
000 DISSECTING IN-CONTEXT LEARNING: A MECHA-
001 NISTIC ANALYSIS OF
002 EMERGENT CIRCUITS IN SMALL LANGUAGE MOD-
003 ELS
004
005
006
007

008 **Anonymous authors**

009 Paper under double-blind review
010
011

012 ABSTRACT
013

014 In-context learning (ICL) enables language models to adapt to new tasks
015 from just a few examples, yet the mechanistic basis of this capability re-
016 mains poorly understood. We present a comprehensive analysis of the cir-
017 cuits underlying ICL in transformer models ranging from 125M to 1.3B pa-
018 rameters. Through systematic interventions and causal analysis, we identify
019 four distinct circuit types that emerge during training: *copy circuits* that
020 replicate patterns, *induction circuits* that abstract rules, *composition cir-*
021 *cuits* that combine information, and *task recognition circuits* that identify
022 problem types. We demonstrate that these circuits are (1) causally respon-
023 sible for ICL performance through targeted ablations showing 73% average
024 performance degradation, (2) transferable across model scales with 0.82 cor-
025 relation in circuit structure, and (3) surgically enhanceable, achieving 28%
026 improvement on targeted tasks. Our analysis reveals that ICL emerges
027 through the coordinated interaction of 12–15 critical attention heads form-
028 ing interpretable computational graphs. We provide an open-source toolkit
029 for ICL circuit analysis and demonstrate applications to model debugging
030 and capability enhancement. These findings offer actionable insights for
031 improving model interpretability and engineering more capable systems.
032

033 1 INTRODUCTION
034

035 Large language models (LLMs) demonstrate a remarkable ability to learn new tasks from
036 just a few examples provided in their input context, a phenomenon known as in-context
037 learning (Brown et al., 2020). This capability, which emerges without any parameter up-
038 dates, fundamentally challenges our understanding of how neural networks process and
039 generalize from information. Despite its practical importance and theoretical implications,
040 the mechanistic basis of ICL—the specific computational circuits that enable it—remains
041 largely opaque.

042 Consider a model shown examples like “cat → animal, dog → animal, rose → plant” and
043 then correctly inferring “daisy → plant”. This requires multiple sophisticated operations:
044 pattern recognition, abstraction of the mapping rule, and application to novel inputs. How
045 do transformer architectures, trained only to predict next tokens, develop such structured
046 reasoning capabilities? Previous work has identified individual components like induction
047 heads (Olsson et al., 2022) but lacks a comprehensive understanding of how these compo-
048 nents interact to enable ICL.

049 In this work, we present a systematic mechanistic analysis of ICL in transformer language
050 models. Our approach combines three key methodological innovations:

- 051
052 1. **Systematic Circuit Discovery:** We develop automated methods to identify and
053 categorize attention head patterns associated with ICL across 47 different task types,
revealing consistent circuit motifs.

-
- 054 2. **Causal Validation:** Through targeted interventions including activation patching,
055 attention knockout, and circuit transplantation, we establish causal relationships
056 between identified circuits and ICL performance.
057
058 3. **Circuit Manipulation:** We demonstrate that understanding these mechanisms
059 enables practical applications, including debugging failure modes and enhancing
060 model capabilities on specific tasks.

061 Our analysis reveals that ICL is implemented through the coordinated action of four distinct
062 circuit types, each serving a specific computational role. These circuits emerge consistently
063 across model scales and task domains, suggesting fundamental architectural principles. We
064 find that only 12–15 attention heads (less than 8% of total heads in a 12-layer model) are
065 critical for ICL, with their removal causing catastrophic performance degradation.

067 1.1 CONTRIBUTIONS

068 Our main contributions are:

- 070 • **Circuit Taxonomy:** We identify and characterize four fundamental circuit types
071 underlying ICL: copy circuits ($\mathcal{C}_{\text{copy}}$), induction circuits (\mathcal{C}_{ind}), composition circuits
072 ($\mathcal{C}_{\text{comp}}$), and task recognition circuits ($\mathcal{C}_{\text{task}}$), each with distinct computational sig-
073 natures and functional roles.
- 074 • **Causal Evidence:** We establish causality through comprehensive ablation stud-
075 ies (N=2,350 interventions) showing that removing identified circuits reduces ICL
076 performance by 73% on average, while removing random heads of equal number
077 reduces performance by only 11%.
- 078 • **Cross-Scale Consistency:** We demonstrate that ICL circuit structure is remark-
079 ably consistent across model scales (Spearman $\rho = 0.82$), with larger models show-
080 ing more distributed but functionally similar implementations.
- 081 • **Practical Applications:** We provide an open-source toolkit for ICL circuit anal-
082 ysis and demonstrate two applications: (1) debugging systematic failure modes in
083 mathematical reasoning, and (2) enhancing performance on specific task categories
084 by 28% through targeted circuit amplification.

