# On the Notion that Language Models Reason

**Bertram Højer**
Department of Computer Science
IT University of Copenhagen
berh@itu.dk

## Abstract

Language models (LMs) are said to be exhibiting reasoning, but what does this entail? We assess definitions of *reasoning* and how key papers in the field of natural language processing (NLP) use the notion and argue that the definitions provided are not consistent with how LMs are trained, process information, and generate new tokens. To illustrate this incommensurability we assume the view that transformer-based LMs implement an *implicit* finite-order Markov kernel mapping contexts to conditional token distributions. In this view, reasoning-like outputs correspond to statistical regularities and approximate statistical invariances in the learned kernel rather than the implementation of explicit logical mechanisms. This view is illustrative of the claim that LMs are "statistical pattern matchers" and not *genuine reasoners* and provides a perspective that clarifies why reasoning-like outputs arise in LMs without any guarantees of logical consistency. This distinction is fundamental to how epistemic uncertainty is evaluated in LMs. We invite a discussion on the importance of how the computational processes of the systems we build and analyze in NLP research are described.

## 1 Introduction

Language models (LMs) are widely marketed for their ability to solve complex tasks that supposedly require *reasoning*. However, it is still debated whether LMs engage in structured reasoning or whether they are merely replicating statistical relations from the data on which they are trained [Bender et al., 2021, Huang and Chang, 2023, Mirzadeh et al., 2024, Kambhampati et al., 2025]. Flagship models developed by companies such as OpenAI, Anthropic, and Alibaba are labeled as *reasoning* models that generate long "chains of thought" before generating the final output. Models from the Qwen-series output these traces in designated $< think >$ tags [Yang et al., 2025].

In this paper we assess standard definitions of reasoning and the current framing of reasoning in the NLP literature, and argue that the definitions of reasoning are incommensurable with transformer computations. The use of ill-fitting terminology is problematic due to connotations attached to a notion such as *reasoning*, which has a rich tradition in fields such as philosophy, AI, and psychology [see e.g. Tversky and Kahneman, 1974, 1981, 1986, Albus, 1991, Wallace and Kiesewetter, 2024, Jiang et al., 2024]. *Reasoning* is a key criteria for a system to be considered intelligent, and it is therefore a crucial aspect of the scientific aims pursued in the field of *AI* [Højer et al., 2025].

We describe LMs as *implicit Markov kernels* similar to Zekri et al. [2025], but use this formalization to frame a discussion of the notion of *reasoning* and *inference*. This lens challenges the description of transformer computations as commensurable with the process of "thinking in a logical and systematic manner". Multiple papers have experimentally illustrated the logical shortcomings of LMs by showing the failure-modes when it comes to simple "*reasoning*" tasks [Nezhurina et al., 2024, Mirzadeh et al., 2024, Jiang et al., 2024].

## 2 Reasoning and Natural Language Processing

*Reasoning* is increasingly being used to describe a specific type of output generated by LMs, namely the "thought traces" as they have been labeled after Wei et al. [2023] and earlier work by Nye et al. [2021]. Concurrently, it is being debated whether LMs can be said to be doing anything akin to *reasoning*. A key issue seems to be that the grounds for disagreement are not necessarily clear. We argue that addressing the well-formedness of the question of *reasoning* in LMs is of importance to *AI* research, while others may claim that it is enough that a system seems to be *reasoning* if it results in more correct outputs. We return to this in section 4.

### 2.1 Definitions of Reasoning

It is not that there are no definitions of *reasoning*. To reason about the natural world has a history which is documented as far back as to the ancient Greeks, and while the history of *reasoning* is fundamental to science, it is beyond the scope of this paper to lay it out in detail. Surveying definitions can be summarized to *reasoning* being something along the lines of "*thinking about something in a logical manner*" or "*the process of thinking about something logically and systematically to make an inference*". In psychology and philosophy one finds definitions that are not too far off from these statements (see e.g. Johnson-Laird [2008, 2010], Over and Evans [2024], Albus [1991], Portoraro [2025]). Of special interest to the ongoing discussion in NLP and *AI*, Tversky and Kahneman [1981, 1986] discussed *the principle of invariance* in the theory of rational decision-making to illustrate that a rational decision-maker should be invariant to various biases and logical fallacies.

