
Subjective Randomness and In-Context Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) exhibit intricate capabilities, often achieving high
2 performance on tasks they were not explicitly trained for. The precise nature
3 of LLM capabilities is often unclear, with different prompts eliciting different
4 capabilities, especially when used with in-context learning (ICL). We propose a
5 “Cognitive Interpretability” framework that enables us to analyze ICL dynamics to
6 understand latent concepts underlying LLMs’ behavioral patterns. This provides
7 a more nuanced understanding than posthoc evaluation benchmarks, but does not
8 require observing model internals as a mechanistic interpretation would require.
9 Inspired by the cognitive science of human randomness perception, we use random
10 binary sequences as context and study dynamics of ICL by manipulating properties
11 of context data, such as sequence length. In the latest GPT-3.5+ models, we find
12 emergent abilities to generate pseudo-random numbers and learn basic formal
13 languages, with striking ICL dynamics where model outputs transition sharply
14 from pseudo-random behaviors to deterministic repetition.

15 1 Introduction

16 Large language models (LLMs), especially when prompted via in-context learning (ICL), demonstrate
17 complex, emergent capabilities [1–12]. Specifically, ICL yields task-specific behaviors in LLMs
18 via use of different prompts (or *contexts*) [1, 5, 13–20]. Although no weight updates occur in ICL,
19 different input contexts can activate, or re-weight, different latent algorithms in an LLM, analogous
20 to how traditional learning methods such as gradient descent use training data to re-weight model
21 parameters to learn representations [21–26]. Two seemingly equivalent prompts can, however, evoke
22 very different behaviors in LLMs [18]. Our central motivation is to interpret emergent capabilities and
23 latent *concepts* underlying complex behaviors in LLMs by analyzing in-context learning behavioral
24 dynamics, without directly observing hidden unit activations or re-training models on varied datasets.

25 Inspired by computational approaches to human cognition [27–31], we model and interpret latent
26 concepts evoked in LLMs by different contexts, without observing or probing model internals.
27 This approach, which we call **Cognitive Interpretability**, is a middle ground between shallow
28 test-set evaluation benchmarks on one hand [17, 32–40] and mechanistic neuron- and circuit-level
29 understanding of pre-trained models’ capabilities on the other [12, 41–53]. Computational cognitive
30 scientists have related algorithmic information theory to human cognition, where mental concepts
31 are viewed as programs, and cognitive hypothesis search over concepts is viewed as Bayesian
32 inference [30, 54–58]. In this vein, Griffiths and Tenenbaum [28] model subjective randomness in
33 human cognition as probabilistic program induction, where a person must search over a space of
34 non-random programs in order to answer the question, “was this sequence generated by a random
35 process?” We argue that ICL can similarly be seen as under-specified program induction, where
36 there is no single “correct” answer; instead, an LLM should appropriately re-weight latent algorithms.
37 The domain of random sequences reflects this framing, in contrast to other behavioral evaluation
38 methodologies, in that there is no correct answer to a random number generation or judgment task
39 (Fig. 1). If the correct behavior is to match a target random process, then the right way to respond to

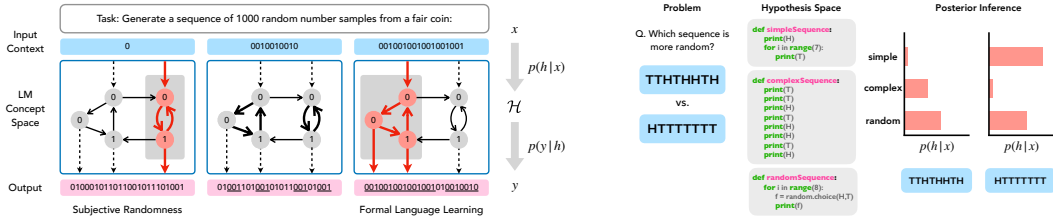


Figure 1: **Overview of our modeling framework.** (Left) Given a pre-trained LLM, we systematically vary input context prompts x . LLM outputs y vary as a function of x , based on some unknown latent concept space embedded in the LLM. With very little context ($x = 0$), GPT-3.5+ generates subjectively random sequences, whereas with adequate context matching a simple formal language ($x = 001001001001$), behavior becomes deterministic $((001)^n)$. (Right) Deciding whether a sequence is random can be viewed as search for a simple program that could generate that sequence. `HTTTTTTT` is described with a short program `simpleSequence` with higher $p(h)$ according to a simplicity prior, compared to `TTHTHHHTH` and `complexSequence`. Both sequences can be generated by `randomSequence`, with lower likelihood $p(x|h)$

40 a prompt *Generate N flips from a fair coin* is at best a uniform distribution over the tokens `Heads` and
 41 `Tails`, instead of a specific sequence, or a more complex algorithm that matches human behavior.

