

How do Language Models Encode Privacy Bias?

Neel Sanjaybhai Faganiya
University of Waterloo
Waterloo, Ontario, Canada

Sajad Rahmanian Ashkezari
University of Waterloo
Waterloo, Ontario, Canada

Lucas Kopp
University of Waterloo
Waterloo, Ontario, Canada

Abstract

Large Language Models (LLMs) increasingly mediate socio-technical systems where privacy judgments are critical, yet they often encode biased privacy norms learned from internet-scale training data. Prior work has predominantly focused on detecting behavioral privacy biases without understanding their mechanistic origins within the model weights. This paper addresses this gap by investigating whether privacy biases are localized in specific circuits of LLMs using the Contextual Integrity (CI) framework.

We use this methodology of combining CI with Mechanistic Interpretability (MI) techniques to identify and analyze these circuits. Our approach constructs controlled vignette pairs that isolate key CI parameters and employs Edge Attribution Patching with Integrated Gradients (EAP-IG) on instruction-tuned LLMs to discover faithful circuits influencing privacy-related decisions. Results reveal specialized mid-to-late layer attention patterns with high fidelity that differentiate appropriate from inappropriate information flows, with low structural overlap indicating modular privacy mechanisms.

Our work bridges behavioral privacy evaluation with internal model interpretability, advancing tools for targeted circuit editing to mitigate privacy biases without costly full model retraining. Our findings provide actionable insights for developing privacy-respecting AI systems.

1 Introduction

During training, Large Language Models (LLMs) can inadvertently learn societal biases and privacy norms embedded in large-scale, non-public datasets. These learned biases may not always align with socially and ethically accepted privacy expectations, sometimes leading to inappropriate disclosures of sensitive information [1, 17, 29]. As LLMs are deployed as interactive agents in socio-technical systems, it becomes imperative to rigorously examine these privacy biases, which is defined as a skewed judgment of the appropriateness of information flows within specific social contexts [24].

Unchecked privacy biases can result in serious ethical, legal, and reputational risks. They may lead to unfair treatment of individuals, violations of data protection laws such as GDPR and HIPAA, and erosion of user trust in AI systems. Moreover, biased privacy judgments can amplify systemic inequalities affecting marginalized groups, potentially causing real-world harms in areas like healthcare, finance, or employment decisions [4, 16]. Therefore, mitigating privacy bias is essential to ensure LLMs uphold privacy standards,

comply with regulations, and foster equitable, trustworthy AI applications.

We ground the notion of privacy bias in the theory of Contextual Integrity, which conceptualizes privacy in terms of appropriate information flows governed by the roles, relationships, datatypes, and transmission principles relevant to a given context [20]. Using this framework, we formulate the novel research problem of localizing privacy biases embedded within LLMs’ internal representations. This approach enables us to move beyond evaluating privacy at a behavioral level, to understanding the mechanistic origins of these biases within the model’s computational circuits.

Evaluating privacy biases is challenged by the high sensitivity of LLM responses to subtle variations in prompt phrasing, which can confound bias assessment if not carefully managed [2, 6, 22, 24]. We address this by developing a robust methodology that identifies consistent privacy judgments through carefully constructed vignette pairs and prompt ensembles.

Our contributions include:

- Generating a dataset to find circuits related to privacy bias.
- Developing Clean and Corrupted prompting strategies for running EAP-IG.
- Finding and Analyzing circuits related to privacy decisions using Mechanistic Interpretability techniques.

2 Background

Contextual Integrity (CI) is a theoretical framework for understanding privacy as the appropriateness of information flows within specific social contexts, governed by norms [20]. These norms are characterized by five key parameters:

- **Roles:** the actors involved, which includes the information sender, recipient and data subject (e.g., *doctor, patient, insurer*).
- **Information Type:** the kind of data being shared, such as *medical history, location, or financial details*.
- **Transmission Principle:** the conditions or constraints under which information is transmitted, e.g., *consent, confidentiality, or legal requirements*.

For example, consider a typical information flow in a healthcare context: **Doctor** (sender) shares **patient’s** (data subject) **medical information** (infotype) with **pharmacist** (recipient) to **prepare prescription** (transmission principle). This flow is aligned to healthcare privacy norms and thus preserves contextual integrity and is an appropriate information

flow. However, if the same medical data were shared with an advertisement agency for marketing purposes, the transmission principle would be violated, constituting a breach of contextual integrity and would be considered an inappropriate information flow [24].

CI has been effectively applied to evaluate privacy expectations in diverse contexts and has recently been adapted to assess privacy compliance and biases in large language models (Section 3.1)[7, 13, 17, 19], yet it lacks insights into the model’s internal mechanisms that drive these judgments.

Privacy Bias in LLMs refers to skewed judgments on the appropriateness of information flows within specific social contexts [24]. That is, the extent to which an LLM deems an information flow to be appropriate or inappropriate relative to established privacy norms learned from training data. For instance, LLMs trained on internet data containing implicit privacy rules may incorrectly suggest sharing patient’s medical records with advertisers or publicly revealing political beliefs as appropriate, reflecting systemic biases in training data.

Mechanistic Interpretability. Neural networks (NNs) are functions whose input is an N -dimensional vector. Thus, it is difficult to predict how they are going to behave on this exponential space. Similar problem arises for many programs that we run on our computers. For these programs, we can take a look at their code and understand their behavior. Perhaps, we could do the same by looking at NNs’ parameters. However, NNs could have tens-hundreds of billions of parameters. The issue here, same as looking at binary representation of a program, is that the parameters themselves are uninterpretable for humans. The goal of mechanistic interpretability is to find smaller descriptions of NNs similar to codes for a binary program [21].

3 State-of-the-Art

3.1 Contextual Integrity in Privacy Assessment

Recent works have advanced the application of Contextual Integrity theory to evaluate privacy in LLMs. Mireshghallah et al. [17] adapts CI through vignette-style probing to compare LLM privacy judgments with human expectations, focusing on measuring adherence to social norms for information sharing without dissecting internal model decision-making. Shao et al. [23] utilizes CI parameters to enable multi-level evaluation of privacy leakage scenarios in Language Model (LM) agents’ actions, aiming to bridge the gap between LM performance in answering the probing questions and the actual behavior when executing client’s instruction in a real-world agentic setup. Yi et al. [27] expands on CI by addressing Contextual Ambiguity, generating disambiguated scenarios that improve prompt stability and privacy judgment accuracy by simulating real-world clarifications.

