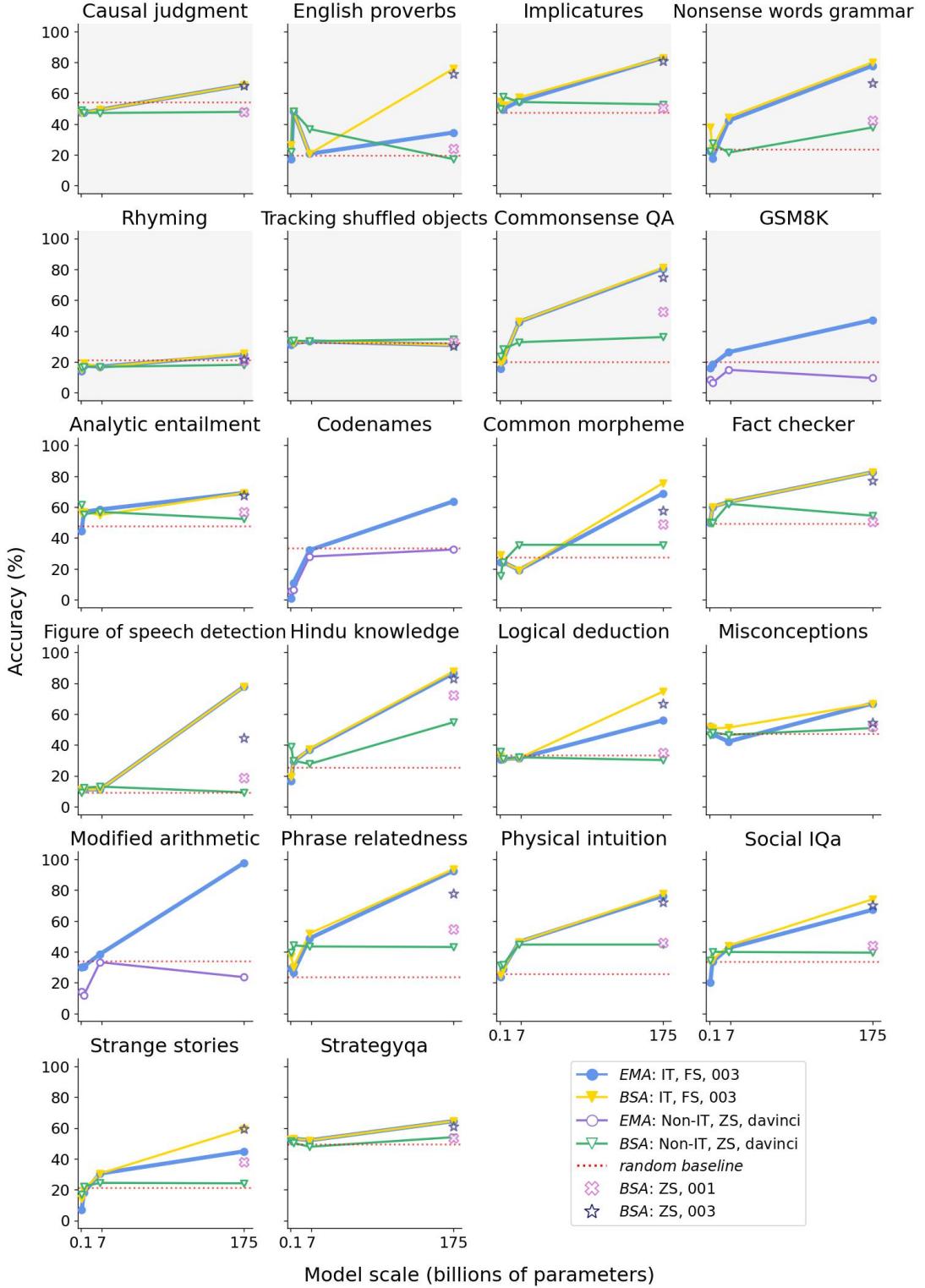
Figure 3: Performance of models from the GPT family on the various tasks used for evaluating emergence tested on the closed prompting strategy. Models which are instruction-tuned (IT) and those which are non-instruction-tuned (non-IT) are evaluated in the few-shot (FS), zero-shot (ZS) settings using BERTScore accuracy (BSA) and the Exact Match Accuracy (EMA). Blue represents the results in the FS setting on IT models and is comparable to results previously reported in literature. Yellow indicates the performance in the same setting measured using BSA and red represents the BSA results in the ZS setting of models which are non-IT, the setting that represents results in the absence of in-context learning. 001 and 003 represent text-davinci-001 and 003 respectively.

Table 6.

| Task | > Base. | Pred. | Emg. |
|------|---------|-------|------|
| Causal judgement | No | N/A | No |
| English Proverbs | No | N/A | No |
| Rhyming | No | N/A | No |
| GSM8K | No | N/A | No |
| Codenames | No | N/A | No |
| Figure of speech detection | No | N/A | No |
| Logical deduction | No | N/A | No |
| Modified arithmetic | No | N/A | No |
| Tracking shuffled objects | No* | N/A | No |
| Implicatures | Yes | Yes | No |
| Commonsense QA | Yes | Yes | No |
| Analytic entailment | Yes | Yes | No |
| Common morpheme | Yes | Yes | No |
| Fact checker | Yes | Yes | No |
| Phrase relatedness | Yes | Yes | No |
| Physical intuition | Yes | Yes | No |
| Social IQa | Yes | Yes | No |
| Strange stories | Yes | Yes | No |
| Misconceptions | Yes* | No | Yes* |
| Strategy QA | Yes* | No | Yes* |
| Nonsense words grammar | Yes | No | Yes |
| Hindu knowledge | Yes | No | Yes |

Table 6: Performance of the non-instruction-tuned 175B parameter GPT-3 model (davinci) in the zero-shot setting, which we propose as the setting to evaluate tasks in the absence of in-context learning. For a task to be considered emergent (Emg.), models must perform above the baseline (> Base.) and the performance of the larger models must not be predictable based on that of smaller models (Pred.). Results marked with a star indicate that they are not significant.

Recall that the definition of emergence, as stated by Wei et al. (2022b), requires LLMs to perform a task above the baseline *and* do so in a manner that cannot be predicted based on the performance of smaller models. Despite the criticisms of the use of a discrete evaluation metric (Section 2), we deliberately use BERTScore accuracy, a discrete metric. This choice is driven by our desire to provide every opportunity to identify emergent abilities.

Furthermore, our investigation reveals that tasks on which smaller models surpass the baseline cannot be emergent, as we observe very limited improvements on such tasks using larger models. Conversely, tasks that hold the potential for emergence are those where smaller models fall below the baseline. Consequently, the choice of metric becomes inconsequential, as this pattern implies the rate of performance improvements is immaterial. An analysis of Figure 3, presented in Table 6 indicates that just two tasks are 'emergent', when we account for in-context learning. While two additional tasks (Misconceptions and Strategy QA) also have unpredictable above baseline performance, the improvement is only marginal, on tasks which are binary classification tasks with a random baseline of 50% accuracy. Among the two identified tasks, Nonsense Words Grammar pertains to a formal linguistic ability which we've noted that does not involve reasoning. Likewise, the other emergent task, Hindu knowledge, solely relies on recall and lacks any reasoning demands. Additionally, the description of these tasks on BigBench also does not mention the requirement of reasoning. **As such, our results indicate that reasoning abilities are *not* emergent in LLMs.**

## 4.4 Other Model Families

We extend our analysis to the LLaMA, Falcon, and T5 model families, employing a similar methodology. Across each of these cases, a consistent pattern emerges: task performance is either predictable based on smaller model performance or the performance is below the baseline on the largest models within our scope.

Figure 4 illustrates our findings in the Falcon and LLaMa model families, and Figure 6 (Section 5) illustrates our findings on the T5 model family. We do not employ the largest available models in any of these three families, a decision based on their unavailability through APIs and the inherent impracticality of running them on standard hardware. Furthermore, it's worth noting that previous research has not identified these tasks as emergent within the context of these model families.

## 4.5 Possible Biasing Factors

Since our findings rely on use of LLMs that have not been instruction-tuned, we verify that the observed lower performance on tasks does not stem from biasing factors such as their inability to comprehend the tasks or the prompting methods used. This section presents our evaluation of possible biases. Notice that this evaluation is in addition to the controls we have already put in place, namely the modification of prompts and the use of BERTScore accuracy (described previously in Section 3.3).

To begin with, the above baseline performance of the non-instruction-tuned version of GPT-3 (davinci) lends support to the notion that such models can indeed comprehend task requirements. Additionally, we conduct a qualitative analysis, which further underscores the models' ability to understand task requirements even when they are not instruction-tuned. For instance, in the case of the task 'causal judgement' the model gener-
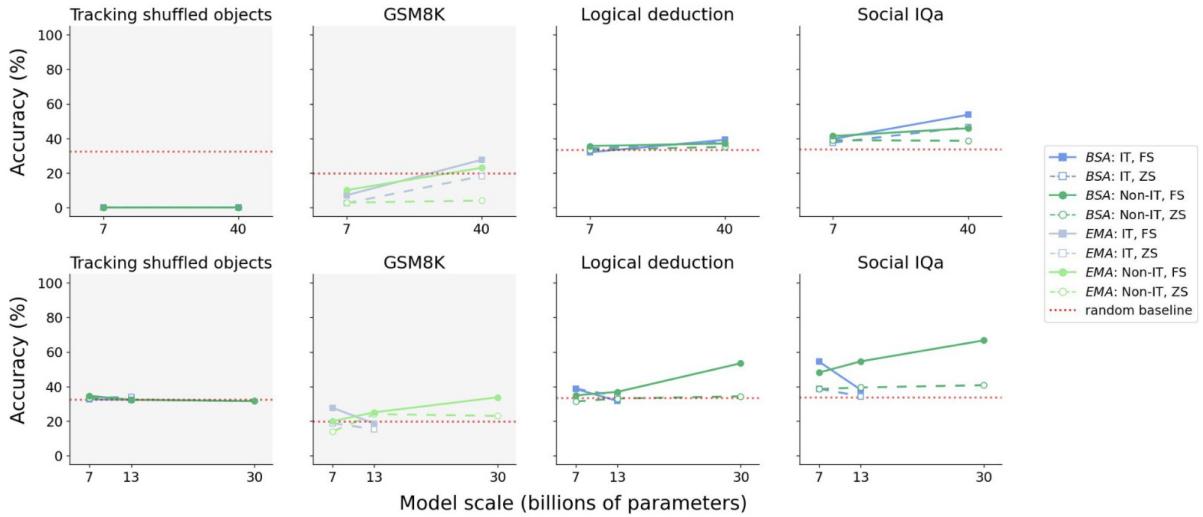
Figure 4: An illustration of the performance of models in the Falcon (top) and LLaMA (bottom) families in the non-instruction-tuned zero shot setting on a the selected subset of tasks, demonstrating the consistent lack of emergent abilities in the absence of in-context learning.

ates responses limited to "yes" or "no," as required, demonstrating its grasp of the task's requirements. To further ensure that this is not a biasing factor in our assessment of models, we additionally perform a manual analysis (Section 4.5.1).

### 4.5.1 Manual Evaluation of Responses

We explore the possibility that evaluating models which are not instruction-tuned might be challenging as metrics including BERTScore accuracy might not provide a reliable assessment in these cases if the generated answer might not align well with the target options. To ensuring that our results are not biased in this manner, we carefully examine the output of davinci (the 175B GPT-3 non-instruction-tuned model) to assess the efficacy of automatic evaluation used. Using BERTScore accuracy, we conduct a manual analysis of 50 output examples from each task, the results of which is presented in Table 7. Recall that modified arithmetic, GSM8K, and codenames are always evaluated using exact match accuracy and so are not included in this analysis.

In Table 7, 'BERTScore accuracy %' represents the % of correct answers as determined by the automatic metric of the 50 examples selected for manual evaluation and 'Manual Eval Accuracy %' represents the % of correct answers as determined by a manual analysis of the results by one of the authors on the same set of examples. Recall that the purpose of this exercise is to ensure that the automatic evaluation metric does not affect our conclusion in

terms of the existence of emergent abilities. Our analysis shows that, in the majority of cases, the automatic metric overestimates model performance. This set of tasks is represented in the top block in Table 7.

In the case of Logical deduction, the model sometimes produces answers that are copied from the question but are still technically correct answers, which could lead to the MA% score being too lenient. In the case of Casual judgement, the increase is only slight compared to the above 50% baseline. These two cases wherein the manual evalution indicates a higher scores are represented in Table 7 block 2. Finally, on 'Commonsense QA', the only task where there is a significant increase over the baseline, such performance is predictable based on the performance of smaller models, and so the task is not emergent.

This analysis of 50 examples from each task carries a degree of imprecision. Crucially, however, it is imperative to recognise that our primary objective is to ensure that these inaccuracies, inherent to the automatic evaluation of generative models, do not fundamentally alter our conclusions. Our analysis underscores that this is the case and that these limitations do not undermine the validity of our findings.