42 2 Background

43 **Bayesian Inference and In-Context Learning** A key methodological tool of cognitive modeling,
 44 recent work has also framed in-context learning as Bayesian inference over models [19, 59, 60].
 45 Specifically, the posterior predictive distribution $p(y|x)$ in these works describes how an LLM
 46 produces output tokens y , given the context, or *prompt*, x . The key assumption is that a context x
 47 will activate latent concepts c within a model according to their posterior probability $p(c|x)$, which
 48 the model marginalizes over to produce the next token y by sampling from the posterior predictive
 49 distribution: $p(y|x) = \int_{c \in \mathcal{C}} p(y|c) p(c|x)$. This model selection takes place in network activation
 50 dynamics, without changing its weights. In our experiments, we assume a hypothesis space \mathcal{H}
 51 that approximates the latent space of LLM concepts \mathcal{C} used when predicting the next token, i.e.,
 52 $p(y|x) = \sum_{h \in \mathcal{H}} p(y|h) p(h|x)$, where \sum_h can be changed to \max_h to represent deterministic
 53 greedy decoding with an LLM temperature parameter of 0. We specifically focus on Bernoulli
 54 processes, regular languages, Markov chains, and a simple memory-constrained probabilistic model
 55 as candidates for the hypothesis space \mathcal{H} for estimating LLM concepts in random binary sequences.
 56 We use a subset of regular languages $(x)^n$, where (x) is a short sequence of values, e.g., $(010)^n$,
 57 where 0 maps to Heads and 1 to Tails.

58 **Algorithmic and Subjective Randomness** Cognitive scientists studying *Subjective Randomness*
 59 model how people perceive randomness, or generate data that is subjectively random but algorithmically
 60 pseudo-random [29, 61–64]. In a bias termed *the Gambler’s Fallacy*, people reliably perceive
 61 binary sequences with long streaks of one value as less random, and judge binary sequences with
 62 higher-than-chance alternation rates as being “more random” than truly random sequences [65, 66].
 63 One way to study subjective randomness is to ask people whether a given data sequence was more
 64 likely to be generated by a Random process or a Non-Random process (Fig. 1). While the posterior
 65 distribution of all non-random processes includes every possible computable function, estimating this
 66 distribution can be simplified to finding the single most probable algorithm to approximate the full
 67 hypothesis space. If the hypotheses are data-generating programs, a natural prior $p(h)$ is to assign
 68 higher probabilities to programs with shorter description lengths, or lower complexity. This opti-
 69 mization problem is equivalent to computing the Kolmogorov complexity of a sequence $K(x)$ [67]
 70 and has motivated the use of “simplicity priors” in a number of domains in computational cognitive
 71 science [30, 54, 56]. Following previous work [28, 29], here we define subjective randomness of
 72 a sequence as the ratio of likelihood of that sequence under a random versus non-random model,
 73 i.e., $\text{randomness}(x) = \log P(x|\text{random}) - \log P(x|\text{non-random})$. The non-random likelihood
 74 $p(x|\text{non-random}) = 2^{-K(x)}$ denotes the probability of the minimal description length program that
 75 generates x , equivalent to Bayesian model selection: $P(x|\text{non-random}) = \max_{h \in \mathcal{H}} p(x|h) p(h)$. In
 76 this work, we study a small subset of \mathcal{H} , which includes formal languages and probabilistic models
 77 inspired by psychological models of human concept learning and subjective randomness [29, 66, 68].

