

Triangulating LLM Progress through Benchmarks, Games, and Cognitive Tests

Filippo Momentè^{1*}, Alessandro Suglia², Mario Giulianelli³, Ambra Ferrari¹,
Alexander Koller⁴, Oliver Lemon², David Schlangen⁵, Raquel Fernández⁶, Raffaella Bernardi¹

¹University of Trento, ²Heriot-Watt University, ³ETH Zürich,
⁴Saarland University, ⁵University of Potsdam, ⁶University of Amsterdam

Abstract

We examine three evaluation paradigms: large question-answering benchmarks (e.g., MMLU and BBH), interactive games (e.g., Signalling Games or Taboo), and cognitive tests (e.g., for working memory or theory of mind). First, we investigate which of the former two—benchmarks or games—is most effective at discriminating LLMs of varying quality. Then, inspired by human cognitive assessments, we compile a suite of targeted tests that measure cognitive abilities deemed essential for effective language use, and we investigate their correlation with model performance in benchmarks and games. Our analyses reveal that interactive games are superior to standard benchmarks in discriminating models. Causal and logical reasoning correlate with both static and interactive tests, while differences emerge regarding core executive functions and social/emotional skills, which correlate more with games. We advocate the development of new interactive benchmarks and targeted cognitive tasks inspired by assessing human abilities but designed specifically for LLMs.

1 Introduction

Evaluating LLMs is critical to track progress, identify blind spots, and ultimately advance towards the kind of language-based AI systems we want as a society (Wooldridge and Jennings, 1995). Currently, the most widespread way to evaluate LLMs is by means of **large benchmarks** made up of miscellaneous question-answering (QA) tasks. Pre-LLM benchmarks such as GLUE and SuperGLUE (Wang et al., 2019b,a) have been replaced by even larger evaluation suites such as MMLU (Measuring Massive Multitask Language Understanding; Hendrycks et al., 2021), GSM8K (Graduate School Math; Cobbe et al., 2021), or BBH (BIG-Bench Hard; Suzgun et al.,

2023; Srivastava et al., 2023). Models with high performance on these benchmarks are taken to possess extensive **world knowledge along with complex problem-solving abilities**.

This trend has promoted standardisation in LLM evaluation protocols, with online leaderboards constantly updated as new models are released. Despite this undeniable benefit, large QA benchmarks like those mentioned above are not without problems. Evaluation results may be inflated by data contamination (see, e.g., Gema et al. 2025 for MMLU and Mirzadeh et al. 2025 for GSM8K) and distorted by model sensitivity to prompt format (Zhuo et al., 2024). Moreover, by design, such benchmarks overlook actual language use in favour of knowledge-intensive tasks where success is measured against gold-standard reference answers provided in a single conversational turn. This contrasts with the view, put forward by philosophers and psycholinguists alike (Wittgenstein, 1953; Austin, 1962; Searle, 1969; Clark, 1996), that the quintessence of language resides in *situated language use*, i.e., using language for a purpose in social and task-based multi-turn interactions (Bisk et al., 2020).

The situated and interactive view underpins a parallel evaluation trend where LLMs are evaluated as **goal-directed language users** by means of **interactive games** (Schlangen, 2023; Suglia et al., 2024).¹ This interactive evaluation paradigm goes beyond single-turn text generation, which is critical for deploying LLMs as agents. Additionally, it is less susceptible to data contamination because the vast space of possible multi-turn interactions is unlikely to be fully represented in the training data. As a result, interactive games provide a more robust framework for evaluating the true generalisation capacity of LLMs (Hupkes et al., 2023).

¹Online leaderboards have started to appear for the interactive games evaluation paradigm; see, e.g., <https://textarena.ai/>, <https://clembench.github.io>.

*Corresponding author.
Email: filippo.momente@studenti.unitn.it

Yet, despite these advantages, it is not easy to pinpoint which specific abilities underpin models’ performance on interactive language games—a difficulty that to some extent also applies to static question-answering benchmarks such as MMLU.

In this paper, we examine these two evaluation paradigms—large QA benchmarks and interactive games—and argue that they can provide complementary perspectives. First, we investigate whether QA benchmarks or games are more effective in gauging qualitative differences between models, e.g., across model families and sizes. We evaluate a selection of current LLMs from four model families and find that games highlight differences between LLMs more strongly than QA benchmarks: While scaling model size leads to systematic improvements on benchmarks, it doesn’t guarantee performance boosts in interactive language use. To shed light on the abilities underlying models’ performance on these two evaluation frameworks, we resort to **targeted cognitive tests**. We propose a taxonomy of cognitive skills motivated by neurocognitive science and compile a list of existing evaluation datasets designed to assess each skill in isolation. We then investigate to what extent increased performance on specific cognitive abilities correlates with performance gain in large QA benchmarks vs. interactive games. Our analysis shows that while causal and logical reasoning correlate with both static and interactive tests, differences emerge regarding core executive functions and social/emotional skills; in particular, working memory and emotional intelligence are only significantly correlated with performance in games.

2 Models

We apply our evaluation framework to a selection of open-weight LLMs ranging from 7B to 72B models. Considering that instruction following capabilities are essential for our analysis, we selected models that have an average performance on the IFEval benchmark (Zhou et al., 2023) higher than 70% (see Figure 1 for details). We evaluate the following models: Olmo-2-1124 with 7 and 13 billion parameters (OLMo-2-1124-* -Instruct) (Walsh et al., 2024); Qwen2.5 with 7B, 32B, and 72B parameters (Qwen2.5-* -Instruct) (Yang et al., 2024; Team, 2024); LLama-3 with 8B (Llama3.1-8B-Instruct) and 70B parameters (Llama3.3-70B-Instruct) (Grattafiori et al., 2024), and Falcon3-10B-Instruct (Falcon Team,

2024). See Appendix A for further model details.

3 Static vs. Interactive Assessments

3.1 Benchmarks

Large QA benchmarks We take MMLU (Hendrycks et al., 2021) and BBH (Suzgun et al., 2023) as representative of large QA benchmarks. MMLU evaluates whether LLMs can apply knowledge from specific domains: it consists of multiple-choice questions spanning 57 academic subjects. BBH assembles diverse tasks drawing problems from linguistics, child development, maths, and common-sense reasoning, among others.

Interactive games We take clembench (Chalamasetti et al., 2023) as a characteristic benchmark to assess LLMs’ gameplay ability in dialogue games. We consider the games 1) *Taboo*, 2) standard *Wordle* and the two variants *Wordle (Clue)* and *Wordle (Critic)*, 3) *Reference Game*, 4) *Image Game*, and 5) *Private/Shared*. See Appendix B.

