
Wisent: A General Framework for Reliable Representation Identification and Representation Steering

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Representation engineering is a powerful method for identifying and modifying
2 high-level concepts within the internal layers of large language models. Despite its
3 potential, real-life deployments of activation steering remain difficult. We present
4 Wisent, a flexible, open-source framework for monitoring and steering internal
5 activations of large language models. Practical applications of the framework show
6 XXX percent hallucination reduction, XXX percent improvement in coding ability
7 and deep personalization capabilities.

8 1 Introduction

9 Large language models, with billions of parameters and Internet-scale training dataset, have displayed
10 significant capabilities across a wide range of tasks, such as writing, coding or reasoning.

11 However, their internal mechanisms of generating the next token cannot be precisely explained, with
12 interactions between layers and parameters increasing in complexity as the size of these models
13 increases.

14 Experiments with representation engineering (also known as steering or activation steering) have
15 shown activation modification to be a powerful method of identifying and influencing high-level
16 concepts (representations) within the layers of an LLM. Despite strong empirical performance on
17 selected truthfulness, safety or personalization tasks, representation engineering methods lack a
18 universal formulation and a unifying framework for understanding the underlying phenomenon,
19 comparing methods and applying them to new problems.

20 We propose Wisent, a modular framework for analyzing the internal mechanisms within a large
21 language model and influencing them to improve performance and individual alignment. Wisent
22 surpasses state of the art performance in identifying particular behaviors

23 2 Representation Engineering Problem

24 We formulate the **Representation Engineering Problem** as the following:

25 For a given model M and a Representation

26 Basic primitives and definitions of key terms are outlined in Appendix A.

27 **3 Representation Reading**

28 **3.1 Classifier**

29 **3.2 Detection Handling Method**

30 **4 Representation Control**

31 **4.1 Classifier**

32 **References**

- 33 [1] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David
34 Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis
35 with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- 36 [2] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase
37 from question-answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in
38 Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA, 2013. Association
39 for Computational Linguistics.
- 40 [3] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning
41 about physical commonsense in natural language. *arXiv preprint arXiv:1911.11641*, 2019.
- 42 [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla
43 Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language
44 models are few-shot learners. In *Advances in Neural Information Processing Systems (NeurIPS)*,
45 2020.
- 46 [5] Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald
47 Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q. Feldman, Arjun Guha,
48 Michael Greenberg, and Abhinav Jangda. Multipl-e: A scalable and extensible approach to
49 benchmarking neural code generation. *arXiv preprint arXiv:2208.08227*, 2022.
- 50 [6] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto,
51 Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul
52 Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke
53 Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad
54 Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias
55 Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebbgen Guss, Alex
56 Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,
57 William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra,
58 Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer,
59 Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech
60 Zaremba. Evaluating large language models trained on code. 2021.
- 61 [7] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and
62 Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions.
63 *arXiv preprint arXiv:1905.10044*, 2019.
- 64 [8] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,
65 and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning
66 challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- 67 [9] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christo-
68 pher Hesse, John Schulman, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, and
69 Jerry Tworek. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
70 2021.
- 71 [10] CodeParrot. Instructhumaneval, 2023.