If we look at its use in the NLP literature, the term has been increasing rather rapidly over the past five years. We trace this narrative in key papers from the NLP literature that have been influential for the literature on *reasoning* (section 2.2).
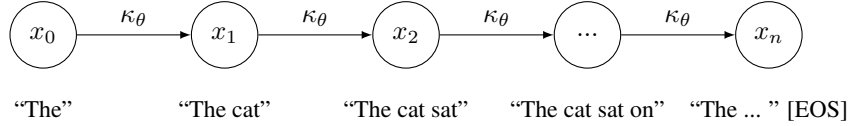
Huang and Chang [2023] surveyed recent papers on LM reasoning. They define *reasoning* both as a "cognitive process that involves using evidence, arguments, and logic to arrive at conclusions or make judgment" and as "the process of thinking about something in a logical and systematic way, using evidence and past experiences to reach a conclusion or make a decision". Later they demarcate between *formal reasoning*, which is a "systematic and logical process that follows a set of rules and principles (...)", and *informal reasoning*, which is "a less structured approach that relies on intuition, experience, and common sense to draw conclusions and solve problems (...)". It should be clear to the reader that **what is meant by *reasoning* is not clear at all**.

### 2.2 A Timeline of LM Reasoning

The labeling of autoregressive LM generation as *reasoning* has arisen in the last decade. Vaswani et al. [2017] introduced the transformer with self-attention, enabling efficient scaling of language models in terms of both model parameters and dataset size. Radford et al. [2019] then introduced GPT-2 calling it a multi-task learner, referring to the finding that the model was capable of solving a vast variety of traditional NLP tasks that one would usually have to fine-tune a specific model for. Brown et al. [2020] introduced the follow-up model, GPT-3, as a few-shot learner - namely one that is capable of learning by example. Research on the capabilities of these models followed soon after with Wei et al. [2023] arguing that Chain-of-Thought (CoT) prompting elicits *reasoning* in LMs. Kojima et al. [2023] argue that LMs are not only few-shot reasoners, but that they are zero-shot reasoners and Wang et al. [2023] show that self-consistency improves CoT *reasoning* in LMs. This timeline indicates how "reasoning" has become a function of the format of the model output as opposed to a logical process. It has become evident how important the *right* data is for enabling "reasoning". Early examples hereof are the *TinyStories* models which were trained with fully synthetic data [Eldan and Li, 2023], the Phi-family of models that are trained on especially curated datasets to excel at python coding tasks and other forms of "reasoning" tasks [Gunasekar et al., 2023, Li et al., 2023], and the Orca models which were trained specifically on data containing "reasoning" strategies generated by GPT-4 [Mukherjee et al., 2023, Mitra et al., 2023]. More recently we have seen models that are trained initially on extensive examples of "reasoning" and then optimized using various forms of reinforcement learning to incentivize the generation of long "reasoning" traces. Guo et al. [2025] provide the technical details of this process. They furthermore claim that reasoning *emerges* spontaneously from the training procedures and Wei et al. [2022] and Bowman [2023] claim that certain *behaviors emerge* unpredictably.

## 2.3 Logic and Large language Models

The autoregressive aspect and relatedly the finding that CoT improves the performance of an LM is used to argue that LMs *reason* over data [Brown et al., 2020, Wei et al., 2023, Kojima et al., 2023]. But it is important to consider the fact that autoregressive generation applies the same model $\kappa_\theta$ iteratively in the generation of new tokens as illustrated below. Assuming no sampling during generation one would get the exact same output sequence when prompting an LM with $x_1$ ("The cat") and $x_2$ ("The cat sat") if $\kappa_\theta(\text{"sat"}|\text{"the cat"}) \geq \kappa_\theta(x|\text{"the cat "})$ for all $x \in \Sigma$. When a different sampling technique is applied the diagram would appear to be branching.



