
Unveiling Linguistic and Mathematical Knowledge: Interpreting Grammar and Arithmetic Embeddings in Small-Scale LLMs

Pratim Chowdhary*
Department of Computer Science
Dartmouth College
cpratim.25@dartmouth.edu

Peter Chin
Department of Engineering
Thayer School of Engineering
pc@dartmouth.edu

Deepernab Chakrabarty
Department of Computer Science
Dartmouth College
deepernab@dartmouth.edu

Abstract

Small transformer language models (≤ 10 B parameters) already solve a surprising range of grammatical and numerical tasks. But *which* internal components drive each capability—and how much circuitry is reused across domains—remains unclear. We study three task families—synthetic arithmetic verification, arithmetic word-problems, and grammatical acceptability—and trace responsibility down to the level of individual attention heads. Using a causal ablation-and-pruning procedure that extracts a *Minimum-Sufficient Head Circuit* (MSHC) for each task, we show that only 10-20 heads (0.4 % of parameters) are needed to recover 90 % of full-model accuracy in the models we analyse (Gemma-9B, Llama-8B, Qwen-8B). The MSHCs for arithmetic verification and word-problems overlap by 40-60 %, revealing a reusable numerical sub-network, whereas grammar circuits are largely disjoint. These findings suggest that small LLMs learn a dedicated "number sense" circuit that generalises from bare arithmetic to text-framed problems, while syntactic competence is carried by a separate set of heads. Our results offer new levers for parameter-efficient fine-tuning and mechanistic interpretability of compact language models.

1 Introduction

Large language models (LLMs) have transformed natural language processing and mathematical reasoning, demonstrating unprecedented capabilities across a spectrum of tasks—from simple question answering to complex linguistic analysis, arithmetic problem solving, and formal reasoning [Brown et al., 2020, Chowdhery et al., 2022, Touvron et al., 2023b, Jiang et al., 2023]. These models, trained through self-supervised learning on vast corpora of text and mathematical data, have increasingly approached human-like performance in generating both grammatically well-formed text and mathematically coherent solutions, despite having no explicit grammatical rules or arithmetic algorithms programmed into their architecture. This emergent dual competence—arising purely from statistical patterns in training data—represents a fascinating case study in how both linguistic and mathematical knowledge can be acquired implicitly through exposure rather than explicit instruction.

*Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

In this work we move beyond aggregate performance metrics and zoom in on the level of individual attention heads. Inspired by circuit-based interpretability analyses [Olah et al., 2020, Elhage et al., 2021], we introduce the *Minimum Sufficient Head Circuit* (MSHC)—the smallest set of heads whose activations suffice to solve a task within a user-specified tolerance. Extracting MSHCs for arithmetic verification, mathematical word-problems, and grammatical acceptability across three open-weight models (Gemma-9B, Llama-8B, Qwen-8B), we uncover a substantial overlap between the arithmetic and word-problem circuits and only minimal intersection with the grammar circuit. These results indicate that small LLMs reuse a shared numerical sub-circuit across superficially different tasks while maintaining a distinct pathway for syntactic reasoning. We quantify this overlap, analyse how it scales with model size, and discuss implications for parameter-efficient fine-tuning and model editing.

The field has witnessed exponential growth in model size, from early transformer models with hundreds of millions of parameters [Vaswani et al., 2017] to modern giants like GPT-4 [OpenAI, 2023] that likely contain trillions of parameters. While these massive models have captured headlines with their impressive capabilities in both language and mathematics, a parallel revolution has been unfolding in the development of smaller, more efficient models in the 1-8 billion parameter range. Models like Llama 2 [Touvron et al., 2023b], Gemma 7B [Jiang et al., 2023], and Qwen [Bai et al., 2023] have demonstrated remarkable performance in both linguistic and numerical tasks despite their relatively modest size, making them particularly valuable for practical applications where computational efficiency and deployment costs are significant concerns.