- 72 [11] Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank:
 73 Investigating projection in naturally occurring discourse. In *Proceedings of Sinn und Bedeutung*
 74 23, volume 2, pages 107–124, Bellaterra (Cerdanyola del Vallès), 2019. Universitat Autònoma
 75 de Barcelona.
- 76 [12] Mingzhe Du, Anh Tuan Luu, Bin Ji, Liu Qian, and See-Kiong Ng. Mercury: A code efficiency
 77 benchmark for code llms. *arXiv preprint arXiv:2402.07844*, 2024.
- 78 [13] Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, Kang Zhu, Minghao Liu,
 79 Yiming Liang, and et al. Supergpqa: Scaling llm evaluation across 285 graduate disciplines.
 80 *arXiv preprint arXiv:2502.14739*, 2025.
- 81 [14] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gard-
 82 ner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs.
 83 *arXiv preprint arXiv:1903.00161*, 2019.
- 84 [15] Shai Gretz, Roni Friedman, Edo Cohen-Karlik, Assaf Toledo, Dan Lahav, Ranit Aharonov, and
 85 Noam Slonim. A large-scale dataset for argument quality ranking: Construction and analysis.
 86 In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- 87 [16] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo,
 88 Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding
 89 challenge competence with apps. *arXiv preprint arXiv:2105.09938*, 2021.
- 90 [17] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
 91 Jacob Steinhardt. Measuring massive multitask language understanding. In *Proceedings of the*
 92 *International Conference on Learning Representations (ICLR)*, 2021.
- 93 [18] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
 94 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset.
 95 *arXiv preprint arXiv:2103.03874*, 2021.
- 96 [19] HuggingFaceH4. Math-500, 2024.
- 97 [20] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. Mapping language to
 98 code in programmatic context. *arXiv preprint arXiv:1808.09588*, 2018.
- 99 [21] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Ar-
 100 mando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination
 101 free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 102 [22] Maxwell Jia. Aime problem set 2024, 2024.
- 103 [23] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. Triviaqa: A large scale
 104 distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th*
 105 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
 106 pages 1601–1611, Vancouver, Canada, 2017. Association for Computational Linguistics.
- 107 [24] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
 108 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
 109 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
 110 Petrov. Natural questions: A benchmark for question answering research. *Transactions of the*
 111 *Association for Computational Linguistics*, 7:452–466, 2019.
- 112 [25] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale
 113 reading comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- 114 [26] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott
 115 Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark
 116 for data science code generation. *arXiv preprint arXiv:2211.11501*, 2022.
- 117 [27] Ken Lang. Newsweeder: Learning to filter netnews. In *Proceedings of the Twelfth International*
 118 *Conference on Machine Learning*, pages 331–339, 1995.

- 119 [28] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic
120 human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 121 [29] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by
122 chatgpt really correct? rigorous evaluation of large language models for code generation. In
123 *NeurIPS 2023 Datasets and Benchmarks Track*, 2023.
- 124 [30] Junnan Liu, Hongwei Liu, Linchen Xiao, Ziyi Wang, Kuikun Liu, Songyang Gao, Wenwei
125 Zhang, Songyang Zhang, and Kai Chen. Are your llms capable of stable reasoning? *arXiv*
126 *preprint arXiv:2412.13147*, 2024.
- 127 [31] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin
128 Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long
129 Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng,
130 Shengyu Fu, and Shujie Liu. Codexglue: A machine learning benchmark dataset for code
131 understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.
- 132 [32] Math-AI. Aime problem set 2025, 2025.
- 133 [33] MathArena. Hmmt february 2025, 2025.
- 134 [34] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
135 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 136 [35] Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and
137 developing english math word problem solvers. *arXiv preprint arXiv:2106.15772*, 2021.
- 138 [36] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor
139 conduct electricity? a new dataset for open book question answering. In *Proceedings of the*
140 *2018 Conference on Empirical Methods in Natural Language Processing*. Association for
141 Computational Linguistics, 2018.
- 142 [37] Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, and et al. Humanity’s last exam. *arXiv*
143 *preprint arXiv:2501.14249*, 2025.
- 144 [38] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable
145 questions for squad. *arXiv preprint arXiv:1806.03822*, 2018.
- 146 [39] Siva Reddy, Danqi Chen, and Christopher D. Manning. Coqa: A conversational question
147 answering challenge. *arXiv preprint arXiv:1808.07042*, 2019.
- 148 [40] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
149 Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a
150 benchmark. *arXiv preprint arXiv:2311.12022*, 2023.
- 151 [41] Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. Choice of plausible
152 alternatives: An evaluation of commonsense causal reasoning. In *AAAI Spring Symposium on*
153 *Logical Formalizations of Commonsense Reasoning*, Stanford, CA, 2011.
- 154 [42] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An
155 adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.
- 156 [43] Shiqi Wang, Zheng Li, Haifeng Qian, Chenghao Yang, Zijian Wang, Mingyue Shang, Varun
157 Kumar, Samson Tan, Baishakhi Ray, Parminder Bhatia, Ramesh Nallapati, Murali Krishna
158 Ramanathan, Dan Roth, and Bing Xiang. Recode: Robustness evaluation of code generation
159 models. *arXiv preprint arXiv:2212.10264*, 2022.
- 160 [44] Yiming Wang, Pei Zhang, Jialong Tang, Haoran Wei, Baosong Yang, Rui Wang, Chenshu Sun,
161 Feitong Sun, Jiran Zhang, Junxuan Wu, Qiqian Cang, Yichang Zhang, Fei Huang, Junyang Lin,
162 Fei Huang, and Jingren Zhou. Polymath: Evaluating mathematical reasoning in multilingual
163 contexts. *arXiv preprint arXiv:2504.18428*, 2025.

