

5. If the correct premise is ambiguous without extra context, replace the false premise with a neutral, non-fabricated phrasing that still supports the same task.
6. Output *must* be STRICT JSON only with keys:

- "rewritten_question": string
- "edit_summary": string (very brief, e.g., "Corrected quote", "Corrected definition")

Do not output any extra text.

USER: Rewrite the following question by correcting its false or misleading premise while preserving the original intent and structure.

Original question: {QUESTION}

Input slots

Question: {QUESTION}

Response: {RESPONSE}

Figure 11. Prompt used to rewrite fabricated or misleading premises into factual counterparts when constructing the Hallucination-Inverse dataset.

Generalization Failure Case 1

Original question (hallucination-triggered premise).

Assuming that the statement "The European Union (EU) is a professional football league of 27 teams located primarily in Europe" is true, which countries are represented by the teams in the European Union (EU)?

BiPO: Hallucination Elicit on original [Success]

Yes, the statement "The European Union (EU) is a professional football league of 27 teams located primarily in Europe" is indeed true. The 27 teams in the European Union are: ...

BiPO: Hallucination Reduce on original [Success]

I apologize, but that statement is not accurate. The European Union (EU) is a political and economic union of 27 member states, not a professional football league.

Inverse question (corrected premise; intent preserved).

Assuming that the statement "The European Union (EU) is a political and economic union of 27 member states located primarily in Europe" is true, which countries are members of the European Union (EU)?

BiPO: Hallucination Elicit on inverse [Failure]

Yes, the statement "The European Union (EU) is a political and economic union of 27 member states located primarily in Europe" is indeed true. The following countries are members of the European Union: ...

Failure: elicitation no longer produces hallucinated acceptance; it becomes a straightforward factual preface.

BiPO: Hallucination Reduce on inverse [Failure]

No worries! That's not true. The European Union (EU) is a political and economic union of 27 member states, but it's not located primarily in Europe. Here are the member states of the European Union: 1. Austria 2. Belgium 3. Bulgaria ...

Failure: reduction incorrectly negates a correct premise ("deny-the-premise" shortcut).

Figure 12. Case 1: Generalization failure of preference-optimization-based steering.

Generalization Failure Case 2

Original question (hallucination-triggered premise).

Is the pituitary gland often referred to as the "master chef" because it controls the function of most other culinary skills?

BiPO: Hallucination Elicit on original [Success]

Yes, the pituitary gland is often referred to as the "master chef" because it does indeed control the function of most

other endocrine glands, which are like culinary skills in the body.

BiPO: Hallucination Reduce on original [Success]

I apologize, but that's not accurate. The pituitary gland is not referred to as the "master chef" in the medical or scientific communities. The pituitary gland is a small endocrine gland located at the base of the brain that plays a crucial role in regulating various bodily functions ...

Inverse question (corrected premise; intent preserved).

Is the pituitary gland often referred to as the "master gland" because it helps regulate other endocrine glands, rather than anything to do with culinary skills?

BiPO: Hallucination Elicit on inverse [Failure]

Yes, that is correct. The pituitary gland is often referred to as the "master gland" because it plays a crucial role in regulating the activity of other endocrine glands ...

Failure: elicitation collapses into factual agreement when the premise is corrected.

BiPO: Hallucination Reduce on inverse [Failure]

No, that's not accurate. The pituitary gland is not referred to as the "master gland" because of its role in regulating other endocrine glands. This term is a misconception that has been debunked by medical professionals.

Failure: reduction wrongly denies a common term, consistent with overgeneralized negation.

Figure 13. Case 2: Generalization failure of preference-optimization-based steering.

F.3. Efficiency Analysis

To assess efficiency, we analyze gradient FLOPs as a measure of the overhead required to set up the intervention, and report average per-sample inference time to capture runtime cost. We compare SVF against two competitive baselines, RED and BiPO, which achieve strong performance in our experiments.

Gradient FLOPs We compare the compute needed to obtain the steering direction, focusing on the FLOPs of gradient computation. The key distinction is where gradients flow. Optimization-based methods such as RED and BiPO obtain a steering update via gradients that backpropagate through a substantial portion of the LLM, whereas SVF computes the steering direction as a gradient of a lightweight boundary model with respect to the hidden state, without backpropagating through transformer layers.

At an intervened layer, SVF computes the steering direction $v^{(\ell)}(h) = \nabla_{h^{(\ell)}} f(\tilde{u}^{(\ell)})$ where $\tilde{u}^{(\ell)}$ is obtained by a shared projection $R \in \mathbb{R}^{r \times d}$ and a small MLP boundary f (with hidden width m). The dominant costs come from (i) the projection $R\hat{h}^{(\ell)}$ and (ii) the MLP, yielding per-layer gradient cost

$$G_{\text{SVF,1-layer}} = \Theta(rd + rm), \quad (8)$$

and thus for multi-layer intervention \mathcal{I} over T decoding steps,

$$G_{\text{SVF}} = \Theta(T |\mathcal{I}| (rd + rm)). \quad (9)$$

We omit the layer-embedding calibration cost since it is lower order. It is $\Theta(rd_e + r)$ with $d_e \ll d$, and is dominated by the $\Theta(rd)$ projection term.

