# Causal Head Gating: A Framework for Interpreting Roles of Attention Heads in Transformers

**Andrew J. Nam**
Princeton Laboratory for AI
Natural and Artificial Minds
Princeton University
andrewnam@princeton.edu

**Henry C. Conklin**
Princeton Laboratory for AI
Natural and Artificial Minds
Princeton University
henry.conklin@princeton.edu

**Yukang Yang**
Department of Electrical and Computer Engineering
Princeton University
yy1325@princeton.edu

**Thomas L. Griffiths**
Department of Psychology
Princeton University
tomg@princeton.edu

**Jonathan D. Cohen**\*
Princeton Neuroscience Institute
Princeton University
jdc@princeton.edu

**Sarah-Jane Leslie**\*
Department of Philosophy
Center for Statistics and Machine Learning
Princeton University
sjleslie@princeton.edu

## Abstract

We present *causal head gating* (CHG), a scalable method for interpreting the functional roles of attention heads in transformer models. CHG learns soft gates over heads and assigns them a causal taxonomy—facilitating, interfering, or irrelevant—based on their impact on task performance. Unlike prior approaches in mechanistic interpretability, which are hypothesis-driven and require prompt templates or target labels, CHG applies directly to any dataset using standard next-token prediction. We evaluate CHG across multiple large language models (LLMs) in the Llama 3 model family and diverse tasks, including syntax, commonsense, and mathematical reasoning, and show that CHG scores yield causal—not merely correlational—insight — validated via ablation and causal mediation analyses. We also introduce *contrastive* CHG, a variant that isolates sub-circuits for specific task components. Our findings reveal that LLMs contain multiple sparse, sufficient sub-circuits, that individual head roles depend on interactions with others (low modularity), and that instruction following and in-context learning rely on separable mechanisms.

## 1 Introduction

Large language models (LLMs) [1, 2, 3] represent the state-of-the art across a wide array of domains, exhibiting remarkable generalization and problem-solving capabilities. Yet, as these models grow in scale and complexity, they become increasingly opaque, making it more difficult to understand, predict, or control their behavior, which raises concerns about safety and misuse [4, 5, 6]. This has motivated a growing body of work on *interpretability*, which seeks to better understand how LLMs learn and represent information, and how their responses can be shaped [7, 8]. Interest has focused

---

\*Equal contribution; authors listed alphabetically

in particular on transformer-based architectures [9] such as GPT [1], LLaMA [3], Gemma [10], and DeepSeek [2], in which the central processing blocks consist of multi-head attention followed by multi-layer perceptrons. Here, there has been considerable research on the roles of individual attention heads, which have been found to exhibit some level of human-interpretability [11, 12, 13].

Two broad categories of approaches dominate research on mechanistic interpretability in LLMs. The first uses a trained mapping from latent representations to human-interpretable concepts, such as syntactic features [7, 14, 15] or identifiable items (e.g., the Golden Gate Bridge [16]). The second uses causal interventions to identify portions of a single weight matrix or individual attention heads responsible for a specific behavior [17, 18]. These approaches often focus on small portions of a model — 'zooming in' [19] in an effort to interpret the role of a single computational subgraph. However, in deep-learning models, computation is often distributed [20] and the role of one component is dependent on another [21, 22, 23], making the behavior of such complex distributed systems difficult to predict from an understanding of their parts alone [24].

To apply a distributed perspective to mechanistic interpretability, we introduce *causal head gating* (CHG) which seeks to identify a parametrically weighted set of heads that contribute to a model's execution of a given task. Given a dataset that defines a task, we fit a set of gating values for each attention head that applies a soft ablation to its output using next-token prediction, so that task-facilitating heads remain unaltered while any task-interfering heads are suppressed. Using a simple regularization procedure that further separates irrelevant heads from those that facilitate or interfere with task performance, CHG assigns meaningful scores to each attention head across an entire model according to its task contribution. We use these scores to define a taxonomy of task relevance according to how individual attention heads contribute to a model's distributed computation of a given task, describing each head as *facilitating, interfering* or *irrelevant*. In this respect, CHG offers an exploratory complement to standard hypothesis-driven approaches to mechanistic interpretablity, assigning causal roles without relying on predefined hypotheses about what each head might be doing.

Beyond its conceptual contribution, CHG also offers several practical methodological advantages over existing mechanistic interpretability tools. First, because CHG operates directly on next-token prediction, it avoids the need for externally-provided labels [7, 14, 15, 16], controlled input-output pairs [7, 14, 15], or rigid prompt templates [25, 12, 13], which are often required for decoding and interventional approaches. Second, CHG naturally accommodates complex target outputs, including chain-of-thought reasoning [26], where the solution spans multiple intermediate steps. Finally, CHG is highly scalable: it introduces only one learnable parameter per attention head and requires no updates to the underlying model weights, so that the CHG parameters can be fitted in minutes using gradient-based optimization, even for LLMs with billions of parameters.

To test its efficacy, we apply CHG across a diverse set of tasks—mathematical, commonsense, and syntactic reasoning—and across LLMs ranging from 1 to 8 billion parameters with varying training paradigms. We use CHG to analyze not only *where* specific computations take place, but also *how distributed* they are across attention heads, and how these patterns vary across different tasks and models. We also validate the causal scores produced by CHG by comparing them against targeted ablations as well as causal mediation analysis [12, 25], showing strong agreement between predicted and observed effects. Finally, we extend CHG to a contrastive setting to identify distinct sub-circuits that support instruction following versus in-context learning, suggesting that even semantically similar tasks can be underpinned by separable mechanisms.