Shvartzshnaider and Duddu [24] apply CI to quantify privacy bias as statistical divergence from contextual norms via multi-prompt ensembles.

3.2 LLM Privacy Bias Detection

Recent studies, such as Mireshghallah et al. [17], reveal that LLM privacy judgments often diverge from human expectations, highlighting significant alignment gaps. Research by Cao et al. [2] demonstrates LLMs’ high sensitivity to prompt templates and a lack of robustness in privacy classification tasks. Investigation into privacy risks, including membership inference attacks, are detailed by Meeus et al. [15], while Carlini et al. [3] expose LLMs’ vulnerability to data reconstruction via training data extraction. However, these works mostly evaluate the models from an external perspective, without pinpointing the internal components behind biased privacy behavior. Identifying such components is a critical step towards effective mitigation.

3.3 Mechanistic Interpretability Techniques

Mechanistic Interpretability methods like Edge Attribute Patching (EAP) by Nanda [18] aim to identify model-internal circuits responsible for specific behaviors by attributing gradients along computational edges. Recent extensions, like EAP-IG by Hanna et al. [10] and EAP-GP by Zhang et al. [28], provide more precise circuit localization and enable lower-impact interventions through improved gradient-based attribution. While these methods have primarily been applied to domains such as factual knowledge editing and bias mitigation, they have rarely been explored in the context of privacy-related behavior. In this work, we employ EAP-IG for its refined attribution capabilities, using it to locate privacy-relevant circuits and then patch those. This bridges the gap between behavioral privacy evaluation and transparent, circuit-level mitigation.

3.4 Existing Privacy Bias Mitigations

Existing mitigation techniques primarily rely on model fine-tuning, which often leads to substantial performance degradation [8]. Other methods like Retrieval-Augmented Generation (RAG) improves context awareness but faces challenges with computational cost and scalability [7]. CI based Chain-of-Thought (CI-COT) prompting helps but it does not alter the underlying model weights to mitigate the issue of privacy bias [12]. These gaps underscore the necessity of interpretable, efficient mechanisms for mitigating privacy bias at the circuit, which this work aims to address.

4 Problem Statement

Recent work has proposed several strategies to mitigate privacy biases in Large Language Models. Fine-Tuning on curated privacy or CI-aligned datasets can shift model behavior

toward safer information-sharing decisions, but it is computationally expensive, requires substantial human annotation, and often degrades unrelated capabilities or leads to over-refusal outside the training distribution [8]. RAG-style systems and external privacy guardrails instead attack policy or CI knowledge bases and use retrieval or secondary models at inference time, reducing direct leakage but introducing latency and infrastructure complexity while leaving the underlying biased model unchanged [23]. CI-based prompting and CI-Reinforcement Learning methods improve adherence to contextual integrity norms at the behavioral level, yet they still operate as black-box adjustments to model outputs and do not reveal or directly modify the internal mechanisms that give rise to privacy biases [12, 17].

At the same time, emerging CI benchmarks show that LLM privacy judgments can diverge substantially from human expectations across contexts, and that these misalignments can be prompt-sensitive and unstable [17, 23, 24]. Because LLMs are increasingly deployed as autonomous agents that act on users’ behalf, such biased or brittle privacy decisions risk inappropriate disclosure or over-withholding of information at scale.

What is missing is a method that (1) **identifies** the internal circuits responsible for privacy decisions under contextual integrity (CI), and (2) **edits** those circuits in a localized way to reduce privacy bias while preserving the model’s broader capabilities. In this work, the goal is to use mechanistic interpretability - specifically, Edge Attribution Patching with Integrated Gradients (EAP-IG) [10] - on ConfAIde, a CI-grounded dataset [17], to discover circuits that mediate privacy judgments. Then, to perform targeted interventions on these circuits to correct privacy bias by patching those circuits without needing to retrain the full model.

Research Questions:

RQ1: Which model components are most important for contextual privacy decisions? We seek to identify the specific parts of the model that are most responsible for judging whether an information flow is appropriate or inappropriate on a CI-grounded benchmark, ConfAIde.

RQ2: To what extent do privacy-relevant components overlap with other behaviors? Here we ask whether components involved in privacy-related behavior rely on shared underlying mechanisms or operate through distinct circuitry. Understanding this overlap is essential for determining how entangled privacy reasoning is with other model functions and whether targeted interventions can be applied safely.

RQ3: How can targeted circuit patching improve the privacy-utility trade-off? Our initial goal was to evaluate whether targeted circuit patching could reduce privacy-related failures, while maintaining or improving task utility. This approach aims to provide more precise control than

broad behavioral methods such as fine-tuning or RAG. However due to time constraint and substantial effort required for RQ1 and RQ2, we were unable to implement this for our course project.

5 Approach

In this section we go over our methodology in detail. We start by describing our experimental setup, then we give a brief overview of EAP-IG, which is the main tool we use for our project. We then explain how we use this method for finding circuits relevant to privacy bias.

5.1 Experimental Setup

Dataset: We construct our evaluation scenarios from the ConfAIde benchmark by Miresghallah et al. [17], which provides contextual-integrity information flows parameterized by **information type**, **recipient**, and **use**. Each scenario instantiates a template of the form: "*Information about {infotype} is collected by {recipient} in order to {use}*".

To elicit model privacy judgments, we ask LLMs to rate how acceptable each scenario is on a 5-point Likert scale ranging from *Strongly Inappropriate* to *Strongly Appropriate*. A prompt example:

Please indicate the extent to which the following scenario is acceptable: {scenario}. The answer must be exactly one of the following options (and nothing else): {Likert options}.

We frame the prompt this way to avoid getting a larger nonsensical response from the LLM. These ratings form the basis for the binary appropriateness labels used in our circuit analysis.

Models: All experiments are conducted on two Instruction-tuned models: **Qwen2.5-7B-Instruct**, and **Llama3.1-8B-Instruct**. We treat both models as black boxes for generation, but later apply EAP-IG to their internal activations for circuit discovery.

5.2 Prompt Paraphrasing to Reduce Wording Bias

LLMs are known to be highly sensitive to prompt wording; logically equivalent questions can yield different answers, especially on nuanced safety and privacy judgments [2, 6, 9, 22, 24]. To reduce this accidental wording bias, we paraphrase each scenario prompt into 11 distinct phrasings (original + 10 paraphrased prompts), varying surface structures (e.g., "Please rate how acceptable this scenario is:", "How acceptable do you find this scenario:", "To what degree is this scenario acceptable:") while preserving the underlying CI information flow [24]. Using multiple paraphrases helps smoothing out wording-specific quirks and gives a more stable and true LLM judgment for each scenario.

