IBCIRCUIT: TOWARDS HOLISTIC CIRCUIT DISCOV ERY WITH INFORMATION BOTTLENECK

Anonymous authors

Paper under double-blind review

ABSTRACT

Circuit discovery has recently attracted attention as a potential research direction to explain the nontrivial behaviors of language model (LM). It aims to find the computational subgraph, also known as *circuit*, that explains LM's behavior on specific tasks. Most studies determine the circuit for a task by performing causal interventions independently on each component. However, they ignored the holistic nature of the circuit, which is an interconnected system of components rather than an independent combination. Additionally, existing methods require redesigning a unique corrupted activation for each task, which are complicated and inefficient. In this work, we propose a novel circuit discovery approach based on the principle of Information Bottleneck, called IBCircuit, to identify the most informative circuit from a holistic perspective. Furthermore, IBcircuit can be applied to any given task without corrupted activation construction. Our experiments demonstrate the ability of IBCircuit to identify the most informative circuit in the model. The results from IBCircuit suggest that **the earlier layers in Transformer-based models are crucial in capturing factual information**.

024 025 026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

027 028

Circuits discovery in transformer-based language models usually involves identifying the subgraph
 (circuits) within the model which are responsible for solving specific tasks. Previous efforts to
 identify circuits within language models have led to the discovery of networks comprising attention
 heads and multi-layer perceptrons (MLPs) that either partially or fully explain the model's behaviors
 on tasks like indirect object recognition, modular arithmetic, docstring completion, and forecasting
 subsequent dates(Wang et al., 2022; Nanda et al., 2023; Hanna et al., 2024). However, the challenge
 is that transformer-based language models primarily operate as complicated black boxes when en gaging in multi-layered nonlinear interactions within high-dimensional spaces (Wang et al., 2022),
 making it exceptionally difficult to understand their behavior.

Work in circuit analysis seeks to decode these models by reverse engineering (Conmy et al., 2023; Wang et al., 2022; Meng et al., 2022; Geva et al., 2021). Recent circuit analysis methods (Räuker et al., 2023), such as activation patching (Meng et al., 2022; Goldowsky-Dill et al., 2023) or knock-040 outs (Wang et al., 2022), ignores the holistic nature of the circuit, which is an interconnected system 041 of components rather than an independent combination. Additionally they scales poorly with the 042 model size. Moreover, when faced with a task that has never been encountered before, most works 043 need to redesign new corrupted activation for patching. Although recent works have proposed at-044 tribution patching (Syed et al., 2023) to estimate the importance of each edge in the computational subgraph without independent patching, they still need to design new corrupted activation for different task, which is inconvenient and complicated. 046

The Information Bottleneck (IB) leverages Shannon mutual information to quantify the compressed and informative nature of data distributions (Yu et al., 2022). This concept, rooted in information theory, aims to balance the trade-off between the complexity of the representation and its predictive power. The primary objective of IB is to distill a compressed yet predictive representation of the input signal (Tishby et al., 2000). By focusing on the most relevant aspects of the data, IB helps in reducing redundancy and noise, thereby enhancing the efficiency of data processing. This technique has been successfully applied in various domains, demonstrating its versatility and effectiveness. In feature selection, IB helps in identifying the most relevant features that contribute to the predictive

power of a model, thereby improving performance and reducing computational costs (Achille & Soatto, 2018b; Kim et al., 2021; Schulz et al., 2020). In representation learning, IB aids in creating compact and meaningful representations of data, which are crucial for tasks such as clustering, classification, and anomaly detection (Luo et al., 2019; Qian et al., 2020; Wu et al., 2020). These applications highlight the broad utility of IB in enhancing model interpretability and performance.

In this work, we propose the Information Bottleneck Circuit (IBCircuit), a novel approach designed 060 to identify critical components within a Transformer-based model that are capable of task execution. 061 The IBCircuit consists of two primary stages: Model Perturbation and Component Selection. 062 During the Model Perturbation stage, the IBCircuit injects Gaussian noise into various components 063 of the model, such as the activations of attention heads and MLP layers. Additionally, the IB weights 064 are learned for the components to control the amount of noise added. The rationale for this noise injection is that it modulates the information flow from the original pretrained model to the perturbed 065 version, with more significant noise leading to greater information distortion. As such, training the 066 IBCircuit encourages the perturbed model to maintain its informativeness, which implies less noise 067 injection. This process effectively approximates the condition of information compression, with less 068 noise injection serving to identify the most informative components. In the Component Selection 069 stage, the circuit is effectively formed by retaining these informative components. The contributions of this work are summarized as follows: 071

- We propose a novel circuit discovery method, IBCircuit, which utilizes information bottleneck to globally identify the most informative circuit of the model for circuit discovery from a holistic perspective.
 - We introduce a model perturbation method that incorporates noise injection and adaptive learning of *IB weights*, which can be applied to any given task without corrupted activation construction.
- We validate the effectiveness of the IBCircuit in identifying the most informative attention heads in the Transformer-based model on the Indirect Object Identification (IOI) and Greater-Than Circuit Discovery task. The results from IBCircuit suggest that **the earlier layers in Transformer-based models are crucial in capturing factual information**.