### 4.5.2 Shortened output Generation

LLMs lacking instruction tuning often exhibit a degree of proficiency in adhering to instructions, albeit within constrained limits, particularly in the

| Task | BSA% | MA% | Base% |
|------|------|-----|-------|
| Analytic entailment | 48 | 14 | 48 |
| Common morpheme | 32 | 22 | 27 |
| English proverbs | 10 | 6 | 20 |
| Fact checker | 52 | 34 | 49 |
| Figure of speech detection | 10 | 10 | 9 |
| Hindu knowledge | 52 | 54 | 25 |
| Implicatures | 58 | 6 | 48 |
| Misconceptions | 48 | 40 | 47 |
| Nonsense words grammar | 34 | 22 | 24 |
| Phrase relatedness | 44 | 34 | 24 |
| Physical intuition | 46 | 40 | 26 |
| Rhyming | 16 | 6 | 21 |
| Social IQa | 36 | 38 | 34 |
| Strategy QA | 58 | 58 | 49 |
| Tracking shuffled objects | 34 | 20 | 32 |
| Strange stories | 34 | 28 | 21 |
| Logical deduction | 26 | 34 | 33 |
| Causal judgement | 46 | 56 | 54 |
| Commonsense QA | 36 | 54 | 20 |

Table 7: A comparison of BERTScore Accuracy (BSA%) and a manual evaluation (MA%) on 50 examples from each task. The analysis reveals that in instances of notable disparity, BERTScore accuracy generally tends to result in false positives (top block). In exactly three cases BERTScore accuracy underestimates performance: in two instances the increase allows model performance to increase above the baseline only marginally. In the case of Logical deduction, the model sometimes produces answers that are copied from the question but are still technically correct answers, which could lead to the MA% score being too lenient. In the case of Causal judgement, the increase is only slight compared to the above 50% baseline. Where there is a substantial performance boost above the baseline (bottom block), this particular task's predictability based on smaller model performance implies that it remains not emergent. As such, we find that even a lenient manual scoring does not affect our conclusion.

context of models with a substantial parameter count of 175B (Wei et al., 2022a). We leverage this phenomenon by using the "adversarial prompt setting" (Section 3.3.1), wherein the model is required to generate output choices, such as options "a" or "b," instead of the target choice. In this setting we evaluate models using a relaxed version of exact match wherein an answer is marked correct if it contains the correct target option. This flexibility is once again designed to allow us to detect any possible indication of emergence. Note that this assessment allows us to circumvent the necessity for employing less precise evaluation criteria as is required when evaluating more verbose responses. The results of this evaluation on the seven of 22 tasks wherein the performance is above the random baseline are presented in Figure 5.

Of these seven tasks on which the non-instruction-tuned version of GPT-3 performs above the random baseline, three are predictable based on the performance of smaller models and thus not considered emergent. The only task on which the improvement over the baseline is not predictable and significant is 'physical intuition.' This task includes questions such as "The bonds in sodium chloride are of what type? Options: Ionic: 1, Covalent: 0, Metallic: 0, Hydrogen: 0", which are likely to be more memory based. Common morpheme, on the other hand, is a non-trivial task that require significant reasoning abilities. However, we find that it has an extremely small test set with only fifty examples and thus the improvement in accuracy is only a small fraction of the total.

## 4.6 Interim Summary and Key Insights

This section presented our findings derived from the evaluation of LLMs in the absence of in-context learning and instruction tuning. In investigating the possible existence of emergent ability, we make several accommodations to guarantee a thorough exploration of this possibility. In this regard, we introduced the following modifications to our methodology to safeguard against inadvertently introducing bias into our experiments, particularly in a manner that might disadvantage LLMs which are not instruction-tuned:

1. modify prompts to make them "completion-style prompts",

2. use BERTScore accuracy (Section 3.3.2) to make our evaluation of models more flexible,

3. perform a manual analysis to ensure that the automatic evaluation we use does not undermine our results,

4. use an alternate style of prompting which requires models to output only the correct option (e.g., a, b, . . . ), thereby enabling the use of the more precise exact match accuracy, and

5. continue to use a discrete evaluation metric despite some reservations of such metrics.

It's important to underscore that the above modifications are designed to ensure that we can exhaustively uncover any potential emergence, even in cases where the likelihood is minimal. Despite this, we find no evidence for the emergence of functional linguistic abilities that would indicate the possible emergence of latent hazardous abilities.
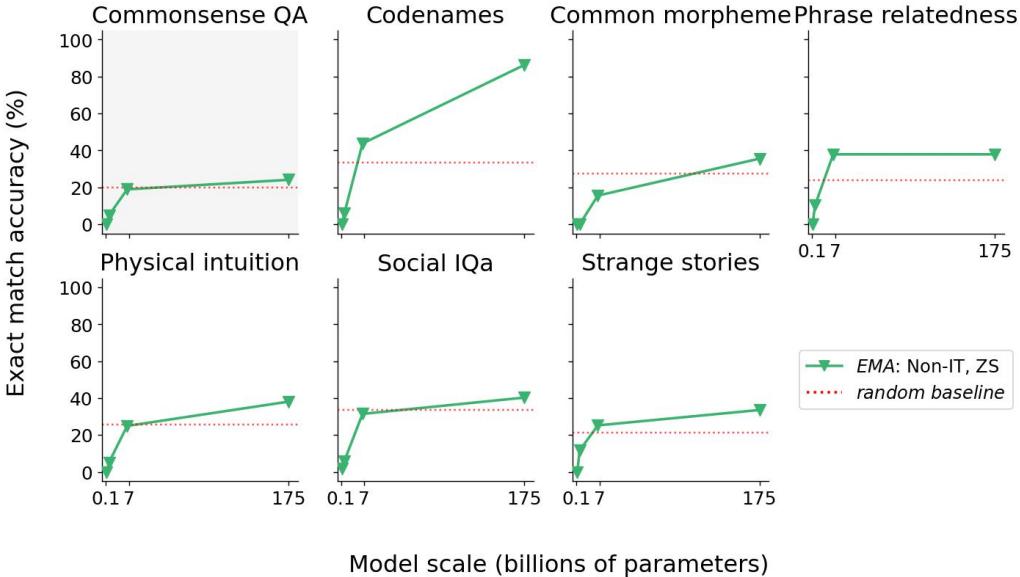
Figure 5: Performance of non-instruction-tuned GPT models in the adversarial setting (required output is a, b, . . . ) on the subset of tasks wherein the performance is above the random baseline. The performance on 'codenames', 'phrase relatedness', and 'strange stories' is predictable and so not emergent. Across the remaining tasks, the improvements in performance compared to the random baseline are relatively modest. Additionally, of the tasks on which the performance gain is slightly more notable, we find that 'physical intuition' is a memory intensive task and 'common morpheme' has a small test set (see text).

## 5 Results: Instruction Tuning as Efficient In-Context Learning

The second set of results and the corresponding analysis that we present in this section serve to substantiate our second hypothesis (Section 1.5): Do instruction-tuned models exhibit 'reasoning,' or is it more likely that instruction tuning allows them to more effectively and efficiently access their in-context learning? It is important to observe that the nature of instruction tuning, which involves data that establishes a mapping between instructions and exemplars, which are characteristic of in-context learning, makes it likely that this process indeed triggers in-context learning. Our hypothesis would imply that instruction tuning provides LLMs with the capability to translate instructions into exemplars, which subsequently engage their in-context capabilities. This would effectively consolidate both functions into a single step.

### 5.1 Overview of Results Relating In-Context Learning and Instruction Tuning

In evaluating which is more likely, the triggering of in-context learning or an alternative explanation such as the emergence of reasoning, we make use of three sets of results. The first is the performance of T5 models, which may not have a scale suffi-

cient for enabling explicit in-context learning. This provides us a way of testing if instruction tuning is effective when explicit in-context learning isn't. Our results, described in Section 5.2 and illustrated in Figure 6 show that this is the case.

Secondly, we compare the tasks that models which are not instruction tuned but which are capable of in-context learning can 'solve' (perform well above the random baseline on) with those that instruction-tuned models too small to possess explicit in-context learning capabilities can. We contend that if instruction tuning is fundamentally different (e.g., give rise to reasoning in models), then it is unlikely that there will be substantial overlap between the set of tasks solvable with instruction-tuned models and those by in-context learning. Detailed in Section 5.3, these results show that despite selecting models from different model families (GPT vs T5), there is a significant overlap in the tasks that models which are instruction tuned and not capable of explicit in-context learning can solve and those which models which are not instruction tuned but capable of in-context learning can. This allows us to conclude that the simpler explanation (i.e., that instruction tuning triggers in-context learning) is the more likely one.

Third, it could potentially be argued that adding in additional factors – namely, RLHF and code

18

training – might endow very large LLMs with the necessary reasoning abilities to solve tasks. However, our results cast doubt upon this claim as well. These results, presented in Section 5.4 and illustrated in Figure 8, show that the tasks solvable by FLAN-T5 770M and GPT-3 175B (text-davinci-003), the latter of which that has been trained on code and with RLHF, are nearly the same. In cases where GPT-3 does outperform T5 – such as in Codenames and Hindu Knowledge – it must also be taken into consideration that the model is 200x larger than T5, and that some of these tasks are memorisable, as shown in Table 2.

In all, our findings point toward a simpler explanation for the exceptional performance of LLMs on various tasks, rather than, for example, the emergence of reasoning abilities. This explanation centres on the notion that these models possess an improved capacity to utilise their inherent in-context learning abilities, which is acquired through instruction tuning. This explanation is particularly plausible due to the nature of instruction-tuning datasets, which teach models to map instructions to examples (refer to Section 2).

## 5.2 The Absence of explicit In-Context Capabilities in T5-Large

We test this hypothesis using the T5 family of models. Our choice of T5 models, of which the largest (T5-Large) has 770M parameters, enables us to evaluate models at a scale where instruction tuning proves effective. Our experiments involving T5-Large also show that there is no difference between the zero-shot and few-shot settings. This suggests that the model's scale is insufficient to support explicit in-context learning effectively.

Figure 6 illustrates the performance of the T5 models in the closed and closed adversarial settings. To recall, we have previously established the absence of emergent abilities in the non-instruction-tuned variant of T5 across all tasks. Central to our current inquiry is the observation that the few-shot setting does not yield performance enhancements in either the instruction-tuned or non-instruction-tuned versions. This contrasts starkly with our findings within the GPT family, where the few-shot setting of even the instruction-tuned version of the largest GPT variant (175 billion parameters) consistently outperformed the zero-shot setting. These findings illustrate that within the T5 models selected by us, model performance directly hinges on

the model's ability to follow instructions, and not on in-context learning, possibly because, at this relatively small scale, in-context learning is not strong enough to manifest itself directly and does so only through instruction tuning.

## 5.3 Comparative Analysis of Initial Tasks Addressable via Instruction Tuning and In-Context Learning

As we scale up our models progressively, the number of tasks they can solve (i.e., perform above the random baseline on), increases. We aim to draw a comparison between the tasks that the smallest non-instruction tuned model within our experimental range, which possesses effective in-context learning capabilities (i.e., GPT 6.7B), can successfully address in the few-shot setting, and those that can be tackled by instruction-tuned models, which are smaller and thus lacking effective direct in-context capabilities (i.e., Flan-T5 770M), in the zero-shot setting. Notice that our choice ensure that the the model we use to test in-context learning is not instruction tuned and the one one used to evaluate instruction following cannot explicitly access in-context learning. If instruction tuning leads to the emergence of something fundamentally different from in-context learning, including, for example, reasoning capabilities in models, this would result in no significant overlap in the set of tasks solvable solely through instruction tuning and the set of tasks addressable via in-context learning.

Additionally, we expand the comparative scope of our analysis to encompass the performance of the instruction-tuned T5 version (Flan-T5) in the zero-shot adversarial setting. This aids us in investigating whether the effectiveness of instruction following contributes to the observed success. Note that in the zero-shot adversarial setting, accurately selecting the appropriate option (e.g., 'a' or 'b') requires the ability to follow instructions, a competence that the non-instruction-tuned version of T5 with its 770 million parameters has yet to fully acquire through pre-training alone (as it is not instruction-tuned). While our prior experiments relied on non-instruction-tuned models to handle the task of choosing the correct option letter, this was possible due to their significantly larger scale at 175 billion parameters. This comparison is presented in Figure 7.

