
Same Story, Different Circuits: Measuring Component Reuse in Transformer Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A central premise of mechanistic interpretability is the existence of task-specific cir-
2 cuits within models. This work investigates the consistency of these circuits by ask-
3 ing how often transformer models *reuse* the same internal parts across different ex-
4 amples of a task. We use edge attribution patching to construct per-example circuits
5 composed of attention heads and MLP layers. We then define a metric, **reuse@P**, to
6 quantify how frequently the same components are reused across different circuits of
7 a task. To verify that shared circuits are actually important for task performance, we
8 knock out these components and compare against a strong baseline, which we call
9 **capacity-conserved control** (C^3). If shared circuit knockout causes larger accuracy
10 drops than the control, this confirms the components' causal specificity for the task.
11 We evaluate this approach across 6 tasks including Addition, Boolean Logic,
12 Indirect Object Identification, CopyColors MCQA, MMLU and AI2 ARC,
13 and 7 models from 4 model families including Qwen3, Llama 3.2, Gemma 2 and
14 QLMo-2-1B. We find that circuit reuse is present in all models but varies by task
15 type: algorithmic tasks exhibit higher reuse than knowledge-intensive ones. Our
16 knockout studies confirm the causal importance of shared circuits, but also sug-
17 gest the presence of functionally redundant pathways within the models. These
18 findings indicate that transformers may solve tasks not with a single, minimal
19 circuit, but with an ensemble of overlapping solutions, opening new directions for
20 understanding generalization and robustness.

21

1 Introduction

22 Neural networks are infamously black box; even when we can elicit strong performance on a task, it
23 is unclear which internal computations are responsible. The field of mechanistic interpretability aims
24 to reverse-engineer the algorithms learned by these models, typically by identifying circuits within a
25 model that explain a behavior. Recent work has reverse-engineered circuits responsible for specific
26 behaviors in transformer language models, such as the Indirect Object Identification (IOI)
27 circuit in GPT-2 (Wang et al., 2023), and documented how apparently similar circuit components can
28 be reused across tasks (Merullo et al., 2024). However, it remains unclear how *within-task circuit*
29 *reuse* behaves.

30 In this paper, we ask a simple question: when a model solves many instances of the same task, does it
31 rely on the same internal parts? Following the Transformer Circuits framework (Elhage et al., 2021),
32 we define circuit reuse as the set of components (attention heads and MLP layers) that recur across
33 inputs within the same task. We map circuits at scale with edge attribution patching (EAP) (Syed
34 et al., 2024), then identify the components appearing in at least $P\%$ of circuits, defined as **reuse@P**.

35 To test whether this shared circuit is functionally important, we knock out its components and measure
36 the resulting drop in task accuracy. To properly calibrate these knockouts, we introduce a new baseline:

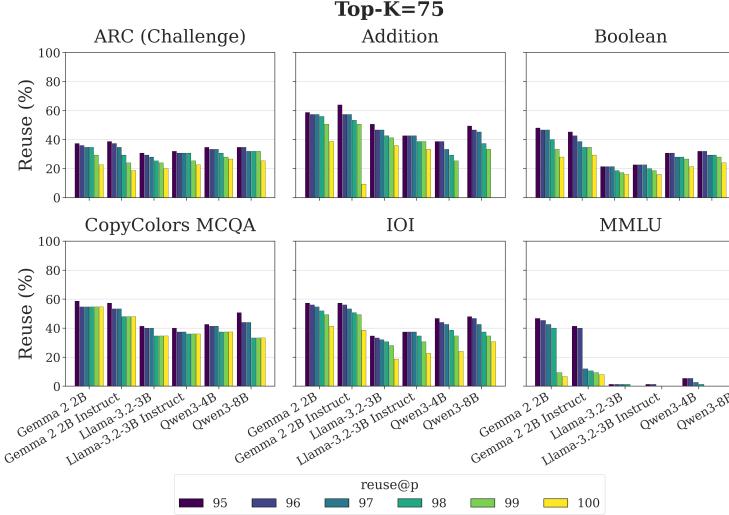


Figure 1: **Circuit reuse varies by task type for top-K=75.** reuse@P measures the percentage of circuit components appearing in at least P% of task example circuits. Reuse is high for algorithmic tasks like addition, and low for knowledge-intensive tasks like MMLU.