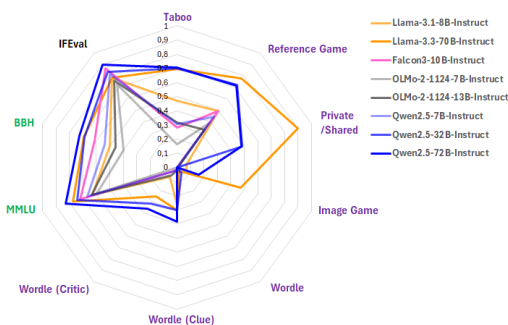


Figure 1: Accuracy across the different model sizes: IFEval, **Static**, and **Interactive** assessments.

3.2 How to Identify Blind Spots in LLMs

LLM evaluation instruments have most practical use when they allow us to track progress by identifying blind spots in models. Here we compare the two evaluation paradigms under study on the extent to which they highlight differences between current models, helping us form hypotheses about possible problem sources and successful mitigation strategies. Figure 1 shows models’ performance on IFEval, the large QA benchmarks, and interactive games. As mentioned in Section 2, all models are reasonably able to follow instructions as measured by IFEval. While the OLMo-2 models are more inconsistent across different model sizes, all the other models exhibit the expected pattern of showcasing better performance on both large QA

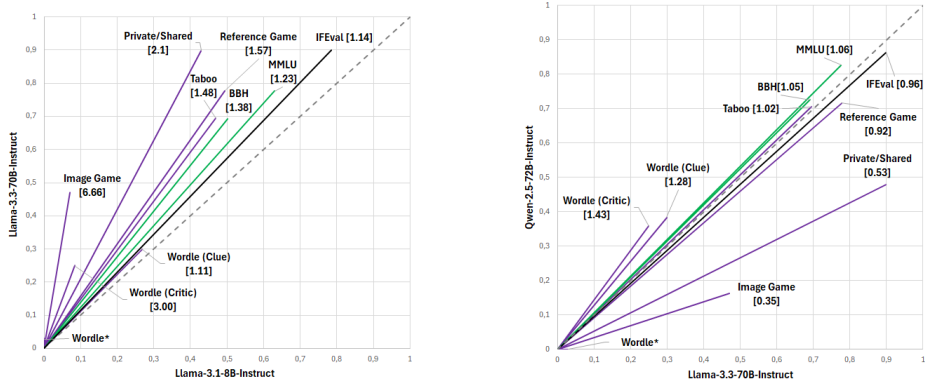


Figure 2: Comparing datasets in their power to discriminate between models of different size but same family (left) and of different families but similarly large (right). The number next to the benchmark’s name indicates the ratio of performance between the two models. The asterisk ‘*’ next to *Wordle* indicates that the ratio is undefined.

benchmarks and interactive games when parameter count increases. At the same time, we observe that most of the interactive games highlight the benefits of large model sizes much more strongly. This can more easily be appreciated in Figure 2 (left) for Llama-3.1-8B vs. Llama-3.3-70B. In this visualisation, the further away a benchmark is from the diagonal, the more affected performance is by model size. While playing *Wordle* is extremely challenging for any model, scaling up the number of parameters appears to be fundamental to succeed at *Private/Shared*, *Image Game*, and *Reference Game*—much more so than for MMLU and BBH.

Is size however all we need? Figure 2 (right) shows that QA benchmarks do not substantially distinguish between large models of comparable size (Llama-3.3-70B-Instruct vs. Qwen2.5-72B-Instruct): scaling on the number of parameters results in performance boosts across model families. Hence, arguably large QA benchmark test for abilities than can be expressed within parametric knowledge. Given that such benchmarks currently are the standard LLM evaluation paradigm, it is not surprising that scaling is high on the agenda of model developers. In contrast, interactive games seem to provide a different picture: models with comparable parametric capacity perform very differently on *Image Game*, *Private/Shared*, and *Wordle (Clue/Critic)*. A similar trend can be observed among the other models we evaluated (see details in Appendix I). This result supports the hypothesis that size is not all there is behind the potential of LLMs to learn inferential strategies for effective language use in interaction.

4 Cognitive Abilities Assessment

We now turn to cognitive tests—a complementary evaluation method that focuses on specific cognitive abilities deemed essential for effective language use in real-world situations. We explore the use of targeted cognitive tests to complement evaluation based on large QA benchmarks and interactive games.

4.1 Taxonomy and Datasets

We present a taxonomy of cognitive abilities involved in human *functional linguistic competence* (Mahowald et al., 2024). It is guided by neurocognitive research (Ward, 2019), and it separates capabilities into two distinct macro-categories known to recruit different brain networks: executive functions and socio-emotional skills. **Executive functions** are broadly defined as the complex processes by which we control and optimise our thoughts and behaviour (Baddeley, 1986) and are divided into *core* and *higher-order* abilities. **Socio-emotional skills** represent the abilities necessary to interact adaptively with other individuals (Higgins, 1987), including the ability to recognize their emotional and cognitive states.

For each cognitive ability, we select an existing evaluation dataset designed to test it in isolation drawing inspiration from human cognitive assessments. We discard datasets that require manual evaluation from the analysis. Table 1 and Table 2 list the abilities in the taxonomy and the datasets we use to evaluate them.² Socio-emotional skills

²We found no dataset to evaluate inhibitory control. The datasets we found for Emotion-regulation, Self-awareness (Liu et al., 2024), Empathy (Chen et al., 2024) and Social Problem-solving (Du et al., 2024) require human evaluation.

	Cognitive Ability	Benchmark
Core	Working Memory	Gong et al. (2024)
	Cognitive Flexibility	Kennedy and Nowak (2024)
	Inhibitory Control	–
HO	Logical Reasoning	Liu et al. (2023)
	Causal Reasoning	Jin et al. (2023)
	Commonsense Reasoning	Sakaguchi et al. (2021)
	Planning	Zheng et al. (2024)

Table 1: Core and Higher-Order Executive Functions.

Cognitive Ability	Benchmark
Pragmatics	Hu et al. (2023)
Theory of Mind	Gu et al. (2025)
Attribution and Judgement	Gu et al. (2025)
Social Commonsense Reasoning	Sap et al. (2019)
Emotional Intelligence	Paech (2023)
Emotion Regulation	–
Self-Awareness	–
Empathy	–
Social Problem-Solving	–

Table 2: Social and Emotional Skills.

have only recently entered the evaluation landscape in NLP, and they have done so with a forceful presence: remarkably, small benchmarks already exist for almost all of the abilities in this category.

