

# Supplementary Materials

Anonymous Author(s)  
 Affiliation  
 Address  
 email

## 1 Appendices

### 2 A Overview and Further Background

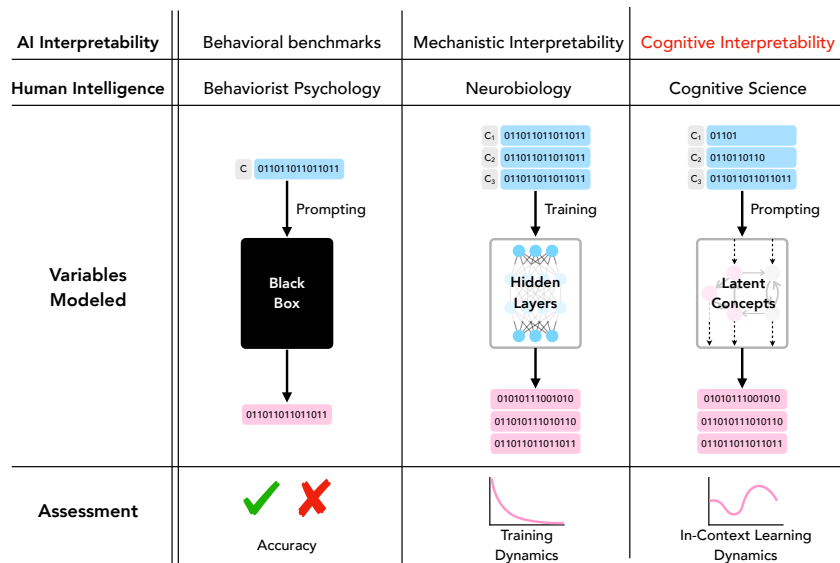


Figure 1: **Cognitive Interpretability in the context of prior work** Cognitive Interpretability studies in-context learning dynamics in LLMs, positing theories of the latent concepts enabling behavioral capabilities. It is a middle ground between behavioral benchmarks, which treat models as black boxes and evaluate hit-or-miss accuracy, and mechanistic interpretability, which studies toy models and training loss dynamics, analogous to how cognitive science is a middle ground between behaviorist psychology and neurobiology in the study of human intelligence.

3 A multi-layer neural networks may possess multiple distributed circuits implementing computational  
 4 primitives, such as basic mathematical operations like addition and sequence copying [1–6]. In state  
 5 of the art LLMs with hundreds of billions (even trillions) of parameters, there may be sub-networks  
 6 implementing various computations, and a wide variety of emergent behaviors corresponding to  
 7 those computations. A similar situation occurs in research focused on understanding the human brain,  
 8 where sub-networks of neuronal cells have been shown to localize specific capabilities. In cognitive  
 9 science, researchers study such aspects of cognition without fully understanding the underlying  
 10 neural circuitry, often modeling behavior without observing brain activity. Analogously, we seek to

11 understand the structure of behaviors in large language models in the wild, without a full mechanistic  
12 understanding of their circuitry (Figure 1).

13 The reliance on benchmarks to evaluate LLM capabilities loosely parallels early behaviorist psy-  
14 chology, when theories of human and animal learning only assumed stimulus-response associations,  
15 without positing theories about mental processes or neural substrates [7]. Circuit-level *mechanistic*  
16 interpretations of neural networks parallel neurobiology, which offers physical models of information  
17 processing in biological neurons, but typically does not account for the structure of high-level be-  
18 havior. Cognitive interpretability, like cognitive science, is aimed at predicting behavioral outcomes  
19 over a potentially infinite space of possible tasks. It defines high-level specifications of behaviors  
20 performed by LLMs, which, we argue, should be the first step before mechanistic understanding of  
21 circuits existent in a model is pursued.

22 Cognitive scientists have used Bayesian predictive and posterior distributions to model learning  
23 dynamics in human adults and children, as well as in non-human animals [8–10]. A discrete  
24 hypothesis space in a probabilistic model can give clear and meaningful explanations to learning  
25 patterns analogous to mode collapse and phase transitions in deep learning. When behavior suddenly  
26 shifts from one pattern, or mode, of behavior to another, this can be understood as one hypothesis  
27 coming to dominate the posterior  $p(h|x)$  as  $x$  grows in scale [11]. Such cognitive analysis of  
28 behavioral shifts as generated from a shift in posterior probability between one hypothesis to another  
29 parallels recent work on learning dynamics in training and prompting LLMs. E.g., sharp phase  
30 changes have been observed when training neural networks, corresponding to the formation of  
31 identifiable mechanisms such as modular addition circuits or induction heads [1, 2, 4, 12]. ICL  
32 research has explored how few-shot prompting can significantly boost LLM performance in various  
33 domains, and how the particular exemplars provided in context determine its overall effectiveness [13–  
34 19]. For further discussion of related work, see Appendix C.

### 35 Learning Dynamics in Model Selection

36 Cognitive scientists have used Bayesian predictive and posterior distributions to model learning  
37 dynamics in human adults and children, as well as in non-human animals [8–10]. A discrete  
38 hypothesis space in a probabilistic model can give clear and meaningful explanations to learning  
39 patterns analogous to mode collapse and phase transitions in deep learning. When behavior suddenly  
40 shifts from one pattern, or mode, of behavior to another, this can be understood as one hypothesis  
41 coming to dominate the posterior  $p(h|x)$  as  $x$  grows in scale [11]. For example, children between the  
42 ages of 3.5–5 years old learning to count undergo a dramatic conceptual shift from knowing the  
43 meanings of only a few number words (“one”, “two”) to a full inductive understanding of counting,  
44 which can be modeled as Bayesian model selection with a simplicity prior over models [20]. Such  
45 cognitive analysis of behavioral shifts as generated from a shift in posterior probability between one  
46 hypothesis to another parallels recent work on learning dynamics in training and prompting LLMs.  
47 E.g., sharp phase changes have been observed when training neural networks, corresponding to the  
48 formation of identifiable mechanisms such as modular addition circuits or induction heads [1, 2, 4, 12].  
49 ICL research has explored how few-shot prompting can significantly boost LLM performance  
50 in various domains, and how the particular exemplars provided in context determine its overall  
51 effectiveness [13–19]. Work on chain-of-thought reasoning in LLMs demonstrates how a few  
52 exemplars of detailed solutions or even a simple prompt like “let’s think this through step-by-step”  
53 can dramatically impact model performance [21–23]. For further discussion of related work, see  
54 Appendix C.

## 55 B Additional Experimental Details

56 All calls were made with the OpenAI API, using default parameters including — important to our  
57 analysis — a temperature parameter of 1.0. We use token-wise log probabilities  $p(y_t|y_{0..t-1})$  from  
58 the OpenAI API where available for cost efficiency and since this is equivalent to drawing repeated  
59 token samples and computing the fraction of samples, e.g.  $N_{\text{Tail}s}/(N_{\text{Head}s} + N_{\text{Tail}s})$ .

60 In the Generation task, the context  $x$  includes the prompt question, as well as an initial set of coin  
61 flips that follow the beginning of the “answer” section of the context. Prompt context in these  
62 experiments includes a specific probability, shown in Fig. 2, where the last \_\_ marks where the model  
63 begins generating tokens  $y$ . In subjective randomness experiments, an initial flip ‘Heads’ is used



90 finite automata. We use a subset of regular languages  $(x)^n$ , where  $(x)$  is a short sequence of values,  
91 e.g.,  $(010)^n$ , where 0 maps to Heads and 1 to Tails.

## 92 C Related Work

93 **Formal languages and transformers.** A number of recent works explore how transformers and other  
94 neural language models learn formal languages [26–36]. One common theme is that neural networks  
95 often learn ‘shortcuts’, degenerate representations of formal languages that fail out-of-samples.

96 **In-Context Learning as Bayesian inference** A number of recent works frame ICL as Bayesian  
97 model selection [12, 37–40]. Two key differences in our work are: first, we analyze state-of-the-art  
98 LLMs based on behaviors alone, whereas prior work trains models from scratch on synthetic data  
99 and analyzes model parameters directly. Second, we consider Bayesian inference as an empirical  
100 modeling framework, as well as a theory, whereas these works only do the latter.

101 **Mechanistic interpretability of transformer models** Prior work has characterized specific circuit-  
102 level implementations of simple high-level behaviors such as sequence copying, modular addition,  
103 and other primitive computational operations [2, 4, 41–52]. Our work differs in that we model  
104 hypothetical algorithms to characterize LM output behavioral patterns, without observing underlying  
105 activation patterns. We see this as analogous to cognitive science complementing neuroscience in  
106 the understanding of human cognition. We characterize the high-level “cognitive” representations in  
107 LLMs as a step towards connecting low-level explanations of neural circuits, such as induction heads,  
108 with sophisticated high-level behaviors that are characteristic of LLMs.

109 **Language model evaluations** Our work resembles evaluation benchmarks such as BIG-Bench [23]  
110 that use behavior alone to evaluate LM understanding and reasoning. However, as described later,  
111 the domain of subjective randomness is fundamentally different in that there is no “correct” answer.  
112 Linguistic probing attempts to characterize the structure of LM representations, but unlike our work,  
113 is a function of hidden unit activations rather than output behavior.