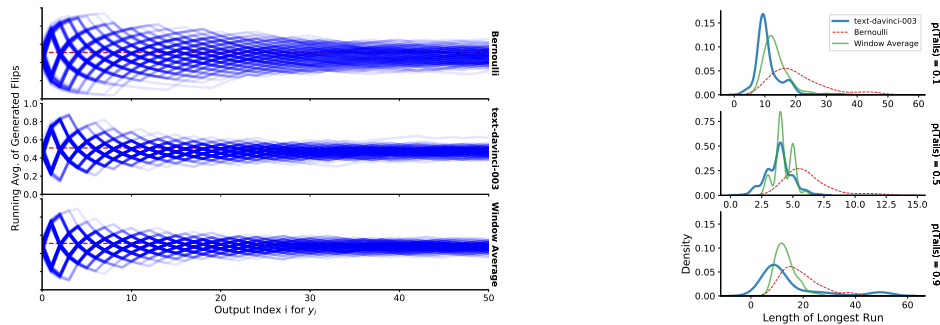


Figure 2: **GPT-3.5 generates pseudo-random binary sequences that deviate from a Bernoulli process.** (Left) Running averages of $p(\text{Tails})$ for flips generated by each model. Compared to a Bernoulli process, sequences generated by GPT and our Window Average model stay closer to the mean. (Right) GPT-3.5 shows a Gambler’s Fallacy bias, avoiding long runs of the same value in a row.

78 3 Experiments

79 **Randomness Generation and Judgment Tasks** In order to assess text generation dynamics and
 80 in-context concept learning, we evaluate LLMs on random sequence **Generation** tasks, analyzing
 81 responses according to simple interpretable models of *Subjective Randomness* and *Formal Language*
 82 *Learning*. In these tasks, the model generates a sequence y of binary values, or *flips*, comma-separated
 83 sequences of `Heads` or `Tails` tokens. We also analyze a smaller set of randomness **Judgment** tasks,
 84 where the prompt includes a sequence of flips, and the model must respond whether the sequence was
 85 generated by Random or Non-Random process. In both cases, y is a distribution over tokens with two
 86 possible values: Random or Non in Judgment tasks, indicating whether the sequence was generated
 87 by a random process with no correlation, or some non-random algorithm. We analyze dynamics in
 88 LLM-generated sequences y simulating a weighted coin with specified $p(\text{Tails})$, with $|x| \approx 0$.

89 **Subjective Randomness Models** We compare LLM-generated sequences to a ground truth “random”
 90 Bernoulli distribution with the same mean ($\mu = \bar{y}_{\text{LLM}}$), to a simple memory-constrained
 91 probabilistic model, and to Markov chains fit to model-generated data y . Hahn and Warren [68]
 92 theorize that the Gambler’s Fallacy emerges as a consequence of human memory limitations, where
 93 ‘seeming biases reflect the subjective experience of a finite data stream for an agent with a limited short-
 94 term memory capacity’. We formalize this as a simple *Window Average* model, which tends towards
 95 a specific probability p as a function of the last w flips: $p(y|x) = \max(0, \min(1, 2p - \bar{x}_{t-w\dots t}))$.

96 **Sub-Sequence Memorization and Complexity Metrics** Bender et al. [69] raise the question
 97 of whether LLMs are ‘stochastic parrots’ that simply copy data from the training set. To measure
 98 memorization, we look at the distribution of unique sub-sequences in y . If an LLM is repeating
 99 common patterns across outputs, potentially memorized from the training data, this should be
 100 apparent in the distribution over length K sub-sequences. Since there are deep theoretical connections
 101 between complexity and randomness [70, 71], we also consider the complexity of GPT-produced
 102 sequences. Compression is a metric of information content, and thus of redundancy over irreducible
 103 complexity [72, 73], and neural language models have been shown to prefer generating low complexity
 104 sequences [21]. As approximations of sequence complexity, we evaluate the distribution of Gzip-
 105 compressed file sizes [74] and inter-sequence Levenshtein distances [75].