Our main contributions are fourfold:

1. We introduce causal head gating (CHG), a parametric, scalable method for identifying potentially distributed, task-relevant sub-circuits in transformer models without requiring prompt templates or labeled outputs, and extend it with contrastive CHG to isolate heads supporting specific sub-tasks.
2. We propose a simple causal taxonomy of heads—facilitating, interfering, and irrelevant—that quantifies the effect of each on task performance using CHG-derived scores.
3. We use CHG to show that models contain multiple low-overlap, task-sufficient sub-circuits, suggesting head roles are not fully modular but depend on interactions with other heads.
4. We use CHG to show that instruction following and in-context learning rely on separable circuits at the head level, and that one can be selectively suppressed through CHG-guided head gating without disrupting the other.

## 2   Related Work

**Representational decoders**   Representational decoders are models trained to map hidden activations to externally labeled properties [7, 14, 15], estimating the mutual information between representations and those properties [27, 28]. However, such probing results are difficult to interpret: simpler decoders may underfit and miss relevant features (false negatives), while complex decoders may overfit and learn spurious correlations (false positives) [29, 30], requiring complexity-accuracy tradeoffs to contextualize results [30]. Moreover, although decodability indicates that a property is encoded in the representation, it does not imply that the model uses that information for its task—highlighting a correlational finding rather than a causal one [31]. Finally, representational decoders require labeled datasets, constraining their use to curated, predefined properties. For a comprehensive review of the probing framework and its limitations, see [27].

Sparse autoencoders can be viewed as a related approach, where the model reconstructs representations through a sparse bottleneck to reveal modular or interpretable features [16, 32]. However, like probing classifiers, their insights remain correlational and still depend on post hoc labeling or interpretation, inheriting the same supervision bottleneck.

**Causal mediation analysis**   Causal mediation analysis (CMA) [33, 34] is used to identify the functional roles of specific attention heads by crafting controlled prompt pairs that isolate a hypothesized behavior, then intervening on model components to measure their causal effect on outputs. For instance, in the indirect-object-identification (IOI) task [25], sentences like "When Alice and John went to the store, John gave a drink to..." are used to identify attention heads responsible for resolving coreference. By patching specific head outputs from a source sentence into a structurally matched target, and checking whether the model changes its prediction (e.g. "Alice" instead of "Mary"), CMA localizes the relevant circuit. It has also uncovered head-level roles in function tracking [12], symbol abstraction [13], and other structured settings [35].

However, CMA relies on manually crafted prompt templates and clear mechanistic hypotheses, which limits its scalability to more complex domains. In open-ended tasks like mathematical reasoning [36, 37, 38], the diversity of required knowledge makes it hard to design effective controlled inputs. A single shared template is unlikely to accommodate even two prompts from the MATH dataset [37], such as: "If $\sum_{n=0}^{\infty} \cos^{2n} \theta = 5$, what is $\cos 2\theta$?" and "The equation $x^2 + 2x = i$ has two complex solutions; determine the product of their real parts." Moreover, LLMs often solve such problems most effectively via chain-of-thought reasoning [26], which unfolds over multiple steps, further complicating the use of a unified prompt structure.

**Head ablations**   Despite the use of multiple heads being commonplace in transformer-based architectures, it has been observed that multiple, and sometimes the majority of, heads can be entirely pruned with minimal impact on model performance [39, 18, 40, 41]. Moreover, entire layers can be pruned while retaining model performance [42, 43, 44]. However, existing works on pruning attention heads have focused primarily on custom-trained small-scale transformers [39, 18, 40] or BERT-based [45] models [41, 43], and the literature is limited for modern causal LLMs such as GPT [46, 1] and Llama [3].

Head pruning has also been used to validate findings from other interpretability methods, such as CMA [25, 13] or attention pattern analysis [18]. In these studies, researchers first identify heads believed to perform specific functions, then ablate them to test their causal impact. Such targeted ablations often lead to disproportionate drops in performance, supporting the hypothesis that those heads are functionally important.

## 3   Our Approach: Causal Head Gating

Causal head gating is based on three ideas: applying multiplicative gates to attention heads to evaluate their roles, using regularization to produce variation in the estimates of the gating parameters, and constructing a taxonomy based on that variation. We introduce these ideas in turn.