5.3 Aggregating Likert Responses into Stable Model Judgments

For each scenario, model, and paraphrased question, we collect a single 5-point likert response. We then map the five options into a binary label (*Inappropriate* vs. *Appropriate*) and aggregate across the 11 paraphrases as follows:

- For each model and scenario, count how many of the 11 responses are labeled **appropriate** and how many **inappropriate**
- If one label appears in more than 66% (i.e., at least 8 out of 11) of the paraphrases, we treat that label as the model’s **true judgment** for that scenario [24]

- Scenarios that do not meet this threshold are considered unstable and are excluded from circuit analysis (grey cells in Figure 1)

This supermajority threshold rule filters out prompts where the model is inconsistent or overly sensitive to wording, leaving a set of stable appropriate/inappropriate judgments (minimizing the survivorship bias) that we use to define clean examples and construct clean-corrupted pairs for EAP-IG [24]. To better understand the global structure of these judgments, we next visualize the aggregated labels over all information-flow combinations as a heatmap.

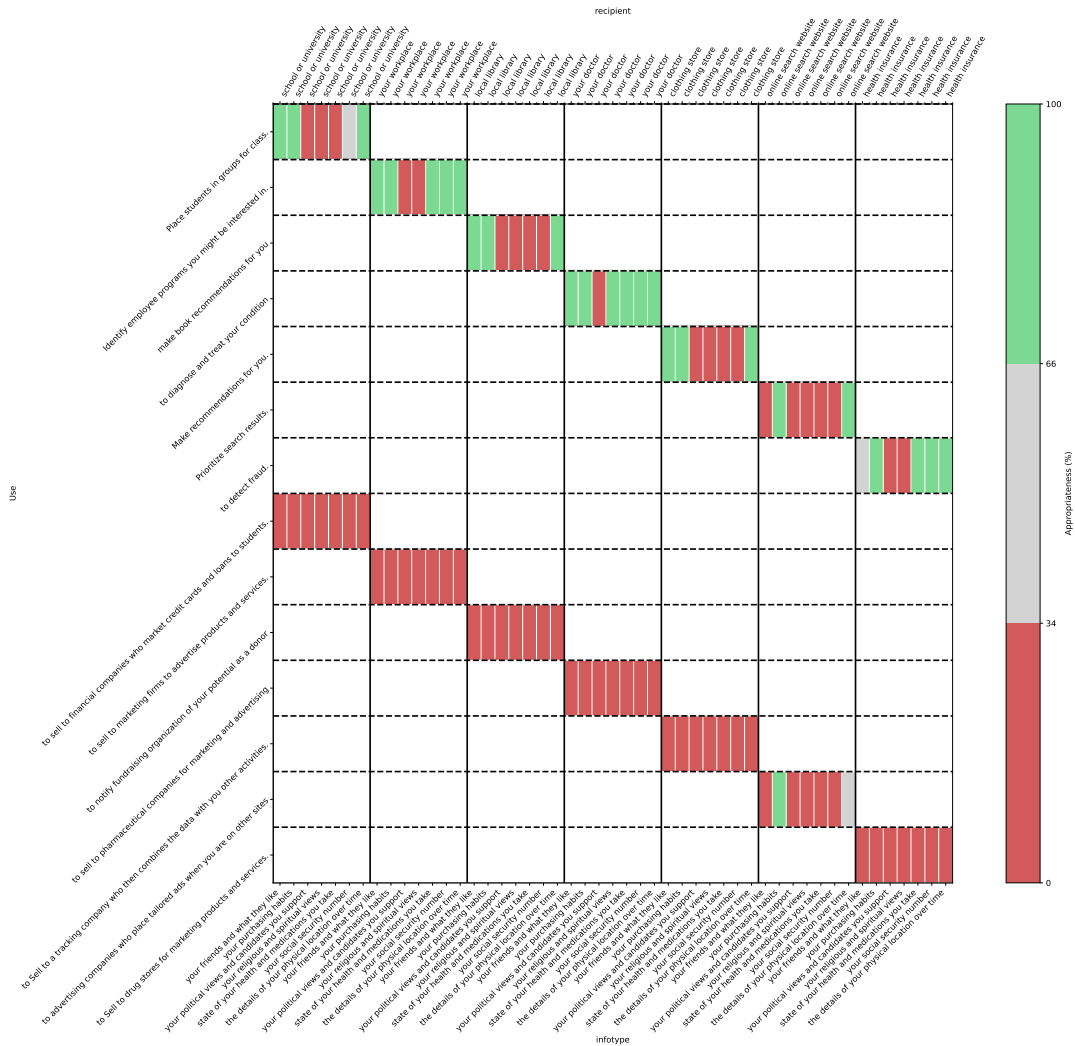


Figure 1. Heatmap of LLAMA-3.1-8B-INSTRUCT’s privacy judgements across CONFAIDE-style information-flow scenarios. Each cell corresponds to a unique combination of *infotype* (bottom x-axis), *recipient* (top x-axis), and *use* (y-axis). Cells are color-coded according to the model’s responses across 11 paraphrased versions of each prompt: green indicates that the model judged the scenario *appropriate* in more than 66% of paraphrases; red indicates that the model judged the scenario *inappropriate* in more than 66% of paraphrases; and gray denotes cases where neither label reached this threshold, reflecting inconsistent or neutral model behaviour.

5.4 Heatmap of Privacy Judgments

Using the stable binary labels for each (infotype, recipient, use) combination, we construct a heatmap that summarizes how often the model judges each information flow as appropriate and inappropriate. This allows us to see, at a glance, which infotype-recipient-use combinations the model generally judges as *appropriate* and which it judges as *inappropriate* (See Figure 1). We later use neighboring cells in this heatmap to define clean-corrupted pairs that differ in a single parameter for our circuit analysis in Section 5.6.

5.5 Circuit analysis and EAP-IG

An LLM can be viewed as a computational graph whose nodes are attention-heads, MLPs, or even neurons (depending how fine-grained we want to analyze it). The edges of this graph are simply connections between these nodes that determine where the output of each node goes. A circuit is a subgraph of this graph that includes the input and output logits. We are interested in finding circuits that are close in performance to the whole graph for a certain task. A task includes a set of clean and corrupted inputs and a metric. For example, $s =$ "The doctor shares your medical records with other doctors for diagnosis. This is appropriate." and $s' =$ "The doctor shares your medical records with other doctors for diagnosis. This is inappropriate." can be clean and corrupted inputs that elicit different responses from LLM (i.e., appropriate and inappropriate). The metric that we're interested in is logit difference (logit(appropriate) - logit(inappropriate)), which captures how much appropriateness changes when we corrupt an input by changing a single parameter in it.