2 RELATED WORK

087

073

075 076

077

078

079

081

082

084 085

Circuit Analysis. Circuit analysis seeks to understand machine learning models by identifying circuits in models that is responsible for given behaviors (Geiger et al., 2021; Wang et al., 2022; 090 Conmy et al., 2023). Most of existing works conduct circuit analysis for language model leveraging 091 activation patching (Zhang & Nanda, 2023) or its variants. Some overwrite activation values with 092 zeros (Cammarata et al., 2021; Olsson et al., 2022), while others erase activation using the mean activation on the dataset (Wang et al., 2022; Hanna et al., 2024). Other works (Geiger et al., 2021; Wu et al., 2024) use the interchange interventions instead, which replaces the activation value of a 094 node on one data point with its value on another data point. However, it has been justified that both 095 zero and mean activations take the model too far away from actually possible activation distributions 096 (Chan et al., 2022). Additionally, methods based on activation patching are inefficient as they require sequential operations on individual activations to assess their impact on task performance.

098 Information Bottleneck in Deep Learning. The principle of the information bottleneck (IB) aims to extract a compressed yet predictive code from the input signal (Tishby et al., 2000). Alemi et 100 al. (Alemi et al., 2016) initially introduced the concept of the variational information bottleneck 101 (VIB) to enhance deep learning. Currently, IB and VIB find applications primarily in representation 102 learning and feature selection in deep learning. In representation learning, researchers aim to learn a 103 compressed representation with the information bottleneck principle (Luo et al., 2019; Goyal et al., 104 2019; Qian et al., 2020; Wu et al., 2020). For feature selection, IB is used to select a subset of input 105 features such as pixels in images or dimensions in vectors (Achille & Soatto, 2018b; Kim et al., 2021; Schulz et al., 2020; Yu et al., 2020), which are maximally predictive to the label of input 106 data. Unlike previous work, we consider a rarely explored perspective: the compressed and relevant 107 information in model components for specific behavior.

¹⁰⁸ 3 BACKGROUNDS

3.1 NEURAL CIRCUITS

112 Circuit discovery seeks to reverse-engineer model behavior by localizing it to subgraphs of the 113 model's computation graph and explaining it. This approach aims to understand how specific parts of a model contribute to its overall functionality. Much research considers models as connected 114 directed computational graphs, denoted as G. These graphs represent the flow of computations 115 within the model, providing a structured way to analyze and interpret its operations. A transformer 116 language model's (LM) computational graph is a directed graph (digraph) that describes the compu-117 tations it performs. This graph flows from the LM's inputs to the unembedding layer, which projects 118 its activations into vocabulary space. The nodes in this digraph are defined to be the LM's attention 119 heads and multi-layer perceptrons (MLPs), though other levels of granularity, such as individual 120 neurons, are also possible. This hierarchical structure allows for detailed analysis at various levels 121 of abstraction. 122

A circuit is a subgraph that connects the inputs to the logits, which are the final outputs before 123 the softmax layer in a language model. In this graph, source nodes represent the model's input, 124 sink nodes represent the model's output, and intermediate nodes represent units of computation. 125 By identifying and analyzing these circuits, researchers can pinpoint which parts of the model are 126 responsible for specific behaviors and how information flows through the network. The concept of 127 *Neural Circuits C* refers to induced subgraphs of G that are responsible for specific behaviors and 128 exhibit distinct functionality. These circuits can be thought of as the building blocks of the model's 129 decision-making process. For instance, in Wang et al. (2022), the authors discover an Indirect Object 130 Identification (IOI) circuit in GPT-2 small based on Activation Patching Meng et al. (2022).

131 132

133

3.2 TRANSFORMER ARCHITECTURE

A transformer model $G: \mathcal{X} \to \mathcal{Y}$ maps a token sequence $x = [x_1, \dots, x_T] \in \mathcal{X}$ to a probability distribution $y \in \mathcal{Y}$. The *i*-th token at layer *l* is embedded as a series of hidden state vectors $h_i^{(l)} \in \mathbb{R}^H$, where *H* is the dimension of hidden state vectors. The input $h_i^{(0)}$ to the transformer is a sum of position pos(i) and token embeddings $emb(x_i)$.