These smaller LLMs offer a compelling balance between capability and efficiency, enabling deployment on consumer hardware, edge devices, and resource-constrained environments. Their reduced inference costs make them attractive for commercial applications, while their smaller memory footprint allows for fine-tuning and adaptation with more modest computational resources. Understanding how these models encode and represent both grammatical and numerical knowledge is therefore not merely an academic exercise but has significant practical implications for developing more capable, efficient, and cognitively robust language technologies that can handle both linguistic and mathematical reasoning.

2 Related Work and Background

2.1 Linguistic and Mathematical Evaluation of Language Models

Research on evaluating neural language models has evolved from early work on LSTM-based architectures [Linzen et al., 2016] to comprehensive assessment of transformer-based models across both linguistic and mathematical domains [Goldberg, 2019, Devlin et al., 2019, Saxton et al., 2019]. For grammatical evaluation, frameworks have progressed from simple agreement tests to comprehensive benchmarks like CoLA [Warstadt et al., 2019] and BLiMP [Warstadt et al., 2020]. For mathematical reasoning, benchmarks include GSM8K [Cobbe et al., 2021], SVAMP [?], and MathQA [?], which test arithmetic computation, word problem understanding, and quantitative reasoning. Recent studies have revealed that while performance in both domains scales with size, challenges remain with complex hierarchical structures in grammar and multi-step reasoning in mathematics [Thrush et al., 2022, Qian et al., 2022, ?].

2.2 Scale and Cognitive Competence

The relationship between model scale and cognitive abilities follows complex patterns beyond the general power-law scaling observed in language modeling [Kaplan et al., 2020]. Research suggests that rare grammatical constructions and complex arithmetic operations require disproportionately more training data [Wei et al., 2021, ?]. Certain phenomena in both domains show nonlinear improvements at specific parameter thresholds [Zhang et al., 2023], with mathematical reasoning often exhibiting steeper scaling curves than linguistic tasks [?]. Studies on compositional generalization indicate that scaling alone may not capture human-like cognitive productivity without architectural innovations that support both symbolic and statistical processing [Hu and Daumé III, 2020, ?].

77 2.3 Probing Language Models for Cognitive Knowledge

78 Researchers have developed various probing techniques to understand how linguistic and mathe-
79 matical knowledge is represented within model parameters. Structural probes have revealed that
80 models implicitly encode both syntactic hierarchies and numerical relationships [Hewitt and Manning,
81 2019, ?], with different types of knowledge appearing at different network depths [Tenney et al.,
82 2019, ?]. Studies show that transformer-based models exhibit emergent capabilities resembling both
83 discrete linguistic rules and arithmetic algorithms [Manning et al., 2020, ?], with individual neurons
84 specializing in specific linguistic features or numerical operations [Geva et al., 2021, Patel and
85 Pavlick, 2022]. Attention patterns often correspond to both syntactic dependencies and mathematical
86 relations, with different attention heads specializing in distinct cognitive phenomena [Clark et al.,
87 2019, ?].

88 2.4 Small LLMs and Comparative Studies

89 The proliferation of open-weight models like OPT [Zhang et al., 2022], Llama [Touvron et al., 2023a],
90 Gemma [Gemma Team, 2024], and Qwen [Bai et al., 2023] has enabled more systematic analyses of
91 how capabilities in both linguistic and mathematical domains scale and how architectural choices
92 affect performance. Studies show that smaller models with innovative architectures can outperform
93 larger ones in specific aspects of language and mathematics [?], and data quality may be as important
94 as scale for robust cognitive representations [?]. Comparative analyses across architectures remain
95 limited but suggest significant variations in performance even at similar parameter scales, with
96 some models showing domain-specific strengths [Talmor et al., 2020, Zhao et al., 2023]. Recent
97 frameworks for quantifying evaluation uncertainties [Roberts et al., 2023] and structured evaluation
98 approaches [Xia et al., 2023] have revealed that architectural design choices impact both grammatical
99 and mathematical phenomena differently, suggesting that model architectures encode cognitive
100 knowledge in fundamentally different ways.