106 **Formal Language Learning Metrics** In our Formal Language Learning analysis, x is a subset of
 107 regular expression repetitions of short token sequences such as $x \in (011)^n$, where longer sequences
 108 x correspond to larger n . This enables us to systematically investigate in-context learning of formal
 109 languages, as $|x|$ corresponds to the amount of data for inducing the correct program (e.g. $(011)^n$) out
 110 of the space of possible algorithms. In Randomness Judgment tasks, we assess formal concept learning
 111 by the dynamics of $p(y = \text{random} | x = C^{|x|})$ as a function of $|x|$. In Randomness Generation tasks,
 112 we assess concept learning according to the language model predictive distribution $p(y|x)$ over output
 113 sequences, inferred from next-token generation data: $p(y_{0\dots T}|x) = p(y_0|x) \prod_{t=1}^T p(y_t|y_{0\dots t-1}, x)$.
 114 Given a space of possible outputs y with length d , $y \in \{0, 1\}^d$, we estimate $p(y|x)$ by enumerating
 115 all y up to some depth d , and computing $\hat{p}(y_d|x, y_1, \dots, y_{d-1}) = \frac{1}{N} \sum_i (y_d = 1)^{(i)}$ as the fraction
 116 of N responses that are “Tails” (or equivalently, by using token-level probabilities directly). We
 117 estimate the predictive probability $p(y_t \in C|x, y_{0\dots t-1})$ assigned to a given regular language by
 118 computing the total probability mass for all trajectories in $y_{0\dots d}$ that exactly match C . For example,
 119 with $C = (011)^n$, there will be 3 trajectories $y_{0\dots d}$ that exactly match C , out of 2^d possible.

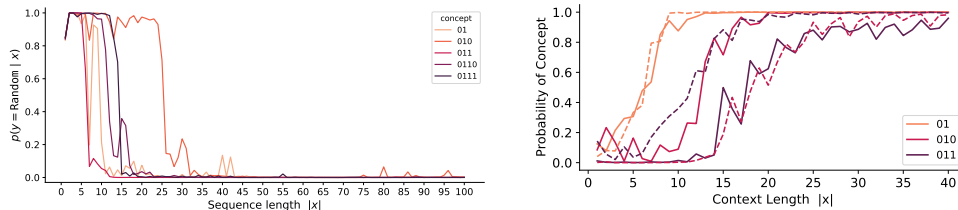


Figure 3: **Sharp transitions in predictive distributions for Randomness Judgment and Generation** (Left) In Randomness Judgment tasks, the predictive distribution $p(y = \text{random}|x)$ for text-davinci-003 transitions from high confidence in x being generated by a random process, to high confidence in a non-random algorithm (Right) in Generation tasks, the predictive $p(y = \text{Tails}|x)$ transitions from pseudo-randomness to deterministic repetition of a particular concept; text-davinci-003 is solid, gpt-3.5-turbo-instruct dashed.

120 4 Results

121 Subjectively Random Sequence Generation

122 In ‘InstructGPT’ models — text-davinci-003, ChatGPT (gpt-3.5-turbo,
 123 gpt-3.5-turbo-instruct) and GPT-4 — we find an emergent behavior of generating
 124 seemingly random binary sequences (Fig. 2). *This behavior is controllable*, where different $p(\text{Tails})$
 125 values leads to different means of generated sequences \bar{y} . However, the distribution of sequence
 126 means, as well as the distribution of the length of the longest runs for each sequence, deviate
 127 significantly from a Bernoulli distribution centered at \bar{y} , analogous to the Gambler’s Fallacy bias
 128 in humans. *Our Window Average model with a window size of $w = 5$ partly explains both biases*,
 129 matching GPT-generated sequences more closely than a Bernoulli distribution. Our cross-LLM
 130 analysis shows that text-davinci-003 is controllable with $P(\text{Tails})$, with a bias towards
 131 $\bar{y} = .50$ and higher variance in sequence means (though lower variance than a true Bernoulli
 132 process). ChatGPT (gpt-3.5-turbo-0301 and 0613) are similar for $P(\text{Tails}) < 50\%$, but
 133 behave erratically with higher $P(\text{Tails})$, with most y repeating ‘Tails’. GPT-4 (0301, 0613) shows
 134 stable, controllable subjective randomness behavior, with lower variances than text-davinci-003.
 135 Earlier models do not show subjective randomness behavior. Also see Appendix.