4.2 Cognitive Ability Analysis

Equipped with our taxonomy and associated cognitive tests, we aim to shed some light on the cognitive abilities involved in interactive games and large QA benchmarks. Figure 3 reports Kendall’s τ correlation coefficients, with asterisks indicating statistical significance ($p < 0.05$); see Appendix I for a detailed correlation matrix between single datasets. The analysis reveals that performance both on static and interactive evaluation correlates with performance on tests measuring higher-order reasoning abilities; while planning is more dominant in static problem-solving tasks, working memory seems to be beneficial for games. Among the social skills, pragmatics appears to be relevant for both static and interactive tests, while emotional intelligence and ToM correlate better with the latter. While these results suggest that interactive tests correlate more strongly with socio-emotional skills than static tests, this analysis remains speculative, as we still lack carefully curated cognitive abilities tests specifically designed for LLMs.

5 Related Work

Waldis et al. (2024) proposes Holmes as a framework to assess the English linguistic competence

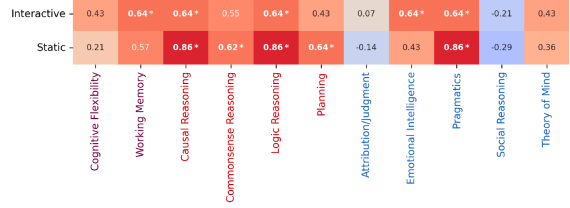


Figure 3: Correlation of cognitive abilities with Static and Interactive assessments (* indicates $p < 0.05$).

of language models. They evaluate models’ competence (morphology, syntax, and semantics) by comparing them across architectures and sizes by probing their internal representations. Moreover, by measuring the correlation between Holmes and downstream tasks results, they observe that morphology highly correlates with reasoning. Rather than on formal linguistic competence, we focus on functional linguistic competences and compare them not just with large QA benchmarks but also with interactive games. Ma et al. (2023) carry out a holistic evaluation of LLMs’ Theory of Mind by inspecting the literature through the competences a model with a ToM should have based on a known taxonomy. Similarly, we take a top-down approach but consider the whole spectrum of cognitive abilities and highlight the importance of connecting them with the complementary benchmarks largely used by the community to monitor LLMs’ progress.

6 Conclusion

Our results show the different discriminating power of interactive games over one-turn static large QA benchmarks. Crucially, we argue that in order to claim that LLMs have emerging abilities, measuring performance on large QA benchmarks or interactive games is not sufficient per se, but should rather be triangulated with controlled tests designed to evaluate such abilities. Furthermore, we highlight the potential value of carefully designed controlled benchmarks inspired by human cognitive ability assessment as a good means for such correlation analyses. While each cognitive assessment test alone does not get us very far in the quest for robust LLM evaluation, we contend that this type of evaluation paradigm has the potential to enhance our understanding of what fundamental abilities LLMs must develop to be able to function effectively as language agents, where multiple skills may be required and possibly interact. Nevertheless, we agree with Millière and Rathkopf (2024) that caution should be exerted before draw-

ing conclusions about LLMs’ abilities from these tests meant for humans. New carefully designed behavioural experiments for LLMs should be proposed, and supplemented with mechanistic studies.

Limitations

Our evaluation prompts models to provide direct answers without employing chain-of-thought (CoT) reasoning or similar capability elicitation techniques. While different elicitation strategies may enhance question-answering, interactive, and cognitive abilities in different ways (Yao et al., 2023; Hao et al., 2023; Li et al., 2024), we opted for an approach that remains agnostic to specific evaluation methods and datasets. This ensures a consistent basis for comparison across models, though future work could explore how alternative prompting strategies influence performance across the three evaluation paradigms. Moreover, for the cognitive abilities assessments, we used currently available datasets; such resources have started to be compiled only very recently, hence the tests we used may not guarantee to evaluate the intended abilities in LLMs. Nevertheless, they help in establishing our message and call for more analysis in such direction.

References

- John Langshaw Austin. 1962. *How to do things with words*. Clarendon Press, London, UK.
- Alan Baddeley. 1986. *Working memory*. Oxford University Press.
- Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020. [Experience grounds language](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735, Online. Association for Computational Linguistics.
- Kranti Chalamalasetti, Jana Götze, Sherzod Hakimov, Brielen Madureira, Philipp Sadler, and David Schlangen. 2023. [clmbench: Using game play to evaluate chat-optimized language models as conversational agents](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11174–11219, Singapore. Association for Computational Linguistics.
- Yuyan Chen, Songzhou Yan, Sijia Liu, Yueze Li, and Yanghua Xiao. 2024. [EmotionQueen: A benchmark for evaluating empathy of large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 2149–2176, Bangkok, Thailand. Association for Computational Linguistics.
- Herbert H Clark. 1996. *Using language*. Cambridge University Press, Cambridge, UK.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Y. Du, P. Rajivan, and C. Gonzalez. 2024. [Large language models for collective problem-solving: Insights into group consensus decision-making](#). In *Proceedings of the Annual Meeting of the Cognitive Science Society*.
- Technology Innovation Institute Falcon Team. 2024. The falcon 3 family of open models.
- Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, Claire Barale, Robert McHardy, Joshua Harris, Jean Kaddour, Emile van Krieken, and Pasquale Minervini. 2025. Are we done with MMLU? In *NAACL 2025*.
- Dongyu Gong, Xingchen Wan, and Dingmin Wang. 2024. [Working memory capacity of chatgpt: an empirical study](#). In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence, AAAI’24/IAAI’24/EAAI’24*. AAAI Press.
- D. A. Grant and E. A. Berg. 1948. Wisconsin card sorting test. *Journal of Experimental Psychology*.
- Aaron Grattafiori, Abhimanyu Dubey, and et al. Abhinav Jauhri. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Yuling Gu, Oyvind Tafjord, Hyunwoo Kim, Jared Moore, Ronan Le Bras, Peter Clark, and Yejin Choi. 2025. [Simpletom: Exposing the gap between explicit tom inference and implicit tom application in llms](#).
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. 2023. [Reasoning with language model is planning with world model](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8154–8173, Singapore. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Tory Higgins. 1987. [Social cognition and social perception](#). *Annual review of psychology*, 38(1):369–425.