The internal computation of hidden states $h_i^{(l)}$ in *G* can be summarized as follows: for each layer, it combines global attention $a_i^{(l)}$ and local MLP $m_i^{(l)}$ contributions computed from previous layers. Additionally, the *residual stream* draws information from previous states. The internal computations are as follows:

$$h_i^{(l)} = h_i^{(l-1)} + a_i^{(l)} + m_i^{(l)}, \tag{1}$$

$$a_i^{(l)} = \operatorname{attn}^{(l)} \left(h_1^{(l-1)}, h_2^{(l-1)}, \dots, h_i^{(l-1)} \right), \tag{2}$$

145 146 147

143 144

$$m_i^{(l)} = W_{proj}^{(l)} \sigma \left(W_{fc}^{(l)} \gamma \left(a_i^{(l)} + h_i^{(l-1)} \right) \right).$$
(3)

Each layer's MLP is a two-layer neural network parameterized by matrices $W_{proj}^{(l)}$ and $W_{fc}^{(l)}$, with rectifying nonlinearity σ and normalizing nonlinearity γ . For each individual attention head, it is parameterized by four matrices W_Q, W_K, W_V and W_O . The QK matrix is used to compute the attention matrix $A = \operatorname{softmax} \left(\frac{QK^T}{\sqrt{d_k}}\right) \in \mathbb{R}^{T \times T}$, where d_k is the dimension of queries Q and keys K, while the OV matrix determines what is written into the residual stream. For further background on transformers, we refer to Vaswani et al. (2017).

154 155

156

3.3 INFORMATION BOTTLENECK

157 The Information Bottleneck (IB) method is a powerful framework in information theory that focuses 158 on finding a compressed representation of data while retaining the most relevant information for a 159 specific task. The IB method aims to balance the trade-off between compressing the input data and 160 preserving the information necessary for predicting an output variable. In essence, the Information 161 Bottleneck method seeks to transform the input data into a new, more compact representation that 162 still contains the critical features needed for accurate prediction or classification. This is particularly



Figure 1: The framework of IBCircuit. During the model perturbation phase, we first inject noise to activation by using the *IB weight* λ , and then optimize λ based on the Information Bottleneck objective. We represent attention components with more noise injected using darker orange squares, and MLP components with more noise injected using darker green rhombuses. The final selection is based on choosing components with the least amount of noise.

useful in scenarios where the original data is high-dimensional or contains a lot of noise, making it challenging to process directly. We denote I(X;Y) the mutual information between two random variables X and Y, it is defined as follows:

$$I(X;Y) = \int_X \int_Y p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \mathrm{d}x \mathrm{d}y.$$
(4)

Given the variable X and its associated label Y, the Information Bottleneck (IB) Tishby et al. (2000) aims to learn the minimal sufficient variable Z by optimizing the following objective: I(X = Z) = I(X = Z)

$$\min_{Z} -I(Y;Z) + \alpha I(X;Z).$$
(5)

187 Here, α is the Lagrange multiplier used to balance the two terms. The first term encourages Z to con-188 tain informative content about the label Y, while the second term minimizes the mutual information 189 between X and Z, ensuring that Z only receives limited information from variable X.

4 Methodology

In this section, we introduce the proposed IBCircuit. In Section 4.1, we present the objective of IBCircuit in the Information Bottleneck framework. We outline the process of implementing Model Perturbation through Noise Injection and IBCircuit Optimization in Section 4.2. Additionally, in Section 4.3, we discuss the applications of Component Selection in Transformer-based models.

196 197

199

190 191

177

178

179

180 181 182

4.1 INFORMATION BOTTLENECK CIRCUIT

Circuit Identification aims to identify the most crucial components within a model, striking a balance between minimizing the number of components and maximizing their ability to perform tasks, which aligns perfectly with the concept of the Information Bottleneck (IB). Inspired by the IB, IBCircuit integrates circuit analysis into the information bottleneck framework to identify key components of pretrained models. Specifically, we denote G as the set of all components of the pretrained model, Y as the output of the pretrained model on a specific task, and C as the circuit composed of critical components. We reformulate Equation equation 5 to obtain the objective of IBCircuit as follows:

206 207

$$\min_{C} -I(Y;C) + \alpha I(G;C).$$
(6)

This objective cannot be directly optimized since the mutual information is intractable to compute. Therefore, we propose an alternative objective to evaluate the compression quality of circuit C by injecting noise into the intermediate activations of the pretrained model. As shown in Fig. 1, the IBCircuit consists of the Model Perturbation and Component Selection.

212

- 213 4.2 MODEL PERTURBATION
- Our approach involves adaptively learning noise and injecting it into each activation, allowing the model to maintain similar performance on a specific task even after perturbation.

