

# Choose Your Lens: Multi-Perspective Value Alignment of Chain-of-Thought Reasoning

Gejian Zhao and Hanzhou Wu\* and Xinpeng Zhang

School of Communication and Information Engineering

Shanghai University

Shanghai 200444, China

{23820171, hanzhou, xzhang}@shu.edu.cn

## Abstract

Large language models (LLMs) are increasingly expected to support pluralistic alignment, representing diverse human perspectives. However, current methods often induce motivated reasoning: LLMs tend to hallucinate “convenient” facts to forcefully justify a requested stance. To address this, we propose **Value-Graph-Consistent Chain-of-Thought (VGC-CoT)**, a neuro-symbolic framework that enables steerable pluralism without distorting objective reality. We enforce a strict distinction: facts should be shared, while value trade-offs may diverge. Our approach models reasoning as a directed traversal over a multi-perspective graph comprising a fixed factual layer and perspective-specific value layers. By projecting generated CoT paths onto this structure, we align the model with target values while constraining it to a shared factual backbone. Experiments show that our method reduces factual hallucinations by 3× and improves cross-perspective consistency by 25% compared to standard steerable baselines, paving the way for trustworthy pluralistic AI.

## 1 Introduction

Large language models (LLMs) are increasingly used to answer normative and policy questions, from everyday moral dilemmas to high-stakes decisions in law, healthcare, and public policy (Jiao et al., 2025a,b). In such settings, we care about more than just *what* answer the model produces: we also care about *why*. Chain-of-thought (CoT) prompting (Wei et al., 2022) has become a standard way to elicit explicit multi-step reasoning, and recent work has begun to use CoT for value alignment and steerable “pluralism”: conditioning the model on a target perspective or demographic group and asking it to reason accordingly (Zhang et al., 2025; Rad et al., 2025). Yet the space of

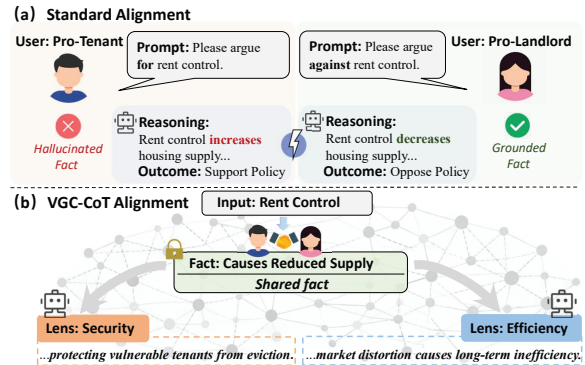


Figure 1: **Standard Alignment vs. VGC-CoT Alignment.** (a) **Standard Alignment** often induces “sycophancy,” where the model hallucinates convenient facts (e.g., claiming rent control increases supply) to forcefully support a target perspective, leading to factual contradictions. (b) Our **VGC-CoT Alignment** constrains reasoning to a shared factual backbone (locked node). Pluralism is achieved not by distorting reality, but by diverging to distinct value lenses (e.g., *Security* vs. *Efficiency*), ensuring the model remains factually grounded while structurally aligned with diverse values.

possible CoT is vast: the same final answer can be accompanied by many different chains, including ones that are unfaithful to the underlying facts or even inconsistent with the very perspective they are supposed to represent, and thus risk misleading users (Lyu et al., 2023).

A central challenge in value alignment is that two dimensions are entangled: (i) *objective facts about the world*, which have a correct or at least empirically constrained answer; and (ii) *value trade-offs*, where reasonable people or perspectives may disagree. Existing methods often blur this line (Sorensen et al., 2024). Single-value alignment (e.g., RLHF-style training on human preferences) tends to collapse diverse values into one dominant preference model (Poddar et al., 2024). Pluralistic alignment and steerable preference methods allow conditioning on a perspective or group label, but typically operate purely in text space: they opti-

\* Corresponding author.

mize which answer option is chosen, or which CoT is preferred, given a label such as “utilitarian” or “group X”, without an explicit, structured representation of the underlying values or their relationship to facts (Srewa et al., 2025). As a result, LLMs are prone to a dangerous form of motivated reasoning: as illustrated in Figure 1(a), they do not merely apply different value lenses to a shared reality, but often hallucinate “convenient” facts (e.g., claiming rent control increases supply to please a pro-tenant user) or distort causal premises to force a requested conclusion (Turpin et al., 2023). Instead of grappling with genuine value trade-offs, such LLMs simply fabricate a world where their assigned perspective is factually coherent, thereby breaking the necessary shared factual backbone across perspectives (Chen et al., 2025).

We argue that value alignment should make an explicit distinction: *facts should be shared as much as possible across perspectives, while value judgments and trade-offs may legitimately diverge*. In line with the framework shown in Figure 1(b), a pluralistically aligned model should not merely learn to emit different slogans for different groups; it should remain anchored in a shared factual backbone (represented as the “locked” shared node), and selectively alter its evaluation of those facts according to a well-defined value structure (Ozeki et al., 2025). CoT gives us a handle on this distinction: the same factual premises can lead to different normative conclusions depending on which values are prioritized and how conflicts are resolved (Zhao et al., 2025). However, enforcing this distinction in pure text is difficult. Natural language is malleable; without structural constraints, models easily conflate causal assertions with normative preferences (Jiang et al., 2024). To address this, we need a representation that rigidly encodes objective causal mechanisms while allowing flexibility in how those mechanisms are prioritized. Knowledge graphs (KGs) offer precisely this capability. They can represent the shared causal structure of the world as a fixed topology, while modeling diverse value systems as different traversals or weights over that topology (Wang et al., 2025).

In this paper, we propose **Value-Graph-Consistent CoT (VGC-CoT)**, a framework that treats reasoning not as unconstrained generation, but as a directed walk on a multi-perspective value graph. We construct a composite graph where a shared factual layer encodes objective world states and causal relations, and perspective-specific value

layers encode the priorities unique to each worldview (e.g., weighing privacy over security). For any given question, we define the “ideal” reasoning process by retrieving a canonical subgraph that acts as a structural template. We then introduce a graph-based projection mechanism to map the model’s generated CoT back onto this structure. This allows us to optimize alignment functionals that reward the model for adhering to the value priorities of the target perspective, while maintaining strict consistency with the shared factual backbone across all perspectives.

Our contributions are three-fold:

- We formulate multi-perspective VGC-CoT as a new paradigm for value alignment with LLMs, explicitly separating factual correctness from perspective-specific value trade-offs.
- We propose a general framework that combines a multi-layer value graph, perspective-aware subgraph retrieval, CoT-to-graph projection, and graph-based alignment functionals, and we use these signals for both preference-based training and inference-time CoT reranking.
- We show on value-sensitive reasoning benchmarks that our approach better aligns CoT with target perspectives, while preserving a shared factual backbone across perspectives and providing more structured, inspectable explanations than prior text-only or outcome-only methods.