Circuit identification. How can we identify good circuits? One way to do this would be to replace an edge activation in the clean forward pass with a corrupted activation, and add the edge to the circuit if the effect of this change is significant. However, this scales poorly with the size of the model. *Edge attribution patching* (EAP) [18] is a method that estimates the effect by a linear approximation. Formally, let $e = (u, v)$ be an edge in the graph, and let z_u, z'_u be activations of node u when we run the model on clean and corrupted inputs, respectively. Then we have:

$$L(z'_u) - L(z_u) \approx (z'_u - z_u)^T \nabla_v(s) \quad (1)$$

where L is our metric and the gradient is computed with respect to the input of v in the run of model on the clean input s . So this only requires two forward passes and one backward pass. A problem with EAP is that for an edge the gradient at the clean input s could be zero (removing attribution of an edge) while the gradient at corrupted input is nonzero. To address this issue, EAP-IG [10] uses integrated gradients, which accumulate the gradient along a line between clean and corrupted activations. That is, the score of an edge becomes:

$$(z'_u - z_u)^T \frac{1}{m} \sum_{k=1}^m \frac{\partial L(z' + \frac{k}{m}(z - z'))}{\partial z_v}, \quad (2)$$

where m is the number of intermediate points between clean and corrupted activations. We set $m = 5$ following the literature.

Faithfulness. The faithfulness of a circuit measures how close it is to the whole computational graph with respect to a metric P . Formally, it is defined as:

$$\frac{P_{\text{circuit}} - P_{\text{corrupted}}}{P_{\text{baseline}} - P_{\text{corrupted}}}$$

We use EAP-IG to identify circuits in our experiments as it was the state-of-the-art method for circuit identification at the start of the project. We refer the readers to [10] and references therein for more details on circuits and their identification methods.

5.6 Identifying circuits for privacy bias

As we discussed in the previous subsection, we need to define a set of clean-corrupted data pairs and a metric and then we can use EAP-IG out of the box. Depending on our target task we can define different clean-corrupted data pairs that differ in different parts. Here, we go over the choices we considered and their intended task.

Appropriate/Inappropriate circuits. Our goal is to identify which parts of the model are responsible for judging an information flow as *appropriate* or *inappropriate*. To do this, we take the information flows that are classified as appropriate or inappropriate (the green and red cells in Figure 1, respectively) and construct clean-corrupted pairs that differ only in the final appropriateness judgment.

For the "appropriate" circuit, we define:

- **Clean:** "{information flow vignette} This is appropriate."
- **Corrupted:** "{information flow vignette} This is inappropriate."

For the "inappropriate" circuit, we swap the labels:

- **Clean:** "{information flow vignette} This is inappropriate."
- **Corrupted:** "{information flow vignette} This is appropriate."

In each pair, the two vignettes are identical; only the final label word differs. The clean label is the majority (correct) label from the LLM over the paraphrases. Our metric is then:

$$\text{logit}(\text{clean label}) - \text{logit}(\text{corrupted label}),$$

so for the appropriate circuit the clean and corrupted labels are "appropriate" and "inappropriate," respectively, and for the inappropriate circuit they are reversed.

Example: For appropriate circuit an example of clean-corrupted pairs is as follows.

- **Clean:** Information about your friends and what they like is collected by school or university in order to Place students in groups for class. This is appropriate.
- **Corrupted:** Information about your friends and what they like is collected by school or university in order to Place students in groups for class. This is inappropriate.

Information flow parameters circuits. The goal here is to find parts of the model responsible for capturing the parameters of an information flow while keeping appropriateness the same. So, our clean and corrupted will be of the following form:

- **Clean:** "Consider the following [appropriate / inappropriate] information flow: {information flow vignette with clean parameter value}."
- **Corrupted:** "Consider the following [appropriate / inappropriate] information flow: {information flow vignette with corrupted parameter value}."

Here, the bracketed phrase "[appropriate / inappropriate]" means that both prompts use the *same* appropriateness label (either "appropriate" or "inappropriate"); only the parameter value differs. The two inputs therefore have different vignettes that differ in a single parameter.

Example: For infotype circuit, an example of clean-corrupted inputs is a follows:

- **Clean:** "It is inappropriate for your workplace to collect information in order to sell it to marketing firms to advertise products and services about your friends and what they like."
- **Corrupted:** "It is inappropriate for your workplace to collect information in order to sell it to marketing firms to advertise products and services about your purchasing habits."

In this case, our metric would be:

$$\text{logit}(\text{clean parameter}) - \text{logit}(\text{corrupted parameter})$$

The input pairs for other parameters are created similarly. For practical purposes, we paraphrase the vignettes so that the differing parameter is at the end of the vignette.

Circuits for both appropriateness and flow parameters. Here the goal is to find circuits that are suitable for both the appropriateness and the parameters. Thus, we look for information flows that differ in a single parameter and their appropriateness differ. We consider two possible choices for phrasing the clean-corrupted pairs: Question and statement phrasing. In question phrasing, the pairs would be "Consider the following information flow: {information flow vignette} Classify the information flow in exactly one word, either appropriate or inappropriate:" where again the information flows for clean and corrupted differ in a single parameter. In the statement phrasing, the pairs would be "Consider the following information flow: {information flow vignette} This is appropriate (inappropriate)". In this setting, the clean and

corrupted prompts differ in both a single information-flow parameter and in their appropriateness label (one is judged appropriate, the other inappropriate), unlike the previous parameter-only circuits subsection where the label is held fixed. The metric used here again is:

$$\text{logit}(\text{clean value}) - \text{logit}(\text{corrupted value})$$

Example: Question phrasing for the parameter "Use":

- **Clean:** "Consider the following information flow: Information about state of your health and medications you take is collected by clothing store in order to Make recommendations for you. Classify the information flow in exactly one word, either appropriate or inappropriate:"
- **Corrupted:** "Consider the following information flow: Information about state of your health and medications you take is collected by clothing store in order to Sell to a tracking company who then combines the data with you other activities. Classify the information flow in exactly one word, either appropriate or inappropriate:"

Example: Statement phrasing for the parameter Infotype:

- **Clean:** "Consider the following information flow: Information about your friends and what they like is collected by clothing store in order to Make recommendations for you. It is inappropriate."
- **Corrupted:** "Consider the following information flow: Information about your purchasing habits is collected by clothing store in order to Make recommendations for you. It is appropriate."