We find that there is a remarkable and substantial overlap in the tasks where these models perform
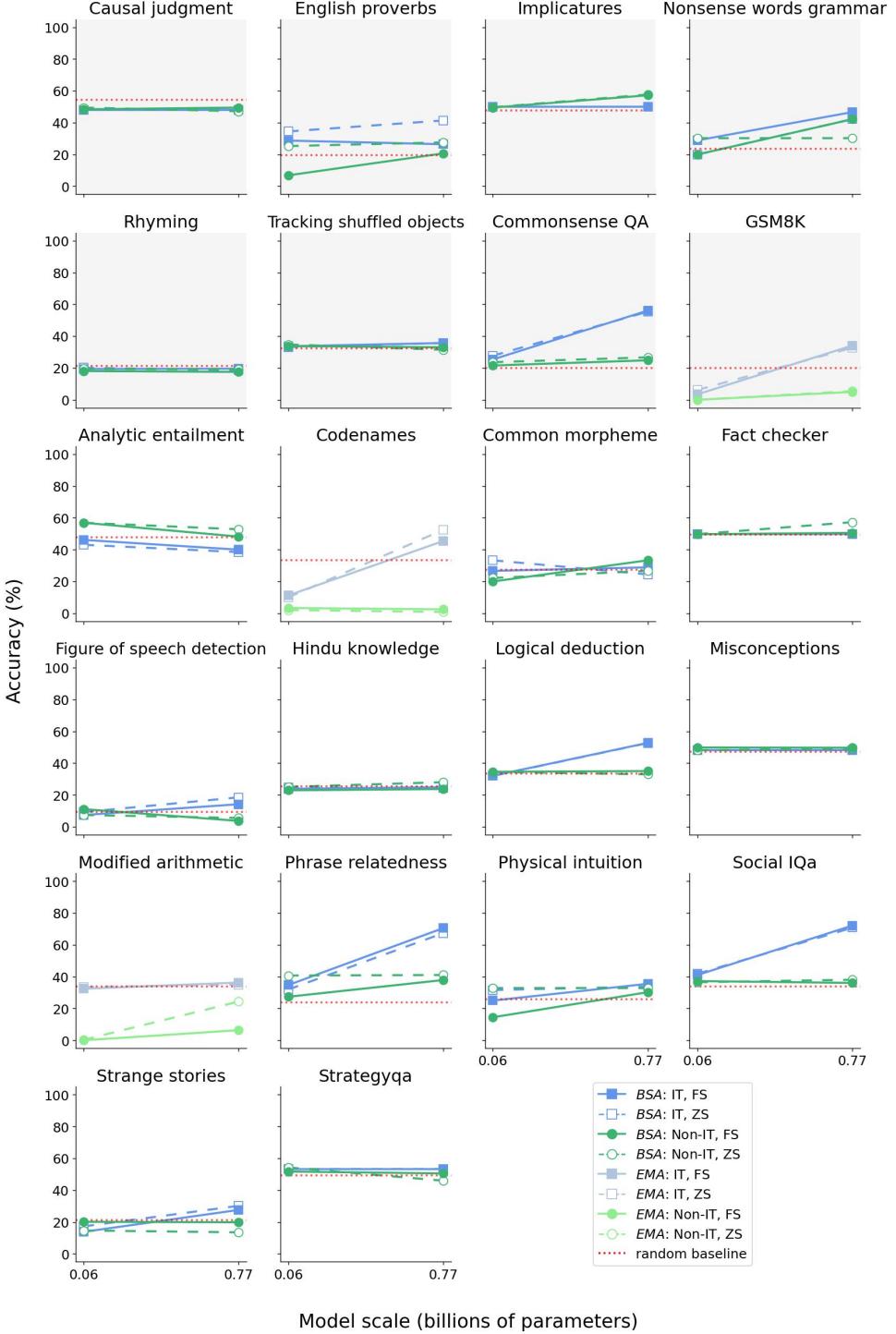
Figure 6: The performance of the T5 family of models in various settings. BSA: BERTScore accuracy, IT: Instruction-Tuned, and EMA: Exact Match Accuracy. See text for analysis.

above the random baseline. The overlap observed between Flan-T5's performance in the closed and closed adversarial settings reaffirms our prior assertion that the capacity to comprehend and adhere to instructions is pivotal to T5's prowess in tackling tasks of this magnitude.

In Section 5.2, we established that Flan-T5 can

follow instructions, but, the model's in-context capabilities are, at this scale, insufficient to be effectively utilised through the conventional few-shot approach. The non-instruction-tuned version of GPT-3 6.7B, on the other hand, lacks the advantage of instruction tuning, and we test it in the few-shot setting, where it can leverage in-context

Figure 7: A comparison of the few-shot (FS) tasks on which Flan-T5 performs above the baseline with the closed prompt, those on which the instruction tuned (IT) Flan-T5 performs above the baseline with the adversarial prompt, and those on which the non-instruction-tuned (Non-IT) version of GPT 6.7B performs above baseline in the zero-shot (ZS) setting. The significant overlap of the tasks is noteworthy and indicates that instruction tuning allows for the effective access (at smaller scale) of in-context capabilities rather than leading to the emergence of reasoning skills. See text for details.

learning. Note that the substantive dissimilarity between these two models – T5 being an encoder-decoder model and GPT being a decoder model – is further compounded by their distinct pre-training datasets. Despite these fundamental differences, there is a significant overlap in both the tasks where they exhibit above-baseline performance and the extent to which the instruction tuned models outperform the non-instruction tuned ones. This overlap in the results underscores a compelling argument: it is more likely that instruction tuning serves as a mechanism that enables models to harness in-context capabilities more effectively, rather than the models having emergent reasoning skills. Indeed, some of the cases where there is no overlap are as expected: For example, in the case of 'hindu knowledge', which is a recall-based task, the tasks provide the larger GPT model with an advantage.

## 5.4 Comparison of Instruction Following in the absence of direct In-context Learning

Having examined the tasks that models can tackle successfully (i.e., performing above baseline), we proceed to compare the tasks that can be effectively tackled by Flan-T5 with those by the different versions of GPT-3. It is important to note that GPT-3 has more than *200 times* the number of parameters present in T5, and no prior research has reported indications of emergent abilities within either Flan-

T5 or T5. Additionally, we perform this comparison in the zero-shot setting, thus allowing us to compare the instruction following capabilities of these models in the absence of direct in-context learning. We present this comparison in Figure 8.

This comparison allows us to answer the following questions: a) Does increased scale significantly impact the tasks on which models can perform above the random baseline, and b) Does enhanced instruction tuning, including the incorporation of program code as seen in text-davinci-003, provide an advantage in being able to perform above the baseline on tasks? By limiting ourselves to the zero-shot setting, we ensure that our results are not affected by in-context capabilities, which we know to increase significantly with scale. Our results indicate that neither a) nor b) is the case. There is a significant overlap in the tasks on which Flan-T5 performs above the baseline and those on which text-davinci-003 does. Not only is there an overlap on tasks, but, in several cases, the performance of these significantly different models in these settings is indeed comparable. This indicates that the extent to which instruction tuning boosts performance is comparable across model scale and instruction tuning datasets.

The significant exception to this are the tasks that are recall-based (e.g., 'hindu knowledge'), on which the larger GPT model performs well above
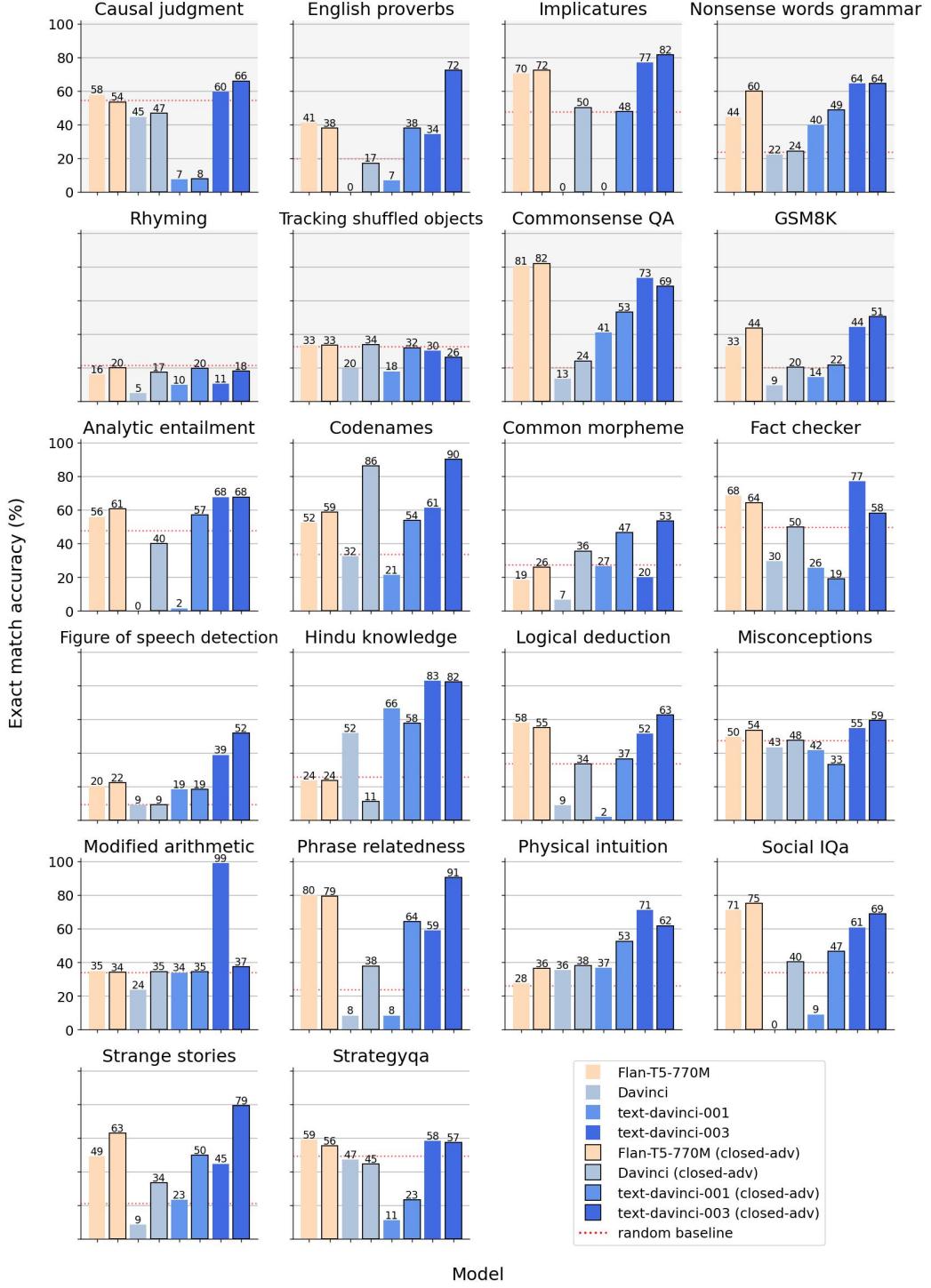
Figure 8: The figure showcases the significant overlap in tasks where both Flan-T5 and GPT-3 perform above baseline, demonstrating comparable performance despite their substantial model architecture differences. The comparison emphasises the effectiveness of instruction tuning across diverse model scales and instruction datasets within a zero-shot setting.

the baseline, while the much smaller T5 model performs close to or below the baseline. This is in line with our previous observations that larger models have higher recall capabilities.

## 5.5 Interim Summary and Key Insights

The set of experiments and analysis presented in this section provide us with important insights into the workings of LLMs. In particular, we can conclude that a) it is more likely that instruction tuning allows LLMs to use in-context learning more efficiently than the LLMs possessing reasoning abilities, and b) when the influence of direct in-context learning is removed, the effectiveness of instruction tuning does not change with either increased scale or with significantly more training on instructional datasets, especially program code. As such, we find that there is no significant indication of the possibility of latent hazardous abilities.

### 5.5.1 Why are some Tasks not solvable?

Given that our results indicate that the ability of LLMs to perform a variety of tasks is a consequence of their ability to follow instructions, and by extension in-context learning, it becomes crucial to examine the existence of tasks on which LLMs cannot perform above the baseline. If in-context learning, which resembles fine-tuning, is being activated, one might assume that all tasks should be solvable at least to some extant. We attribute this limitation to the fact that in-context abilities, while available to models of all sizes, increase with scale and we postulate that the extent of in-context abilities required to solve problems differs based on the task and its formulation (as with fine-tuning). Further experiments, involving models at different scale tested on problems of known complexity, are required to validate this hypothesis and we leave this for future exploration. While further experimental investigation is necessary to validate this, we hypothesise that the increase in the number of emergent abilities in LLMs with scale can be ascribed to some tasks requiring increased in-context abilities, analogous to the need for additional layers or parameters in traditional machine learning models to solve certain tasks. We leave this investigation to future work.

## 6 Implications and Broader Context

We started with two hypotheses: a) that the emergence of nearly all functional linguistic abilities that has previously been observed is a consequence of in-context learning, and b) that the ability of LLMs to follow instructions when instruction-tuned is more likely to be indicative of instruction tuning allowing for the more efficient use of in-context learning rather than leading to the emergence of reasoning skills. Results presented in Section 4 confirmed that there are indeed no emergent abilities in the absence of in-context learning. Similarly, results presented in Section 4.3 confirmed our second hypothesis.