37 the **Capacity-Conserved Control (C³)**, which removes an equal number of randomly sampled
 38 components. This control helps verify that any observed effects reflect functional specialization,
 39 not just generic capacity reduction. We evaluate circuit reuse on a variety of tasks (e.g. Addition,
 40 Boolean Logic, IOI, CopyColors MCQA, AI2 ARC, MMLU) with a wide range of open-weight
 41 models (e.g. QLMo-2-1B, Gemma 2 2B, Llama 3.2 3B, Qwen3 8B).

42 Our paper makes several contributions:

- 43 1. We propose **reuse@P**, a useful metric for quantifying the degree of circuit reuse across a
 44 distribution of inputs for a specific task.
- 45 2. We introduce the **Capacity-Conserved Control (C³)**, a strong baseline for circuit knockouts
 46 that helps calibrate whether observed performance drops are actually due to the specific choice
 47 of components, not just the reduction in model capacity.
- 48 3. We conduct a large-scale empirical study of circuit reuse trends across seven distinct tasks and
 49 seven transformer models (differing in size, architecture and post-training).

50 2 Background

51 **Transformer Circuits.** The Transformer Circuits framework is widely used for mechanistic inter-
 52 pretability to reason about model computations and component-level interventions (Elhage et al.,
 53 2021). Within this framework, a decoder-only transformer can be viewed as a directed computation
 54 over layers operating on a central residual stream. The residual stream acts as a high-dimensional
 55 communication channel where each layer reads information via linear projections and writes results
 56 back through additive contributions. Attention heads and MLP layers are treated as the primary
 57 computational components. A circuit is then defined as a sparse subgraph of components and their
 58 connecting edges that is sufficient, and minimal, for implementing a given behavior.

59 **Edge Attribution Patching.** To identify which components form a circuit for a given behavior,
 60 researchers typically use causal intervention methods – systematically ablating or patching activations
 61 to measure each component’s contribution to the model’s output. Edge Attribution Patching
 62 (EAP) is an efficient automated circuit discovery method that approximates the effect of activation
 63 patching (a computationally expensive technique) using gradient information (Syed et al., 2024).
 64 EAP estimates the importance of all edges using just two forward passes and one backward pass by
 65 computing attribution scores based on activation differences and gradients. Edges are then ranked by

66 their absolute attribution scores and those below a given top-K threshold are pruned to identify the
67 minimal circuit.

68 **Tasks.** We evaluate circuit reuse across a diverse set of tasks designed to probe model capabilities.
69 They can be grouped into three main categories:

- 70 • **Algorithmic Tasks:** These tasks follow a clear, learnable algorithm. We use 3-digit Addition
71 and Boolean Logic evaluation. High circuit reuse is expected as the underlying algorithm is
72 consistent across all inputs.
- 73 • **Linguistic & Reasoning Tasks:** These require understanding sentence structure and context. We
74 use Indirect Object Identification (IOI) (Wang et al., 2023) and CopyColors MCQA.
- 75 • **Knowledge-Intensive Tasks:** These tasks require recalling factual knowledge, potentially stored
76 heterogeneously across the model’s parameters. We use the high school European history subset
77 from MMLU (Hendrycks et al., 2021) and AI2 ARC (Clark et al., 2018).

78 3 Methodology

79 To investigate within-task circuit reuse, we construct shared circuits for each task to identify important
80 components across multiple examples. For each task, we sample $n = 1000$ input/output pairs from the
81 training set and extract per-example circuits using edge attribution patching (EAP), which produces
82 importance scores for every attention head and MLP layer. We define the per-input circuit as the
83 top-K components by score for each example.

84 Given these per-input circuits from the training data, we construct a shared circuit at threshold P
85 containing components that appear in at least P% of training examples. We quantify the extent
86 of this sharing using **reuse@P**: the fraction of a typical circuit covered by components appearing
87 in at least P% of examples (formal definition in Appendix B). For example, if 75% of the top-
88 K components appear in at least 95% of task examples, then reuse@95 = 75%. This measures
89 circuit consistency across the task distribution - higher values indicate that models rely on the same
90 components repeatedly.

91 To validate that shared circuits are causally important rather than just frequent, we perform knockout
92 experiments on held-out validation data, comparing shared circuit removal against two controls:

- 93 • **Capacity-Conserved with parity (C^3 -parity).** Match the head/MLP counts of the shared circuit
94 and sample from the entire pool (including shared circuit components). This control maintains
95 capacity and component mix while allowing overlap with the shared circuit.
- 96 • **Capacity-Conserved without parity (C^3 -exclude).** Match the head/MLP counts but explicitly
97 exclude all shared circuit components from sampling. This maximizes contrast with the shared
98 circuit while still controlling for capacity and component type.