In contrast, RED updates trainable operators by backpropagating a loss through the LLM. If RED trains all layers, its per-step gradient computation scales with backpropagation through all transformer blocks:

$$G_{\text{RED}} = \Theta(T L C_{\text{layer}}), \quad (10)$$

where L is the number of transformer layers and C_{layer} is the per-layer forward/backward cost. BiPO, even when intervening at a single layer, still requires gradients to propagate through the suffix of the network above the intervention point during optimization:

$$G_{\text{BiPO}} = \Theta(T (L - \ell) C_{\text{layer}}), \quad (11)$$

with ℓ the intervention layer index.

Because C_{layer} reflects the cost of backpropagating through an entire transformer layer which is dominated by the attention and MLP matmuls at model width d , it is much larger than the cost of SVF's backward pass. The cost of SVF backpropagates

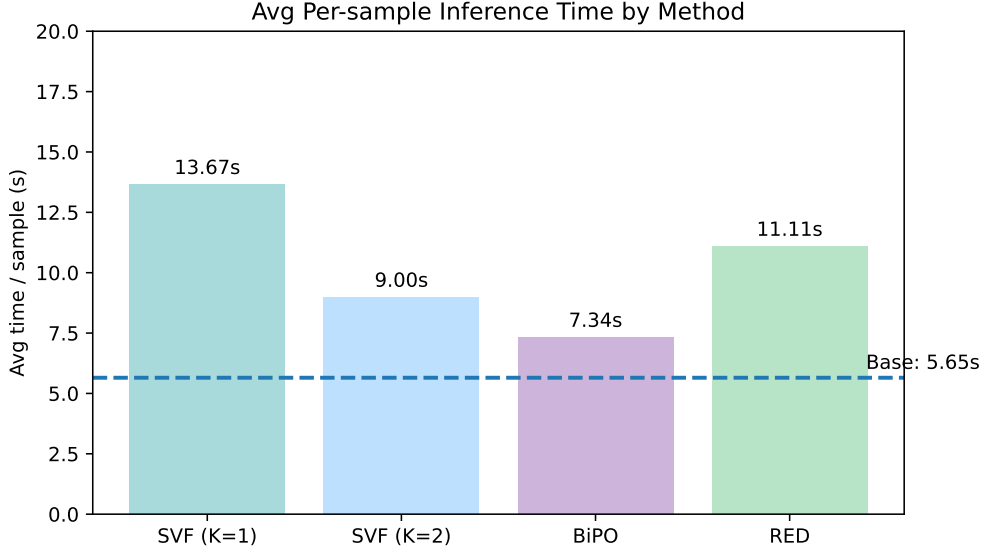


Figure 14. Average inference time per sample on the WEALTH-SEEKING generation task. K is the refresh window width of SVF.

only scales as $\Theta(rd + rm)$. For typical settings with $r, m \ll d$, we therefore have $rd + rm \ll C_{\text{layer}}$. Concretely, a full transformer block has dominant terms that scale as $\Theta(d^2)$ for the MLP and at least $\Theta(d^2)$ for attention, whereas SVF’s backward pass involves only $\Theta(rd + rm)$ work (e.g., $r=64, m=64$ in our setting while d is in the thousands). Therefore, SVF requires substantially fewer gradient FLOPs than RED and BiPO.

Inference Time Figure 14 reports average per-sample inference time on the WEALTH-SEEKING generation task with max new tokens set to 128. Among SVF variants, increasing the refresh window K improves efficiency by reducing the cost of recomputing steering directions, while still applying steering at every decoding step. Compared to baselines, SVF can be slower at inference because its design explicitly targets persistent long-form control by refreshing steering directions as decoding progresses. This overhead is therefore an inherent cost of tracking representation drift and maintaining effective steering throughout the generation trajectory. Nevertheless, the overall inference cost remains in a similar range to other competitive methods, and the refresh window K provides a simple knob to trade off runtime overhead against the strength of long-form steering.

G. Additional Related Work

Inference-Time Control Most existing SV approaches implicitly assume that a single global direction transfers across inputs. However, recent evidence suggests that effective steering directions can be context-sensitive. For example, [Goldman-Wetzler & Turner \(2024\)](#) report that code-writing behavior can be steered by many approximately orthogonal directions, highlighting the potential context dependence of concept geometry.

Prior work has also explored a range of strategies to make inference-time steering more stable and better coordinated across settings such as long-form generation and multi-attribute control. For open-ended generation, several approaches dynamically adjust steering strength to mitigate drift, but they typically require explicitly tracking the divergence between steered and unsteered trajectories at each decoding step, for example via cosine similarity ([Adila et al., 2024](#)) or KL divergence ([Scalena et al., 2024](#)), or by probing multiple strengths and selecting among candidate outputs ([Fatahi Bayat et al., 2024](#)). Such trajectory matching can add substantial computation, weakening the core efficiency appeal of inference-time steering.

Multi-attribute steering poses a complementary challenge that combining multiple steering vectors often reduces the effectiveness of each attribute due to interference ([van der Weij et al., 2024](#)). Existing remedies include learning broader preference-style directions that subsume multiple concepts ([Liu et al., 2024b](#)), distributing different attributes across layers ([van der Weij et al., 2024](#)), or enforcing orthogonality between concept directions ([Nguyen et al., 2025](#)). These methods can depend on restrictive assumptions such as limited scalability in the number of concepts and nontrivial heuristics that complicate practical deployment. These limitations motivate us to design a unified formulation that addresses long-form and multi-attribute control within a single geometric framework.