A Modular Dataset to Demonstrate LLM Abstraction Capability

Anonymous ACL submission

Abstract

1

Large language models (LLMs) exhibit 2 impressive capabilities but struggle with 3 reasoning errors due to hallucinations and 4 flawed logic. To investigate their internal 5 representations of reasoning, we introduce 6 ArrangementPuzzle, a novel puzzle 7 dataset with structured solutions and 8 automated stepwise correctness 9 verification. We trained a classifier model 10 on LLM activations on this dataset and 11 found that it achieved over 80% accuracy 12 in predicting reasoning correctness, 13 implying that LLMs internally distinguish 14 between correct and incorrect reasoning 15 steps, with the strongest representations in 16 middle-late Transformer layers. Further 17 analysis reveals that LLMs encode abstract 18 reasoning concepts within the middle 19 activation layers of the transformer 20 architecture, distinguishing logical from 21 semantic equivalence. These findings 22 provide insights into LLM reasoning 23 mechanisms and contribute to improving 24 AI reliability and interpretability, thereby 25 offering the possibility to manipulate and 26 refine LLM reasoning. 27

28 1 Introduction

²⁹ Recently, large language models (LLMs) based on
³⁰ the Transformer architecture (Vaswani et al., 2017)
³¹ have demonstrated competence across a wide
³² range of domains, from reading comprehension to
³³ coding to mathematics. However, in domains such
³⁴ as these, LLMs can often generate incorrect
³⁵ responses due to hallucinations and incorrect
³⁶ reasoning (Rawte et al., 2023). Large reasoning
³⁷ models (LRMs) explicitly trained to produce
³⁸ accurate chains of thought such as o1 (OpenAI et al., 2024) and DeepSeek-R1 (DeepSeek-AI et al., 40 2025) promise to increase the effectiveness of

⁴¹ LLM reasoning. Even so, hallucinations and
⁴² reasoning inaccuracies remain, preventing these
⁴³ models from exceling at more complex, multi-step
⁴⁴ tasks such as PlanBench (Valmeekam et al., 2023;
⁴⁵ Valmeekam et al., 2024) designed to evaluate the
⁴⁶ planning and reasoning capabilities of LLMs.

As LLM adoption increases and they begin to 48 be used in increasingly critical applications, it 49 becomes more necessary to detect and prevent 50 such mistakes. The field of Explainable AI 51 (reviewed in Ferrando et al., 2024) seeks to tackle 52 this problem by trying to understand the inner 53 workings of LLMs, and using that information to 54 gain insight into how and why they go wrong.

For example, recent work suggests that LLMs 55 56 understand the difference between truth and 57 falsehood in factual statements (Azaria and 58 Mitchell, 2023) internally. Specifically, the LLM 59 contains representations of truthfulness that are 60 strongest in its intermediate layers. Probing these 61 layers for this representation can outperform 62 directly prompting the LLM about the truthfulness 63 of a statement (Liu et al., 2023). Interestingly, 64 larger LLMs appear to have more capacity for 65 abstraction, as their representations of truthfulness 66 generalize better across different data modalities 67 (Marks and Tegmark, 2024). Furthermore, 68 understanding how truthfulness and other similar 69 concepts are represented in the activations of 70 LLMs can allow us to manipulate those 71 representations to affect LLM behavior – for ⁷² example, causing it to be more honest or less angry 73 (Zou et al., 2023).

The recent advances in reasoning LLMs pose To the analogous questions for reasoning: do LLMs To have an internal concept of reasoning, and if so To how much abstraction are they capable of in Reasoning tasks? These are challenging questions to address with most publicly available reasoning datasets such as those used to train LRMs, which



Figure 1: Diagram of how stepwise verification of LLM reasoning is used to train a classifier on the LLM's activations over those steps. The classifier is trained to predict whether a particular reasoning step is accurate.

⁸¹ typically contain a variety of solutions with ⁸² different structures that makes analysis difficult.

In this paper, we introduce ArrangementPuzzle, a dataset suitable for analyzing LLM internal representations of reasoning. We use it to discover that LLMs have an innate understanding of when their reasoning steps are accurate, and they also have an internal representation of abstraction that separates logical from semantic equivalence.