99 For each knockout, we zero-out the outputs of the selected components and measure the resulting
100 accuracy on the validation set. We measure relative impact using *lift*, defined as:

$$\text{lift} = \frac{\text{acc}_{\text{shared}} - \text{acc}_{\text{control}}}{\text{acc}_{\text{base}}}$$

101 where acc_{base} is baseline accuracy on the validation set, $\text{acc}_{\text{shared}}$ is validation accuracy after
102 knocking out S_P , and $\text{acc}_{\text{control}}$ is validation accuracy after knocking out a control set. If shared
103 circuit knockout causes larger accuracy drops than controls, this confirms the components’ causal
104 specificity for the task.

105 We hypothesize that, holding K and P fixed, higher within-task reuse will correlate with more
106 negative lift under C^3 -exclude; in other words, we think that highly reused components are more
107 likely to be the task-specific circuit. This is a directional hypothesis and not assumed by our
108 evaluation.

109 4 Results

110 Our experiments show that circuit reuse is a consistent phenomenon, though its magnitude varies with
111 task type. As shown in Figure 1, algorithmic tasks and linguistic tasks (Addition, IOI, CopyColors

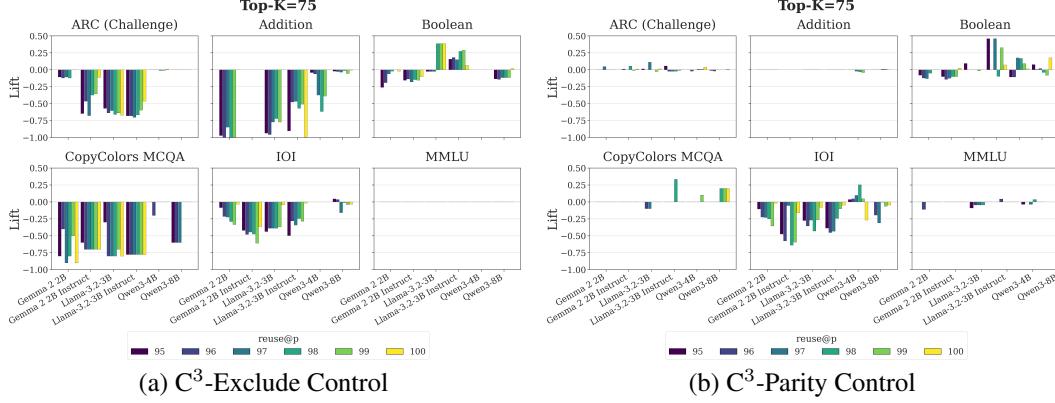


Figure 2: Lift in task accuracy when knocking out shared circuits vs. two controls for top-K=75.
(a) The C³-exclude control knocks out a disjoint set of components, yielding a strong negative lift.
(b) The C³-parity control knocks out a random set that can overlap, resulting in a near-zero lift.

MCQA) exhibit high reuse across all models, often exceeding 50% at P=95. Meanwhile, knowledge-intensive and complex reasoning tasks like MMLU and ARC show lower reuse. For MMLU, reuse is notably low for all models except the Gemma family, suggesting that different models may adopt distinct strategies for storing and accessing factual knowledge. We observe that reuse increases with the circuit size K, though not always monotonic (see Appendix C for full results). Additionally, unlike the other models we studied, OLMo-2-1B released its pretraining checkpoints. This allowed us to extend our analysis and observe that circuit reuse increases as the model undergoes more training steps, especially for larger circuit sizes K (see Figure 7 and Figure 6).

Our knockout studies yield two key findings. As seen in Figure 2, knocking out the shared circuit against the *C³-exclude* control produces a strong negative lift across nearly all tasks and models. This indicates that the shared components are more critical for task performance than a disjoint set of the same size. However, with the *C³-parity* control, shown in Figure 2. In this condition, the lift is consistently near-zero across all tasks and models. This implies that knocking out the shared circuit is no more damaging than knocking out a randomly selected set of components of the same size, provided that the random set is allowed to overlap with the shared circuit.

The C³-exclude result confirms the shared circuit's causal role, while the C³-parity result suggests its components are not unique in this role. We interpret this by hypothesizing that models develop functionally redundant pathways to solve tasks. If so, components in the shared circuit aren't unique; other components can perform similar functions. Then a random knockout is likely to damage these alternative pathways, leading to a performance drop comparable to knocking out the shared circuit itself. This suggests transformer circuits function less like a single, minimal algorithm and more like an ensemble of overlapping solutions.