114 **LLM Text Generation Dynamics** Work on chain-of-thought reasoning in LLMs demonstrates  
115 how a few exemplars of detailed solutions or even a simple prompt like “let’s think this through  
116 step-by-step” can dramatically impact model performance [21–23], but typically only the model’s  
117 final answer is analyzed, not the trajectory of its intermediate steps. Our memory-constrained Window  
118 Average model, inspired by Hahn and Warren [53], is similar in spirit to the claim of Prystawski and  
119 Goodman [54], that ‘[chain-of-thought] reasoning emerges from the locality of experience’. Zhang  
120 et al. [55] demonstrate that invalid reasoning can snowball in LLMs, where hallucinations during  
121 intermediate steps lead to hallucinations in the final answer.

122 **Random number generation in LLMs** Renda et al. [56] explore random number generation in  
123 LLMs, in addition to cursory explorations by [57, 58]. These investigations do not analyze dynamics  
124 of sequence generation, nor do they ground their analysis, as we do, in theories of ICL as Bayesian  
125 model selection and the cognitive science of subjective randomness. Ortega et al. [59] uses a similar  
126 domain as ours with random binary sequences and has a similar binary tree visualization over possible  
127 sequences, but they train models from scratch and analyze model hidden states, rather than behavioral  
128 trajectories as we do.

129 **Bayesian program learning in cognitive science** Our work is inspired by computational cognitive  
130 science work that theoretically treats concepts as programs, and empirically uses structured Bayesian  
131 models to understand human cognition in various domains [8, 10, 11, 60]. We use models based on  
132 the cognitive science of subjective randomness [61], drawing particularly on the Bayesian program  
133 induction definitions of subjective randomness in Griffiths and Tenenbaum [62, 63], Griffiths et al.  
134 [64]. Our method of studying learning as probabilistic inference over formal languages with varying  
135  $|x|$  is also similar to Goodman et al. [9], Piantadosi et al. [20], Yang and Piantadosi [65], Bigelow  
136 and Piantadosi [66], who use more sophisticated grammar-based models of concept learning.



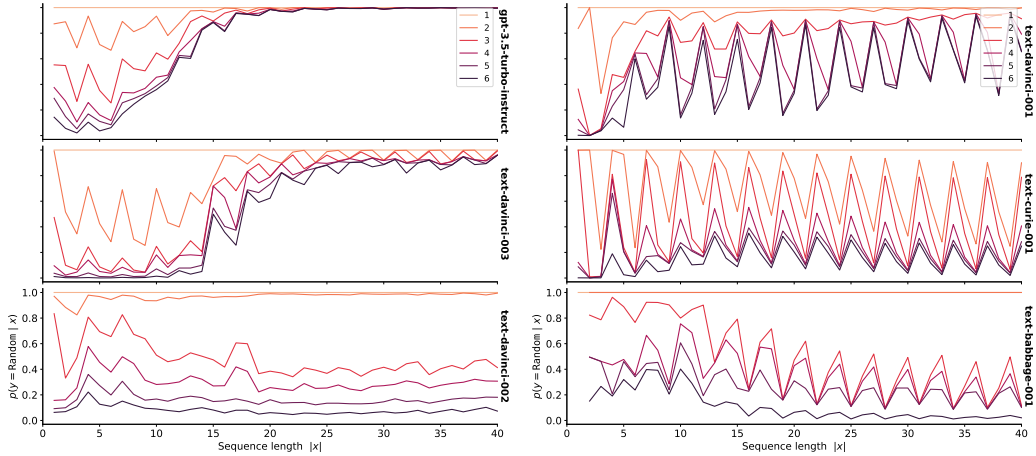


Figure 5: **Predictive distributions  $p(y|x)$  by each LLM for Concept  $(011)^n$ , at each prediction depth  $d$ .** Colors correspond to different prediction depths, also refer to Figure 3.

152 **Additional Predictive  $p(y|x)$  Trees**

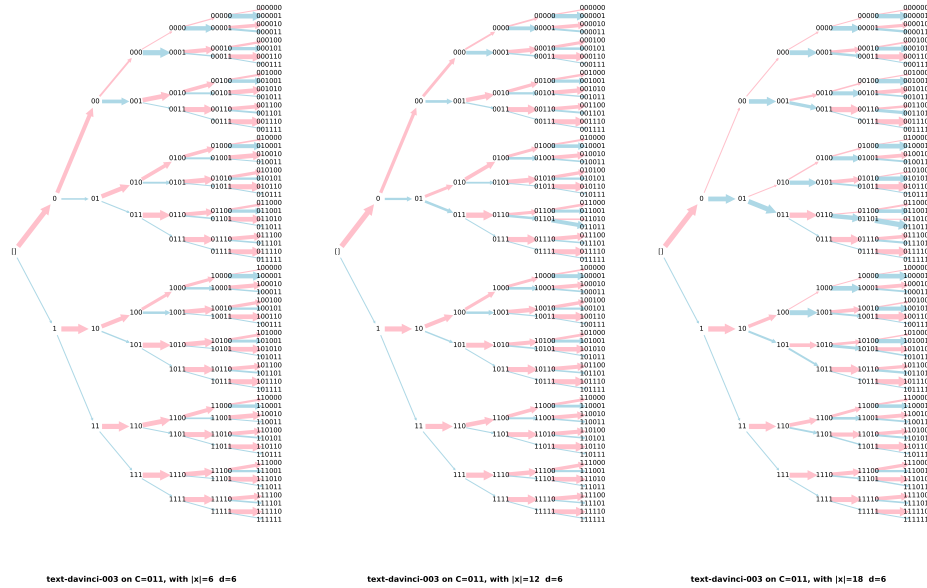


Figure 6: **Predictive distribution  $p(y|x)$  trees with  $d = 6$  for concept  $C = (011)^n$  with  $|x| \in \{6, 12, 18\}$ .** Since  $|x|$  is increasing by the same depth as the tree  $\Delta_{|x|} = d = 6$ , the transition from generating pseudo-random numbers to deterministically repeating 011 is visibly apparent. Also see Figure 3.

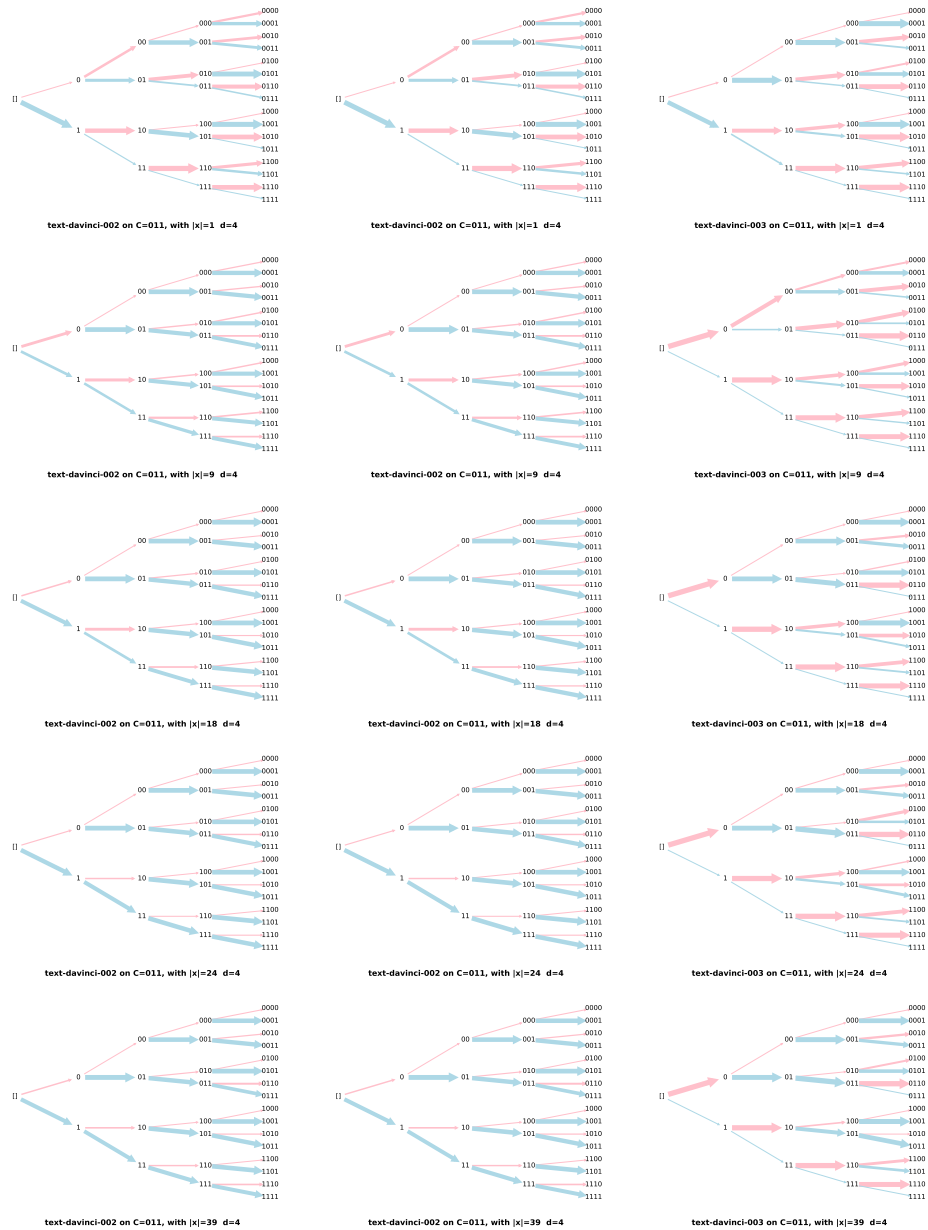


Figure 7: Predictive distribution  $p(y|x)$  trees with  $d = 4$  for concept  $C = (011)^n$  with  $|x| \in \{1, 9, 18, 24, 39\}$ . Models shown are text-davinci-002, text-davinci-003, and gpt-3.5-turbo-instruct. Also see Figure 3.