101 3 Methodology

102 Our analysis proceeds in three strands:

- 103 **(1) Low-dimensional linear separability (LS).** We measure how well a D -dimensional projection
104 of the hidden state at each layer separates correct from incorrect examples ($D \leq 3$).
- 105 **(2) Minimum Sufficient Head Circuit (MSHC).** Guided by the LS curves we isolate the smallest
106 set of attention heads whose activations suffice for the task.
- 107 **(3) Controlled datasets.** All tasks are cast as *minimal pairs* so that a single factual change flips the
108 label, letting us attribute separability (or its absence) to that fact alone.

109 Throughout, let \mathcal{V} denote the vocabulary and $m : \mathcal{V}^* \rightarrow \mathbb{R}^d$ a causal transformer with L layers and
110 H heads per layer. For a sequence $x = (x_1, \dots, x_n)$ we write $\mathbf{h}_{x,\ell} \in \mathbb{R}^d$ for the activation at the
111 end-of-sequence (EOS) token at layer ℓ , and $\mathbf{a}_{x,\ell,h} \in \mathbb{R}^d$ for the contribution of head $h \in \{1, \dots, H\}$
112 in that layer.

113 **Activation collection.** For every item we construct the full prompt by concatenating a single
114 one-shot demonstration, the query sentence or equation, and the end-of-sequence marker $\langle /s \rangle$
115 (or its model-specific analogue). The model is executed in teacher-forcing mode. At each layer
116 $\ell \in \{0, \dots, L\}$ we record the hidden state $\mathbf{h}_{x,\ell}$ at the position of that final EOS token—thus collecting
117 a compact $L+1$ -vector trace that summarises the entire computation leading to the model’s output
118 logit. We omit intermediate attention projections and all key/value tensors to minimise I/O overhead;
119 subsequent probes operate solely on this last-token trace. Because all prompts fit within the context
120 window, no padding is required, and teacher forcing guarantees determinism across repeated runs.

121 3.1 Task families and evaluation sets

122 We consider three evaluation corpora, each built from *minimal pairs* (x_c, x_i) in which the two
123 members differ in exactly one fact that determines correctness. Every item appears in three prompt

124 variants: the raw text, a two-choice question (“A” or “B”), and a single-candidate acceptability probe,
125 each preceded by a one-shot demonstration so that prompting remains uniform across models.

126 **Grammar (G).** The 67k sentence pairs from BLiMP [Warstadt et al., 2020] cover twelve syntactic
127 phenomena.

Good Sentence: *Who should Derek hug after shocking Richard?*
Bad Sentence: *Who should Derek hug Richard after shocking?*

128 **Arithmetic verification (A).** We generate 10^3 addition pairs with addends $n_1, n_2 \in [1, 10^3]$;
129 incorrect results satisfy the deviation constraint in Eq. (1).

Good Equation: $1338 + 88 = 1426$
Bad Equation: $1338 + 88 = 2139$

$$.5 * (n1 \pm n2) \leq \text{noisy}(n1 \pm n2) \leq 1.5 * (n1 \pm n2) \quad (1)$$

130 **Word-problem arithmetic (W).** The 100 story problems are produced by numerically perturbing
131 the ending of a template narrative:

Good Expression: Tim has 5 apples and eats 2, leaving him with 3 apples.
Bad Expression: Tim has 5 apples and eats 2, leaving him with 10 apples.

132 where the perturbation is applied by taking a sentence generated with a template and replacing the
133 numbers with perturbed numbers from the 10000-item dataset of addition/subtraction pairs generated
134 as in (1).

135 3.2 Low-dimensional linear separability metric

136 A well-known pathology of very high-dimensional spaces is that almost *any* two finite clouds are
137 linearly separable with overwhelming probability [see ?]. Raw accuracy of a linear probe in the full
138 residual space \mathbb{R}^d therefore over-states how “easy” a task is. To obtain a more conservative—and
139 hence more informative—measure of representational structure, we evaluate separability after first
140 collapsing the hidden states onto just *two* principal directions. We call the resulting statistic the
141 *low-dimensional linear separability* score, $\text{LS}_{t,\ell}$.