6 Results

6.1 Faithfulness

We measured the faithfulness of our circuits¹ across each circuit size and across three different circuit-finding methods: EAP, EAP-IG, and Clean-Corrupt. EAP-IG tended to yield the most faithful circuits. In Figure 2, you can see that the faithfulness of EAP-IG scores are marginally, but consistently, better than the other circuit finding methods. We tested two types of metrics when running EAP-IG. The first was logit difference, as discussed in the previous section, testing on both the difference between the median logit over multi-token labels, and the maximal logit over the label². We also tested using Kullback-Leibler (KL) divergence between the model's clean and corrupted distributions (EAP-IG-KL),

¹For the sake of space, here we do not explicitly discuss both the "Statement" and "Question" formulations as they saw similarly nonsensical results. Conclusions for "Statement" apply equally to "Question."

²Note that this was only for Use and Infotype parameters since "Appropriate" and "Inappropriate" are single tokens, so we do not have to take median or maximum in this case.

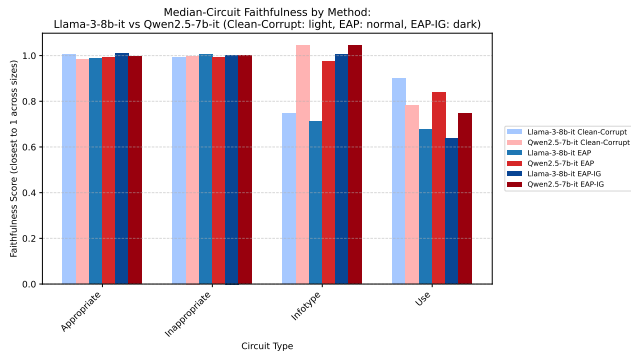


Figure 2. This bar graph contains the faithfulness scores for Appropriate, Inappropriate, Infotype, and Use circuits along the x-axis (in the listed order). Since we tested at various circuit sizes, the score on the y-axis is the the faithfulness (from 0 to 1) for whichever circuit size performed the best (closest to 1). For simplicity, the alternative formulations for Use and Infotype (Statement and Question) are not included as their results are not within the expected range. Llama depicted in blue and Qwen in red. Shade indicates the circuit-finding method, where the lightest is clean-corrupted, darkest is EAP-IG, and the middle is EAP.

which measures how much the full predicted token distribution shifts under the intervention. Unless stated otherwise, the scores are reported for median logit difference as we tended to see the best, consistent scores (see Figures 4 and 5).

Looking at Figure 3, we got high faithfulness scores for the non-statement formulations. Additionally, we find that some circuits have faithfulness scores above 1, which means that they are more faithful than the full model. We hypothesize that because large language models have a very large number of components, there may be some components that are negatively impacting the model’s performance, so when we run the model on a subset of the most influential nodes and edges, we can actually gain performance³. This should be explored in future work. Finally, the scores seen for the “Statement” formulation of Use and Infotype are very bad. Faithfulness scores should never be negative or significantly above 1, which indicates to us that our prompting strategy for the Clean and Corrupted prompts were not effective. EAP-IG could not identify consistent circuits across all the prompts for these datasets. One reason for this is that there were two changes in these prompts. Both the parameter and appropriateness were different between the clean and corrupted pair. Having two differences did not help to identify the appropriateness or inappropriateness of the parameter itself, rather it caused EAP-IG to fail to identify faithful circuits.

³Thank you to Anthony Hughes for helping us formulate this hypothesis.

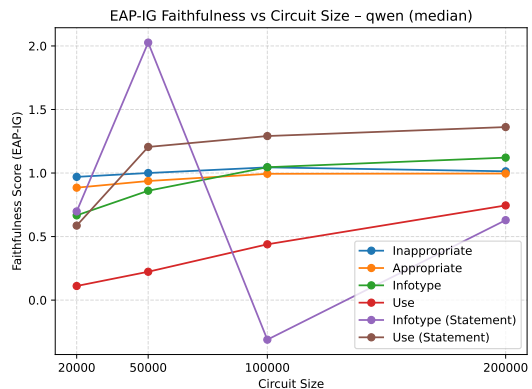


Figure 3. This plots the faithfulness of circuits for Qwen, with circuit size on the x-axis and faithfulness on the y-axis. Different colours indicate different circuits.

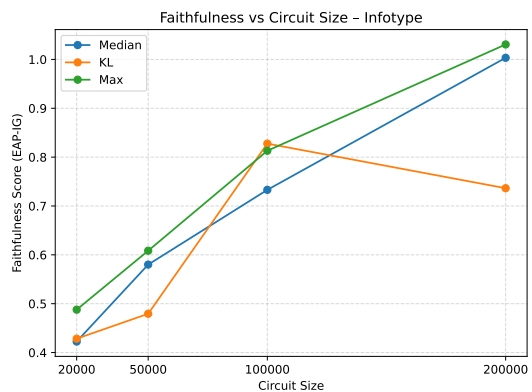


Figure 4. As a metric, KL performs worse than both logit difference methods. Results here are for Llama.

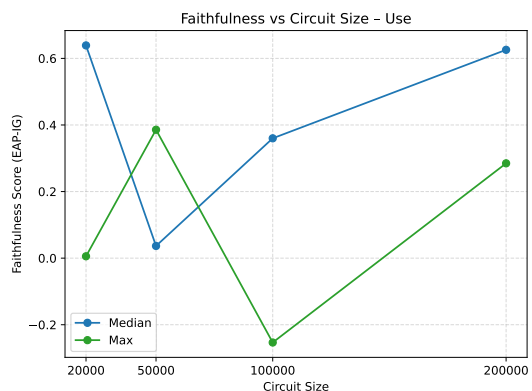


Figure 5. On the Use circuit, the max logit difference method saw negative faithfulness, indicating that it did not function as well as median. Results here are for Llama.

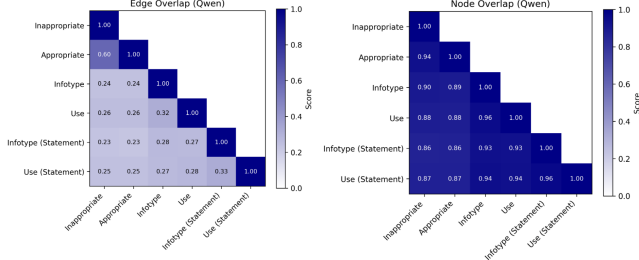


Figure 6. Heatmap of overlap between nodes and edges in Qwen. Scored using the Jaccard index (Intersection over Union) where 0 (white) is the no overlap and 1 (dark) is complete overlap.