90 2 ArrangementPuzzle Dataset

⁹¹ To test LLM reasoning representations, we ⁹² constructed ArrangementPuzzle, a customizable ⁹³ puzzle dataset. This dataset and the code for ⁹⁴ generating it are publicly available in our GitHub ⁹⁵ repository¹. Each puzzle contains some clues as to ⁹⁶ how a certain number of people are arranged, and ⁹⁷ what colors they are wearing. The LLM (Llama-⁹⁸ 3.1-8B-Instruct in our case) is then tasked with ⁹⁹ determining the full arrangement. An example ¹⁰⁰ puzzle is shown on the left of Figure 1.

The puzzles are guaranteed to have a unique solution given the available clues, and they are guaranteed to not have any redundant clues (that is, removing any single clue would result in a nonus unique solution). Inspired by (Mirzadeh et al., 106 2024) and (Jiang et al., 2024), the puzzles 107 randomize the exact names and colors used each
108 time, as well as the correct arrangement and clues
109 given. Our results differ from previous work in
110 that our dataset allows statement-level (rather than
111 solution-level) accuracy checking, and it contains
112 a deterministic solution generator capable of
113 generating full reasoning traces as it makes logical
114 deductions to solve the puzzle (see Figure 5 for
115 examples and Solution Generator for more
116 details). Additionally, our focus is on using LLM
117 activations derived from our dataset to understand
118 LLM internal representations of reasoning, rather
119 than performance benchmarking.

120 2.1 Reasoning Dataset

121 To this end, we evaluate the LLM on 10.000 122 prompts from puzzles with n = 2 people, of 123 which it gets 67.6% correct. We save its 124 activations for each generated token in each puzzle 125 to disk. Simultaneously, we use a regular 126 expression-based parser (also available on 127 GitHub) to extract reasoning statements from the 128 LLM's text output as it constructs partial claims 129 about the correct arrangement (for example, 130 Alice is not sitting on the right). This 131 works because our prompting approach 132 encourages the LLM to use very specific phrasing

¹ Link removed to preserve anonymity.

for its reasoning in line with the phrasing of thesolution generator (see LLM Prompting).

Once the parser has extracted such statements, the it evaluates them for correctness against the ground-truth solution to the puzzle. This then the creates a dataset of model activations labeled by whether they came from a correct reasoning step the or incorrect reasoning step.

141 3 Reasoning Classifier

¹⁴² We then use this dataset to train a classifier model ¹⁴³ (Figure 1). The model takes as input the LLM ¹⁴⁴ activations at a particular set of layers (typically 1 ¹⁴⁵ layer) at the last five token positions in a reasoning ¹⁴⁶ statement. It is trained to predict whether that ¹⁴⁷ reasoning statement was correct or incorrect. The ¹⁴⁸ classifier is a feedforward architecture that ¹⁴⁹ contains a single convolutional layer over token ¹⁵⁰ position followed by several fully-connected ¹⁵¹ layers (more details in Classifier Architecture). It ¹⁵² was trained on 7544 training puzzles and 1130 ¹⁵³ validation puzzles for 50 epochs and evaluated on ¹⁵⁴ a withheld set of 1141 test puzzles. The source ¹⁵⁵ code is available in our GitHub repository.

156 3.1 Isomorphic puzzles

157 One important property of our dataset is that 158 distinct puzzles can be generated by permuting ¹⁵⁹ some of the puzzle details (eg: names and colors). 160 This, together with LLM output randomness, 161 allow generating a variety of different data from a 162 smaller set of distinct logic puzzles. However, to 163 ensure our reasoning classifier learned 164 information about actual logical reasoning and 165 was not able to memorize patterns based on the 166 exact clues given, we took steps to ensure that the ¹⁶⁷ validation and test datasets contained distinct sets 168 of logical puzzles. To this end, we define two 169 puzzles to be isomorphic if there exists a 170 permutation of the clues, and substitutions of 171 names and colors, to transform one puzzle into the 172 other. Then we ensure our training, validation, and 173 testing datasets contain disjoint isomorphism 174 classes of puzzles.

175 3.2 Classifier Performance

The performance of the classifier as a function of the transformer layer it was trained on is shown in Figure 2. This high level of performance (>80% for most layers) demonstrates that the LLM does in fact contain distinct representations of correct and incorrect reasoning patterns. Additionally, these representations of reasoning appear to be strongest in the middle-late attention layers. This echoes previous findings (e.g. (Azaria and Mitchell, 2023)), which indicate that these layers also encode abstract representations of truth. We additionally ran an analysis where we trained a classifier on all of these top 5 performing layers (15, 17, 23, 25, and 30), but performance did not substantially increase (dashed red line), suggesting that these layers contain similar representations of the reasoning information.