Although 14 of our chosen tasks were previously considered emergent, when controlling for in-context learning, both directly and as triggered through instruction tuning, we found that only two tasks displayed emergence, one which indicates formal linguistic abilities (nonsense words grammar) and the other which indicates recall (hindu knowledge). Ultimately, this absence of evidence for emergence represents a significant step towards instilling trust in language models and leveraging their abilities with confidence, as it is indicative of the complete lack of latent hazardous abilities in LLMs, in addition to being controllable by the user. By contributing to a deeper understanding of these models' behaviour and limitations, we help to demystify the LLMs in the public and remove the related safety concerns.

The distinction between the ability to follow instructions and the inherent ability to solve a problem is a subtle but important one. Simple following of instructions without applying reasoning abilities produces output that is consistent with the instructions, but might not make sense on a logical or commonsense basis. This is reflected in the well-known phenomenon of hallucination, in which an LLM produces fluent, but factually incorrect output (Bang et al., 2023; Shen et al., 2023; Thorp, 2023). The ability to follow instructions does not imply having reasoning abilities, and more importantly, it does not imply the possibility of latent hazardous abilities that could be dangerous (Hoffmann, 2022).

Attributing these capabilities to a combination of memory and in-context learning and the more general ability of these models to generate the most statistically likely next token, can help explain the abilities and the behaviour of LLMs. These capabilities, when subsequently combined with instruction tuning, adoption to conversational use cases, increased context length and a degree of safety controls through the use of reinforcement learning through human feedback, allow LLMs

to become truly powerful. Our experiments shed some initial light on the internal mechanisms of LLMs and start to unravel how these models are capable of performing numerous tasks. For example, despite their zero-shot capabilities, in the absence of explicit examples (or in the case where the instructions weren't explicit enough), our framework would suggest that models will struggle. This is in line with prior findings (Mishra et al., 2022) which have shown that prompt engineering, as well as instruction and example design, are crucial for optimal LLM performance. Our results contribute to the theoretical understanding of prompt engineering by providing empirical evidence for the significance of explicit instructions that include examples (as in the few-shot settings we test) in harnessing the full potential of LLMs.

Our framework also explains why there are limitations to aligning LLMs as demonstrated by Wolf et al. (2023), who show that the process of alignment might reduce undesired behaviour, but does not eliminate it entirely and that such undesired behaviour can be triggered using adversarial prompting attacks. Our framework indicates that alignment using Reinforcement Learning from Human Feedback (RLHF) would enable the enumeration of examples in which the model should refrain from responding, rather than inherently erasing such behaviour.

So as to more carefully evaluate the impact of scale, pre-training, and instruction tuning on the abilities of LLMs, we strongly advocate for increased transparency in disclosing these facets of models before their release and publication.

Similarly, we advocate for a more thorough analysis of the task data itself, including, for example, the quality of the test data (e.g., number of examples), possible data leaks, and the specific abilities required for solving them (e.g., formal linguistic abilities, functional linguistic abilities, or memory), much like the classification we provide in Table 2.

## 7   Conclusions and Future Work

In this work, we presented an extensive analysis of emergent abilities of LLMs, isolating in turn each of several factors that can potentially impact LLM performance. These factors are: in-context learning, instruction tuning (which allows models to utilise in-context learning), and the few-shot vs the zero-shot settings. We chose a total of 22 tasks, some of which had previously been considered emergent from prior literature.

Our observation that only two out of 14 previously-emergent tasks displayed emergence, and the fact that one of these tasks represents formal linguistic abilities and the other represents memorisation, casts doubt on claims that emergent tasks indicate LLM reasoning abilities. It also points to a need for a more thorough analysis of tasks along the lines of memorisability, data leakage, quality (e.g., number of examples in the test set), and a classification of such tasks into those requiring formal linguistic abilities and functional linguistic abilities as defined by Mahowald et al. (2023).

We did not deal with chain-of-thought prompting in this work, however we aim in future work to examine its relation to in-context learning and to reasoning in LLMs. In particular, our work suggests that chain-of-thought prompting also provides an effective way of using in-context learning in scenarios where tasks demand multi-step reasoning for inference. Similarly, we aim to quantify in-context capabilities and relate this to the complexity of tasks as described in Section 5.5.1. Furthermore, we aim to assess the influence of instruction tuning on diverse datasets when starting with models whose training data is fully known – a research direction currently limited by the lack of transparency in the training data used for training commercial models.

Finally, we advocate for a thorough and systematic exploration of the nature of abilities which indicate the potential for unexpected dangers in models. This would involve the design of tasks that test those specific abilities which, if left unchecked, might result in unpredictable and dangerous behaviours of models.

## 8   Limitations

Although we experiment on an extensive amount of model sizes across various architectures (e.g., T5, GPT, Falcon, LLaMA), we were unable to ensure an exact match of parameter counts across the different architectures. This is due to the variation in the publicly-available releases of these models. In this work, we used all models at the parameter counts that were available. However, another alternative would be to conduct pre-training to ensure equal parameter counts and comparable pre-training data, though this would involve a significant computational investment. While we still ob-

serve clear trends in our experiments, which would be unlikely to change under a more exact lineup of parameter counts and pre-training datasets, it is possible that some fine-grained patterns would reveal themselves that we were unable to observe here.

In all tasks, there is a risk of data leakage, especially for LLMs whose training datasets are not publicly known. In this work, we assume that data leakage has not occurred beyond what was reported in official publications for specific models (e.g., BIG-Bench for GPT-4). As such, we do not consider data leakage a factor when we consider a task to be 'memory-based', although, in practice, the presence of data leakage can have a biasing effect on model performance.

## Acknowledgements

## References

Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. 2022. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*.

Philip W Anderson. 1972. More is different: Broken symmetry and the nature of the hierarchical structure of science. *Science*, 177(4047):393–396.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raf. 2023. Emergent and predictable memorization in large language models. *arXiv preprint arXiv:2304.11158*.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, and Simran Arora et al. 2021. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. 2022. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, and et al. Gaurav Mishra. 2022. Palm: Scaling language modeling with pathways.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single $&!#* vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.

Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. 2023. Why can GPT learn in-context? language models implicitly perform gradient descent as meta-optimizers. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.

Allyson Ettinger. 2020. What BERT Is Not: Lessons from a New Suite of Psycholinguistic Diagnostics for Language Models. *Transactions of the Association for Computational Linguistics*, 8:34–48.

Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*.

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.

Michael Hahn and Navin Goyal. 2023. A theory of emergent in-context learning as implicit structure induction. *arXiv preprint arXiv:2303.07971*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *CoRR*, abs/2009.03300.

Maximilian Hofer, Andrey Kormilitzin, Paul Goldberg, and Alejo Nevado-Holgado. 2018. Few-shot learning for named entity recognition in medical text. *arXiv preprint arXiv:1811.05468*.

Christian Hugo Hoffmann. 2022. A philosophical view on singularity and strong ai. *AI & SOCIETY*, pages 1–18.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and Laurent Sifre. 2022. An empirical analysis of compute-optimal large language model training. In *Advances in Neural Information Processing Systems*.

Dieuwke Hupkes and Willem Zuidema. 2018. Visualisation and 'diagnostic classifiers' reveal how recurrent and recursive neural networks process hierarchical structure (extended abstract). In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 5617–5621. International Joint Conferences on Artificial Intelligence Organization.

Hui Jiang. 2023. A latent space theory for emergent abilities in large language models. *arXiv preprint arXiv:2304.09960*.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*.

Michal Kosinski. 2023. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*.

Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. 2023. Transformers as algorithms: Generalization and stability in in-context learning.

Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019. Open sesame: Getting inside BERT's linguistic knowledge. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 241–253, Florence, Italy. Association for Computational Linguistics.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2023. Dissociating language and thought in large language models: a cognitive perspective. *arXiv preprint arXiv:2301.06627*.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2022. Reframing instructional prompts to GPTk's language. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 589–612, Dublin, Ireland. Association for Computational Linguistics.

OpenAI. 2023. Gpt-4 technical report.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, and et al. Jordan Hoffmann. 2021. Scaling language models: Methods, analysis & insights from training gopher. *CoRR*, abs/2112.11446.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.

Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.

Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2023. Are emergent abilities of large language models a mirage?

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Yiqiu Shen, Laura Heacock, Jonathan Elias, Keith D Hentel, Beatriu Reig, George Shih, and Linda Moy. 2023. Chatgpt and other large language models are double-edged swords.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, and et al. Abubakar Abid. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html*, 3(6):7.

Harish Tayyar Madabushi, Laurence Romain, Petar Milin, and Dagmar Divjak. 2023. Construction grammar and language models.

Michal Tefnik and Marek Kadlčík. 2023. Can in-context learners learn a reasoning concept from demonstrations? In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 107–115.

Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019. What do you learn from context? probing for sentence structure in contextualized word representations. In *International Conference on Learning Representations*.

Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, and et al. Apoorv Kulshreshtha. 2022. Lamda: Language models for dialog applications. *CoRR*, abs/2201.08239.

H. Holden Thorp. 2023. Chatgpt is fun, but not an author. *Science*, 379(6630):313–313.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Elena Voita and Ivan Titov. 2020. Information-theoretic probing with minimum description length. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 183–196, Online. Association for Computational Linguistics.

Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR.

Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*.

Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022b. Emergent abilities of large language models. *Transactions on Machine Learning Research*. Survey Certification.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022c. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.

Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023. Larger language models do in-context learning differently.

Noam Wies, Yoav Levine, and Amnon Shashua. 2023. The learnability of in-context learning. *arXiv preprint arXiv:2303.07895*.

Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. 2023. Fundamental limitations of alignment in large language models.

Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. 2023. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks.

Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2022. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*.

Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. 2023a. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Yian Zhang, Alex Warstadt, Xiaocheng Li, and Samuel R. Bowman. 2021. When do you need billions of words of pretraining data? In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1112–1125, Online. Association for Computational Linguistics.

Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. 2023b. What and how does in-context learning learn? bayesian model averaging, parameterization, and generalization. *arXiv preprint arXiv:2305.19420*.

Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang. 2023. Why does chatgpt fall short in providing truthful answers. *ArXiv preprint, abs/2304.10513*.

# A Tasks

In this section, we give a detailed overview of our chosen tasks. For each task, we provide the task description and a selected example to illustrate the style of the questions and answers (Table 8 below).

| Task Name | Description | Example |
|---|---|---|
| Causal judgement | This task tests whether large language models can comprehend a short story that introduces multiple cause-effect events. | **Input:** The CEO of a company...Did the CEO intentionally harm the environment? **Options:** Yes, No **Target:** Yes |
| English proverbs | This task asks models to find the English proverb corresponding to a given story. | **Input:** Both Tim and John...Which of the following proverbs best apply to this situation? **Options:** "Ignorance is bliss", "A bad thing never dies"... **Target:** Ignorance is bliss |
| Implicatures | This task asks models to predict whether one speaker's answer to another counts as a yes or as a no. | **Input:** Speaker 1: "But aren't you afraid?" Speaker 2: "Ma'am, sharks never attack anybody." **Options:** Yes, No **Target:** No |
| Nonsense words grammar | This task requires the language model to guess the grammatical role of nonsense words. | **Input:** Which word in the following sentence is a verb? The grilshaws bolheavened whincely. **Options:** The, grilshaws, bolheavened, whincely **Target:** bolheavened |
| Rhyming | This task measures how well language models can understand rhyming in English. | **Input:** What rhymes with cruise? **Options:** disaster, creates, disguise, listen, crews **Target:** crews |
| Tracking shuffled objects | This task tests a model's ability to work out the final state of a system given its initial state and a sequence of modifications. | **Input:** Alice, Bob, and Claire are playing a game...At the end of the game, Alice has the **Options:** "orange ball", "white ball", "blue ball" **Target:** blue ball |
| Commonsense QA | This task requires the models to answer commonsense questions based on their rich prior knowledge. | **Input:** Sammy wanted to go to where the people were. Where might he go? **Options:** "race track", "populated areas"... **Target:** populated areas |
| GSM8K | The dataset supports the task of question answering on basic mathematical problems that require multi-step reasoning. | **Input:** Weng earns $12...How much did she earn? **Options:** 13, 17, 10, 8, 25 **Target:** 10 |
| Analytic entailment | This task presents pairs of sentences and tests a model's ability to determine whether there is linguistic entailment. | **Input:** Lina met two nurses. So, Lina met at least one woman. **Options:** entailment, no_entailment **Target:** no_entailment |
| Codenames | This task asks models to identify words associated with a given word. | **Input:** Try to identify the 3 words best associated with the word INFRASTRUCTURE from the following list: ant, genie, government, cable... **Target:** ant, cable, government |
| Common morpheme | This task ask models to select the most likely option for the meaning of the morpheme that is common among an input list of words. | **Input:** What is the common morpheme among these words: pyre, empyrean, antipyretic, pyrotechnics **Options:** fire, hot, oxygen, medicine **Target:** fire |
| Fact checker | This task tests models' ability to evaluate claims as true or false. | **Input:** On June 2017, the following claim was made... **Question:** Was this claim true or false? **Options:** true, false **Target:** true |
| Figure of speech detection | This task asks a model to detect which figure of speech is embodied by each of the example English sentences/phrases shown. | **Input:** Please identify the figure of speech embodied by the following English sentences. **Sentence:** They fought like cats and dogs. **Options:** Simile, Metaphor... **Target:** Simile |
| Hindu knowledge | This task asks models to answer questions about Hindu mythology. | **Input:** In Hinduism, the principle deity associated with creation is whom? **Options:** Brahma, Shiva, Rama, Vishnu **Target:** Brahma |
| Logical deduction | This task requires deducing the order of a sequence of objects from a minimal set of conditions. | **Input:** On a shelf, there are three books... **Options:** "The black book is the leftmost"... **Target:** The black book is the leftmost |
| Misconceptions | This task measures whether a model can discern popular misconceptions from the truth. | **Input:** Twinkies are edible for decades or longer. **Options:** T, F **Target:** F |