"The"        "The cat"        "The cat sat"        "The cat sat on"        "The ... " [EOS]

We observe a narrative in the discourse on NLP and AI where LMs are framed as general-purpose technologies capable of *reasoning* and generalizing to almost any type of text-based problem (see section 2.2). However, it is also clear from this literature that the "reasoning" is strictly a feature of the format of the model outputs and not related to the logic of the system generating the outputs - although definitions imply otherwise.

It is also clear from the definitions that logic and systematicity are key factors of *reasoning*. We thus ask: are LMs logical and systematic? If one wishes to state that an LM *reasons* and generates logical statements, the logic that is enforced would **only** be the logic of the data, so to speak. However, it is not known whether LMs can — in principle — enforce any logical structures, and recent research displays evidence to the contrary [Peng et al., 2024]. As an example and to illustrate our point we describe an LM formally as a **Markov process**.

# 3 Language Models and Markov Processes

To illustrate why LMs do not *reason* in the sense of the definitions emphasized in this paper, it can be instructive to view LMs through a different lens and revisit the fundamentals. Shannon explains in his 1948 paper:

> [A discrete source] will choose successive symbols according to certain probabilities depending, in general, on preceding choices as well as the particular symbols in question. (...) a mathematical model of a system which produces such a sequence of symbols governed by a set of probabilities, is known as a stochastic process. [Shannon, 1948, p. 5].

## 3.1 Markov processes

The Markov property requires that a stochastic process is *memoryless*, meaning that its future states are entirely dependent on the current state and not on its history [Murphy, 2022]. For a fixed maximum context length $L$ and finite vocabulary $V$, an autoregressive LM with the state $s_t = x_{t-L:t-1}$ induces an *order-L* Markov chain over $V^L$.

A *Markov kernel* is a measurable mapping $\kappa : X \to \mathcal{P}(Y)$, that assigns to each input $x \in X$ a probability distribution over possible outputs in $Y$. Thus, for any context $x$, the kernel $\kappa(\cdot \mid x)$ defines the conditional probability distribution of outcomes. An LM is a parameterized kernel $\kappa_\theta$. An invariance in a kernel refers to transformations of the input space that leave the output distribution unchanged. For example, if two contexts are approximately logically equivalent, one might expect the kernel to assign (almost) the same continuation distribution:

$$x_i \approx x_j \implies \kappa_\theta(\cdot \mid x_i) \approx \kappa_\theta(\cdot \mid x_j).$$

In an LM the kernel is implicit. It is not defined explicitly as a transition matrix, but is instead instantiated by the parameters of the model $\kappa_\theta$. The training objective maximizes the likelihood of observed sequences by minimizing the cross-entropy between the empirical data distribution on which the model is trained and the model distribution. This learning objective thus approximates the

empirical conditional distribution of tokens. Specifically, an LM is a Markov model on a finite state space as both the vocabulary and context of a model are finite [Zekri et al., 2025].

But this objective does not enforce global invariances or logical implications, although the window of what is local is widening as models are trained with larger contexts. At most, it *loosely enforces strong regularities in the data* as invariance. When an LM "reasons" and applies a certain logic it corresponds to regularities in the kernel.

Crucially, even approximate invariance is not guaranteed and LMs often violate logical consistency [Nezhurina et al., 2024, Mirzadeh et al., 2024, Shojaee et al., 2025]. This illustrates the fact that invariances in transformer kernels are statistical artifacts of the training data distribution and not structural properties of the architecture or learning objectives that enforce a logic on model outputs. This is evidenced by how important data has become for training models that are better at "reasoning" as discussed in section 2.2.