142 **Notation.** Fix a task t and a transformer layer ℓ . For every prompt x we write $\mathbf{h}_{x,\ell} \in \mathbb{R}^d$ for the
143 centred hidden state at the EOS position: $\mathbf{h}_{x,\ell} := \mathbf{h}_{x,\ell}^{(n)} - \bar{\mathbf{h}}_\ell$, where $\bar{\mathbf{h}}_\ell = \frac{1}{|\mathcal{D}_t|} \sum_{x \in \mathcal{D}_t} \mathbf{h}_{x,\ell}^{(n)}$ is the
144 layer mean computed over the complete training split \mathcal{D}_t .

145 **Step A: variance-maximising projection.** Let $\Sigma_\ell = \frac{1}{|\mathcal{D}_t|} \sum_{x \in \mathcal{D}_t} \mathbf{h}_{x,\ell} \mathbf{h}_{x,\ell}^\top$ be the empirical covari-
146 ance matrix. We seek an orthonormal matrix $\mathbf{W}_\ell \in \mathbb{R}^{d \times 2}$ that captures the greatest possible variance
147 under a rank-2 constraint:

$$\max_{\mathbf{W}^\top \mathbf{W} = \mathbf{I}_2} \text{Tr}(\mathbf{W}^\top \Sigma_\ell \mathbf{W}). \quad (2)$$

148 The optimal columns are the two leading eigenvectors of Σ_ℓ . The corresponding 2-D coordinates are
149 $\tilde{\mathbf{h}}_{x,\ell} = \mathbf{W}_\ell^\top \mathbf{h}_{x,\ell} \in \mathbb{R}^2$.

150 **Step B: linear decision boundary.** Assign labels $y_i = +1$ for *correct* members of a minimal pair
151 and $y_i = -1$ for *incorrect* ones. We train a soft-margin support-vector machine in the projected space:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^{N_t} \max(0, 1 - y_i (\mathbf{w}^\top \tilde{\mathbf{h}}_{x_i,\ell} + b)), \quad C = 10. \quad (3)$$

152 Because $\tilde{\mathbf{h}}_{x,\ell} \in \mathbb{R}^2$, the resulting classifier depends on at most three parameters, precluding the kind
153 of pathological over-fitting that haunts full-dimensional probes.

154 **Definition of the metric.** Let $\hat{y}_i = \text{sign}(\mathbf{w}^\top \tilde{\mathbf{h}}_{x_i, \ell} + b)$ denote the SVM’s prediction on a held-out
 155 validation example. We define

$$\text{LS}_{t, \ell} = \frac{1}{N_t^{\text{val}}} \sum_{i=1}^{N_t^{\text{val}}} \mathbf{1}[\hat{y}_i = y_i], \quad \text{LS}_{t, \ell} \in [0, 1]. \quad (4)$$

156 **Interpretation.** Random guessing gives $\text{LS}_{t, \ell} = 0.5$. A score close to 1 implies that the first two
 157 principal axes already support a linear decision boundary, i.e. the class-conditional embeddings differ
 158 in a *low-codimension* direction. Conversely, scores near 0.5 indicate either geometric entanglement
 159 or dispersion of the signal across many dimensions.

160 In what follows we plot $\ell \mapsto \text{LS}_{t, \ell}$ for each task. Peaks in these curves spotlight layers where the
 161 model’s hidden states make the good/bad distinction in the simplest possible way: by shifting along
 162 just two orthogonal directions.