153 **E Randomness Judgments**

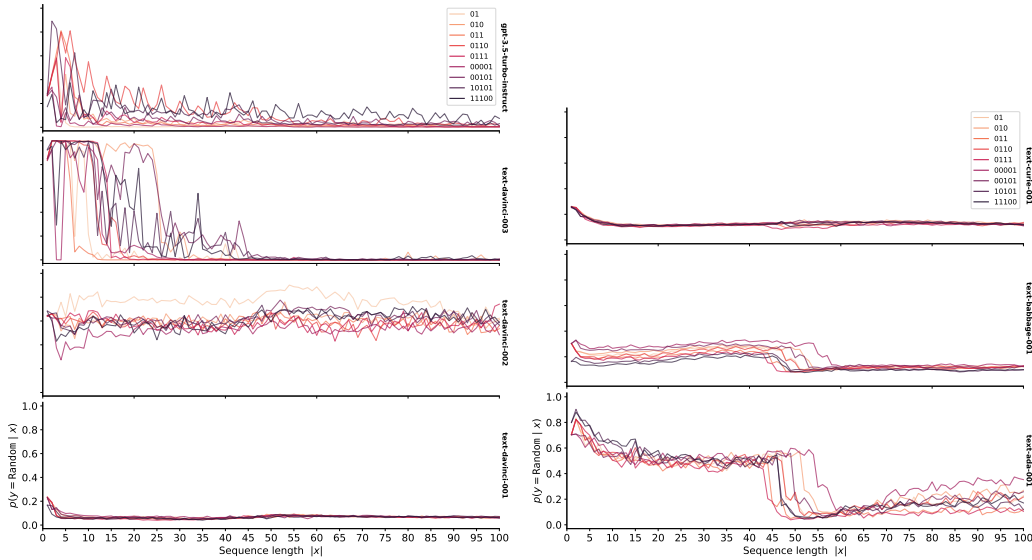


Figure 8: Randomness judgments across GPT models for 9 concepts.

154 (Fig 8) `text-davinci-003` shows a stable pattern of being highly confident (high token probability)  
 155 in the process being random up to some amount of context  $|x|$ , at which point it rapidly transitions to  
 156 being highly confident in the process being non-random, with transition points varying substantially  
 157 between concepts. `chat-gpt-3.5-instruct` does not go through a stable high-confidence random  
 158 period like `text-davinci-003`, and stable high-to-low confidence dynamics are observed for only a  
 159 subset of concepts. The majority of earlier GPT models (`text-davinci-002`, `text-davinci-001`,  
 160 `text-curie-001`, `text-babbage-001`) show no ‘formal language learning’, at all. However,  
 161 surprisingly OpenAI’s smallest available GPT model `text-ada-001` shows S-shaped in-context  
 162 learning dynamics, with the peak close to .5 instead of 1.0 as in `text-davinci-003`. Additionally,  
 163 the learning dynamics and transition points for all concepts appear nearly identical, approximately at  
 164  $|x| = 50$ , and some concepts show less stable “non-random” patterns for larger  $|x|$ .



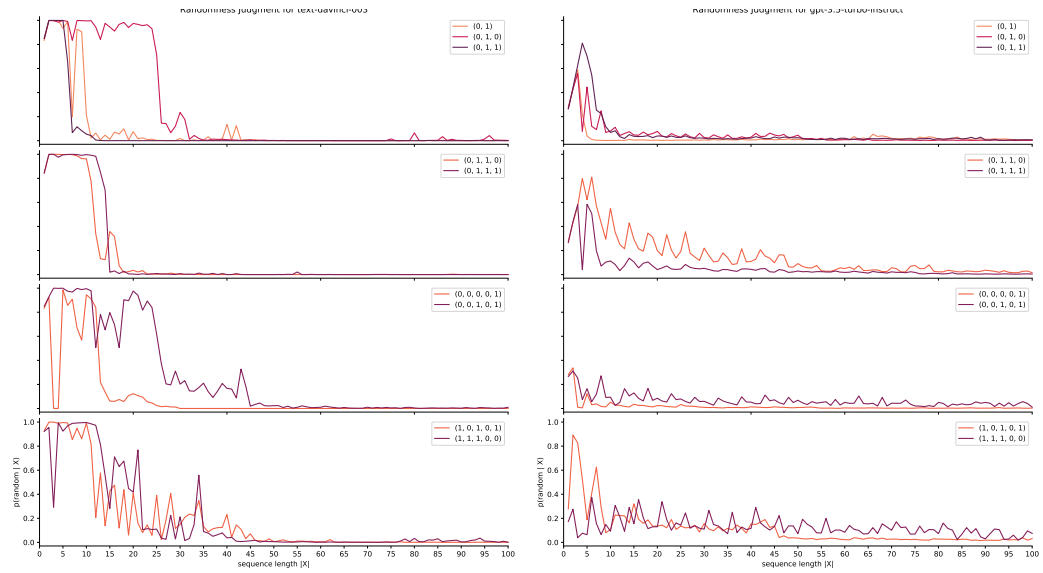


Figure 9: Randomness Judgment  $p(y = \text{random}|x)$  dynamics for each concept tested, for text-davinci-003 and gpt-3.5-turbo-instruct

165 **F Random Sequence Generation by GPT Model**

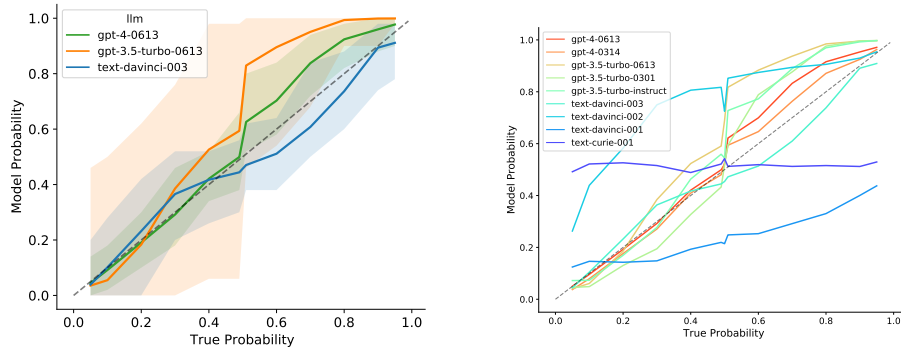


Figure 10: **Probability  $p(\text{Tails})$  bias across LLMs** text-davinci-003 and GPT-4 models are least biased relative to the specified  $p(\text{Tails})$  (x-axis). In the left figure, error bars represent the maximum and minimum sequence means  $\bar{y}$  for each  $p(\text{Tails})$ .

166 Our cross-LLM analysis (Fig. 10, 11) shows that text-davinci-003 is controllable with  $P(\text{Tails})$ ,  
 167 with a bias towards  $\bar{y} = .50$  and higher variance in sequence means (though lower variance than a  
 168 true Bernoulli process). ChatGPT (gpt-3.5-turbo-0301 and 0613) demonstrate similar behavior  
 169 for  $P(\text{Tails}) < 50\%$ , but behave erratically with higher  $P(\text{Tails})$  and the majority of sequences  $y$   
 170 converge to repeating ‘Tails’. GPT-4 (0301, 0613) show stable, controllable subjective randomness  
 171 behavior, but with lower variances than sequences generated by text-davinci-003. Earlier models  
 172 do not show subjective randomness behavior, with text-davinci-002 and text-davinci-001  
 173 being heavily biased and uncontrollable, and text-curie-001 generates sequences with  $\bar{y} = .50$   
 174 regardless of  $P(\text{Tails})$ .

175 Fig. 11 and the left side of Fig. 10 demonstrate that text-davinci-003 and GPT-4 models not  
 176 only are more controllable, following the correct probability more closely on average, but also have  
 177 substantially lower variance than ChatGPT, which is both less controllable and has more variability  
 178 in its distribution of responses. Further, GPT-4 is lower variance than text-davinci-003, with  
 179 sequences staying even closer to their means  $\bar{y}$ .