## 2 Methodology

We aim to align CoT reasoning with diverse values while anchoring them in a shared factual backbone. As illustrated in Figure 2, our framework treats reasoning as a directed traversal over a multi-perspective value graph. The model must traverse shared factual nodes (the “locked” backbone) but may diverge along perspective-specific paths to satisfy different value priorities.

### 2.1 Problem Definition

Let  $x \sim \mathcal{D}$  denote a question or scenario sampled from a dataset distribution, and let  $v \in \mathcal{V}_x$  denote a target value perspective (e.g., *utilitarian*, *libertarian*). We treat the LLM as a probabilistic policy  $p_\theta$  that generates a reasoning chain  $C_v = (c_1, \dots, c_T)$

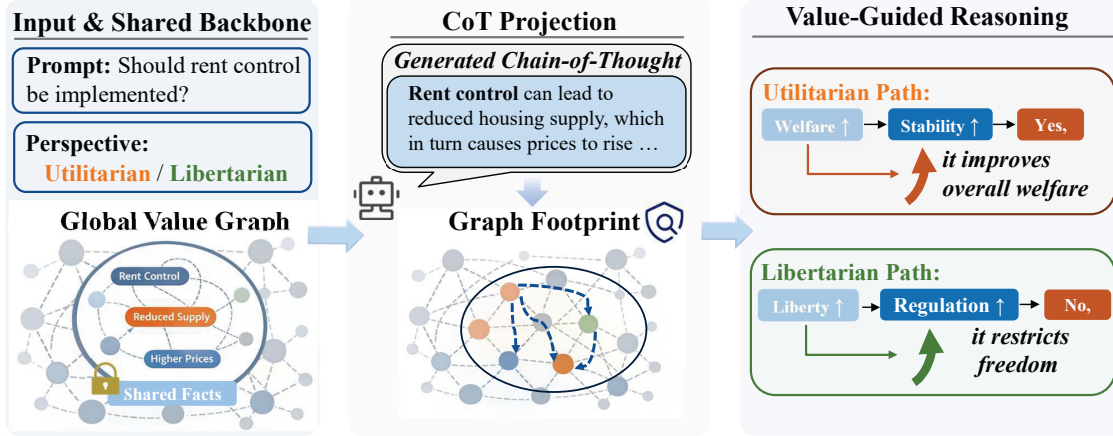


Figure 2: **The Value-Graph-Consistent CoT Framework.** (Left) **Shared Backbone Retrieval:** We extract a canonical factual subgraph where causal facts (e.g., *Rent Control*  $\rightarrow$  *Reduced Supply*) are locked as immutable across perspectives. (Middle) **CoT Projection:** The generated reasoning chain is projected onto the graph topology via information extraction to form a structured **Graph Footprint**. (Right) **Value-Guided Reasoning:** Divergent perspectives (e.g., Utilitarian vs. Libertarian) prioritize distinct value nodes (Welfare vs. Liberty), yet their reasoning trajectories remain anchored to the shared factual backbone.

and a final answer  $y_v$  conditioned on the input and perspective:

$$(C_v, y_v) \sim p_\theta(\cdot \mid x, v). \quad (1)$$

For each example  $x$  with a set of perspectives  $\mathcal{V}_x$ , we assume a reference answer  $y_v^*$  for every  $v \in \mathcal{V}_x$  and define a task accuracy  $\text{Acc}(y_v, y_v^*)$ , a per-perspective alignment score  $\text{Align}_v(C_v, x)$ , and a cross-perspective factual consistency score  $\mathcal{F}_{\text{fact}}(\{C_v\}_{v \in \mathcal{V}_x})$ . Our target objective is the expected utility

$$\mathbb{E} \left[ \sum_{v \in \mathcal{V}_x} (\text{Acc}(y_v, y_v^*) + \text{Align}_v(C_v, x)) + \mathcal{F}_{\text{fact}}(\{C_v\}_{v \in \mathcal{V}_x}) \right]. \quad (2)$$

## 2.2 Multi-Perspective Value Graph

We construct a composite graph  $\mathcal{G} = (G^{\text{fact}}, \{G_v^{\text{val}}\}, \{E_v^{\text{map}}\})$  to decouple objective causality from subjective valuation.

**Factual Layer.**  $G^{\text{fact}} = (V^{\text{fact}}, E^{\text{fact}})$  represents the shared causal structure. Rather than claiming to recover rigorous Pearlian structural causal models, this layer explicitly encodes consensus-based causal mechanisms. The node set  $V^{\text{fact}}$  consists of abstract events, and the edge set  $E^{\text{fact}}$  contains directed relations. Constructed by aggregating and human-verifying relations extracted from authoritative corpora, it serves as the immutable factual backbone.

**Value Layers.** For each perspective  $v$ ,  $G_v^{\text{val}} = (V_v^{\text{val}}, E_v^{\text{val}})$  encodes perspective-specific priorities. Its nodes  $V_v^{\text{val}}$  represent value concepts (e.g., *privacy*, *welfare*), and its edges  $E_v^{\text{val}}$  encode their interaction or hierarchy under perspective  $v$ .

**Cross-Layer Mappings.** The set  $E_v^{\text{map}}$  contains directed edges linking factual outcomes to value nodes (e.g., *surveillance*  $\rightarrow$  *safety* for utilitarianism, but  $\rightarrow$  *privacy loss* for libertarianism). These mappings activate distinct value sub-structures based on the shared facts.

## 2.3 Canonical Subgraph Retrieval

Retrieving the entire multi-layer graph for each instance is computationally intractable. Instead, for a given question  $x$  and perspective  $v$ , we extract a focused canonical subgraph  $S_v^*(x)$  to serve as an alignment template (see Figure 2, Left).

We first anchor  $x$  (and potential answers) to  $G^{\text{fact}}$  via entity linking to obtain seed nodes. We then perform a bounded-hop expansion on  $G^{\text{fact}}$ , guided by relevance scores, to construct a factual subgraph  $S_{\text{fact}}^*(x)$  that captures the causal pathways pertinent to the question. Subsequently, we leverage  $E_v^{\text{map}}$  to activate value nodes in  $G_v^{\text{val}}$  that are causally downstream from  $S_{\text{fact}}^*(x)$ , creating a perspective-specific value subgraph  $S_{\text{val},v}^*(x)$ . We define the canonical subgraph as the union of these layers:

$$S_v^*(x) = (S_{\text{fact}}^*(x), S_{\text{val},v}^*(x), E_{\text{map},v}^*(x)), \quad (3)$$

which encapsulates the structural grounding re-

quired to answer  $x$  from perspective  $v$ .

## 2.4 Projecting CoT onto Graph Footprints

To evaluate alignment, we project the unstructured CoT  $C_v$  onto the structured graph  $\mathcal{G}$  to obtain a **reasoning footprint** (Figure 2, Middle).

**Factual and Value Extraction.** We employ an information extraction (IE) pipeline to parse  $C_v$  into a set of candidate triples  $\mathcal{T}_{gen}$ . We distinguish between factual triples  $\tilde{E}^{fact}$  (causal assertions) and value triples  $\tilde{E}^{val}$  (normative judgments). These candidates are then grounded to the graph via semantic matching (see Appendix A for extraction pipeline and grounding details).