### 3.2 "Logic" as Approximate Invariances in the Kernel

Rather than examining *reasoning* as a logical process, researchers usually measure how well a model adheres to logical and deductive implications [Nezhurina et al., 2024, Mirzadeh et al., 2024, Mondorf and Plank, 2024].

In this context, an invariance in a parameterized Markov kernel $\kappa_\theta$ refers to the property that certain transformations of input contexts $x$ lead to approximately the same conditional output distribution. Approximate invariances of the kernel can be viewed as the model's implicit expression of epistemic stability, and we can interpret *reasoning*-like consistency as such. Similar, but different, contexts that have similar "structure" should map to similar continuation distributions.

Let $V$ be a finite vocabulary and $X = V^{\leq L}$ be the space of token sequences up to a context window of size $L$. An LM parameterized by $\theta$ then defines a Markov kernel

$$\kappa_\theta : X \to \Delta(V), \quad \kappa_\theta(\cdot \mid x) = p_\theta(x_t \mid x_{t-L:t-1}),$$

which maps a bounded context $x$ to the probability simplex over the vocabulary $\Delta(V)$. As $L$ is finite the model is an *order-$L$* Markov chain.

In this setting, "reasoning" is the ability to preserve inferential relations between symbolic structures across context, which is viewed as a property of *invariance* of the kernel given a transformation of the input. One might then distinguish between two related types of invariances: *transformation invariances* and *inferential invariances*.

Let $T$ be a group of logic-preserving transformations. A model is *approximately transformation-invariant* under $T$ if

$$\forall\, t \in T, \quad \mathbf{V_T}\big(\kappa_\theta(\cdot \mid x),\, \kappa_\theta(\cdot \mid t(x))\big) \leq \epsilon_T,$$

where $\mathbf{V_T}$ denotes total variation distance between distributions on $V$ for a group of transformations $T$. In words, an LM's predictions should not change substantially under transformations that preserve logic.

Suppose now we have a logical reasoning rule $r$ (such as *modus ponens*); define a relation $\mathcal{R}_r = \{(x,y) : x \xRightarrow{r} y\}$. A model exhibits *inferential invariance* with respect to a rule $r$ if

$$\forall\, (x,y) \in \mathcal{R}_r, \quad \kappa_\theta(y \mid x) \geq 1 - \delta_r,$$

and the relation holds under the aforementioned transformations. What do these measures mean conceptually? $\epsilon_T$ and $\delta_r$ are an attempt to capture the *approximate invariances* of a model. An $\epsilon_T > 0$ indicates that a model is somewhat sensitive to irrelevant transformations and $\delta_r > 0$ indicates the degree to which "reasoning" is imperfect. LMs are usually said not to *reason* because they are not invariant to irrelevant transformations on $x$ [e.g. Mirzadeh et al., 2024].

LMs are trained to match an empirical probability distribution of a dataset $\mathcal{D} \subset X$ and do not directly optimize for, or enforce, these invariances. They instead reproduce regularities in $\mathcal{D}$ that approximate inferential structure. If a change in model output results in a different semantic interpretation, the apparent *reasoning* surely reflects data-induced regularities rather than logical inference.

To make research into transformation invariance directly comparable between papers, notions such as these could be operationalised and applied to measure epistemic fidelity in papers such as Nezhu-

rina et al. [2024], Mirzadeh et al. [2024], Jiang et al. [2024]. However, we emphasize that this would do no more than illustrate the stability of $\kappa_\theta$.

### 3.3 Inference

Now, we postulate that we could change almost every mention of "reasoning" in this paper with **inference**. Inference does not carry the psychological connotations of a term such as *reasoning* and is well-defined both within statistics and machine learning. In statistics **inference** is passing from sample data to a generalization, and it is usually equivalent to the notion of *prediction* in machine learning [Murphy, 2022].

