# ALICE IN WONDERLAND: SIMPLE TASKS REVEAL SEVERE GENERALIZATION AND BASIC REASONING DEFICITS IN STATE-OF-THE-ART LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Large Language Models (LLMs) are often described as being instances of foundation models - that is, models that possess strong generalization and therefore transfer robustly across various tasks and conditions in few-show or zero-shot manner, while exhibiting scaling laws that predict generalization improvement when increasing the pre-training scale. These claims of strong generalization and advanced reasoning function enabling it rely on measurements by various standardized benchmarks where state-of-the-art (SOTA) models score high. We demonstrate here a dramatic breakdown of generalization and basic reasoning of all SOTA models which claim strong function, including advanced models like GPT-4 or Claude 3 Opus trained at the largest scales, using a simple, short common sense problem formulated in concise natural language, easily solvable by humans (AIW problem). The breakdown is dramatic as it manifests in both low average performance and strong performance fluctuations on natural problem variations that change neither problem structure nor its difficulty, while also often expressing strong overconfidence in the wrong solutions, backed up by plausible sounding explanation-like confabulations. Various standard interventions in an attempt to get the right solution, like chain-of-thought prompting, or urging the models to reconsider the wrong solutions again by multi step re-evaluation, fail. We take these observations to the scientific and technological community to stimulate re-assessment of the capabilities of current generation of LLMs as claimed by standardized benchmarks. Such re-assessment also requires common action to create standardized benchmarks that would allow proper detection of such deficits in generalization and reasoning that obviously remain undiscovered by current state-of-the-art evaluation procedures, where SOTA LLMs obtain high scores.<sup>1</sup>.

037

007

009

010 011 012

013

015

016

017

018

019

021

023

025

026

027

028

029

030

031

032

034

#### 1 INTRODUCTION

In the recent breakthroughs in transferable learning that were achieved in various classical domains 040 of machine learning like visual recognition (Radford et al., 2021) or language understanding (Devlin 041 et al., 2018; Raffel et al., 2020; Brown et al., 2020), large language models (LLMs) have played a very 042 prominent role. The generic form and scalability of autoregressive language modelling (Brown et al., 043 2020) allowed to push towards training scales not achievable before with conventional supervised 044 label-based learning. Scaling laws derived from experiments on smaller scales hinted on strong 045 function and generalization capability appearing at larger scales (Kaplan et al., 2020; Hoffmann et al., 046 2022), which was then confirmed by training models at the large scales, measuring their performance 047 on set of standardized benchmarks (MMLU, HellaSwag, ARC, MATH, GSM8k, etc) where they 048 scored high on few- and zero-shot transfer across various tasks (Kojima et al., 2022), following 049 accurately the predictions (Brown et al., 2020; Kaplan et al., 2020; Achiam et al., 2023; Touvron 050 et al., 2023a;b; Jiang et al., 2023).

There were however observations made by various works that questioned the claimed strong generalization, transfer and reasoning capabilities attributed to LLMs (Mitchell, 2023). These works pointed

<sup>1</sup>Code for reproducing experiments in the paper and raw experiments data can be found at AIW repo

054 out various function failures that were seemingly incompatible with postulated strong capabilities as 055 measured by standardized benchmarks (Wu et al., 2023; Golchin & Surdeanu, 2023; Li & Flanigan, 056 2024; Frieder et al., 2024). However, it has also been noted that observed failures can frequently 057 be addressed through simple adjustments to the prompts or by repeated execution and evaluation 058 using majority voting, or by requesting the model to perform self-verification (Kadavath et al., 2022; Wang et al., 2023; Zhou et al., 2023; Zhang et al., 2024; Pan et al., 2024). It remained thus unclear where those observations of failures are pointing to some fundamental deficits in core model capa-060 bilities affecting generalization and reasoning, or whether those are just symptoms of minor issues 061 easily resolvable by simple interventions, leaving claim of strong core function as put forward by 062 standardized benchmarks unaffected. 063

064 To shed light on current situation, we study whether the claim of SOTA LLMs possessing strong functions across various complex tasks can be put to test by using tasks that are very simple, in 065 contrast to those employed by standardized benchmarks. We introduce a short conventional common 066 sense problem that is formulated without any ambiguities in concise natural language and can be 067 easily solved by humans. The problem (in following Alice in Wonderland, AIW problem) has 068 following template: "Alice has N brothers and she also has M sisters. How many sisters does 069 Alice's brother have?". Crucially, instantiating natural numbers  $N, M \leq 7$  allows us to naturally introduce systematic controlled variations that do not change problem structure and difficulty and 071 thus should not affect ability to solve it. We use then this technique of creating problem structure and 072 difficulty preserving variations to measure models' sensitivity to problem irrelevant perturbations 073 across multiple repetitive trials, testing models' generalization ability.

074 Surprisingly, when confronted with AIW problem and its structure preserving variations, all SOTA 075 models including most advanced large-scale ones (eg GPT-4 (OpenAI, a), Claude 3 Opus (Anthropic, 076 2024c)) suffer severe function breakdown. This breakdown manifests (i) in average correct response 077 rates that are unexpectedly low for such a simple problem and in (ii) strong fluctuations in correct response rates across AIW problem variations despite those being entirely irrelevant for coping with 079 the problem. Strong fluctuations remain despite using various standard interventions to improve model function like chain-of-thought prompting. By creating further control versions of AIW problem 081 and observing that models are successfully coping with those, we are able to rule out that observed failures might be rooted in minor low level issues of tokenization/natural language parsing, handling the basic family structure, binding attributes, or executing arithmetic operations necessary to solve the 083 problem. Despite a specific form of simple AIW problem, we can thus conclude that observed failure 084 has generic character. The lack of robustness revealed in all SOTA models by problem irrelevant 085 variations of a simple problem points to severe generic deficits in generalization and reasoning. 086

087 The observed breakdown of function and generalization is in strong contrast to scores on standardized 088 benchmarks, which contain problems of higher difficulty. Many tested models that score high on such benchmarks show correct response rates close to zero across simple AIW problem variations. Claim 089 put forward by standardized benchmarks to properly reflect model capabilities such as generalization 090 and reasoning cannot be upheld in face of the evident failure to detect such severe function deficits 091 as revealed in the simple AIW problem setting. Our study highlights necessity to re-assess current 092 capabilities of SOTA LLMs by creating novel benchmarks that properly reflect their true abilities to 093 generalize and reason. Such benchmarks will be able to correctly spot deficits overlooked so far and 094 thus show the path for improvement of current still unsatisfactory state. 095

096

098

#### 2 METHODS & EXPERIMENT SETUP

#### 2.1 SIMPLE COMMON SENSE REASONING PROBLEMS AND THEIR VARIATIONS

100 AIW Problem. To measure models' sensitivity to problem irrelevant variations and thus probe the 101 zero-shot generalization, we use following problem template: "Alice has N brothers and she also 102 has M sisters. How many sisters does Alice's brother have?". The problem has a simple common 103 sense solution which assumes all sisters and brothers in the problem setting share the same parents. 104 The correct response C - number of sisters - is easily obtained by calculating M + 1 (Alice and 105 her sisters), which gives the number of sisters Alice's brother has. To create problem variations, we choose to vary natural numbers  $N, M \leq 7$ , obtaining AIW variations 1-4 (see Suppl. Tab. 2) 106 that all pose same problem using variations irrelevant for problem solving. We further use 3 prompt 107 types, RESTRICTED, STANDARD and THINKING, to ensure we measure models across various

prompt formulations for making general conclusions about their behavior (see Suppl. Sec. B for full overview)

Control AIW Light problems To control for models struggling either with basic family relations structure handling or with executing arithmetic operations in frame of the posed AIW problem, we make various versions of AIW problem - AIW Light Family, AIW Light Arithmetic Siblings and AIW Light Arithmetic Total Girls. The AIW Light problems keep problem template close to the original, changing only the final question part such that the posed modified question tests particular operations. The variations 1-4 are created in the same way like in AIW original by varying natural numbers of brothers and sisters, while ensuring that the natural numbers for final correct answers in AIW original and AIW Light are matched across variations 1-4 (see also Suppl. Sec. B)

118 **AIW Light Arithmetic Siblings**. AIW Light Arithmetic Siblings has following problem template: 119 "Alice has N brothers and she also has M sisters. How many siblings does Alice have?". Compared 120 to AIW original, only question part is modified. To solve the problem, summing up already given 121 numbers of brothers and sisters is sufficient - the correct answer is C = N + M. This requires basic 122 grasping of relational family structure (realizing Alice's siblings are her sisters and brothers) and 123 selection and execution of elementary arithmetic sum operation. In contrast to AIW original, it does not require execution of set operations nor binding sex attribute to Alice to properly assign her to 124 125 correct sets. Should the issues with solving AIW original be rooted in selection and execution of elementary arithmetic operations in family frame, we should see models also failing here. Again, we 126 create variations 1-4 by varying natural numbers N, M, such that correct responses C are matched 127 with AIW original variations 1-4 (Suppl. Tab. 3) 128

129 AIW Light Family. AIW Light Family has following problem template: "Alice has N brothers 130 and she also has M sisters. How many brothers does Alice's sister have?". Compared to AIW 131 original, only question part is modified. To solve the problem, reporting already given number of brothers is sufficient - the correct answer is C = N. This requires only basic grasping of relational 132 family structure (understanding entity "Alice's sister", binding female attribute to Alice and realizing 133 Alice and her sisters share same brothers). It does NOT require execution of any arithmetic or set 134 operations, in contrast to AIW original. Should the issues with solving AIW original be rooted in 135 handling basic family structure, we should see models also failing here. Again, we create AIW Light 136 Family variations 1-4 by varying natural numbers N, M, such that correct responses C are matched 137 with AIW original variations 1-4. (Suppl. Tab. 4) 138

AIW Light Arithmetic Total Girls. AIW Light Arithmetic Total Girls has following problem 139 template: "Alice has N brothers and she also has M sisters. How many girls are there in total?". 140 Compared to AIW original, only question part is modified. To solve the problem, it is necessary to 141 bind female attribute to Alice via the pronoun "she", to assign correct female attributes to the sisters 142 and to execute the correct arithmetic sum operation adding all the obtained girls - the correct answer 143 is C = M + 1. This requires basic grasping of family structure (realizing who are the girls in the 144 family) and selection and execution of elementary arithmetic sum operation. In contrast to AIW 145 original, it does not require execution of set operations to properly assign Alice to sisters set. Should 146 the issues with solving AIW original be rooted in binding correct sex attributes or counting total 147 members of particular sex in family frame given its structure, we should see models also failing here. Again, we create variations 1-4 by varying natural numbers N, M, such that correct responses C are 148 matched with AIW original variations 1-4. (Suppl. Tab. 5) 149

150 151

152

#### 2.2 PROMPT TYPES AND RESPONSE PARSING

153 **Model prompt types.** It is well known that so-called prompt engineering can heavily influence 154 the model behavior and model response quality (Arora et al., 2022; Wei et al., 2022; White et al., 155 2023). To check that our observations reflect model sensitivity to controlled problem structure 156 preserving variations in same manner independent of particular prompt type, we employed 3 various 157 prompt types to provide model's input: STANDARD (prompt with instruction to format final answer 158 output as a natural number), THINKING (prompt that in addition encourages thinking in spirit 159 of CoT) and RESTRICTED (prompt with instruction to output nothing else but final answer as a natural number). THINKING v2 prompt type is a minor variation of THINKING type that just 160 adds "step by step" after already existing "think carefully" phrasing (control experiments show that 161 THINKING and THINKING v2 are equivalent in terms of observed performance, so we use both

interchangeably). STANDARD and THINKING prompt types allow models to generate any text
 output before delivering the final answer, while RESTRICTED is used as control with restricted
 output to measure model behavior when the only output allowed is the final answer (Suppl. Tab. 2)

Furthermore, we make use of other prompt types (see Suppl. Sec.B for overview) to demonstrate various important properties and the different success or failure modes of the model behavior for the AIW problem. In those prompts, we re-use the main problem formulation as introduced in Sec. 2.1, while adding various modifications. This allows us for instance to observe confabulations that contain clearly broken statements with reasoning-like convincing sound backing up wrong final answers or responses showing model overconfidence.

171 **Parsing model responses.** To perform evaluations of model performance, it is necessary to parse 172 and extract the final answer from the responses provided by the models. Each input to the model 173 is combination of a AIW problem variation, followed by one of prompt types as described before. 174 To keep the parsing procedure simple, we add to each problem prompt following output format 175 instruction: "provide the final answer in following form: "### Answer: "". We observed that all 176 models we have chosen to test were able to follow such an instruction, providing a response that 177 could be easily parsed. We also ran control experiments without such formatting instruction in the 178 problem formulation, ensuring that behavior does not depend on it.

- 179
- 180 181

#### 2.3 Selecting models for evaluation and conducting experiments

We are interested in testing current state-of-the-art models that claim strong function, especially in generalization and reasoning, backed up by high scores shown on standardized benchmarks that are assumed to measure generalization and reasoning capabilities to solve problems. We therefore select models widely known and used in the ML community that also appear in the top rankings of the popular LLM leaderboards, like openLLM leaderboard by HuggingFace or ELO leaderboard by LMsys. We provide the overview of the selected models in Suppl. Tab. 1 and list in Suppl. Tab. 7 the corresponding standardized benchmarks where they obtain strong scores.

We expose selected SOTA LLMs, including most advanced models at largest scales (see Suppl. Tab. 1) to AIW problem variations 1-4 (Suppl. Tab. 2) and AIW Light control problems (Suppl. Tab. 3, 4, 5), using different prompt types as described above. For each combination of model, AIW problem variation and prompt type, at least 30 trials are collected to compute correct response rates, Suppl. Fig. 21. For details on correct response rates estimation procedure, see Suppl. Sec. A.2

We use hosting platforms that offer API access or local deployment via vLLM (Kwon et al., 2024) for testing the models, and automatize the procedure by scripting the routines necessary to prompt models with our prompts set. The routines are simple and can be used by anybody with access to the APIs (we used liteLLM and TogetherAI for our experiments) or to locally hosted models to reproduce and verify our results. We protocol all the data from interactions with the models to enable community checking. We release all the collected raw response data, correct response rates estimates and routines used to conduct experiments as open-source for reproducibility and further usage.

201 202

203 204

205

206

#### 3 Results

## 3.1 HUMPTY DUMPTY SAT ON A WALL: BREAKDOWN OF SOTA LLMS ON THE SIMPLE AIW PROBLEM

AIW reveals severe generalization and reasoning deficits in SOTA LLMs. Following the procedures described in Sec 2, we expose the selected models that claim strong function and reasoning capabilities (Suppl. Tab. 1) and measure their correct response rate performance across and for each AIW variations 1-4 using various prompt types, executing > 30 trials for each combination (see also Suppl. Tab. 2 and Suppl. Fig. 21). The results suggest that confronted with the AIW problem, models suffer a severe function breakdown. This breakdown has two main manifestations:

**1. Low correct response rates.** Despite evident problem's simplicity, many models are not able to deliver a single correct response, and the majority stay well below correct response rate of p = 0.2. We summarize the main results in the Fig. 1. The only major exceptions from the observation of very low correct response rates are the largest scale closed models GPT-4 and Claude 3 Opus. These two



Figure 1: Collapse of SOTA LLMs on AIW problem. (main) Models with non-zero AIW correct response rate, average over STANDARD, THINKING, RESTRICTED prompt types and AIW variations 1-4. Omitted models score 0. (inlay) Strong fluctuations on AIW variations 1-4, despite problem structure and difficulty remaining entirely unchanged across variations.

238

239

240

models at largest scales obtain correct response rates well above p = 0.3, leaving the remaining large and smaller scales open-weights (e.g., Mistral-7B, Mixtral, Qwen, Command R+, and Dbrx Instruct) and closed-weights models (e.g., Gemini Pro, Mistral Large) far behind. Remarkably, many models that claim high scores on standardized benchmarks, show very low correct response rates close to 0, eg. Llama-3-8B, Mixtral-8x22B, Qwen1.5-110B, or exhibit even complete breakdown on AIW with correct response rate of zero across all variations, eg Command R+ or Qwen1.5-72B (Suppl. Tab. 7)

249 The results presented in the Fig. 1 show estimates for correct response rates averaged across 250 RESTRICTED, STANDARD and THINKING prompt types (Suppl. Tab. 2, prompt IDs provided for reproducibility; Suppl. Fig. 8 with models scoring 0). RESTRICTED prompt type was used as 251 further control that forces models into short outputs, restricting the compute for providing a solution and thus serving as low baseline for the performance (see Suppl. Sec. C and Suppl. Fig. 11). Among 253 the 4 models that are able to cross p = 0.3, two clear winners are the GPT-40 (p = 0.649) and Claude 254 3 Opus (p = 0.431). The only open-weights model in this set of better performers is the rather older Llama-2 70B Chat (p = 0.3). For these better performers, when inspecting the responses with correct 256 final answers, we see also correct reasoning backing up the final answers. For the poor performing 257 models with low correct response rates, by inspecting those rare responses with correct answers we 258 also in some cases still can see correct reasoning. In the poor performers, among the responses with 259 a correct final answer we see however often responses where final answer, after careful inspection, 260 turns out to be an accident of executing entirely wrong reasoning with various mistakes leading 261 coincidentally to the final output number corresponding to the right answer. Such responses are encountered in models with low correct performance rates (p < 0.3) (see Suppl. Sec. D for response 262 examples), and we correct via manual inspection the status of correct response for such cases. 263

264
 2. Strong performance fluctuations across irrelevant AIW problem variations. Importantly, we
 also observe strong fluctuation of correct response rates across AIW variations 1-4 as introduced
 in Sec. 2. Such fluctuations strongly affect better performers with higher average correct response
 rates like GPT-4/40 and Claude 3 Opus. As shown in the Fig. 1 (inlay) for the STANDARD and
 Fig. 2 for the THINKING prompt type, the correct response rates can fluctuate between being close
 to being close to 0, depending on AIW variation. Remarkable is that such fluctuations appear
 despite AIW variations being all instances of the very same simple problem, as changes in numbers



Figure 2: Strong fluctuations across AIW problem variations, THINKING prompt. Also for better performers, eg GPT-40, GPT-4 and Claude Opus 3, correct response rates vary strongly from close to 1 to close to 0, despite AIW variations being irrelevant for problem structure (a color per each variation 1-4). This shows clear lack of model robustness, revealing generalization and basic reasoning deficits.



Figure 3: Correct response rates across AIW Light Arithmetic Siblings control problem variations 1-4 (THINKING v2 prompt type). Strong performance is observed across problem variations (a color per each variation 1-4; prompt IDs in the legend, Suppl. Tab. 3). Models that entirely collapse on AIW, like Command R Plus and Dbrx Instruct, are clearly able to solve this version, with correct response rates going up to 1 or close to 1 across all problem variations. This shows that executing arithmetic operations or handling basic family setting is not an issue for the tested models.

used across AIW variations do not change either the problem structure or its difficulty at all. This
 lack of robustness on such a simple problem hints on severe deficits in generalization. The strong
 fluctuations across variations appear independent of employed prompt types (Suppl. Fig. 9), while
 correct response rate averaged across all variations also varies across prompt types, showing in
 addition expected prompt type dependency (Suppl. Fig. 11, 12)

311

284

285

287

299

300

301

302

303

304 305

- 312 3.1.1 CONTROL EXPERIMENTS USING AIW LIGHT PROBLEMS
- In all following experiments, for each AIW variation, 60 trials were executed to estimate correct response rate and its variance.