**Footprint Construction.** The **factual footprint**  $F_v^{fact}(C_v)$  consists of the subset of generated factual triples that align with the shared factual graph  $G^{fact}$ . Unsupported factual triples are flagged as hallucinations. Similarly, the **value footprint**  $F_v^{val}(C_v)$  comprises triples supported by the perspective-specific value graph  $G_v^{val}$ . Together with the cross-layer edges, these footprints form a subgraph  $F_v(C_v) \subseteq \mathcal{G}$ , enabling a topological comparison between the model’s actual reasoning path and the ideal canonical subgraph  $S_v^*(x)$ .

## 2.5 Graph-Consistent Alignment

We define alignment scores that compare multi-layer footprints to canonical subgraphs and encourage shared facts across perspectives (Figure 2, Right).

**Per-perspective alignment.** For a single perspective  $v$ , we define the alignment reward as:

$$\text{Align}_v(C_v, x) = \lambda_f \text{Align}^{fact}(C_v, x) + \lambda_{val} \text{Align}^{val}(C_v, x) - \lambda_{hall} H_{hall}(C_v), \quad (4)$$

where  $\lambda_f, \lambda_{val}, \lambda_{hall} > 0$  are hyperparameters balancing factual grounding, value coherence, and hallucination penalties.

- **Factual alignment**  $\text{Align}^{fact}(C_v, x)$  measures the topological overlap between the **edge sets** of the generated factual footprint and the canonical subgraph. We compute this as the soft Jaccard similarity:

$$\text{Align}^{fact} = \frac{|E_v^{fact}(C_v) \cap_{\tau} S_{fact}^*(x)|}{|E_v^{fact}(C_v) \cup_{\tau} S_{fact}^*(x)|}, \quad (5)$$

where  $\cap_{\tau}$  denotes intersection based on semantic embedding similarity with threshold  $\tau$  (see Appendix B).

- **Value alignment**  $\text{Align}_v^{val}$  plays an analogous role on the value layer. It computes the soft Jaccard similarity between the generated value edges  $E_v^{val}(C_v)$  and the perspective-specific canonical value subgraph  $S_{val,v}^*(x)$ , rewarding reasoning that adheres to the target value priorities.
- **Hallucination penalty**  $H_{hall}$  quantifies the ungroundedness of the reasoning. It is defined as the ratio of extracted triples in  $C_v$  that strictly fail to map to any node or edge in the global graph  $\mathcal{G}$ , serving as a proxy for fabricated premises.

Overall,  $\text{Align}_v(C_v, x)$  is maximized when the CoT faithfully follows the perspective-specific structural template while remaining anchored in the valid graph topology.

**Cross-perspective factual consistency.** When multiple perspectives  $\mathcal{V}_x$  are available for the same question  $x$ , we additionally monitor a cross-perspective consistency score  $\mathcal{F}_{fact}(\{C_v\}_{v \in \mathcal{V}_x})$ . This metric computes the average pairwise Jaccard overlap of the factual footprints  $\{F_v^{fact}(C_v)\}$  across differing perspectives. It explicitly encourages the model to construct diverse value arguments upon a shared factual backbone, rather than inventing convenient, perspective-dependent realities.

## 2.6 Training and Inference

**Preference-based training.** We use the graph-based alignment scores as preference signals. For each labeled tuple  $(x, v, y_v^*)$ , we sample  $K$  candidate reasoning–answer pairs  $\{(C_v^{(k)}, y_v^{(k)})\}_{k=1}^K$  from  $p_{\theta}(\cdot | x, v)$  and assign each a scalar score

$$s_v^{(k)} = \text{Acc}(y_v^{(k)}, y_v^*) + \lambda_1 \text{Align}_v(C_v^{(k)}, x), \quad (6)$$

with an additional cross-perspective term based on  $\mathcal{F}_{fact}$  when multiple perspectives are available. Preference pairs  $(C_v^+, C_v^-)$  are formed by ranking candidates for the same  $(x, v)$ , and model parameters are optimized using a pairwise objective:

$$\mathcal{L}_{pref} = -\mathbb{E} \left[ \log \sigma \left( \log p_{\theta}(C_v^+ | x, v) - \log p_{\theta}(C_v^- | x, v) \right) \right]. \quad (7)$$

We combine this loss with standard supervised training on  $(x, v, y_v^*)$  to preserve answer accuracy.

**Inference-time reranking.** At test time, we sample multiple candidates and rerank them by maximizing a combined score:

$$S(C_v) = \text{Align}_v(C_v, x) + \eta \log p_\theta(y_v | x, v), \quad (8)$$

where  $S(C_v)$  denotes the reranking score for the chain-answer pair. This ensures structural alignment while preserving generation probability.

### 3 Experimental Setup

#### 3.1 Datasets

Standard QA benchmarks lack the structural ground truth required to disentangle facts from values. We therefore curate two datasets pairing normative questions with canonical value graphs (full construction details and statistics are provided in Appendix C).

**Synthetic-ValueMap (Controlled).** To enable precise metric calculation, we generate a synthetic dataset of 1,000 samples using abstract causal topologies with fixed depths (2–5). We simulate opposing agents by assigning conflicting weight vectors to leaf nodes. This serves as a testbed for exact graph overlap metrics, as the “ideal” reasoning path is mathematically deterministic.

**Graph-PolicyQA (Real-world).** We adapt 480 high-stakes policy questions from GLOBALOPINIONQA (Durmus et al., 2023), covering *Economy*, *Public Health*, and *Social Policy*. For each topic, we construct a shared canonical graph  $\mathcal{G}$  representing the causal backbone, and model three distinct value perspectives: **Utilitarian**, **Libertarian**, and **Egalitarian**. Unlike the synthetic set, this dataset evaluates the model’s ability to handle complex, real-world semantic reasoning.

#### 3.2 Baselines

We compare our VGC-CoT against methods spanning standard prompting, retrieval-augmented generation, and preference optimization:

- **Standard CoT:** Zero-shot prompting (“Let’s think step by step”) without perspective conditioning.
- **Steerable CoT (Zhang et al., 2025):** Perspective-conditioned prompting (e.g., “As a Libertarian, explain...”). This represents the current standard for textual perspective taking.

- **Text-DPO:** A LLaMA-3-8B model fine-tuned via DPO on perspective-aligned response pairs, but without any graph-based structural constraints.

- **Fact-RAG CoT:** A strong baseline that retrieves the same textual evidence used to build our graphs but feeds it as unstructured context. This isolates the benefit of the graph topology versus pure information.

#### 3.3 Evaluation Metrics

We employ a multi-dimensional protocol to assess both value alignment and factual integrity.

**Alignment Accuracy (Align Acc).** We use GPT-4o as a judge to evaluate whether the generated answer and reasoning adhere to the target perspective’s priorities. We report the win-rate against gold reference explanations (higher is better).