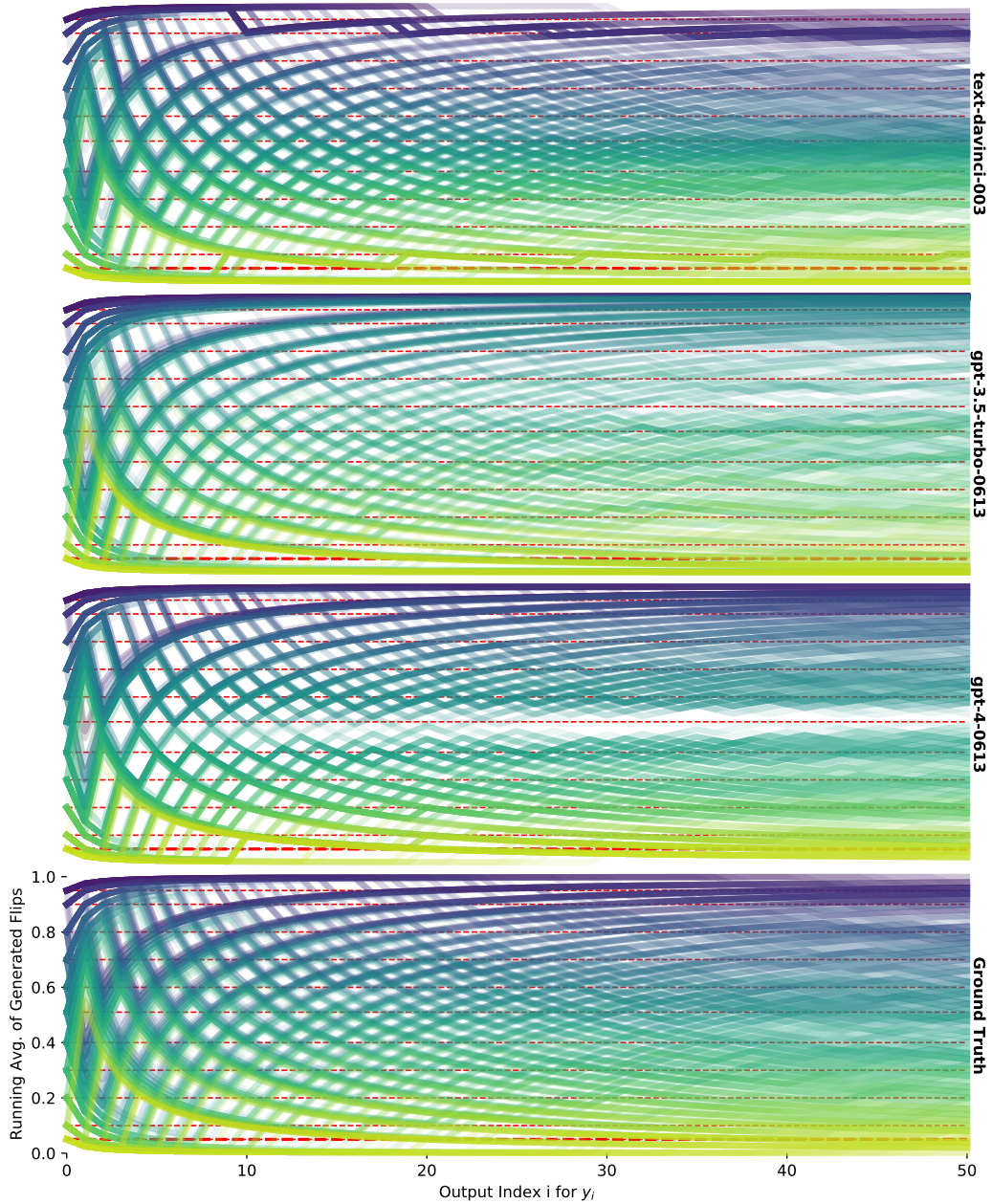


Figure 11: 50 sequences sampled by each GPT model, for each  $p(Tails)$ . Color is assigned according to specified  $p(Tails)$ . Red dotted lines are drawn for each  $p(Tails)$ .

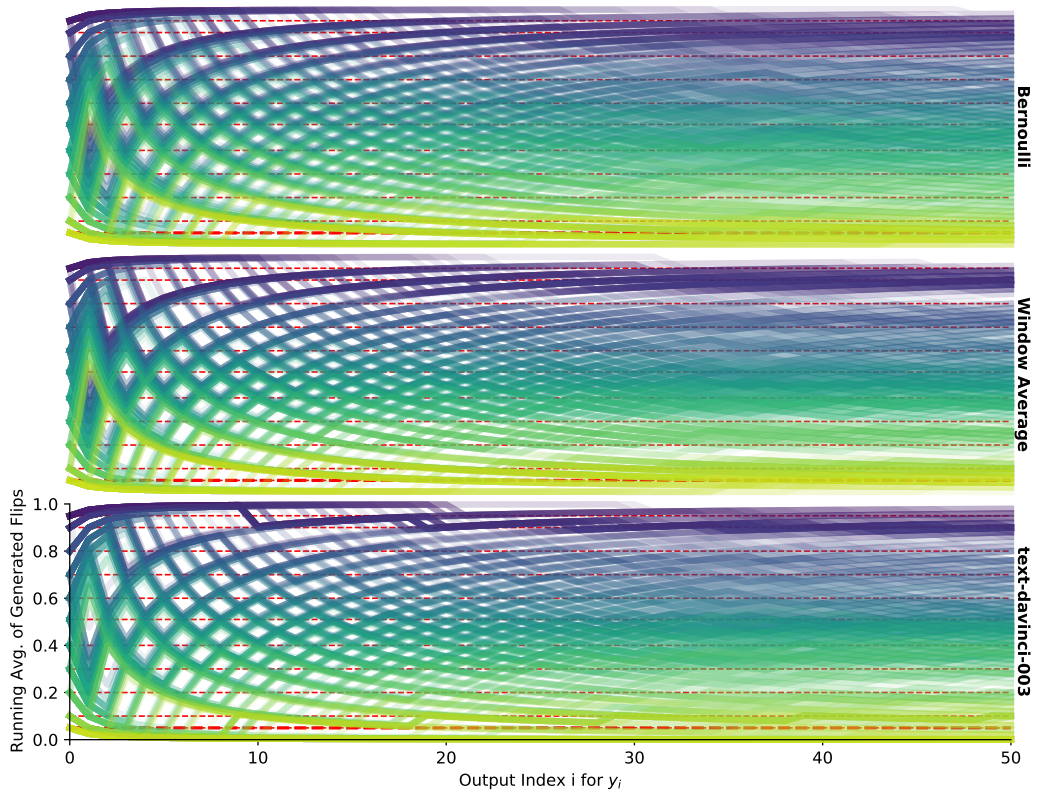


Figure 12: 50 sequences sampled by `text-davinci-003`, for each  $p(Tails)$ , compared with samples from Bernoulli and Window Average models fit to  $y_{LLM}$  for each  $p(Tails)$ . Color is assigned according to specified  $p(Tails)$ .

180 **G Gambler’s Fallacy Metrics by GPT Model**

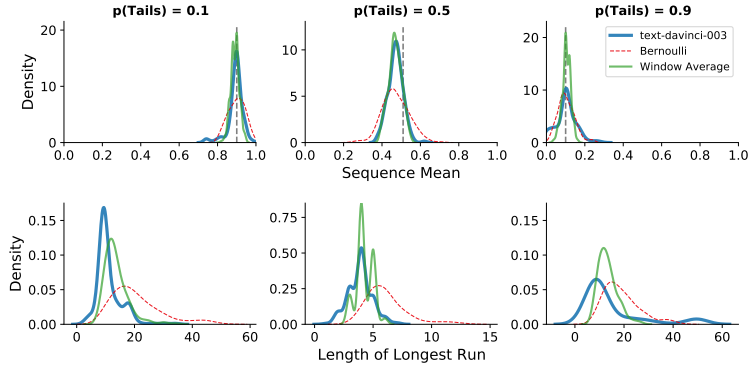


Figure 13: **GPT-3.5 shows a Gambler’s fallacy bias of avoiding long runs.** (Top) Distribution of mean values of flip sequences ( $\mu = \frac{1}{T} \sum_t y_t$ ) generated by GPT-3.5 (`text-davinci-003`) with the specified  $p(Tails)$ , compared with a Bernoulli process and our Window Average model with the same mean as the GPT-3.5 flips. Flips generated by GPT approximately follow the expected mean  $p(Tails)$ , but have lower variance than a Bernoulli distribution. (Bottom) Length of the longest run for each sequence, where a run is a sub-sequence of the same value repeating. In this case, we see a clear bias in GPT-3.5 to avoid long runs, with a similar pattern across all values of  $p(Tails)$  despite the x-axis changing in scale.

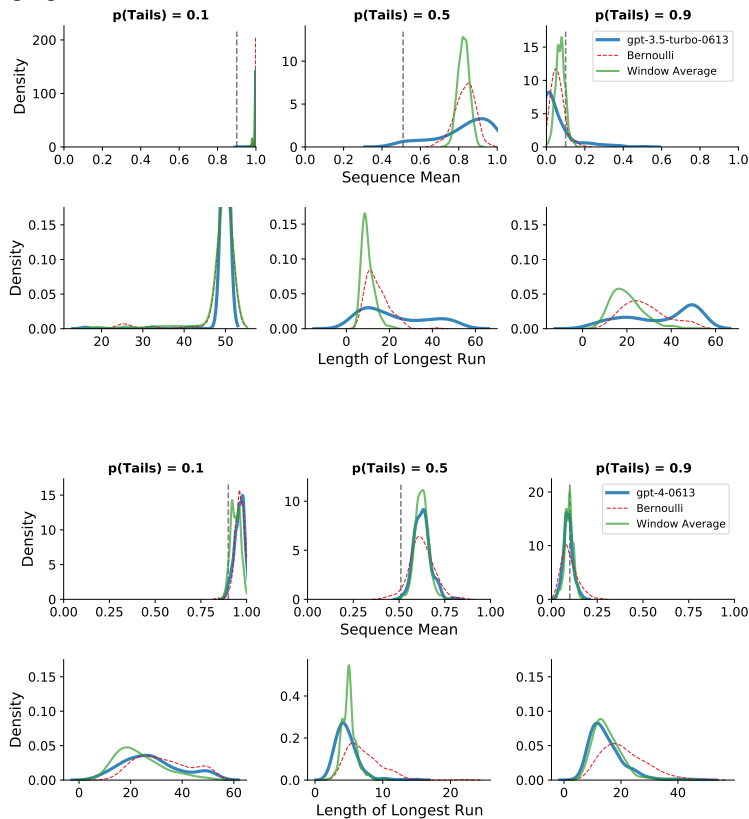


Figure 14: **Gambler’s Fallacy histograms for ChatGPT (Top) and GPT-4 (Bottom).** Also see Fig. 13.