316 AIW Light Arithmetic Siblings. We show tested models' performance in Fig. 3. While all 317 tested models clearly have struggled with AIW original (Fig. 1, Suppl. Fig. 8), we observe them 318 successfully solving AIW Light Arithmetic Siblings. Correct response rates go high up close to 319 1 for most tested models across all variations 1-4. This is also the case for the models that show 320 very low correct response rates close to 0 or 0 on AIW original, like Command R+ or Dbrx Instruct 321 (Suppl. Fig. 8, Suppl. Tab. 7). Strong fluctuations we observe across variations on AIW original (Fig. 1, 2) also disappear. This clearly demonstrates that models neither struggle with basic grasping of 322 relational family structure - realizing Alice's siblings are her sisters and brothers, nor with selection 323 and execution of elementary arithmetic sum operation.



Figure 4: Correct response rates across AIW Light Family control problem variations 1-4 (THINKING v2 prompt type). Strong performance is observed across problem variations (a color per each variation 1-4). Models that entirely collapse on AIW, like Command R Plus and Dbrx Instruct, are clearly able to solve this version, with correct response rates going up to 1 or close to 1 across all problem variations. This shows that handling basic family relations is not an issue for the tested models.



Figure 5: Correct response rates across AIW Light Arithmetic Total Girls control problem variations 1-4 (THINKING v2 prompt type). Strong performance is observed across problem variations (a color per each variation 1-4; prompt IDs in the legend, Suppl. Tab. 5). Models that entirely collapse on AIW, like Command R Plus and Dbrx Instruct, are clearly able to solve this version, with correct response rates going up to 1 or close to 1 across all problem variations. This rules out that either binding of female attributes to Alice and the sisters entities or selection and execution of arithmetic operations necessary to count total females is an issue for the tested models.

AIW Light Family. We show tested models' performance in Fig. 4. Also here we observe all the tested models that are struggling with AIW original successfully solving AIW Light Family. Correct response rates go high up close to 1 for most tested models across all variations 1-4. This is also the case for the models that show very low correct response rates close to 0 or 0 on AIW original. like
 Command R+ or Dbrx Instruct (Suppl. Fig. 8 & Tab. 7). Also strong fluctuations that we observe across variations on AIW original (Fig. 1, 2) disappear. This clearly demonstrates that models handle well basic grasping of relational family structure - understanding entity "Alice's sister", binding female attribute to Alice and realizing Alice and her sisters share same brothers.

AIW Light Arithmetic Total Girls. We show tested models' performance in Fig. 5. Again, we observe also here strong performance for all tested models that clearly have struggled with AIW original. Correct response rates go high up close to 1 for most tested models across all variations 1-4. This is also the case for the models that show very low correct response rates close to 0 or 0 on AIW original. like Command R+ or Dbrx Instruct. Also strong fluctuations that we observe across variations on AIW original (Fig. 1, 2) are gone. This clearly demonstrates that models successfully cope with binding female attribute to entity of Alice, handle assignment of correct female attributes to the sisters and select and execute the correct arithmetic sum operation adding all the girls together. 

From these control experiments, we are thus able to obtain strong evidence that all tested models do
not suffer from low-level issues with tokenization and natural language or natural numbers parsing and
can handle well basic family relations structure and selection and execution of elementary arithmetic
operations necessary to solve AIW problem. This further strengthen the hypothesis that observed
failures and strong fluctuations in all tested SOTA models on AIW problem (Fig. 1, 2) are rooted in
problem unspecific, generic deficits in generalization and basic reasoning.



Figure 6: Altering model performance by fully redundant information. Adding fully redundant information "Alice is female" leads to increase of average correct response rates in (a) compared to AIW original (b) (see also Suppl. Fig. 10). For some models, eg Llama 3 70B or Qwen 2 72B, this boost via redundant info is especially pronounced and happens across all variations, resulting in clear overall improvement from (b) to (a). Strong fluctuations across variations 1-4 persist. This again shows lack of model robustness, hinting on severe generalization and basic reasoning deficits.

- 407
- 408
- 409 410

#### 3.2 CURIOUSER AND CURIOSER: FURTHER PROPERTIES OF OBSERVED BREAKDOWN

411 Boosting by reduntant information and persisting fluctuations: Alice female power boost. One 412 clear signature of generalisation and reasoning breakdown are the strong fluctuations we observe across AIW problem variations 1-4 that differ only in instantiated numbers (Fig. 2). We investigate a 413 further AIW problem version by adding "Alice is female" to the original AIW problem formulation 414 (see Suppl. Tab. 6 and Suppl. Sec. C.1). This is a fully redundant information, as Alice's gender is 415 already unambiguously specified by the "she" pronoun used in original AIW problem. As evident 416 from Fig. 6 and Suppl. Fig. 10, the average correct response rates are increasing, despite the provided 417 "female boost" information being entirely redundant and not revealing anything new necessary for 418 AIW problem solution. Altering performance by fully redundant information that should not affect 419 problem solving reveals again deficits in generalization and basic reasoning. While average correct 420 response rates increase, the strong fluctuations across AIW variations 1-4 remain (Fig. 6a). For 421 instance, GPT-40 has on AIW variations 2,4 correct response rate close to 1, while dropping heavily 422 for AIW variations 1,3, showing same lack of robustness despite the average boost.

423 Standardized benchmarks failure. We observe failure of standardized reasoning benchmarks 424 to properly reflect generalization and basic reasoning skills of SOTA LLMs by noting significant 425 disparity between the model's performance on the AIW problem and the scores on conventional 426 standardized benchmarks. All of the tested models report high scores on various standardized 427 benchmarks that claim to test problem solving via reasoning, e.g. MMLU, ARC, Hellaswag. Our 428 observations of SOTA models breaking down on the simple AIW problem hint that the benchmarks do 429 not reflect deficits in generalization and basic reasoning of those models properly. We visualize this failure by plotting scores tested models obtain on wide-spread and accepted standardized benchmarks 430 like MMLU versus the performance we observe on our proposed AIW problem. As strikingly evident 431 from Fig. 7, there is a strong mismatch between high scores on MMLU reported by the models and



Figure 7: Failure of standardized benchmark MMLU to properly reflect and compare model basic reasoning capabilities as shown by strong discrepancy between AIW correct response rate vs MMLU average score. Many models, eg. Command R+, score 0 on AIW, but have high MMLU score.

- the correct response rates they obtain on AIW. This mismatch and lack of differentiation makes it impossible for a given model to predict from its score on MMLU whether it will suffer breakdown on a simple problem like AIW, making the score unreliable for measuring core capabilities. Also model ranking fails, as models claiming higher scores can be strongly outperformed by models with lower scores when looking at their correct response rates on simple AIW problem. For instance, Llama-2-70B with lower MMLU score clearly outperforms on AIW problem models (eg Mistral-Large, Command R+, Dbrx Instruct) that are crowded in high MMLU - low AIW score region (left upper part of Fig. 7). For similar evidence on other standardized benchmarks, see Suppl. Sec. C.3
- Further relevant observations. 1. Dominance of wrong responses We measure distribution of natural numbers responses on output, showing that for AIW variations with low correct response rate, peaks are on wrong answers, excluding majority voting methods as a fix. (Suppl. Sec. C.2) 2. Confabulations and overconfident tone We observe that wrong responses are often accompanied by persuasive explanation-like confabulations and overconfident tone about correctness of the wrong solutions provided by the models, which can further mislead model users (Suppl. Sec. E) 3. Inability to revise wrong responses. Models show failure to properly detect mistakes and to revise wrong solutions when encouraged to do so in experiments with multi-turn AIW problem interaction and self-verification. (Suppl. Sec. F). 4. Reformulation of AIW as relational SQL database problem. We make use of relational logic underlying the AIW problem structure and prompt models to reformulate AIW into a correct relational SOL database format, to test their ability to extract formal problem structure. We observe that smaller scale models, and also some larger scale ones, consistently fail to generate a correct relational SQL form. Some models that are able to do so more frequently, e.g., Mistral/Mixtral, still fail to provide correct final answer most of the time (Suppl. Sec. G)

#### **RELATED WORK & LIMITATIONS**

Measuring LLMs capabilities. Since the seminal breakthroughs in language modelling (Devlin et al., 2018; Raffel et al., 2020; Brown et al., 2020), measuring LLM capabilities became indispensable for evaluations and model comparison. To measure how well a language model performs on reasoning, there exists a plethora of different standardized reasoning benchmarks. These benchmarks can be roughly divided into categories by what exact reasoning capability we want to test such as ARC (Clark et al., 2018), PIQA (Bisk et al., 2020), GSM8K (Cobbe et al., 2021), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2020) or WinoGrande (Sakaguchi et al., 2019). Multiple works aim on improving reasoning performance of LLMs as measured by those standardized benchmarks in various ways (Wei et al., 2022; Yao et al., 2024; Zhou et al., 2022; Wang et al., 2022; Pfau et al., 2024).

486 Stress-testing LLMs' weaknesses. Paralleling impressive progress shown by LLM research, cautious 487 voices have been raising concern about discrepancy between claimed capabilities as measured by 488 standardized benchmarks and true LLM reasoning skills by presenting carefully selected evidence for 489 model failures (Mitchell, 2023). In response, the research community has been undertaking attempts 490 to create more challenging benchmarks like HELM (Liang et al., 2023) or BIG-bench (Srivastava et al., 2023). These benchmarks also aimed at properly testing generalization capabilities beyond 491 memorization, in line with recent works that pointed out high test dataset contamination due to 492 large-scale pre-training on web-scale data (Golchin & Surdeanu, 2023; Li & Flanigan, 2024). 493

494 Similar in spirit to our work, multiple studies (Wu et al., 2023; Dziri et al., 2024; Lewis & Mitchell, 495 2024; Berglund et al., 2023; Moskvichev et al., 2023; Huang et al., 2023) have shown breakdowns of 496 language models reasoning capabilities in different scenarios and also lack of robustness to variation of problem formulation (Zong et al., 2024; Zheng et al., 2024). Other works were looking into 497 particular reasoning failures like deficits in causality inference (Jin et al., 2023b;a). These works 498 operate often with formalized, rather complex problems that does not have simple common sense 499 character. Here we show breakdown on a common sense problem with very simple structure using 500 natural, controlled variations that keep problem structure and difficulty unchanged, which emphasizes 501 a generic deficiency in generalization and basic reasoning about problem structure. A key limitation 502 of our current approach is the lack of sufficient diversity in AIW problem variations. This can be addressed in future work by systematic procedural instance generation for broader response 504 evaluation. 505

- 505
- 506 507

508 509

#### 5 DISCUSSION & CONCLUSION

In our work, using a very simple AIW problem (Sec. 2) that can be easily solved by adults and arguably 510 even children, we observe a striking breakdown of SOTA LLMs performance when confronted with 511 the AIW problem and its variations (Suppl. Tab. 2). The breakdown is manifested in (i) Low 512 average correct response rates (Fig. 1) and (ii) Strong performance fluctuation across structure and 513 difficulty preserving natural variations of the same problem, which hints at fundamental issues with 514 the generalization capability of the models (Fig. 2). The observed breakdown is in dramatic contrast 515 with claims about strong core functions of SOTA LLMs. Specifically, the claim of strong reasoning 516 cannot hold, as any system claiming even basic reasoning should be able to obtain 100% correct 517 response rates on problems as simple as AIW. The evidence also falsifies the claim of strong zero-shot 518 generalization - as in such a simple problem, strong performance fluctuations across variations that 519 keep problem structure and difficulty unchanged reveal severe generalization deficits in all tested 520 SOTA LLMs. By executing control experiments, we provide evidence that the observed failures are not specific to the problem type we study and thus hint on generic deficits in generalization and 521 basic reasoning (Sec. 3.1.1). Our study also clearly points to failure of standardized benchmarks to 522 properly measure core model functionality such as generalization or reasoning (Suppl. Sec. C.3, Fig. 523 7, Tab. 7). Standardized benchmarks assigning high scores to SOTA LLMs fail to reveal severe model 524 weaknesses made evident by breakdown on simple AIW problem. It has to be noted that despite 525 observed breakdowns with low average correct response rates, the reasoning is not entirely absent and 526 better performer larger scale models like GPT-4 or Claude 3 Opus do show examples of fully correct 527 reasoning (see Suppl. Fig. 25, 26). As our results show, this reasoning capability is however fragile 528 and cannot be accessed robustly, even in such a simple scenario as posed by AIW problem variations.

529 The observations urge re-assessment of the claimed capabilities of current generation of LLMs, 530 with evidence suggesting that current SOTA LLMs are not capable of strong generalization and 531 robust reasoning, and enabling those is still subject of basic research. Such re-assessment also 532 requires common action to create standardized benchmarks that would allow proper detection of 533 such generalization and basic reasoning deficits as observed in our study that obviously manage to 534 remain undiscovered by current state-of-the-art evaluation procedures and benchmarks. Variations built into problem templates can serve as technique to create new benchmarks that are, in contrast to 536 current common benchmarks, no longer static and can serve as better measurement tool for properly 537 testing model generalization and reasoning. New benchmarks should follow Karl Popper's principle of falsifiability (Popper, 1934), attempting everything to break model's function, highlighting its 538 deficits, and thus showing possible directions for model improvement, which is the way of scientific method, also offering protection from overblown claims about models' core functions.

### 540 REFERENCES

547

551

556

567

568

569

570

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
  Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Alibaba. Introducing qwen1.5, 2024a. URL https://qwenlm.github.io/blog/qwen1.
   5/.
- 548 Alibaba. Hello qwen2, 2024b. URL https://qwenlm.github.io/blog/qwen2/.
- 549 Anthropic. Introducing the next generation of claude, 2024a. URL https://www.anthropic. com/news/claude-3-family.
- 552 Anthropic. Claude 3.5 sonnet, 2024b. URL https://www.anthropic.com/news/ claude-3-5-sonnet.
- Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024c. URL https://www.
   anthropic.com/claude-3-model-card.
- Simran Arora, Avanika Narayan, Mayee F Chen, Laurel Orr, Neel Guha, Kush Bhatia, Ines Chami, and Christopher Re. Ask me anything: A simple strategy for prompting language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
   Yu Han, Fei Huang, et al. Qwen technical report. arXiv preprint arXiv:2309.16609, 2023.
- Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. The reversal curse: Llms trained on" a is b" fail to learn" b is a". *arXiv preprint arXiv:2309.12288*, 2023.
- 566 Berri.AI. Litellm, 2024. URL https://github.com/BerriAI/litellm.
  - Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
   Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
   few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
  Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 581 Cohere. Command r+ documentation, 2024a. URL https://docs.cohere.com/docs/
   582 command-r-plus.
- 583 584 584 585 Cohere. Model card for c4ai command r+, 2024b. URL https://huggingface.co/ CohereForAI/c4ai-command-r-plus.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
   bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. Faith and fate: Limits of transformers on compositionality. *Advances in Neural Information Processing Systems*, 36, 2024.
- Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz,
   Philipp Petersen, and Julius Berner. Mathematical capabilities of chatgpt. Advances in Neural Information Processing Systems, 36, 2024.

626

627

- 594 Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large 595 language models. arXiv preprint arXiv:2308.08493, 2023. 596
- Google. Gemma: Introducing new state-of-the-art open models, 2024a. URL https://blog. 597 google/technology/developers/gemma-open-models/. 598
- Google. Gemma model card, 2024b. URL https://ai.google.dev/gemma/docs/ 600 model card. 601
- 602 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In International Conference on 603 Learning Representations, 2020. 604
- 605 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Ruther-606 ford, Diego de las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, 607 Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon 608 Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack William Rae, and Laurent Sifre. An 609 empirical analysis of compute-optimal large language model training. In Alice H. Oh, Alekh Agar-610 wal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing 611 Systems, 2022. URL https://openreview.net/forum?id=iBBcRUIOAPR.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei 613 Han. Large language models can self-improve. In The 2023 Conference on Empirical Methods 614 in Natural Language Processing, 2023. URL https://openreview.net/forum?id= 615 uuUQraD4XX. 616
- 617 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 618 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023. 619
- 620 Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, Zhiheng LYU, Kevin Blin, Fer-621 nando Gonzalez Adauto, Max Kleiman-Weiner, Mrinmaya Sachan, and Bernhard Schölkopf. CLad-622 der: A benchmark to assess causal reasoning capabilities of language models. In Thirty-seventh 623 Conference on Neural Information Processing Systems, 2023a. URL https://openreview. 624 net/forum?id=e2wtjx0Yqu. 625
- Zhijing Jin, Jiarui Liu, LYU Zhiheng, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona T Diab, and Bernhard Schölkopf. Can large language models infer causation from correlation? In The Twelfth International Conference on Learning Representations, 2023b. 628
- 629 Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas 630 Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) 631 know what they know. arXiv preprint arXiv:2207.05221, 2022. 632
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott 633 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. 634 arXiv preprint arXiv:2001.08361, 2020. 635
- 636 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large 637 language models are zero-shot reasoners. Advances in neural information processing systems, 35: 638 22199-22213, 2022.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, 640 Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. vllm, 2024. URL https://github.com/ 641 vllm-project/vllm. 642
- 643 Martha Lewis and Melanie Mitchell. Using counterfactual tasks to evaluate the generality of 644 analogical reasoning in large language models. arXiv preprint arXiv:2402.08955, 2024. 645
- Changmao Li and Jeffrey Flanigan. Task contamination: Language models may not be few-shot 646 anymore. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 647 18471-18480, 2024.

648 649 650	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re.
651	Diana Acosta-Navas Drew Arad Hudson Eric Zelikman Esin Durmus Faisal Ladhak Frieda
652	Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshay Santhanam, Laurel Orr, Lucia Zheng, Mert
653	Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter
654	Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli,
655	Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li,
656	Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models. Transactions
657	on Machine Learning Research, 2023. ISSN 2835-8856. URL https://openreview.net/
658	forum?id=i04LZibEqW. Featured Certification, Expert Certification.
659	
660	//about_meta_com/news/2023/07/llama=2/
661	// about . meta . com/ news/ 2023/ 07/ 11ama 2/.
662 663	Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024a. URL https://ai.meta.com/blog/meta-llama-3/.
664	
665	Ilama3/blob/main/MODEL CARD.md.
666	
667	Mistral-AI-Team. Announcement mistral ai 7b, 2023. URL https://mistral.ai/news/
668	announcing-mistral-7b/.
669	
670	Mistral-Al-Team. Announcement mixtral ai /x22b, 2024a. URL https://mistral.ai/news/
671	mixtral-8x22b/.
672	Mistral-AI-Team Announcement mixtral ai 7x8b 2024b JIRL https://mistral_ai/news/
673	mixtral-of-experts/.
674	mano and the trip of the first state of the
675	Mistral-AI-Team. Mistral ai models api versioning, 2024c. URL https://docs.mistral.ai/
676	getting-started/models/#api-versioning.
677	Misteral ALTrans. As lange Insistent ail frontion of in source hands 2024d LIDL hat have a finance in the second
678	pi/news/mistral_large/
679	ai/news/mistiai iaige/.
680 681	Melanie Mitchell. How do we know how smart ai systems are? Science, 381(6654):eadj5957, 2023.
682	Mosaic. Introducing dbrx: A new state-of-the-art open llm. URL https://www.databricks.
683	com/blog/introducing-dbrx-new-state-art-open-llm.
684	
685	Arsenii Kirillovich Moskvichev, Victor Vikram Odouard, and Melanie Mitchell. The conceptarc
686	machine learning research 2023
687	machine learning research, 2025.
688	<b>OpenAI.</b> Gpt-4 turbo and gpt-4 model docs., a. URL https://platform.openai.com/
689	docs/models/gpt-4-turbo-and-gpt-4.
690	
691	OpenAI. Models - openai gpt 3.5 turbo docs, b. URL https://platform.openai.com/
692	docs/models/gpt-3-5-turbo.
693	OpenAL Models - openai ant 3.5 turbo undate a LIRI https://openai.com/index/
694	new-embedding-models-and-api-updates/
695	non embedding modele and api apadees/.
696 697	OpenAI. Introducing chatgpt, 11 2022. URL https://openai.com/blog/chatgpt.
698	OpenAL Cost to model does 2024a LIPL https://platform.openai.com/deac/
699	models/gpt=/o
700	MOROTO, Abr in.
701	OpenAI. Announcement: Hello gpt-40, 2024b. URL https://openai.com/index/ hello-gpt-40/.