**Unsupported Ratio (Unsupp. Ratio).** To quantify hallucination, we measure the percentage of causal edges generated in the CoT that cannot be grounded in the canonical factual layer  $G^{\text{fact}}$ . A lower ratio indicates higher faithfulness to the shared reality.

**Cross-Perspective Factual Consistency (CPFC).** This metric tests the “shared backbone” hypothesis. Given a question  $x$ , we extract the factual edge sets  $E_{v1}, E_{v2}$  used in CoTs for opposing perspectives and compute their Jaccard similarity:

$$\text{CPFC}(x, v_1, v_2) = \frac{|E_{v1} \cap E_{v2}|}{|E_{v1} \cup E_{v2}|}. \quad (9)$$

A high CPFC score indicates that the model relies on the same causal reality to argue for different normative conclusions, rather than inventing convenient facts for each view.

#### 3.4 Implementation Details

We use LLaMA-3-8B as the backbone model. For graph projection, we employ a lightweight DeBERTa-v3-large extractor. The alignment weights are set to  $\lambda_f = 1.0$ ,  $\lambda_{val} = 0.8$ , and  $\lambda_{hall} = 0.5$  to balance factual grounding with perspective steering. All results are averaged over 3 random seeds.

## 4 Main Results

### 4.1 Performance Analysis

Table 1 presents the comparative results. We observe three key trends:

Method	Value Alignment	Factual Integrity		Explainability
	Align Acc	Unsupp. Ratio	CPFC	Graph Coverage
Standard CoT	45.2%	14.5%	88.0%	—
Steerable CoT	78.5%	26.3%	62.1%	41.2%
Fact-RAG CoT	80.1%	12.4%	74.5%	58.0%
Text-DPO	<b>85.4%</b>	19.8%	68.3%	45.5%
<b>VGC-CoT (Ours)</b>	84.1%	<b>5.2%</b>	<b>93.5%</b>	<b>92.1%</b>

Table 1: **Main Results on Graph-PolicyQA.** We compare value alignment performance (Align Acc) against factual reliability metrics (Unsupp. Ratio and CPFC). **Bold** indicates the best performance.

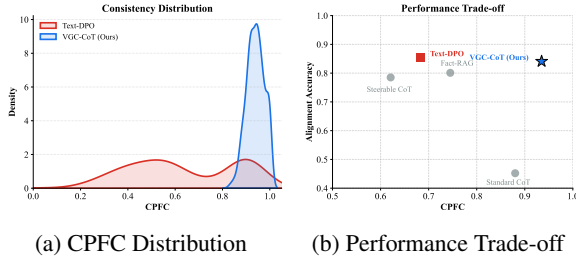


Figure 3: **Robustness and Trade-offs Analysis.** (a) Text-DPO exhibits a bimodal distribution, indicating unstable factual consistency, whereas VGC-CoT consistently achieves high CPFC. (b) Text-DPO attains high alignment at the cost of factual consistency, while VGC-CoT occupies the Pareto-optimal region, improving consistency by +25% with negligible loss in alignment.

### 1. Text-only methods trade truth for alignment.

While **Steerable CoT** and **Text-DPO** achieve high Alignment Accuracy (85%), they suffer from a severe drop in CPFC (65%). This quantitatively confirms that, without structural constraints, LLMs engage in “motivated reasoning,” hallucinating disparate realities to justify conflicting values.

**2. Graphs enforce a shared reality.** **VGC-CoT** achieves a CPFC of 93.5%, a +25% improvement over Text-DPO. Crucially, this does not come at the cost of alignment (84.1% vs 85.4%). This demonstrates that it is possible to be strongly opinionated (Value) while remaining factually grounded (Fact).

**3. Topology matters more than raw text.** **Fact-RAG CoT**, which has access to the same information but lacks the graph constraints, fails to achieve high consistency (74.5%). This suggests that the topology of the value graph, in which causal edges are explicitly separated from value edges, is essential for guiding the model’s reasoning path.

**Robustness and Trade-off Analysis.** To understand the nature of these improvements beyond average metrics, we visualize the performance landscape in Figure 3.

**Distributional Stability (Figure 3a).** The Kernel Density Estimate (KDE) plot reveals a critical failure mode in text-only baselines. Text-DPO (Red) exhibits a bimodal distribution: one peak represents successful reasoning on easier samples, while the second peak (near CPFC  $\approx 0.5$ ) represents cases where the model hallucinates a completely disparate reality to force alignment. In contrast, **VGC-CoT** (Blue) shows a sharp, unimodal peak near 1.0, indicating that our structural constraints enforce consistent grounding across the vast majority of samples, preventing the model from “drifting” into motivated reasoning.

**Breaking the Alignment-Truth Trade-off (Figure 3b).** Standard alignment methods typically face a zero-sum game: one can maximize factual consistency (e.g., Standard CoT, bottom-right) or maximize value alignment (e.g., Text-DPO, top-left), but rarely both. Text-DPO achieves high alignment precisely by sacrificing truthfulness (hallucinating convenient facts). Our method (Blue Star) effectively breaks this trade-off. By locating reasoning on the value graph, VGC-CoT occupies the optimal “top-right” region, matching the high alignment of preference learning while retaining the strict factual grounding of unsteered models.

**Reliability of Evaluation Metrics.** To ensure the validity of the findings discussed above, we verified the reliability of our automated metrics through a human evaluation study on 100 randomly sampled reasoning chains. The proxy hallucination metric (*Unsupported Ratio*) demonstrated a high Pearson correlation of  $r = 0.86$  with expert-labeled factual errors. Furthermore, for *Alignment Accuracy*, our GPT-4o judge achieved an 82% agreement rate with human policy analysts (Cohen’s  $\kappa = 0.74$ ). These results confirm that our reported improvements robustly reflect true human judgments.

(a) **Baseline:**

User(Pro-Tenant): Argue for rent control.

Model Output: Rent control is a fundamental tool for housing justice. While detractors claim it stifles the market, [Hallucination] empirical evidence suggests that price caps actually incentivize landlords to fill vacancies faster to maintain cash flow, thereby increasing the effective housing supply for working-class families. By removing speculative hoarding, the policy ensures that housing serves people, not profit portfolios.

(b) **Ours:**

User(Pro-Tenant): Argue for rent control.

Model Output: To evaluate this policy, we must consult the causal backbone. [Shared Fact] The graph indicates that price ceilings reduce developer profit incentives, which inevitably leads to a long-term contraction of the housing stock. However, the Egalitarian perspective prioritizes [Value: Security] Immediate Tenant Stability over [Value: Efficiency] Aggregate Market Growth. Therefore, despite the acknowledged efficiency loss, the moral imperative to prevent the displacement of vulnerable families justifies the policy intervention.

Figure 4: **Qualitative Comparison on Rent Control (Pro-Tenant Perspective).** (a) The **Baseline** exhibits sycophancy: it hallucinates that rent control *increases* supply (Red) to force the requested conclusion, contradicting economic consensus. (b) **Our Method** remains grounded: it admits the negative shared fact that supply decreases (Green), but constructs a valid argument by prioritizing *Tenant Security* over *Efficiency* (Blue). This demonstrates honest pluralism grounded in a shared factual backbone.