5 Conclusion

In this paper, we presented a methodology (**reuse@P** and the **Capacity-Conserved Control (C³)**) for quantifying and validating circuit reuse in transformer models. Our analysis covered seven models and six tasks, suggesting that circuit reuse is a task-dependent phenomenon, with algorithmic tasks showing higher consistency than knowledge-intensive ones. Our knockout studies, however, show mixed results. While we see clear evidence that shared circuits are important for task performance, we cannot rule out the possibility that models develop functionally redundant pathways that complicate the isolation of a single, task-specific circuit. To resolve this ambiguity, future work focuses on training small transformers on synthetic tasks with controlled settings to better understand and accurately characterize circuit reuse and redundancy in transformer models.

6 Limitations

Our approach has several limitations. First, Edge Attribution Patching (EAP) is a gradient-based approximation, less faithful than causal methods. More precise variants like EAP-IG (Hanna et al.,

147 2024) could improve circuit identification but are more expensive. Second, zero-ablation may
148 yield out-of-distribution activations, potentially confounding results. Alternatives such as mean-
149 ablation (Miller et al., 2024) may better preserve activation distributions. Finally, our study focuses
150 on specific models and tasks, so results may not generalize across architectures or domains.

151 **References**

- 152 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
153 Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning
154 challenge. *ArXiv*, abs/1803.05457.
- 155 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
156 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac
157 Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse,
158 and 6 others. 2021. A mathematical framework for transformer circuits. *Transformer Circuits*
159 *Thread*. <Https://transformer-circuits.pub/2021/framework/index.html>.
- 160 Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. 2024. Have faith in faithfulness: Going
161 beyond circuit overlap when finding model mechanisms. In *ICML 2024 Workshop on Mechanistic*
162 *Interpretability*.
- 163 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
164 Steinhardt. 2021. Measuring massive multitask language understanding. *Proceedings of the*
165 *International Conference on Learning Representations (ICLR)*.
- 166 Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2024. Circuit component reuse across tasks in
167 transformer language models. In *The Twelfth International Conference on Learning Representa-*
168 *tions*.
- 169 Joseph Miller, Bilal Chughtai, and William Saunders. 2024. Transformer circuit evaluation metrics
170 are not robust. In *First Conference on Language Modeling*.
- 171 Aaron Mueller, Atticus Geiger, Sarah Wiegreffe, Dana Arad, Iván Arcuschin, Adam Belfki, Yik Siu
172 Chan, Jaden Fiotto-Kaufman, Tal Haklay, Michael Hanna, Jing Huang, Rohan Gupta, Yaniv
173 Nikankin, Hadas Orgad, Nikhil Prakash, Anja Reusch, Aruna Sankaranarayanan, Shun Shao,
174 Alessandro Stolfo, and 4 others. 2025. MIB: A mechanistic interpretability benchmark. In
175 *Forty-second International Conference on Machine Learning*.
- 176 Aaquib Syed, Can Rager, and Arthur Conmy. 2024. Attribution patching outperforms automated
177 circuit discovery. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting*
178 *Neural Networks for NLP*, pages 407–416, Miami, Florida, US. Association for Computational
179 Linguistics.
- 180 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023.
181 Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The*
182 *Eleventh International Conference on Learning Representations*.

183 A Task Details

184 This section provides descriptions and examples for each task used in our circuit reuse analysis.

Addition. The Addition task requires models to perform 3-digit arithmetic addition presented in natural language format. For example, given the input "What is $247 + 163$?", the model should output "410". We expect this task to be relatively simple for large language models to complete.

Boolean Logic. The Boolean task involves evaluating logical expressions with AND, OR, and NOT operations over boolean variables. An example input "True AND (False OR True)" should yield "True". This task requires systematic logical reasoning and an understanding of operator order.

191 Indirect Object Identification (IOI). Following Wang et al. (2023), the IOI task requires models
192 to identify the indirect object in sentences with specific templates involving two names.

193 *Input:* “As Carl and Maria left the consulate, Carl gave a fridge to”
194 *Expected Output:* “Maria”

195 This task is taken from the Mechanistic Interpretability Benchmark (MIB) (Mueller et al., 2025).

196 CopyColors MCQA. The CopyColors task presents multiple-choice questions requiring models to
197 identify and copy color information from a given context. It is similar to the IOI task.

198 *Input:* “Coconuts are brown. What color are coconuts? (A) red (B) orange (C)
199 brown (D) purple”
200 *Expected Output:* “(C) brown”

²⁰¹ This task is taken from the MIB (Mueller et al., 2025).