216 4.2.1 NOISE INJECTION 217

218 We define a pretrained model G as a composition of n components, denoted as $G = [v_1, v_2, \cdots, v_n]$, with the corresponding intermediate activations $\mathbf{h} = [h_1, h_2, \cdots, h_n]$. To introduce perturbations 219 into G, we inject noise into intermediate activations with learnable IB weights $\lambda = [\lambda_1, \lambda_2, \dots, \lambda_n]$. 220 We sample noise ϵ from a parametric noise distribution. We insert Gaussian noise to the intermediate activations. In high-dimensional space, random vectors tend to concentrate on a spherical surface, 222 and their distribution can be approximated by a Gaussian distribution. For each component v_i , we

224

225

$$\hat{h}_i = \lambda_i h_i + (1 - \lambda_i)\epsilon. \tag{7}$$

226 Here, \hat{h}_i represents the perturbed activation. In order to avoid introducing the noise that devia-227 tion from the original distribution, we calculate the mean and variance of the whole dataset 228 as the variance of our injected noise. The learnable *IB* weight λ_i acts as a transmission probabil-229 ity, controlling the amount of information extracted from h_i to \hat{h}_i . To ensure that $\lambda_i \in (0, 1)$, we 230 define $\lambda_i = \text{Sigmoid}(\omega_i)$, where $\omega_i \sim \mathcal{N}(0,1)$ is the learnable parameter. When $\lambda_i \to 1$, all the 231 information from h_i is transferred to \hat{h}_i without loss. Conversely, when $\lambda_i \to 0$, \hat{h}_i contains no in-232 formation from h_i but only noise. Unlike baseline methods that iteratively perturb each intermediate 233 activation, the proposed method allows for simultaneous adjustment of the information flow from 234 all activations **h** to $\mathbf{h} = [h_1, h_2, \dots, h_n]$ by learning and updating all the *IB* weights λ together. We 235 denote the perturbed model as $\hat{G} = [\hat{v}_1, \hat{v}_2, \cdots, \hat{v}_n]$, we can then select the informative components 236 of G into C with learned IB weights. The selection process will be elaborated in Section 4.3. 237

perturb the intermediate activation h_i by combining it with noise ϵ using the *IB weight* $\lambda_i \in (0, 1)$:

4.2.2 Optimization of Information Bottleneck Circuit

The perturbed model G is learned by extracting information from the pretrained model G to achieve 240 the same performance Y on a specific task. On the one hand, we compress the effective information 241 in G by injecting noise. On the other hand, we aim to maximize the information content of the 242 perturbed model \hat{G} to achieve the performance Y. Therefore, the reconstructed IBCircuit objective 243 is as follows: 244

256

262

238

239

$$\min_{\hat{G}} -I(Y;\hat{G}) + \alpha I(G;\hat{G}).$$
(8)

246 The first term encourages \hat{G} to be sufficient for predicting Y, and the second term constrains the 247 information that \hat{G} learns from G. These two terms require us to inject noise into G selectively so 248 that \hat{G} maintains valuable information as much as possible. The intuition is that injecting noises 249 into the critical components of G is more harmful to the functionality of G than that into the irrel-250 evant components. In that sense, the critical components are less likely to be injected with noise. 251 Therefore, we can select the circuit C from \hat{G} by this criterion after training the IBCircuit. 252

We justify the above formulation through the following derivation. Let G_s be a subset of G, which 253 is independent to Y. Denote G_{ϵ} as the noisy subset of G determined by injected noise, if we select 254 the circuit C by dropping G_{ϵ} in G, the following inequality holds: 255

$$I(G_s; C) \le I(G_s; \hat{G}) \le I(G; \hat{G}) - I(Y; \hat{G}).$$
 (9)

257 This equation indicates when setting $\alpha = 1$ in Eq. equation 8, the IBCircuit objective upper bounds 258 the mutual information of G_s and C. Hence, optimizing the IBCircuit objective encourages C to 259 be less related to components in G_s which are irrelevant to Y. The detailed proof is provided in 260 Appendix A. 261

Minimizing $-I(Y;\hat{G})$. We first examine the first term $-I(Y;\hat{G})$ in Eq. equation 8, which encourages \hat{G} to be informative of the output Y. We derive the upper bound of $-I(Y; \hat{G})$ as follows:

$$-I(Y; \hat{G}) = \mathbb{E}_{Y}[\log p(Y)] - \mathbb{E}_{Y,\hat{G}}[\log p(Y|\hat{G})]$$

$$\leq -\mathbb{E}_{Y,\hat{G}}[\log p(Y|\hat{G})]$$

$$:= \mathcal{L}_{CE}(q_{\theta}(Y|\hat{G}))$$
(10)

where $q_{\theta}(Y|\hat{G})$ is the variational approximation to the true posterior $p(Y|\hat{G})$. This inequality 269 demonstrates that the minimization of $-I(Y; \hat{G})$ can be achieved by minimizing the training loss of the model. Since we validate IBCircuit on the Transformer-based language model in this paper, we use the Cross Entropy Loss for the next token prediction training, denoted as \mathcal{L}_{CE} .