Table 1: Causal taxonomy for head roles and corresponding gating patterns.

| Role | Description | $G^+$ | $G^-$ | Metric | Ablation Effect |
|------|-------------|-------|-------|--------|-----------------|
| Facilitating | Supports task performance | High | High | $G^-$ | Decreases task performance |
| Interfering | Interferes with task performance | Low | Low | $1 - G^+$ | Increases task performance |
| Irrelevant | Negligible impact on performance | High | Low | $G^+ \times (1 - G^-)$ | No effect on task performance |

## 3.1 Applying gates to attention heads

For a transformer with $L$ layers and $H$ attention heads, we define a gating matrix $G \in [0,1]^{L \times H}$, where $G_{\ell,h}$ scales the output of head $h$ in layer $\ell$, just before the output projection matrix $W_\ell^O$ (shown in red for an example head in Figure 1a). Given input hidden states $X \in \mathbb{R}^{\text{seq} \times d_{\text{model}}}$, each head computes:

$$A_{\ell,h} = \text{softmax}\left(\frac{XW_Q^{\ell,h}(XW_K^{\ell,h})^\top}{\sqrt{d_k}}\right), \quad V_{\ell,h} = XW_V^{\ell,h}, \quad Z_{\ell,h} = G_{\ell,h} \cdot (A_{\ell,h}V_{\ell,h})$$

where $W_Q^{\ell,h}, W_K^{\ell,h}, W_V^{\ell,h} \in \mathbb{R}^{d_{\text{model}} \times d_k}$ are learned projection matrices for queries, keys, and values.

The gating coefficient $G_{\ell,h}$ modulates the contribution of head $h$ by scaling its output $Z_{\ell,h}$ after attention is applied but before the heads are combined (see Figure 1a). The gated outputs are then concatenated and projected:

$$\text{Output}_\ell = \text{Concat}(Z_{\ell,1}, \ldots, Z_{\ell,H})W_\ell^O, \quad W_\ell^O \in \mathbb{R}^{H d_k \times d_{\text{model}}}$$

We fit $G$ by freezing the parameters of the model $\mathcal{M}_\theta$ and minimizing the negative log-likelihood (NLL) on a next-token prediction task with a regularization term specified below.

## 3.2 Producing variation through regularization

We add a regularization term to the objective that introduces a small but consistent gradient—clipped to ensure NLL remains the dominant term—that nudges the gates for task-irrelevant heads toward 1 or 0 while leaving task-relevant ones relatively unaffected. The NLL optimizes towards improving task performance, and tunes the heads by either increasing the gating values for task-facilitating heads or decreasing the gating values for task-interfering heads. However, if a head does not affect task performance, i.e. is task-irrelevant, then the expected gradient from the NLL is 0, which confounds interpretation of task relevance when evaluating the tuned gating values: a gate $G_{l,h}$ may be close to 1 either because it is important for performing the task (causal), or because gating it has no effect (incidental). We address this limitation by introducing an $L_1$-regularization term in our objective
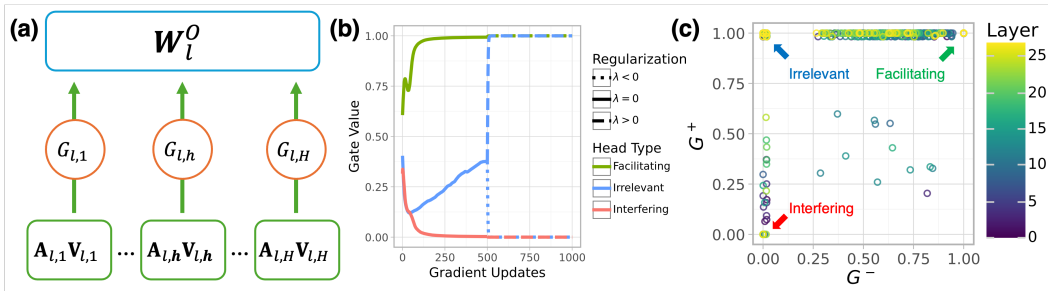


Figure 1: (a) Schematic of a single multihead attention block with CHG-determined gating attenuation (in red). (b) Gate fitting trajectories for three heads. When fitting with $\lambda < 0$ and $\lambda > 0$, $G^+$ and $G^-$ both stay near 1 for facilitating heads and near 0 for interfering heads, but bifurcate to 1 and 0 respectively for irrelevant heads. (c) Gate values after fitting.

function, with weight $\lambda$ that either nudges gates toward 1 for maximal density ($\lambda > 0$) or toward 0 for maximal sparsity ($\lambda < 0$):

$$\mathcal{L}(G; \mathcal{M}_\theta, \mathcal{D}, \lambda) = \underbrace{- \sum_{(x,y)\in\mathcal{D}} \log P(y \mid x; \mathcal{M}_\theta, G)}_{\text{Negative log-likelihood (NLL)}} \underbrace{- \lambda \sum_{i,j} \sigma^{-1}(G_{l,h})}_{\text{Regularization}} \tag{1}$$

where $\mathcal{M}_\theta$ is the model being analyzed, $y$ is the target text sequence for a given prompt $x$ in dataset $\mathcal{D}$, and $\sigma^{-1}$ is the clipped inverse-sigmoid function.

We fit $G$ twice: once with $\lambda > 0$ to encourage retention ($G^+$), and once with $\lambda < 0$ to encourage removal ($G^-$). To ensure that the heads are aligned across both optimizations, we first fit $G$ with $\lambda = 0$ to establish a shared initialization (see Figure 1), so that any differences between $G^+$ and $G^-$ reflect only the effect of the regularization and not divergent optimization paths.

### 3.3 Constructing a taxonomy of task relevance

The $G^+$ and $G^-$ matrices allow us to interpret the functional role of each head. To formalize this, we introduce a causal taxonomy (Table 1) in which each head is assigned one of three roles—*facilitating*, *interfering*, or *irrelevant*—based on its predicted impact on model performance under ablation. Facilitating heads positively contribute to performance, while ablating them degrades it. Conversely, interfering heads negatively contribute to performance, while ablating them improves it. Finally, irrelevant heads have negligible effect, with ablation leaving performance effectively unchanged.