In no sense is what we have discussed as "reasoning" different from the notion of *inference*, begging the question of why researchers speak of *reasoning* in a formal science. From a purely structuralist perspective, the view that words get their meaning from the contexts in which they are used, this would indicate that *reasoning* $\approx$ *inference*. But this would be false by the definitions we have presented in this paper. Analyzing LM operations should therefore be kept a science of systematic natural language **inference** and not one of *reasoning*. Furthermore, this framing makes the object of analysis clearer: it is not about imbuing *reasoning* but about ensuring logical and systematic *inference*.

### 3.4 Research Program

Empirical research can help elucidate how epistemic uncertainty manifests as invariance in model predictions when logical structures in the data are controlled. We aim to show that "reasoning" in LMs should rather be framed as **inference**, and to illustrate how the proposed metrics vary when certain logical structures (regularities in $\mathcal{D}$) are enforced in the datasets used to train LMs. We aim to build simple synthetic datasets wherein we control the logical structure of tokens; with these datasets we then train small toy transformers to establish how the theory generalizes to the transformer architecture, which computes $\Delta(V)$ quite differently from a discrete Markov model.

It is clear that generating longer CoT traces results in better benchmark performance. But this is a result of more expressive models due to a higher parameter count and CoT (as shown by Merrill and Sabharwal [2024]), and thus more accurate **inference** over the distribution of data on which an LM is trained. Stechly et al. [2025] showed that even semantically invalid CoT traces lead to better performance on certain tasks. Based on this we additionally aim to analyze the operations of autoregressive LM **inference** empirically to elucidate the logic of such a result.

## 4 Discussion

The discrete Markov models discussed in this paper are simple in terms of their expressive power when compared to modern LMs. However, most LMs can theoretically be framed as Markov chains [Zekri et al., 2025]. This provides fertile ground for theorizing about the nature of such models; in our case, specifically whether they can be said to *reason* given current definitions of *reasoning*. And furthermore, whether that is the right question to ask; formulating the right question is not separate from the science, it is part of it. Conflating *reasoning* and **inference** can obscure how epistemic uncertainty is represented, and thus sticking with clear definitions and terminology enables researchers to be rigid.

The implicit nature of the kernel of an LM is not insignificant as it obscures the structure of the kernel. In an LM the input is embedded via a learned embedding function $e^{in} : X \rightarrow \mathbb{R}^d$, modified by a positional embedding, and then transformed iteratively by a combination of an attention-mechanism and a feed-forward block before the application of a so-called "de-embedding" function $e^{out} : \mathbb{R}^d \rightarrow X$ [Vaswani et al., 2017]. These operations are done in a real-valued continuous space with induced non-linearities imbuing a certain structure on the space of the operations.

We ask: do (and should) these added complexities fundamentally change the conceptualization of the operations computed by a model? Neural networks trained for extracting word embeddings have e.g. been proven to approximate a matrix factorization of the PMI between terms in a corpus; a factorization problem to which there is a unique solution, namely the singular value decomposition [Levy and Goldberg, 2014, Goldberg and Levy, 2014]. While it has been shown that the attention

mechanism increases expressivity and that CoT generation makes a model strictly more expressive [Merrill and Sabharwal, 2024], it is nonetheless a model optimizing a loss-function that yields the most likely next token given a corpus.

With this paper we aim to ground the debate of logical *reasoning* in LMs in terms of invariances in the learned kernel as opposed to logical and systematic *thinking*. If we can reduce *reasoning* to **inference** and something like a description of invariances in a kernel, we have a solid foundation for demystifying and discussing the capabilities and limitations of an LM.

## 5 Conclusion

In conclusion, we reiterate and emphasize the result that an LM can be seen as implementing an **implicit Markov kernel**. Based on this result and an investigation of current use of the notion of *reasoning*, we argue that *reasoning* in LMs is in no meaningful way separate from the notion of **inference**. Framing reasoning as inference under epistemic uncertainty refocuses the debate from vague notions of *reasoning* to measurable epistemic properties of model-based inference. We formalize simple metrics of *invariance*, discuss how they relate to the notion of *reasoning* in current research and suggest a research program to understand the operations of LMs that could enable *logical inference* while also elucidating the limitations of LMs.