136 **Sub-Sequence Memorization and Complexity** We find significant differences between the
 137 distributions of sub-sequences for GPT-3.5-generated sequences and sequences sampled from a
 138 Bernoulli distribution (see Appendix for figures). This difference is partly accounted for with a
 139 Window Average model with a window size $w = 5$, although GPT repeats certain longer sub-
 140 sequences, for example length-20 sub-sequences, that are far longer than 5. However, the majority
 141 of sub-sequences have very low frequency, and though further experiments would be required
 142 to conclude that all sub-sequences are not memorized from training data, it seems unlikely that
 143 these were in the training set, since we find thousands of unique length- k (with varying k) sub-
 144 sequences generated at various values of $P(\text{Tails})$. *This indicates that GPT-3.5 combines dynamic,*
 145 *subjectively random sequence generation with distribution-matched memorization.* Across three
 146 metrics of sequence complexity — number unique sub-sequences, Gzip file size, and inter-sequence
 147 Levenshtein distance — we find that *GPT-3.5+ models, with the exception of ChatGPT, generate*
 148 *low complexity sequences*, showing that structure is repeated across sequences and supporting prior
 149 work [21, 73].

150 4.1 Distinguishing Formal Languages from Randomness

151 *GPT-3.5 sharply transitions between behavioral patterns, from generating pseudo-random values to*
 152 *generating non-random sequences that perfectly match the formal language (Fig. 3).* We observe
 153 a consistent pattern of formal language learning in GPT-3.5 generating random sequences where
 154 predictions $p(y|x)$ of depth $d \geq 4$ are initially random with small $|x|$, and have low $p(y \in C|x)$
 155 where C is a given concept. This follows whether the prompt describes the process as samples
 156 from “a weighted coin” or “a non-random-algorithm”. We also find *sharp phase changes in GPT-3.5*
 157 *behavioral patterns in Randomness Judgment tasks across 9 binary concepts (Fig. 3).* These follow
 158 a stable pattern of being highly confident in that the sequence is Random (high $p(y = \text{random}|x)$)
 159 when x is low, up to some threshold of context at which point it rapidly transitions to being highly
 160 confident in the process being non-random. Transition points vary between concepts, but the pattern
 161 is similar across concepts (see additional figures in Appendix).

References

- 162
- 163 [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
164 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
165 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 166 [2] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
167 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 168 [3] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani
169 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large
170 language models. *arXiv preprint arXiv:2206.07682*, 2022.
- 171 [4] Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao
172 Liu, Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently.
173 *arXiv preprint arXiv:2303.03846*, 2023.
- 174 [5] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
175 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
176 *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- 177 [6] Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych.
178 Are emergent abilities in large language models just in-context learning? *arXiv preprint*
179 *arXiv:2309.01809*, 2023.
- 180 [7] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von
181 Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the
182 opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 183 [8] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan
184 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv*
185 *preprint arXiv:2109.01652*, 2021.
- 186 [9] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
187 language models are zero-shot reasoners. *Advances in neural information processing systems*,
188 35:22199–22213, 2022.
- 189 [10] Alexander Pan, Jun Shern Chan, Andy Zou, Nathaniel Li, Steven Basart, Thomas Woodside,
190 Hanlin Zhang, Scott Emmons, and Dan Hendrycks. Do the rewards justify the means? measuring
191 trade-offs between rewards and ethical behavior in the machiavelli benchmark. In *International*
192 *Conference on Machine Learning*, pages 26837–26867. PMLR, 2023.
- 193 [11] Vishal Pallagani, Bharath Muppasani, Keerthiram Murugesan, Francesca Rossi, Biplav Srivas-
194 tava, Lior Horesh, Francesco Fabiano, and Andrea Loreggia. Understanding the capabilities of
195 large language models for automated planning. *arXiv preprint arXiv:2305.16151*, 2023.
- 196 [12] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin
197 Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic
198 task. *arXiv preprint arXiv:2210.13382*, 2022.
- 199 [13] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing
200 Xu, and Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 201 [14] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale
202 Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables
203 complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.
- 204 [15] Xinyi Wang, Wanrong Zhu, and William Yang Wang. Large language models are implicitly
205 topic models: Explaining and finding good demonstrations for in-context learning. *arXiv*
206 *preprint arXiv:2301.11916*, 2023.
- 207 [16] Seongjin Shin, Sang-Woo Lee, Hwijeen Ahn, Sungdong Kim, HyoungSeok Kim, Boseop
208 Kim, Kyunghyun Cho, Gichang Lee, Woomyoung Park, Jung-Woo Ha, et al. On the effect
209 of pretraining corpora on in-context learning by a large-scale language model. *arXiv preprint*
210 *arXiv:2204.13509*, 2022.