ARC. The AI2 Reasoning Challenge (Clark et al., 2018) contains science exam questions requiring reasoning over natural language. We use the Challenge set, which filters out questions answerable by simple retrieval.

Input: “George wants to warm his hands quickly by rubbing them. Which skin surface will produce the most heat? (A) dry palms (B) wet palms (C) palms covered with oil (D) palms covered with lotion”
Expected Output: “(A) dry palms”

²⁰⁹ This task is taken from the MIB (Mueller et al., 2025).

MMLU. The Massive Multitask Language Understanding benchmark (Hendrycks et al., 2021) spans 57 subjects from elementary to professional levels.

Input: “This question refers to the following information. Read the following quotation to answer questions. The various modes of worship which prevailed in the Roman world were all considered by the people as equally true; by the philosopher as equally false; and by the magistrate as equally useful. Edward Gibbon, The Decline and Fall of the Roman Empire, 1776–1788 Gibbon’s interpretation of the state of religious worship in ancient Rome could be summarized as (A) In ancient Rome, religious worship was decentralized and tended to vary with one’s social position. (B) In ancient Rome, religious worship was the source of much social tension and turmoil. (C) In ancient Rome, religious worship was homogeneous and highly centralized. (D) In ancient Rome, religious worship was revolutionized by the introduction of Christianity.”

Expected Output: “(A) In ancient Rome, religious worship was decentralized and tended to vary with one’s social position.”

225 **B Methodology Details**

226 For each task, we sample $n = 1000$ input/output pairs from the training set and extract a circuit for
227 every example using edge attribution patching (EAP). We define the per-input circuit \mathcal{C}_i as the top- K
228 components by score for example i (we sweep $K \in \{25, 50, 75, 100\}$; main figures show $K = 75$).

229 The shared circuit at threshold P is formally defined as:

$$S_P = \left\{ c : c \in \mathcal{C}_i \text{ for at least } \left\lceil \frac{P}{100} \cdot n \right\rceil \text{ training examples} \right\}.$$

230 The reuse@P metric is:

$$\text{reuse}@P = \frac{\min \{|S_P|, K\}}{K} \times 100$$

231 **C Additional Results**

232 This section provides lift and reuse results across full sweeps for all models, as well as lift and reuse
233 results for OLMo-2-1B pretraining checkpoints.