703 Automatically correcting large language models: Surveying the landscape of diverse automated 704 correction strategies. Transactions of the Association for Computational Linguistics, 12:484–506, 705 2024. 706 Jacob Pfau, William Merrill, and Samuel R Bowman. Let's think dot by dot: Hidden computation in transformer language models. arXiv preprint arXiv:2404.15758, 2024. 708 709 Sundar Pichai and Demis Hassabis. Introducing gemini: Google's most capable ai model yet, 2023. 710 URL https://blog.google/technology/ai/google-gemini-ai/. 711 Sundar Pichai Hassabis. and Demis Our next-generation model: Gem-712 https://blog.google/technology/ai/ 1.5, 2024. URL ini 713 google-gemini-next-generation-model-february-2024/. 714 Karl Raimund Popper. The logic of scientific discovery. 1934. 715 716 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 717 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 718 models from natural language supervision. In International Conference on Machine Learning, pp. 719 8748-8763. PMLR, 2021. 720 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi 721 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text 722 transformer. Journal of machine learning research, 21(140):1-67, 2020. 723 724 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste 725 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 726 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530, 2024. 727 728 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WINOGRANDE: an 729 adversarial winograd schema challenge at scale, 2019. 730 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam 731 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, 732 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. 733 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda 734 Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan 735 Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, 736 Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, 738 Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut 739 Erdem, Ayla Karakas, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, 740 Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk 741 Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, Cesar Ferri, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta 742 Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt, Christopher D 743 Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, 744 Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, 745 Dan Roth, C. Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle 746 Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David 747 Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz 748 Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho 749 Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad 750 Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà, 751 Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan 752 Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio 754 Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-Lopez, Gregor Betz, Guy Gur-Ari, Hana 755 Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar,

Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang.

756 Henry Francis Anthony Shevlin, Hinrich Schuetze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon 758 Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocon, 759 Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason 760 Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan 761 Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg 762 Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, 764 Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel, 765 Kevin Omondi, Kory Wallace Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, 766 Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, 767 Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, 768 Ludwig Schmidt, Luheng He, Luis Oliveros-Colón, Luke Metz, Lütfi Kerem Senel, Maarten 769 Bosma, Maarten Sap, Maartje Ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, 770 Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez-Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás 771 Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael Andrew 772 Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michael 773 Swędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac 774 Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, 775 Mukund Varma T, Nanyun Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover, 776 Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas 777 Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah 778 Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, 779 Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter W Chang, Peter Eckersley, Phu Mon Htut, 781 Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing 782 Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe 783 Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan Le Bras, Rosanne 784 Liu, Rowan Jacobs, Rui Zhang, Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan 785 Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, 786 Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel Stern Schoenholz, 787 Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, 788 Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, 789 Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane 790 Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima Shammie Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, 793 Stella Biderman, Stephanie Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, 794 Tariq Ali, Tatsunori Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, 796 Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu, 798 Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout 799 Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, 800 Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, 801 Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, 802 Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating 803 the capabilities of language models. Transactions on Machine Learning Research, 2023. ISSN 804 2835-8856. URL https://openreview.net/forum?id=uyTL5Bvosj.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

809

TogetherAI. Togetherai api, 2024. URL https://docs.together.ai/.

834

835

836

847

848

849

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
  Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
  efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay 814 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-815 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, 816 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, 817 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel 818 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, 819 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, 820 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, 821 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh 822 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, 823 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. 824 arXiv, 7 2023b. URL https://arxiv.org/abs/2307.09288. 825
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
   Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
   models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=1PL1NIMMrw.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*, 2023.
- Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek, Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas, and Yoon Kim. Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks. *arXiv e-prints*, pp. arXiv–2307, 2023.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
  Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
  - Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019. URL https://arxiv.org/abs/1905.07830.
- Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming
  Lu. Self-contrast: Better reflection through inconsistent solving perspectives. *arXiv preprint arXiv:2401.02009*, 2024.
- Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models
   are not robust multiple choice selectors. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=shr9PXz7T0.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.
- Benny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=WZH7099tgfM.

864 865 866	Yongshuo Zong, Tingyang Yu, Ruchika Chavhan, Bingchen Zhao, and Timothy Hospedales. Fool your (vision and) language model with embarrassingly simple permutations. In <i>Forty-first International Conference on Machine Learning</i> , 2024
867	Conjerence on Machine Learning, 2024.
868	
869	
870	
871	
872	
873	
874	
875	
876	
877	
878	
879	
880	
881	
882	
883	
884	
885	
886	
887	
888	
889	
890	
891	
892	
893	
894	
895	
896	
897	
898	
899	
900	
901	
902	
903	
904	
905	
906	
907	
908	
909	
910	
911	
912	
913	
914	
915	
916	
917	

#### Supplementary.

#### A ADDITIONAL DETAILS ON PERFORMED EXPERIMENTS

Here we give further details on the procedures around the executed experiments.

A.1 MODELS SELECTED FOR EXPERIMENTS

To provide overview over origin of core tested models used for the AIW experiments, we list those in Suppl. Tab. 1. All tested models use same default inference hyperparameters, T = 0.1, top-p = 1.0 (we executed control experiments to check that various settings do not change the main pattern in the observed behavior). The output was limited to 2048 tokens, and as evident from Suppl. Fig. 22, most observed responses stayed well below this limit.

Table 1: Names, origin and versioning of core test models used in the experiments.

937 938	Name	Origin	Released	Open Weights	Sources
939	GPT-40-2024-05-13	OpenAI	13.05.2024	No	(Achiam et al., 2023; OpenAI, 2024a;b)
941	GPT-4-turbo-2024-04-09	OpenAI	09.04.2024	No	(Achiam et al., 2023; OpenAI, a)
942 943	GPT-4-0125-preview	OpenAI	25.01.2024	No	(Achiam et al., 2023; OpenAI, a)
944	GPT-4-0613	OpenAI	13.06.2023	No	(Achiam et al., 2023; OpenAI, a)
945	GPT-3.5-turbo-0125	OpenAI	24.01.2024	No	(OpenAL 2022: b:c)
946	Claude-3-5-sonnet-20240620	Anthropic	21.06.2024	No	(Anthropic, 2024b)
947	Claude-3-opus-20240229	Anthropic	04.03.2024	No	(Anthropic, 2024a;c)
948	Claude-3-sonnet-20240229	Anthropic	04.03.2024	No	(Anthropic, 2024a;c)
949	Claude-3-haiku-20240307	Anthropic	04.03.2024	No	(Anthropic, 2024a;c)
950	Gemini 1.0 Pro	Google	06.12.2023	No	(Pichai & Hassabis, 2023; Team et al., 2023)
951 952	Gemini 1.5 Pro	Google	16.02.2024	No	(Pichai & Hassabis, 2024; Reid et al., 2024)
953	gemma-7b-it	Google	05.04.2024 (v1.1)	Yes	(Google, 2024a;b)
954	gemma-2b-it	Google	05.04.2024 (v1.1)	Yes	(Google, 2024a;b)
055	Mistral-large-2402	Mistral AI	26.02.2024	No	(Mistral-AI-Team, 2024c;d)
955	Mistral-medium-2312	Mistral AI	23.12.2023	No	(Mistral-AI-Team, 2024c;d)
956	Mistral-small-2402	Mistral AI	26.02.2024	No	(Mistral-AI-Team, 2024c;d)
957	open-mixtral-8x22b-instruct-v0.1	Mistral AI	17.04.2024	Yes	(Mistral-AI-Team, 2024c;a)
958	open-mixtral-8x7b-instruct-v0.1	Mistral AI	11.12.2023	Yes	(Mistral-AI-Team, 2024c;b)
959 960	open-mistral-7b-instruct-v0.2	Mistral AI	11.12.2023	Yes	(Jiang et al., 2023; Mistral-AI- Team, 2024c; 2023)
900	Command R+	Cohere	04.04.2024	Yes	(Cohere, 2024a;b)
961	Dbrx Instruct	Mosaic	27.03.2024	Yes	(Mosaic)
962 963	Llama 2 70B Chat	Meta	18.07.2023	Yes	(Meta, 2023; Touvron et al., 2023b)
964	Llama 2 13B Chat	Meta	18.07.2023	Yes	(Meta, 2023; Touvron et al., 2023b)
965 966	Llama 2 7B Chat	Meta	18.07.2023	Yes	(Meta, 2023; Touvron et al., 2023b)
967	Llama 3 70B Chat	Meta	18.04.2024	Yes	(Meta, $2024a$ :b)
968	Llama 3 8B Chat	Meta	18.04.2024	Yes	(Meta, 2024a;b)
969	Qwen 1.5 1.8B - 72B Chat	Alibaba	04.02.2024	Yes	(Bai et al., 2023; Alibaba, 2024a)
970 971	Qwen 2 72B Instruct	Alibaba	07.06.2024	Yes	(Alibaba, 2024b)

#### 972 A.2 EVALUATING MODEL RESPONSES 973

974 The formatting instruction makes it possible to extract for each prompting trial whether a model has 975 provided a correct answer to the AIW problem posed in the input. We can interpret then any number 976 n of collected responses as executing n trials given a particular prompt for a given model (n - number 977 of Bernoulli trials), observing in each *i*-th trial a Bernoulli variable  $X_i = \{0, 1\}$ . We interpret the number of correct responses  $X = \sum_{i} X_{i}$  as random variable following a Beta-Binomial distribution 978 with unknown probability p of correct response that we also treat as random variable that comes from 979 980 a Beta distribution, i.e.  $p \sim Beta(\alpha, \beta)$ , where  $\alpha$  and  $\beta$  are parameters of the Beta distribution. To obtain plots showing correct response ratios, we would like to estimate Beta distribution underlying 981 p, and for that, we first estimate the mean of p and its variance from the collected observations. To 982 estimate  $\hat{p}$ , we use the formula for estimating the mean of p for a binomial distribution:  $\hat{p} = X/n$  (i.e. 983 as a proportion of successes). We can report the estimate  $\hat{p}$  as the estimate of the correct response rate 984 of a given model and also, compare the correct response rates of various tested models. Moreover, 985 we can estimate the variance of the probability of a correct response by using the following formula: 986

987

$$\operatorname{var}\left(\frac{1}{n}\sum_{i=1}^{n}X_{i}\right) = \frac{1}{n^{2}}\sum_{i=1}^{n}\operatorname{var}(X_{i}) = \frac{n\operatorname{var}(X_{i})}{n^{2}} = \frac{\operatorname{var}(X_{i})}{n} = \frac{p(1-p)}{n}$$
(1)

993

994 995

996

997

998

The estimates of the variance and the standard deviation of p can be thus obtained by using  $\hat{p}$  as  $\frac{\hat{p}(1-\hat{p})}{n}$  and  $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$  respectively. Using the estimated variance and mean of p, we can use the following relations for the variance:  $\left(\sigma^2 = \frac{n\alpha\beta(\alpha+\beta+n)}{(\alpha+\beta)^2(\alpha+\beta+1)}\right)$  and the mean  $\left(\mu = \frac{\alpha}{\alpha+\beta}\right)$  in order to obtain  $\alpha$  and  $\beta$  parameters for the Beta distribution. To simulate data for the plots, we draw N random samples corresponding to correct and incorrect responses using the estimated distribution of p and obtain the plots showing performance on the task for various models of interest as a full distribution of the respective p.

999 1000 1001

#### **B PROMPT TYPES AND VARIATIONS**

1002 1003

For testing the model dependence on input prompt type as well as robustness against problem variations when solving AIW and AIW Light problems, we used three main prompt types - STANDARD (original prompt with answer formatting instructions), THINKING (prompt that encourages thinking with answer formatting instructions) and RESTRICTED (prompt that instructs model to output only formatted answer and nothing else). THINKING v2 prompt type is a minor variation of THINKING type that just adds "step by step" after already existing "think carefully" phrasing (control experiments show that THINKING and THINKING v2 are equivalent in terms of observed performance, so we use both interchangeably, Suppl. Fig. 12b).

For testing the models' robustness to problem perturbations, we try different variations of main AIW problem (AIW Variations 1-4, see Sec. 2, Suppl. Tab. 2), where we keep the same problem structure while varying numbers of brothers and sisters and their mentioning order within the sentence. Those variations are made intentionally in such a way that they do not affect problem structure or its difficulty and thus should not affect how models cope with the problem.

We employ a further AIW version - AIW Light - as control to test whether models are able to deal
with various aspects of original AIW problem, eg handling the specific relational family structure
frame, or executing elementary arithmetic operations necessary to solve AIW problem. See Sec. 2
for details on the AIW Light design.

See Suppl. Tab. 2 (for AIW problem) and Suppl. Tab. 3, 4, 5 (for AIW Light problems) for examples with full prompt versions for each presented problem and its variations<sup>2</sup>.

<sup>1024</sup> 1025

<sup>&</sup>lt;sup>2</sup>All prompts and their IDs available at https://anonymous.4open.science/r/AITW\_ anonymous-69A6/prompts/prompts.json

Table 2: AIW main variations, prompt types and correct answers overview.

Var.	Prompt	Type/Answer	ID
1	Alice has 3 brothers and she also has 6 sisters. How many sisters does Alice's	STANDARD /	55
	brother have? Solve this problem and provide the final answer in following	7	
	form: "### Answer: ".		
1	Alice has 3 brothers and she also has 6 sisters. How many sisters does Alice's	THINKING /	57
	brother have? Before providing answer to this problem, think carefully and	7	
	double check the path to the correct solution for any mistakes. Provide then the		
	final answer in following form: "### Answer: ".		
1	Alice has 3 brothers and she also has 6 sisters. How many sisters does Alice's	RESTRICTED	53
-	brother have? To answer the question, DO NOT OUTPUT ANY TEXT EX-	17	
	CEPT following format that contains final answer: "### Answer: ".		
2	Alice has 2 sisters and she also has 4 brothers. How many sisters does Alice's	STANDARD /	56
-	brother have? Solve this problem and provide the final answer in following	3	
	form: "### Answer: "	5	
2	Alice has 2 sisters and she also has 4 brothers. How many sisters does Alice's	THINKING /	55
-	brother have? Refore providing answer to this problem think carefully and	3	50
	double check the path to the correct solution for any mistakes. Provide then the	5	
	final answer in following form: "### Answer: "		
2	Alice has 2 sisters and she also has 4 brothers. How many sisters does Alice's	RESTRICTED	5/
4	hrother have? To answer the question DO NOT OUTDUT AND TEVT EV	/3	54
	CEPT following format that contains final answer: "### Answer: "	15	
3	Alice has A sisters and she also has 1 brother. How many sisters does Alice's		6
5	hrother have? Solve this problem and provide the final answer in following	5 IANDARD /	0.
	form: "### A newer: "	J	
3	Alice has A sisters and she also has 1 brother. How many sisters does Alice's	THINKING /	6,
3	After have? Refere providing answer to this problem, think corefully and	5	04
	double sheak the path to the correct solution for any mistelies. Provide there the	J	
	final answer in fallowing forms "#### A reserver"		
2	Inal answer in following form: "### Answer:".	DECTDICTED	~
3	Affice has 4 sisters and sne also has 1 prother. How many sisters does Affice's	KESIKICIED	03
	oroner nave? To answer the question, DU NUT OUTPUT ANY TEXT EX-	13	
4	CEPT TOHOWING TOFINAL UNAL CONTAINS TINAL ANSWER: "### Answer:".		61
4	Affice has 4 brothers and she also has 1 sister. How many sisters does Affice's	STANDARD /	65
	brother nave? Solve this problem and provide the final answer in following	2	
	form: "### Answer: ".		_
4	Alice has 4 brothers and she also has 1 sister. How many sisters does Alice's	THINKING /	/(
	brother have? Before providing answer to this problem, think carefully and	2	
	double check the path to the correct solution for any mistakes. Provide then the		
	tinal answer in following form: "### Answer: ".		_
4	Alice has 4 brothers and she also has 1 sister. How many sisters does Alice's	RESTRICTED	71
	brother have? To answer the question, DO NOT OUTPUT ANY TEXT EX-	/2	
	CEPT following format that contains final answer: "### Answer: ".		
4	Alice has 4 brothers and she also has 1 sister. How many sisters does Alice's	CONFIDENCE	11
	brother have? Solve the problem by taking care not to make any mistakes. Ex-	/ 2	
	press your level of confidence in the provided solution as precisely as possible.		
3	Alice has 4 sisters and she also has 1 brother. How many sisters does Alice's	SCIENTIST /	40
	brother have? To solve the problem, approach it as a very intelligent, accurate	5	
	and precise scientist capable of strong and sound reasoning. Provide the solution		
	to the problem by thinking step by step, double checking your reasoning for		
	any mistakes, and based on gathered evidence, provide the final answer to the		
	problem in following form: "### Answer: ".		

Var.	Prompt	Type/Answer	ID
1	Alice has 3 brothers and she also has 4 sisters. How many siblings does Alice have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer:"	THINKING v2/7	277
2	Alice has 2 sisters and she also has 1 brother. How many siblings does Alice have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer:"	THINKING v2/3	278
3	Alice has 4 sisters and she also has 1 brother. How many siblings does Alice have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer: "	THINKING v2 / 5	279
4	Alice has 1 brother and she also has 1 sister. How many siblings does Alice have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer:".	THINKING v2/2	280
C N P	IODEL PERFORMANCE AND BEHAVIOR ON AIW AND AIW LIGH ROBLEM	ΙT	
2024). For the Fig. 8. Suppl. Suppl. Compa For the see Su	e full overview of average correct response rate including models that score zero. For the statistics on number of trials conducted for each model and each prom Fig. 21. For the statistics on the average output length across models and promp Fig. 22. For models' behavior on RESTRICT prompt types, see Suppl. Fig. 11. rison of THINKING v2 prompt type to THINKING and STANDARD, see see Sup e control of observed strong fluctuations being the same independent of employed p ppl. Fig. 9	o, see Suppl. pt type, see ot types, see For control ppl. Fig. 12. prompt type,	
C.1	BOOSTING BY REDUNDANT INFORMATION AND PERSISTING FLUCTUATIONS: FEMALE POWER BOOST	ALICE	
We rep Tab. 6 wherea Suppl. AIW o being solving	ort in Sec. 3.2, Fig. 6 how introducing fully redundant information "Alice is female for full prompts) causes increase of average correct response rates across AIW va as strong fluctuations across variations remain. Here we visualize the observed Fig. 10. We see that for most models that had some non-negligible correct respo riginal average correct response rate os significantly boosted, despite the provided fully redundant. This change in performance caused by information irrelevant f g again hints on deficits in generalization and basic reasoning across the models.	" (see Suppl. riations 1-4, increase in nse rates on information For problem	
C.2	FREQUENCY DISTRIBUTION OF NATURAL NUMBERS ON OUTPUT AND DOMIN WRONG RESPONSES.	ANCE OF	
To she with A for AΓ	d more light on modes of correct or wrong responses provided by the models when IW problem variations, we show here frequency distribution for natural numbers o W variations with higher and lower correct response rates.	n confronted n the output	
As evi are oft Fig. 13	dent from the plots, in higher performance AIW variations (Suppl. Fig. 14), domi en positioned on correct answer C=M+1, while for lower performance AIW variat 3), dominant peaks fall on wrong answer M. Further, for weaker models, distribution	nants peaks ions (Suppl. on broadens.	