- 211 [17] Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and
212 Luke Zettlemoyer. Rethinking the Role of Demonstrations: What Makes In-Context Learning
213 Work?, October 2022. URL <http://arxiv.org/abs/2202.12837>. arXiv:2202.12837 [cs].
- 214 [18] Chenglei Si, Dan Friedman, Nitish Joshi, Shi Feng, Danqi Chen, and He He. Measuring
215 inductive biases of in-context learning with underspecified demonstrations. *arXiv preprint*
216 *arXiv:2305.13299*, 2023.
- 217 [19] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An Explanation of
218 In-context Learning as Implicit Bayesian Inference, July 2022. URL <http://arxiv.org/abs/2111.02080>. arXiv:2111.02080 [cs].
- 220 [20] Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does
221 in-context learning learn? bayesian model averaging, parameterization, and generalization.
222 *arXiv preprint arXiv:2305.19420*, 2023.
- 223 [21] Micah Goldblum, Marc Finzi, Keefer Rowan, and Andrew Gordon Wilson. The no free lunch
224 theorem, kolmogorov complexity, and the role of inductive biases in machine learning. *arXiv*
225 *preprint arXiv:2304.05366*, 2023.
- 226 [22] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can
227 gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers.
228 In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*,
229 2023.
- 230 [23] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to imple-
231 ment preconditioned gradient descent for in-context learning. *arXiv preprint arXiv:2306.00297*,
232 2023.
- 233 [24] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander
234 Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by
235 gradient descent. In *International Conference on Machine Learning*, pages 35151–35174.
236 PMLR, 2023.
- 237 [25] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What
238 learning algorithm is in-context learning? investigations with linear models. *arXiv preprint*
239 *arXiv:2211.15661*, 2022.
- 240 [26] Yingcong Li, M Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as
241 algorithms: Generalization and implicit model selection in in-context learning. *arXiv preprint*
242 *arXiv:2301.07067*, 2023.
- 243 [27] Joshua Tenenbaum. Bayesian modeling of human concept learning. *Advances in neural*
244 *information processing systems*, 11, 1998.
- 245 [28] Thomas L Griffiths and Joshua B Tenenbaum. From Algorithmic to Subjective Randomness.
246 *Neural Information Processing Systems*, 2003.
- 247 [29] Thomas L. Griffiths, Dylan Daniels, Joseph L. Austerweil, and Joshua B. Tenenbaum. Subjective
248 randomness as statistical inference. *Cognitive Psychology*, 103:85–109, June 2018. ISSN
249 00100285. doi: 10.1016/j.cogpsych.2018.02.003. URL <https://linkinghub.elsevier.com/retrieve/pii/S0010028517302281>.
- 251 [30] Steven T Piantadosi, Joshua B Tenenbaum, and Noah D Goodman. Bootstrapping in a language
252 of thought: A formal model of numerical concept learning. *Cognition*, 123(2):199–217, 2012.
- 253 [31] Michael J Spivey, Sarah E Anderson, and Rick Dale. The phase transition in human cognition.
254 *New Mathematics and Natural Computation*, 5(01):197–220, 2009.
- 255 [32] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid,
256 Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al.
257 Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.
258 *arXiv preprint arXiv:2206.04615*, 2022.