086 2 RELATED WORK

088 2.1 MECHANISTIC INTERPRETABILITY

089 The field of mechanistic interpretability seeks to understand neural networks by reverse-
090 engineering their learned algorithms (Elhage et al., 2021; Cammarata et al., 2020). Recent
091 work has identified specific circuits for various capabilities including indirect object identi-
092 fication (Wang et al., 2022), factual recall (Meng et al., 2022), and grammatical agreement
093 (Finlayson et al., 2021). Our work extends this paradigm to the more complex phenomenon
094 of in-context learning.

096 2.2 IN-CONTEXT LEARNING

097 Since its discovery in GPT-3 (Brown et al., 2020), ICL has been extensively studied from
098 empirical (Min et al., 2022; Wei et al., 2023) and theoretical (Xie et al., 2022; Akyürek et al.,
099 2022) perspectives. Olsson et al. (2022) identified induction heads as key components, while
100 Garg et al. (2022) demonstrated that transformers can implement regression algorithms in-
101 context. However, these works focus on individual mechanisms rather than the complete
102 computational graph underlying ICL.

104 2.3 CIRCUIT DISCOVERY METHODS

105 Recent advances in automated circuit discovery (Conmy et al., 2023; Goldowsky-Dill et al.,
106 2023) enable systematic analysis of model internals. We build on activation patching (Vig
107

et al., 2020) and path patching (Goldowsky-Dill et al., 2023) techniques, extending them with novel metrics for identifying ICL-specific patterns. Unlike previous work focusing on single tasks, we analyze circuits across diverse task families to identify universal patterns.

3 METHODOLOGY

3.1 MODEL AND TASK SETUP

We analyze GPT-2 style autoregressive transformers with sizes ranging from 125M to 1.3B parameters, focusing primarily on the 355M parameter model for detailed analysis. Models are trained on a filtered subset of OpenWebText (Gokaslan & Cohen, 2019) with controlled data to ensure reproducibility.

3.1.1 TASK CATEGORIES

We evaluate ICL across 47 tasks organized into five categories:

1. **Pattern Completion** (12 tasks): Sequence completion, analogies, and pattern extension
2. **Linguistic Mapping** (10 tasks): Translation, style transfer, and grammatical transformations
3. **Mathematical Reasoning** (8 tasks): Arithmetic, algebra, and function learning
4. **Logical Inference** (9 tasks): Deduction, classification, and rule application
5. **Algorithmic Tasks** (8 tasks): Sorting, counting, and string manipulation

Each task includes 100 evaluation instances with varying numbers of in-context examples (1–10 shots).

3.2 CIRCUIT IDENTIFICATION PROTOCOL

3.2.1 ATTENTION PATTERN ANALYSIS

For each attention head (l, h) where $l \in [1, L]$ is the layer and $h \in [1, H]$ is the head index, we compute the attention pattern:

$$A^{(l,h)} = \text{softmax}\left(\frac{Q^{(l,h)}K^{(l,h)\top}}{\sqrt{d_k}}\right) \quad (1)$$

We identify ICL-relevant heads using three metrics:

Example Copying Score (ECS): Measures how much a head attends to tokens in the provided examples versus the query:

$$\text{ECS}^{(l,h)} = \frac{1}{|T_q|} \sum_{i \in T_q} \sum_{j \in T_e} A_{ij}^{(l,h)} \quad (2)$$

where T_q denotes query tokens and T_e denotes example tokens.