#### Table 3: AIW Light Arithmetic Siblings variations





Figure 9: Strong fluctuations on AIW problem variations (a color per each variation 1-4) appear to same extent independent of employed prompt type, on example of GPT-4 (gpt-4-0613). AIW variations 1-4 correspond to different instantiations of numbers N, M for brothers and sisters in the same AIW problem template. Varying numbers should not affect problem solution at all, as it does not affect problem structure and its difficulty. However, correct response rate varies strongly depending on the variation. E.g., it drops close to 0 for variation 3, while going up to 1 for variation 4. This observation is consistent for different prompt types - STANDARD, THINKING and RESTRICTED (from left to right). Full input for each single trial has a form <instantiated-template> <prompt-type>, where <instantiated-template> is template with substituted numbers instantiating one of AIW Variations 1-4 and *<prompt-type>* contains instructions corresponding to one of 3 prompt types. Lack of robustness to irrelevant variations of such a simple problem points to severe generalization deficits. 

Table 4: AIW Light Family variations

Var.	Prompt	Type/Answer	ID
1	Alice has 7 brothers and she also has 3 sisters. How many brothers does Alice's sister have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer:".	THINKING v2 / 7	271
2	Alice has 4 sisters and she also has 3 brothers. How many brothers does Alice's sister have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer:".	THINKING v2/3	272
3	Alice has 2 sisters and she also has 5 brothers. How many brothers does Alice's sister have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer: ".	THINKING v2 / 5	273
4	Alice has 2 brothers and she also has 3 sisters. How many brothers does Alice's sister have? Before providing answer to this problem, think carefully step by step and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: "### Answer:".	THINKING v2/2	274

Tuote of the angle the internet of the office the angle of the office of	Table 5: AIW	Light	Arithmetic	Total	Girls	variations
--	--------------	-------	------------	-------	-------	------------

Var.	Prompt	Type/Answer	ID
1	Alice has 6 sisters and she also has 3 brothers. How many girls are there in	THINKING	343
	total? Before providing answer to this problem, think carefully step by step and	v2/7	
	double check the path to the correct solution for any mistakes. Provide then the		
	final answer in following form: "### Answer:".		
2	Alice has 2 sisters and she also has 4 brothers. How many girls are there in	THINKING	344
	total? Before providing answer to this problem, think carefully step by step and	v2/3	
	double check the path to the correct solution for any mistakes. Provide then the		
_	final answer in following form: "### Answer:".		
3	Alice has 4 sisters and she also has 1 brother. How many girls are there in total?	THINKING	345
	Before providing answer to this problem, think carefully step by step and double	v2/5	
	check the path to the correct solution for any mistakes. Provide then the final		
	answer in following form: "### Answer: ".		
4	Alice has 1 sister and she also has 4 brothers. How many girls are there in	THINKING	346
	total? Before providing answer to this problem, think carefully step by step and	v2/2	
	double check the path to the correct solution for any mistakes. Provide then the		
	final answer in following form: "### Answer:".		

Table 6: AIW Alice Female Power Boost and AIW Original, variations 1-4, THINKING v2 prompt type

1	263	
1	264	

Var.	Prompt	Type/Answer	ID
1	Alice is female and has 3 brothers and she also has 6 sisters. How many sisters	FEMALE	193
	does Alice's brother have? Before providing answer to this problem, think	BOOST / 7	
	carefully step by step and double check the path to the correct solution for any		
	mistakes. Provide then the final answer in following form: "### Answer: ".		
1	Alice has 3 brothers and she also has 6 sisters. How many sisters does Alice's	NO BOOST /	205
	brother have? Before providing answer to this problem, think carefully step by	7	
	step and double check the path to the correct solution for any mistakes. Provide		
•	then the final answer in following form: "### Answer: ".		
2	Alice is female and has 2 sisters and she also has 4 brothers. How many sisters	FEMALE	197
	does Alice's brother have? Before providing answer to this problem, think	BOOST/3	
	carefully step by step and double check the path to the correct solution for any		
n	Alice has 2 sisters and she also has 4 brothers. How many sisters does Alice's	NO DOOST /	206
Ζ	Ance has 2 sisters and she also has 4 brothers. How many sisters does Ance s	NO BOOST /	200
	step and double check the path to the correct solution for any mistakes. Provide	5	
	then the final answer in following form: "### Answer: "		
3	Alice is female and has 4 sisters and she also has 1 brother. How many sisters	FEMALE	189
5	does Alice's brother have? Before providing answer to this problem think	BOOST / 5	107
	carefully step by step and double check the path to the correct solution for any	2000170	
	mistakes. Provide then the final answer in following form: "### Answer: ".		
3	Alice has 4 sisters and she also has 1 brother. How many sisters does Alice's	NO BOOST /	187
	brother have? Before providing answer to this problem, think carefully step by	5	
	step and double check the path to the correct solution for any mistakes. Provide		
	then the final answer in following form: "### Answer: ".		
4	Alice is female and has 4 brothers and she also has 1 sister. How many sisters	FEMALE	190
	does Alice's brother have? Before providing answer to this problem, think	BOOST / 2	
	carefully step by step and double check the path to the correct solution for any		
	mistakes. Provide then the final answer in following form: "### Answer: ".		
4	Alice has 4 brothers and she also has 1 sister. How many sisters does Alice's	NO BOOST /	188
	brother have? Before providing answer to this problem, think carefully step by	2	
	step and double check the path to the correct solution for any mistakes. Provide		
	then the final answer in following form: "### Answer: ".		



Figure 10: AIW "Alice Female Power Boost" version. Average correct response rate (measured across AIW variations 1-4) increases after addition of entirely redundant information "Alice is female" (pronoun "she" already fully indicates the gender in original AIW). Thinking v2 prompt type is used for both AIW versions. See also Fig. 6 for persisting strong fluctuations across variations 1-4.

1325

covering more numbers (eg in Llama 3 8b), while for better performers, responses concentrate on M and M+1, peaking on correct or wrong answer on depending on AIW variation. Remarkably, for lower performance AIW variations (Suppl. Fig. 13), performance cannot be rescued by major voting or by similar ensemble like strategies, as peaks on wrong response numbers dominate clearly peaks on numbers for correct responses, which would still correspond to committing wrong answer when performing majority voting.

For the AIW Light problem versions used in control experiments, we observe as expected clear dominant peaks on the numbers corresponding for correct responses across all tested models (Suppl. Fig. 15, 16), as AIW Light problems are successfully solved across all their variations.

We note that distribution characteristics, eg concentration on numbers around the correct answer, height of the peaks, can be a further signature that reflects model's capability to handle the problem. More capable models retain dominant peaks on number corresponding to correct answer with smaller peaks on neighboring numbers, while weak models have large peaks on numbers corresponding to wrong answers or in general broad distribution across all natural numbers below 10. Computing scores from distribution shape can thus also enable model ranking.

1341 1342

C.3 STANDARDIZED BENCHMARKS FAILURE.

In Section 3.2, we observe failure of standardized reasoning benchmarks to properly reflect generalization and basic reasoning skills of SOTA LLMs by noting significant disparity between the model's
performance on the AIW problem and the outcomes on conventional standardized benchmarks,
taking MMLU as representative examples. Here, we confirm this finding on further standardized
reasoning benchmarks like MATH, ARC-c, GSM8K and Hellaswag (Suppl. Tab. 7). We provide
plots visualizing failure of these standardized benchmarks, reflected in strong mismatch between high
benchmark scores reported by many models and the low correct response rates they obtain on AIW



(a) Correct response rates for RESTRICTED prompt (b) Strong fluctuations across AIW variations 1-4, RE-type, averaged across AIW var 1-4
 STRICTED prompt type

Figure 11: Correct response rates on RESTRICTED prompt type. The prompt type enforcing to output only final answer without any further text was used as further control. (a) Correct response rates averaged over variations 1-4 resemble behavior with STANDARD and THINKING types, while looking at fluctuations across variations 1-4 in (b) reveals stronger models' lack of robustness compared to other prompt types (see for comparison Fig. 2). We thus used THINKING prompt types across main experiment not to put models into disadvantage on AIW testing.



1392 1393

1394

1375

(a) Correct response rates THINKING v2 vs. STAN-(b) Correct response rates THINKING v2 vs. THINK-DARD prompt type, averaged across AIW var 1-4 ING prompt type, averaged across AIW var 1-4

Figure 12: Control comparison of correct response rates averaged across AIW variations 1-4. (a) THINKING v2 vs. STANDARD, (b) THINKING v2 vs. THINKING prompt types. THINKING provides better average correct response rates for tested models. We thus used THINKING prompt types for main and control experiments to ensure tested models are not disadvantaged on AIW problem. THINKING and THINKING v2 show highly similar behavior across tested models (b) and can be used interchangeably (THINKING v2 only difference to THINKING is the explicit phrasing "step by step", Suppl. Tab. 6)

1402



Figure 13: Frequency distribution of output numbers in models' responses. Shown are numerical 1415 outputs for AIW Variation 3, THINKING prompt type (prompt ID 64), that has correct answer 1416 C=M+1=5, with M=4 number of sisters of Alice. For this AIW variation, models have low perfor-1417 mance (see also Figure D.). Correspondingly, peaks are on the dominant wrong response, R=M=4. 1418 For this low performance variation, performance cannot be rescued by majority voting or other simple 1419 ensembling strategies, as also for better performing models like GPT-40, there are dominant peaks on 1420 wrong numbers that would overrule less dominant peaks for correct numbers. Weaker models, eg 1421 Llama 3 8B, show also broader distribution. Distributions were computed over 60 trials executed for 1422 each model, taken from original collected responses data. 1423



1436 Figure 14: Frequency distribution of output numbers in models' responses. Shown are numerical 1437 outputs for AIW Variation 4, THINKING prompt type (prompt ID 70), that has correct answer 1438 C=M+1=2, with M=1 number of sisters of Alice. For this AIW variation, models have higher 1439 performance (see also Figure D.). Correspondingly, peaks for better performing models (eg GPT-40, 1440 GPT-4, Claude Opus 3) are on the dominant correct response, R=M+1=2. For models with worse 1441 performance, peaks are on the dominant wrong response, R=M=1. For weaker models, eg Llama 3 8B, also broader distribution over numbers appears, with further wrong clear peaks that are further 1442 away from C=M+1 (eg M=4). The distribution shape and peaks nature can be thus used as signature 1443 of model's capability to handle the problem, also allowing model ranking dependent on peak types 1444 and distribution sharpness. Distributions were computed over 60 trials executed for each model, taken 1445 from original collected responses data. 1446

1424

1440

(which in some cases is 0 for models with high standardized benchmark scores), in Figures 17, 20, 18, 19.

We see thus that standardized benchmarks fail to properly reflect true model capabilities to generalize
and reason - the majority of the tested models score high on standardized benchmarks, suggesting
strong function, while showing extreme low correct response rates on simple AIW problem. Many
of the models with high scores on standardized benchmarks cannot solve AIW problem a single
time (e.g. Command R+ is unable to solve a single AIW problem instance, see Suppl. Tab. 7). This
discrepancy refutes the claim of standardized benchmarks to measure correctly current models' core
functionality.

Table 7: Performance of tested models on MMLU, Hellaswag, ARC-c, GSM8k and AIW problems.
 Correct response rate averaged across AIW variations 1-4, across STANDARD and THINKING prompt types.

Model	MMLU	Hellaswag	ARC-c	GSM8k	Correct av- erage resp. rate
gpt-4o-2024-05-13	0.89	-	-	-	0.65
claude-3-opus-20240229	0.87	95.40	96.40	95.00	0.43
gpt-4-0613	0.86	95.30	96.30	92.00	0.37
llama-2-70b-chat	0.64	85.90	64.60	56.80	0.30
llama-2-7b-chat	0.55	77.10	43.20	25.40	0.13
dbrx-instruct	0.74	88.85	67.83	67.32	0.11
gpt-4-turbo-2024-04-09	0.80	-	-	-	0.10
llama-3-8b-chat	0.67	78.55	60.75	79.60	0.05
llama-3-70b-chat	0.80	85.69	71.42	93.00	0.05
qwen1.5-1.8b-chat	0.46	46.25	36.69	38.40	0.05
gemma-2b-it	0.38	71.40	42.10	17.70	0.04
llama-2-13b-chat	0.66	80.70	48.80	77.40	0.03
qwen1.5-4b-chat	0.56	51.70	40.44	57.00	0.02
claude-3-sonnet-20240229	0.79	89.00	93.20	92.30	0.01
mistral-large-2402	0.81	89.20	94.20	81.00	0.01
gpt-3.5-turbo-0125	0.70	85.50	85.20	57.10	0.01
gemini-pro	0.72	84.70	-	77.90	0.01
open-mixtral-8x22b	0.78	89.08	72.70	82.03	0.01
open-mistral-7b	0.64	84.88	63.14	40.03	0.01
qwen1.5-7b-chat	0.62	59.38	52.30	62.50	0.01
claude-3-haiku-20240307	0.75	85.90	89.20	88.90	0.00
open-mixtral-8x7b	0.72	87.55	70.22	61.11	0.00
command-r-plus	0.76	88.56	70.99	70.74	0.00
qwen1.5-14b-chat	0.69	63.32	54.27	70.10	0.00
gemma-7b-it	0.54	81.20	53.20	46.40	0.00
qwen1.5-72b-chat	0.77	68.37	65.36	79.50	0.00
qwen1.5-32b-chat	0.75	66.84	62.97	77.40	0.00



Figure 15: Frequency distribution of output numbers in models' responses. Shown are numerical outputs for AIW Light Family, Variation 3, THINKING prompt type (prompt ID 273), that has correct answer C=5 (number of Alice's brothers). For this AIW Light version, all models have high performance. Correspondingly, peaks are on the dominant correct response, R=5. However also here, weaker models like Llama 3 8B show broader distribution with non-vanishing peaks besides the correct response (eg R=0, R=2) hinting on their weaker capabilities to deal robustly with the problem. Distributions were computed over 60 trials executed for each model.



Figure 16: Frequency distribution of output numbers in models' responses. Shown are numerical outputs for AIW Light Arithmetic, Variation 3, THINKING prompt type (prompt ID 279), that has correct answer C=5 (total number of Alice's siblings). For this AIW Light version, all models have high performance. Correspondingly, peaks are on the dominant correct response, R=5. However also here, weaker models like Llama 3 8B show broader distribution with non-vanishing peaks besides the correct response (eg R=4, R=6) hinting on their weaker capabilities to deal robustly with the problem. Distributions were computed over 60 trials executed for each model.

1548 1549

1550 1551

1552

1553

1554 1555

1556

1557 1558

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

#### D EXAMPLES OF CORRECT AND FAILED RESPONSES

We provide all collected model responses we obtained during this study in the collected\_responses folder in the AIW repo. Here we also showcase some correct and incorrect answers as an example (see Figs. 23, 26, 24, 25).

## E CONFABULATIONS AND OVERCONFIDENT TONE ACCOMPANYING WRONG ANSWERS

Overconfident tone. In ideal scenario, if LLM cannot correctly solve the AIW problem, it should at least be capable of expressing high uncertainty about the provided incorrect solution to the user. We used CONFIDENCE prompt type (see Suppl. Tab. 2) for AIW problem to see how confident tested models are in their wrong solutions.

From our experiments we can see that LLMs most of the time express high certainty even if their answers are completely wrong, thus mediating strong confidence (see Fig. 27). The models also use highly persuasive tone to argue for the expressed certainty and correctness of the provided wrong solutions, using words like "highly confident", "definitive answer", or "accurate and unambiguous".



Figure 17: Discrepancy between the AIW correct response rate and the MATH average score, indicating the limitation of standardized benchmark MATH in accurately assessing and comparing basic reasoning capabilities of models. Numerous models, such as Command R+, exhibit a stark contrast in performance, scoring zero on AIW while achieving high scores on MATH.



Figure 18: Limitation of the standardized benchmark GSM8k in accurately reflecting and comparing basic reasoning capabilities of models, as illustrated by the stark discrepancy between the AIW correct response rate and the GSM8k average score. Notably, the majority of tested models exhibit low performance on AIW problems while achieving relatively high scores on GSM8k, a graduate-level math benchmark for large language models. Among models with slightly better calibration are Claude Opus and GPT 4 that outperform other models on AIW, which coincides with their high GSM8k scores. Llama 2 70b also shows better calibration, where its modest AIW performance matches its modest GSM8k score. In contrast, models like Mistral Large, Gemini Pro, Dbrx Instruct, or Command R+, while scoring high on GSM8k, show breakdown on AIW (Command R+ has 0 correct response rate, Mistral Large and Gemini Pro 0.01, Dbrx Instruct 0.11, see also Suppl. Tab. 7) 



Figure 19: Limitation of the standardized benchmark Hellaswag in accurately assessing and comparing basic reasoning capabilities of models, as evidenced by the significant discrepancy between the AIW correct response rate and the Hellaswag average score.

We see also strong overconfidence expressed in multi-turn interactions with models, where user is insisting on solution provided being incorrect, and observe there high resistance of models to revise their decisions, which was already referred to as "stubbornness" in other works (Zhang et al., 2024) (see Suppl. Sec. F and also data provided in the AIW repo)

Confabulations. In our experiments we observe frequent tendency of those tested models that show strong reasoning collapse and produce frequent wrong answers for AIW problem to generate at the same time persuasive sounding pseudo-explanations to back up their incorrect answers. We term here such pseudo-explanations confabulations, and present a selection of those as examples.

Such confabulations can contain mathematical calculations or other logic-like expressions and operations that make little or absolutely no sense given the problem to be solved, see examples for Olmo-7B, Fig. 28 and Command R+, Fig. 30.

Further confabulations make use of various social and cultural norm specific context to argue for the posed problem to be inappropriate to solve or to provide non-sense arguments for various incorrect answers. There are many such examples that we have observed, we present here only a small selection.

1719 CodeLlama-70B-instruct for instance seems to be specifically prone to claim ethical or moral reasons
1720 for not addressing the problem correctly, in the presented example inventing out of nowhere a person
1721 with Down syndrome and then pointing out that question has to be modified to be addressed due to
1722 potential perpetuation of harm towards individuals or groups, which has nothing to do with original
1723 task, Fig. 29.

Another example are confabulations provided by Command R Plus. These confabulations use
concepts of gender identity such as non-binary gender or concepts related to inclusion or to cultural
context dependent family identification in the provided wrong reasoning leading to incorrect answers.
In the attempt to solve the problem, the model first fails to provide obvious common sense solution
and then goes on to describe potential scenarios where brothers and sisters may self-identify as



Figure 20: Failure of standardized benchmark ARC-c to properly reflect and compare model basic reasoning capabilities as shown by strong discrepancy between AIW correct response rate vs ARC-c average score.

1756 non-binary, although providing information on brothers and sisters in the problem usually means 1757 via common sense that those persons self-identify correspondingly to their known status as brother 1758 or sister (while Alice is clearly identified via "she" pronoun). Model thus clearly fails to grasp that 1759 problem structure has nothing to do with the social and cultural norms. The solutions derived by the 1760 model from considering those factors that are far beyond Occam's razor and common sense inherent 1761 to the simple AIW problem all lead to wrong answers and generate more confusion, while again 1762 keeping the persuasive tone that suggests that model is on some right path to provide the correct solutions (Fig. 31) 1763

For more illustrative examples, see the raw data on interactions with the models collected in AIW repo)

#### 1766 1767

#### 1768 F INABILITY TO REVISE WRONG SOLUTIONS

We look into ability of the models to verify and revise their solution in two ways.