Figure 2: Performance of the reasoning classifier trained on a specific layer evaluated on testing data.

194 4 Reasoning Information Abstraction

¹⁹⁵ Since our classifier is a neural network, it does not ¹⁹⁶ readily yield information on *how* the reasoning ¹⁹⁷ information is stored. However, it did reveal that ¹⁹⁸ the reasoning representation was strongest in the ¹⁹⁹ middle-late attention layers. Based on this, we ²⁰⁰ hypothesized that these layers might store more ²⁰¹ abstract representations of reasoning and logic.

To test this hypothesis, we used our dataset to develop an abstraction test. Specifically, we algorithmically generated solutions for all our puzzles using our solution generator (see Solution Generator) and evaluated the LLM on these solutions and stored its layer activations. Notably, the LLM itself was not used to generate the text. We then used this dataset to compare the LLM's activations across solutions.

211 4.1 Information Abstraction in LLMs

212 Running the solution generator on our set of
213 10,000 puzzles, we identified two sets of puzzles
214 of interest to studying abstraction in LLMs:

- 1. Logically distinct puzzles with solutions that contain identical lines of text, but at different places in the logical sequence.
- 2. Puzzles with isomorphic solutions that is, where one solution can be transformed into

215

216

217

218

219

220 and colors. 221



Figure 3: Examples of puzzles with distinct logical structure but identical text (Puzzles 1 and 2), and isomorphic puzzles with identical logical structure but distinct text (Puzzles 1 and 3). LLM activations on highlighted text are compared via correlation.

Activations

Activations

222

We randomly sampled 10,000 pairs of lines of 223 text from each of these categories - that is, identical lines of text from logically distinct 225 puzzles or corresponding but non-identical lines of 226 text from isomorphic puzzle solutions (examples 227 shown in Figure 3). For each pair of tokens in each 228 229 pair of lines of text sampled in this way, we 230 computed the correlation coefficient across hidden 231 activations for a given layer between the two 232 corresponding tokens. We then averaged this 233 together across tokens in each line and then 234 averaged it together across lines. We excluded

the other by replacing details like names 235 pairs of lines with different numbers of tokens (eg: 236 isomorphic lines where the token lengths of the ²³⁷ substituted fields were not the same).



Figure 4: Correlation of attention layer activations between lines of text in different puzzles.

In this manner, we were able to compute the 239 "abstraction level" of each layer in the Llama model by comparing the concordance of its activations between the two conditions (Figure 4). 242 Layers with high correlation in the "Identical" 243 condition and low correlation in the "Isomorphic" 244 condition, such as the 0th embedding layer, contain 245 246 mostly token-specific information. On the other 247 hand, layers where the reverse is true, such as 248 layers 10-20, contain more abstract information ²⁴⁹ about higher-level logical features of the puzzle.

250 5 Discussion

251 By leveraging ArrangementPuzzle, we trained a classifier that accurately distinguishes correct 252 from incorrect reasoning steps, confirming that ²⁵⁴ LLMs internally encode logical consistency. 255 Specifically, our study demonstrates that LLMs 256 possess internal representations of reasoning correctness, with the strongest signals emerging in 257 258 middle-late Transformer layers. Additionally, our 259 analysis of abstraction in model activations 260 suggests that LLMs differentiate between logical semantic equivalence. These findings 261 and 262 contribute to a deeper understanding of LLM ²⁶³ reasoning processes and may inform future efforts 264 to enhance model reliability, interpretability, and 265 trustworthiness. In particular, the fact that LLMs 266 already encode an innate representation of ²⁶⁷ reasoning may explain their ability to gain massive ²⁶⁸ improvements in reasoning capability via 269 distillation from LRMs, even without additional 270 reinforcement learning (DeepSeek-AI et al., 2025) 271 and even with as few as 1000 SFT samples for 272 distillation (Muennighoff et al., 2025).