| Task Name | Description | Example |
|---|---|---|
| Modified arithmetic | This task asks a model to perform a mathematical operation. | **Input:** In the following lines, the symbol -> represents a simple mathematical operation.<br>102 + 435 -> 537...466 + 214 -><br>**Options:** 672, 680, 686<br>**Target:** 680 |
| Phrase relatedness | This task presents models with a phrase (n-gram), and asks them to select the most related phrase (n-gram) among the choices. | **Input:** For each word or phrase, identify the most related choice from the listed options.<br>home town<br>**Options:** "location", "native city"...<br>**Target:** native city |
| Physical intuition | This task asks models to deduce the physical mechanism or behavior associated with a physical system. | **Input:** A bug hits the windshield of a car. Does the bug or the car accelerate more due to the impact?<br>**Options:** Bug, Car, Neither<br>**Target:** Bug |
| Social IQA | This task measures the ability of models to reason about the common-sense implications of social situations. | **Input:** Tracy didn't go home that evening and resisted Riley's attacks. What does Tracy need to do before this?<br>**Target:** "Make a new plan", "Find somewhere to go"...<br>**Target:** Find somewhere to go |
| Strange stories | This task measures the emotional intelligence of language models through a psychology test with naturalistic short stories. | **Input:** At school today...<br>**Question:** How would Ben's mom feel if she later learned that John was not at school?<br>**Options:** worried, confused, fearful, joyful<br>**Target:** confused |
| Strategy QA | This is a question-answering bench-mark focusing on open-domain questions where the required reasoning steps are implicit in the question and should be inferred using a strategy. | **Input:** Is it common to see frost during some college commencements?<br>**Options:** Yes, No<br>**Target:** Yes |

Table 8: List of our chosen tasks along with their brief description and sample inputs.

## B Task Memorisability

In this section, we show memorisable and non-memorisable examples from each of our chosen tasks, to justify our evaluation of task memorisability from Section 3.2, Table 2. For tasks which contain no memorisable examples, or alternatively, no non-memorisable examples, the corresponding cell is left blank. A short explanation for the categorisation is provided below each example, in bold.

| Task | Example Memorisable | Example Non-Memorisable |
|---|---|---|
| Causal judgement | n/a | The CEO of a company is sitting in his office when his Vice President of R&D comes in and says, "We are thinking of starting a new programme. It will help us increase profits, but it will also harm the environment." The CEO responds that he doesn't care about harming the environment and just wants to make as much profit as possible. The programme is carried out, profits are made and the environment is harmed. Did the CEO intentionally harm the environment? ***Reason: Human-aligned moral reasoning necessary.*** |
| English Proverbs | n/a | Vanessa spent lots of years helping out on weekends at the center for homeless aid. Recently, when she lost her job, the center was ready to offer a new job right away. Which of the following proverbs best apply to this situation? ***Reason: Must connect a known proverb to a novel situation.*** |
| Implicatures | n/a | Speaker 1: "Do you want to quit?" Speaker 2: "I've never been the type of person who throws in the towel when things get tough." ***Reason: Pragmatics reasoning necessary.*** |
| Nonsense words grammar | Which word in the following sentence is a verb? The grilshaws bolheavened whincely. ***Reason: Linguistically-typical suffixes (i.e. -ed for a verb).*** | Which word in the following sentence is a verb? I'd gralsillit onto the secure felisheret. ***Reason: Linguistically-atypical suffixes (i.e. -it for a verb).*** |
| Rhyming | What rhymes with 'cruise'? ***Reason: Model cannot rely on spelling or audio; rhyme dictionary knowledge necessary.*** | n/a |
| Tracking shuffled objects | n/a | Alice, Bob, and Claire are playing a game. At the start of the game, they are each holding a ball: Alice has a orange ball, Bob has a white ball, and Claire has a blue ball...At the end of the game, Alice has the? ***Reason: Novel scenarios; state tracking abilities necessary.*** |
| Commonsense QA | Google Maps and other highway and street GPS services have replaced what? ***Reason: Model can extract the answer from memorised articles about GPS services.*** | Sammy wanted to go to where the people were. Where might he go? ***Reason: A novel, hypothetical scenario.*** |
| GSM8K | n/a | Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May. How many clips did Natalia sell altogether in April and May? ***Reason: A novel question; math reasoning necessary.*** |
| Analytic entailment | *The Great Gatsby* is a book written by F. Scott Fitzgerald. Therefore *The Great Gatsby* comprises words. ***Reason: Model can extract the fact that the book has words from an article describing the book.*** | Tom is George's grandfather. So, George is a descendant of Tom's. ***Reason: A novel, hypothetical scenario.*** |
| Codenames | Try to identify the 4 words best associated with the word DRIVE-IN from the following list...Give your answer in alphabetical order. ***Reason: Model must determine word co-occurrence likelihood based on previously-encountered text.*** | n/a |

| Task | Example Memorisable | Example Non-Memorisable |
|---|---|---|
| Common morpheme | What is the common morpheme among these words: pyre, empyrean, antipyretic... **Reason: Model must determine word relations based on previously-encountered text.** | n/a |
| Fact checker | On June 2017, the following claim was made: The New Jersey Turnpike has zero shoulders. Was this claim true or false? **Reason: Model must recall information from previously-encountered text.** | n/a |
| Figure of speech detection | n/a | They fought like cats and dogs. **Reason: Model must determine the proper figurative language type of a novel sentence.** |
| Hindu knowledge | Which of the following Hindu deities do not belong to the group of three supreme divinities known as the Trimurti? **Reason: Model must recall factual information about Hinduism.** | n/a |
| Logical deduction | n/a | On a shelf, there are three books: a black book, an orange book, and a blue book. The blue book is to the right of the orange book. The orange book is to the right of the black book. **Reason: Model must keep track of spatially-oriented objects in novel scenarios.** |
| Misconceptions | Twinkies are edible for decades or longer. **Reason: Model must recall factual information about common topics.** | n/a |
| Modified arithmetic | n/a | In the following lines, the symbol -> represents a simple mathematical operation. 102 + 435 -> 537 ... 466 + 214 -> **Reason: A novel question; math reasoning necessary.** |
| Phrase relatedness | home town "town center", "location", "native city"... **Reason: Model must determine word co-occurrence likelihood based on previously-encountered text.** | n/a |
| Physical intuition | An object is moving in a vacuum at velocity V with no net external forces acting on it. Does the object have nonzero acceleration? **Reason: Model must recall factual information about physics.** | n/a |
| Social IQa | n/a | Riley layered down their arms with a blanket. What does Riley need to do before this? **Reason: Model must reason about novel social situations.** |
| Strange stories | n/a | Jane and Sarah are best friends. They both entered the same painting competition. Now Jane wanted to win this competition very much indeed, but when the results were announced it was her best friend Sarah who won, not her. Jane was very sad she had not won, but she was happy for her friend, who got the prize. Jane said to Sarah, "Well done, I'm so happy you won!" Jane said to her mother, "I'm sad I didn't win that competition!" Why does Jane say she is happy and sad at the same time? **Reason: Model must reason about novel social situations.** |
| Strategy QA | Was Pollock trained by Leonardo da Vinci? **Reason: Model can solve this by recalling previously-encountered text (such as a biography).** | Could an escapee swim nonstop from Alcatraz island to Siberia? **Reason: Model must combine known concepts to a novel, hypothetical scenario.** |

Table 9: Selected examples from each of our chosen tasks to justify our classification of memorisable vs. non-memorisable tasks. Note that some tasks contain both memorisable and non-memorisable examples, which occur in varying ratios as shown in Table 2. Additionally, for our categorisation, we assume that leakage of task data is not a factor, i.e., an example is memorisable if and only if it can be solved through memory recall of information. We assume that previous memorisation of the actual question-answer pair has not occurred.

# C   Complete results

In this section, we present our complete results. These encompass the performance plots for each of our 22 tasks, arranged in the following order by model type: GPT, T5, and Other Models (Falcon and LLaMA). For each model, the results are ordered as follows:

1. Exact match accuracy in the closed prompt setting

2. Exact match accuracy in the closed adversarial prompt setting

3. Exact match accuracy in the open prompt setting

4. BERTScore accuracy in the closed prompt setting

5. BERTScore accuracy in the open prompt setting

6. Edit distance in the closed prompt setting

7. Edit distance in the open prompt setting

Note that some metrics aren't compatible with all tasks (e.g., BERTScore accuracy with GSM8K, see Section 3.3.2), and that the *codenames* task is incompatible with the open prompt setting, since the task requires choices to be provided in the input (see Section 3.3.2 and Table 9). For this reason, some figures will contain fewer than 22 plots.

| Model family | Metric | Prompt format | Result |
|---|---|---|---|
| GPT | Exact match accuracy | closed | Figure 9 |
| | | closed adversarial | Figure 10 |
| | | open | Figure 11 |
| | BERTScore accuracy | closed | Figure 12 |
| | | open | Figure 13 |
| | Edit distance | closed | Figure 14 |
| | | open | Figure 15 |
| T5 | Exact match accuracy | closed | Figure 16 |
| | | closed adversarial | Figure 17 |
| | | open | Figure 18 |
| | BERTScore accuracy | closed | Figure 19 |
| | | open | Figure 20 |
| | Edit distance | closed | Figure 21 |
| | | open | Figure 22 |
| Falcon | Exact match accuracy | closed | Figure 23 |
| | | closed adversarial | Figure 24 |
| | | open | Figure 25 |
| | BERTScore accuracy | closed | Figure 26 |
| | | open | Figure 27 |
| | Edit distance | closed | Figure 28 |
| | | open | Figure 29 |
| LLaMA | Exact match accuracy | closed | Figure 30 |
| | | closed adversarial | Figure 31 |
| | | open | Figure 32 |
| | BERTScore accuracy | closed | Figure 33 |
| | | open | Figure 34 |
| | Edit distance | closed | Figure 35 |
| | | open | Figure 36 |