1771 First, we observe in the collected data responses that contain examples of self-verification. Those 1772 can arise following from THINKING prompt that encourages to double-check the solution, or they 1773 appear by following customized prompts that request to produce different solutions and check which one is to prefer, or those that appear entirely unprompted (An example of a customized prompt that 1774 encourages to produce various solutions and evaluate those is "Look at the problem step by step and 1775 formulate 3 different solutions that come to different results. Then evaluate which solution seems 1776 to be the best and then come to a definitive final statement.", see also Fig. 30. In all those cases, 1777 we see only poor ability of the models the provide proper self-checks. In the examples we observed, 1778 self-verification provides longer narration, but does not lead to successful revision of wrong answers. 1779 Second, we looked into multi-turn interactions with the user and model, where it might be arguably 1780

second, we looked into multi-turn interactions with the user and model, where it might be arguably
 easier for the model to check if solution is right or wrong by looking at the full previous history
 of interaction and use the user's feedback. In such interactions, the model is prompted with AIW

High       High       Finite       Part of the second se	1783					
Treshold/P         Treshold/P         Treshold/P         Treshold/P         Treshold/P           1766         queri 1.5 / Xub /	1784		AIW Nu	umber of responses pe	r model	
No.         No. <th>1785</th> <th></th> <th>RESTRICTED</th> <th></th> <th>THINKING</th> <th></th>	1785		RESTRICTED		THINKING	
1100         0	1796	gwen1 5-72h-chat	5.70+02	3.6e+02	4.80+02	
1/16/1         gn 3.3 starb 0.12         3.4 m 0         2.7 m 0         2.7 m 0           1786         gp 4.4 m 0.2 9 4 4 6 0         3.6 m 0<	1700	open-mixtral-8x7b	3.4e+02	2.7e+02	3e+02	
1786         gp1-4eb/302-000         3.8-02         2.6-02         5.6-02           1789         cemara 6 tab         3.8-02         2.4-02         3.8-02         1.8-02           1790         cemara 6 tab         3.8-02         2.4-02         3.8-02         1.8-02           1791         fam3 2-786-784         3.8-02         3.8-02         3.8-02         3.8-02           1792         daub 3-cone 3000228         3.8-02         2.1-02         3.8-02         3.8-02           1795         calack 3-bab 2000228         3.8-02         2.1-02         2.8-02         3.8-02           1795         calack 3-bab 2000228         3.8-02         2.1-02         2.8-02         3.8-02           1795         calack 3-bab 200028         3.8-02         1.8-02         2.4-02         2.4-02           1796         calack 3-bab 200028         3.8-02         1.8-02         2.4-02         2.4-02           1796         calack 3-bab 2000028         3.8-02         1.8-02         2.4+02         2.4+02           1799         calack 3-bab 200002         1.8-02         2.4+02         2.4+02         2.4+02           1800         calack 3-bab 200002         1.8-02         2.4+02         2.4+02           180	1/8/	gpt-3.5-turbo-0125	3.3e+02	2.6e+02	2.7e+02	
1786         gen intra-b.         32-8-2         2.4-9/2         3-9-2         3-9-2           1790         conventioned in the second intervent in the second intervent inte	1788	gpt-4-turbo-2024-04-09	3.3e+02	2.6e+02	2.6e+02	
1790         convent-trains         3.44-02         2.44-02         2.84-92           1791         kmax-2-70-cm         3.34-70         2.8-72         3.34-78           1792         cdack-2-00-cm         3.34-70         2.8-72         3.34-78           1793         mono-kapp-402         3.4-70         3.34-78         1.4-72         3.34-78           1793         cdack-2-00-cm         3.34-70         3.34-78         1.4-72         3.34-78           1795         cdack-2-00-cm         3.40-70         3.34-78         2.34-70         3.34-78           1795         cdack-2-00-cm         3.40-70         3.34-70         3.34-70         3.34-70           1796         cdack-2-00-cm         3.40-70         3.44-70         3.44-70         3.44-70           1797         cdack-2-00-cm         3.40-70         1.46-70         2.44-70         3.44-70           1798         codellam-15-70-cm         2.7-70         1.46-70         2.44-70         3.44-70           18001         opert-1.50-cm         2.7-70         1.66-70         2.44-70         3.44-70           1801         opert-1.50-cm         2.7-70         1.66-70         2.44-72         3.44-70           18001         opert-1.66-70<	1789	open-mistral-7b	3.2e+02	2.4e+02	3e+02	
1991         Uma2-700-xml         12x00         22xm02         13x02         22xm02         33x02         35x02         33x02         3x02         3x02         3x02         3x02         3x02	1790	command-r-plus	3.4e+02	2.4e+02	2.8e+02	
memory and stands         22-4-02         33-4-02         33-4-02         33-4-02           results large-302         3-4-03         22-1-62         33-4-02         33-4-02           results large-302         3-4-03         22-1-62         33-4-02         33-4-02           results large-302         3-4-03         22-1-62         33-4-02         33-4-02           results large-302         3-4-02         21-1-62         22-4-02         34-4-02           results large-302         5-1-60         2-1-62         24-4-02         24-4-02           results large-302         5-1-60         2-1-62         24-4-02         24-4-02           results large-302         1-1-6-62         24-4-02         24-4-02         24-4-02           results large-302         1-1-6-62         24-4-02         24-4-02         24-4-02           results large-304         2-2-1-02         1-1-6-62         24-4-02         24-4-02           results large-304         2-2-1-02         1-1-6-62         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02         24-4-02	1791	llama-2-70b-chat	1.2e+02	2.3e+02	1.2e+02	500
1924         cdude 3 some 30(2020)         3 46-02         2 14-02         1 34-03           1933         meistal-sep-202         3 36-02         2 14-02         3 36-02         3 36-02           1934         deales 3-ope 32(2022)         3 36-02         2 4-102         3 36-02         3 46-02<	1702	mistral-small	3.4e+02	2.2e+02	3.2e+02	- 500
1733         meta-alga-22         3.84-62         2.16-02         3.86-62         2.8-02         3.66-02           1795         dedde-5-200202         3.66-62         2.8-02         3.66-02         2.6-02         3.66-02           1795         dedde-5-30020         3.66-02         2.6-02         3.26-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02         2.6-02         3.6-02	1752	claude-3-sonnet-20240229	3.6e+02	2.1e+02	3.2e+02	
1794         0.0000 (-)puic (2.000)         3.6.4.0.         2.4.4.0.         2.6.4.0.         2.6.4.0.           1795         0.0000 (-)puic (2.000)         3.6.4.0.2         2.4.4.0.2         3.2.4.0.2         3.4.4.0.2         2.4.4.0.2         3.4.4.0.2         1.6.4.0.2         2.4.4.0.2	1793	mistral-large-2402	3.4e+02	2.1e+02 2.1e+02	3.1e+02 2.3e+02	
1795         dudb 3-halto 2020000         44402         244-02         324-02         324-02         164-02         244-02           1796         coolelane 7b-matu         34-02         168-02         244-02         244-02           1797         durch statut         34-02         168-02         244-02         244-02           1799         quent 5-7b-chat         276-02         168-02         24-02         24-02           1800         quent 5-7b-chat         276-02         168-02         24-02         24-02           1801         quent 5-7b-chat         276-02         168-02         24-02         24-02           1801         quent 5-7b-chat         276-02         168-02         24-02         26-02           1802         coolelane 3-4b-chat         34-02         168-02         24-02         26-02           1803         ccoslere 3-7b-chat         34-02         168-02         24-02         24-02           1804         genma 7b-b         26-02         168-02         24-02         24-02           1805         gin-4705-dt         3-902         168-02         24-02         24-02           1806         git-4715         3-902         168-02         24-02         24-02 <td>1794</td> <td>open-mixtral-8x22b</td> <td>3.5e+02</td> <td>2:10+02 2e+02</td> <td>2.3e+02</td> <td></td>	1794	open-mixtral-8x22b	3.5e+02	2:10+02 2e+02	2.3e+02	
1796         codelama 7b estand         3+62         1 8+62         2 4+62           1797         din A and din A and quent 5-7b -bit quent 2-7b -bit quent 3-7b	1795	claude-3-haiku-20240307	4e+02	2e+02	3.2e+02	
1797         dbs/minder         9+92         1.6+02         2.4+92           1798         odellema 13b-mitude         3+92         1.6+02         2.4+92           1799         quent 5.5b-bitte         2.4+92         1.6+02         2.4+92           1800         quent 5.4b-bitte         2.4+92         1.6+02         2.4+92           1801         quent 5.4b-bitte         2.4+92         1.6+02         2.4+92           1802         codellema 70b-bitte         2.4+92         1.6+02         2.4+92           1803         quent 5.16b-bitte         2.4+92         1.6+02         2.4+92           1804         genmo 7b-a         3+92         1.6+02         2.4+92           1805         yi 34b-bitt         2.9+92         1.6+02         2.4+92           1806         gg4-6224-65-13         3+92         1.6+02         2.4+92           1807         gg4-6224-65-13         3+92         1.6+02         2.4+92           1806         gg4-6224-65-13         3+92         1.6+02         2.4+92           1807         gg4-613         2.9+92         1.6+02         2.4+92           1808         gg4-612         3+92         1.6+02         2.4+92           1810	1796	codellama-7b-instruct	3e+02	1.8e+02	2.4e+02	
1.1.1         contained         2.8+r02         1.8+r02         2.4+r02         2.4+r02           1799         contained 3bindted         2.4+r02         1.8+r02         2.4+r02         2.4+r02           1800         quert.15-84-brit         2.4rr02         1.8+r02         2.4+r02         2.4+r02           1801         quert.15-84-brit         2.4rr02         1.8+r02         2.4+r02         2.4+r02           1802         codelimes-34binstric         3-4r02         1.8+r02         2.4+r02         2.4+r02           1803         codelimes-34binstric         3-4r02         1.8+r02         2.4+r02         2.4+r02           1804         gemms-7bin         2.8-r02         1.8+r02         2.4+r02         2.4+r02           1806         gif-4-022-6bind         3-4r02         1.8+r02         2.4+r02         2.4+r02           1806         gif-4-022-6bind         3-4r02         1.8+r02         2.2+r02         2.4+r02           1807         gif-4-0125-proving         3-8+r02         1.8+r02         2.2+r02         2.4+r02           1808         gif-4-0125-proving         3-8+r02         1.8+r02         2.2+r02         2.4+r02           1810         gif-4-0125-proving         3-8+r02         1.8+r02	1797	dbrx-instruct	3e+02	1.8e+02	2.4e+02	
1750         codelimen :3b:initiation         3e-02         1 8e-02         2 4e+02           1800         qwen1.5-32b-drd         2 7e+02         1 8e+02         2 4e+02           1801         qwen1.5-32b-drd         2 4e+02         1 8e+02         2 4e+02           1802         codelimen:34b-initiation         2 4e+02         1 8e+02         2 4e+02           1803         codelimen:34b-initiation         3 e+02         1 8e+02         2 4e+02           1804         gumms.7b-4         2 8e+02         1 8e+02         2 4e+02           1805         gumms.7b-4         2 8e+02         1 8e+02         2 4e+02           1806         gumms.7b-4         2 8e+02         1 8e+02         2 4e+02           1806         gumms.7b-4         2 8e+02         1 8e+02         2 4e+02           1806         gumms.7b-4         2 8e+02         1 8e+02         2 4e+02           1807         gumms.7b-4         3 e+02         1 8e+02         2 4e+02           1808         gumms.7b-4         3 e+02         1 8e+02         2 4e+02           1809         gumms.7b-4         3 e+02         1 8e+02         2 4e+02           1810         gumms.7b-4         3 e+02         1 8e+02         2 4e+0	1700	qwen1.5-7b-chat	2.9e+02	1.8e+02	2.4e+02	
1799       quent.54bc.htt       27402       1.8+02       2.4+02       2.4+02         1800       quent.5.14bc.htt       2.24+02       1.8+102       2.4+02         1801       quent.5.14bc.htt       2.24+02       1.8+102       2.4+02         1802       codelinma-70-initate       38+02       1.8+102       2.4+02         1803       codelinma-70-initate       38+02       1.8+102       2.4+02         1804       gemma-71-4       2.8+02       1.8+102       2.4+02         1805       dim-70-initate       38+02       1.8+102       2.4+02         1806       gi+0-2024-051       38+02       1.8+102       2.4+02         1806       gi+0-2024-051       38+02       1.8+102       2.4+02         1807       lima-3-70-tht       38+02       1.8+102       2.4+02         1808       ggi+0.013 priviev       38+02       1.8+102       2.4+02         1810       ggi+0.013 priviev       38+02       1.8+102       2.4+02         1811       ggi+0.013 priviev       38+02       1.8+102       2.4+02         1810       ggi+0.013 priviev       38+02       1.8+102       2.4+02         1811       ggi+0.4013 priviev       1.8+102       2	1790	codellama-13b-instruct	3e+02	1.8e+02	2.4e+02	
1800         gwn1 5 4b child         2.44+02         1.84+02         2.44+02           1801         gwn1 5 4b child         2.44+02         1.84+02         2.44+02           1802         codellama 4b child         3a+02         1.84+02         2.44+02           1803         codellama 4b child         3a+02         1.84+02         2.44+02           1804         germa 7b+         3a+02         1.84+02         2.44+02           1805         y-4b-child         3a+02         1.84+02         2.44+02           1806         germa 7b+         3a+02         1.84+02         2.44+02           1806         y-4b-child         3a+02         1.84+02         2.44+02           1807         germa 3-70b-child         3a+02         1.84+02         2.44+02           1808         ggh-40-22:A0+03         38+02         1.84+02         2.44+02           1809         ggh-40-22:A0+03         38+02         1.84+02         2.44+02           1809         ggh-40-22:A0+03         38+02         1.84+02         2.44+02           1810         gh-12:X0+04         38+02         1.84+02         2.44+02           1811         gh-22         38+02         1.54+02         2.44+02	1799	qwen1.5-4b-chat	2.7e+02	1.8e+02	2.4e+02	- 400
1801         quart 51-8b-chi         2.24-02         1.8+02         2.4+02           1802         quart 51-8b-chi         2.27+02         1.8+02         2.4+02           1803         codilamo-3b-instrut         3x+02         1.8+102         2.4+02           1804         gemma-7b-instrut         3x+02         1.8+102         2.4+02           1805         gemma-7b-instrut         3x+02         1.8+102         2.4+02           1806         gpt-4-0204-0513         3x+02         1.8+102         2.4+02           1806         gpt-4-0204-0513         3x+02         1.8+102         2.4+02           1807         Isama-38-chi         3x+02         1.8+102         2.4+02           1806         gpt-4-025-proive         3x+02         1.8+102         2.4+02           1807         Isama-37b-chi         3x+02         1.8+102         2.4+02           1809         gpt-4-025-proive         3x+02         1.8+102         2.4+02           1810         gpt-4-025-proive         3x+02         1.8+102         2.4+02           1811         phi2         3x+02         1.8+102         2.4+02           1811         gpt-4-025-proive         3x+02         1.8+102         2.4+02	1800	qwen1.5-32b-chat	2.4e+02	1.8e+02	2.4e+02	
1802	1801	qwen1.5-14b-chat	2.4e+02	1.8e+02	2.4e+02	
coolaima-34-manded         38402         1.84102         2.84102           1803         coolaima-34-manded         38402         1.84102         2.44102           1804         oims-7b-instruct         2.84102         1.84102         2.44102           1806         oims-7b-instruct         384102         1.84102         2.44102           1806         oims-7b-instruct         384102         1.84102         2.44102           1806         oims-7b-instruct         384102         1.84102         2.44102           1807         iiama-3b-chat         384102         1.84102         2.44102           1808         oims-7b-chat         384102         1.84102         2.484102           1809         opt-4.0125 proview         384102         1.84102         2.484102           1810         opt-4.013         2.88402         1.84102         2.84402           1811         opt-2         38402         1.84102         2.84402           1814         wardim-2.83.bctt         3.8402         1.84402         2.48402           1811         opt-2         3.8402         1.84402         2.84402           1814         wardim-2.83.bctt         3.8402         1.84402         2.48402	1802	qwen1.5-1.8b-chat	2.7e+02	1.8e+02	2.4e+02	
1003         000000000000000000000000000000000000	1803	codellama-34b-instruct	30+02	1.80+02	2.60+02	
1804         u         u         1.0 <th1.0< th=""> <th1.0< th=""> <th1.0< th=""></th1.0<></th1.0<></th1.0<>	1005	gemma-7b-it	2 8e+02	1.8e+02	2.4e+02	