- 259 [33] Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal
260 analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
- 261 [34] Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny
262 Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. Measuring
263 faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*, 2023.
- 264 [35] Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas Liao, Kamilė Lukošiuūtė, Anna Chen,
265 Anna Goldie, Azalia Mirhoseini, Catherine Olsson, Danny Hernandez, et al. The capacity for
266 moral self-correction in large language models. *arXiv preprint arXiv:2302.07459*, 2023.
- 267 [36] Samuel R Bowman, Jeeyoon Hyun, Ethan Perez, Edwin Chen, Craig Pettit, Scott Heiner, Kamile
268 Lukosuite, Amanda Askell, Andy Jones, Anna Chen, et al. Measuring progress on scalable
269 oversight for large language models. *arXiv preprint arXiv:2211.03540*, 2022.
- 270 [37] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez,
271 Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language
272 models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- 273 [38] Ethan Perez, Sam Ringer, Kamilė Lukošiuūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig
274 Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model
275 behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*, 2022.
- 276 [39] Ben Prystawski and Noah D Goodman. Why think step-by-step? reasoning emerges from the
277 locality of experience. *arXiv preprint arXiv:2304.03843*, 2023.
- 278 [40] Miles Turpin, Julian Michael, Ethan Perez, and Samuel R Bowman. Language models don't
279 always say what they think: Unfaithful explanations in chain-of-thought prompting. *arXiv
280 preprint arXiv:2305.04388*, 2023.
- 281 [41] Neel Nanda, Lawrence Chan, Tom Liberum, Jess Smith, and Jacob Steinhardt. Progress
282 measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- 283 [42] Gabriel Goh, Nick Cammarata, Chelsea Voss, Shan Carter, Michael Petrov, Ludwig Schubert,
284 Alec Radford, and Chris Olah. Multimodal neurons in artificial neural networks. *Distill*, 6(3):
285 e30, 2021.
- 286 [43] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers
287 are key-value memories. *arXiv preprint arXiv:2012.14913*, 2020.
- 288 [44] Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational
289 Linguistics*, 48(1):207–219, 2022.
- 290 [45] Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.
291 Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv
292 preprint arXiv:2211.00593*, 2022.
- 293 [46] Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse
294 engineering how networks learn group operations. *arXiv preprint arXiv:2302.03025*, 2023.
- 295 [47] Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris
296 Bertsimas. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint
297 arXiv:2305.01610*, 2023.
- 298 [48] Alex Foote, Neel Nanda, Esben Kran, Ioannis Konstas, Shay Cohen, and Fazl Barez. Neuron to
299 graph: Interpreting language model neurons at scale. *arXiv preprint arXiv:2305.19911*, 2023.
- 300 [49] Ekdeep Singh Lubana, Eric J Bigelow, Robert P Dick, David Krueger, and Hidenori Tanaka.
301 Mechanistic mode connectivity. In *International Conference on Machine Learning*, pages
302 22965–23004. PMLR, 2023.
- 303 [50] Tom Lieberum, Matthew Rahtz, János Kramár, Geoffrey Irving, Rohin Shah, and Vladimir
304 Mikulik. Does circuit analysis interpretability scale? evidence from multiple choice capabilities
305 in chinchilla. *arXiv preprint arXiv:2307.09458*, 2023.