181 ChatGPT shows no clear Gambler’s Fallacy bias, whereas GPT-4 does show this pattern, but is less  
 182 pronounced than `text-davinci-003` (Fig. 14).

183 In both plots of Fig. 15, we observe that `text-davinci-003` shows a Gambler’s Fallacy  
 184 bias across  $p(Tails)$ , of higher-than-chance alternation rates and shorter runs; ChatGPT



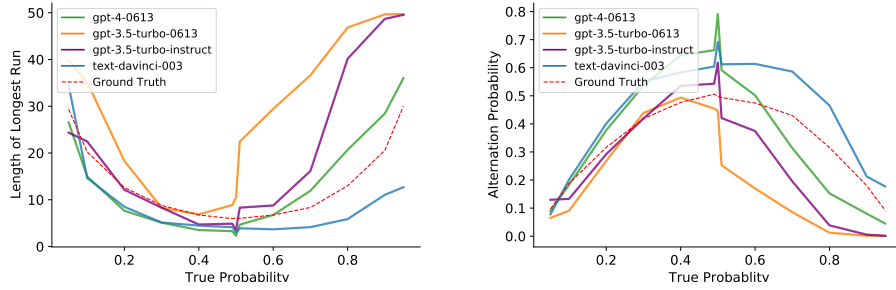


Figure 15: **Comparing metrics of Gambler’s Fallacy across probabilities and LLMs** (Left) The mean longest run for each sequence  $y$ , at each specified probability  $p(\text{Tail}s)$ , where a run is a consecutive sub-sequence of the same flip repeating multiple times in a row. (Right) The mean alternation rate for each LLM, where alternation rate is the fraction of consecutive flips that are not equal  $p(y_t \neq y_{t-1})$ .

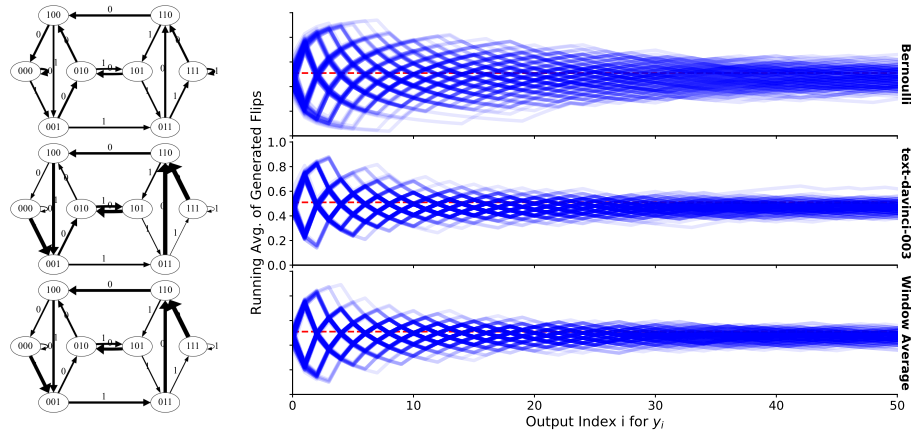


Figure 16: **GPT-3.5 generates pseudo-random binary sequences that deviate from a Bernoulli process.** (Left) Empirical conditional probabilities for a third-order Markov Chain fit to sequences  $y$  generated by GPT-3.5 text-davinci-003, a Bernoulli process centered at the mean of GPT sequence  $\bar{y}$ , and our Window Average model ( $w = 5$ ). In the simulated Bernoulli process, edges are fairly uniform; the conditional probabilities for GPT-3.5 and the Window Average model demonstrate a similar non-uniform bias. (Right) Running averages for flip sequences generated by each model, where 0 denotes ‘Heads’ and 1 denotes ‘Tails’. Compared to a Bernoulli process (top), sequences generating using GPT (middle) and those of our Window Average model (bottom) stay closer to the mean, repeating the same patterns more often.

185 (gpt-3.5-turbo-0613) produces more tails-biased and higher-variance sequences  $y$  when  
 186  $p(\text{Tail}s) > 50\%$ ; GPT-4 and gpt-3.5-turbo-instruct interpolate between the two distinct  
 187 trends of text-davinci-003 and ChatGPT. The red dotted line represents a Bernoulli process with  
 188 mean  $p(\text{Tail}s)$ .

189 It is unclear how the capabilities we identify are implemented at a circuit level, or why they only  
 190 seem to emerge in the most powerful and heavily tuned GPT models. For the latter, one hypothesis  
 191 is that internet corpora contain text with human-generated or human-curating subjectively random  
 192 binary sequences, and fine-tuning methods such as instruction fine-tuning, supervised fine-tuning, and  
 193 RLHF make LLMs more controllable, enabling them to apply previously inaccessible capabilities in  
 194 appropriate circumstances. Another hypothesis is that these fine-tuning methods bias LLMs towards  
 195 non-repetitiveness, or induce some other general bias that plays a role in the in-context learning  
 196 dynamics we observe in our particular domain. We hope that future work in cognitive and mechanistic  
 197 interpretability will shed further light on these questions.

198 **H Memorization, Compression, and Complexity**

199 Across three metrics of sequence complexity — number unique sub-sequences, Gzip file size, and  
 200 inter-sequence Levenshtein distance (see Fig. 20, 19 in Appendix) — we find that *GPT-3.5+ models*,  
 201 *with the exception of ChatGPT, generate low complexity sequences*, showing that structure is repeated  
 202 across sequences and supporting Goldblum et al. [67], Delétang et al. [68]. By the metrics of mean  
 203 Levenshtein distance and number of unique sub-sequences, ChatGPT generates higher complexity  
 204 sequences than chance. We speculate that this phenomenon might explained by a cognitive model  
 205 that avoids sampling with replacement.

206 For the Generation task, we note that with a specification of  $P(Tails) = 50%$ , but not 49%, 51%  
 207 or other values, sequences  $y$  generated by GPT-3.5+ are dominated by repeating ‘Heads, Tails,  
 208 Heads, Tails, ...’. This pattern is consistent across variations of the prompts listed in Fig. 2,  
 209 including specifying ‘fair’ or ‘unweighted’ instead of a ‘weighted coin’, and produces a visible kink  
 210 in many cross- $p(Tails)$  metrics (Fig. 20, 19, 15). For this reason, in Fig. 13 we show results for  
 211  $P(Tails) = 51%$ .

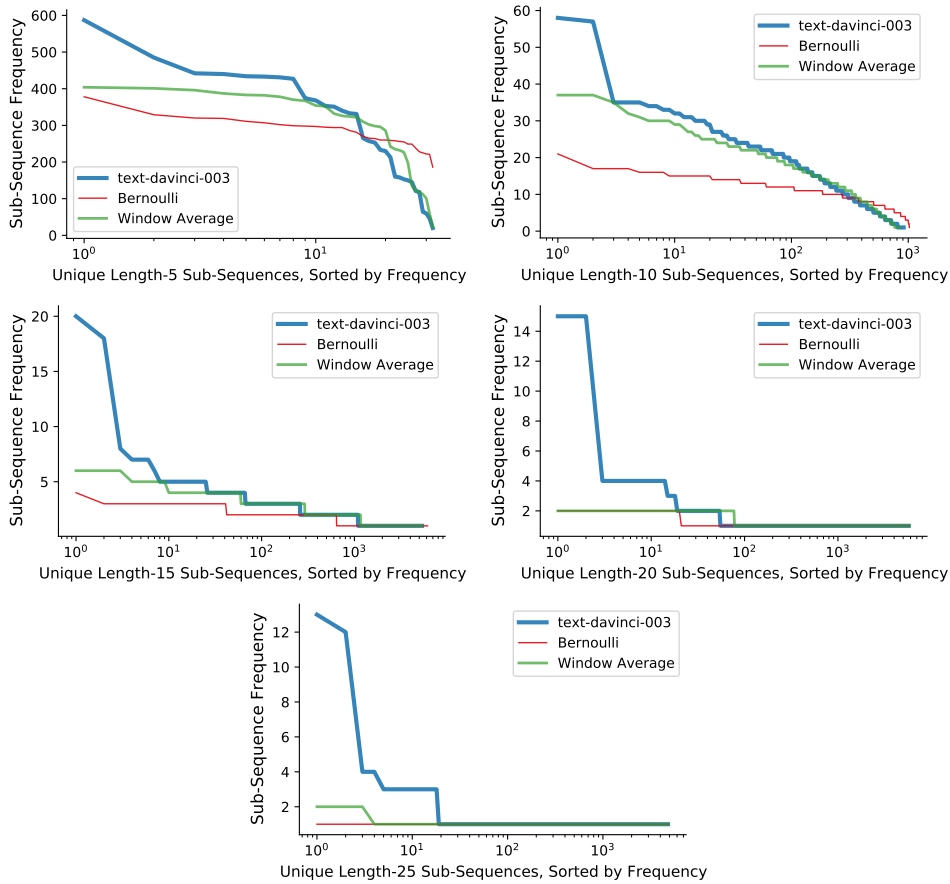


Figure 17: **Distribution of unique sub-sequences for text-davinci-003 for varying sub-sequence lengths**