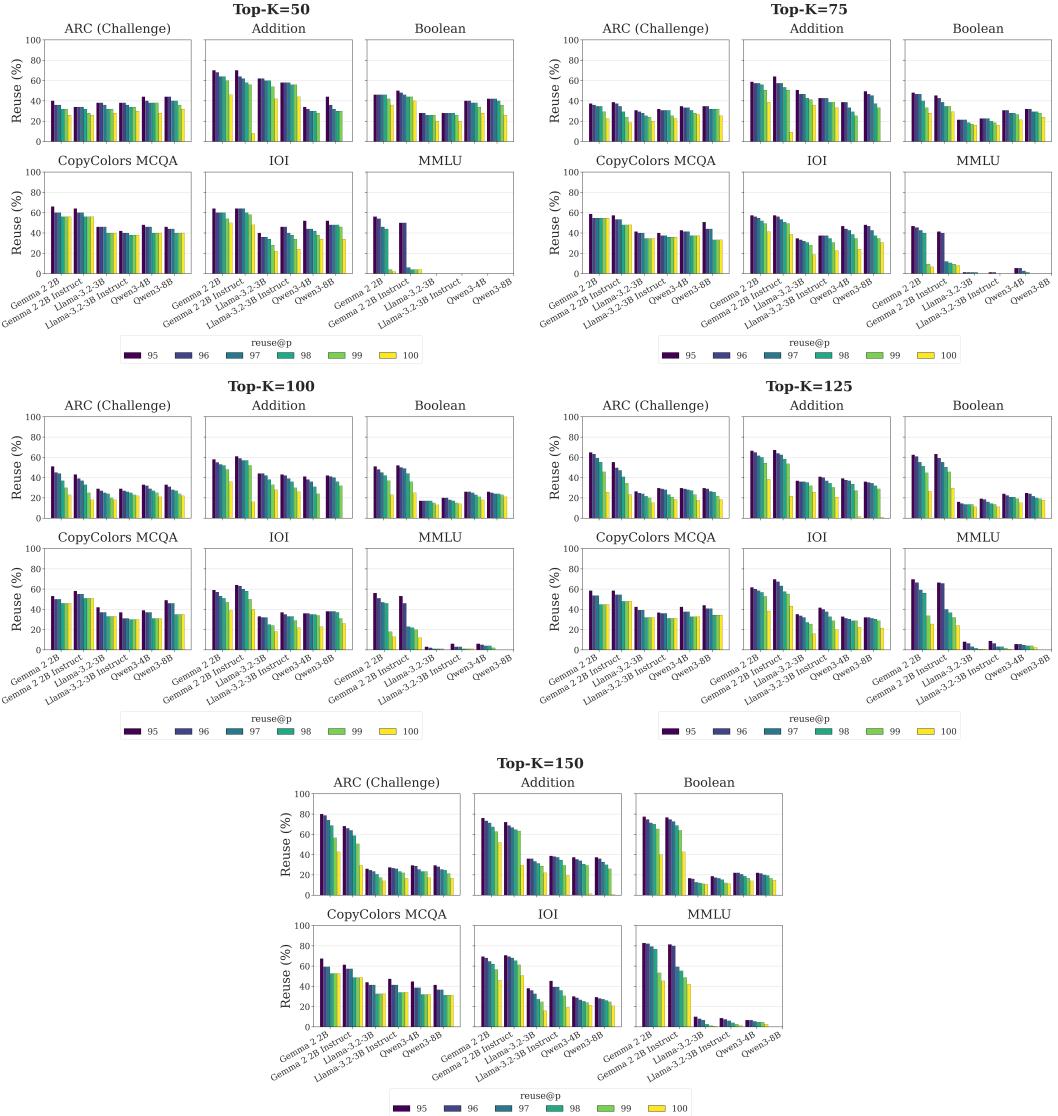


Figure 3: Reuse@P across all K.



Figure 4: Lift under C^3 -parity for all K.