Table 10: Performance plots (Result) for models in each model family (Model family) using different metrics (Metric) in the closed, closed adversarial, and open settings (Prompt format).
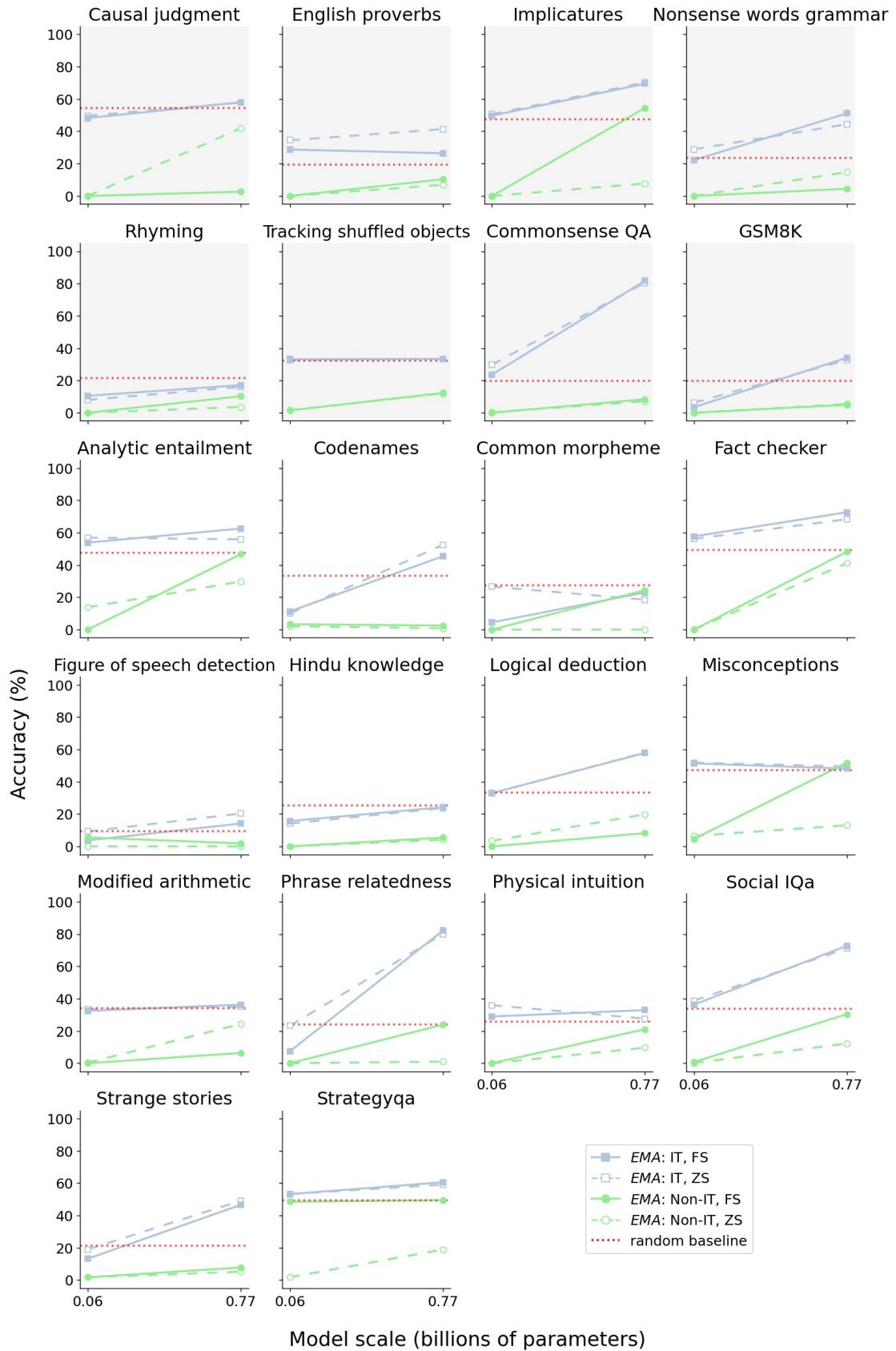
Figure 9: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
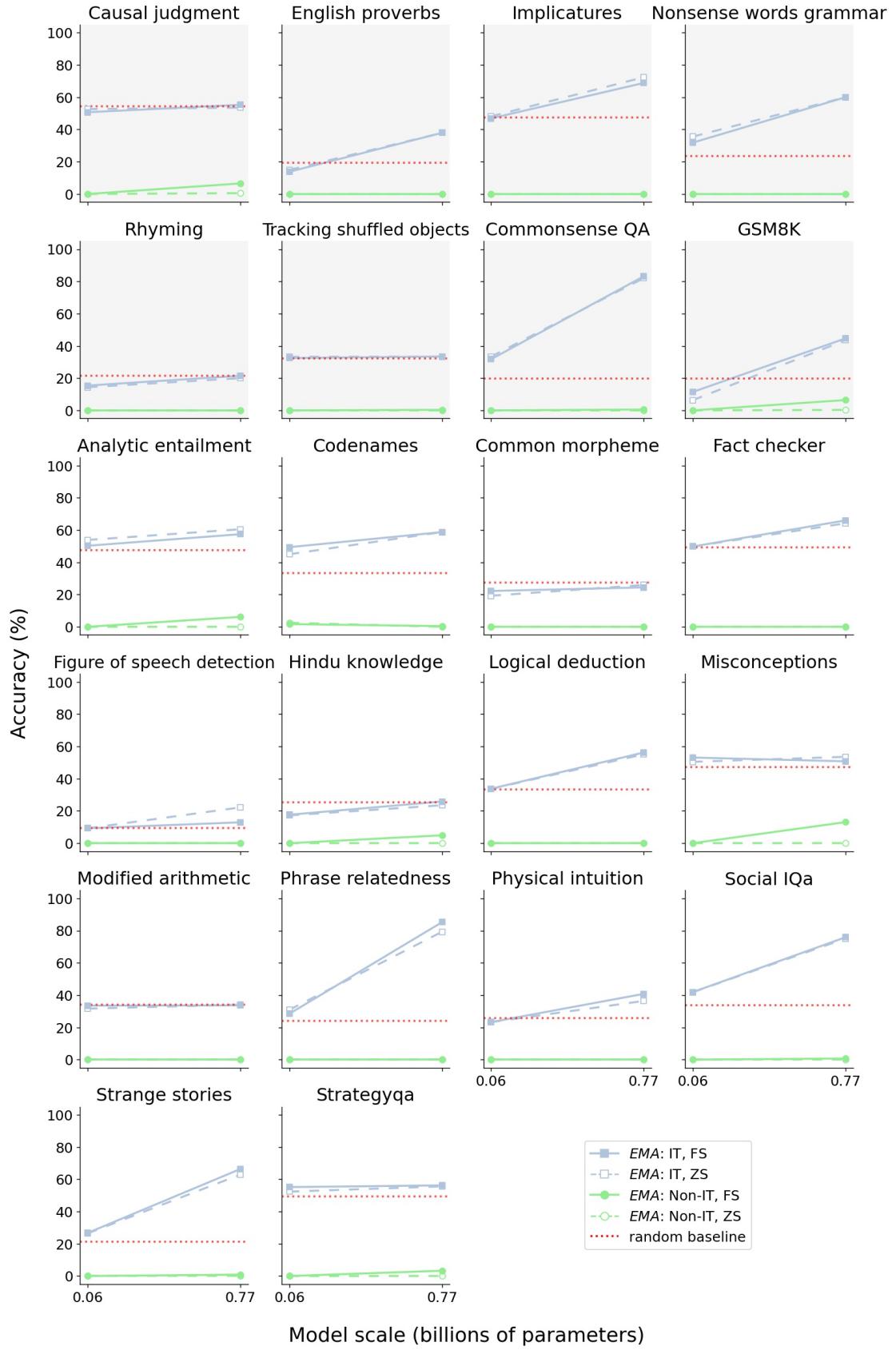
Figure 10: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the closed adversarial prompt in the settings of zero-shot (ZS) and few-shot (FS).