273 6 Limitations and Ethical Concerns

274 6.1 Limitations

275 While our study provides insights into the internal 276 representations of reasoning in large language 277 models (LLMs), it has several limitations. First, 278 our analysis is restricted to a specific class of 279 structured logic puzzles, which may not fully 280 capture the complexity of reasoning required in 281 more real-world scenarios. diverse The ²⁸² constrained nature of our dataset, where solutions 283 follow deterministic patterns, may not generalize 284 to open-ended reasoning tasks that require 285 commonsense knowledge, probabilistic inference, 286 or multi-modal reasoning. Additionally, we did not 287 examine whether our classifier's success in 288 distinguishing correct from incorrect reasoning 289 steps generalizes to other types of reasoning 290 problems, preventing us from making claims 291 about the generalizability of the reasoning 292 representations we uncovered.

Another key limitation lies in our reliance on 293 probing techniques to analyze model activations. 294 While we demonstrate that certain Transformer 295 296 layers encode representations of reasoning 297 correctness and abstraction, our approach does not ²⁹⁸ provide a mechanistic explanation of how these 299 representations emerge or how they influence 300 downstream reasoning behavior. Furthermore, our 301 classifier is trained on activations from a single ³⁰² LLM architecture, and it remains unclear whether 303 these findings generalize across different model ³⁰⁴ families, sizes, or training paradigms. Future work 305 should explore more diverse reasoning 306 benchmarks, conduct broader cross-model 307 analyses including LRM models, and develop 308 methods for directly steering LLM reasoning 309 based on these learned representations.

310 7 Ethical Concerns

We do not anticipate any immediate ethical impact arising from our work. However, our work does highlight the potential for LLMs to prioritize pattern-matching over accuracy, as we have demonstrated that LLMs have an internal representation of reasoning accuracy yet output incorrect reasoning anyway. Additionally, latent LLM abstraction capabilities highlighted in this work suggest that it may be relatively easy to "jailbreak" open-weight LLMs with a small amount of fine-tuning into producing potentially dangerous output.

323 8 Use of AI Assistants

We employed the use of AI assistants, primarily ChatGPT (versions 40, 01, and 03-mini), to help enerate some of the code and text of this manuscript. The authors have examined all output of these assistants to ensure accuracy.

329 References

- 330 Amos Azaria and Tom Mitchell. 2023. The Internal
- 331 State of an LLM Knows When It's Lying. In
- 332 Findings of the Association for Computational
- 333 Linguistics: EMNLP 2023, pages 967–976,
- ³³⁴ Singapore. Association for Computational
- 335 Linguistics.
- 336 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei
- 337 Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
- 338 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi,
- 339 Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
- 340 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, et
- al. 2025. DeepSeek-R1: Incentivizing Reasoning
- 342 Capability in LLMs via Reinforcement Learning.
- 343 arXiv:2501.12948 [cs].
- Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and
 Marta R. Costa-jussà. 2024. A Primer on the Inner
 Workings of Transformer-based Language Models.
 arXiv:2405.00208 [cs].
- 348 Bowen Jiang, Yangxinyu Xie, Zhuoqun Hao,
- 349 Xiaomeng Wang, Tanwi Mallick, Weijie J Su,
- 350 Camillo Jose Taylor, and Dan Roth. 2024. A Peek
- 351 into Token Bias: Large Language Models Are Not
- 352 Yet Genuine Reasoners. In Proceedings of the 2024
- 353 Conference on Empirical Methods in Natural
- 354 Language Processing, pages 4722-4756, Miami,
- 355 Florida, USA. Association for Computational
- 356 Linguistics.
- 357 Kevin Liu, Stephen Casper, Dylan Hadfield-Menell,
- 358 and Jacob Andreas. 2023. Cognitive Dissonance:
- 359 Why Do Language Model Outputs Disagree with
- 360 Internal Representations of Truthfulness?. In
- 361 Proceedings of the 2023 Conference on Empirical
- 362 Methods in Natural Language Processing, pages
- 363 4791–4797, Singapore. Association for
- 364 Computational Linguistics.
- 365 Samuel Marks and Max Tegmark. 2024. The
- 366 Geometry of Truth: Emergent Linear Structure in
- 367 Large Language Model Representations of
- 368 True/False Datasets. arXiv:2310.06824 [cs].

369 Iman Mirzadeh, Keivan Alizadeh, Hooman

- 370 Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad
- 371 Farajtabar. 2024. GSM-Symbolic: Understanding the
- 372 Limitations of Mathematical Reasoning in Large
- 373 Language Models. arXiv:2410.05229 [cs].