- Jennifer Hu, Sammy Floyd, Olessia Jouravlev, Evelina Fedorenko, and Edward Gibson. 2023. [A fine-grained comparison of pragmatic language understanding in humans and language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4194–4213, Toronto, Canada. Association for Computational Linguistics.
- Dieuwke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, et al. 2023. [A taxonomy and review of generalization research in NLP](#). *Nature Machine Intelligence*, 5(10):1161–1174.
- Zhijing Jin, Yuen Chen, Felix Leeb, Luigi Gresele, Ojasv Kamal, Zhiheng Lyu, Kevin Blin, Fernando Gonzalez, Max Kleiman-Weiner, Mrinmaya Sachan, and Bernhard Schölkopf. 2023. [CLadder: Assessing causal reasoning in language models](#). In *NeurIPS*.
- Sean M Kennedy and Robert D Nowak. 2024. [Cognitive flexibility of large language models](#). In *ICML 2024 Workshop on LLMs and Cognition*.
- Zaijing Li, Gongwei Chen, Rui Shao, Yuquan Xie, Dongmei Jiang, and Liqiang Nie. 2024. [Enhancing emotional generation capability of large language models via emotional chain-of-thought](#). *Preprint*, arXiv:2401.06836.
- Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, Ming Zhou, and Yue Zhang. 2023. [Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2947–2962.
- Ziyi Liu, Abhishek Anand, Pei Zhou, Jen-tse Huang, and Jieyu Zhao. 2024. [InterIntent: Investigating social intelligence of LLMs via intention understanding in an interactive game context](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6718–6746, Miami, Florida, USA. Association for Computational Linguistics.
- Ziqiao Ma, Jacob Sansom, Run Peng, and Joyce Chai. 2023. [Towards a holistic landscape of situated theory of mind in large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1011–1031, Singapore. Association for Computational Linguistics.
- Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2024. [Dissociating language and thought in large language models: a cognitive perspective](#). *Trends in Cognitive Sciences*, 28.
- Raphaël Millière and Charles Rathkopf. 2024. [Anthropocentric bias and the possibility of artificial cognition](#). In *ICML 2024 Workshop on LLMs and Cognition*.
- Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2025. [GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models](#). In *The Thirteenth International Conference on Learning Representations*.
- Leora Morgenstern and Charles L. Ortiz. 2015. The winograd schema challenge: evaluating progress in commonsense reasoning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, AAAI’15, page 4024–4025. AAAI Press.
- Samuel J. Paech. 2023. [Eq-bench: An emotional intelligence benchmark for large language models](#). *Preprint*, arXiv:2312.06281.
- R. D. Rogers and S. Monsell. 1993. Costs of a predictable switch between simple cognitive tasks. *Journal of Experimental Psychology*, 124:207–231.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. [Winogrande: an adversarial winograd schema challenge at scale](#). *Commun. ACM*, 64(9):99–106.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. [Social IQa: Commonsense reasoning about social interactions](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.
- David Schlangen. 2023. [What a situated language-using agent must be able to do: A top-down analysis](#). *Preprint*, arXiv:2302.08590.
- John R Searle. 1969. *Speech acts: An essay in the philosophy of language*. Cambridge University Press, Cambridge, UK.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew Kyle Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubaranjan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden,

Benno Stein, Berk 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 Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt, Christopher D Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, C. Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan 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 Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-Lopez, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, 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 Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocon, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason 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 Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg 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, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel, Kevin Omondi, Kory Wallace Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonnell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros-Colón, Luke Metz, Lütfi Kerem Senel, Maarten Bosma, Maarten Sap, Maartje Ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez-Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás Schu-

bert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael Andrew Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdah Gheini, Mukund Varma T, Nanyun Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, 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, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing 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 Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Roman Le Bras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel Stern Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane 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, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, 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, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadhollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Sophie Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi

- Wu. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *Transactions on Machine Learning Research*. Featured Certification.
- Alessandro Suglia, Ioannis Konstas, and Oliver Lemon. 2024. Visually grounded language learning: a review of language games, datasets, tasks, and models. *Journal of Artificial Intelligence Research*, 79:173–239.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. 2023. [Challenging BIG-bench tasks and whether chain-of-thought can solve them](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13003–13051, Toronto, Canada. Association for Computational Linguistics.
- Qwen Team. 2024. [Qwen2.5: A party of foundation models](#).
- Raphael Vallat. 2018. [Pingouin: statistics in python](#). *Journal of Open Source Software*, 3(31):1026.
- A. Waldis, Y. Perlitz, L. Choshen, Y. Hou, and I. Gurevych. 2024. [Holmes: A benchmark to assess the linguistic competence of language models](#).
- Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. [2 olmo 2 furious](#).
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. [Superglue: A stickier benchmark for general-purpose language understanding systems](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *International Conference on Learning Representations*.
- Jamie Ward. 2019. *The student’s guide to cognitive neuroscience*. Routledge.
- Ludwig Wittgenstein. 1953. *Philosophical Investigations*. Wiley-Blackwell, New York, NY, USA.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Michael Wooldridge and Nicholas R Jennings. 1995. Intelligent agents: Theory and practice. *The knowledge engineering review*, 10(2):115–152.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Peng, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 11809–11822. Curran Associates, Inc.
- Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang, Xinyun Chen, Minmin Chen, Azade Nova, Le Hou, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and Denny Zhou. 2024. [Natural plan: Benchmarking llms on natural language planning](#).
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sidhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.
- Jingming Zhuo, Songyang Zhang, Xinyu Fang, Haodong Duan, Dahua Lin, and Kai Chen. 2024. Prosa: Assessing and understanding the prompt sensitivity of llms. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 1950–1976.

Appendix

A Models

The Olmo-2-1124 series (Walsh et al., 2024) includes models with 7 billion and 13 billion param-

eters (OLMo-2-1124-*-Instruct). Both models are designed for a variety of tasks, including chat, mathematics, and reasoning. They have undergone supervised fine-tuning on the Tulu 3 dataset and further training using DPO techniques.

The Qwen2.5 series (Yang et al., 2024; Team, 2024) includes models with 7B, 32B, and 72B parameters (Qwen2.5-*-*Instruct). They are multi-lingual, supporting over 29 languages, and excel in coding, mathematics, and instruction following.

The Llama-3 series (Grattafiori et al., 2024) includes models with 8B (Llama3.1-8B-Instruct) and 70B parameters (Llama3.3-70B-Instruct). These models are optimized for multilingual dialogue and support various languages. They use an optimized transformer architecture and are fine-tuned for instruction following.