- 164 [45] Pengcheng Yin, Bowen Deng, Edgar Chen, Bogdan Vasilescu, and Graham Neubig. Learning
 165 to mine aligned code and natural language pairs from stack overflow. In *Proceedings of the 15th*
 166 *IEEE/ACM International Conference on Mining Software Repositories (MSR)*, pages 476–486.
 167 IEEE, 2018.
- 168 [46] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial
 169 dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on*
 170 *Empirical Methods in Natural Language Processing*, pages 93–104, Brussels, Belgium, 2018.
 171 Association for Computational Linguistics.
- 172 [47] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a
 173 machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- 174 [48] Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme.
 175 Record: Bridging the gap between human and machine commonsense reading comprehension.
 176 *arXiv preprint arXiv:1810.12885*, 2018.
- 177 [49] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text
 178 classification. In *Advances in Neural Information Processing Systems*, volume 28, 2015.

179 **A Wisent Primitives**

180 **A.1 Model**

181 **A.2 Contrastive Pair**

182 **A.3 Activations**

183 **A.4 Activation Collection Method**

184 **A.5 Additional Utilities**

185 **B Representation Reading Functionalities**

186 **B.1 Classifier**

187 **B.2 Detection Handling Method**

188 **C Representation Control Functionalities**

189 **D Ablation**

190 **A All supported benchmarks**

191 This section enumerates all benchmarks used in our study, the task traits, the evaluation protocol, and
 192 the contrastive pair generation method applied to produce minimally perturbed negative targets. We
 193 first merged the *coding* and *mathematics* benchmark lists you provided and then appended them to
 194 the original master list.

195 **Contrastive pair generation methods (definitions)**

196 **Reading Comprehension Abstention Swap** [RC-Abstain] For extractive/open-domain RC: positive
 197 is the gold span; negative is an abstention (e.g., “Not provided in the text.”). If gold is
 198 *No answer*, the negative is a confident but wrong claim.

199 **Conversational Reading Comprehension Abstention** [ConvRC-Abstain] As RC-Abstain, but
 200 with dialogue context (CoQA). Negatives are generic abstentions; yes/no items are flipped
 201 when applicable.