374 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang

375 Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke

376 Zettlemoyer, Percy Liang, Emmanuel Candès, and

377 Tatsunori Hashimoto. 2025. s1: Simple test-time

³⁷⁸ scaling. arXiv:2501.19393 [cs].

379 OpenAI: Aaron Jaech, Adam Kalai, Adam Lerer,

380 Adam Richardson, Ahmed El-Kishky, Aiden Low,

381 Alec Helyar, Aleksander Madry, Alex Beutel, Alex

382 Carney, Alex Iftimie, Alex Karpenko, Alex Tachard

³⁸³ Passos, Alexander Neitz, Alexander Prokofiev,

³⁸⁴ Alexander Wei, Allison Tam, Ally Bennett, et al.

385 2024. OpenAI o1 System Card. arXiv:2412.16720

386 [cs].

³⁸⁷ Vipula Rawte, Amit Sheth, and Amitava Das. 2023.

388 A Survey of Hallucination in Large Foundation

389 Models. arXiv:2309.05922 [cs].

390 Karthik Valmeekam, Matthew Marquez, Alberto

³⁹¹ Olmo, Sarath Sreedharan, and Subbarao

392 Kambhampati. 2023. PlanBench: An Extensible

393 Benchmark for Evaluating Large Language Models

394 on Planning and Reasoning about Change. Advances

395 in Neural Information Processing Systems,

396 36:38975-38987.

³⁹⁷ Karthik Valmeekam, Kaya Stechly, and Subbarao

398 Kambhampati. 2024. LLMs Still Can't Plan; Can

399 LRMs? A Preliminary Evaluation of OpenAI's o1 on

⁴⁰⁰ PlanBench. arXiv:2409.13373 [cs].

401 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob

402 Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz

403 Kaiser, and Illia Polosukhin. 2017. Attention is All

404 you Need. In Advances in Neural Information

405 Processing Systems, volume 30. Curran Associates,

406 Inc.

407 Andy Zou, Long Phan, Sarah Chen, James Campbell,

⁴⁰⁸ Phillip Guo, Richard Ren, Alexander Pan, Xuwang

409 Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,

410 Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan

411 Wang, Alex Mallen, Steven Basart, Sanmi Koyejo,

412 Dawn Song, Matt Fredrikson, et al. 2023.

413 Representation Engineering: A Top-Down Approach

414 to AI Transparency. arXiv:2310.01405 [cs].

415 A Solution Generator

416 The algorithm begins by initializing all possible 417 assignments for people, positions, and colors, 418 effectively considering every permutation. For 419 example, at the beginning it would initialize 420 Andrew as being at one of positions left or 421 right. It then iteratively applies the clues to 422 eliminate invalid combinations, updating the sets 423 of possibilities for each entity based on the 424 constraints provided by the clues. Through 425 constraint propagation, the algorithm refines 426 these possibilities by intersecting sets and 427 removing options when only one remains for a 428 given entity. This iterative process continues until 429 a unique solution is found, solving the puzzle. 430 Importantly, every time the algorithm updates its ⁴³¹ internal possibilities, it outputs a reasoning step in 432 text format.

433

434 **B** LLM Prompting

⁴³⁵ To prompt the LLM to solve our puzzles, we use ⁴³⁶ a 3-shot approach where we append the actual ⁴³⁷ puzzle (to be solved) to three example ⁴³⁸ puzzle/solution pairs. These solutions were ⁴³⁹ generated by our solution generator. This primes ⁴⁴⁰ the LLM to reason through the puzzles using



441 similar logic. The full 3-shot prompt is available
442 on our GitHub repository in the file
443 prompt.txt.

444

445 C Classifier Architecture

⁴⁴⁶ Our classifier uses a convolutional layer with 128 ⁴⁴⁷ output channels and kernel size 3 that convolves ⁴⁴⁸ over the five token positions, treating each ⁴⁴⁹ individual activation (of the 4096 hidden units) ⁴⁵⁰ and each layer as a different channel. These ⁴⁵¹ outputs are then fed through three fully-⁴⁵² connected layers with hidden sizes 256 and 128, ⁴⁵³ before finally producing a single logit which is ⁴⁵⁴ then passed through a sigmoid to produce the ⁴⁵⁵ final prediction.

456

Figure 5: Example puzzles with solutions from our deterministic solution generator. The left two puzzles are isomorphic to each other.