The Falcon3 series (Falcon Team, 2024) includes a model with 10 billion parameters. It achieves state-of-the-art results in reasoning, language understanding, instruction following, code, and mathematics tasks. It supports four languages (English, French, Spanish, Portuguese) and a context length of up to 32K.

B Interactive Games

We leverage *clmbench* (Chalamalasetti et al., 2023), a benchmark that assesses models’ game-play ability in well-defined dialogue games such as:

1) *Taboo*: a game where one player tries to get the other player to guess a word without using certain ‘taboo’ words in their clues;

2) *Wordle*: a game where one player thinks of a word and the other player tries to guess it by suggesting words that are similar or related;

3) *Wordle (Clue)*: a variant of Wordle where the guesser is given a clue to help them guess the target word;

4) *Wordle (Critic)*: A variant of Wordle where the guesser’s suggestions are evaluated by a third player, who provides feedback on their accuracy;

5) *Reference Game*: a game where one player is presented with three grids and is tasked to make the other—who is also presented with the same three grids, but potentially in a different order—identify a pre-specified one;

6) *Private/Shared*: a game where a customer agent goes through a form with a service agent and, after each turn, a third agent,³ probes the customer

on what they believe that the service agent already knows; hence, which information is “private” and which is “shared”.

C Taxonomy of Cognitive Abilities

- Executive Functions:

- Core abilities

Working Memory: Hold and manipulate information in mind over short periods;

Inhibitory Control Suppress automatic, inappropriate, or impulsive responses and resist distractions;

Cognitive Flexibility: Adapt to new situations, switch between tasks, and think about multiple concepts simultaneously;

- Higher-order abilities

Planning: Set goals, develop steps to achieve them, and anticipate potential obstacles;

Causal Reasoning: Understand cause-and-effect relationships;

Logical Reasoning : Deductive and inductive reasoning;

Commonsense Reasoning: Apply general common knowledge to everyday scenarios, including understanding basic physical properties, such as gravity, solidity, and object interaction;

- Socio-emotional skills:

Social Commonsense Reasoning:

Understand social norms and expectations;

Social Problem-Solving: Analyze social situations, generate solutions, and make decisions that foster positive interactions;

Emotional Intelligence: Recognize, interpret, and manage one’s own and others’ emotions.

Emotion Regulation: Manage and modify one’s emotional responses in appropriate ways;

Self-Awareness: Recognize and understand one’s own emotions, thoughts, and behaviors;

Master”. This agent plays a particularly active role in *Private/Shared*.

³In *clmbench*, all interactions are mediated by a “Game

Empathy: Share and understand the feelings of others, both emotionally and cognitively;

Theory of Mind: Understand that others have thoughts, beliefs, desires, and intentions different from one’s own;

Attribution and Judgment: Interpret the causes of others’ behavior, distinguishing between intentional and unintentional actions.

Pragmatics: Aspects of communication that go beyond formal language competence: considering communicative intentions, the communicative context of the utterance, shared knowledge between speakers, manners, social and cultural norms.

D Benchmarks for Cognitive Abilities

Working Memory (Gong et al., 2024) (referred as *WM* in this work) is a set of verbal and spatial n -back tasks presented with three levels of difficulties from $n = 1$ to $n = 3$. The model has to identify whether the current stimulus (a letter in a string or a spatial location in a grid) is the same as the n back stimulus or not.

Cognitive Flexibility (Kennedy and Nowak, 2024) (referred as *LLM-Cognitive-Flexibility* in this work) aims to test to what degree LLMs can rapidly switch tasks within a single context window. To this end, it employs two neuropsychological tests, the Wisconsin Card Sorting Test (WCST) (Grant and Berg, 1948) and the Letter-Number Test (LNT) (Rogers and Monsell, 1993) commonly used to measure cognitive flexibility in humans.

Logical Reasoning *LogiQA 2.0* (Liu et al., 2023) This dataset evaluates logical reasoning using the same data both in NLI and Machine Reading Comprehension format (text, question, multiple-choice) for each of the following (deductive) reasoning types: categorical, sufficient condition, necessary condition, disjunctive, conjunctive reasoning.

Causal Reasoning *CLADDER* (Jin et al., 2023) focuses on formal causal reasoning (causal inference), as opposed to commonsense causal reasoning. The dataset is constructed from formal logic-based templates that are then verbalised into natural language as binary questions.

Commonsense Reasoning *WinoGrande* (Sakaguchi et al., 2021) A large-scale dataset of 44k commonsense reasoning problems consisting of pairs of nearly identical questions with two answer choices (as in the original Winograd Schema Challenge (Morgenstern and Ortiz, 2015) but without its bias).

Planning *NATURAL PLAN* (Zheng et al., 2024) is a realistic planning benchmark consisting of three tasks expressed in natural language: Trip Planning, Meeting Planning and Calendar Scheduling. Models are given a situation and a problem to solve (e.g. find a trip plan that satisfies some constraints given the situation described.) Each task contains problems of different levels of complexity based on the number of cities, people or days involved. The problems are often based on numerical reasoning too. We evaluate models on the Trip Planning and the Calendar Scheduling tasks.

Emotional Intelligence *EQ-Bench* (Paech, 2023) the model is given an emotionally charged short dialogue (generated by GPT-4) and has to score the four possible emotions felt by a given character. Scores are compared against a reference score.

Pragmatics (Hu et al., 2023) (referred as *LM-Pragmatics* in this work) is a benchmark evaluating LLMs’ understanding of seven pragmatics phenomena: deceit, indirect speech, irony, maxims, metaphor, humour, and coherence. Scenarios are grounded into social situations, requiring LLMs to interpret utterances. The task is designed as a multi-choice questionnaire with 2-5 questions based on the subtask.

Social Commonsense *SOCIAL IQA* (Sap et al., 2019) a dataset for evaluating social commonsense reasoning and emotional intelligence. Each sample includes a short scenario and three multiple-choice questions across six categories: intentions, reactions, descriptions, motivations, needs, and consequences. Transfer learning on this dataset has shown strong performance on other commonsense reasoning benchmarks.

Attribution and Judgment/Theory of Mind

SimpleToM (Gu et al., 2025) contains concise, diverse stories each with questions that ask models to predict behavior ("Will Mary

pay for the chips or report the mold?"), judgment ("Mary paid for the chips. Was that reasonable?") or mental states ("Is Mary likely to be aware that 'The can of Pringles has moldy chips in it.'? Yes or No?") The first two subtasks have been taken as a reference for the Attribution and Judgment cognitive ability, while the last as a reference for Theory of Mind.