Minimizing $I(G; \hat{G})$. For the second term $I(G; \hat{G})$ in Eq. equation 8, which aims to extract the informative components from G that contains minimal information about G. We minimize it by training the IBCircuit. By injecting more noise into insignificant components, while injecting less noise into more informative ones. Following Yu et al. (2022), by choosing the distribution of noise $\epsilon \sim \mathcal{N}(\mu_h, \sigma_h^2)$, where μ_h, σ_h^2 are mean and variance of intermediate activations h in G. $I(G; \hat{G})$ has a tractable variational upper bound as follows:

282 283

284

286 287 288

289

290

296

297

298 299

300 301

302 303

305 306

307

308

310

311

312

313

314

315

 $I(G;\hat{G}) \leq \mathbb{E}_{G}(-\frac{1}{2}\log A_{G} + \frac{1}{2n}A_{G} + \frac{1}{2n}B_{G}^{2}) =: \mathcal{L}_{MI}(G;\hat{G}),$ where $A_{G} = \sum_{i=1}^{n} (1-\lambda_{i})^{2}$ and $B_{G} = \frac{\sum_{i=1}^{n} \lambda_{i}(h_{i}-\mu_{h})}{\sigma_{h}}.$

Final Objectives. Finally, the overall loss is defined as follows:

$$\mathcal{L} = \mathcal{L}_{CE} + \alpha \mathcal{L}_{MI},\tag{12}$$

(11)

where α is hyperparameter used to adjust the weights of the loss.

4.3 COMPONENT SELECTION

We identify the critical components with less perturbation in the Component Selection stage.

Attention Heads Selection. To identify circuit attention heads, we perturb the attention matrix A in each attention head of the Transformer module as $\hat{A}_i = \lambda_i A_i + (1 - \lambda_i)\epsilon$. We then define a threshold δ and select significant attention patterns between tokens based on the condition $\lambda_i > \delta$. Then we examine whether the corresponding attention patterns meet the definitions of critical attention heads to capture the model's circuit.

MLP Layers Selection. We perturb the hidden state m of each MLP layer defined in Eq. equation 3 to identify critical MLP layers, i.e., $\hat{m}_i^{(l)} = \lambda_i m_i^{(l)} + (1 - \lambda_i)\epsilon$. Similarly, we can find critical MLP layers by defining a threshold δ and identifying MLP layers where $\lambda_i > \delta$, or identify commonly occurring layers with larger λ across multiple examples as critical MLP layersMeng et al. (2022).

5 EVALUATING IBCIRCUIT

To investigate the effectiveness of IBCircuit, we conduct the evaluation to answer the following research questions (RQs):

- **RQ1** (Grounded in Previous Work): Can IBCircuit effectively reproduce canonical circuits taken from previous works that found an end-to-end circuit explaining behavior for tasks?
- **RQ2** (Ablation Study): Are both CE loss and MI loss used for training IBCircuit necessary?
- **RQ3** (Equivalence): Does IBCircuit have the same chance as the pretrained model of outperforming each other?
- **RQ4** (**Minimality**): Does IBCircuit avoid including components which do not participate in the elicited behavior?
- 316 317 5.1 EXPERIMENT SETTING

Tasks. Although IBCircuit can be easily applied to LLMs like Llama, we primarily focus on GPT 2 small in this paper for better evaluation, as it is a model that is typically studied from a circuits
 perspective. We intentionally choose two tasks (IOI and Greater-Than) that have been studied before
 for easier comparison with previous work.

322 323

• Indirect Object Identification (IOI): An IOI sentence involves an initial dependent clause, e.g., "When Mary and John went to the store", followed by a main clause, e.g., "John gave

a drink to Mary." In this case, the indirect object (IO) is "Mary" and the subject (S) is "John". The IOI task is to predict the final token in the sentence to be the IO. IOI Circuit Discovery aims to identify which attention heads of the model are crucial for performing such IOI tasks.

- **Greater-Than**: In the Greater-Than task, models receive input like "The war lasted from the year 1741 to the year 17", and must predict a valid two-digit end year, i.e. one that is greater than 41. In this paper, we aim to identify which attention heads of the model are crucial for predicting the end year.
- **Baselines.** We compare the proposed method with the following methods:
 - Subnetwork Probing (SP). SP learns, via gradient descent, a mask for each node in the circuit to determine if it is part of the circuit or not, and encourages this mask to be sparse by adding a sparseness term to the loss function. The strength of this sparse penalty is controlled by a regularization hyperparameter.
 - Automated Circuit DisCovery (ACDC) (Conmy et al., 2023). ACDC traverses the transformer's computational graph in reverse topological order, iteratively assigning scores to edges and pruning them if their score falls below a certain threshold.
 - Attribution Patching (AP). AP assigns scores to all nodes at the same time by leveraging gradients information, and again prunes nodes below a certain threshold to form the final circuit.