We instantiate this taxonomy using the fitted CHG matrices $G^+$ and $G^-$, which reflect head behavior under opposing regularization pressures. Facilitation is measured by $G^-$: heads that remain active despite pressure to suppress are likely necessary for the task. Interference is measured by $1 - G^+$: heads that are suppressed even under encouragement to remain are likely harmful. Irrelevance is measured via $G^- \odot (1 - G^+)$, identifying heads that vary in gate values based on regularization.

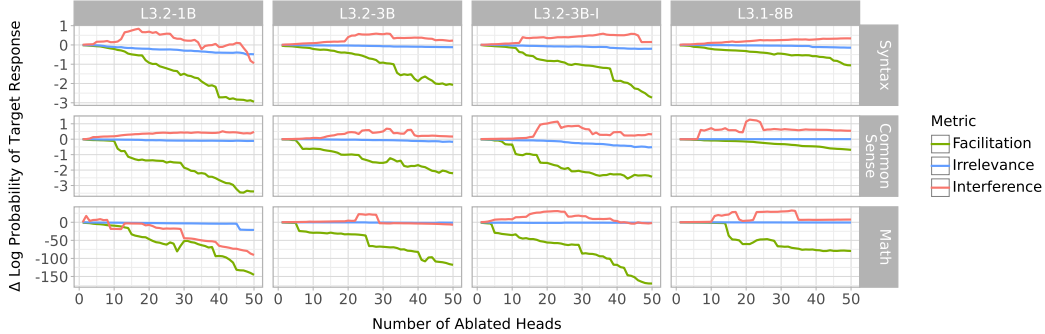## 4 Experiments and analyses



Figure 2: Difference in target log-probability when sequentially setting individual gates in $G^+$ to 1 and 0 in order of facilitation, irrelevance, and interference scores. The horizontal axis shows the number of heads ablated in descending score order. Positive values indicate task improvement, negative values indicate degradation, and values near zero indicate no effect. Note that not all heads in the top 50 necessarily have high absolute scores.

### 4.1 Causal roles of attention heads

We begin by reporting experiments that evaluated the causal taxonomy presented in Table 1 across four variants of the Llama 3 LLM [3]: L3.1-8B, a pre-trained 8B-parameter model; L3.2-3B, a 3B-parameter model distilled from Llama-3.1-70B (not used in this paper); L3.2-3BI, an instruction-tuned version of Llama-3.2-3B; and L3.2-1B, a 1B-parameter model distilled from L3.1-8B. For each model, we fit CHG matrices on three task types performed over distinct datasets: mathematical reasoning from OpenMathInstruct2 [38], syntactic reasoning from the subset labeled "syntax" in

BIG-Bench [47], and commonsense reasoning from CommonsenseQA [48]. We fit CHG matrices independently for each model-dataset pair across 10 random seeds.

We first test whether the causal scores align with the taxonomy's predictions about performance. Specifically, the taxonomy predicts that, when ablated, attention heads scoring highly on facilitation, irrelevance, or interference should decrease, leave unchanged, or increase the model's task performance, respectively. To test this, we sort heads in descending order by each causal metric and evaluate the model using the $G^+$ matrix while toggling each head to 0 or 1 in order of its score. While both $G^+$ and $G^-$ match the context in which scores were computed, we use $G^+$ as it retains more heads, providing a more interpretable baseline for ablation. We then compare the retained and ablated masks by the model's log-probability of the target sequence, expecting the resulting change in log-probability to follow the predicted pattern. As shown in Figure 2, these interventions match the predicted patterns: the difference in target log-probability is negative when progressively ablating facilitating heads, near 0 when ablating irrelevant heads, and positive when ablating interfering heads, up until the set of interfering heads is exhausted.
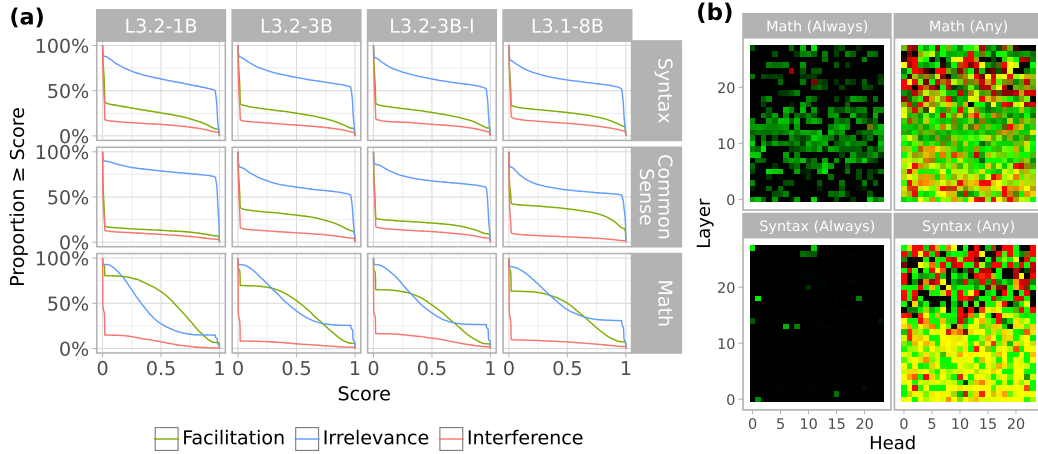
## 4.2 Distribution of causal roles