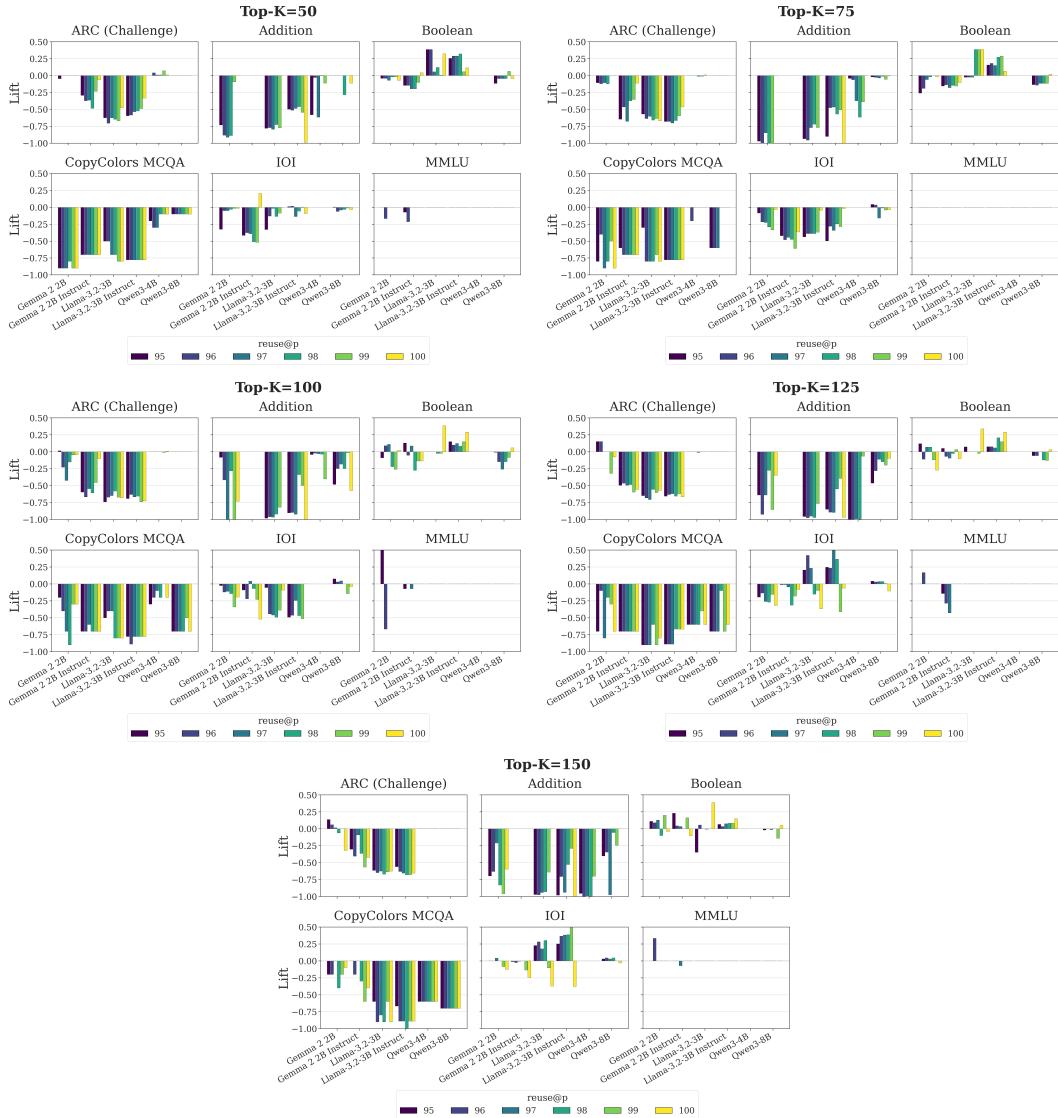


Figure 5: Lift under C^3 -exclude for all K.

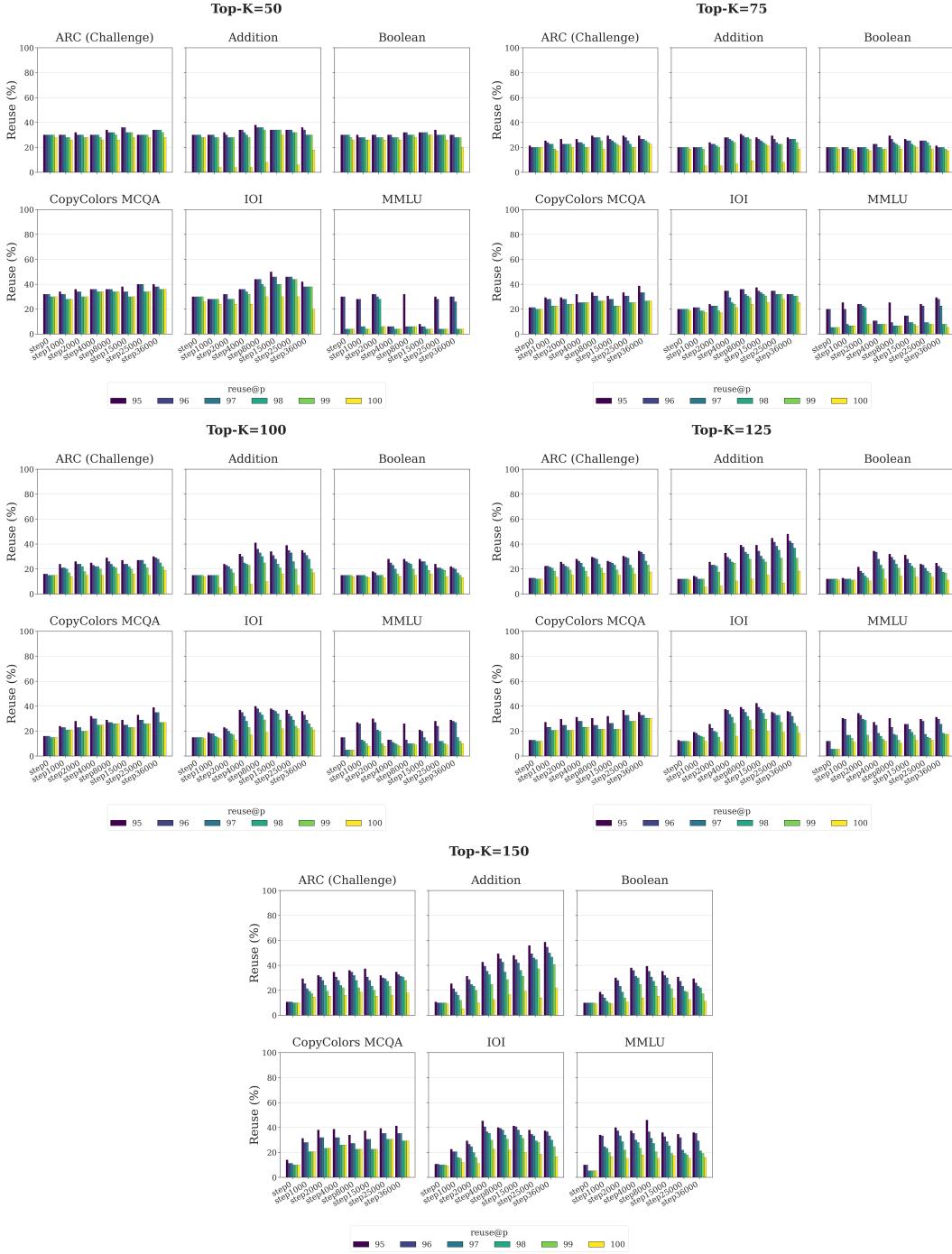


Figure 6: Reuse across OLMo-2 1B pretraining checkpoints.