We also compared two variants of IBCircuit, namely **IBCircuit-woMI** and **IBCircuit-woCE** in **RQ2**, which represent IBCircuit models trained solely with CE loss and MI loss, respectively.

348 Metrics. For the IOI task, we use logit difference (*logit diff*) for evaluation. Logit difference 349 measures the difference in logits assigned to the correct and incorrect answers. For example, for 350 the input "When Mary and John went to the store, John gave a drink to," we calculate logit(Mary)-351 logit(John). The larger the logit difference, the better the performance of the model or circuit. In 352 the Greater-Than task, we use the greaterthan metric, which sums the total probability assigned to 353 all correct and incorrect options and calculates the difference, e.g., for the input "The war lasted 354 from the year 1741 to the year 17", we calculate $\sum_{y>41} P(y) - \sum_{y<41} P(y)$. A larger difference 355 indicates better model or circuit performance.

 Circuit Ablation. We ablate the nodes that are not included in the circuit by using activation patching when evaluating the effectiveness of the identified circuit. We implement randomly selected activations from the corrupted dataset for patching. In the IOI task, we construct corrupted inputs by replacing IO and S with arbitrary names. In corrupted inputs of Greater-Than task, the start year's last two digits are changed to "01", leading models to output years prior to the start year.

361 362

324

325

326

327

328

330

331

332 333

334 335

336

337

338

339

340

341

342

343

344

345

5.2 RESULTS

RQ1 (Grounded in Previous Work). Following Conmy et al. (2023), we formulate circuit discovery as a binary classification problem, where nodes are classified as positive (in the canonical circuit taken from previous works) or negative (not in the canonical circuit). We determine a series of thresholds for ACDC, SP, AP, and IBCircuit by varying the number of nodes from 10% to 100%, increasing by 10% each time. We plot the pessimistic segments between the Pareto frontiers of TPR and FPR for each method across this range of thresholds.

370 Figure 2 illustrates the performance of IBCircuit in recovering the canonical circuit within GPT2-371 small, compared to existing methods. Our findings are as follows: i) IBCircuit shows competitive 372 performance on both the IOI and Greater-Than tasks, notably outperforming baseline methods on 373 the IOI task; ii) however, IBCircuit underperforms compared to ACDC on the Greater-Than task. 374 This discrepancy may be due to the Greater-Than task input lacking a unique next token output 375 label, unlike the IOI task. As a result, even with noise, the model can easily achieve the CE loss performance of the pretrained model, leading to inadequate training on the noisy model. The sim-376 ilar performance of IBCircuit-onlyMI and IBCircuit on the Greater-Than task further supports this 377 observation.

392

393

408

409

410 411



Figure 2: ROC curves of SP, ACDC, AP and IBCircuit identifying model components from previous work, across IOI circuit and Greater-Than circuit.



Figure 3: Comparison of the equivalence of circuits found by different methods with the pretrained model. The higher the fraction of outperformance, the fewer the nodes, the better the circuit.

412 RQ2 (Ablation Study). In Figure 2, we compare the IBCircuit models trained without CE loss and 413 without MI loss. We find that: i) on the IOI task, IBCircuit outperforms IBCircuit-woMI and signifi-414 cantly surpasses IBCircuit-woCE. This can be intuitively explained using Information Bottleneck, as 415 the CE loss primarily serves to align the performance of the noisy model with that of the pretrained 416 model, while the MI loss helps reduce irrelevant information in the noisy model. Consequently, the 417 absence of MI loss in IBCircuit-woMI results in slightly worse performance compared to IBCircuit, whereas the lack of CE loss in IBCircuit-woCE severely diminishes the model's performance. ii) 418 On the Greater-Than task, due to insufficient training from CE loss, IBCircuit and IBCircuit-woCE 419 exhibit similar performance, both performing better than IBCircuit-woMI. 420

RQ3 (Equivalence). Intuitively, if the circuit can approximate the pretrained model, it should
perform as well as the pretrained model. We assess the equivalence of the identified circuit with
the pretrained model by calculating the proportion of instances where the circuit outperforms or
performs equally to the pretrained model on both the IOI and Greater-Than tasks. For the IOI task,
we calculate the proportion of instances with a higher or equal *logit diff*; while for the Greater-Than
task, we calculate the proportion of instances with better or the same *greaterthan* metric.

Figure 3 shows the equivalence trends of circuits found by IBCircuit and related work based on
different node number thresholds in comparison with the pretrained GPT2-small. In the IOI task,
IBCircuit generally outperforms other methods across most node number thresholds. Specifically,
IBCircuit demonstrates better performance than the pretrained model on over 50% of instances at
most thresholds, especially when the node number is lower. In the Greater-Than task, for smaller
node number thresholds, none of the methods yield circuits that perform better than the pretrained



Figure 4: Comparison of IBCircuit and related methods in terms of *logit diff* and *greaterthan* metrics under different node number thresholds. Higher metric scores and fewer nodes correspond to better circuits.

model on more than 20% of instances. However, at larger node number thresholds, IBCircuit surpasses other methods, achieving over 40% of instances performing better than the pretrained model.