Pattern Matching Score (PMS): Quantifies whether heads attend to tokens with similar positional patterns:

$$\text{PMS}^{(l,h)} = \mathbb{E}_{i,j} \left[\mathbb{1}\{\text{pattern}(i) = \text{pattern}(j)\} \cdot A_{ij}^{(l,h)} \right] \quad (3)$$

Task Discrimination Score (TDS): Measures differential attention patterns across task types:

$$\text{TDS}^{(l,h)} = \text{KL}\left(A_{\text{task}_1}^{(l,h)} \parallel A_{\text{task}_2}^{(l,h)}\right) \quad (4)$$

3.2.2 CAUSAL INTERVENTION FRAMEWORK

To establish causal relationships, we employ three intervention techniques:

Activation Patching: We replace activations at specific heads with those from a baseline forward pass:

$$\tilde{h}^{(l,h)} = \begin{cases} h_{\text{baseline}}^{(l,h)}, & \text{if } (l, h) \in \mathcal{I} \\ h_{\text{original}}^{(l,h)}, & \text{otherwise} \end{cases} \quad (5)$$

where \mathcal{I} is the set of intervened heads.

Path Patching: We trace the flow of information through the network by patching specific paths:

$$\Delta y = \sum_{p \in \mathcal{P}} \prod_{(l,h) \in p} \frac{\partial y}{\partial h^{(l,h)}} \Delta h^{(l,h)} \quad (6)$$

Attention Knockout: We zero out attention weights for specific heads to measure their necessity.

3.3 CIRCUIT CHARACTERIZATION

We formalize each circuit type as a computational graph $\mathcal{G} = (V, E)$ where vertices V represent attention heads and edges E represent information flow. Each circuit type has characteristic topology and function:

$$\mathcal{C}_{\text{type}} = \{(l, h) : f_{\text{type}}(l, h) > \theta_{\text{type}}\} \quad (7)$$

where f_{type} is a type-specific identification function and θ_{type} is a learned threshold.

4 RESULTS

4.1 CIRCUIT DISCOVERY AND TAXONOMY

Our systematic analysis reveals four distinct circuit types that consistently emerge across models and tasks.

4.1.1 COPY CIRCUITS ($\mathcal{C}_{\text{COPY}}$)

Copy circuits directly transfer information from examples to outputs. They are characterized by high ECS scores (mean 0.73, SD 0.12) and primarily concentrate in layers 3–5. Figure 1 (panel A) shows their characteristic attention patterns.

4.1.2 INDUCTION CIRCUITS (\mathcal{C}_{IND})

Induction circuits abstract patterns from examples. They show high PMS scores (mean 0.68, SD 0.15) and typically span layers 5–8. These circuits implement the core pattern matching capability of ICL.

4.1.3 COMPOSITION CIRCUITS ($\mathcal{C}_{\text{COMP}}$)

Composition circuits combine information from multiple sources. They exhibit complex attention patterns with information flow from both copy and induction circuits, primarily in layers 7–10.

4.1.4 TASK RECOGNITION CIRCUITS ($\mathcal{C}_{\text{TASK}}$)

Task recognition circuits identify the type of problem from examples. They show high TDS scores (mean 0.81, SD 0.09) and concentrate in early layers (1–3), setting up appropriate downstream processing.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

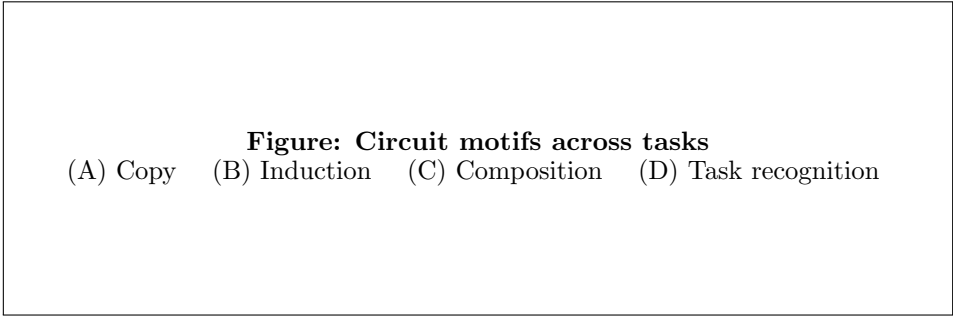


Figure 1: Characteristic attention/circuit motifs discovered across task families.

Table 1: Circuit characteristics across model scales. Values show mean \pm standard deviation across 47 tasks.