E Benchmark Implementations

For the majority of the static benchmarks evaluated in this work we relied on the popular framework for the evaluation of LLMs *lm-eval*⁴ (ver. 0.4.7), which already made available many of the selected benchmarks, and enabled a common interface for the implementation of most of the remaining ones.

The benchmarks which were already present within the framework are: SOCIAL IQA (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021), EQ-Bench (Paech, 2023), LogiQA 2.0 (Liu et al., 2023), MMLU (Hendrycks et al., 2021). The benchmarks which have been implemented in the framework over the course of the study are: CLADDER (Jin et al., 2023), LM-Pragmatics (Hu et al., 2023), SimpleToM (Gu et al., 2025), NATURAL PLAN (Zheng et al., 2024).

As for BBH and IFEval, they were also available in the framework and potentially usable, however we decided to rely on the scores made available by Huggingface in their Open Leaderboard 2⁵ as they also rely on the *lm-eval* framework for evaluation. This has been possible for all of the models except one (*OLMo-2-1124-13B-Instruct*), for which the scores were not available in the leaderboard at the time of this study. In this case, we have used the code made available by Huggingface and followed their instruction for reproducing the results (in particular, given the instruction-tuned nature of the model, the settings *apply_chat_template* and *fewshot_as_multiturn* were applied).

As for the interactive games, we have used the implementation provided by version 1.5 of the *clembench* (Chalamalasetti et al., 2023). The remaining benchmarks (WM, LLM-Cognitive-Flexibility) have been implemented outside of the framework, as *lm-eval* did not provide support for the multi-turn nature of the tasks.

⁴<https://github.com/EleutherAI/lm-evaluation-harness>

⁵<https://huggingface.co/spaces/open-llm-leaderboard/blog>

E.1 Zero-shot and Few-shot Tasks

The majority of the tasks have been evaluated in a zero-shot setting with the exception of MMLU (5-shot), BBH (3-shot) (following common practices in model evaluation, e.g. in the Open Leaderboard 2 for BBH) and NATURAL PLAN (5-shot). In the case of NATURAL PLAN, our models performed really poorly when evaluated in a zero-shot fashion—resulting in scores close to 0. Given that the task relies on the models producing answers in a strict format for parsing, we opted for using the 5-shot version provided by the benchmark’s authors.

E.2 Metrics

E.2.1 Evaluation

In the evaluation of models, we followed the original works’ implementations as well as associated metrics. However, it may be the case that for a certain benchmark more metrics were defined, or that the original work did not aggregate results across subtasks. For this reason, we report here the metrics we used for evaluating models.

In the case of *Clembench games*, we computed performance by computing the ratio between the quality score (a number from 0 to 100) and the percentage of played games (a number between 0 and 1) divided by 100.

In the case of *IFEval*, following what was done in the Open Leaderboard 2, we averaged the results obtained on prompt-level and instruction-level strict accuracy.

As for *EQ-Bench*, we computed the task-specific score as it was implemented in the *lm-eval*.

Regarding *WM*, we only considered the subtask *Verbal N-3*, and we computed the accuracy for the results obtained across the 50 trials defined in the original work.

In the case of *LLM-Cognitive-Flexibility*, we ran each subtask 8 times with 25 trials each, and computed the average of the accuracy obtained in each run. In this case, the accuracy was computed only on the trials for which response parsing was successful. We then averaged the accuracy obtained on both subtasks to compute the final score.

In the case of CLADDER, we followed the original work which treated the task as generative and probed for the presence of the substrings "yes"/"no" at the beginning of the model’s answer.

In NATURAL PLAN, the original work defined a rule-based procedure to parse specific data from the generated plan (e.g., dates). We reuse their

parsing procedure and verify whether the expected elements are all present in the parsed plan.

For the remaining tasks (LogiQA 2.0, WinoGrande, LM-Pragmatics, SOCIAL IQA, MMLU, BBH, SimpleToM), we treated them as a multiple-choice question answering task that is evaluated based on the likelihood of the correct answer for the task.

In the case of BBH, the Open Leaderboard 2’s evaluation code excludes three of the original tasks from the overall score’s computing: *dyck languages*, *navigate* and *word sorting*. The performances on these subtasks are therefore also ignored in the performance reported in this study.

In the case multiple subtasks were present (LM-Pragmatics, MMLU, BBH, NATURAL PLAN, LLM-Cognitive-Flexibility), we computed the micro-average over the results achieved on each subtask. In the specific case of SimpleToM, since the subtasks were associated with two different Cognitive Abilities, we’ve aggregated the score of the subtasks *behaviour* and *judgment* into a single score (under Attribution and Judgment), and considered the *mental state* subtask separately (under Theory of Mind).

E.2.2 Correlation

For measuring the pair-wise correlation between benchmarks, we’ve computed the Kendall rank correlation coefficient (or Kendall’s Tau) (Tau-b version). It measures rank correlation according to this formula:

$$\tau_b = \frac{P - Q}{\sqrt{(P + Q + T_x)(P + Q + T_y)}}$$

where:

P = number of concordant pairs,

Q = number of discordant pairs,

T_x = tie correction for variable X ,

T_y = tie correction for variable Y .

This method was preferable over the others given its robustness in case of few data points, as it was in our case. We have also experimented with the Pearson correlation coefficient and observed that in the majority of the cases, the correlation patterns were similar, however with larger positive as well as negative correlations compared to Kendall. We’ve relied on the implementation provided by the *pingouin* Python package (Vallat, 2018).

E.3 Generation Settings

The tasks which required the models to generate text are: EQ-Bench, WM, MMLU, IFEval, the clembench games, LLM-Cognitive-Flexibility, NATURAL PLAN, CLADDER. With the exception of Working Memory and LLM-Cognitive-Flexibility, all tasks have been evaluated by applying a temperature of 0. Following the original implementation, we have applied a temperature of 1 to WM and 0.7 to LLM-Cognitive-Flexibility. In these cases, however, the increased randomness caused by the higher temperature was mitigated by averaging the results obtained over multiple trials.

As for the other generation settings, we also have followed what was prescribed in the original works regarding the tokens for the termination of the generation, the maximum or minimum number of tokens. In the case of NATURAL PLAN, the original work did not provide specific information regarding the settings they have adopted for the evaluation. Given the highly challenging nature of the task, we have set the minimum and maximum number of tokens to 90 and 350 respectively. This was derived based on the minimum and maximum number of tokens in the gold plans.