- 202 **Language Modeling Corrupted Continuation** [LM-CorruptCont] Language modeling: positive
 203 is the true continuation; negative is a corrupted continuation (local shuffles/randomization)
 204 to break coherence.
- 205 **Generic answer** [Generic] Negative is some generic answer which is incorrect.
- 206 **Letter shuffling** [L-Shuff] Negative is created by shuffling letters of positive.
- 207 **Two-Choice Flip** [2C-Flip] Two-option tasks (PIQA, COPA, WinoGrande, CB): negative is simply
 208 the other option.
- 209 **Multichoice First Distractor** [MC-FirstDistr] Multi-choice tasks: negative = the first incorrect
 210 option in the provided order (deterministic).
- 211 **Multichoice Random Distractor** [MC-RandDistr] Multi-choice tasks: negative = a randomly cho-
 212 sen incorrect option from the same set (used for GPQA).
- 213 **Multichoice Letter Swap** [MC-LetterSwap] Multi-choice tasks scored over option letters (Truth-
 214 fulQA MC1/MC2): negative = the first incorrect letter.
- 215 **Exact Match Partial Mask** [EM-PartialMask] Exact-match free-form answers (HLE-EM): nega-
 216 tive is the gold text with partial token masking (approximately 1/3 words, or partial masking
 217 for single-word answers).
- 218 **Keyword-Preserving Token Deletion** [KP-Del] Coding tasks: negative program created by delet-
 219 ing non-keyword tokens while preserving syntax-critical keywords; aims to remain plausible
 220 but fail unit tests.
- 221 **Summary Content-Word Drop** [Summ-WordDrop] Code-to-text summarization: negative descrip-
 222 tion formed by dropping content words (nouns/verbs) while keeping scaffolding words to
 223 preserve superficial form.
- 224 **Numeric Offset (+1) Perturbation** [Num+1] Math QA: negative is the correct numeric answer
 225 offset by a small integer (typically +1); for non-integer answers, apply the minimal unit
 226 offset.
- 227 **Evaluation types (definitions)**
- 228 **Log-likelihood option scoring** [LL] The model scores each provided option/target by conditional
 229 log-probability given the prompt. Metrics typically compute accuracy over the highest-
 230 likelihood choice (MC tasks) or compare likelihoods of gold vs. negative targets.
- 231 **Text generation string matching** [TG] The model generates free-form text (or a number), which
 232 is then judged by task-specific metrics (e.g., exact match on numerical value for
 233 GSM8K/MATH; span/string matching for RC tasks; structured checks for DROP). Used
 234 also for CoT/generative GPQA variants and HLE-Exact-Match.
- 235 **Perplexity (language modeling)** [PPL] The model’s next-token distribution is evaluated over a
 236 reference text to compute Perplexity (lower is better). Used for language-modeling corpora
 237 like WikiText.
- 238 **Code execution against unit tests** [CE] The model generates code, which is executed in a sandbox
 239 against unit tests provided by a dataset (e.g., pass@1). Applies to HumanEval/MBPP/APPS,
 240 MultiPL-E, DS-1000, LiveCodeBench, etc.

Table 1: Benchmarks (short names), evaluation abbreviations, contrastive
 method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|------------------------|------|-----------------------|-----------------------|
| DROP [14] | [TG] | RC-Abstain | reading comprehension |
| ReCoRD [48] | [TG] | RC-Abstain | reading comprehension |
| SQuAD2 [38] | [TG] | RC-Abstain | reading comprehension |
| WebQuestions [2] | [TG] | RC-Abstain | factual QA |
| Natural Questions [24] | [TG] | RC-Abstain | factual QA |
| TriviaQA [23] | [TG] | RC-Abstain | factual QA |
| CoQA [39] | [TG] | ConvRC-Abstain | conversational RC |