Figure 3: CHG score distributions and consistency. (a) Empirical cumulative distribution of CHG scores across all attention heads, showing the proportion of heads with scores below a given threshold for facilitation, irrelevance, and interference. (b) Aggregated CHG scores on L3.2-3BI, where red and green color channels represent interference $(1 - G^+)$ and facilitation $(G^-)$, respectively. Colors are combined using RGB rules: black indicates irrelevance (low in both), and yellow indicates both facilitation and interference (high in both). *Always* aggregates using the minimum across seeds (highlighting consistent effects); *Any* uses the maximum (highlighting any effect across seeds).

Having validated the causal scores using targeted ablations, we next analyze how they are distributed across models and tasks. Figure 3a shows that for each task, the distribution of head roles is highly consistent across all four model variants. This holds despite large differences in model size (1B to 8b) and training setup (pretraining, distillation, instruction tuning). We quantify these similarities by computing Pearson correlations between head scores across all model pairs for each task and causal metric, yielding 54 model pairs, all of which show high agreement with a minimum correlation of 94.92% and an average of 99.2%. Across tasks, however, we observe notable differences, with the math dataset standing out in particular. For syntax and commonsense reasoning, most heads are irrelevant—63.0% and 64.6% have irrelevance scores $\geq 0.5$, respectively—with only a sparse subset of facilitating heads (25.6% and 27.4% with facilitation scores $\geq 0.5$), suggesting that compact, redundant circuits are sufficient for these tasks. In contrast, mathematical reasoning activates a much larger fraction of facilitating heads: 52.6% have facilitation scores $\geq 0.5$, while only 39.0% are irrelevant, likely reflecting the task's higher complexity and need for broader sub-circuitry to support multi-step, latent computations.

It is also worth noting that, across all tasks, 84.0% of heads are marked as facilitating or interfering (score $\geq 0.5$) in at least one seed, yet only a small fraction are consistently facilitating or interfering

Table 2: Percent of heads with facilitation (F) or interference (N) scores $\geq 0.5$ across all seeds (always) or in at least one seed (any).

| Task | Agg. | L3.2-1B | | L3.2-3B | | L3.2-3BI | | L3.1-8B | |
|---|---|---|---|---|---|---|---|---|---|
| | | F | N | F | N | F | N | F | N |
| Syntax | Always | 1.2 | 0.2 | 1.5 | 0.1 | 0.7 | 0.0 | 1.4 | 0.0 |
| | Any | 72.1 | 57.2 | 67.9 | 51.3 | 72.8 | 56.1 | 68.5 | 59.2 |
| Common Sense | Always | 3.9 | 0.0 | 4.5 | 0.0 | 3.0 | 0.0 | 18.7 | 0.6 |
| | Any | 56.6 | 41.0 | 75.4 | 52.4 | 68.2 | 55.7 | 60.3 | 22.2 |
| Math | Always | 38.3 | 0.4 | 24.6 | 1.3 | 18.3 | 0.1 | 25.3 | 1.0 |
| | Any | 81.1 | 26.0 | 75.1 | 13.8 | 74.4 | 47.2 | 75.0 | 21.2 |

across all seeds (Figure 3b). In syntax and commonsense tasks, most models have fewer than 5% of heads that are always facilitating and virtually none that are always interfering (Table 2). In contrast, math reveals more rigid and consistent circuitry, with up to 38.3% of heads consistently facilitating and 1.3% consistently interfering. These patterns suggest that individual attention heads may not have modular, context-independent roles, but instead participate in a flexible ensemble of overlapping sub-circuits, in which their function depends on the configuration of others.

## 4.3 Comparison with causal mediation analysis

CMA, like CHG, aims to identify attention heads that facilitate task execution, though it does so in a more hypothesis-driven manner. Framed in signal detection terms, CMA and CHG are complementary. CMA exhibits high precision but relatively low sensitivity: while many facilitating heads may go undetected (false negatives), those it does identify are reliably task-relevant (few false positives) Conversely, CHG is biased toward sensitivity over precision. This suggests that that heads identified by CMA should also be identified (as showing strong facilitation) under CHG. We test this by comparing CHG to the results of two former studies using CMA, replicating their methods to identify attention heads with specific computations: heads that encode task information in function vectors [12] and heads that perform symbolic reasoning [13].

For function vectors, we use the six in-context learning tasks used in [12]: 'antonym', 'capitalize', 'country-capital', 'English-French', 'present-past', and 'singular-plural'. Each prompt is presented in an in-context learning (ICL) [46] format consisting of 10 input-output examples using a "Q: X\n A: Y" template, followed by a query to be answered. To perform CMA, we corrupt the prompt by randomly shuffling example outputs to induce mismatched pairs, then patch individual head outputs with clean activations to identify which heads recover performance—interpreting high recovery as evidence of causal mediation.