212 In Figure 17, we find that GPT repeats specific sub-sequences more often than chance (Bernoulli  
 213 with  $\mu = \bar{y}$ ), or what is predicted by our Window Average model. While the Window Average model  
 214 (green) generates fewer unique sub-sequences than a Bernoulli process (red), this does not account  
 215 for the bias in GPT-3.5 (*text-davinci-003*, in blue) to repeat many of the same sub-sequences.  
 216 This disparity increases with longer sub-sequences



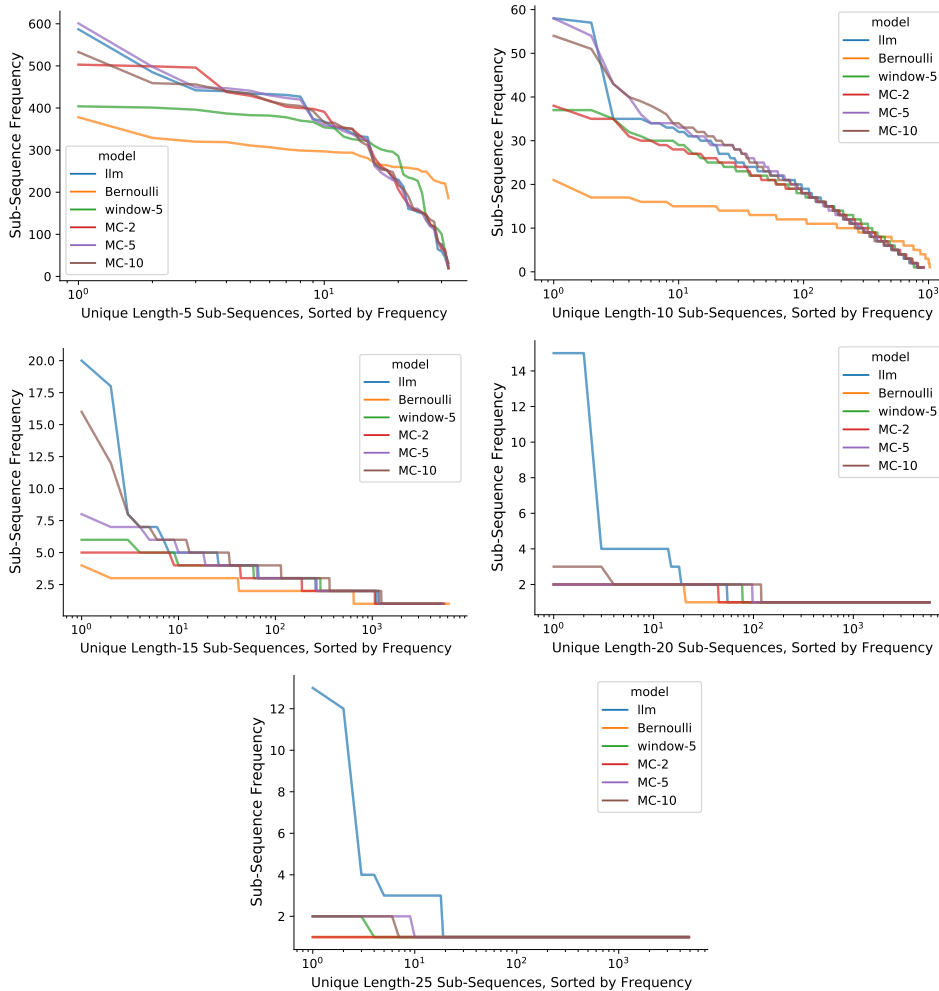


Figure 18: **Distribution of unique sub-sequences for text-davinci-003, with additional models, varying sub-sequence lengths** MC-2, MC-5, and MC-10 are Markov Chain models fit to GPT-3.5 flips, with orders  $k = \{2, 5, 10\}$

217 In Fig. 18, we show that Markov chains of high order  $k$  can account for the sub-sequence distribution,  
 218 but this only applies when  $k \leq w$  where  $w$  is the sub-sequence length, and the Markov chains can  
 219 effectively memorizing the sub-sequence distribution of  $y$ .

220 Across both unnormalized and normalized distributions of unique sub-sequences (Fig. 19), we find  
 221 that GPT-4 repeats the same length-10 sub-sequences significantly more than the other models, and  
 222 both ChatGPT-based models (gpt-3.5-turbo-0613, gpt-3.5-turbo-instruct) follow different  
 223 patterns for  $p(\text{Tail}s) < 50\%$  and  $p(\text{Tail}s) > 50\%$ , even when controlling for sequence bias (Right).  
 224 The only model that generates more unique sub-sequences than chance (above dotted line) is ChatGPT  
 225 (gpt-3.5-turbo-0613).

226 As a coarse approximation of sequence complexity, we use Gzip file size of appended sequences  
 227  $\text{gzip}(y : y' : y'' : \dots)$  and mean Levenshtein distance between sequences  $d(y, y')$ . Gzip [69],  
 228 a common algorithm for file compression that is highly optimized for compressing strings with  
 229 redundancy into small file sizes, and Gzip file size has been found to be an effective feature extractor  
 230 for NLP [70]. Levenshtein distance [71] is a measure of edit distance between two strings.

231 Since sequence compression is highly correlated with probability, e.g. all sequences with  $\bar{y} = 0.99$   
 232 will be highly compressible, we normalize the distribution of both plots in Fig. 20 by dividing by  
 233 the same metric (appended Gzip size, or mean Levenshtein distance) for a Bernoulli distribution centered

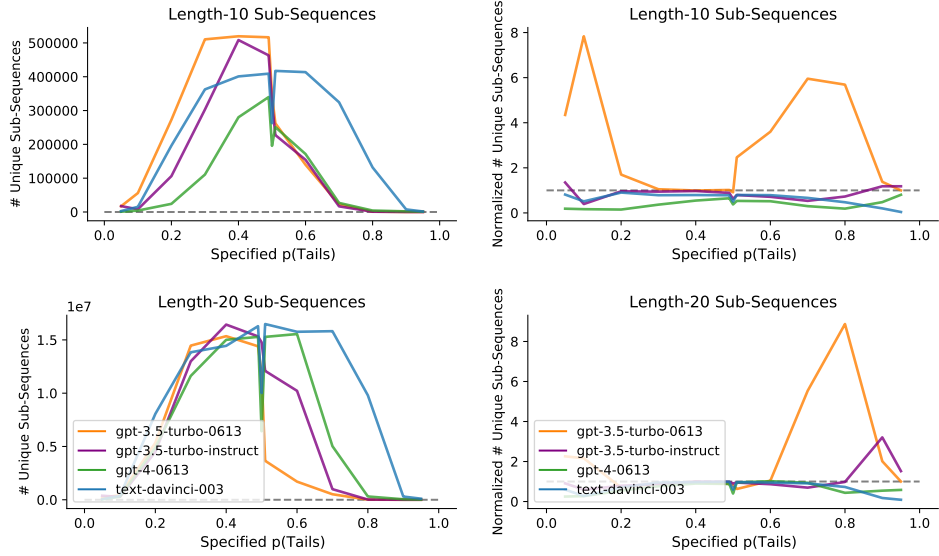


Figure 19: **GPT-4 repeats the same sub-sequences more often than other GPT models** (Left) Number of unique length-10 and length-20 sub-sequences as a function of specified probability  $p(\text{Tails})$ , across all sequences  $y$  (note:  $|y| = 50$ ) generated by each GPT model. (Right) The same distributions, with the y-axis normalized by dividing by the same metric (appended Gzip size, or mean Levenshtein distance) for a Bernoulli distribution centered at  $\bar{y}$ , to control for sequence compression being correlated with probability, e.g. with  $\bar{y} = 0.99$ , the same sub-sequences of only ‘Tails’ flips will appear many times.

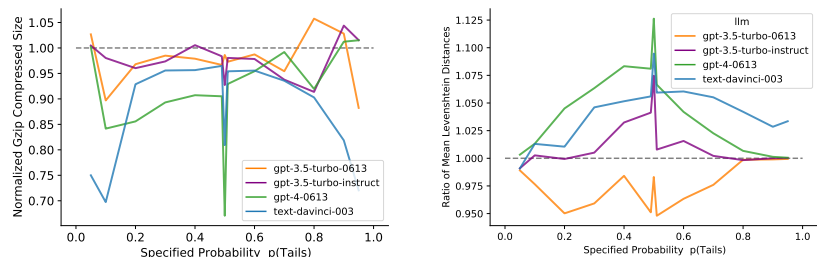


Figure 20: **GPT-generated sequences have lower complexity than Bernoulli sequences**

234 at  $\bar{y}$ . For all GPT models except ChatGPT, generated sequences have smaller Levenshtein distance  
 235 than a Bernoulli process. This is evidence that these LLMs are using memorized sub-sequences  
 236 (‘parrotting’), since sequences have repeated structure. On the other hand, ChatGPT produces more  
 237 dissimilar sequences than chance, suggesting *higher* complexity. In Gzip file size, however, we see a  
 238 lower-complexity bias in all LLMs (except for a few higher values of  $p(\text{Tails})$ ), to varying degrees,  
 239 produce data  $Y$  that is more compressible than data from an equal probability Bernoulli process.