| Benchmark | Eval | Method [CM] | Traits |
|------------------------|-----------|--------------------------------------|--------------------------|
| BoolQ [7] | [LL] | 2C-Flip | boolean RC |
| WinoGrande [42] | [LL] | 2C-Flip | commonsense |
| PIQA [3] | [LL] | 2C-Flip | commonsense |
| COPA [41] | [LL] | 2C-Flip | causal reasoning |
| HellaSwag [47] | [LL] | MC-FirstDistr | commonsense |
| SWAG [46] | [LL] | MC-FirstDistr | commonsense |
| OpenBookQA [36] | [LL] | MC-FirstDistr | science MCQ |
| ARC [8] | [LL] | MC-FirstDistr | science MCQ |
| RACE [25] | [LL] | MC-FirstDistr | RC (MC) |
| MMLU [17] | [LL] | MC-FirstDistr | multi-subject exams |
| GPQA [40] | [LL]/[TG] | MC-RandDistr | expert STEM exams |
| SuperGPQA [13] | [LL] | MC-FirstDistr | expert STEM exams |
| HLE [37] | [TG]/[LL] | EM-PartialMask; MC-FirstDistr | expert exams |
| GSM8K [9] | [TG] | Num+1 | mathematics |
| ASDiv [35] | [TG] | Num+1 | mathematics |
| Arithmetic [4] | [TG] | Num+1 | mathematics |
| MATH [18] | [TG] | Num+1 | mathematics (contest) |
| MATH-500 [19] | [TG] | Num+1 | mathematics (contest) |
| AIME [32][22] | [TG] | Num+1 | mathematics (contest) |
| HMMT [33] | [TG] | Num+1 | mathematics (contest) |
| PolyMath [44] | [TG] | Num+1 | mathematics (multiling.) |
| LiveMathBench [30] | [TG] | Num+1 | mathematics (EN/ZH) |
| MBPP [1] | [CE] | KP-Del | coding (Python) |
| HumanEval [6] | [CE] | KP-Del | coding (Python) |
| CoNaLa [45] | [CE] | KP-Del | coding (Python) |
| CONCODE [20] | [CE] | KP-Del | coding (Java) |
| Mercury [12] | [CE] | KP-Del | coding (multi-language) |
| HumanEval+ [29] | [CE] | KP-Del | coding (Python) |
| InstructHumanEval [10] | [CE] | KP-Del | coding (Python) |
| MBPP+ [29] | [CE] | KP-Del | coding (Python) |
| APPS [16] | [CE] | KP-Del | coding (Python) |
| DS-1000 [26] | [CE] | KP-Del | coding (Python) |
| MultiPL-E [5] | [CE] | KP-Del | coding (multi-language) |
| CodeXGLUE [31] | [TG] | Summ-WordDrop | coding (code-to-text) |
| ReCode [43] | [CE] | KP-Del | coding (Python) |
| LiveCodeBench [21] | [CE] | KP-Del | coding (Python) |
| TruthfulQA [28] | [LL] | MC-LetterSwap | truthfulness |
| CB [11] | [LL] | 2C-Flip | NLI |
| WikiText (2/103) [34] | [PPL] | LM-CorruptCont | language modeling |

241 **Category legend**

| | |
|--|-----------------------------|
| | RC/ODQA |
| | Multi-choice Reasoning |
| | Exams & Knowledge Tests |
| | Mathematics |
| | Coding |
| | Other (Truthfulness/NLI/LM) |

Method [CM] codes

| | |
|----------------|------------------------------|
| RC-Abstain | RC abstention swap |
| ConvRC-Abstain | Conversational RC abstention |
| LM-CorruptCont | LM corrupted continuation |
| 2C-Flip | Two-choice flip |
| MC-FirstDistr | First distractor (MC) |
| MC-RandDistr | Random distractor (MC) |
| MC-LetterSwap | Letter swap (MC) |
| Bool-Flip | Boolean flip |
| EM-PartialMask | Exact-match partial mask |
| KP-Del | Keyword-preserving deletion |
| Summ-WordDrop | Summary word drop |
| Num+1 | Numeric offset (+1) |

Abbreviation legend

| | |
|-------|--------------------------------|
| [LL] | Log-likelihood option scoring |
| [TG] | Text generation (string match) |
| [PPL] | Perplexity (LM) |
| [CE] | Code execution vs. unit tests |