We apply a similar logic to symbolic reasoning tasks from [13], where the goal is to generalize abstract identity rules such as ABA ("flow^Started^flow") or ABB ("flow^Started^Started"). We deploy the same CMA procedure used in [13] to identify the three-stage symbolic processing mechanism that was reported: (1) *symbol abstraction* heads that abstract symbols ("A" or "B") away from the actual tokens in the in-context examples; (2) *symbolic induction* heads that operate over the abstracted symbols to induce the symbol for the missing token in the query; (3) *retrieval* heads that retrieve the actual token based on the induced symbol to complete the query. To screen heads of each type, we construct prompt pairs in which either the same token is assigned to different symbols ("A" or "B") or tokens are swapped while preserving the same rule, and patch activations at certain token positions between them. Attention heads that steer model behavior towards specific hypotheses about the three head types after patching (either converting the abstract rule or altering the actual token) are labeled as mediating. We conduct all experiments on the Llama-3.2-3B-Instruct model.

As predicted, CMA-identified heads tend to exhibit high facilitation scores under CHG in both domains (Figure 4). To quantify this, we compare the CHG facilitation scores of CMA-identified heads—those with three standard deviations above the mean in function vector tasks or with statistical significance in ABA/ABB tasks [13]—to the remaining ones. Since facilitation and irrelevance depend on the specific sufficient circuit identified by CHG, a head may appear irrelevant in one run but facilitating in another if multiple circuits exist. To account for this, we fit 10 CHG masks per function vector task and 20 per ABA/ABB task, and compute each head's maximum facilitation score
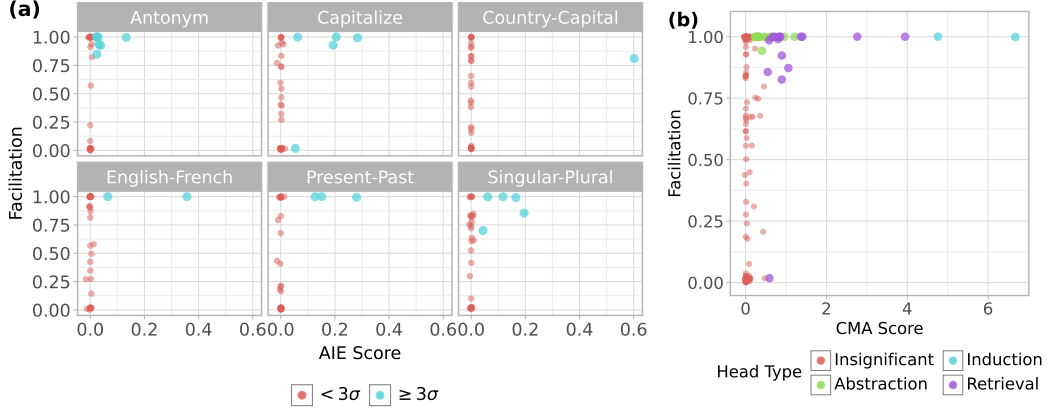
Figure 4: Task-facilitation scores versus (a) average indirect effect for function vector tasks and (b) CMA scores for symbolic reasoning tasks, showing significant heads by type (*abstraction*, *induction*, *retrieval*) and using the maximum CMA score across types for insignificant heads.

across runs—capturing whether it participates in any sufficient circuit. We find significantly greater facilitation among mediating heads in both the function vector tasks ($t(23.05) = 8.52$, $p < 10^{-8}$) and the ABA/ABB tasks ($t(53.77) = 11.18$, $p < 10^{-15}$), supporting the relationship between CMA and CHG-identified task relevance.

## 4.4 Contrastive Causal Head Gating

The results above indicate that CHG effectively distinguishes among facilitating, irrelevant, and interfering attention heads. However, as an exploratory method, it lacks the granularity to characterize the specific functions of these subnetworks. For instance, consider the 'antonym' task from Section 4.3, presented in an in-context learning (ICL) format with 10 examples and a single-word response, as defined in [12]. To perform this task successfully, the model must not only generate the appropriate antonym of a given word, but also infer the task itself from the 10 input-output pairs in the prompt. Thus, a minimal circuit of task-facilitating heads will contain both those involved in task inference and those involved in antonym production, and CHG cannot distinguish between the two. This becomes more pronounced as task complexity increases, as in the OpenMathInstruct2 dataset, where the minimal circuit must jointly support diverse sub-tasks, including English comprehension, mathematical reasoning, chain-of-thought processing, and LaTeX generation.

To address this, we introduce a simple extension of CHG that not only identifies facilitating heads for a given task but also isolates the sub-circuit responsible for a particular sub-task. We generate parallel variants of the same task that share all features except for a controlled difference in the required operation, allowing us to isolate the corresponding sub-circuits. In doing so, we take a step toward a hypothesis-driven approach, decomposing the task into sub-steps while remaining agnostic to the mechanistic implementations For example, the antonym task can be constructed as an ICL task using the default format from [12], or as an instruction-following task where the model is presented with the task description "Given an input word, generate the word with opposite meaning". By comparing the resulting attention circuits, we can disentangle components responsible for task inference from those involved in antonym generation.