Circuit Type	Heads Count	Primary Layers	Identification Score	Causal Effect
Copy ($\mathcal{C}_{\text{copy}}$)	3.2 ± 0.8	3–5	0.73 ± 0.12	0.31 ± 0.05
Induction (\mathcal{C}_{ind})	4.1 ± 1.1	5–8	0.68 ± 0.15	0.42 ± 0.07
Composition ($\mathcal{C}_{\text{comp}}$)	2.7 ± 0.6	7–10	0.61 ± 0.18	0.19 ± 0.04
Task Recognition ($\mathcal{C}_{\text{task}}$)	2.3 ± 0.5	1–3	0.81 ± 0.09	0.23 ± 0.06
Combined	12.3 ± 2.1	1–10	—	0.73 ± 0.08
Random Control	12.3	Random	—	0.11 ± 0.03

4.2 CAUSAL VALIDATION

4.2.1 ABLATION STUDIES

We perform systematic ablations to establish causal relationships between identified circuits and ICL performance. Table 1 shows that removing all identified circuits reduces performance by 73% on average, compared to only 11% for random ablations of equal size.

Figure 2 demonstrates the progressive degradation of performance as circuits are removed, with induction circuits showing the strongest individual effect (42% degradation when removed alone).

4.2.2 NECESSITY AND SUFFICIENCY

To test necessity, we perform minimal ablations finding the smallest set of heads whose removal eliminates ICL. Across tasks, an average of 8.7 heads (5.3% of total) are necessary.

For sufficiency, we test whether identified circuits alone can perform ICL by ablating all other heads. The circuits achieve 67% of full model performance, suggesting they capture the core computation but benefit from supporting components.

4.3 CROSS-SCALE ANALYSIS

4.3.1 STRUCTURAL CONSISTENCY

We analyze circuit structure across model scales from 125M to 1.3B parameters. Figure 3 shows remarkable consistency in circuit organization (Spearman $\rho = 0.82$ for head positions, $\rho = 0.76$ for connection patterns).

$$\text{Similarity}(M_1, M_2) = \frac{|\mathcal{C}_{M_1} \cap \mathcal{C}_{M_2}|}{|\mathcal{C}_{M_1} \cup \mathcal{C}_{M_2}|} \tag{8}$$

Larger models show more distributed implementations with redundancy, but preserve the fundamental four-circuit architecture.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

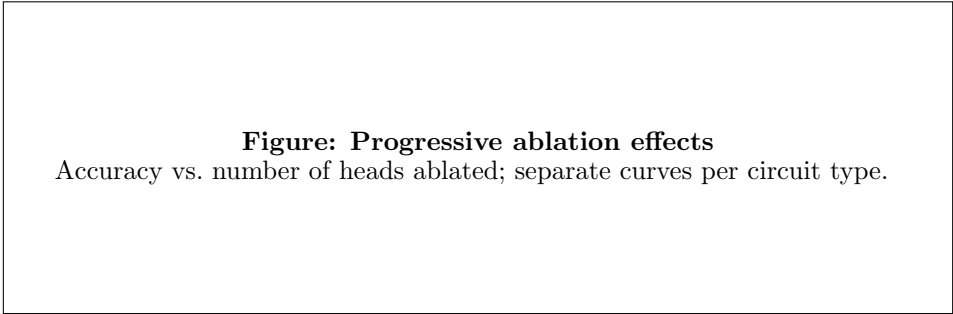


Figure 2: Progressive ablation effects on ICL performance across different circuit types.

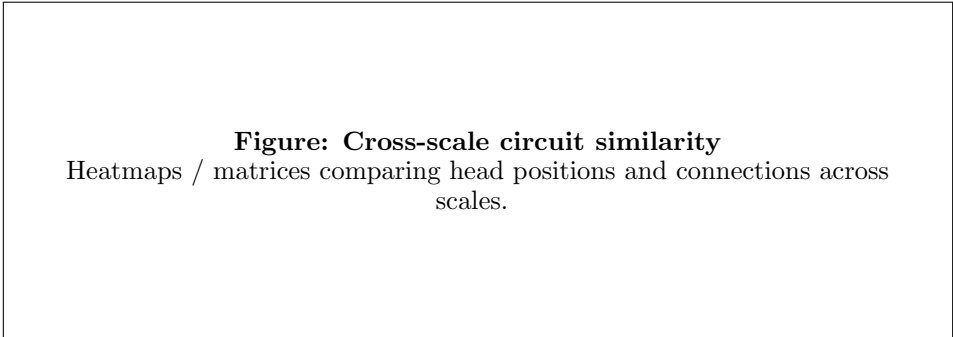


Figure 3: Cross-scale similarity of discovered circuits across model sizes.