F Limitations in the Evaluations

In certain cases, results have not been computed on all the subtasks available for that benchmark. In addition to BBH, we also made special arrangements for Working Memory and NATURAL PLAN. In the case of NATURAL PLAN, we have not considered results coming from the *meeting* subtask, while for WM we have only considered those coming from the *Verbal (Base) N-3* subtask. In the first case, the high amount of resources required for evaluating the task, especially for the larger models (60% of the prompts contained above 14k space-separated words (up to 38.3k) vs 100% below 15.1k for the *travel* subtask and below 7.09k for the *calendar* subtask prevented us from doing so. As for the second, we’ve only considered the base version of the verbal subtask and excluded its variations as they would not provide meaningful information for this study.

G Model Implementations

All the models used in this study have been made available by Huggingface, and have been accessed through their *transformers* (Wolf et al., 2020) library. For text generation, we have been applying

the default chat template specified within the same library.

H Computational Resources

As a reference, the time required for running through all the benchmarks for the Llama-3.1-8B-Instruct model on 1 A100 GPU with batch size set to 'auto' in the *lm-eval* (i.e. it automatically fits into the memory the maximum batch size possible for each task). For the Clembench games, LLM-Cognitive-Flexibility and WM, the batch size is 1. The time also includes time required for procedures performed by the *lm-eval* prior to the actual evaluation (only for those datasets evaluated through the framework) and loading the model into the memory (all tasks).

LLM-Cognitive-Flexibility: ~1:50 min

LogiQA 2.0: ~5 min

CLADDER: ~19:30 min

WinoGrande: ~1 min

NATURAL PLAN: ~4:50 hours

WM ~2:40 min

EQ-Bench: ~3 min

LM-Pragmatics: ~6:30 min

SOCIAL IQA: ~1:30 min

SimpleToM: ~2:40 min

MMLU: ~14 min

Taboo: ~3:30 min

Reference Game: ~3:00 min

Image Game: ~2.40 min

Wordle: ~7:50 min

Wordle (Critic): ~2:50 min

Wordle (Clue): ~2:15 min

Private/Shared: ~17:30 min

I Supplementary Plots

Figure 4 and Figure 5 show the supplementary plots for the results in Section 3 (comparing models of different size but same family and models of similar size respectively). Moreover, we report the supplementary plots for the results in Section 4. Figure 6 presents a direct comparison of models based on our selected cognitive tests. Figure 7 reports an extended version of Figure 3. Finally, Figure 8 reports by means of example two scatter plots representing respectively situations of high and low correlation between two benchmarks (a game and a cognitive ability).

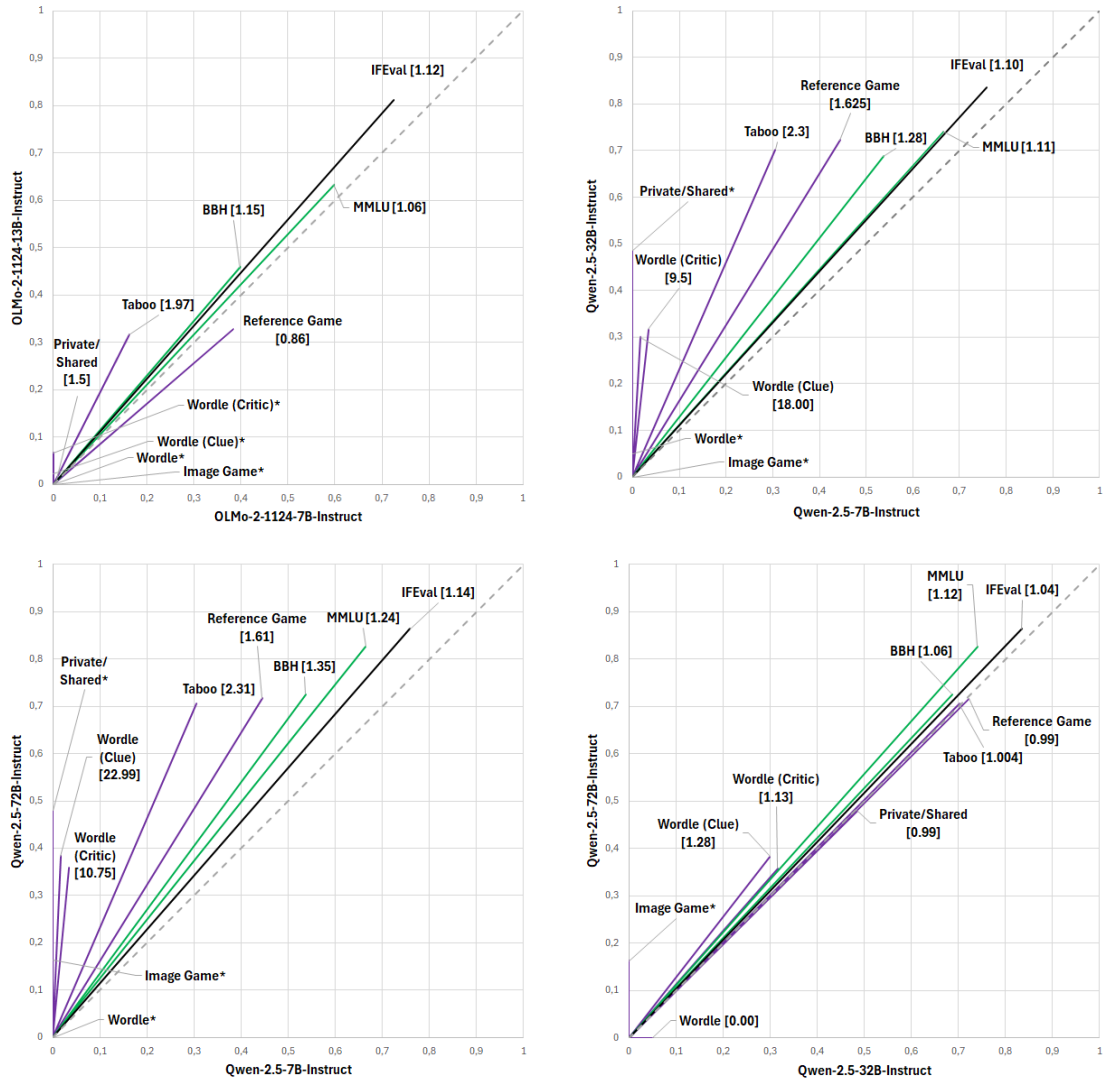


Figure 4: Comparing datasets in their power to discriminate models across size.