Furthermore, rather than simply applying CHG to each version and directly comparing the results, we propose a combined approach that fits a single mask with a joint objective to forget one variant of the task while retaining the other, so that the resulting gate matrix suppresses heads uniquely necessary for one variant but dispensable for the other:

$$\mathcal{L}_{\text{contrastive}}(G; \mathcal{M}_\theta, \mathcal{D}_R, \mathcal{D}_F, \lambda) = \mathcal{L}(G; \mathcal{M}_\theta, \mathcal{D}_R, \lambda) - \mathcal{L}(G; \mathcal{M}_\theta, \mathcal{D}_F, \lambda) \qquad (2)$$

where $\mathcal{D}_R$ is the dataset to retain, $\mathcal{D}_F$ is the dataset to forget, $\lambda < 0$, and $\mathcal{L}$ is the original CHG loss defined in Equation 1.

We evaluate this method using the six function vector tasks from Section 4.3, leveraging the natural language task descriptions provided in [12] to construct instruction-based variants. For each problem,

we replace the 10-shot word-pair examples with a prompt containing the task instruction and a single example. We then fit the *contrastive* causal head gating (CCHG) mask to forget the ICL variant of five tasks while retaining the instruction-based format, holding out the sixth task for evaluation. If task inference from examples, instruction-following, and task execution are indeed mediated by separable circuits, this analysis should disable example-based generalization while preserving instruction-based performance. We perform our experiments in both directions (forgetting ICL while retaining instruction-following, and vice versa), using each of the six tasks as the held-out evaluation task. All experiments were conducted on the LLaMA-3.2-3B-Instruct model.
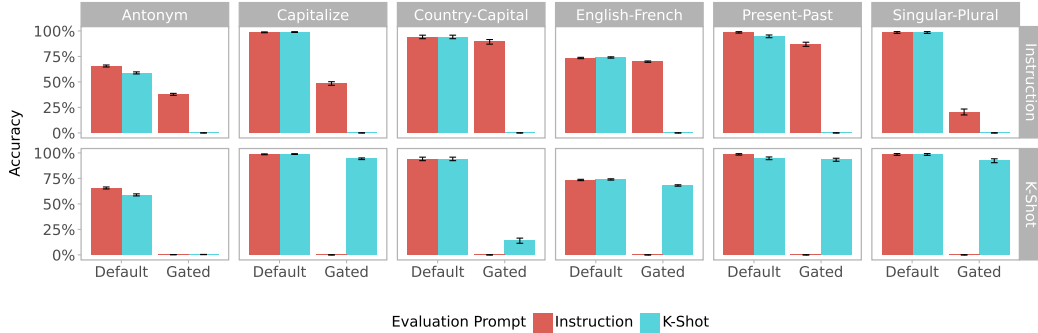


Figure 5: Task accuracy under CCHG. Columns indicate held-out evaluation tasks and rows indicate the retained prompt format. Bar color shows the evaluation prompt format. "Default" and "gated" indicate whether CCHG is applied during evaluation. Error bars indicate 95% CI.

As shown in Figure 5, the CCHG masks generalize to the held-out task. When the model is induced to forget task inference from ICL examples across five tasks, its target task accuracy drops to zero on the ICL variant of the held-out task while in most cases remaining well above zero—and often close to the unablated baseline—on the instruction-based variant. A similar pattern emerges when forgetting is applied using the instruction-based format: performance collapses on instruction prompts while generally remaining intact for example-based ones.

Interestingly, while degradation is often small for the retained prompt format, this pattern is not consistent across all tasks. For example, when the gating matrix is fitted to retain ICL and forget instruction-following, the 'singular-plural' task shows only a small drop in ICL accuracy ($98\% \rightarrow 92\%$) but a complete failure on instruction prompts ($98\% \rightarrow 0\%$). When this setup is reversed—fitted to retain instruction-following and forget ICL—accuracy on ICL drops from 98% to 0%, while instruction accuracy drops more modestly ($98\% \rightarrow 21\%$). Across the 6 tasks, 3 ('country-capital', 'English-French', 'present-past') remain robust as held-out tasks under instruction prompts, and 4 ('capitalize', 'English-French', 'present-past', 'singular-plural') do so under ICL prompts.

Thus, our results indicate that the circuits for instruction following and ICL may be separable at the head level. However, this separability also depends on the task, suggesting that task execution circuits may share heads with those used for task understanding and representation.

# 5 Discussion

In this work, we introduced Causal Head Gating (CHG), a flexible and scalable method for identifying causally relevant attention heads in large language models. CHG assigns each head a graded score for facilitation, interference, or irrelevance based on its causal effect on task performance, going beyond correlational or observational analyses. Crucially, it does so using next-token prediction alone, thereby avoiding reliance on labeled data or handcrafted prompts, making it broadly and easily applicable. Moreover, CHG requires no finetuning or auxiliary decoder model, and introduces only one parameter per head, allowing it to run in minutes even on billion-scale models.

We validated CHG by showing that its scores predict performance changes under targeted ablations, confirming that facilitation, interference, and irrelevance scores capture causal impact. Heads identified by CMA also receive high CHG facilitation scores.

9

Using CHG, we analyzed a range of models and tasks, revealing that attention heads form task-sufficient sub-circuits with low overlap. Interestingly, we found that head roles are not fixed: the same head may facilitate or interfere depending on which other heads are active, underscoring the distributed and context-dependent nature of computation in LLMs. Lastly, we extended CHG to a contrastive setting, showing that heads supporting instruction following and in-context learning can be selectively suppressed, indicating separable functional mechanisms.