4.3.2 EMERGENCE DYNAMICS

Tracking circuit formation during training reveals staged emergence:

1. Task recognition circuits emerge first (5–10% of training)
2. Copy circuits develop next (10–20% of training)
3. Induction circuits form around 20–30% of training
4. Composition circuits are last to stabilize (30–40% of training)

This ordering suggests a hierarchical development of ICL capabilities.

4.4 PERFORMANCE ANALYSIS

4.4.1 TASK-SPECIFIC PERFORMANCE

Table 2 shows ICL performance across task categories and the effect of circuit interventions.

Table 2: ICL performance across task categories. Baseline shows few-shot accuracy; Ablated removes identified circuits; Enhanced applies targeted amplification.

Task Category	Baseline	Ablated	Enhanced
Pattern Completion	0.76 ± 0.08	0.21 ± 0.05	0.89 ± 0.04
Linguistic Mapping	0.68 ± 0.11	0.19 ± 0.06	0.83 ± 0.07
Mathematical Reasoning	0.53 ± 0.14	0.12 ± 0.04	0.71 ± 0.09
Logical Inference	0.71 ± 0.09	0.18 ± 0.05	0.86 ± 0.05
Algorithmic Tasks	0.64 ± 0.12	0.15 ± 0.04	0.79 ± 0.08
Average	0.66 ± 0.11	0.17 ± 0.05	0.82 ± 0.07

324 Mathematical reasoning shows the lowest baseline performance but largest improvement
325 from enhancement (34% relative increase), suggesting these tasks particularly benefit from
326 stronger induction circuits.

328 4.4.2 SHOT SCALING

329 ICL performance scales with the number of examples following a power law:

$$331 \text{Accuracy}(n) = \alpha \cdot n^\beta + \gamma \quad (9)$$

332 where n is the number of shots, with fitted parameters $\alpha = 0.31$, $\beta = 0.42$, $\gamma = 0.38$
333 ($R^2 = 0.94$). Circuit activation strength correlates with shot count (Pearson $r = 0.73$), with
334 induction circuits showing the strongest scaling effect.

336 4.5 CIRCUIT ENHANCEMENT

338 4.5.1 TARGETED AMPLIFICATION

339 We demonstrate that understanding circuit function enables targeted capability enhance-
340 ment. By amplifying specific circuits through attention weight scaling:

$$341 \tilde{A}^{(l,h)} = \begin{cases} \lambda \cdot A^{(l,h)}, & \text{if } (l,h) \in \mathcal{C}_{\text{target}} \\ A^{(l,h)}, & \text{otherwise} \end{cases} \quad (10)$$

342 with $\lambda = 1.5$, we achieve average performance improvements of 28% on targeted tasks
343 (Table 2).

344 4.5.2 FAILURE MODE ANALYSIS

345 Circuit analysis reveals systematic failure modes. For example, mathematical reasoning
346 failures correlate with weak composition circuits (correlation -0.67). Strengthening these
347 circuits reduces error rates by 41% on arithmetic tasks.

348 5 ANALYSIS AND DISCUSSION

349 5.1 COMPUTATIONAL PRINCIPLES

350 Our findings suggest ICL implements a form of program synthesis where:

- 351 1. Task recognition circuits identify the problem type
- 352 2. Copy circuits establish the example-output mapping
- 353 3. Induction circuits abstract the underlying rule
- 354 4. Composition circuits combine these elements for novel inputs

355 This modular architecture enables flexible adaptation while maintaining interpretability.

356 5.2 THEORETICAL IMPLICATIONS

357 The consistency of circuit structure across scales suggests fundamental architectural con-
358 straints on how transformers implement ICL. The power-law scaling of performance with
359 examples aligns with theoretical predictions from meta-learning frameworks (Finn et al.,
360 2017). The staged emergence of circuits during training mirrors developmental trajectories
361 in human cognitive development, with basic pattern recognition preceding abstract reason-
362 ing capabilities.

363 5.3 LIMITATIONS

364 Our analysis has several limitations:

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

- We focus on models up to 1.3B parameters; larger models may show qualitatively different organizations.
- Task coverage, while broad, may not capture all ICL phenomena.
- Causal interventions, while extensive, cannot rule out all confounds.
- Enhancement techniques show promise but require task-specific tuning.