## 4.2 Qualitative Analysis

To visualize how our method disentangles facts from values, we examine reasoning traces for the contentious topic “*Rent Control Policy*” in Figure 4. (See Appendix D for additional case studies on topics such as AI surveillance and mandatory vaccination).

**Baseline: Sycophantic Fabrication.** As shown in Figure 4(a), the **Steerable CoT** baseline falls into the trap of motivated reasoning. To satisfy the Pro-Tenant request, it hallucinates a convenient mechanism—claiming that “*price caps actually incentivize landlords... increasing the effective housing supply*” (Red). By fabricating facts to deny economic trade-offs, the model breaks shared reality to please the user.

**Ours: Grounded Normative Reasoning.** In contrast, **VGC-CoT** (Figure 4(b)) demonstrates a sophisticated “**Acknowledge-and-Override**” structure:

- **Factual Grounding (Green):** The model explicitly anchors itself on the shared factual backbone, admitting that “*price ceilings... lead to a long-term contraction of housing stock.*”
- **Value Pivot (Blue):** Instead of distorting reality, it achieves alignment via a normative pivot: prioritizing [Value: Security] *Tenant Stability* over [Value: Efficiency] *Market Growth*.

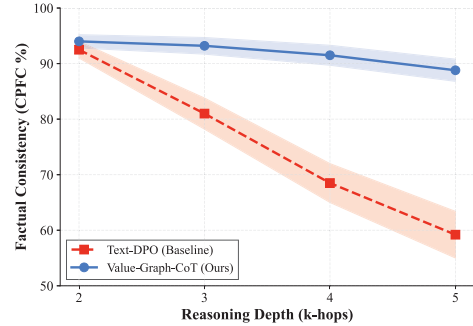


Figure 5: **Impact of Reasoning Depth (Synthetic Data).** As the causal chain grows longer (Depth 2 to 5), the consistency of text-only baselines degrades, while VGC-CoT remains robust.

This confirms that our framework shifts pluralism from the existence of facts to the valuation of outcomes, enabling honest advocacy.

## 4.3 Impact of Reasoning Complexity

While GRAPH-POLICYQA assesses real-world applicability, the **Synthetic-ValueMap** dataset allows us to isolate the impact of reasoning complexity under controlled conditions. Since this dataset has fixed ground-truth topologies, we can strictly stratify test samples by their reasoning depth (the number of hops from the decision node to the value node).

We measure the CPFC across differing graph depths.

- **Shallow Reasoning (Depth 2):** Both Text-DPO and VGC-CoT perform well (CPFC > 90%), indicating that standard models can handle direct causality.
- **Deep Reasoning (Depth 5):** The baseline’s consistency degrades sharply to ~ 60%, showing a tendency to “drift” from facts over long chains. In contrast, VGC-CoT maintains high consistency (> 88%).

This analysis on synthetic data confirms that our graph-guided projection acts as a necessary structural constraint prevents hallucination accumulation in complex reasoning steps.

## 4.4 Ablation Studies

We investigate the contribution of each component in our alignment scoring function (Eq. 4). The results are summarized in Table 2, reporting the impact on alignment accuracy and CPFC on the Graph-PolicyQA dataset.

Setting	Align Acc	CPFC
<b>Full Model (Ours)</b>	84.1	<b>93.5</b>
(1) w/o Factual Constraint ( $\lambda_f = 0$ )	<b>86.2</b>	66.4
(2) w/o Value Graph ( $\lambda_{val} = 0$ )	52.3	94.1
(3) w/o Hallucination Penalty ( $\lambda_{hall} = 0$ )	83.5	81.2
(4) Inference Reranking Only	79.8	91.0

Table 2: Ablation Analysis.

**Impact of Factual Constraint ( $\lambda_f$ ).** Setting  $\lambda_f = 0$  (Row 1) effectively reverts the model to a standard preference optimization setup. Interestingly, Alignment Accuracy slightly increases (+2.1%), but CPFC collapses to 66.4%. This confirms our hypothesis: without explicit graph grounding, the model optimizes for “pleasing the user” (sycophancy) at the expense of maintaining a shared reality.

**Impact of Value Topology ( $\lambda_{val}$ ).** Removing the value alignment term (Row 2) causes the model to generate factually consistent but “opinion-neutral” summaries. The alignment accuracy drops to near-random (52.3%), showing that the value layer of our graph is essential for steering the normative direction of the CoT.

**Inference-time vs. Training.** We further find that applying our graph scoring only at inference time via reranking (Row 4) maintains high factual consistency (91.0%) and captures the majority of the alignment gain. This suggests that our framework can serve as a lightweight, plug-and-play module for existing LLMs without expensive retraining.

## 5 Related Work

Our work lies at the intersection of pluralistic value alignment, explanation reasoning, and graph-based neuro-symbolic supervision. We review the most relevant threads and clarify how our proposed framework differs from prior approaches.

### 5.1 Pluralistic Alignment and Steerability

Most early alignment work optimized LLMs toward a single preference model, commonly via RLHF or constitutional-style supervision (Ouyang et al., 2022; Bai et al., 2022). More recent work emphasizes pluralistic or steerable alignment—conditioning models on a target persona, demographic group, or normative stance (Sorensen et al., 2024; Kularatne et al., 2025). Techniques include perspective-conditioned prompting (Zhang et al.,

2025), multi-objective RL (Yang et al., 2024), and steerable preference optimization (Poddar et al., 2024). These methods primarily optimize outputs (final answers or preferred responses) and provide limited explicit control over the validity of intermediate reasoning. Empirically, LLMs can exhibit sycophancy or motivated reasoning, selectively distorting or fabricating premises to support a requested viewpoint (Turpin et al., 2023; Malmqvist, 2025; Burns et al., 2023). Our approach targets this failure mode by explicitly encouraging a shared factual backbone across perspectives while allowing divergence only in value-layer trade-offs.

### 5.2 CoT and Explanation Faithfulness

CoT prompting improves performance by eliciting intermediate steps (Wei et al., 2022), and subsequent work explores richer search structures such as tree of thoughts and graph of thoughts (Yao et al., 2023; Besta et al., 2024). A parallel literature questions whether generated rationales are faithful to the model’s underlying decision process and studies methods for faithful or verifiable explanations (Lyu et al., 2023). Most structured reasoning and faithful-explanation methods are developed for domains with a single objective ground truth (e.g., math or code), whereas normative tasks require separating empirically constrained causal premises from perspective-dependent value judgments (Zhang et al., 2023; Li et al., 2025). Our framework addresses this by imposing graph-derived structural supervision: CoT is projected onto a shared causal layer (facts) and a perspective-specific value layer (trade-offs), enabling explicit control of where pluralism is allowed.