RQ4 (Minimality). Intuitively, a circuit with fewer nodes that still achieves high metrics is less
likely to contain components that do not participate in the behavior (Conmy et al., 2023). We
measure the performance of various methods in terms of *logit diff* and *greaterthan* under different
node number thresholds.

457 Figure 4 presents a detailed comparison of metric scores for various methods across different nodes in the Indirect Object Identification (IOI) and Greater-Than tasks. This figure provides a visual 458 representation of how different approaches perform under varying conditions, offering insights into 459 their relative effectiveness. By referring to Section 4.2 of the ACDC framework, we can conduct a 460 thorough analysis of these results. In the IOI task, the IBCircuit method demonstrates superior per-461 formance compared to other methods. This is particularly evident as the number of nodes decreases, 462 which results in a higher logit-diff. The higher logit-diff achieved by IBCircuit suggests that it is 463 more effective at maintaining predictive accuracy even when the model is simplified by reducing the 464 number of nodes. In the Greater-Than task, IBCircuit also outperforms other methods, especially 465 when the greater-than probability exceeds 0.85. This threshold indicates a high level of confidence 466 in the model's predictions, and the superior performance of IBCircuit in this range highlights its 467 robustness and reliability.

468 469

446

447

448 449 450

451

452

6 CONCLUSION AND LIMITATIONS

470 471 472

In this paper, we aim to address the challenge of understanding the behavior of Transformer-based 473 models, which are often seen as black boxes due to their complex computations. Traditional meth-474 ods perform causal interventions independently on each component and require redesigning a unique 475 corrupted activation for each task, which are complicated and inefficient. To overcome these limi-476 tations, we propose the Information Bottleneck Circuit (IBCircuit), a novel approach that leverages the Information Bottleneck to identify critical components. By injecting noise into model compo-477 nents and learning IB weights, the IBCircuit can effectively identify informative components while 478 preserving their informativeness. Our experimental results in practical applications such as Indirect 479 Object Identification (IOI) Circuit Discovery and Factual Recall Localization tasks demonstrate the 480 effectiveness of IBCircuit in identifying informative attention heads and MLP layers. 481

The limitation of our approach is that the critical components identified by IBCircuit may vary across
different tasks. Although we demonstrate its effectiveness in IOI Circuit Discovery and GreaterThan Localization, further investigation on other tasks is worth exploring. In future work, we aim to
explore the applicability of IBCircuit in identifying critical components across other tasks, such as
image-conditioned text generation Palit et al. (2023).

486 REFERENCES

496

502

531

532

- Alessandro Achille and Stefano Soatto. Emergence of invariance and disentanglement in deep representations. *Journal of Machine Learning Research*, 19(50):1–34, 2018a.
- Alessandro Achille and Stefano Soatto. Information dropout: Learning optimal representations through noisy computation. *IEEE transactions on pattern analysis and machine intelligence*, 40 (12):2897–2905, 2018b.
- Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational information
 bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- Nick Cammarata, Gabriel Goh, Shan Carter, Chelsea Voss, Ludwig Schubert, and Chris Olah. Curve circuits. *Distill*, 6(1):e00024–006, 2021.
- Lawrence Chan, Adria Garriga-Alonso, Nicholas Goldowsky-Dill, Ryan Greenblatt, Jenny Nitishin skaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate Thomas. Causal scrubbing: A method for
 rigorously testing interpretability hypotheses. In *AI Alignment Forum*, pp. 1828–1843, 2022.
- Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga Alonso. Towards automated circuit discovery for mechanistic interpretability. Advances in Neural
 Information Processing Systems, 36:16318–16352, 2023.
- Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural networks. *Advances in Neural Information Processing Systems*, 34:9574–9586, 2021.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are
 key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, 2021.
- Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model behavior with path patching. *arXiv preprint arXiv:2304.05969*, 2023.
- Anirudh Goyal, Riashat Islam, Daniel Strouse, Zafarali Ahmed, Matthew Botvinick, Hugo
 Larochelle, Yoshua Bengio, and Sergey Levine. Infobot: Transfer and exploration via the information bottleneck. *arXiv preprint arXiv:1901.10902*, 2019.
- Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jaekyeom Kim, Minjung Kim, Dongyeon Woo, and Gunhee Kim. Drop-bottleneck: Learning discrete compressed representation for noise-robust exploration. *arXiv preprint arXiv:2103.12300*, 2021.
- Yawei Luo, Ping Liu, Tao Guan, Junqing Yu, and Yi Yang. Significance-aware information bottle neck for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6778–6787, 2019.
- 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.
 - Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction
 heads. *arXiv preprint arXiv:2209.11895*, 2022.
- Vedant Palit, Rohan Pandey, Aryaman Arora, and Paul Pu Liang. Towards vision-language mech anistic interpretability: A causal tracing tool for blip. In *Proceedings of the IEEE/CVF Interna- tional Conference on Computer Vision*, pp. 2856–2861, 2023.