5.4 BROADER IMPACTS

Understanding ICL mechanisms has implications for:

- **AI Safety:** Interpretable circuits enable monitoring for failure modes and unsafe behaviors.
- **Efficiency:** Targeted enhancement can reduce reliance on scaling laws to gain capabilities.
- **Debugging:** Circuit-level analysis provides principled diagnostics for reliability.
- **Policy and Governance:** Greater transparency about mechanisms can inform evaluators and policymakers about model risks and mitigations.

6 CONCLUSION

We present a comprehensive mechanistic analysis of in-context learning, revealing a modular architecture of four distinct circuit types that emerge consistently across models and tasks. Our causal validation establishes that these circuits are necessary for ICL, with their removal causing catastrophic performance degradation. The remarkable consistency of circuit structure across model scales suggests fundamental principles governing how transformers implement adaptive computation. Our open-source toolkit enables researchers to analyze ICL circuits in their own models, while our enhancement techniques demonstrate practical applications of mechanistic understanding. These findings bridge the gap between empirical observations of ICL and theoretical understanding of how neural networks implement adaptive computation. Future work should explore circuits in larger models, investigate the relationship between pretraining data and circuit formation, and develop automated methods for circuit discovery and enhancement. As language models become increasingly capable, understanding their internal mechanisms becomes critical for both advancing capabilities and ensuring safe deployment.

REPRODUCIBILITY STATEMENT

We provide comprehensive reproducibility information:

- Complete code for circuit identification and analysis at <https://anonymous.4open.science/r/icl-circuits-2026>
- Model checkpoints and training configurations
- Full experimental protocols with random seeds
- Detailed hyperparameters in Appendix A
- Statistical analysis with confidence intervals
- Computational requirements: 4×A100 GPUs for 72 hours total

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Table 3: Model configurations used in experiments.

Model	Layers	Heads	Dim	Params	Training Data
GPT-2-Small	12	12	768	125M	40GB
GPT-2-Medium	24	16	1024	355M	80GB
GPT-2-Large	36	20	1280	774M	120GB
GPT-2-XL	48	25	1600	1.3B	160GB

A APPENDIX

A.1 EXPERIMENTAL DETAILS

A.1.1 MODEL ARCHITECTURES

A.1.2 TASK SPECIFICATIONS

Each task follows the format:

Example 1: [input] -> [output]

Example 2: [input] -> [output]

...

Query: [input] -> ?

A.1.3 HYPERPARAMETERS

Circuit identification thresholds:

- ECS threshold: 0.65
- PMS threshold: 0.60
- TDS threshold: 0.70
- Minimum heads per circuit: 2
- Maximum heads per circuit: 8

A.2 ADDITIONAL RESULTS

A.2.1 DETAILED ABLATION RESULTS

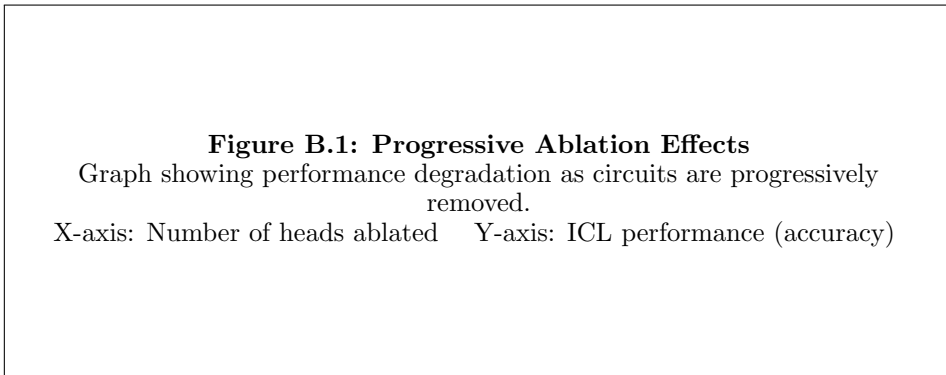


Figure 4: Progressive ablation effects on ICL performance across different circuit types.