6.2 Overlap

Following previous work [10, 11], our goal is to discover how the circuits that correspond to the appropriateness and inappropriateness of information flows interact with one another. This permits greater precision when performing the ablation study to rectify misaligned model behaviour. We find a circuit C_i for each of “Appropriate”, “Inappropriate”, “Use”, and “Infotype” and compute the overlap between the high importance edges and nodes for each pair of circuits. The high important edges for each C_i are those that have an EAP-IG score above τ_i , which is a dynamic threshold to capture those nodes and edges with scores above a certain percentile. In our study, following previous work, we use $\tau = 95\%$ and $\tau = 99\%$ [11]. The overlap between two circuits C_i and C_j was calculated using the Jaccard index (Intersection over Union):

$$\text{Overlap}(C_i, C_j, \tau) = \frac{|E_i^{\text{high}} \cap E_j^{\text{high}}|}{|E_i^{\text{high}} \cup E_j^{\text{high}}|} \times 100\% \quad (3)$$

We can see in Figure 6 that we have low edge overlap but high node overlap. This indicates that the model has specialized circuits for each of the CI parameters for which we identified circuits. The high node overlap, however, shows that CI seems to “live in” a subset of attention heads within the model, and the model relies on a subset of attention patterns, rather than completely different mechanisms.

6.3 Attention Patterns

Using the 95% and 99% thresholds, we can see the most important nodes (attention heads) in our circuits. This can help us identify which heads at which layers to intervene on in future patching work.

In Qwen, averaged across all circuit types, we see layers 17-19 are the most influential with several attention heads seeing high importance scores, and an early concentration in layer 0 (see Figure 7). In Llama, we see a numerous attention heads with higher activation scores, with a generally high level of activation in the early to middle layers (see Figure 8).

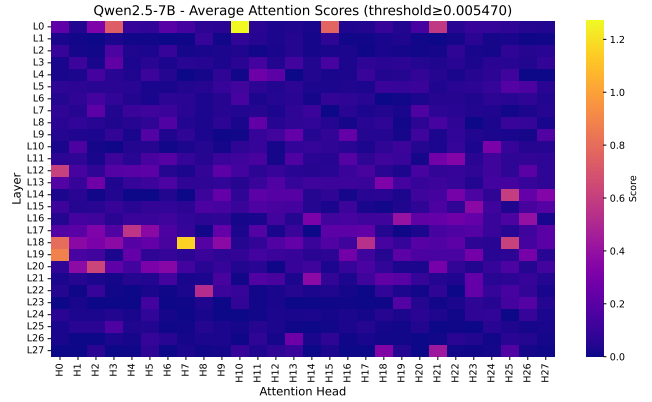


Figure 7. Heatmap of average attention scores for Qwen-2.5-7B across all layers and attention heads. Colour intensity encodes the average attention score for each layer–head pair, computed over the dataset and visualized after applying a threshold of 99.

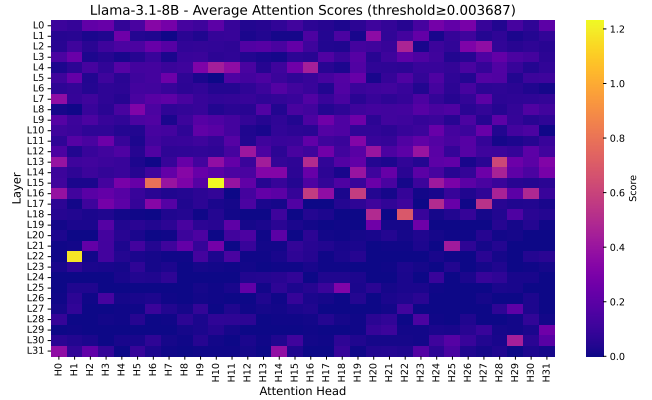


Figure 8. Heatmap of average attention scores for Llama-3.1-8B across all layers and attention heads. Colour intensity encodes the average attention score for each layer–head pair, computed over the dataset and visualized after applying a threshold of 99%.

The important thing here, is that although attention patterns differ between models, both have faithful and relevant circuits related to the parameters of interest.

7 Patching Strategy

Due to time constraints, we were unable to perform patching but we have a general idea for how we would proceed. Although we do not have an exact algorithm, looking at the overlap between circuits gives us a good idea of where to start. For example if we are concerned with the inappropriateness of infotypes, and want to intervene on the intersection of these circuits, we can verify they have overlap by looking at the heatmap (i.e., Figure 6), and notice that Use and Inappropriate have high node overlap. We can then check the

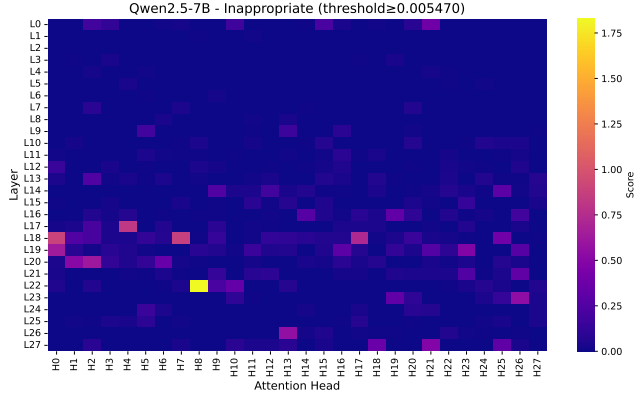


Figure 9. Heatmap of attention scores for Inappropriate circuit on Qwen using a 99% threshold.