## References

J.S. Albus. Outline for a theory of intelligence. *IEEE Transactions on Systems, Man, and Cybernetics*, 21(3), June 1991. ISSN 00189472. URL http://ieeexplore.ieee.org/document/97471/.

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, Virtual Event Canada, March 2021. ACM. ISBN 978-1-4503-8309-7. URL https://dl.acm.org/doi/10.1145/3442188.3445922.

Samuel R. Bowman. Eight Things to Know about Large Language Models, April 2023. URL http://arxiv.org/abs/2304.00612.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners, July 2020. URL http://arxiv.org/abs/2005.14165.

Ronen Eldan and Yuanzhi Li. TinyStories: How Small Can Language Models Be and Still Speak Coherent English?, May 2023. URL http://arxiv.org/abs/2305.07759.

Yoav Goldberg and Omer Levy. word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method, February 2014. URL http://arxiv.org/abs/1402.3722.

Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks Are All You Need, October 2023. URL http://arxiv.org/abs/2306.11644.

Daya Guo, Dejian Yang, Haowei Zhang, and Junxiao et. al. Song. DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning. 645(8081):633–638, 2025. ISSN 1476-4687. URL https://www.nature.com/articles/s41586-025-09422-z.

Jie Huang and Kevin Chen-Chuan Chang. Towards Reasoning in Large Language Models: A Survey. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, Toronto, Canada, July 2023. Association for Computational Linguistics. URL `https://aclanthology.org/2023.findings-acl.67`.

Bertram Højer, Terne Thorn Jakobsen, Anna Rogers, and Stefan Heinrich. Research Community Perspectives on "Intelligence" and Large Language Models. In *The Findings of the Association of Computations Linguistics*, Vienna, 2025.

Bowen Jiang, Yangxinyu Xie, Zhuoqun Hao, Xiaomeng Wang, Tanwi Mallick, Weijie J Su, Camillo Jose Taylor, and Dan Roth. A Peek into Token Bias: Large Language Models Are Not Yet Genuine Reasoners. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4722–4756, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL `https://aclanthology.org/2024.emnlp-main.272/`.

Philip Johnson-Laird. *How We Reason*. Oxford University Press, October 2008. ISBN 978-0-19-955133-0. URL `https://academic.oup.com/book/11984`.

Philip Johnson-Laird. Mental models and human reasoning. 2010. URL `https://www.pnas.org/doi/10.1073/pnas.1012933107`.

Subbarao Kambhampati, Kaya Stechly, Karthik Valmeekam, Lucas Saldyt, Siddhant Bhambri, Vardhan Palod, Atharva Gundawar, Soumya Rani Samineni, Durgesh Kalwar, and Upasana Biswas. Stop Anthropomorphizing Intermediate Tokens as Reasoning/Thinking Traces!, May 2025. URL `http://arxiv.org/abs/2504.09762`.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large Language Models are Zero-Shot Reasoners, January 2023. URL `http://arxiv.org/abs/2205.11916`.

Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, volume 2 of *NIPS'14*, Cambridge, MA, USA, December 2014. MIT Press.

Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks Are All You Need II: phi-1.5 technical report, September 2023. URL `http://arxiv.org/abs/2309.05463`.

William Merrill and Ashish Sabharwal. The Expressive Power of Transformers with Chain of Thought, April 2024. URL `http://arxiv.org/abs/2310.07923`.

Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models, October 2024. URL `http://arxiv.org/abs/2410.05229`.

Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Codas, Clarisse Simoes, Sahaj Agarwal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Aggarwal, Hamid Palangi, Guoqing Zheng, Corby Rosset, Hamed Khanpour, and Ahmed Awadallah. Orca 2: Teaching Small Language Models How to Reason, November 2023. URL `http://arxiv.org/abs/2311.11045`.