1805         yi-34b-that         2x+02         1.8e+02         2.4e+02         2.4e+02           1806         gpl-4-0.2024.05-13         3e+02         1.8e+02         2.4e+02         2.4e+02           1807         ilama-3-bb-that         3e+02         1.8e+02         2.4e+02         2.4e+02           1809         gpl-4-0155 previor         3e+02         1.8e+02         2.4e+02         2.4e+02           1810         gpl-4-0155 previor         3e+02         1.8e+02         2.4e+02         2.4e+02           1811         gpl-4-0155 previor         3e+02         1.8e+02         2.4e+02         2.4e+02           1811         gpl-4-0155 previor         3e+02         1.8e+02         2.4e+02         2.4e+02           1812         gemin-150         3e+02         1.7e+02         2.4e+02         3e+02         1.7e+02         2.4e+02           1813         gwen1.5-10b-that         1.2e+02         1.2e+02         1.2e+02         3e+02         1.8e+02         3e+02         1.8e+02         3e+02         1.8e+02         3e+02         1.8e+02         2.4e+02         3e+02         1.8e+02         2.4e+02         3e+02         1.8e+02         1.8e+02         2.4e+02         3e+02         1.8e+02         1.8e+02         1.	1804	olmo-7b-instruct	3e+02	1.8e+02	2.4e+02	
1806         gpt40-2024-05-13         3e+02         1 8e+02         2.4e+02           1807         Iama-3-05-hat         3e+02         1 8e+02         2.4e+02           1809         gpt4-0125-prover         3e+02         1 8e+02         2.4e+02           1800         gpt4-0125-prover         3e+02         1 8e+02         2.4e+02           1810         gpt4-013         2.8e+02         1 8e+02         2.4e+02           1811         gpt-40         3e+02         1 8e+02         2.4e+02           1811         gpt-40         3e+02         1 8e+02         2.4e+02           1811         gpt-40         3e+02         1 8e+02         2.4e+02           1811         gpt-40-61         2.9e+02         1 7e+02         2.4e+02           1813         qvert1-5+105-hat         1 2e+02         1 2e+02         1 3e+02           1814         wizzdm-2-4bi         1 2e+02         1 2e+02         1 3e+02           1816         command-nightly         60         60         80           1817         sovefalse-actic-instruct         14         14         14           1818         gpt3-5-bit-05         2.20         2.20         2.20           1820         cd	1805	yi-34b-chat	2e+02	1.8e+02	1.8e+02	
1807         IIma-3-9b-chat         3e+02         1 8e+02         2.2e+02         -300           1808         0         18e+02         2.4e+02         2.4e+02           1809         0         18e+02         1.8e+02         2.4e+02           1809         0         2.8e+02         1.8e+02         2.8e+02         2.8e+02           1810         0         2.8e+02         1.8e+02         2.8e+02         2.8e+02           1811         0         0         3e+02         1.7e+02         2.8e+02           1811         0         0         3e+02         1.7e+02         2.4e+02           1812         0         0         1.5e+02         2.4e+02         1.8e+02         2.4e+02           1814         wizadmin-2.8e>20         3e+02         1.5e+02         2.4e+02         1.8e+02         1.3e+02         1.3e+02           1815         IIma-2-13b-chat         1.2e+02         1.2e+02         1.3e+02         2.4e+02         1.8e+02         2.4e+02         1.8e+0	1806	gpt-4o-2024-05-13	3e+02	1.8e+02	2.4e+02	
1808         0         0         18402         18402         24+02           1809         0         18402         18402         24+02           1800         0         18402         28+02         28+02           1810         0         28+02         18+02         28+02           1811         0         0         28+02         18+02         28+02           1811         0         0         28+02         18+02         24+02           1812         0         0         70*02         17*02         24+02           1813         0         9         370*02         1.5*02         24+02           1813         0         9         3*02         1.5*02         24+02           1814         12*02         1.2*02         1.3*02         1.2*02           1816         0         0         60         80           1817         sovak=-arctinatest         4         4         4         4           1818         0         0         0         0         0         10           1820         0         0         0         20         182         0         20         182         18	1807	llama-3-8b-chat	3e+02	1.8e+02	2.2e+02	- 300
00         deepsek-lim-67b-chat         38+02         18+02         24s+02           1810         gpt-40125-preview         38+02         1.8s+02         2.8s+02         2.4s+02         2.8s+02         2.8s+02         2.4s+02         2.4s+02         2.8s+02         2.8s+02         2.4s+02         2.4s+02         2.8s+02         1.8s+02         2.4s+02         2.4s+02         2.4s+02         2.8s+02         1.8s+02         2.4s+02         2.4s+02         2.4s+02         2.8s+02         1.8s+02         2.4s+02         2.4s+02         1.8s+02         2.4s+02	1808	B llama-3-70b-chat	3e+02	1.8e+02	2.4e+02	
1605         gpt-4.0125-preview         3a+02         1.8e+02         2.8e+02           1810         gpt-4.0125-preview         2.6e+02         1.8e+02         2.6e+02           1811         ph-2         3a+02         1.8e+02         2.6e+02           1811         gemini-pro         3.7e+02         1.7e+02         2.4e+02           1813         qwen1.5-110-chat         3e+02         1.7e+02         2.4e+02           1814         wizardin2-8x2b         3e+02         1.5e+02         2.4e+02           1814         wizardin2-8x2b         3e+02         1.2e+02         2.4e+02           1815         Ilama-2-7b-chat         1.2e+02         1.2e+02         1.2e+02         1.2e+02           1816         command-nighty         60         60         80         80           1817         snowflake-arctic-instruct         14         14         14           1818         gemini-1-5 pro-lates         4         4         4         4           1819         chronos-hermes-13b         20         20         20         20           1820         claude-2t         20         20         20         20         20           1821         claude-instant-12	1900	deepseek-llm-67b-chat	3e+02	1.8e+02	2.4e+02	
1810       ggl-40613       2.86+02       1.86+02       2.66+02       2.66+02         1811       phi-2       36+02       1.86+02       2.46+02         1811       gemm-2b-it       2.96+02       1.76+02       2.96+02         1812       gemm-2b-it       2.96+02       1.76+02       2.46+02         1813       quen1.5-100-chat       38+02       1.56+02       2.46+02         1814       wizardm-24x2b       36+02       1.56+02       2.46+02         1815       Iama-2-130-chat       1.20+02       1.36+02       2.46+02         1816       comman-highty       60       60       80         1817       snowtake-artic-instruct       14       14       14         1818       gelmin-1.5-pro-latest       4       4       4         1819       chaude-21       20       20       20         1820       claude-15pro-latest       4       4       4       4       4         1819       chaude-21       20	1009	gpt-4-0125-preview	3e+02	1.8e+02	2.8e+02	
1811	1810	gpt-4-0613	2.8e+02	1.8e+02	2.6e+02	
1812         gemm-2b-ii gemm-2b-ii gemm-2b-ii         2.9e+02         1.7e+02         2.4e+02           1813         qwen1.5-110b-chat         3e+02         1.5e+02         2.4e+02           1814         wizardim-2-8x22b         3a+02         1.5e+02         2.4e+02           1815         iiiama-2-13b-chat         1.2e+02         1.2	1811	pni-2	30+02	1.8e+02	2.4e+02	
1813         quent.5-10b-chat         38-02         1.5-02         2.4e+02           1814         wizardm-2-8x2b         3e+02         1.5e+02         2.4e+02           1815         Ilama-2-13b-chat         1.2e+02         1.2e+02         1.3e+02         1.2e+02           1816         command-nighty         60         60         80         1.2e+02	1812	gemma-2h-it	2.9e+02	1.7e+02	2.3e+02	
1814       wizardin-2-8x2b       3e+02       1.5e+02       2.4e+02         1815       Ilama-2-13b-chat       1.2e+02       1.3e+02       1.2e+02         1816       command-nightly       60       60       80         1817       snowflake-arcie-instrut       14       14       14         1818       gemini-1.5-pro-lates       4       4       4         1819       reka-core-20240501       12       4       5         1820       claude-2       20       20       20         1821       claude-2       20       20       20         1822       claude-instant-1       20       20       20         1823       codellama-34b-python       20       20       20         1824       gpt-3.5-turbo-6061       20       20       20         1825       gpt-3.5-turbo-6061       20       20       20         1826       gpt-3.5-turbo-6061       20       20       20         1826       gpt-3.5-turbo-6061       20       20       20         1826       gpt-3.5-turbo-6061       20       20       20         1827       gpt-4.1106-preview       20       20       20	1813	gwen1.5-110b-chat	3e+02	1.5e+02	2.4e+02	
Iama-2-13b-chat         1.2e+02         1.3e+02         1.3e+02         1.2e+02         1.3e+02         1.2e+02	1814	wizardlm-2-8x22b	3e+02	1.5e+02	2.4e+02	
1613         Iama-2-7b-chat         1.2e+02	1015	llama-2-13b-chat	1.2e+02	1.2e+02	1.3e+02	200
1816         command-nightly         60         60         80           1817         snowflake-arctic-instruct         14         14         14           1818         gemini-1.5-pro-latest         4         4         4           1819         chronos-hermes-13b         12         4         5           1820         claude-2         20         20         20           1821         claude-2.1         20         20         20           1822         claude-instant-1.2         20         20         20           1823         codeliama-34b-python         20         20         20           1824         gpt-3.5-turbo         20         20         20           1825         gpt-3.5-turbo-0613         20         20         20           1826         gpt-3.5-turbo-0614         20         20         20           1827         gpt-4.1106-preview         20         20         20           1828         mistral-tiny         20         20         20           1829         nous-hermes-Ilama2-13b         20         20         20           1820         wizardcoder-15b-v10         20         20         20	1010	llama-2-7b-chat	1.2e+02	1.2e+02	1.2e+02	200
1817         snowllake-arctic-instruct         14         14         14         14           1818         gemini-1.5-pro-latest         4         4         4           1819         cka-core-2024050         12         4         5           1820         chronos-hermes-13b         20         20           1821         claude-2.1         20         20           1822         claude-1.12         20         20           1823         codellama-34b-python         20         20           1824         gpt-3.5-turbo-0301         20         20           1825         gpt-3.5-turbo-0301         20         20           1826         gpt-3.5-turbo-0301         20         20           1827         gpt-3.5-turbo-0301         20         20           1828         mistral-medium         20         20           1829         mous-hermes-llama2-13b         20         20           1820         wizardcoder-15b-v1.0         20         20           1831         wizardcoder-15b-v1.0         20         20	1010	command-nightly	60	60	80	
1818     gemini 1.5-pro-latest     4     4     4       1819     reka-core-2024050     12     4     5       1820     chronos-hermes-13b     20     20       1821     claude-2.1     20     20       1822     claude-instant-1.2     20     20       1823     codellama-34b-python     20     20       1824     gpt-3.5-turbo-0301     20     20       1825     gpt-3.5-turbo-0301     20     20       1828     mistral-medium     20     20       1828     mistral-medium     20     20       1829     nous-hermes-ilama2-13b     20     20       1831     wizardcoder-15b-v1.0     20     20	1817	snowflake-arctic-instruct	14	14	14	
1819     chronos-hermes-13b     12     4     5       1820     chronos-hermes-13b     20     20       1821     claude-2.1     20     20       1822     claude-instant-1.2     20     20       1823     codellama-34b-python     20     20       1824     gpt-3.5-turbo-0301     20     20       1825     gpt-3.5-turbo-0301     20     20       1826     gpt-3.5-turbo-04301     20     20       1828     mistral-medium     20     20       1829     mistral-medium     20     20       1820     wizardcoder-15b-v1.0     20     20       1831     ************************************	1818	gemini-1.5-pro-latest	4	4	4	
1820       claude-2       20         1821       claude-2.1       20         1822       claude-instant-1.2       20         1823       codellama-34b-python       20         1824       gpt-3.5-turbo-0301       20         1825       gpt-3.5-turbo-0301       20         1826       gpt-3.5-turbo-0401       20         1827       gpt-3.5-turbo-0401       20         1828       mistral-medium       20         1829       mistral-medium       20         1820       wizardcoder-15b-v1.0       20         1831       wizardcoder-15b-v1.0       20	1819	chronos hormos 13h	12	4	20	
1821       claude-2.1       20         1822       claude-instant-1.2       20         1823       codellama-34b-python       20         1824       gpt-3.5-turbo-0301       20         1825       gpt-3.5-turbo-0301       20         1826       gpt-3.5-turbo-0400       20         1827       gpt-3.5-turbo-0400       20         1828       mistral-medium       20         1829       mistral-medium       20         1829       nous-hermes-llama2-13b       20         1831       wizardcoder-15b-v1.0       20	1820	claude-2			20	
102.1       claude-instant-1       20         1822       claude-instant-1.2       20         1823       codellama-34b-python       20         1824       gpt-3.5-turbo-0301       20         1825       gpt-3.5-turbo-0301       20         1826       gpt-3.5-turbo-0613       20         1827       gpt-3.5-turbo-1016       20         1828       mistral-medium       20         1829       mistral-medium       20         1820       wizardcoder-15b-v1.0       20         1831       wizardcoder-15b-v1.0       20	1821	claude-2.1			20	
1822       claude-instant-1.2       20       20         1823       codellama-34b-python       20       20         1824       gpt-3.5-turbo-0301       20       20         1825       gpt-3.5-turbo-04301       20       20         1826       gpt-3.5-turbo-04301       20       20         1826       gpt-3.5-turbo-04301       20       20         1826       gpt-3.5-turbo-04301       20       20         1827       gpt-4.1106-preview       20       20         1828       mistral-medium       20       20         1829       mous-hermes-ilama2-13b       20       20         1830       wizardcoder-15b-v1.0       20       20         1831	1021	claude-instant-1			20	
1823       codellama-34b-python       20       -100         1824       gpt-3.5-turbo       20       20         1825       gpt-3.5-turbo-0613       20       20         1826       gpt-3.5-turbo-0613       20       20         1827       gpt-4.1106-preview       20       20         1828       mistral-medium       20       20         1829       nous-hermes-llama2-13b       20       20         1831       wizardcoder-15b-v1.0       20       20	1822	claude-instant-1.2			20	
1824       gpt-3.5-turbo-0301       20         1825       gpt-3.5-turbo-0301       20         1826       gpt-3.5-turbo-0610       20         1826       gpt-3.5-turbo-0610       20         1827       gpt-4.1106-preview       20         1828       mistral-medium       20         1829       mistral-turb       20         1830       wizardcoder-15b-v1.0       20         1831       20       20	1823	codellama-34b-python			20	- 100
1825       gpt-3.5-turbo-0301       20         1826       gpt-3.5-turbo-0610       20         1827       gpt-3.5-turbo-0610       20         1827       gpt-4.1106-preview       20         1828       mistral-medium       20         1829       mistral-tithy       20         1830       wizardcoder-15b-v1.0       20         1831       20       20	1824	gpt-3.5-turbo			20	
1826       gpt-3.5-turbo-0613       20         1826       gpt-3.5-turbo-1106       20         1827       gpt-4.1106-preview       20         1828       mistral-medium       20         1829       mistral-imity       20         1830       wizardcoder-15b-v1.0       20         1831       3       3	1825	gpt-3.5-turbo-0301			20	
gpt-3.5-turbo-1106     20       1827     gpt-4.1106-preview     20       1828     mistral-medium     20       1829     mistral-tiny     20       1830     wizardcoder-15b-v1.0     20       1831     3     3	1826	gpt-3.5-turbo-0613			20	
Instal     gpt-4-1100-preview     20       1828     mistral-medium     20       1829     mistral-tiny     20       1830     wizardcoder-15b-v1.0     20       1831     20     20	1827	gpt-3.5-turbo-1106			20	
1820         mistral-tiny         20           1829         mistral-tiny         20           1830         wizardcoder-15b-v1.0         20           1831         20         20	1021	gpt-4-1106-preview			20	
1829         nous-hermes-lama2-13b         20           1830         wizardcoder-15b-v1.0         20           1831          20	1020	mistral-tiny			20	
1830 wizardcoder-15b-v1.0 20	1829	nous-hermes-llama2-13b			20	
1831	1830	wizardcoder-15b-v1.0			20	
	1831		ł			