240 **I Background on Algorithmic and Subjective Randomness**

241 We focus on cognitive interpretability of LLMs in the domain of random sequences of binary values.  
 242 Random binary sequences are a minimal domain that have been studied extensively in statistics,  
 243 formal language theory, and algorithmic information theory [24, 72, 73]. We can use this domain  
 244 to systematically test few-shot learning as a function of context length  $|x|$  by testing different input  
 245 sequences  $x$ . We can also test zero-shot learning by having models generate sequences with no context  
 246 ( $|x| = 0$ ), without relying on alternate prompt formats such as chain-of-thought reasoning [21, 22].  
 247 Moreover, language generation trajectories over binary sequences can also be analyzed and visualized  
 248 much more easily than typical user-chatbot interaction trajectories [55, 74], since the token-by-token  
 249 branching factor is only two. Random binary sequences have also been a target domain in cognitive  
 250 science (specifically, *subjective randomness*), where researchers have studied the mechanisms and  
 251 concepts that underlie how people generate random binary sequences or evaluate the randomness of  
 252 given sequences [61, 64, 75].

253 Randomness of a sequence  $x$ , defined in terms of Bayesian model comparison between the class of  
 254 non-random models with the class of random models, can be translated to be the difference between  
 255 the sequence length  $|x|$  and the algorithmic complexity, or *Kolmogorov complexity* of the sequence  
 256  $K(x)$ .

$$\begin{aligned} \text{randomness}(x) &= \log P(x|\text{random}) - \log P(x|\text{non-random}) \\ &= \log 2^{-|x|} - \log 2^{-K(x)} \\ &= K(x) - |x| \end{aligned}$$

257 The likelihood given a truly random Bernoulli process  $p(x|\text{random}) = 2^{-|x|}$  since sequences of equal  
 258 length have equal probability and there are  $2^{|x|}$  binary sequences of length  $|x|$ . This can be thought  
 259 of as a uniform prior over programs, where every program is an exact copy of the output string.

260 The likelihood of  $x$  given the space of non-random processes marginalizes over the posterior of all  
 261 non-random programs (hypotheses)  $\mathcal{H}$ :

$$p(x|\text{non-random}) = \sum_{h \in \mathcal{H}} p(h) p(x|h)$$

262 .

263 A natural prior for programs  $p(h)$  is the description length of that program, where common metrics  
 264 used in software engineering such as *lines of code* or *number of functions* can be seen as practical  
 265 estimations of program description length.

266 If we assume  $p(x|h)$  is a binary likelihood, that is:

$$p(x|h) = \begin{cases} 1 & \text{if } h \text{ generates } x \\ 0 & \text{otherwise} \end{cases}$$

267 and we simplify the problem to finding the maximum a-priori hypothesis  $h$ , and set a prior over  
 268 hypotheses (programs) proportional to their length  $p(h) = 2^{-|x|}$ , this equates to finding the program  
 269 with lowest Kolmogorov complexity  $K(x)$ :

$$P(x|\text{non-random}) \approx \max_h p(h) p(x|h) = 2^{-K(x)}$$

270 where Kolmogorov complexity  $K(x)$  is defined as the description length of the shortest program that  
 271 generates  $x$  as output:

$$K(x) = \operatorname{argmin}_{\{p \in \Sigma^* \mid \text{Evaluate}(p) = x\}} |p|$$

272 The notation  $p \in \Sigma^*$  is analogous to  $h \in \mathcal{H}$ , but refers to a formal alphabet  $\Sigma$  that programs are  
 273 comprised of. In the general case, Kolmogorov complexity  $K(x)$  is uncomputable due to the halting  
 274 problem, since the expression  $\text{Evaluate}(p) = x$  might run forever if  $p$  has an infinite loop.

## 275 References

- 276 [1] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom  
277 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning  
278 and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- 279 [2] Neel Nanda, Lawrence Chan, Tom Liberum, Jess Smith, and Jacob Steinhardt. Progress  
280 measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- 281 [3] Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-  
282 than?: Interpreting mathematical abilities in a pre-trained language model. *arXiv preprint*  
283 *arXiv:2305.00586*, 2023.
- 284 [4] Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two  
285 stories in mechanistic explanation of neural networks. *arXiv preprint arXiv:2306.17844*, 2023.
- 286 [5] Stephanie CY Chan, Adam Santoro, Andrew K Lampinen, Jane X Wang, Aaditya Singh,  
287 Pierre H Richemond, Jay McClelland, and Felix Hill. Data distributional properties drive  
288 emergent few-shot learning in transformers. *arXiv preprint arXiv:2205.05055*, 2022.
- 289 [6] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a  
290 transformer: A memory viewpoint. *arXiv preprint arXiv:2306.00802*, 2023.
- 291 [7] Howard Gardner. *The mind’s new science: A history of the cognitive revolution*. Basic books,  
292 1987.
- 293 [8] Joshua Tenenbaum. Bayesian modeling of human concept learning. *Advances in neural*  
294 *information processing systems*, 11, 1998.
- 295 [9] Noah D Goodman, Joshua B Tenenbaum, Jacob Feldman, and Thomas L Griffiths. A rational  
296 analysis of rule-based concept learning. *Cognitive science*, 32(1):108–154, 2008.
- 297 [10] Tomer D Ullman and Joshua B Tenenbaum. Bayesian models of conceptual development:  
298 Learning as building models of the world. *Annual Review of Developmental Psychology*, 2:  
299 533–558, 2020.
- 300 [11] Tomer D Ullman, Noah D Goodman, and Joshua B Tenenbaum. Theory learning as stochastic  
301 search in the language of thought. *Cognitive Development*, 27(4):455–480, 2012.
- 302 [12] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisti-  
303 cians: Provable in-context learning with in-context algorithm selection. *arXiv preprint*  
304 *arXiv:2306.04637*, 2023.
- 305 [13] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
306 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
307 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 308 [14] Mengjie Zhao, Yi Zhu, Ehsan Shareghi, Ivan Vulić, Roi Reichart, Anna Korhonen, and Hinrich  
309 Schütze. A closer look at few-shot crosslingual transfer: The choice of shots matters. *arXiv*  
310 *preprint arXiv:2012.15682*, 2020.
- 311 [15] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use:  
312 Improving few-shot performance of language models. In *International Conference on Machine*  
313 *Learning*, pages 12697–12706. PMLR, 2021.
- 314 [16] Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and  
315 Luke Zettlemoyer. Rethinking the Role of Demonstrations: What Makes In-Context Learning  
316 Work?, October 2022. URL <http://arxiv.org/abs/2202.12837>. arXiv:2202.12837 [cs].
- 317 [17] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically  
318 ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv*  
319 *preprint arXiv:2104.08786*, 2021.

- 320 [18] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale  
321 Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables  
322 complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.
- 323 [19] Miles Turpin, Julian Michael, Ethan Perez, and Samuel R Bowman. Language models don't  
324 always say what they think: Unfaithful explanations in chain-of-thought prompting. *arXiv*  
325 *preprint arXiv:2305.04388*, 2023.
- 326 [20] Steven T Piantadosi, Joshua B Tenenbaum, and Noah D Goodman. Bootstrapping in a language  
327 of thought: A formal model of numerical concept learning. *Cognition*, 123(2):199–217, 2012.
- 328 [21] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,  
329 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.  
330 *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- 331 [22] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
332 language models are zero-shot reasoners. *Advances in neural information processing systems*,  
333 35:22199–22213, 2022.
- 334 [23] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid,  
335 Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al.  
336 Beyond the imitation game: Quantifying and extrapolating the capabilities of language models.  
337 *arXiv preprint arXiv:2206.04615*, 2022.
- 338 [24] Michael Sipser. Introduction to the theory of computation. *ACM Sigact News*, 27(1):27–29,  
339 1996.
- 340 [25] Noam Chomsky. Three models for the description of language. *IRE Transactions on information*  
341 *theory*, 2(3):113–124, 1956.
- 342 [26] Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot  
343 Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, et al. Neural networks and the  
344 chomsky hierarchy. *arXiv preprint arXiv:2207.02098*, 2022.
- 345 [27] Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International*  
346 *Conference on Machine Learning*, pages 11080–11090. PMLR, 2021.
- 347 [28] Hui Shi, Sicun Gao, Yuandong Tian, Xinyun Chen, and Jishen Zhao. Learning bounded context-  
348 free-grammar via lstm and the transformer: Difference and the explanations. In *Proceedings of*  
349 *the AAAI Conference on Artificial Intelligence*, volume 36, pages 8267–8276, 2022.
- 350 [29] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar.  
351 *arXiv preprint arXiv:2305.13673*, 2023.
- 352 [30] Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of trans-  
353 formers to recognize formal languages. *arXiv preprint arXiv:2009.11264*, 2020.
- 354 [31] Kaiyue Wen, Yuchen Li, Bingbin Liu, and Andrej Risteski. (un) interpretability of transformers:  
355 a case study with dyck grammars. 2023.
- 356 [32] Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Exposing  
357 attention glitches with flip-flop language modeling. *arXiv preprint arXiv:2306.00946*, 2023.
- 358 [33] Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers  
359 learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022.
- 360 [34] William Merrill and Ashish Sabharwal. The parallelism tradeoff: Limitations of log-precision  
361 transformers. *Transactions of the Association for Computational Linguistics*, 11:531–545, 2023.
- 362 [35] William Merrill, Nikolaos Tsilivis, and Aman Shukla. A tale of two circuits: Grokking as  
363 competition of sparse and dense subnetworks. *arXiv preprint arXiv:2303.11873*, 2023.
- 364 [36] William Merrill and Ashish Sabharwal. Transformers implement first-order logic with majority  
365 quantifiers. *arXiv preprint arXiv:2210.02671*, 2022.