163 3.3 The Minimum Sufficient Head Circuit

164 A transformer’s computation is frequently dominated by a surprisingly small subset of its attention
 165 heads. We formalise this with the *Minimum Sufficient Head Circuit* (MSHC). Fix an accuracy
 166 tolerance $\epsilon \in (0, 0.5)$. Let $\mathcal{H} = \{(\ell, h) : 1 \leq \ell \leq L, 1 \leq h \leq H\}$ be the set of all heads. A subset
 167 $\mathcal{C} \subseteq \mathcal{H}$ is an **MSHC $_\epsilon$** if, *with high probability (WHP)*, enabling any single head from \mathcal{C} is already
 168 sufficient to lift accuracy above chance:

$$(\forall h \in \mathcal{C}) \quad \Pr_{x \sim \mathcal{D}} [\text{Acc}(m; \{h\}) > 0.5 + \epsilon] \geq 1 - \delta,$$

169 where δ is a user-specified failure probability (we use $\delta = 0.05$ in all experiments), and \mathcal{C} is
 170 inclusion-minimal with this property.

171 **Hunting the circuit** Our search still employs a sliding window but now ranks layers by the *size* of the
 172 accuracy loss they incur. We slide a window of width $w = \lfloor xL \rfloor$ layers from the bottom to the top of
 173 the network, disable *all* heads in that window, and measure the resulting accuracy Acc_{off} . For every
 174 layer ℓ covered by the current window we store $\text{ACCLAYER}[\ell] = \max(\text{ACCLAYER}[\ell], \text{Acc}_{\text{off}})$ —the
 175 *best* accuracy ever observed when *any* window that contains ℓ is switched off. After scanning all
 176 windows we compute a per-layer drop score $\Delta_\ell = \text{Acc}_{\text{full}} - \text{ACCLAYER}[\ell]$ and keep those layers
 177 whose Δ_ℓ lies in the top 25th percentile. These high-impact layers seed the head-level search:

178 **Step 1.** Disable every head *in the selected top-quartile layers*, establishing a “dark-start” baseline.

179 **Step 2.** Run a stochastic pruning loop (Alg. 1) that begins with bundle size $k_0 = \lceil 2\sqrt{|\mathcal{C}|} \rceil$. For
 180 each bundle size k , we draw N random bundles of k heads, measure their accuracies, and
 181 remove the bundle that attains the *lowest* accuracy whenever that value is $\leq 0.5 + \epsilon$. When
 182 the worst of the N bundles exceeds the threshold, we tighten the constraint by updating
 183 $k \leftarrow \max(1, \lfloor k/2 \rfloor)$ and continue.

184 **Step 3.** Stop when $k = 1$ and every random bundle succeeds—by definition, the remaining heads
 185 form an MSHC.

186 The procedure terminates because the candidate set shrinks monotonically and can be pruned at most
 187 $|\mathcal{H}|$ times.

188 The narrative interpretation is simple: we first find *where* knowledge lives, then keep only those
 189 filaments that reliably relight the lamp.

190 Empirically, the circuit coalesces in fewer than 300 iterations on all three models studied—Gemma-
 191 9B, Llama-8B, and Qwen-8B.

192 3.4 Theoretical analysis of MSHC

193 We aim to show that the MSHC is a good approximation of the minimum number of heads that are
 194 needed to solve the task.

Algorithm 1 Sliding-window percentile localisation followed by stochastic discovery of an MSHC_ϵ

