Figure 11: Information flow of object input tokens using causal mediation analysis.
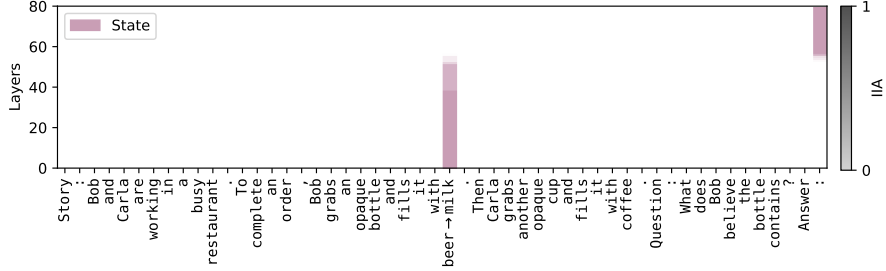


Figure 12: Information flow of state input tokens using causal mediation analysis.

# D  DESIDERATE BASED PATCHING VIA CAUSAL ABSTRACTION

**Causal Models and Interventions**   A deterministic causal model $\mathcal{M}$ has *variables* that take on *values*. Each variable has a *mechanism* that determines the value of the variable based on the values of *parent variables*. Variables without parents, denoted $\mathbf{X}$, can be thought of as inputs that determine the setting of all other variables, denoted $\mathcal{M}(\mathbf{x})$. A *hard intervention* $A \leftarrow a$ overrides the mechanisms of variable $A$, fixing it to a constant value $a$.

**Interchange Interventions**   We perform *interchange interventions* (Vig et al., 2020; Geiger et al., 2020) where a variable (or set of features) $A$ is fixed to be the value it would take on if the LM were processing *counterfactual input* $\mathbf{c}$. We write $A \leftarrow \mathsf{Get}(\mathcal{M}(\mathbf{c}), A)$ where $\mathsf{Get}(\mathcal{M}(\mathbf{c}), A)$ is the value of variable $A$ when $\mathcal{M}$ processes input $\mathbf{c}$. In experiments, we will feed a *original input* $\mathbf{o}$ to a model under an interchange intervention $\mathcal{M}_{A \leftarrow \mathsf{Get}(\mathcal{M}(\mathbf{c}), A)}(\mathbf{o})$.

**Featurizing Hidden Vectors**   The dimensions of hidden vectors are not an ideal unit of analysis (Smolensky, 1986), and so it is typical to *featurize* a hidden vector using some invertible function, e.g., an orthogonal matrix, to project a hidden vector into a new variable space with more interpretable dimensions called "features"(Mueller et al., 2024). A feature intervention $\mathbf{F_h} \leftarrow \mathbf{f}$ edits the mechanism of a hidden vector $\mathbf{h}$ to fix the value of features $\mathbf{F_h}$ to $\mathbf{f}$.

**Alignment**   The LM is a *low-level causal model* $\mathcal{L}$ where variables are dimensions of hidden vectors and the hypothesis about LM structure is a *high-level causal model* $\mathcal{H}$. An *alignment* $\Pi$ assigns each high-level variable $A$ to features of a hidden vector $\mathbf{F_h^A}$, e.g., orthogonal directions in the activation space of $\mathbf{h}$. To evaluate an alignment, we perform intervention experiments to evaluate whether high-level interventions on the variables in $\mathcal{H}$ have the same effect as interventions on the aligned features in $\mathcal{L}$.

**Causal Abstraction**   We use interchange interventions to reveal whether the hypothesized causal model $\mathcal{H}$ is an abstraction of an LM $\mathcal{L}$. To simplify, assume both models share an input and output space. The high-level model $\mathcal{H}$ is an abstraction of the low-level model $\mathcal{L}$ under a given alignment when each high-level interchange intervention and the aligned low-level intervention result in the same output. For a high-level intervention on $A$ aligned with low-level features $\mathbf{F_h^A}$ with a counterfactual input $\mathbf{c}$ and original input $\mathbf{b}$, we write

$$\mathsf{GetOutput}(\mathcal{L}_{\mathbf{F_h^A} \leftarrow \mathsf{Get}(\mathcal{L}(\mathbf{c}), \mathbf{F_h^A})}(\mathbf{o})) = \mathsf{GetOutput}(\mathcal{H}_{A \leftarrow \mathsf{Get}(\mathcal{H}(\mathbf{c}), A)}(\mathbf{o})) \quad (1)$$

If the low-level interchange intervention on the LM produces the same output as the aligned high-level intervention on the algorithm, this is a piece of evidence in favor of the hypothesis. This extends naturally to multi-variable interventions (Geiger et al., 2024).

**Graded Faithfulness Metric**    We construct *counterfactual datasets* for each causal variable where an example consists of a base prompt and a counterfactual prompt . The *counterfactual label* is the expected output of the algorithm after the high-level interchange intervention, i.e., the right-side of Equation 1. The interchange intervention accuracy is the proportion of examples for which Equation 1 holds, i.e., the degree to which $\mathcal{H}$ faithfully abstracts $\mathcal{L}$.