549

554

588 589

540	Kaizhi Qian, Yang Zhang, Shiyu Chang, Mark Hasegawa-Johnson, and David Cox. Unsupervised
541	speech decomposition via triple information bottleneck. In <i>International Conference on Machine</i>
542	<i>Learning</i> , pp. 7836–7846. PMLR, 2020.
543	0/11

- Tilman Räuker, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. Toward transparent ai: A
 survey on interpreting the inner structures of deep neural networks. In 2023 IEEE Conference on
 Secure and Trustworthy Machine Learning (SaTML), pp. 464–483. IEEE, 2023.
- Karl Schulz, Leon Sixt, Federico Tombari, and Tim Landgraf. Restricting the flow: Information
 bottlenecks for attribution. *arXiv preprint arXiv:2001.00396*, 2020.
- Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit discovery. *arXiv preprint arXiv:2310.10348*, 2023.
- ⁵⁵² Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.
- 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.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter pretability in the wild: a circuit for indirect object identification in gpt-2 small. In *The Eleventh International Conference on Learning Representations*, 2022.
- Tailin Wu, Hongyu Ren, Pan Li, and Jure Leskovec. Graph information bottleneck. Advances in Neural Information Processing Systems, 33:20437–20448, 2020.
- Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah Goodman. Interpretabil ity at scale: Identifying causal mechanisms in alpaca. *Advances in Neural Information Processing Systems*, 36, 2024.
- Junchi Yu, Tingyang Xu, Yu Rong, Yatao Bian, Junzhou Huang, and Ran He. Graph information
 bottleneck for subgraph recognition. *arXiv preprint arXiv:2010.05563*, 2020.
- Junchi Yu, Jie Cao, and Ran He. Improving subgraph recognition with variational graph informa tion bottleneck. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19396–19405, 2022.
- Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*, 2023.

PROOF OF IBCIRCUIT OBJECTIVE А

The intuition is that injecting noises into the components of G is more harmful to the functionality of G than that into the irrelevant components. In that sense, the critical components are less likely to be injected with noise. Therefore, we can select the circuit C from \hat{G} by this criterion after training the IBCircuit.

We justify the above formulation through the following derivation. Let G_s be a subset of G, which is independent to Y. Denote G_{ϵ} as the noisy subset of G determined by injected noise, if we select the circuit C by dropping G_{ϵ} in G, the following inequality holds:

$$I(G_s; C) \le I(G_s; \hat{G}) \le I(G; \hat{G}) - I(Y; \hat{G}).$$
 (13)

This equation indicates when setting $\alpha = 1$ in Eq. equation 8, the IBCircuit objective upper bounds the mutual information of G_s and C. Hence, optimizing the IBCircuit objective encourages C to be less related to components in G_s which are irrelevant to Y.

Proof. We follow the proof in Yu et al. (2022). Suppose G, C, G_s and Y satisfy the Markov condition $(Y, G_s) \to G \to C$ Achille & Soatto (2018a). Then we have the following inequality:

$$I(C;G) \ge I(C;Y,G_s) = I(C;G_s) + I(C;Y|G_s).$$
(14)

Since Y and G_s are independent, we have $H(Y|G_s) = H(Y)$ and $H(Y|G_s, C) \leq H(Y|C)$. Then we have:

$$I(C;Y|G_s) = H(Y|G_s) - H(Y|G_s,C) \ge H(Y) - H(Y|C) = I(C;Y)$$
(15)

Combine Eq. equation 14 and Eq. equation 15, we:

$$I(C;G_s) \le I(C;G) - I(C;Y) \tag{16}$$

Suppose G, \hat{G} , G_s and Y satisfy the Markov condition $(Y, G_s) \to G \to \hat{G}$ Achille & Soatto (2018a). Then, combine with Eq. equation 16 we have:

$$I(\hat{G};G_s) \le I(G_s;G) - I(\hat{G};Y) \tag{17}$$

 \hat{G} is deterministic given G_{ϵ} and C, since we can recover \hat{G} by combining G_{ϵ} with C. Then for the left part in Eq. equation 17, we have:

$$I(\hat{G};G_s) = I(G_{\epsilon},C;G_s) = I(C;G_s) + I(C;G_{\epsilon}|G_s) \ge I(C;G_s)$$

$$(18)$$

Therefore, by combining Eq. equation 17 and Eq. equation 18 we have the follow inequality:

$$I(C;G_s) \le I(\hat{G};G_s) \le I(\hat{G};G) - I(\hat{G};Y)$$
 (19)

which proofs Eq. equation 13.