- 306 [51] Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang.
307 Hidden progress in deep learning: Sgd learns parities near the computational limit. *Advances in*
308 *Neural Information Processing Systems*, 35:21750–21764, 2022.
- 309 [52] Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams.
310 Towards understanding grokking: An effective theory of representation learning. *Advances in*
311 *Neural Information Processing Systems*, 35:34651–34663, 2022.
- 312 [53] Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two
313 stories in mechanistic explanation of neural networks. *arXiv preprint arXiv:2306.17844*, 2023.
- 314 [54] Nick Chater and Paul Vitányi. Simplicity: a unifying principle in cognitive science? *Trends in*
315 *cognitive sciences*, 7(1):19–22, 2003.
- 316 [55] Stanislas Dehaene, Fosca Al Roumi, Yair Lakretz, Samuel Planton, and Mathias Sablé-Meyer.
317 Symbols and mental programs: a hypothesis about human singularity. *Trends in Cognitive*
318 *Sciences*, 2022.
- 319 [56] Noah D Goodman, Joshua B Tenenbaum, Jacob Feldman, and Thomas L Griffiths. A rational
320 analysis of rule-based concept learning. *Cognitive science*, 32(1):108–154, 2008.
- 321 [57] Tomer D Ullman, Noah D Goodman, and Joshua B Tenenbaum. Theory learning as stochastic
322 search in the language of thought. *Cognitive Development*, 27(4):455–480, 2012.
- 323 [58] Tomer D Ullman and Joshua B Tenenbaum. Bayesian models of conceptual development:
324 Learning as building models of the world. *Annual Review of Developmental Psychology*, 2:
325 533–558, 2020.
- 326 [59] Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers
327 as Algorithms: Generalization and Stability in In-context Learning, February 2023. URL
328 <http://arxiv.org/abs/2301.07067>. arXiv:2301.07067 [cs, stat].
- 329 [60] Michael Hahn and Navin Goyal. A Theory of Emergent In-Context Learning as Implicit Struc-
330 ture Induction, March 2023. URL <http://arxiv.org/abs/2303.07971>. arXiv:2303.07971
331 [cs].
- 332 [61] An T. Oskarsson, Leaf Van Boven, Gary H. McClelland, and Reid Hastie. What’s next? Judging
333 sequences of binary events. *Psychological Bulletin*, 135(2):262–285, 2009. ISSN 1939-1455,
334 0033-2909. doi: 10.1037/a0014821. URL [http://doi.apa.org/getdoi.cfm?doi=10.](http://doi.apa.org/getdoi.cfm?doi=10.1037/a0014821)
335 [1037/a0014821](http://doi.apa.org/getdoi.cfm?doi=10.1037/a0014821).
- 336 [62] Giorgio Gronchi and Steven A. Sloman. Regular and random judgements are not two sides
337 of the same coin: Both representativeness and encoding play a role in randomness perception.
338 *Psychonomic Bulletin & Review*, 28(5):1707–1714, October 2021. ISSN 1069-9384, 1531-
339 5320. doi: 10.3758/s13423-021-01934-9. URL [https://link.springer.com/10.3758/](https://link.springer.com/10.3758/s13423-021-01934-9)
340 [s13423-021-01934-9](https://link.springer.com/10.3758/s13423-021-01934-9).
- 341 [63] Florent Meyniel, Maxime Maheu, and Stanislas Dehaene. Human inferences about sequences:
342 A minimal transition probability model. *PLoS computational biology*, 12(12):e1005260, 2016.
- 343 [64] Samuel Planton, Timo van Kerkoerle, Leila Abbih, Maxime Maheu, Florent Meyniel, Mariano
344 Sigman, Liping Wang, Santiago Figueira, Sergio Romano, and Stanislas Dehaene. A theory of
345 memory for binary sequences: Evidence for a mental compression algorithm in humans. *PLOS*
346 *Computational Biology*, 17(1):e1008598, January 2021. ISSN 1553-7358. doi: 10.1371/journal.
347 [pcbi.1008598](https://dx.plos.org/10.1371/journal.pcbi.1008598). URL <https://dx.plos.org/10.1371/journal.pcbi.1008598>.
- 348 [65] Ruma Falk and Clifford Konold. Making sense of randomness: Implicit encoding as a basis for
349 judgment. 1997.
- 350 [66] Raymond S Nickerson. The production and perception of randomness. *Psychological review*,
351 109(2):330, 2002.
- 352 [67] Ming Li, Paul Vitányi, et al. *An introduction to Kolmogorov complexity and its applications*.
353 Springer, 1997.

- 354 [68] Ulrike Hahn and Paul A Warren. Perceptions of randomness: why three heads are better than
355 four. *Psychological review*, 116(2):454, 2009.
- 356 [69] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On
357 the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021*
358 *ACM conference on fairness, accountability, and transparency*, pages 610–623, 2021.
- 359 [70] Gregory J Chaitin. Algorithmic information theory. *IBM journal of research and development*,
360 21(4):350–359, 1977.
- 361 [71] Gregory J Chaitin. *Information, randomness & incompleteness: papers on algorithmic informa-*
362 *tion theory*, volume 8. World Scientific, 1990.
- 363 [72] David JC MacKay. *Information theory, inference and learning algorithms*. Cambridge university
364 press, 2003.
- 365 [73] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christo-
366 pher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, et al.
367 Language modeling is compression. *arXiv preprint arXiv:2309.10668*, 2023.
- 368 [74] Zhiying Jiang, Matthew YR Yang, Mikhail Tsirlin, Raphael Tang, and Jimmy Lin. Less is more:
369 Parameter-free text classification with gzip. *arXiv preprint arXiv:2212.09410*, 2022.
- 370 [75] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and
371 reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union, 1966.