Despite its strengths, CHG has key limitations. First, it identifies causally relevant heads but not why they matter or what they compute. Even contrastive CHG offers limited insight into underlying mechanisms, making it best suited for exploratory analysis alongside hypothesis-driven tools like CMA. Additionally, $G^+$ and $G^-$ do, albeit rarely, diverge, possibly reflecting richer head interactions. Finally, head roles can vary across runs, so multiple CHG fits may be needed.

We hope that our work encourages further exploration of causal structure in language models as a foundation for more mechanistic understanding. Future work may build on these tools to develop circuit-level explanations of how models implement complex behaviors.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[2] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

[3] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

[5] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does LLM safety training fail? *Advances in Neural Information Processing Systems*, 36:80079–80110, 2023.

[6] Laura Weidinger, Joslyn Barnhart, Jenny Brennan, Christina Butterfield, Susie Young, Will Hawkins, Lisa Anne Hendricks, Ramona Comanescu, Oscar Chang, Mikel Rodriguez, et al. Holistic safety and responsibility evaluations of advanced ai models. *arXiv preprint arXiv:2404.14068*, 2024.

[7] Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. *arXiv preprint arXiv:1905.05950*, 2019.

[8] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2, 2023.

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

[10] Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.

[11] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformer-circuits.pub/2021/framework/index.html.

[12] Eric Todd, Millicent L Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. Function vectors in large language models. *arXiv preprint arXiv:2310.15213*, 2023.

[13] Yukang Yang, Declan Campbell, Kaixuan Huang, Mengdi Wang, Jonathan Cohen, and Taylor Webb. Emergent symbolic mechanisms support abstract reasoning in large language models. *arXiv preprint arXiv:2502.20332*, 2025.

[14] Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R Bowman, Dipanjan Das, et al. What do you learn from context? Probing for sentence structure in contextualized word representations. *arXiv preprint arXiv:1905.06316*, 2019.

[15] John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.

[16] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from Claude 3 Sonnet. *Transformer Circuits Thread*, 2024.

[17] Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K Kummerfeld, and Rada Mihalcea. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. *arXiv preprint arXiv:2401.01967*, 2024.

[18] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint arXiv:1905.09418*, 2019.

[19] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.

[20] Geoffrey E Hinton. Learning distributed representations of concepts. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 8, 1986.

[21] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. `https://transformer-circuits.pub/2022/toy_model/index.html`.

[22] Kayson Fakhar and Claus C Hilgetag. Systematic perturbation of an artificial neural network: A step towards quantifying causal contributions in the brain. *PLOS Computational Biology*, 18(6):e1010250, 2022.

[23] Tyler Giallanza, Declan Campbell, Jonathan D Cohen, and Timothy T Rogers. An integrated model of semantics and control. *Psychological Review*, 2024.

[24] Melanie Mitchell. *Complexity: A guided tour*. Oxford University Press, 2009.

[25] 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. *arXiv preprint arXiv:2211.00593*, 2022.

[26] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

[27] Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022.

[28] Tiago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell. Information-theoretic probing for linguistic structure. *arXiv preprint arXiv:2004.03061*, 2020.

[29] John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. *arXiv preprint arXiv:1909.03368*, 2019.

[30] Elena Voita and Ivan Titov. Information-theoretic probing with minimum description length. *arXiv preprint arXiv:2003.12298*, 2020.

[31] Abhilasha Ravichander, Yonatan Belinkov, and Eduard Hovy. Probing the probing paradigm: Does probing accuracy entail task relevance? *arXiv preprint arXiv:2005.00719*, 2020.

[32] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.

[33] 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. *Advances in Neural Information Processing Systems*, 33:12388–12401, 2020.

[34] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.

[35] Zifan Zheng, Yezhaohui Wang, Yuxin Huang, Shichao Song, Mingchuan Yang, Bo Tang, Feiyu Xiong, and Zhiyu Li. Attention heads of large language models. *Patterns*.

[36] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

[37] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.

[38] Shubham Toshniwal, Wei Du, Ivan Moshkov, Branislav Kisacanin, Alexan Ayrapetyan, and Igor Gitman. Openmathinstruct-2: Accelerating ai for math with massive open-source instruction data. *arXiv preprint arXiv:2410.01560*, 2024.

[39] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in Neural Information Processing Systems*, 32, 2019.

[40] Jiaoda Li, Ryan Cotterell, and Mrinmaya Sachan. Differentiable subset pruning of transformer heads. *Transactions of the Association for Computational Linguistics*, 9:1442–1459, 2021.

[41] Mengzhou Xia, Zexuan Zhong, and Danqi Chen. Structured pruning learns compact and accurate models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1513–1528, 2022.

[42] Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured dropout. *arXiv preprint arXiv:1909.11556*, 2019.

[43] Hassan Sajjad, Fahim Dalvi, Nadir Durrani, and Preslav Nakov. On the effect of dropping layers of pre-trained transformer models. *Computer Speech & Language*, 77:101429, 2023.

[44] Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li. What matters in transformers? not all attention is needed. *arXiv preprint arXiv:2406.15786*, 2024.

[45] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.

[46] 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*, 33:1877–1901, 2020.

[47] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022.

[48] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *arXiv preprint arXiv:1811.00937*, 2018.