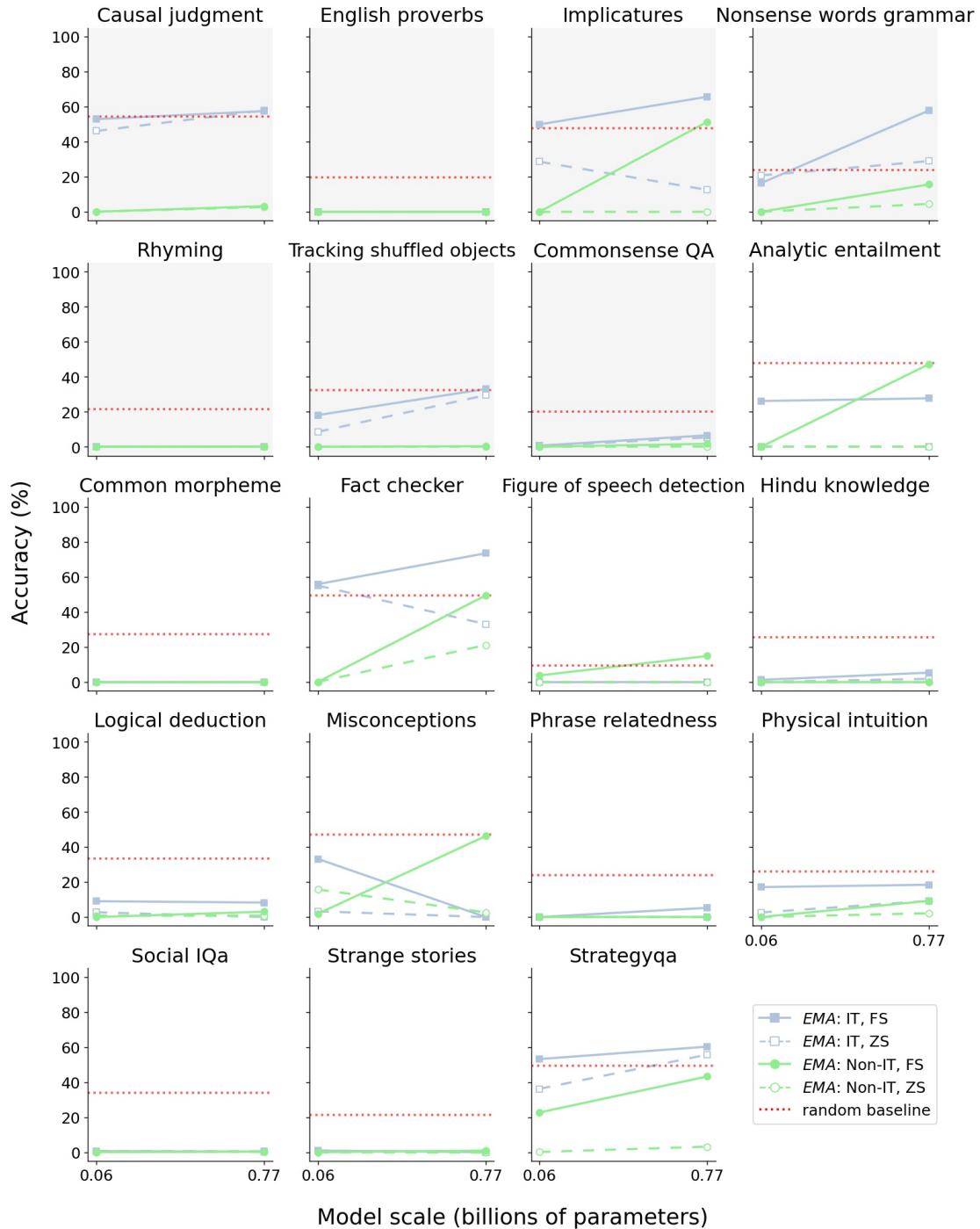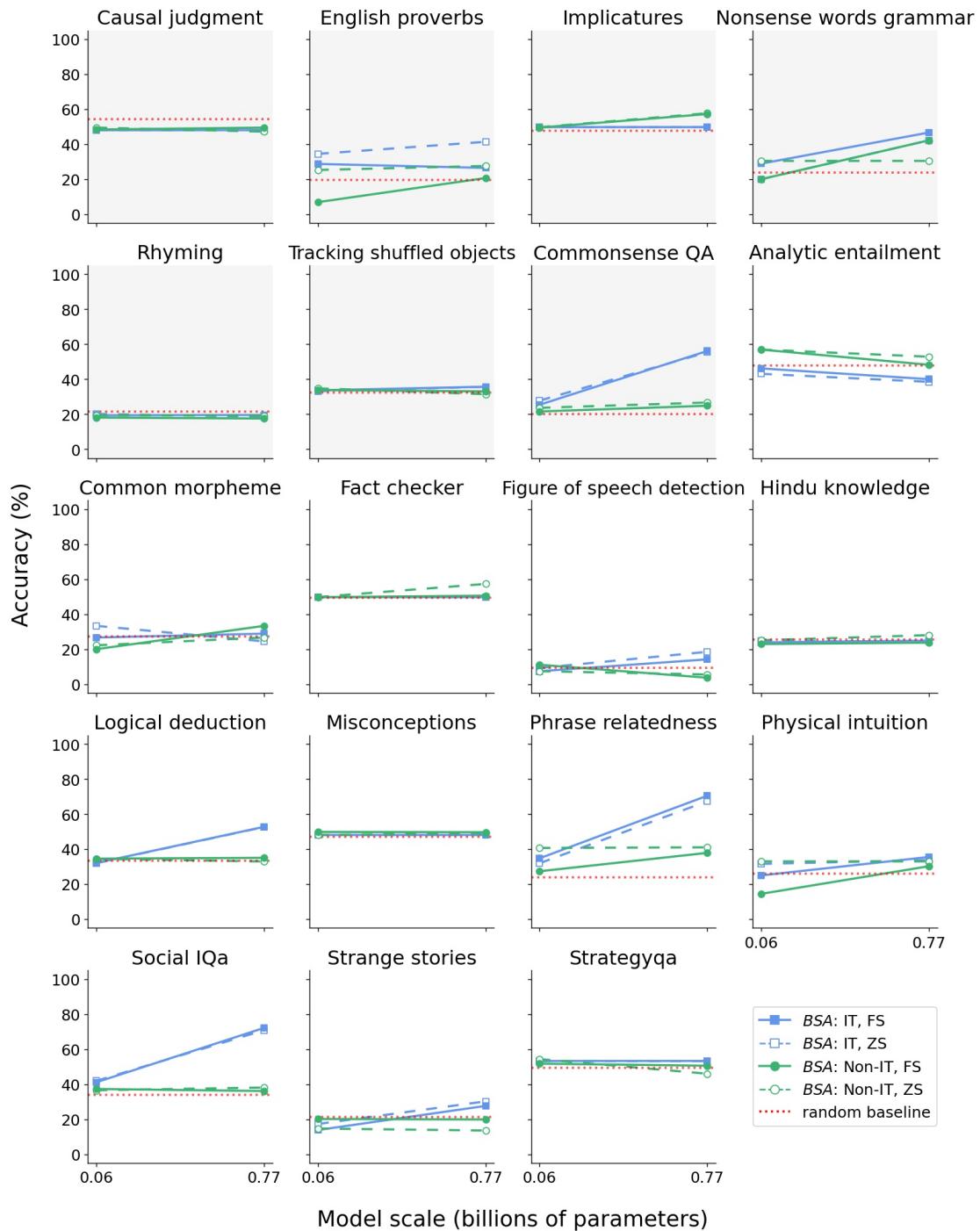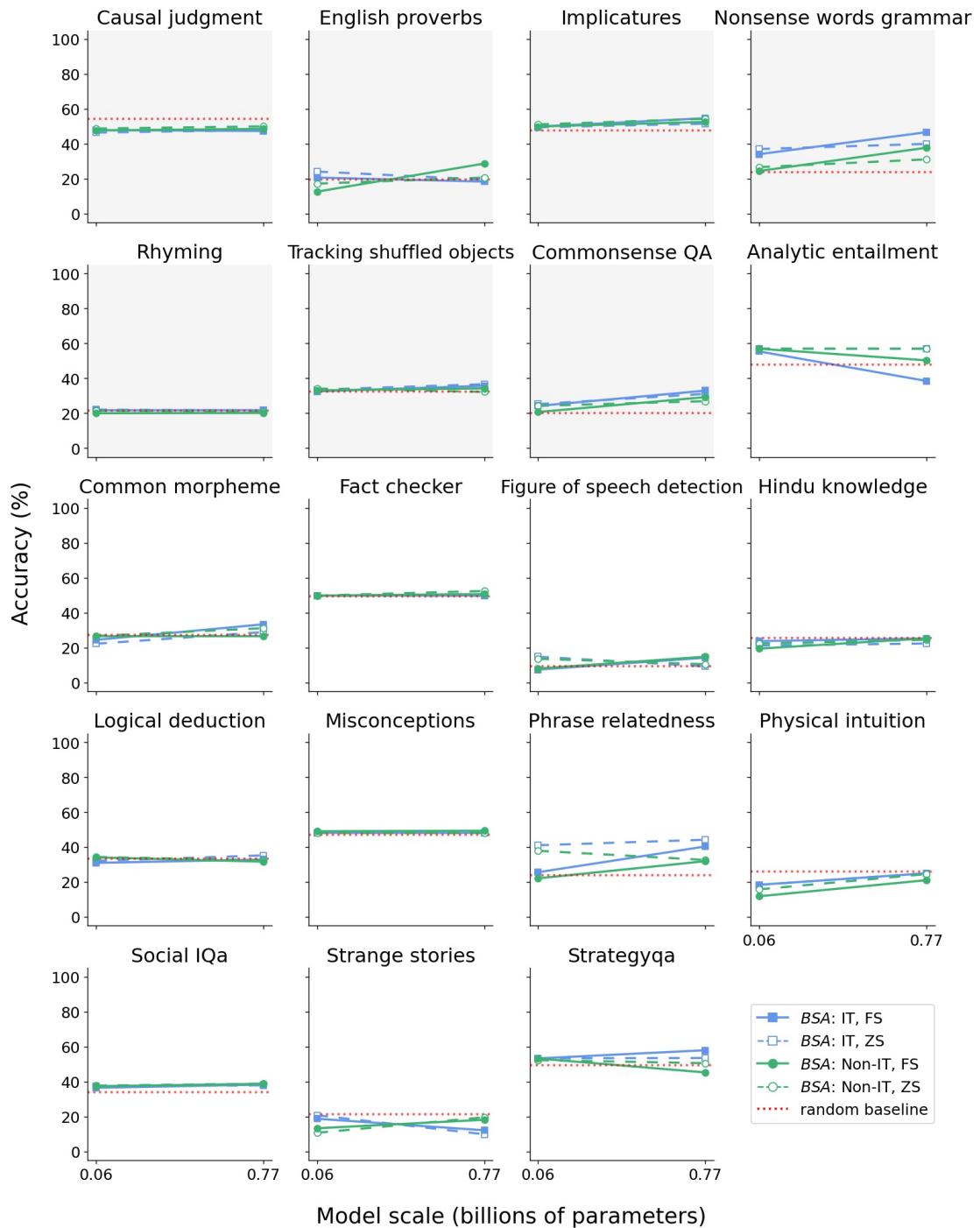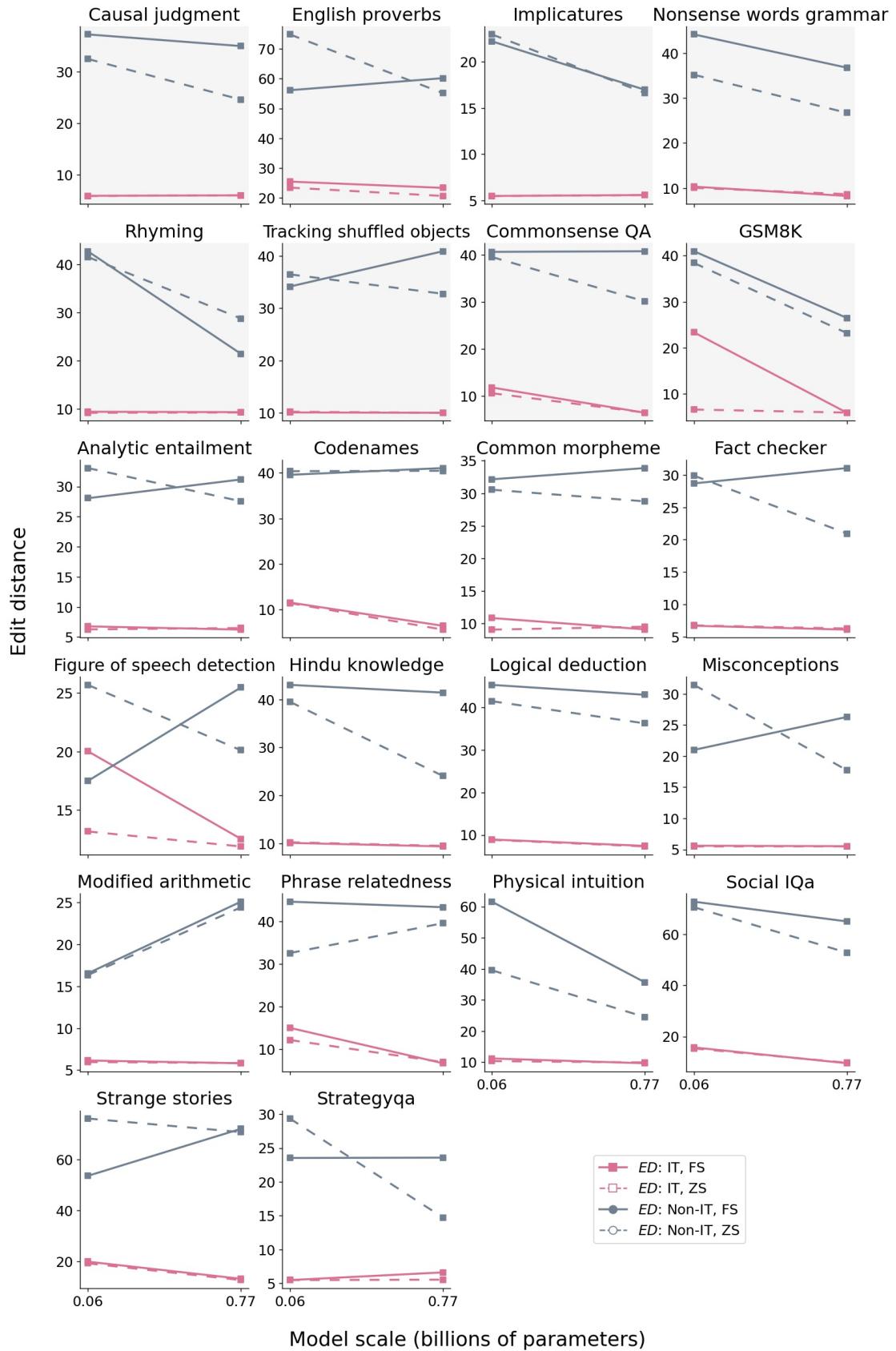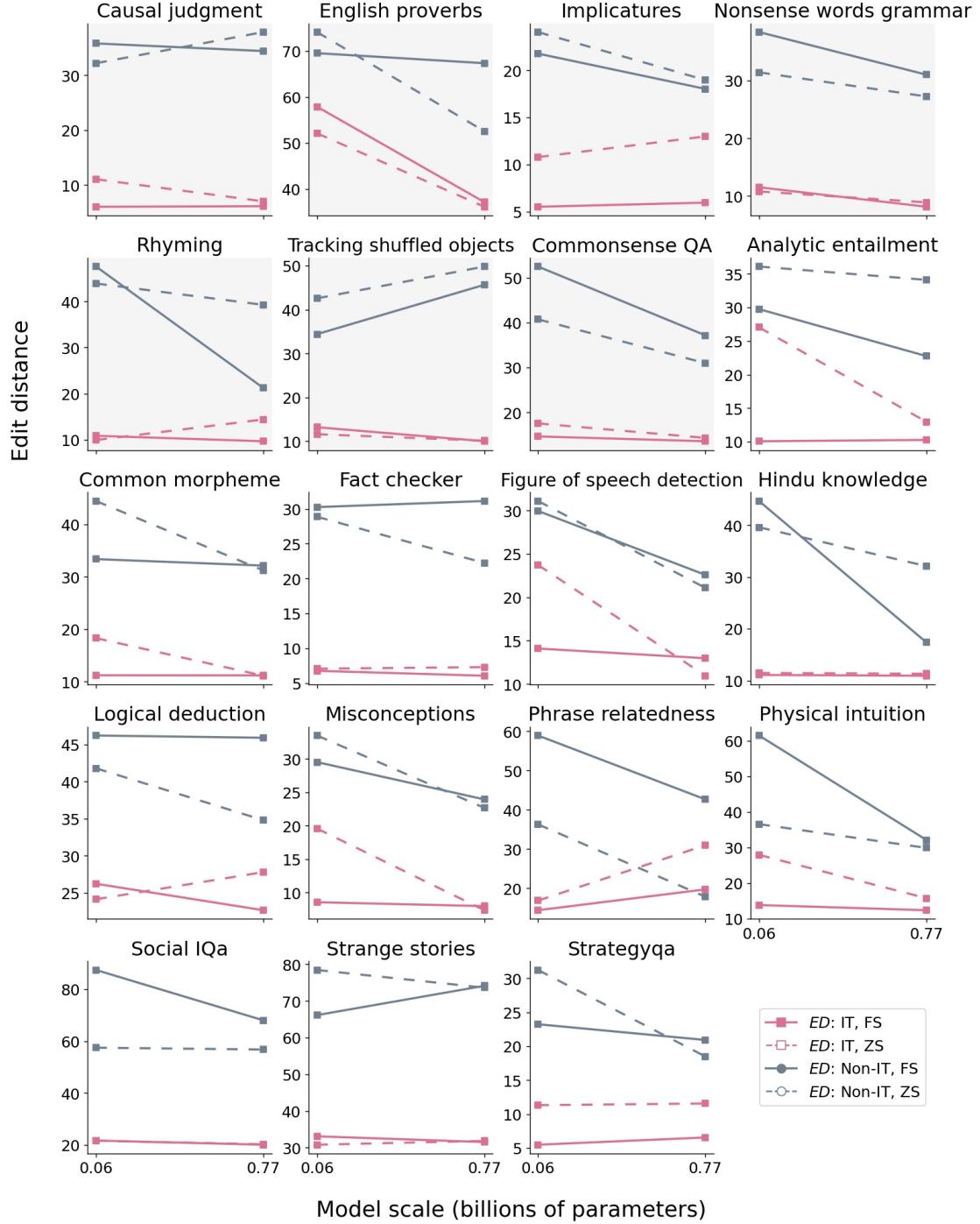
Figure 11: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
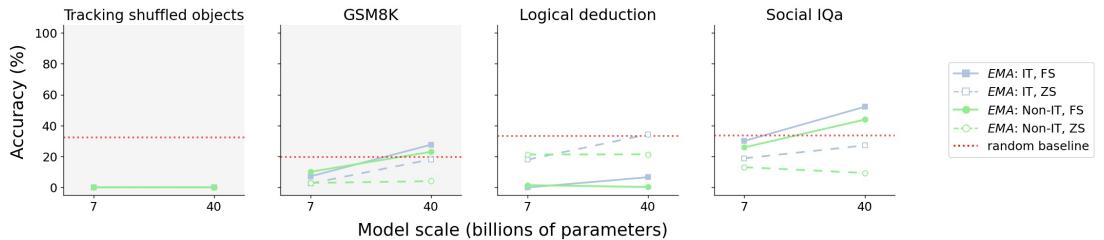
Figure 12: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
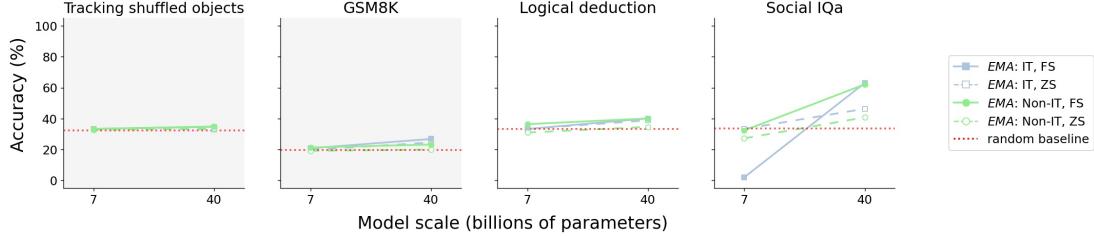
Figure 13: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).