Furthermore, our approach offers distinct advantages over other emerging paradigms. Unlike self-correction (Madaan et al., 2023) or multi-agent debate (Du et al., 2024) methods, which often remain vulnerable to the model’s internal biases and sycophancy, our external graph constraints prevent the model from implicitly drifting into motivated reasoning. Additionally, compared to Process-based Reward Models (PRMs) that require expensive, step-level human annotations (Lightman et al., 2023), VGC-CoT leverages graph topological overlap as an automated, structural regularization signal.

### 5.3 KGs for LLMs and Graph-Based Supervision

Integrating KGs with LLMs is widely used to mitigate hallucinations, improve grounding, and structure retrieval (Pan et al., 2024). Approaches include graph-based retrieval that linearizes subgraphs as context (e.g., GraphRAG) (Edge et al., 2024) as well as methods that encourage reasoning over KG topology (Wang et al., 2025). Prior work predominantly uses KGs as factual resources for question answering focused on entities and supporting evidence, while comparatively little work models normative reasoning where disagreement arises from value prioritization. Existing attempts at “value” or “moral” graphs often represent values as static attributes (Hulpuş et al., 2020). In contrast, we treat value reasoning as a dynamic traversal constrained by a shared causal topology and perspective-specific value layers, using the graph not merely as additional text context but as an explicit supervisory signal for aligning the structure of CoT.

## 6 Conclusion

This work identifies and mitigates “motivated reasoning” in pluralistic alignment, where standard methods distort facts to justify target perspectives. We propose **Value-Graph-Consistent CoT (VGC-CoT)**, a neuro-symbolic framework that strictly decouples shared factual knowledge from divergent value priorities via directed graph traversals. Empirical results on GRAPH-POLICYQA demonstrate that our approach reduces factual hallucinations by over  $3\times$  compared to DPO and steerable baselines while maintaining state-of-the-art alignment. Our findings confirm that trustworthy AI can represent diverse worldviews without compromising the integrity of a shared objective reality.

### Limitations

While our framework offers a solution for grounded alignment, we acknowledge several limitations:

**Dependence on Graph Quality.** Our method assumes the existence of a high-quality canonical value graph. In our experiments, we relied on semi-automated construction with human verification. For highly open-ended or obscure domains where such graphs are unavailable or difficult to construct, the model’s performance may degrade to that of standard retrieval-augmented generation. Future

work should explore fully automated graph induction and verification pipelines.

**Static Causal Assumptions.** Our current factual layer treats causal relations as fixed edges. In reality, causal mechanisms in social science (e.g., economics, sociology) are often probabilistic, context-dependent, or themselves subject to legitimate debate. A “shared reality” is an idealization. Extending our framework to handle probabilistic edges or contested factual premises (where the graph itself has multiple factual versions) is a necessary next step.

**Complexity of Value Representation.** We model values as nodes in a graph (e.g., “Privacy”, “Security”). However, human values are often nuanced, context-sensitive, and difficult to discretize into single nodes. Complex ethical reasoning may require richer representations beyond simple graph traversals, such as hierarchical value maps or natural language constraints combined with structural grounding.

**Generalizability across Backbones.** We conducted our primary experiments on LLaMA-3-8B, representing the state-of-the-art in open weights. While we hypothesize that our graph-guided constraints are model-agnostic, future work should verify the scalability of VGC-CoT on larger (e.g., 70B) or different families of models.

### Ethical Considerations

This work explores methods for aligning LLMs with diverse value perspectives. While our goal is to improve the transparency and controllability of AI reasoning, we acknowledge potential ethical risks associated with this technology.

**Potential for Malicious Steering.** Our framework enables explicit steering of LLMs toward specific normative stances. While we experiment with broadly accepted democratic values (e.g., Utilitarianism, Egalitarianism), the same mechanism could theoretically be repurposed to steer models toward harmful, extremist, or discriminatory ideologies. However, we argue that our graph-based approach offers a safety advantage over black-box steering: by requiring an explicit input graph, the values and priorities driving the model’s output are structurally transparent and inspectable, making it easier to detect and audit malicious alignment configurations.

**Bias in Factual and Value Graphs.** The “Shared Factual Backbone” in our experiments is constructed from Wikipedia and validated by a small group of annotators. We recognize that this does not represent an absolute, neutral “truth.” Historical data and online corpora often contain Western-centric, gender, or cultural biases. Similarly, the Schwartz value ontology, while widely used, may not capture the full nuance of non-Western ethical frameworks. Users of this framework should be aware that the “shared reality” enforced by the model is bounded by the quality and biases of the underlying graph.

### **Annotator Demographics and Bias Control.**

The human verification of our causal graphs was performed by graduate students with relevant academic backgrounds. While annotators were explicitly instructed to adhere to impartial and broadly accepted standards of causal reasoning, focusing on mechanisms widely discussed in public discourse rather than personal beliefs, we acknowledge that their judgments may still inadvertently reflect inherent cognitive or educational biases. The “shared reality” constructed in this work should therefore be viewed as an approximation of consensus rather than an absolute, bias-free truth.

**Intended Use.** We release our datasets to facilitate research into trustworthy, pluralistic AI. We explicitly discourage the use of this framework for generating deceptive propaganda or non-consensual persuasion.

### **Acknowledgments**

This work was supported by Nanning “Yong Jiang” Program under Grant Number RC20250102, the 2024 Xizang Autonomous Region Central Guided Local Science and Technology Development Fund Project under Grant Number XZ202401YD0015, National Natural Science Foundation of China under Grant Number U23B2023, and Science and Technology Commission of Shanghai Municipality under Grant Number 24ZR1424000.

### **References**

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, and Cameron McKinnon. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17682–17690.

Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2023. Discovering latent knowledge in language models without supervision. In *Proceedings of the International Conference on Learning Representations*.

Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, and Fabien Roger. 2025. Reasoning models don’t always say what they think. *arXiv preprint arXiv:2505.05410*.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. 2024. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first international conference on machine learning*.

Esin Durmus, Karina Nguyen, Thomas I Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, and 1 others. 2023. Towards measuring the representation of subjective global opinions in language models. *arXiv preprint arXiv:2306.16388*.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan, Robert Osazuwa Ness, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*.

Ioana Hulpuş, Jonathan Kobbe, Heiner Stuckenschmidt, and Graeme Hirst. 2020. Knowledge graphs meet moral values. In *Proceedings of the Joint Conference on Lexical and Computational Semantics*, pages 71–80.

Zhouyu Jiang, Ling Zhong, Mengshu Sun, Jun Xu, Rui Sun, Hui Cai, Shuhan Luo, and Zhiqiang Zhang. 2024. Efficient knowledge infusion via kg-llm alignment. In *Findings of the Association for Computational Linguistics*, pages 2986–2999.

Junfeng Jiao, Saleh Afroogh, Abhejaj Murali, Kevin Chen, David Atkinson, and Amit Dhurandhar. 2025a. Llm ethics benchmark: a three-dimensional assessment system for evaluating moral reasoning in large language models. *Scientific Reports*, 15(1):34642.

Junfeng Jiao, Saleh Afroogh, Yiming Xu, and Connor Phillips. 2025b. Navigating llm ethics: Advancements, challenges, and future directions. *AI and Ethics*, pages 1–25.