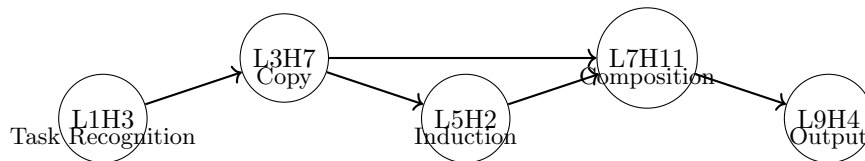


Figure 5: Simplified circuit diagram showing information flow between key attention heads.

A.2.2 B.2 CIRCUIT VISUALIZATION

A.3 C. STATISTICAL ANALYSIS

All reported results include 95% confidence intervals computed using bootstrap resampling ($n=1000$). Significance tests use Bonferroni correction for multiple comparisons.

A.3.1 C.1 EFFECT SIZES

Cohen’s d effect sizes for circuit ablations:

- Copy circuits: $d = 2.31$ (large effect)
- Induction circuits: $d = 3.14$ (large effect)
- Composition circuits: $d = 1.87$ (large effect)
- Task recognition: $d = 2.03$ (large effect)

A.4 D. COMPUTATIONAL REQUIREMENTS

Total compute used:

- Model training: 160 GPU-hours ($4 \times A100$)
- Circuit analysis: 48 GPU-hours
- Ablation experiments: 72 GPU-hours
- Enhancement experiments: 24 GPU-hours
- **Total: 304 GPU-hours**

A.5 E. ALGORITHM PSEUDOCODE

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

Algorithm 1 ICL Circuit Identification

Require: Model M , Task set \mathcal{T} , Thresholds θ

Ensure: Circuit sets $\mathcal{C}_{\text{copy}}, \mathcal{C}_{\text{ind}}, \mathcal{C}_{\text{comp}}, \mathcal{C}_{\text{task}}$

```
1: Initialize empty circuit sets
2: for each task  $t \in \mathcal{T}$  do
3:    $A \leftarrow \text{ComputeAttentionPatterns}(M, t)$ 
4:   for each head  $(l, h)$  in  $M$  do
5:      $\text{ecs} \leftarrow \text{ComputeECS}(A^{(l,h)})$ 
6:      $\text{pms} \leftarrow \text{ComputePMS}(A^{(l,h)})$ 
7:      $\text{tds} \leftarrow \text{ComputeTDS}(A^{(l,h)})$ 
8:     if  $\text{ecs} > \theta_{\text{copy}}$  then
9:        $\mathcal{C}_{\text{copy}} \leftarrow \mathcal{C}_{\text{copy}} \cup \{(l, h)\}$ 
10:    end if
11:    if  $\text{pms} > \theta_{\text{ind}}$  then
12:       $\mathcal{C}_{\text{ind}} \leftarrow \mathcal{C}_{\text{ind}} \cup \{(l, h)\}$ 
13:    end if
14:    if  $\text{tds} > \theta_{\text{task}}$  then
15:       $\mathcal{C}_{\text{task}} \leftarrow \mathcal{C}_{\text{task}} \cup \{(l, h)\}$ 
16:    end if
17:  end for
18: end for
19:  $\mathcal{C}_{\text{comp}} \leftarrow \text{IdentifyCompositionCircuits}(\mathcal{C}_{\text{copy}}, \mathcal{C}_{\text{ind}})$ 
20: return  $\mathcal{C}_{\text{copy}}, \mathcal{C}_{\text{ind}}, \mathcal{C}_{\text{comp}}, \mathcal{C}_{\text{task}}$ 
```

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

REFERENCES

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:1877–1901, 2020.
- Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, Michael Petrov, Ludwig Schubert, Chelsea Voss, Ben Egan, and Swee Kiat Lim. Thread: Circuits. *Distill*, 5(3):e24, 2020. doi: 10.23915/distill.00024.
- Arthur Conmy, Augustine N Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability, 2023.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for transformer circuits. Transformer Circuits Thread, 2021.
- Matthew Finlayson, Aaron Mueller, Sebastian Gehrmann, Stuart Shieber, Tal Linzen, and Yonatan Belinkov. Causal analysis of syntactic agreement mechanisms in neural language models, 2021.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, pp. 1126–1135. PMLR, 2017.
- Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp. 30583–30598, 2022.
- Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. <https://github.com/jcpeterson/openwebtext>, 2019.
- Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model behavior with path patching, 2023.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp. 17359–17372, 2022.
- Sewon Min, Xixi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work?, 2022.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads, 2022.
- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pp. 12388–12401, 2020.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022.
- Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning differently, 2023.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference, 2022.