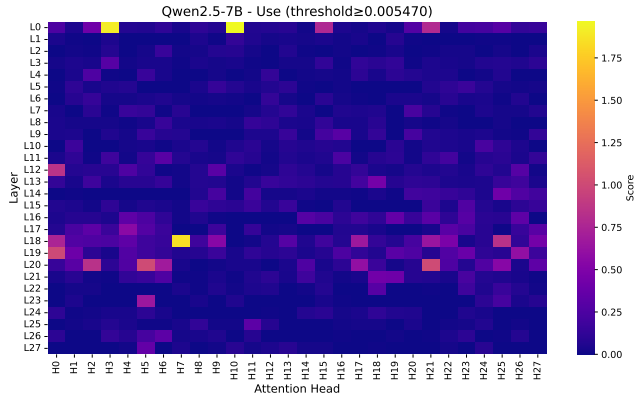


Figure 10. Heatmap of attention scores for Use circuit on Qwen using a 99% threshold.

attention head patterns in Figures 9 and 10 to identify which attention heads to ablate. For example, heads 1, 7, and 17 at layer 18 all have relatively high activation in both circuits, so we could test both zero and mean ablation strategies on these heads to start (Hughes et al. [11] suggests that the optimal strategy seems to be model dependent). Additionally, we can see in Figure 6 that Appropriate and Inappropriate have high edge and node overlap, so we can see which attention heads they diverge in, and we can test zero and mean ablation on heads that are not in their intersection, to fix instances where we see that the LLM is biased towards answering that an information flow is appropriate.

8 Future Work

There remains a lot of work to reach our ultimate goal of rectifying privacy bias in LLMs. Some additional work to be done is to run our existing pipeline on additional models. We also need to create an algorithm that will enable us to efficiently and effectively patch the models. Ideally, we would

see the patched models achieve high accuracy on the ConFAIde dataset, and we would be able to achieve this without a significant loss to the overall model utility. One measure that we would use for utility, which is common in related works [5, 11, 14, 25, 26], is to measure perplexity. We expect that, as seen in Hughes et al. [11], the perplexity would be higher over the baseline model (lower is better), but ideally we would identify an optimal patching algorithm that would minimize the increase over the baseline. Moreover, we would need to compare with other existing approaches for reducing the privacy bias in LLMs.

9 Conclusion

This work takes a first step toward opening the “black box” of how LLMs internalize privacy norms. By combining the Contextual Integrity framework with mechanistic interpretability, we constructed controlled pairs of clean and corrupted prompts from the ConFAIde dataset and used them to probe Llama 3.1-8B-it and Qwen 2.5-7B-it. Applying EAP-IG to these datasets, we identified faithful circuits corresponding to key CI parameters (RQ1), characterized their attention patterns, and quantified the overlap of each pair of circuits’ edges and nodes (RQ2). Our findings reveal faithful, specialized circuits that reliably distinguish appropriate from inappropriate information flows, with relatively low overlap between circuit nodes, which suggest specialized privacy mechanisms within these models.

Together, these results bridge behavioural privacy evaluations with an internal, circuit-level view of model decision-making. Beyond characterizing where privacy-relevant behaviour “lives” in the network, our circuit analyses and importance scores lay a concrete foundation for future circuit-editing algorithms that selectively patch privacy biases without retraining entire models. Moving forward, extending this methodology to additional architectures and richer real-world scenarios, and empirically evaluating circuit-level interventions, against less fine-grained interventions, on downstream privacy outcomes, will be crucial steps toward developing AI systems that are aligned with human privacy-expectations.

10 Novelty

Our work introduces some methodological contributions to the study of privacy-related behavior in Large Language Models. First, we construct a new dataset specifically designed to find and analyze circuits associated with privacy bias in information-flow judgments. To support mechanistic interpretability methods, we developed tailored clean-corrupted prompting strategies for running EAP-IG. We also demonstrate, to our knowledge, the first identification of circuits linked to privacy bias across two model families. Finally, we extend EAP-IG to operate effectively with multi-token labels by allowing the attribution objective to aggregate logits

over all label tokens using either the median or maximum token logit to better suit our use-case.

11 Contributions & Acknowledgments

The dataset creation was primarily performed by Neel and Sajad. Circuit finding was completed by Sajad and Lucas. Results and analysis were performed by Lucas and Neel. Neel, Sajad, and Lucas each completed 33.33% of the work.

Finally, we would like to acknowledge the contributions of Vasishth Duddu and Anthony Hughes, without whom we would not have been able to complete this project. Their advice and insights were invaluable throughout this entire process.