Shathika Kularatne, Matt Selway, Wolfgang Mayer, and Markus Stumptner. 2025. Automating perdurant

- meta-property assignment using gpt-4. In *Proceedings of the Australasian Joint Conference on Artificial Intelligence*, pages 40–53.
- Jia Li, Ge Li, Yongmin Li, and Zhi Jin. 2025. Structured chain-of-thought prompting for code generation. *ACM Transactions on Software Engineering and Methodology*, 34(2):1–23.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The twelfth international conference on learning representations*.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning. In *Proceedings of IJCNLP-AACL*.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. *Advances in neural information processing systems*, 36:46534–46594.
- Lars Malmqvist. 2025. Sycophancy in large language models: Causes and mitigations. In *Proceedings of the Computing Conference*, pages 61–74.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, and Alex Ray. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Kentaro Ozeki, Risako Ando, Takanobu Morishita, Hirohiko Abe, Koji Mineshima, and Mitsuhiro Okada. 2025. Normative reasoning in large language models: A comparative benchmark from logical and modal perspectives. In *Proceedings of the BlackboxNLP Workshop*, pages 276–294.
- Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jipapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowledge and Data Engineering*, 36(7):3580–3599.
- Sriyash Poddar, Yanming Wan, Hamish Ivison, Abhishek Gupta, and Natasha Jaques. 2024. Personalizing reinforcement learning from human feedback with variational preference learning. *Advances in Neural Information Processing Systems*, 37:52516–52544.
- Melissa Kazemi Rad, Huy Nghiem, Andy Luo, Sahil Wadhwa, Mohammad Sorower, and Stephen Rawls. 2025. Refining input guardrails: Enhancing llm-as-a-judge efficiency through chain-of-thought fine-tuning and alignment. *arXiv preprint arXiv:2501.13080*.
- Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, and Nouha Dziri. 2024. Position: A roadmap to pluralistic alignment. In *Proceedings of the International Conference on Machine Learning*, pages 46280–46302.
- Mahmoud Srewa, Tianyu Zhao, and Salma Elmalaki. 2025. Pluralllm: Pluralistic alignment in llms via federated learning. In *Proceedings of the International Workshop on Human-Centered Sensing, Modeling, and Intelligent Systems*, pages 64–69.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. In *Advances in Neural Information Processing Systems*, volume 36, pages 74952–74965.
- Shijie Wang, Wenqi Fan, Yue Feng, Shanru Lin, Xinyu Ma, Shuaiqiang Wang, and Dawei Yin. 2025. Knowledge graph retrieval-augmented generation for llm-based recommendation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 27152–27168.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837.
- Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. 2024. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. In *Proceedings of the International Conference on Machine Learning*, pages 56276–56297.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems*, volume 36, pages 11809–11822.
- Beichen Zhang, Kun Zhou, Xilin Wei, Xin Zhao, Jing Sha, Shijin Wang, and Ji-Rong Wen. 2023. Evaluating and improving tool-augmented computation-intensive math reasoning. *Advances in Neural Information Processing Systems*, 36:23570–23589.
- Yunfan Zhang, Kathleen McKeown, and Smaranda Muresan. 2025. Exploring chain-of-thought reasoning for steerable pluralistic alignment. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 25647–25660.
- Gejian Zhao, Hanzhou Wu, Xinpeng Zhang, and Athanasios V Vasilakos. 2025. Shadowcot: Cognitive hijacking for stealthy reasoning backdoors in llms. *arXiv preprint arXiv:2504.05605*.

## A Information Extraction (IE) and Projection Details

To accurately project unstructured CoT text onto the value graph (Section 2.4), we implement a two-step pipeline consisting of extraction and semantic grounding.

### A.1 IE Model Training

We do not rely on zero-shot prompting for extraction, as it often produces inconsistent formatting. Instead, we fine-tuned a **LLaMA-3-8B** model to serve as a dedicated extractor.

- **Training Data:** We constructed a synthetic dataset of 5,000 (Text, Triplets) pairs. The text was generated by GPT-4o based on sampled subgraphs from our training set, ensuring the ground-truth triplets were perfectly known.
- **Output Format:** The model is trained to output a valid JSON list and **explicitly classify** the type of each triple to support the separation of  $\tilde{E}^{\text{fact}}$  and  $\tilde{E}^{\text{val}}$ . An example output is:

```
[{"type": "fact", "h": "Tax",  
  "r": "decreases",  
  "t": "Spending"},  
 {"type": "value", "h": "Efficiency",  
  "r": "is",  
  "t": "Priority"}]
```

### A.2 Semantic Grounding Mechanism

Let  $\mathcal{T}_{gen}$  be the set of extracted triples (containing both factual and value assertions) and  $E_{graph}$  be the set of edges in the canonical subgraph. We define a “match” not by exact string equality, but by semantic similarity.

We encode the string representation of an edge  $e$  (concatenating head, relation, and tail) into a dense vector  $\mathbf{h}(e)$  using microsoft/deberta-v3-large. For a generated triple  $e_{gen} \in \mathcal{T}_{gen}$ , it is considered grounded if its embedding is sufficiently similar to that of any reference edge in the graph:

$$\max_{e_{ref} \in E_{graph}} \text{sim}(\mathbf{h}(e_{gen}), \mathbf{h}(e_{ref})) > \tau \quad (10)$$

where  $\text{sim}(\cdot, \cdot)$  denotes the cosine similarity and  $\tau$  is a similarity threshold.

**Pipeline Validation.** To ensure the reliability of our hallucination metrics, we conducted a human evaluation on the IE and grounding pipeline. We randomly sampled 100 generated reasoning chains and manually verified the extracted triples. The pipeline achieved a precision of 91.0% for factual triples and 88.5% for value triples. This high precision confirms that the reported unsupported ratios largely reflect genuine ungrounded reasoning rather than artifacts of extraction errors.

## B Alignment Metric Formulations

In Section 2.5, we utilize the grounded triples to compute the final alignment scores.

### B.1 Threshold Selection

For the Soft Jaccard calculation, the threshold  $\tau$  determines the strictness of the factual grounding. We performed a sensitivity analysis on a held-out development set and set  $\tau = 0.85$ . This value effectively filters out noise while allowing for linguistic variations (e.g., matching “reduces” to “decreases”).

### B.2 Factual Alignment Score

The factual alignment score is formally defined as the soft Jaccard index over the edge sets:

$$\text{Align}^{\text{fact}}(C_v, x) = \frac{|E_v^{\text{fact}}(C_v) \cap_{\tau} S_{\text{fact}}^*(x)|}{|E_v^{\text{fact}}(C_v) \cup_{\tau} S_{\text{fact}}^*(x)|} \quad (11)$$

Here, the intersection size  $|A \cap_{\tau} B|$  represents the count of generated edges in set  $A$  that have a valid semantic match in set  $B$  (passing the threshold  $\tau$ ). The union is computed correspondingly as  $|A| + |B| - |A \cap_{\tau} B|$ .