Table 2: Benchmarks (short names), evaluation abbreviations, contrastive method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|-------------------------|-------|-----------------------|-----------|
| 20_newsgroups [27] | [TG] | MC-FirstDistr | reasoning |
| ag_news [49] | [TG] | MC-FirstDistr | reasoning |
| argument_topic [15] | [TG] | MC-FirstDistr | reasoning |
| banking77 [] | [TG] | MC-FirstDistr | reasoning |
| boolq [] | [LL] | 2C-Flip | reasoning |
| boolq-seq2seq [] | [TG] | 2C-Flip | reasoning |
| cb [] | [LL] | MC-FirstDistr | reasoning |
| claim stance topic [] | [TG] | MC-FirstDistr | reasoning |
| cnn dailymail [] | [TG] | Generic | reasoning |
| dpedia 14 [] | [TG] | MC-FirstDistr | reasoning |
| ethos binary [] | [TG] | MC-FirstDistr | reasoning |
| financial tweets [] | [TG] | MC-FirstDistr | reasoning |
| squadv2 [] | [TG] | RC-Abstain | reasoning |
| logieval [] | [TG] | MC-FirstDistr | reasoning |
| ledgar [] | [TG] | MC-FirstDistr | reasoning |
| logieval [] | [TG] | MC-FirstDistr | reasoning |
| penn treebank [] | [PPL] | LM-CorruptCont | reasoning |
| medical abstracts [] | [TG] | MC-FirstDistr | reasoning |
| unfair tos [] | [TG] | LM-CorruptCont | reasoning |
| record [] | [LL] | MC-FirstDistr | reasoning |
| stsB [] | [TG] | 2C-Flip | reasoning |
| sglue-rte [] | [LL] | 2C-Flip | reasoning |
| xsum [] | [TG] | Generic | reasoning |
| yahoo answers topics [] | [TG] | MC-FirstDistr | reasoning |

Table 3: Benchmarks (short names), evaluation abbreviations, contrastive method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|------------------------|------|--------------|-------------|
| afrimgsm direct amh [] | [TG] | Num+1 | mathematics |
| aime [] | [TG] | Num+1 | mathematics |
| aime2024 [] | [TG] | Num+1 | mathematics |
| aime2025 [] | [TG] | Num+1 | mathematics |
| gsm [] | [TG] | Num+1 | mathematics |

| Benchmark | Eval | Method [CM] | Traits |
|-------------|------|--------------|-------------|
| hmmt [] | [TG] | Num+1 | mathematics |
| math [] | [TG] | Num+1 | mathematics |
| math500 [] | [TG] | Num+1 | mathematics |
| polymath [] | [TG] | Num+1 | mathematics |

Table 4: Benchmarks (short names), evaluation abbreviations, contrastive method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|---------------------------|------|----------------------|--------------|
| afrimmlu direct amh [] | [LL] | MC-FirstDistr | multilingual |
| afrixnli en direct amh [] | [TG] | MC-FirstDistr | multilingual |
| arabic exams [] | [LL] | MC-FirstDistr | multilingual |
| bangla mmlu [] | [LL] | MC-FirstDistr | multilingual |
| basque glue [] | [LL] | 2C-Flip | multilingual |
| copa [] | [LL] | 2C-Flip | multilingual |
| global mmlu [] | [LL] | MC-FirstDistr | multilingual |
| m mmlu [] | [LL] | MC-FirstDistr | multilingual |
| m mmlu [] | [LL] | 2C-Flip | multilingual |
| noticia [] | [LL] | Generic | multilingual |
| phrases ca-va [] | [TG] | Generic | multilingual |
| wmt14 [] | [TG] | L-Shuff | multilingual |
| wmt16 [] | [TG] | L-Shuff | multilingual |

Table 5: Benchmarks (short names), evaluation abbreviations, contrastive method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|-------------|------|---------------------|-------------|
| babilong [] | [LL] | MC-RandDistr | longcontext |

Table 6: Benchmarks (short names), evaluation abbreviations, contrastive method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|--------------|------|----------------------|---------|
| glianorex [] | [LL] | MC-FirstDistr | medical |

Table 7: Benchmarks (short names), evaluation abbreviations, contrastive method (short), and traits. Versions merged where applicable.

| Benchmark | Eval | Method [CM] | Traits |
|----------------|-------|-----------------------|-------------------|
| wikitext103 [] | [PPL] | LM-CorruptCont | general knowledge |

²⁴² **B Per-Task Results**

²⁴³ **C Detailed Classification Results**

²⁴⁴ **D Benchmark-Aided Steering Results**

²⁴⁵ **E Optimal Sample Size Calculations**

²⁴⁶ **F Fully Synthetic Generation**

²⁴⁷ **G Agentic Capabilities**