Figure 21: AIW Average number of responses per model for each prompt type (4 AIW variations per prompt type.). Models with less than 100 responses per prompt type are excluded from further analysis. All those models have negligible correct response rates, either 0 or close to 0.

Hind of the probability of th	1837					
RESULCIED       PERMINE       PERMINE       PERMINE         READ       0001       0010       00	1838		AIW A	Average length of resp	onses	
1840         prod         1 hadra         2 hadra         2 hadra         2 hadra         2 hadra           1841         codian-Shall-200007         4 hadra         4 hadra         6 hadra         6 hadra           1842         dauk-Shall-200007         1 hadra	1839		RESTRICTED	prompt_type STANDARD	THINKING	
A. M. Wellins 7/1-rature         4.5-62         4.4-62         4.4-62         4.6-62	1840	phi-2	1.3e+03	2.5e+03	2.8e+03	
0.64 (mathematical mathematical m	10/1	codellama-70b-instruct	4.7e+02	4.9e+02	6.9e+02	
0142         cubbs 3 oper 30202020         1.98-02         3.8-02         5.8-02	1041	claude-3-haiku-20240307	1e+02	4.2e+02	4.1e+02	
1414       tando       tando       1.4-02       1.3.4-02 <td>1842</td> <td>claude-3-opus-20240229</td> <td>1.5e+02</td> <td>3.9e+02</td> <td>7.8e+02</td> <td></td>	1842	claude-3-opus-20240229	1.5e+02	3.9e+02	7.8e+02	
1844         cdusk3 some 2004020         00         3 June 20         00         3 June 20         00         3 June 20         00         0 June 20         1 June 20         1 June 20         1 June 20         0 June 20	1843	llama-3-8b-chat	1e+02	3.8e+02	5.3e+02	
1445	1844	claude-3-sonnet-20240229	86	3.7e+02	6.3e+02	0500
1846	1845	wizardIm-2-8x22b	2.7e+02	3.2e+02	2.6e+02	- 2500
1847         dim-7b-147bc         2.64-92         2.84-92         3.44-92           1848         gen-7b-147bc         2.012         2.24-92         3.44-92           1849         gen-7b-147bc         7b         7b         7b         7b           1849         gen-7b-147bc         7b	1846	gnt-40-2024-05-13	87	3e+02	6.20+02 4 7e+02	
1.1.1         2x42         2x402         34402           1848         gen15.75.041         76         2x402         34402           1849         gen2.15.75.041         75         2x402         44402           1850         Imma 2.18.041         39         2x402         44402           1851         Imma 2.18.041         48e02         2x402         44e02           1852         coodian-15-battod         48e02         2x402         35602           1853         mitral-larp-2402         2x402         35602         36602           1854         open-mitral         2x402         2x402         35602           1855         dis-Antitral         12862         2x402         32802           1856         genma 5.70-batt         12862         24802         33802           1857         Imma 5.70-batt         12862         24802         38802           1858         genma 5.70-batt         12862         28802         18802         39802           1859         gen1.5.70-batt         12862         18802         39802         24e02         18802         13866         18802         39802         24e02         18802         18802         18802         18802 <td>1847</td> <td>olmo-7b-instruct</td> <td>2.6e+02</td> <td>2.9e+02</td> <td>6.1e+02</td> <td></td>	1847	olmo-7b-instruct	2.6e+02	2.9e+02	6.1e+02	
1040         open mids 40%         76         224402         324402         344072           1850         1am-2/13-040         73         274402         444072           1851         1am-2/13-040         73         274402         444072           1852         oodellam-15-33-040         14402         284402         444072           1853         oodellam-15-34040         14402         284402         336402           1854         oopen midrah 70         14402         284402         2.84402         2.84402           1855         daw-intato         2.24402         2.24402         3.96402           1856         daw-intato         2.24402         2.84402         2.84402           1856         daw-intato         2.24402         2.84402         2.84402           1856         daw-intato         2.24402         2.84402         2.84402           1857         daw-intato         2.24402         2.84402         2.84402           1856         daw-intato         2.24402         2.84402         2.84402           1857         daw-intato         2.24402         2.84402         2.84402           1856         dawnin 5-16-041         1.84402         2.84402         1.	1010	reka-core-20240501	2e+02	2.9e+02	3.4e+02	
1249         open-mdd:BA0;         2.2en/2         2.7en/2         4.1en/2           1850         liama 2.7b-dat         59         2.7en/2         4.4en/2           1851         liama 2.7b-dat         59         2.7en/2         4.4en/2           1852         cooffailes actic/introl         4.4en/2         2.8en/2         4.1en/2           1853         cooffailes 15b-ettrod         4.4en/2         2.8en/2         3.8en/2           1854         open-model Abor         2.1en/2         2.8en/2         3.8en/2           1855         data-strong         2.1en/2         2.8en/2         2.8en/2         2.8en/2           1855         data-strong         2.1en/2         2.8en/2         2.8en/2         2.8en/2           1856         data-strong         2.2en/2         2.8en/2         2.8en/2         2.8en/2           1857         data-strong         2.2en/2         1.8en/2         2.8en/2         1.8en/2         2.8en/2           1856         data-strong         66         1.8en/2         2.8en/2         1.8en/2         2.8en/2           1861         get-1.5en/2         1.8en/2         2.8en/2         2.8en/2         1.8en/2         2.8en/2           1862         get-1.5en/2	1040	qwen1.5-7b-chat	76	2.9e+02	3.4e+02	
1850       lima-2-10-oht       75       2.74-02       3.64-02         1851       uma-2-70-oht       39       2.24-02       3.64-02         1852       ucolalian-3/b-intiti       1.44-02       2.84-02       3.64-02         1853       maint-3/b-intiti       1.44-02       2.84-02       3.64-02         1854       open-maint-7b       2.14-02       2.84-02       3.64-02         1855       dow-inititii       2.84-02       2.84-02       3.64-02         1856       open-maint-7b       2.14-02       2.84-02       3.64-02         1857       lima-3.7b-oht       15       2.24-02       2.84-02       3.84-02         1856       deepseekins 7b-oht       2.52-02       2.84-02       3.84-02       3.84-02         1857       lima-3.7b-oht       2.52-02       2.84-02       3.84-02       3.84-02         1860       open-minit-15-0-intat       1.84-02       2.84-02       3.84-02       3.84-02         1861       gpit-4-tito-202-04-00       01       1.86+02       2.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02       3.84-02	1849	open-mixtral-8x7b	2.2e+02	2.7e+02	4.1e+02	
1851         Immo 27-book         39         2.7-m02         3.6-m02           1852         scooldelam-102-instruct         4.16-102         2.6-m02         3.6-m02           1853         minital-lang-2012         7.4         2.5-m02         3.6-m02           1854         oper-minital-N         7.1         2.5-m02         3.6-m02           1855         diox-minital-N         2.1-m02         2.4-m02         3.8-m02           1856         diox-minital-N         2.1-m02         2.4-m02         3.8-m02           1855         diox-minital-N         2.2         2.4-m02         3.8-m02           1856         diox-minital-N         2.2         2.4-m02         3.8-m02           1857         lima-3.70-bal         2.9         2.1-m02         3.8-m02           1856         degestack-Im-57-bat         2.5         1.8-m02         3.8-m02           1860         opf-4-mino-202-bal         649         1.8-m02         3.8-m02           1861         opf-4-mino-202-bal         649         1.8-m02         3.8-m02           1862         opf-1.100-bal         649         1.8-m02         3.8-m02           1864         oper-1.5-100-bal         649         1.8-m02         2.8-m02	1850	llama-2-13b-chat		2.7e+02	4.4e+02	
1852         soudilate ardici-inition         14:e-02         2.6er02         3.6er02         3.6er02           1853         codelinition-31-bindtud         14:e-02         2.6er02         3.6er02         3.6er02           1854         open-mistab-70         2.1er02         2.4er02         3.6er02         3.6er02           1855         down-inition         2.5er02         2.4er02         3.2er02         4.1er02           1856         down-inition         1.2er02         2.1er02         3.2er02         3.8er02           1856         deepseek-Im-67b-bat         1.2er02         2.1er02         3.2er02         3.8er02           1857         deepseek-Im-67b-bat         1.2er02         1.8er02         1.8er02         2.6er02           1859         gemma 7b-         2.5er02         1.8er02         2.6er02         3.8er02           1861         geb-4-ubb-3026/4040         61         1.8er02         2.6er02         3.8er02           1862         gewant.5-14b-bat         66         1.6er02         2.6er02         3.8er02           1863         open-mistab-822         1.3er02         1.6er02         2.6er02         3.6er02           1864         owent.5-32b-bat         66         1.6er02         3.	1851	llama-2-7b-chat	39	2.7e+02	3.3e+02	
codeliman :18: maintained         14: m22         2.6 m22         3.8 m22           1853         maintaineg-202         74         2.6 m22         3.8 m22           1854         open-maintaineg-202         2.4 m22         2.6 m22         3.8 m22           1855         dinx-instruct         2.6 m22         2.4 m22         2.6 m22         3.8 m22           1856         dinx-instruct         15         2.2 m22         2.8 m22         3.8 m22           1857         dinma-370-maint         62         2.1 m22         2.8 m22         3.8 m22           1856         deepsed-in-m7h-maint         5.2 m22         1.9 m22         2.8 m22         1.9 m22           1859         amma-370-bant         2.7 m22         2.6 m22         1.9 m22         1.9 m22           1861         open-minitain-822         2.6 m22         1.9 m22         1.9 m22         1.9 m22           1861         open-minitain-822         1.9 m22         1.9 m22         2.8 m22         1.9 m22           1862         open-minitain-822         1.9 m22         2.8 m22         1.9 m22         2.8 m22           1863         open-minitain-823         1.9 m22         2.8 m22         1.9 m2         1.9 m2           1864         open-mi	1852	snowflake-arctic-instruct	4.9e+02	2.6e+02	4.1e+02	
1555         mistal-arge-2402         74         2.24+02         3.24+02         2.24+02         2.24+02         3.24+02           1855         dor-mistal-7         2.14+02         2.24+02         3.24+02         3.24+02           1856         dor-mistal-7         12.04+02         2.14+02         2.14+02         3.24+02           1856         dor-mistal-7         12.04+02         2.14+02         2.84+02         3.84+02           1857         lisma-3-70-bate         92         2.14+02         2.84+02         3.84+02           1858         despose-lism-7b-rbate         92         2.14+02         3.84+02         3.84+02           1859         gemma-7b-te         2.66+02         1.86+02         3.24+02         3.84+02           1860         codelame-7b-intute         09         1.66+02         3.28+02         3.84+02           1861         gen-mistal-802         1.38+02         1.86+02         3.28+02         3.84+02           1862         gen-mistal-1.54+0-04         64         1.86+02         2.84+02         3.84+02           1864         gen-mistal-1.54+0-04         68         1.84+02         3.84+02         3.84+02           1864         gen-mistal-1.54+0-04         68         9	1853	codellama-13b-instruct	1.4e+02	2.6e+02	3.6e+02	- 2000
1854         open-minit for open-m	1055	mistral-large-2402	74	2.5e+02	3.8e+02	
1855	1854	open-mistral-/b	2.1e+02	2.4e+02	2.5e+02	
1856         Link of 15/2b         12-02         22-02         23-02	1855	Ilama-2-70b-chat	2.50+02	2.40+02	3.20+02 4 1e+02	
1857         Imma-3-70-chtal         27m-62         2.1m+02         3.3m+02           1858         despeaek.Im-70-chtal         27m-62         3.2m+02         3.2m+02           1859         gemini-15-pro-lates         16         1.8m+02         3.2m+02           1860         codeliama-34b-natrod         50         1.8m+02         2.4m+02         3.2m+02           1861         gemini-15-pro-lates         16         1.8m+02         2.4m+02         3.2m+02           1861         gemini-15-bro-lates         60         1.8m+02         3.2m+02         3.2m+02           1862         gemini-50-shato         68         1.8m+02         3.2m+02         3.2m+02           1863         gemini-51-bro-brato         68         1.8m+02         3.2m+02         3.2m+02           1864         gemini-51-bro-brato         68         1.8m+02         3.2m+02         3.2m+02           1865         gemini-51-bro-brato         68         1.8m+02         3.2m+02         3.2m+02           1866         gemini-51-bro-brato         68         1.8m+02         3.2m+02         3.2m+02           1866         gemini-51-bro-brato         68         1.8m+02         3.2m+02         3.2m+02           1866 <t< td=""><td>1856</td><td>gwen1.5-72b-chat</td><td>1.2e+02</td><td>2.1e+02</td><td>2.8e+02</td><td></td></t<>	1856	gwen1.5-72b-chat	1.2e+02	2.1e+02	2.8e+02	
1858         despesek-lim/67b-chat         2.7e+02         2.9e+02         3.2e+02         3.2e+02           1859         gemm-7b-i         2.5e+02         1.9e+02         3.2e+02         1.9e+02           1860         codellam-3-th-instrue         59         1.8e+02         2.7e+02         3.2e+02           1861         gent-1-th/o-10400         61         1.8e+02         3.2e+02         3.8e+02           1862         gent-15-14b-chat         64         1.8e+02         3.2e+02         3.8e+02           1864         qeent-15-14b-chat         68         1.8e+02         2.8e+02         3.8e+02           1864         qeent-15-14b-chat         68         1.8e+02         2.8e+02         3.8e+02           1865         qeent-15-14b-chat         73         1.6e+02         3.8e+02         3.8e+02           1866         qeent-15-14b-chat         73         1.6e+02         2.8e+02         3.8e+02           1866         gentma-2b+1         33         78         1.3e+02         3.8e+02           1867         quad-15-11b-chat         14         14         13         3           1871         gpt-40125-preview         14         18         15           1874         cdu	1857	llama-3-70b-chat	52	2.1e+02	3.3e+02	
15.50         gemmin.7b-it         2.5e+02         1 9e+02         3.2e+02         9e-02         9e-02 <td>1858</td> <td>deepseek-llm-67b-chat</td> <td>2.7e+02</td> <td>2e+02</td> <td>2.9e+02</td> <td></td>	1858	deepseek-llm-67b-chat	2.7e+02	2e+02	2.9e+02	
1559         gemin 1-5 pro-latest         15         1.8e+02         1.9e+02         2.4e+02           1860         opi-4-tubo 224-04-09         61         1.8e+02         2.7e+02         2.8e+02           1861         opi-4-tubo 224-04-09         61         1.8e+02         3.5e+02         3.5e+02           1862         opi-motal-15-14b-04         64         1.8e+02         3.2e+02         3.5e+02           1864         qeen-15-14b-04         68         1.8e+02         2.8e+02         3.6e+02           1865         qeen15-16b-04         68         1.8e+02         2.5e+02         3.6e+02           1866         qeen15-16b-04         68         1.4e+02         2.5e+02         3.6e+02         1.8e+02         3.6e+02	1050	gemma-7b-it	2.5e+02	1.9e+02	3.2e+02	
1860       codellama-34-instruct       99       1.8e-02       2.4e-02       7e-02         1861       gpt-4-tubo-2024-04-09       61       1.8e+02       3.5e-02       3.5e-02         1862       ope       diata       1.3e+02       1.8e+02       3.2e+02       3.5e+02         1863       ope       1.3e+02       1.8e+02       2.5e+02       3.5e+02       3.5e+02 </td <td>1859</td> <td>gemini-1.5-pro-latest</td> <td>15</td> <td>1.8e+02</td> <td>1.9e+02</td> <td></td>	1859	gemini-1.5-pro-latest	15	1.8e+02	1.9e+02	
1861         gp1-4.ubc.2022-04.09         61         1.8e+02         2.7e+02         7.000           1862         pg0         codellama-7b-instrut         89         1.8e+02         3.5e+02           1863         open-mixtral-8x22b         1.3e+02         1.8e+02         2.8e+02           1864         open-mixtral-8x22b         1.3e+02         1.8e+02         2.8e+02           1865         open-mixtral-8x22b         1.3e+02         1.8e+02         3.6e+02           1866         qwen1.5-10b-chat         66         1.4e+02         2.5e+02           1866         qwen1.5-13b-chat         60         1.4e+02         2.5e+02           1866         qwen1.5-14b-chat         36         91         1.1e+02           1866         gpt3.5-turb-0125         30         37         78           1869         gpt4.4-0613         13         14         13           1870         command-rplus         14         13         13           1871         gpt4.4-0613         13         13         13           1872         command-rplus         13         13         13           1873         durons-hermes-1b         2.8e+02         2.8e+02           1876	1860	codellama-34b-instruct	59	1.8e+02	2.4e+02	1500
1862         0         1.8e+02         3.5e+02           1863         qven1.5-14b-chat         64         1.8e+02         2.8e+02           1864         qven1.5-110b-chat         68         1.8e+02         2.8e+02           1865         mistral-small         1.3e+02         1.6e+02         2.8e+02           1866         qven1.5-18b-chat         60         1.4e+02         2.5e+02           1867         gven1.5-18b-chat         30         78         1.8e+02           1869         gp1.3-5-tu+b-chat         30         37         78           1869         gp1.3-5-tu+b-chat         31         13         13           1870         command-rupus         14         14         43           1871         gp4.4o+03         33         13         13           1871         chaude-21         0         2.8e+02           1873         chaude-21         0         2.8e+02 <td< td=""><td>1861</td><td>gpt-4-turbo-2024-04-09</td><td>61</td><td>1.8e+02</td><td>2.7e+02</td><td>- 1500</td></td<>	1861	gpt-4-turbo-2024-04-09	61	1.8e+02	2.7e+02	- 1500
E         qwn15.14b-chaf         64         1.8e+02         3.2e+02           1864         open-mixtra <sup>1</sup> 8/22b         1.8e+02         2.8e+02           1865         mistra <sup>1</sup> -small         1.3e+02         1.8e+02         2.8e+02           1866         qwen1.5-10b-chat         68         1.8e+02         1.8e+02         1.8e+02           1866         qwen1.5-12b-chat         73         1.6e+02         3.8e+02         2.8e+02           1866         qwen1.5-14b-chat         73         1.6e+02         2.5e+02         3.8e+02           1867         qwen1.5-4b-chat         36         91         1.1e+02         3.8e+02           1868         gemma-2b+1         33         78         1.3e+02         3.8e+02           1869         gp1.4-0125         30         37         78         3.8e+02           1870         command-rptus         14         14         43           1871         gp4-0-013         13         14         13           1872         command-rptus         9.8         10         26           1874         claude-2          1.8e+02         1.8e+02           1875         claude-2          1.8e+02	1862	codellama-7b-instruct	89	1.8e+02	3.5e+02	
1364         open-mittral-8/22         136+02         136+02         2.56+02           1864         qiven1.5-10b-bal         68         1.66+02         1.86+02           1865         qiven1.5-18b-bal         73         1.66+02         3.66+02           1866         qiven1.5-18b-bal         73         1.66+02         3.66+02           1867         qiven1.5-4b-bal         36         91         1.16+02           1868         gemma-2b-ll         33         78         1.36+02           1869         gpt-40125-preview         14         18         15           1870         command-rplus         14         14         43           1871         gpt-40125-preview         14         14         13           1872         command-rplus         14         14         43           1873         dchronzs-hermes-12b	1863	gwen1.5-14b-chat	64	1.8e+02	3.2e+02	
1604         dywn, Srindshat         08         Leero2         2.59402           1865         mistaismall         1.3e+02         1.6e+02         3.6e+02           1866         qwen1.5-1.8b-chat         60         1.4e+02         2.5e+02           1867         qwen1.5-2b-chat         60         1.4e+02         2.5e+02           1867         qwen1.5-4b-chat         36         91         1.1e+02           1868         gemma-2b-t         33         76         1.3e+02           1869         gpt-3.5-turb-0/125         30         37         78           1870         command-repus         14         16         15           1871         gpt-4-0125-preview         14         13         13           1872         command-repus         14         13         13           1873         gemin-pro         9.8         10         2.6e+02           1874         claude-2         1.6e+02         1.6e+02         1.6e+02           1875         claude-instan-1         2.6e+02         1.6e+02         1.6e+02           1876         claude-instan-12         2.6e+02         1.6e+02         1.6e+02           1877         ccodelam-34b-python	1067	open-mixtral-8x22b	1.3e+02	1.8e+02	2.8e+02	
1865         Industry of the second seco	1004	qwent.5-1105-chat	1 36+02	1.60+02	2.50+02	
1866         quent.5-32b-chat         60         1.4e+02         2.5e+02           1867         quent.5-32b-chat         36         91         1.1e+02           1868         gemma-2b-it         33         78         1.3e+02           1869         gp14-3.5-turb-0125         30         37         78         1.3e+02           1869         gp14-0125-preview         14         18         15         16           1870         command-r-plus         14         14         13         14         13         14         13         14         13         14         14         14         14         14         14         14         14         14         14         14         14         14         14	1865	gwen1.5-1.8b-chat	73	1.6e+02	3.6e+02	
1867         qwen1.5-4b-chat         36         91         1.1e+02           1868         gemma-2b-it         33         78         1.3e+02           1869         gpl-3.5-turbo-0125         30         37         78           1870         command-rplus         14         18         15           1871         gpt-4.0125-preview         14         13         13           1871         gpt-4.0613         13         14         13           1872         command-rightty         13         13         13           1873         gemin-pro         9.8         10         26           1874         claude-21         2.0         1.0e+02           1875         claude-21         2.0e+02         1.0e+02           1875         claude-1nstant-1         3.0e+02         1.0e+02           1876         claude-instant-12         2.0e+02         1.0e+02           1877         codellama-34b-python         2.0e+02         1.0e+02           1879         gp1-3.5-turbo-0613         3.0e+02         1.0e+02           1880         gp1-3.5-turbo-0166         60         1.0e+02           1881         gp1-4.106-preview         0         1.0e+02	1866	gwen1.5-32b-chat	60	1.4e+02	2.5e+02	
1868       gema-2b-it       33       78       1.3e+02         1869       gpt-3.5-turb-0125       30       37       78         1870       command-r-plus       14       18       15         1871       gpt-4-013       13       14       13         1872       command-righty       13       13       13         1873       gemin-pro       9.8       00       26         1874       claude-2       2       28e+02       28e+02         1875       claude-1stant-1       28e+02       28e+02       28e+02         1875       claude-1stant-1       3.6e+02       3.6e+02       3.6e+02         1876       claude-instant-1       2.8e+02       3.6e+02       3.6e+02         1876       claude-instant-1       2.8e+02       3.6e+02       3.6e+02         1876       claude-instant-1       2.8e+02       3.6e+02       3.6e+02       3.6e+02         1877       codellama-34b-python       3.6e+02       80       3.6e+02       3.6e+02 <t< td=""><td>1867</td><td>qwen1.5-4b-chat</td><td>36</td><td>91</td><td>1.1e+02</td><td></td></t<>	1867	qwen1.5-4b-chat	36	91	1.1e+02	
1869         gpt-3.5-turbo-0125         30         37         78         -1000           1870         gpt-4.0125-preview         14         18         15           1871         gpt-4.0613         13         14         13         14         13         13         13         13         13         13         13         13         13         13         13         13	1868	gemma-2b-it	33	78	1.3e+02	
1000         gpt-4-0125-preview         14         18         15           1870         command-r-plus         14         14         43           1871         gpt-4-0613         13         14         13           1872         command-rightty         13         13         13         13           1872         command-rightty         13         13         13         13         13           1873         chronos-hermes-109         9.8         10         26         28+02         28+0	1869	gpt-3.5-turbo-0125	30	37	78	- 1000
1670         command-r-plus         14         14         43           1871         gpt-4-0613         13         14         13           1872         command-righty         13         13         13           1873         command-righty         13         13         13           1873         chronos-hermes-13b         9.8         10         26           1874         claude-2         16         1.6e+02           1875         claude-2.1         180         3.6e+02           1876         claude-instant-1         2.6e+02         3.6e+02           1876         claude-instant-1         2.6e+02         3.6e+02           1877         codellama-34b-python         2.6e+02         3.6e+02           1878         gpt-3.5-turbo-0301         10         3.6e+02           1879         gpt-3.5-turbo-0301         58         80           1881         gpt-4.1106-preview         6         80           1882         mistral-medium         2.5e+02         3.6e+02           1883         mistral-timp         2.5e+02         3.6e+02           1884         wizardcoder-15b-10.0         2.5e+02         3.6e+02           1885         1	1070	gpt-4-0125-preview	14	18	15	
18/1       gpt.4-0613       13       14       13         1872       command-nightly       13       13       13       13         1872       command-nightly       13       13       13       13         1873       chronos-hermes-13b       9.8       10       26         1874       claude-2       2.8e+02       1.8e+02         1875       claude-2.1       1.8e+02       1.8e+02         1876       claude-instant-1       2.6e+02       3.6e+02         1877       codellama-34b-python       2.6e+02       2.6e+02         1878       gpt-3.5-turbo       11       14       13         1879       gpt-3.5-turbo-0301       11       14       16+02         1880       gpt-3.5-turbo-0403       58       80       13       16+02         1881       gpt-4.1106-preview       1.3e+02       80       13       16       2.6e+02       13       16       2.6e+02       188       13       13       13       13       14       16       16       16       16       16       16       16       16       16       16       16       16       16       16       16       16       16 <t< td=""><td>1070</td><td>command-r-plus</td><td>14</td><td>14</td><td>43</td><td></td></t<>	1070	command-r-plus	14	14	43	
1872         command-nginiy gemini-pro         13 <t< td=""><td>1871</td><td>gpt-4-0613</td><td>13</td><td>14</td><td>13</td><td></td></t<>	1871	gpt-4-0613	13	14	13	
1873         definition         30         10         20         20           1873         chronos-hermes-13b         2.8e+02         2.8e+02         1.6e+02           1874         claude-2.1         1.8e+02         3.6e+02         1.8e+02           1876         claude-instant-1         2.8e+03         3.6e+02         1.8e+02           1876         claude-instant-1.2         2.6e+02         3.6e+02         1.8e+03           1877         codellama-34b-python         2.3e+03         1.8e+03         1.8e+03           1878         gpt-3.5-turbo         1.1e+02         1.8e+02         1.8e+03           1879         gpt-3.5-turbo-0813         68         80         1.3e+02           1880         gpt-3.5-turbo-0106         80         1.3e+02         1.8e+02           1881         gpt-4.1106-preview         2.6e+02         1.3e+02         1.8e+02           1882         mistral-imedium         2.6e+02         1.8e+02         1.8e+02           1883         mistral-imedium         2.5e+02         2.5e+02         1.8e+02           1884         wizardcoder.15b-1,0         71         71         1.8e+02	1872	command-nightly	13	13	13	
1874       claude-2       1.6e+02         1875       claude-2.1       1.6e+02         1876       claude-instant-1       3.6e+02         1876       claude-instant-1.2       2.6e+02         1877       codellama-34b-python       2.3e+03         1878       gpt-3.5-turbo       11e+02         1879       gpt-3.5-turbo-0301       117         1880       gpt-3.5-turbo-0103       68         1881       gpt-4.1106-preview       60         1882       mistral-medium       2.6e+02         1883       mistral-tiny       2.5e+02         1884       wizardcoder-15b-v1.0       71	1873	chronos-bermes-13b	9.8	10	20	
1875       claude-2.1       18e+02         1876       claude-instant-1       3.6e+02         1877       codellama-34b-python       2.6e+02         1878       gpt-3.5-turbo       11e+02         1879       gpt-3.5-turbo-0301       117         1880       gpt-3.5-turbo-0103       68         1881       gpt-4.1106-preview       80         1882       mistral-medium       2.6e+02         1883       mistral-ting       2.5e+02         1884       wizardcoder.15b-0.0       71	1874	claude-2			1.6e+02	
1876       claude-instant-1       3.6e+02       3.6e+02         1877       codellama-34b-python       2.6e+02       2.6e+02         1878       gpt-3.5-turbo       1e+02       1e+02         1879       gpt-3.5-turbo-0813       58       68         1880       gpt-3.5-turbo-0106       80       1.3e+02         1881       gpt-4.1106-preview       2.6e+02       1.3e+02         1882       mistral-medium       2.6e+02       1.3e+02         1883       mistral-itruy       2.5e+02       2.5e+02         1884       wizardcoder.15b-v1.0       71	1875	claude-2.1			1.8e+02	
1876       claude-instant-1.2	1075	claude-instant-1			3.6e+02	
1877       codellama-34b-python       2.3e+03         1878       gpt-3.5-turbo       1e+02         1879       gpt-3.5-turbo-0301       17         1880       gpt-3.5-turbo-0103       58         1881       gpt-4.1106-preview       80         1882       mistral-medium       2.6e+02         1883       nous-hermes-llama2-13b       2.5e+02         1884       wizardcoder-15b-v1.0       71	1876	claude-instant-1.2			2.6e+02	- 500
1878       gpt-3.5-turbo       1e+02         1879       gpt-3.5-turbo-0301       17         1880       gpt-3.5-turbo-0613       58         1880       gpt-3.5-turbo-106       80         1881       gpt-4.1106-preview       1.3e+02         1882       mistral-medium       2.6e+02         1883       mistral-tirup       2.5e+02         1884       wizardcoder-15b-v1.0       71	1877	codellama-34b-python			2.3e+03	500
1879       gpt-3.5-turbo-0301       17         1880       gpt-3.5-turbo-0613       58         1881       gpt-3.5-turbo-106       80         1881       gpt-4.1106-preview       1.3e+02         1882       mistral-medium       2.6e+02         1883       mistral-tiny       2.5e+02         1884       wizardcoder-15b-v1.0       71	1878	gpt-3.5-turbo			1e+02	
1880     gpt-3.5-turbo-0613     58       gpt-3.5-turbo-106     80       1881     gpt-4-1106-preview     1.3e+02       1882     mistral-medium     2.6e+02       1883     mistral-tiny     2.5e+02       1884     wizardcoder-15b-v1.0     71	1879	gpt-3.5-turbo-0301			17	
gpt3.5-turbo-1106     80       1881     gpt4-1106-preview     1.3e+02       1882     mistral-medium     2.6e+02       1883     mistral-tiny     2.5e+02       1884     wizardcoder-15b-v1.0     71	1880	gpt-3.5-turbo-0613			58	
Install         gpt-4-110b-preview         I.3e+02           1882         mistral-medium         2.6e+02           1883         mistral-liny         2.5e+02           1884         wizardcoder-15b-v1.0         2.5e+02           1885         71	1881	gpt-3.5-turbo-1106			80	
1882     misual-incoduit     2.0e+02       1883     mistral-tiny     2.5e+02       1884     wizardcoder-15b-v1.0     2.5e+02       1885     71	1001	gpt-4-1106-preview			1.30+02	
1883         nous-hermes-llama2-13b         2.5e+02           1884         wizardcoder-15b-v1.0         71           1885         1885         1885	1882	mistral-mealum			2.00+02 2.5e+02	
1884 wizardcoder-15b-v1.0 71	1883	nous-hermes-llama2-13h			2.5e+02	
1885	1884	wizardcoder-15b-v1.0			71	
	1885		1			