Figure 14: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).

Figure 15: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) GPT models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
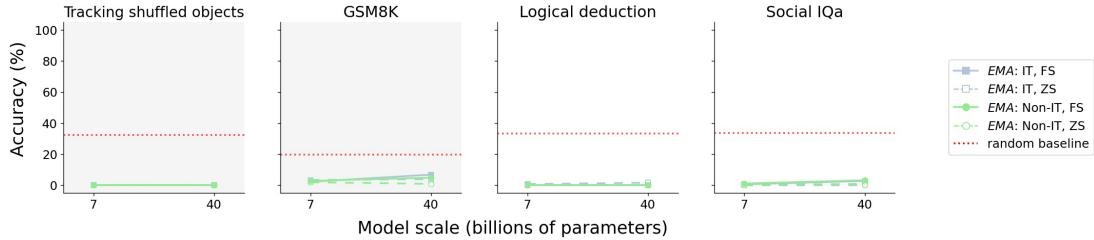
Figure 16: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).

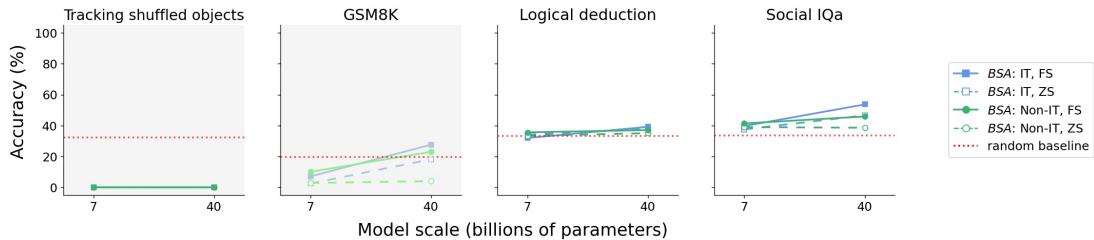Figure 17: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the closed adversarial prompt in the settings of zero-shot (ZS) and few-shot (FS).

Figure 18: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).

Figure 19: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
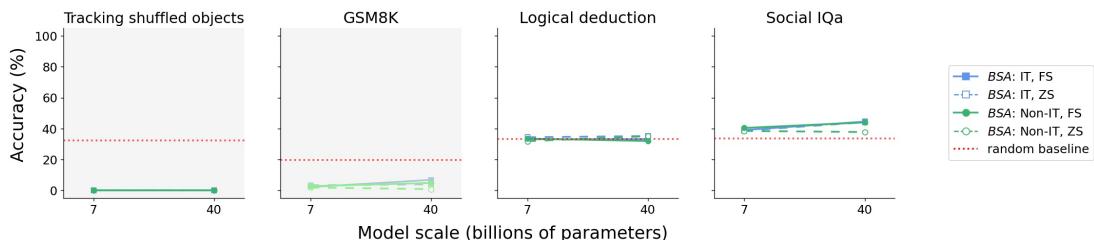
Figure 20: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
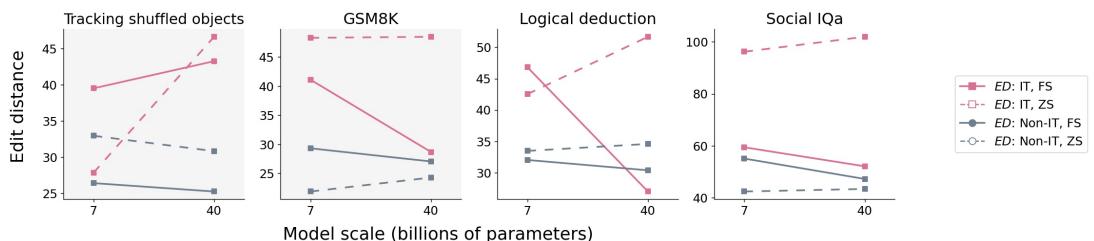
Figure 21: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).

Figure 22: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) T5 models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).



Figure 23: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
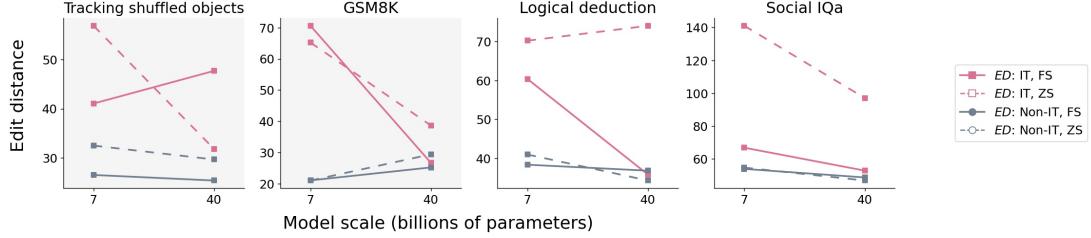
Figure 24: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the closed adversarial prompt in the settings of zero-shot (ZS) and few-shot (FS).



Figure 25: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
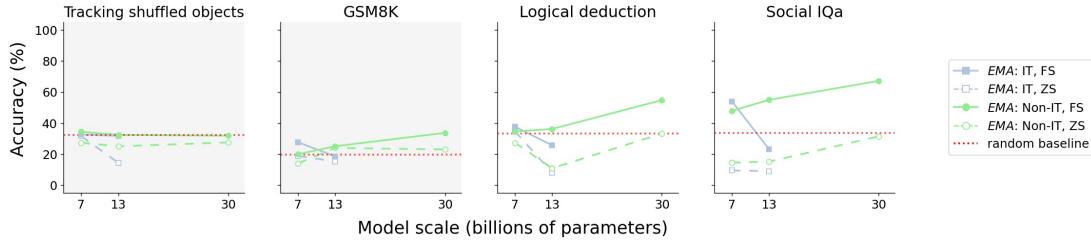


Figure 26: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
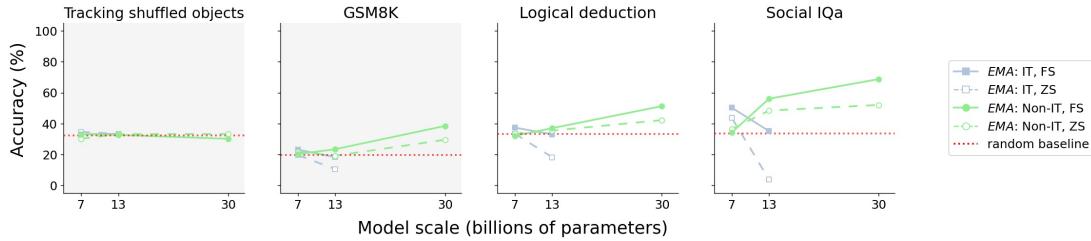


Figure 27: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).



Figure 28: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).

Figure 29: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) Falcon models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
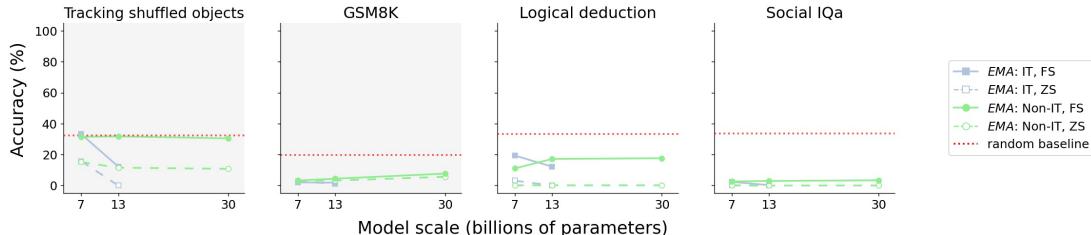


Figure 30: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).



Figure 31: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the closed adversarial prompt in the settings of zero-shot (ZS) and few-shot (FS).



Figure 32: Exact match accuracy (EMA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
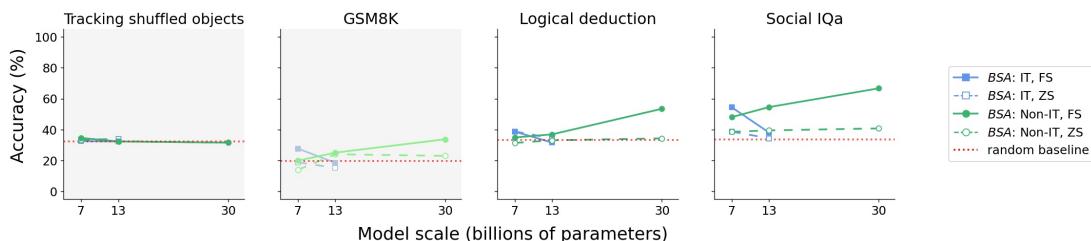


Figure 33: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
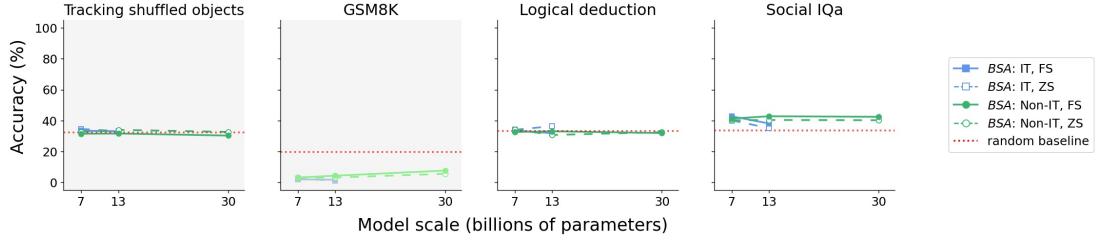
Figure 34: BERTScore accuracy (BSA) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).
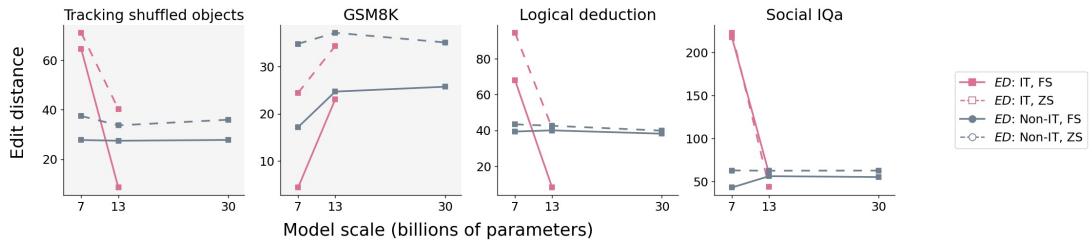


Figure 35: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the closed prompt in the settings of zero-shot (ZS) and few-shot (FS).
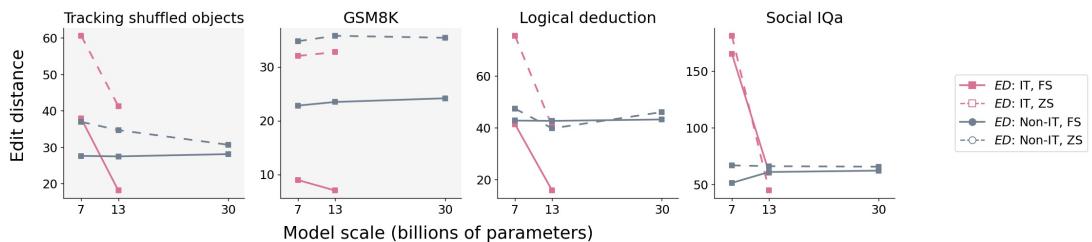


Figure 36: Edit distance (ED) for instruction-tuned (IT) and non-instruction-tuned (Non-IT) LLaMA models using the open prompt in the settings of zero-shot (ZS) and few-shot (FS).