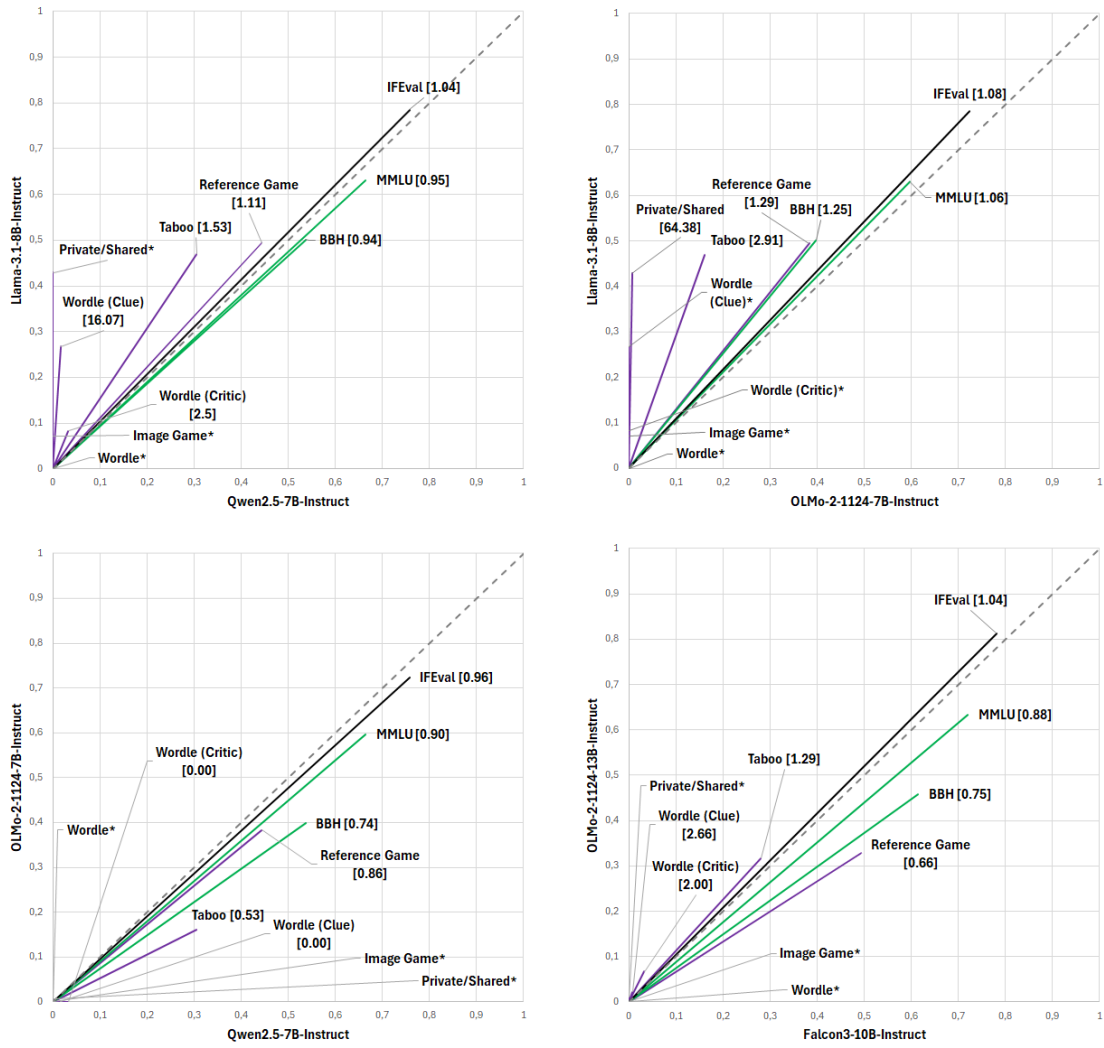


Figure 5: Comparing datasets in their power to discriminate models with similar size across families.

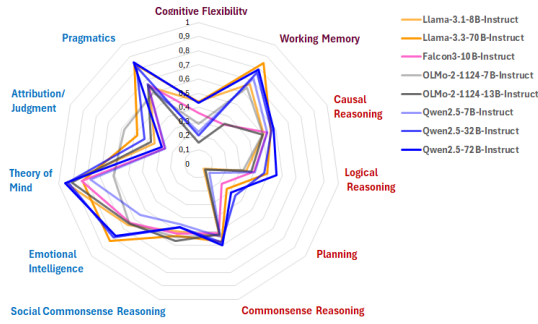


Figure 6: Cognitive Abilities Spectrum of LLMs

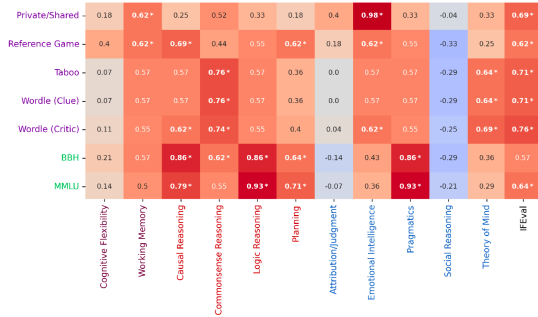


Figure 7: Correlation of cognitive abilities with Large QA benchmarks and Interactive Games. The correlation matrix does not include results on Wordle and Image Game as model performances' were too low.

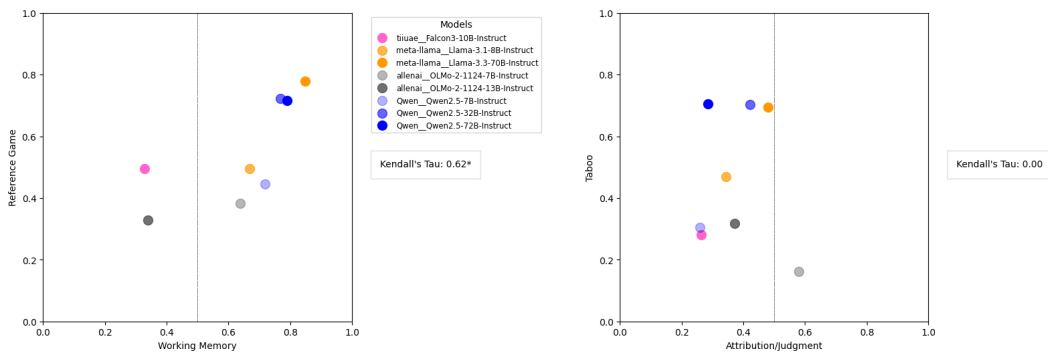


Figure 8: High (left) and low (right) correlation