Figure 22: Average length (on characters) of responses per model for each AIW prompt variation.
We see that phi-2 has the highest average length of responses (probably because it is not a classical instruction tuned model, but a base model, capable of following instructions).



#### Human

Alice has 4 sisters and she also has 1 brother. How many sisters does Alice's brother have? Before providing answer to this problem, think carefully and double check the path to the correct solution for any mistakes. Provide then the final answer in following form: \"### Answer: \"

#### **Claude 3 Opus**

Let's think through this problem step by step: 1. Alice has 4 sisters and 1 brother. 2. Alice's brother is part of the same family as Alice. 3. All of Alice's sisters are also sisters to her brother. 4. Therefore, Alice's brother has the same number of sisters as Alice has, which is 4. 1've double-checked the reasoning and can't find any mistakes in the logic. ### Answer: Alice's brother has 4 sisters

Figure 23: Example of a wrong response for Claude 3 Opus on THINKING prompt type, featuring a failure of self-verification.

#### Human

Alice has four brothers and she also has a sister. How many sisters does Alice's brother have?

#### LLaMA-3-70B-Instruct

A classic lateral thinking puzzle!

Let's break it down:

\* Alice has four brothers.

\* Alice also has a sister (herself).

So, each of Alice's brothers has... (drumroll) ... one sister, which is Alice herself!

Therefore, the answer is: each of Alice's brothers has one sister.

Figure 24: Example of a wrong response for Llama 3 70B Instruct. The model counts Alice as her own sister therefore produces wrong answer.

Human

Alice has 4 brothers and she also has 1 sister. How many sisters does Alice's brother have? Solve this problem and provide the final answer in following form: \"### Answer: \".

#### GPT-4-turbo-2024-04-09

Alice has 1 sister. Since they are siblings, **each of Alice's brothers would have the same number of** sisters, which includes Alice and her other sister. Thus, each of Alice's brothers has 2 sisters. ### Answer: 2

1941 1942 1943

Figure 25: Example of a correct final answer with fully correct reasoning for GPT-4 Turbo.

#### Human Alice has three brothers and she also has three sisters. How many sisters does Alice's brother have? Claude 3 Opus Okay, let's think about this step-by-step: 1) We know that Alice has three brothers and three sisters. This means that there are a total of 7 children in Alice's family (Alice + 3 brothers + 3 sisters). 2) Now, let's consider one of Alice's brothers. Let's call him Bob. 3) How many sisters does Bob have? Well, all of Alice's sisters are also Bob's sisters. Alice has 3 sisters, so Bob has 3 sisters too. 4) But wait! Alice is also Bob's sister, but we didn't count her yet. So we need to add Alice to the count of Bob's sisters. 5) Therefore, Bob has 3 sisters (Alice's sisters) + 1 sister (Alice herself) = 4 sisters in total. Since all of Alice's brothers are in the same situation as Bob, each of Alice's brothers has 4 sisters. Figure 26: Example of a correct final answer with fully correct, detailed reasoning for Claude 3 Opus. Human Alice has four brothers and she also has a sister. How many sisters does Alice's brother have? Solve the problem taking care about not doing mistakes. Command-R+ This problem can be solved accurately and unambiguously. Alice has four brothers, which clearly indicates that she has four male siblings. Additionally, it is stated that she also has a sister, which designates a female sibling. Therefore, Alice's brothers have one sister. This interpretation is straightforward and avoids any potential mistakes or ambiguities. It is based solely on the information provided, and there is no need for further interpretation or speculation. So, the definitive answer is: Alice's brothers have one sister.

Figure 27: Overconfidence: example of a wrong response for Command R Plus, containing highly persuasive wording to back up its confidence into incorrect solution

#### Human

Alice has X brothers and she also has Y sisters. How many sisters does Alice's brother have? Solve this problem and provide the final answer in following form: \"### Answer: \"

#### OLMo-7B-Instruct

The number of Alice's brothers is X number of Alice's sisters
The number of X brothers is Y, and the total number of brothers and sisters is $(X + Y)$ . how, we can
represent X's brothers by x and Y's sisters by y, so $(X + Y) = x + y$ .
Since the total number of brothers and sisters is $X + Y, X = (X + Y) - Y$ .
Now, divide both sides by $(X + Y)$ to find X: $X = (X + Y) - Y = (x + y) - y$ - In other words: $X = x - y$ . Subtract y
from both sides: X = x - y.Now, scale the -y to the left side by adding 5 to both sides: X = ((x + 5) - y).
So, the number of Alice's brothers is either x, x + 1, x + 2, x + 3, or x + 4.\\### Answer: x - 4

Figure 28: Confabulation: example of a wrong response for OLMo-7b-Instruct for AIW-param problem. The model generates text that contains different algebraic manipulations involving X and Y and has little to do with reasoning.

#### Human

Alice has four brothers and she also has a sister. How many sisters does Alice's brother have?

#### CodeLlama-70b-Instruct

I apologize, but as a responsible AI language model, I cannot provide an answer that would potentially perpetuate harm or discrimination towards any individual or group, **including a person of living with Down syndrome**. It is important to recognize that every individual has inherent dignity and worth, regardless of their abilities or characteristics. Instead, I **suggest rephrasing the question to focus on Alice's relationship with the person living with Down syndrome** in a respectful and inclusive manner. I can assist you in modifying the question to promote understanding and empathy towards individuals with diverse abilities and characteristics. If you have any other questions or concerns, please feel free to ask, and I will do my best to provide helpful and responsible responses.

Figure 29: Example of a confabulation response for CodeLlama-70b-Instruct. In this example, the model produces wrong reasoning and in addition distracts from the actual main problem to solve by bringing up that it is disrespectful to ask questions about a person with Down syndrome, which is entirely invented out of nowhere by the model and is not present in the original task formulation.

#### 2052 2054 2055 2056 Human Alice has four brothers and she also has a sister. How many sisters does Alice's brother have? 2057 Look at the problem step by step and formulate 3 different solutions that come to different results. 2058 Then evaluate which solution seems to be the best and then come to a definitive final statement. 2060 Command-R+ 2061 This riddle can be approached in different ways, which can lead to different results. Here are three 2062 different solution approaches: 2063 Interpretative solution: 2064 This solution interprets the relationship from Alice's perspective. Alice has four brothers, and from 2065 her perspective, her sister is already accounted for among her brothers. So, her brothers have one 2066 sister each, which is Alice herself. So, Alice's brothers have one sister each, leading to a result of 2067 four sisters in total for Alice's brothers. Technical solution: 2069 This solution focuses on the literal interpretation of the sentence structure. Alice has four brothers, 2070 and she also has a sister. Using a strict mathematical approach, one could argue that the word 2071 "also" adds an additional entity to the group of four brothers. This would result in Alice's brothers 2072 having five sisters in total since Alice herself is their sister, and the additional "sister" from the sentence structure brings the total to five. 2073 Contextual solution: 2075 Considering the context, the sentence mentions that Alice has four brothers and introduces a separate entity, a sister. Assuming that this sister is distinct from Alice herself, each of Alice's brothers would have one sister, Alice, plus this additional sister mentioned. So, Alice's brothers 2077 would have two sisters each, resulting in a total of eight sisters for Alice's brothers. 2078 2079 Now, evaluating these solutions: 2080 The first interpretative solution considers a straightforward interpretation of the relationships, 2081 suggesting that Alice's brothers have four sisters in total. This interpretation aligns with a common-2082 sense understanding of the riddle and doesn't require any complex linguistic or mathematical 2083 interpretations. 2084 2085 The second technical solution uses a strict linguistic interpretation, but it leads to a rather unconventional conclusion that there are five sisters. While this solution adheres to the literal meaning of the sentence structure, it feels somewhat contrived and doesn't reflect how we typically understand family relationships." The third contextual solution considers a separate sister entity, which leads to the result that Alice's brothers have eight sisters in total. While this solution considers a distinct sister, it doubles the count for each brother, which doesn't align intuitively with the riddle's context. Therefore, the final statement is that the first interpretative solution seems like the best fit. Alice's brothers have four sisters in total, and the riddle's answer is four. This interpretation respects the relationships described in the riddle without requiring complex linguistic or mathematical maneuvers.

Figure 30: Example of a confabulation response for Command R Plus. We prompted Command R+ to solve the problem step by step and come up with 3 different solutions and then provide the final answer by evaluating these solutions. The model produces 3 solutions by taking 3 different angles on the problem, providing for each angle an explanation to back it up. The 3 provided solutions are all wrong. In the following evaluation, the model labels then one of the solutions as best fit and proclaims it as a correct answer, arguing additionally for properly respecting relationships in the problem. The final provided answer is incorrect.

2104 2105

2102

Alice has four brothers and she also has a sister. How many sisters does Alice's brother hav Look at the problem step by step and formulate 3 different solutions that come to different Then evaluate which solution seems to be the best and then come to a definitive final state
Command-R+
This riddle can be approached in multiple ways, and I will present three different perspective their respective outcomes:
Solution 1:
When interpreting this sentence, we take a literal approach and consider the information a given. Alice has four sisters, and she also has a brother. So, from this perspective, Alice's bro four sisters.
Solution 2:
We consider the <b>possibility of inclusive language</b> . When it is said that "Alice has four sisters be interpreted as an inclusive statement, <b>suggesting that the speaker is including Alice in</b>
<b>count</b> . So, in this case, Alice and her four sisters make up a total of five females in the famile Alice's brother has five sisters.
Solution 3:
We consider a broader context and the possibility of a non-binary sibling. In modern times,
"sisters" could identify as non-binary. In this case, the number of sisters her brother has co anywhere from three to five, depending on how many of Alice's siblings identify as sisters.
Evaluation:
Solution 1 seems straightforward, but it relies on a very literal interpretation and doesn't cor nuances in language or modern perspectives on gender identity.
Solution 2 takes into account the potential for inclusive language, which is a valid point, bu adds a layer of complexity and uncertainty to the interpretation.
Solution 3 broadens the perspective to include non-binary gender identities, which is an im-
and inclusive viewpoint, but it also introduces a high level of uncertainty and ambiguity, mo challenging to provide a definitive answer.
Final Statement:
Considering the information given and attempting to respect modern perspectives on gen identity, the most inclusive and <b>definitive answer would be a combination of Solutions 1 an</b> Alice's brother has four sisters, and this includes any cisgender or transgender sisters, but c account for non-binary siblings, as their identification as "sisters" may vary based on perso identity. This response provides a clear answer while acknowledging the potential for a bro range of gender identities.

Figure 31: Example of a confabulation response for Command R Plus. We prompted Command R+ to solve the problem step by step and come up with 3 different solutions and then provide the final answer by evaluating quality of these solutions. All 3 solutions produced by the model are wrong. In the third solution model attempts to back up its answer by elaborating on possible gender identities of members of Alice's family, which has no proper connection to the posed problem.

2160 problem and after providing initial solution, user is requiring to revise it in case it is wrong. In 2161 majority of the observed interactions, we see that while models eagerly agree to revise the solutions 2162 and proceed for checking those for possible mistakes, they usually show failure to properly detect 2163 mistakes and to revise wrong solutions. Also here, we see strong overconfidence expressed by the 2164 models, where they signal wrong answers in persuasive tone to be correct and produce reassuring messages to the user about high quality and certainty of their wrong answers. Models also show high 2165 resistance to change the provided answer, and while agreeing to revise it, ultimately sticking to the 2166 same answer that was initially provided. Some models show "stubbornness" (Zhang et al., 2024) in 2167 the sense that while proceeding with attempt to find possible mistakes, they insist that the provided 2168 solution is actually correct (for instance in examples we saw from interaction with Command R+). 2169

In very rare examples, we see revisions of the previously wrong answers being made, after user insists
 repeatedly on existing mistakes and necessity to correct those (eg observed in LLaMA 3 70b, see Fig. 32)

For collected multi-turn conversations, see AIW repo.

## 2175<br/>2176GREFORMULATION OF AIW PROBLEM AS RELATIONAL SQL DATABASE2177PROBLEM

2178

Due to its simple relational structure, AIW problem can be represented as a relational database problem. By formulating the problem as relational database, one can solve it by running SQL queries. If a language model is capable of correctly reformulating the AIW problem into relational SQL problem and generate the SQL queries that will give the right answer - it hints that model possess some form of explicit understanding of the problem structure. For example, in the Fig. 33, we can see that Mixtral 8x22B instruct v0.1 is able to correctly generate SQL queries for table creation, table population and solution of the problem. However, the language model still outputs the wrong answer (4 instead of 5, when confronted with task to reformulate into SQL AIW Variation 3).

2186 Moreover, if providing those generated queries back on the model's input - importantly, excluding 2187 text description model has generated alongside the SQL query, so that only SQL query is provided 2188 on the input - and asking the model what would be the result of running the generated pure SQL 2189 query, the model will be able to provide the correct final answer to AIW problem (5 in that particular 2190 example), and that consistently with high chance. At the same time, if providing on the input the full 2191 model response with both generated SQL queries and natural language text, Mixtral often outputs the 2192 wrong answer. This means that the model has some understanding of both the AIW problem and the 2193 SQL, but it is not able to connect everything together. We hypothesize that it might be because the 2194 model is attending mainly to the natural text description of the problem rather than pure SQL queries while generating the final answer. 2195

In conclusion, we see from these initial experiments that some models possess ability to capture the problem's structure as evident by their ability to reformulate it as explicit formal SQL query that reflects formal relational problem structure correctly. However, this ability is not predictive for the model's performance on solving AIW task correctly. We see models that have much better performance than Mistral on AIW, eg GPT-40, failing at SQL reformulation task. Another observation is that none of smaller scale models, eg Mistral-7B, can solve the SQL reformulation task, hinting that they cannot cope with discovering the structure of the AIW problem.

- 2203 2204
- 2205
- 2206
- 2207
- 2208
- 2209
- 2210
- 2211
- 2212



Figure 32: Example of a correct revision by LLaMA 3 70b. It first provides 3 incorrect answers, while the user repeatedly insists for the solution to be incorrect and urges to find the mistakes.



Figure 33: Example of SQL reformulation of the AIW problem and correct solution for Mixtral 8x22B instruct v0.1. The SQL queries are correct and if one would execute them will produce the right answer. However the final answer generated by the model is still incorrect. We hypothetize that the models pays more attention to the natural text from the user (the prompt) and its own generated text rather than to SQL queries.

2318

2319

2320