Require: Window fraction x , tolerance ϵ , sample count N

```
1:  $\text{Acc}_{\text{full}} \leftarrow \text{Acc}(m)$ 
2:  $w \leftarrow \lfloor xL \rfloor$ 
3: initialise array  $\text{ACCLAYER}[1:L] \leftarrow 0$   $\triangleright$  best accuracy seen with each layer disabled
4: for  $s = 1$  to  $L - w + 1$  do
5:   disable all heads in layers  $s$  to  $s + w - 1$ 
6:    $\text{Acc}_{\text{off}} \leftarrow \text{Acc}(m)$ 
7:   for  $\ell = s$  to  $s + w - 1$  do
8:      $\text{ACCLAYER}[\ell] \leftarrow \max(\text{ACCLAYER}[\ell], \text{Acc}_{\text{off}})$ 
9:   re-enable the disabled layers
10: for  $\ell = 1$  to  $L$  do
11:    $\text{DROP}[\ell] \leftarrow \text{Acc}_{\text{full}} - \text{ACCLAYER}[\ell]$   $\triangleright$  compute per-layer drop scores
12:  $\tau \leftarrow 75\text{th percentile of DROP}$ 
13:  $\mathcal{L} \leftarrow \{\ell \mid \text{DROP}[\ell] \geq \tau\}$   $\triangleright$  top-quartile layers
14:  $\mathcal{C} \leftarrow$  all heads in  $\mathcal{L}$   $\triangleright$  initial candidate circuit
15:  $k \leftarrow \lfloor |\mathcal{C}|/2 \rfloor$ 
16: while  $k \geq 1$  do
17:   repeat
18:      $\text{minAcc} \leftarrow 1; \mathcal{K}_{\text{min}} \leftarrow \emptyset$ 
19:     for  $i = 1$  to  $N$  do
20:       draw  $\mathcal{K} \sim \text{Unif}\{\mathcal{S} \subseteq \mathcal{C} : |\mathcal{S}| = k\}$ 
21:        $\text{acc} \leftarrow \text{Acc}(m; \mathcal{K})$ 
22:       if  $\text{acc} < \text{minAcc}$  then
23:          $\text{minAcc} \leftarrow \text{acc}; \mathcal{K}_{\text{min}} \leftarrow \mathcal{K}$ 
24:       if  $\text{minAcc} \leq 0.5 + \epsilon$  then
25:          $\mathcal{C} \leftarrow \mathcal{C} \setminus \mathcal{K}_{\text{min}}$   $\triangleright$  prune the worst bundle
26:     until  $\text{minAcc} > 0.5 + \epsilon$ 
27:      $k \leftarrow \max(1, \lfloor k/2 \rfloor)$ 
28: return  $\mathcal{C}$ 
```

195 3.5 Experimental protocol

196 Our analysis focuses on three open-weight checkpoints that occupy a comparable size class yet
197 differ architecturally: *Gemma-9B*, *Llama-8B* and *Qwen-8B*. Each model is evaluated on exactly the
198 same train/validation/ test splits (60 : 20 : 20) of the G, A, and W corpora described above. For
199 a given model we first trace the accuracy-vs-layer curves (§??), then feed the validation set to the
200 sliding-window search with $(\epsilon, k) = (0.1, 3)$ and a window fraction $x=0.2$. The resulting MSHC is
201 frozen before we inspect test performance, and uncertainty bands are computed via 1 000-sample
202 bootstrap over minimal pairs. All experiments run on a single A100 80 GB GPU; wall-clock times
203 are listed in Appendix C.

204 4 Experiments

205 5 Discussion

206 6 Conclusion

207 References

- 208 Jinze Bai, Shuai Lv, Sheng Peng, Yida Wang, Xingjian Zhang, Ziyue Yang, Beilei Yang, Haotian
209 Gong, Zhiyu Fu, Kongming Liu, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*,
210 2023.
- 211 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
212 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
213 few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- 214 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
215 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
216 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 217 Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does bert look at?
218 an analysis of bert’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing*
219 *and Interpreting Neural Networks for NLP*, pages 276–286, 2019.
- 220 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, et al. Training verifiers to solve math word
221 problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 222 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
223 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2019.
- 224 Nelson Elhage, Neel Nanda, Catherine Olsson, et al. A mathematical framework for transformer
225 circuits. *Transformer Circuits Thread*, 2021. [https://transformer-circuits.pub/2021/](https://transformer-circuits.pub/2021/framework/index.html)
226 [framework/index.html](https://transformer-circuits.pub/2021/framework/index.html).
- 227 Gemma Team. Gemma: Lightweight open models for language understanding. *arXiv preprint*
228 *arXiv:2402.19155*, 2024.
- 229 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are
230 key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural*
231 *Language Processing*, pages 5484–5495, 2021.
- 232 Yoav Goldberg. Assessing bert’s syntactic abilities. In *arXiv preprint arXiv:1901.05287*, 2019.
- 233 John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations.
234 In *Proceedings of the 2019 Conference of the North American Chapter of the Association for*
235 *Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*,
236 pages 4129–4138, 2019.
- 237 Weihs Hu and Hal Daumé III. Systematic evaluation of causal discovery in visual model based
238 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:12578–12590,
239 2020.
- 240 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
241 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
242 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 243 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
244 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
245 *arXiv preprint arXiv:2001.08361*, 2020.
- 246 Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. Assessing the ability of lstms to learn syntax-
247 sensitive dependencies. *Transactions of the Association for Computational Linguistics*, 4:521–535,
248 2016.