- 366 [37] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An Explanation of  
367 In-context Learning as Implicit Bayesian Inference, July 2022. URL [http://arxiv.org/  
368 abs/2111.02080](http://arxiv.org/abs/2111.02080). arXiv:2111.02080 [cs].
- 369 [38] Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers  
370 as Algorithms: Generalization and Stability in In-context Learning, February 2023. URL  
371 <http://arxiv.org/abs/2301.07067>. arXiv:2301.07067 [cs, stat].
- 372 [39] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What  
373 learning algorithm is in-context learning? investigations with linear models. *arXiv preprint*  
374 *arXiv:2211.15661*, 2022.
- 375 [40] Michael Hahn and Navin Goyal. A Theory of Emergent In-Context Learning as Implicit Struc-  
376 ture Induction, March 2023. URL <http://arxiv.org/abs/2303.07971>. arXiv:2303.07971  
377 [cs].
- 378 [41] Gabriel Goh, Nick Cammarata, Chelsea Voss, Shan Carter, Michael Petrov, Ludwig Schubert,  
379 Alec Radford, and Chris Olah. Multimodal neurons in artificial neural networks. *Distill*, 6(3):  
380 e30, 2021.
- 381 [42] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers  
382 are key-value memories. *arXiv preprint arXiv:2012.14913*, 2020.
- 383 [43] Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational*  
384 *Linguistics*, 48(1):207–219, 2022.
- 385 [44] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin  
386 Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic  
387 task. *arXiv preprint arXiv:2210.13382*, 2022.
- 388 [45] Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.  
389 Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv*  
390 *preprint arXiv:2211.00593*, 2022.
- 391 [46] Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse  
392 engineering how networks learn group operations. *arXiv preprint arXiv:2302.03025*, 2023.
- 393 [47] Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris  
394 Bertsimas. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint*  
395 *arXiv:2305.01610*, 2023.
- 396 [48] Alex Foote, Neel Nanda, Esben Kran, Ioannis Konstas, Shay Cohen, and Fazl Barez. Neuron to  
397 graph: Interpreting language model neurons at scale. *arXiv preprint arXiv:2305.19911*, 2023.
- 398 [49] Ekdeep Singh Lubana, Eric J Bigelow, Robert P Dick, David Krueger, and Hidenori Tanaka.  
399 Mechanistic mode connectivity. In *International Conference on Machine Learning*, pages  
400 22965–23004. PMLR, 2023.
- 401 [50] Tom Lieberum, Matthew Rahtz, János Kramár, Geoffrey Irving, Rohin Shah, and Vladimir  
402 Mikulik. Does circuit analysis interpretability scale? evidence from multiple choice capabilities  
403 in chinchilla. *arXiv preprint arXiv:2307.09458*, 2023.
- 404 [51] Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang.  
405 Hidden progress in deep learning: Sgd learns parities near the computational limit. *Advances in*  
406 *Neural Information Processing Systems*, 35:21750–21764, 2022.
- 407 [52] Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams.  
408 Towards understanding grokking: An effective theory of representation learning. *Advances in*  
409 *Neural Information Processing Systems*, 35:34651–34663, 2022.
- 410 [53] Ulrike Hahn and Paul A Warren. Perceptions of randomness: why three heads are better than  
411 four. *Psychological review*, 116(2):454, 2009.
- 412 [54] Ben Prystawski and Noah D Goodman. Why think step-by-step? reasoning emerges from the  
413 locality of experience. *arXiv preprint arXiv:2304.03843*, 2023.

- 414 [55] Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A Smith. How language model  
415 hallucinations can snowball. *arXiv preprint arXiv:2305.13534*, 2023.
- 416 [56] Alex Renda, Aspen Hopkins, and Michael Carbin. Can llms generate random numbers?  
417 evaluating llm sampling in controlled domains llm sampling underperforms expectations, 2023.  
418 <http://people.csail.mit.edu/renda/llm-sampling-paper>.
- 419 [57] janus. Mysteries of mode collapse. LessWrong, 2022. URL [https://www.lesswrong.com/  
420 posts/t9svvNPNmFf5Qa3TA/mysteries-of-mode-collapse](https://www.lesswrong.com/posts/t9svvNPNmFf5Qa3TA/mysteries-of-mode-collapse).
- 421 [58] Andrej Karpathy. A baby GPT with two tokens 0/1 and context length of 3, view[ed] as  
422 a finite state markov chain., 2023. URL [https://twitter.com/karpathy/status/  
423 1645115622517542913](https://twitter.com/karpathy/status/1645115622517542913).
- 424 [59] Pedro A Ortega, Jane X Wang, Mark Rowland, Tim Genewein, Zeb Kurth-Nelson, Razvan  
425 Pascanu, Nicolas Heess, Joel Veness, Alex Pritzel, Pablo Sprechmann, et al. Meta-learning of  
426 sequential strategies. *arXiv preprint arXiv:1905.03030*, 2019.
- 427 [60] Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Goodman. How to grow  
428 a mind: Statistics, structure, and abstraction. *science*, 331(6022):1279–1285, 2011.
- 429 [61] Ruma Falk and Clifford Konold. Making sense of randomness: Implicit encoding as a basis for  
430 judgment. 1997.
- 431 [62] Thomas L Griffiths and Joshua B Tenenbaum. From Algorithmic to Subjective Randomness.  
432 *Neural Information Processing Systems*, 2003.
- 433 [63] Thomas L Griffiths and Joshua B Tenenbaum. Probability, algorithmic complexity, and subjective  
434 randomness. *Proceedings of the Cognitive Science Society*, 2004.
- 435 [64] Thomas L. Griffiths, Dylan Daniels, Joseph L. Austerweil, and Joshua B. Tenenbaum. Subjective  
436 randomness as statistical inference. *Cognitive Psychology*, 103:85–109, June 2018. ISSN  
437 00100285. doi: 10.1016/j.cogpsych.2018.02.003. URL [https://linkinghub.elsevier.  
438 com/retrieve/pii/S0010028517302281](https://linkinghub.elsevier.com/retrieve/pii/S0010028517302281).
- 439 [65] Yuan Yang and Steven T Piantadosi. One model for the learning of language. *Proceedings of  
440 the National Academy of Sciences*, 119(5):e2021865119, 2022.
- 441 [66] Eric Bigelow and Steven T Piantadosi. Inferring priors in compositional cognitive models. In  
442 *CogSci*, 2016.
- 443 [67] Micah Goldblum, Marc Finzi, Keefer Rowan, and Andrew Gordon Wilson. The no free lunch  
444 theorem, kolmogorov complexity, and the role of inductive biases in machine learning. *arXiv  
445 preprint arXiv:2304.05366*, 2023.
- 446 [68] Grégoire Delétang, Anian Ruoss, Paul-Ambroise Duquenne, Elliot Catt, Tim Genewein, Christo-  
447 pher Mattern, Jordi Grau-Moya, Li Kevin Wenliang, Matthew Aitchison, Laurent Orseau, et al.  
448 Language modeling is compression. *arXiv preprint arXiv:2309.10668*, 2023.
- 449 [69] Peter Deutsch. Gzip file format specification version 4.3. Technical report, 1996.
- 450 [70] Zhiying Jiang, Matthew YR Yang, Mikhail Tsirlin, Raphael Tang, and Jimmy Lin. Less is more:  
451 Parameter-free text classification with gzip. *arXiv preprint arXiv:2212.09410*, 2022.
- 452 [71] Vladimir I Levenshtein et al. Binary codes capable of correcting deletions, insertions, and  
453 reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union, 1966.
- 454 [72] Gregory J Chaitin. On the length of programs for computing finite binary sequences. *Journal of  
455 the ACM (JACM)*, 13(4):547–569, 1966.
- 456 [73] Ming Li, Paul Vitányi, et al. *An introduction to Kolmogorov complexity and its applications*.  
457 Springer, 1997.



- 458 [74] Lukas Berglund, Asa Cooper Stickland, Mikita Balesni, Max Kaufmann, Meg Tong, Tomasz  
459 Korbak, Daniel Kokotajlo, and Owain Evans. Taken out of context: On measuring situational  
460 awareness in llms. *arXiv preprint arXiv:2309.00667*, 2023.
- 461 [75] Amos Tversky and Thomas Gilovich. The “hot hand”: Statistical reality or cognitive illusion?  
462 *Chance*, 2(4):31–34, 1989.