References

- [1] Eugene Bagdasarian, Ren Yi, Sahra Ghalebikesabi, Peter Kairouz, Marco Gruteser, Sewoong Oh, Borja Balle, and Daniel Ramage. 2024. AirGapAgent: Protecting Privacy-Conscious Conversational Agents. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security* (Salt Lake City, UT, USA) (CCS '24). Association for Computing Machinery, New York, NY, USA, 3868–3882. doi:10.1145/3658644.3690350
- [2] Bowen Cao, Deng Cai, Zhisong Zhang, Yuexian Zou, and Wai Lam. 2024. On the Worst Prompt Performance of Large Language Models. arXiv:2406.10248 [cs.CL] <https://arxiv.org/abs/2406.10248>
- [3] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2020. Extracting Training Data from Large Language Models. In *USENIX Security Symposium*. <https://api.semanticscholar.org/CorpusID:229156229>
- [4] Irene Chen, Peter Szolovits, and Marzyeh Ghassemi. 2019. Can AI Help Reduce Disparities in General Medical and Mental Health Care? *AMA journal of ethics* 21 (02 2019), E167–179. doi:10.1001/amajethics.2019.167
- [5] Ruizhe Chen, Tianxiang Hu, Yang Feng, and Zuozhu Liu. 2024. Learnable Privacy Neurons Localization in Language Models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 256–264. doi:10.18653/v1/2024.acl-short.25
- [6] Federico Errica, Davide Sanvito, Giuseppe Siracusano, and Roberto Bifulco. 2025. What Did I Do Wrong? Quantifying LLMs' Sensitivity and Consistency to Prompt Engineering. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Association for Computational Linguistics, 1543–1558. doi:10.18653/v1/2025.naacl-long.73
- [7] Wei Fan, Haoran Li, Zheyang Deng, Weiqi Wang, and Yangqiu Song. 2024. GoldCoin: Grounding Large Language Models in Privacy Laws via Contextual Integrity Theory. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 3321–3343. doi:10.18653/v1/2024.emnlp-main.195
- [8] Kathleen C. Fraser, Hillary Dawkins, Isar Nejadgholi, and Svetlana Kiritchenko. 2025. Fine-Tuning Lowers Safety and Disrupts Evaluation Consistency. In *Proceedings of the The First Workshop on LLM Security (LLMSEC)*, Leon Derczynski, Jekaterina Novikova, and Muhao Chen (Eds.). Association for Computational Linguistics, Vienna, Austria, 129–141. <https://aclanthology.org/2025.llmsec-1.10/>
- [9] Chengguang Gan and Tatsunori Mori. 2023. Sensitivity and Robustness of Large Language Models to Prompt Template in Japanese Text Classification Tasks. In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, Chu-Ren Huang, Yasunari Harada, Jong-Bok Kim, Si Chen, Yu-Yin Hsu, Emmanuele Chersoni, Pranav A, Winnie Huiheng Zeng, Bo Peng, Yuxi Li, and Junlin Li (Eds.). Association for Computational Linguistics, Hong Kong, China, 1–11. <https://aclanthology.org/2023.paclic-1.1/>
- [10] Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. 2024. Have faith in faithfulness: Going beyond circuit overlap when finding model mechanisms. *arXiv preprint arXiv:2403.17806* (2024).
- [11] Anthony Hughes, Vasishth Duddu, N. Asokan, Nikolaos Aletras, and Ning Ma. 2025. PATCH: Mitigating PII Leakage in Language Models with Privacy-Aware Targeted Circuit Patching. arXiv:2510.07452 [cs.CR] <https://arxiv.org/abs/2510.07452>
- [12] Guangchen Lan, Huseyin A. Inan, Sahar Abdelnabi, Janardhan Kulkarni, Lukas Wutschitz, Reza Shokri, Christopher G. Brinton, and Robert Sim. 2025. Contextual Integrity in LLMs via Reasoning and Reinforcement Learning. arXiv:2506.04245 [cs.AI] <https://arxiv.org/abs/2506.04245>
- [13] Haoran Li, Wei Fan, Yulin Chen, Jiayang Cheng, Tianshu Chu, Xuebing Zhou, Peizhao Hu, and Yangqiu Song. 2025. Privacy Checklist: Privacy Violation Detection Grounding on Contextual Integrity Theory. arXiv:2408.10053 [cs.CL] <https://arxiv.org/abs/2408.10053>
- [14] Nils Lukas, Ahmed Salem, Robert Sim, Shruti Tople, Lukas Wutschitz, and Santiago Zanella-Beguelin. 2023. Analyzing Leakage of Personally Identifiable Information in Language Models. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, Los Alamitos, CA, USA, 346–363. doi:10.1109/SP46215.2023.10179300
- [15] Matthieu Meeus, Igor Shilov, Shubham Jain, Manuel Faysse, Marek Rei, and Yves-Alexandre de Montjoye. 2025. SoK: Membership Inference Attacks on LLMs are Rushing Nowhere (and How to Fix It). arXiv:2406.17975 [cs.CL] <https://arxiv.org/abs/2406.17975>
- [16] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2022. A Survey on Bias and Fairness in Machine Learning. arXiv:1908.09635 [cs.LG] <https://arxiv.org/abs/1908.09635>
- [17] Niloofar Mireshghallah, Hyunwoo Kim, Xuhui Zhou, Yulia Tsvetkov, Maarten Sap, Reza Shokri, and Yejin Choi. 2024. Can LLMs Keep a Secret? Testing Privacy Implications of Language Models via Contextual Integrity Theory. In *The Twelfth International Conference on Learning Representations*.
- [18] Neel Nanda. 2023. Attribution Patching: Activation Patching At Industrial Scale. <https://www.neelnanda.io/mechanistic-interpretability/attribution-patching>. Accessed: YYYY-MM-DD.
- [19] Iviline C. Ngong, Swanand Ravindra Kadhe, Hao Wang, Keerthiram Murugesan, Justin D. Weisz, Amit Dhurandhar, and Karthikeyan Natesan Ramamurthy. 2025. Protecting Users From Themselves: Safeguarding Contextual Privacy in Interactions with Conversational Agents. In *Findings of the Association for Computational Linguistics: ACL 2025*, Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (Eds.). Association for Computational Linguistics, Vienna, Austria, 26196–26220. doi:10.18653/v1/2025.findings-acl.1343
- [20] Helen Nissenbaum. 2009. Privacy in Context: Technology, Policy, and the Integrity of Social Life. *American Behavioral Scientist* 58 (2009).
- [21] Chris Olah. 2022. Mechanistic Interpretability, Variables, and the Importance of Interpretable Bases. <https://www.transformer-circuits.pub/2022/mech-interp-essay>. Accessed: 2025-12-04.
- [22] Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying Language Models' Sensitivity to Spurious Features in Prompt Design or: How I learned to start worrying about prompt formatting. arXiv:2310.11324 [cs.CL] <https://arxiv.org/abs/2310.11324>
- [23] Yijia Shao, Tianshi Li, Weiyan Shi, Yanchen Liu, and Diyi Yang. 2025. PrivacyLens: Evaluating Privacy Norm Awareness of Language Models in Action. arXiv:2409.00138 [cs.CL] <https://arxiv.org/abs/2409.00138>

- [24] Yan Shvartzshnaider and Vasisht Duddu. 2025. Investigating Privacy Bias in Training Data of Language Models. arXiv:2409.03735 [cs.LG] <https://arxiv.org/abs/2409.03735>
- [25] Xinwei Wu, Weilong Dong, Shaoyang Xu, and Deyi Xiong. 2024. Mitigating Privacy Seesaw in Large Language Models: Augmented Privacy Neuron Editing via Activation Patching. In *Findings of the Association for Computational Linguistics: ACL 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 5319–5332. doi:10.18653/v1/2024.findings-acl.315
- [26] Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. 2023. DEPN: Detecting and Editing Privacy Neurons in Pretrained Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 2875–2886. doi:10.18653/v1/2023.emnlp-main.174
- [27] Ren Yi, Octavian Suci, Adria Gascon, Sarah Meiklejohn, Eugene Bagdasarian, and Marco Gruteser. 2025. Privacy Reasoning in Ambiguous Contexts. arXiv:2506.12241 [cs.AI] <https://arxiv.org/abs/2506.12241>
- [28] Lin Zhang, Wenshuo Dong, Zhuoran Zhang, Shu Yang, Lijie Hu, Ninghao Liu, Pan Zhou, and Di Wang. 2025. EAP-GP: Mitigating Saturation Effect in Gradient-based Automated Circuit Identification. arXiv:2502.06852 [cs.LG] <https://arxiv.org/abs/2502.06852>
- [29] Shikun Zhang, Yan Shvartzshnaider, Yuanyuan Feng, Helen Nissenbaum, and Norman Sadeh. 2022. Stop the Spread: A Contextual Integrity Perspective on the Appropriateness of COVID-19 Vaccination Certificates. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (Seoul, Republic of Korea) (FAccT ’22)*. Association for Computing Machinery, New York, NY, USA, 1657–1670. doi:10.1145/3531146.3533222