249 Christopher D Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. Emergent
250 linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the*
251 *National Academy of Sciences*, 117(48):30046–30054, 2020.

252 Chris Olah, Nick Cammarata, Ludwig Schubert, et al. Zoom in: An introduction to circuits. *Distill*, 5
253 (3):e00024, 2020. doi: 10.23915/distill.00024.

254 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

255 Krishna Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. *arXiv*
256 *preprint arXiv:2210.02539*, 2022.

257 Peng Qian, Tahsina Huang, Reza Firoozi, Zhiyuan Wang, Qiyang Zhou, Eric Wong, Kevin Chen,
258 Shaopeng Pan, Zhou Yu, Yang Xiang, et al. Limitations of language models in arithmetic and
259 symbolic reasoning. *arXiv preprint arXiv:2208.05051*, 2022.

260 Adam Roberts, Albert Webson, Colin Larson, Leo Gao, Niket Tandon, Kai-Wei Tai, Hyung Won
261 Chung, Colin Raffel, and Gaurav Mishra. Quantifying language models’ sensitivity to spurious
262 features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv*
263 *preprint arXiv:2310.11324*, 2023.

264 David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical
265 reasoning abilities of neural models. *arXiv preprint arXiv:1904.01557*, 2019.

266 Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. olmpics-on what language model
267 pre-training captures. In *Transactions of the Association for Computational Linguistics*, volume 8,
268 pages 743–758. MIT Press, 2020.

269 Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline. In *Pro-*
270 *ceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages
271 4593–4601, 2019.

272 Tristan Thrush, Sanjay Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and
273 Candace Ross. Winoground: Probing vision and language models for visio-linguistic composition-
274 ality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
275 pages 5238–5248, 2022.

276 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
277 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
278 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

279 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
280 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
281 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

282 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
283 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing*
284 *Systems*, 30, 2017.

285 Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments.
286 In *Transactions of the Association for Computational Linguistics*, volume 7, pages 625–641. MIT
287 Press, 2019.

288 Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and
289 Samuel R Bowman. Blimp: The benchmark of linguistic minimal pairs for english. *Transactions*
290 *of the Association for Computational Linguistics*, 8:377–392, 2020.

291 Jerry Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
292 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Frequency effects on syntactic rule learning
293 in transformers. *arXiv preprint arXiv:2109.07020*, 2021.

294 Jiacheng Xia, Songbo Li, Haozhao Xu, Danny Chen, Yang Liu, Bill Cohen, and Leyang Zhang. Struc-
295 tured prompting: Scaling in-context learning to 1,000 examples. *arXiv preprint arXiv:2212.06713*,
296 2023.

297 Hugh Zhang, Amy Webb, Saujas Petryk, Yiheng Han, Jason Lei, and Chelsea Finn. Language
298 modeling with reduced spurious correlations. *arXiv preprint arXiv:2306.01708*, 2023.

299 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
300 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language
301 models. *arXiv preprint arXiv:2205.01068*, 2022.

302 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
303 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*
304 *preprint arXiv:2303.18223*, 2023.

305 **A Appendix / supplemental material**

306 Optionally include supplemental material (complete proofs, additional experiments and plots) in
307 appendix.