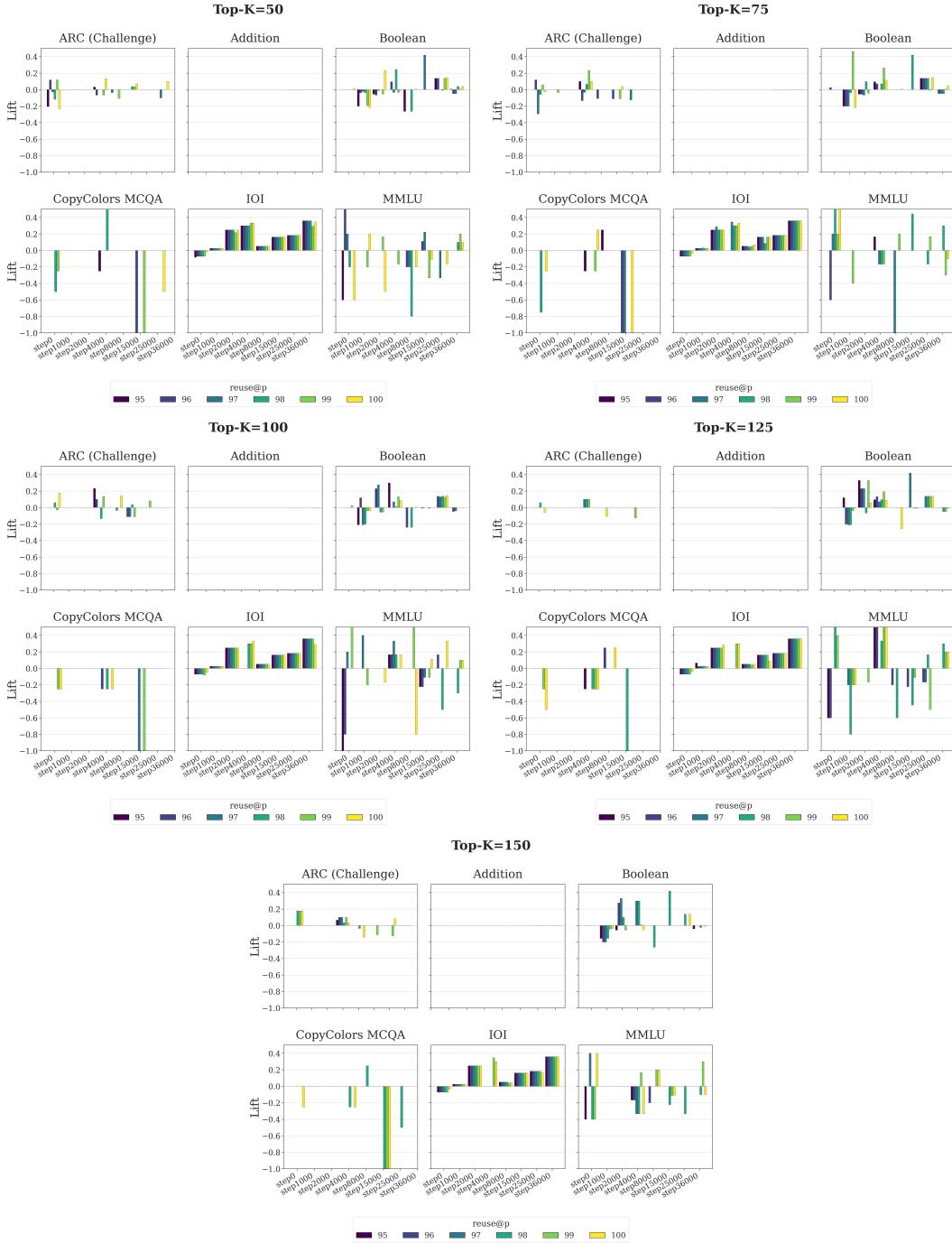


Figure 7: Lift under C^3 -parity for OLMo-2 1B checkpoints.