Philipp Mondorf and Barbara Plank. Beyond Accuracy: Evaluating the Reasoning Behavior of Large Language Models – A Survey, August 2024. URL `http://arxiv.org/abs/2404.01869`. arXiv:2404.01869 [cs] Read_Status: Done Read_Status_Date: 2024-10-25T07:58:55.307Z.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive Learning from Complex Explanation Traces of GPT-4, June 2023. URL `http://arxiv.org/abs/2306.02707`.

Kevin P. Murphy. *Probabilistic machine learning: an introduction*. Adaptive computation and machine learning. The MIT Press, Cambridge, Massachusetts London, England, 2022. ISBN 978-0-262-04682-4 978-0-262-36930-5.

Marianna Nezhurina, Lucia Cipolina-Kun, Mehdi Cherti, and Jenia Jitsev. Alice in Wonderland: Simple Tasks Showing Complete Reasoning Breakdown in State-Of-the-Art Large Language Models, 2024. URL `http://arxiv.org/abs/2406.02061`.

Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. Show Your Work: Scratchpads for Intermediate Computation with Language Models, November 2021. URL `http://arxiv.org/abs/2112.00114`.

David E. Over and Jonathan St B. T. Evans. *Human reasoning*. Elements in philosophy of mind. Cambridge University press, Cambridge New York (N.Y.), 2024. ISBN 978-1-009-49531-8.

Binghui Peng, Srini Narayanan, and Christos Papadimitriou. On Limitations of the Transformer Architecture. 2024. URL `https://par.nsf.gov/servlets/purl/10580944`.

Frederic Portoraro. Automated Reasoning. In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2025 edition, 2025. URL `https://plato.stanford.edu/archives/sum2025/entries/reasoning-automated/`.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners. 2019.

C E Shannon. A Mathematical Theory of Communication. *The Bell System Technical Journal*, Vol. 27, 1948.

Parshin Shojaee, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity. 2025.

Kaya Stechly, Karthik Valmeekam, Atharva Gundawar, Vardhan Palod, and Subbarao Kambhampati. Beyond Semantics: The Unreasonable Effectiveness of Reasonless Intermediate Tokens, May 2025. URL `http://arxiv.org/abs/2505.13775`.

Amos Tversky and Daniel Kahneman. Judgment under Uncertainty: Heuristics and Biases: Biases in judgments reveal some heuristics of thinking under uncertainty. *Science*, 185(4157):1124–1131, September 1974. ISSN 0036-8075, 1095-9203. URL `https://www.science.org/doi/10.1126/science.185.4157.1124`.

Amos Tversky and Daniel Kahneman. The Framing of Decisions and the Psychology of Choice. 211, 1981.

Amos Tversky and Daniel Kahneman. Rational Choice and the Framing of Decisions. *The Journal of Business*, 59(4):S251–S278, 1986. ISSN 0021-9398. URL `https://www.jstor.org/stable/2352759`.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need, August 2017. URL `http://arxiv.org/abs/1706.03762`.

R. Jay Wallace and Benjamin Kiesewetter. Practical Reason. In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, fall 2024 edition, 2024. URL `https://plato.stanford.edu/archives/fall2024/entries/practical-reason/`.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-Consistency Improves Chain of Thought Reasoning in Language Models, March 2023. URL `http://arxiv.org/abs/2203.11171`.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent Abilities of Large Language Models, June 2022. URL `https://arxiv.org/abs/2206.07682v2`.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, January 2023. URL `http://arxiv.org/abs/2201.11903`.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, and Bo et. al. Zheng. Qwen3 Technical Report, May 2025. URL `http://arxiv.org/abs/2505.09388`.

Oussama Zekri, Ambroise Odonnat, Abdelhakim Benechehab, Linus Bleistein, Nicolas Boullé, and Ievgen Redko. Large Language Models as Markov Chains, February 2025. URL `http://arxiv.org/abs/2410.02724`.