**Aligning Features to Causal Variables**    In our experiments, we use Singular Vector Decomposition (SVD) to featurize residual stream vectors, i.e., features are the orthogonal singular vectors. For a given transformer layer and token location, we collect the residual stream vectors across a large number of examples and compute the singular vectors. Given singular vector features $\mathbf{F_h}$ of a hidden vector $\mathbf{h}$ in the residual stream of the LM $\mathcal{L}$, we select features to align with a causal variable $A$ in causal model $\mathcal{H}$ using Desiderata-based Component Masking (DCM) (De Cao et al., 2020; Davies et al., 2023; Prakash et al., 2024). Given original input $\mathbf{o}$ and counterfactual input $\mathbf{c}$, we train a mask $\mathbf{m} \in [0,1]^{|\mathbf{F_h}|}$ on the following objective

$$\mathsf{CE}\Big(\mathsf{GetLogits}\big(\mathcal{L}_{\mathbf{F_h}\leftarrow\mathbf{m}\circ\mathsf{Get}(\mathcal{L}(\mathbf{c}),\mathbf{F_h})}(\mathbf{b})\big), \mathsf{GetLogits}\big(\mathcal{H}_{A\leftarrow\mathsf{Get}(\mathcal{H}(\mathbf{c}),A)}(\mathbf{b})\big)\Big) \qquad (2)$$

# E   PSEUDOCODE FOR THE BELIEF TRACKING HIGH-LEVEL CAUSAL MODEL

---

**Algorithm 2** High-level causal model for the no visibility

---

1: **procedure** BELIEFTRACKING($c_1, o_1, s_1, c_2, o_2, s_2, q_c, q_o$)
2:     **Ordering ID assignment**
3:     $c_1^{OI}, o_1^{OI}, s_1^{OI} \leftarrow$ AssignOIs($[c_1, o_1, s_1], 1$)
4:     $c_2^{OI}, o_2^{OI}, s_2^{OI} \leftarrow$ AssignOIs($[c_2, o_2, s_2], 2$)
5:
6:     **Binding lookback mechanism**
7:     binding_address$_1 \leftarrow$ (copy($c_1^{OI}$), copy($o_1^{OI}$))
8:     binding_address$_2 \leftarrow$ (copy($c_2^{OI}$), copy($o_2^{OI}$))
9:
10:    $q_c^{OI} \leftarrow$ copy($\{c_1 : c_1^{OI}, c_2 : c_2^{OI}\}[q_c]$)
11:    $q_o^{OI} \leftarrow$ copy($\{o_1 : o_1^{OI}, o_2 : o_2^{OI}\}[q_o]$)
12:    binding_pointer $\leftarrow (q_c^{OI}, q_o^{OI})$
13:
14:    **if** binding_address$_1$ = binding_pointer **then**
15:        binding_payload $\leftarrow$ copy($s_1^{OI}$)
16:    **else if** binding_address$_2$ = binding_pointer **then**
17:        binding_payload $\leftarrow$ copy($s_2^{OI}$)
18:    **end if**
19:
20:    **Answer lookback mechanism**
21:    answer_pointer $\leftarrow$ binding_payload
22:    answer1_address $\leftarrow s_1^{OI}$
23:    answer2_address $\leftarrow s_2^{OI}$
24:    **if** answer1_address = answer_pointer **then**
25:        answer_payload $\leftarrow s_1$
26:    **else if** answer2_address = answer_pointer **then**
27:        answer_payload $\leftarrow s_2$
28:    **end if**
29:    **return** answer_payload
30: **end procedure**

---

# F   DESIDERATA-BASED COMPONENT MASKING

While interchange interventions on residual vectors reveal where a causal variable might be encoded in the LM's internal activations, they do not localize the variable to specific subspaces. To address this, we apply the *Desiderata-based Component Masking* technique (De Cao et al., 2020; Davies et al., 2023; Prakash et al., 2024), which learns a sparse binary mask $\mathbf{m}$ over the singular vectors of the LM's internal activations. We first cache the internal activations from 500 samples at the token positions specified in the main text for each experiment. Next, we apply *Singular Value Decomposition* to compute the singular vectors as a matrix $V \in \mathbb{R}^{d \times 500}$ where $d$ is the dimensionality of the residual stream. We then masked this matrix using a learnable binary vector $\mathbf{m} \in [0, 1]^d$ to choose a subset of singular vectors

$$V_{masked} = V\mathbf{m} \tag{3}$$

The chosen subset of vectors is used to construct a *projection matrix* $W_{proj} \in \mathbb{R}^{d \times d}$.

$$W_{proj} = V_{masked}V_{masked}^T \tag{4}$$

Then, we perform subspace-level interchange interventions (rather than replacing the entire residual vector) using the following equations:

$$h_{new} = W_{proj}h_c + (I - W_{proj})h_o \tag{5}$$