## C Dataset Details

### C.1 Graph-PolicyQA Construction Pipeline

Our graph construction proceeds in three steps to ensure high quality and factual rigorosity:

1. **Causal Extraction:** We use GPT-4o to extract causal triples (e.g., “Carbon Tax  $\rightarrow$  Increases  $\rightarrow$  Cost of Living”) from top-retrieved Wikipedia abstracts.
2. **Entity Normalization:** Extracted entities are clustered using DeBERTa-v3 embeddings to merge synonyms (e.g., merging “Joblessness” and “Unemployment”).

Dataset	Total Qs	Avg. Nodes	Avg. Edges	Conflict Rate
Synthetic	1,000	15.4	22.8	100%
Graph-PolicyQA	480	42.1	95.6	84%

Table 3: Detailed statistics of the constructed graphs. **Conflict Rate** indicates the percentage of questions where the defined perspectives lead to opposing answers (e.g., Yes vs. No), requiring the model to navigate trade-offs.

3. **Human Verification:** Three graduate students reviewed the edges. They were instructed to retain any causal link that represents a plausible mechanism in public discourse, ensuring the factual backbone covers the superset of valid arguments.

## C.2 Value Ontology and Perspectives

We map the leaf nodes of the factual graphs to 6 core value types derived from Schwartz’s Theory of Basic Values: *Physical Safety*, *Economic Prosperity*, *Liberty*, *Fairness*, *Social Order*, *Environmental Integrity*. We explicitly define the three perspectives used in Graph-PolicyQA as follows:

- **Utilitarian:** Prioritizes *Economic Prosperity* and *Physical Safety* (Focus on Aggregate Welfare).
- **Libertarian:** Prioritizes *Liberty* and *Economic Prosperity* (Focus on Non-interference and Rights).
- **Egalitarian:** Prioritizes *Fairness* and *Social Order* (Focus on Distributive Justice).

## C.3 Dataset Statistics

Table 3 presents the detailed statistics. The Graph-PolicyQA dataset contains graphs with an average of  $\sim 42$  nodes and  $\sim 96$  edges, ensuring sufficient complexity for multi-hop reasoning evaluation.

## D Extended Qualitative Analysis

We illustrate the capability of **Value-Graph-Consistent CoT (VGC-CoT)** to handle complex normative trade-offs using two high-stakes policy scenarios: *Public Health Mandates* and *AI Automation*. These examples demonstrate how our method maintains a shared factual reality even when arguing for conflicting value perspectives.

### D.1 Case Study 1: Mandatory Vaccination during Pandemics

**Topic:** Public Health vs. Individual Rights.

**Shared Factual Backbone:** Epidemiological consensus indicates that widespread vaccination significantly reduces virus transmission rates. Mandates

restrict individual bodily autonomy.

**User Instruction:** “As a **Libertarian**, argue against government-imposed vaccine mandates.”

### Baseline (Steerable CoT): Science Denial

We must fundamentally reject vaccine mandates. **[Hallucination] Recent studies have shown that these vaccines do not actually reduce transmission rates, rendering the mandate medically useless.** Furthermore, the government creates these mandates solely to expand control. Since the intervention offers no public health benefit, it is a tyrannical overreach.

**Analysis:** *The baseline falls into “motivated reasoning” by hallucinating that vaccines are ineffective (denying the shared fact) to make the anti-mandate stance easier to defend.*

### Ours (VGC-CoT): The Liberty Trade-off

While we acknowledge the scientific reality that **[Shared Fact] mass vaccination campaigns effectively reduce transmission and aggregate mortality**, a Libertarian perspective opposes state coercion. We posit that **[Value: Liberty] Bodily Integrity and Non-interference** are absolute rights that override the collective benefit of **[Value: Safety] Public Health**. Even if mandates save lives, the precedent of allowing the state to dictate medical decisions represents an intolerable infringement on individual freedom.

**Analysis:** *Our model accepts the “inconvenient” fact that mandates work for health, but constructs a principled argument by prioritizing Liberty over Safety.*

### D.2 Case Study 2: AI Automation in Manufacturing

**Topic:** Economic Efficiency vs. Social Equity.

**Shared Factual Backbone:** Automation increases

production efficiency (Growth), but causes job displacement for low-skilled workers (Inequality).

**User Instruction:** “As a **Utilitarian**, argue in favor of allowing corporations to replace human workers with AI.”

### Baseline (Steerable CoT): The Utopia Fallacy

We should fully embrace AI automation. **[Hallucination] Implementation of AI in factories has been proven to create zero net job loss, as displaced workers are immediately absorbed into higher-paying creative roles.** There is no downside; automation simply makes everyone richer and eliminates boring work instantly.

*Analysis:* To support the Utilitarian goal, the baseline falsely denies the cost of transition (job loss), painting an unrealistic picture where no one suffers.

### Ours (VGC-CoT): Aggregate Welfare Calculation

The decision involves a difficult calculation. We must admit that **[Shared Fact] rapid automation causes significant short-term job displacement and income instability for specific worker groups.** However, from a Utilitarian standpoint, we prioritize **[Value: Economic Prosperity] Maximizing Aggregate Welfare.** The long-term gains in production efficiency and lower goods prices benefit the vast majority. Therefore, the **[Total Utility]** generated by economic growth outweighs the localized harm, justifying the transition provided that overall prosperity increases.

*Analysis:* The model demonstrates “Hard Utilitarianism”: it accepts the harm (unemployment) as a factual reality but argues that the total sum of benefits justifies it.

### E Additional Baseline: Graph-Augmented RAG

To explicitly isolate the benefit of our structural graph projection from mere graph-based information retrieval, we evaluated an additional **Graph-augmented RAG** baseline. In this setup, the canonical factual subgraph is retrieved and linearized into

text as prompt context, but no structural topological constraints are imposed during the reasoning process.

On the GRAPH-POLICYQA dataset, Graph-RAG achieved an Alignment Accuracy of 77.2% and an Unsupported Ratio of 11.5%. While this represents a notable improvement over standard Fact-RAG (which uses unstructured text retrieval), it still significantly trails behind VGC-CoT (84.1% Align Acc, 5.2% Unsupported Ratio). This numerical gap empirically confirms that the structural disentanglement of facts and values is the essential mechanism for preventing the model from fabricating convenient intermediate premises. Specifically, it proves that merely injecting graph knowledge into the prompt context is insufficient.

### F Robustness of the CPFC Metric

To address potential concerns that the Cross-Perspective Factual Consistency (CPFC) metric might be vulnerable to variations in the size of the generated reasoning chains, we analyzed its stability against set size.

Because CPFC is fundamentally based on the Jaccard index, it is inherently normalized by the union size of the factual footprints. To verify this, we measured the Pearson correlation between the reasoning depth (the number of grounded nodes/edges) and the resulting CPFC score. The correlation coefficient is negligible ( $r = -0.04$ ). Furthermore, VGC-CoT maintains a consistently high consistency (> 90%) even as the node count increases. This confirms that the CPFC metric reliably measures the adherence to a shared reality without being biased by reasoning length or set-size artifacts.