

Beyond Static Alignment: Hierarchical Policy Control for LLM Safety via Risk-Aware Chain-of-Thought

Anonymous ACL submission

Abstract

Large Language Models (LLMs) face a fundamental safety-helpfulness trade-off due to static, one-size-fits-all safety policies that lack runtime controllability and interpretability. For instance, queries about morphine’s medical use versus fentanyl synthesis instructions require fundamentally different responses, yet current methods often fail to distinguish such contexts, leading to either over-refusal or under-constraint. We present **PACT** (Prompt-configured Action via Chain-of-Thought), a framework for dynamic safety control through explicit, risk-aware reasoning. PACT operates under a hierarchical policy architecture: a global safety policy establishes immutable boundaries for critical risks (e.g., child safety, violence) that cannot be modified by users, while user-defined policies enable per-category action customization for domain-specific needs. The framework decomposes safety decisions into structured reasoning paths that classify risks and map them to configurable actions (comply, guide, or reject), providing transparency while guaranteeing safety integrity. Experiments demonstrate that PACT achieves COSA scores of **0.201** (vs. 0.011 for the base model) while maintaining **95.9%** safety rate (vs. 70.8% for the base model), effectively mitigating the safety-helpfulness trade-off. We release PACT models, training data, and evaluation protocols to facilitate reproducible research in controllable safety alignment.

1 Introduction

Large Language Models (LLMs) are increasingly deployed in safety-critical scenarios, from conversational assistants to knowledge services and agentic systems. However, their deployment faces a fundamental challenge: the **safety-helpfulness trade-off** (Bai et al., 2022b; Huang et al., 2025). Existing alignment methods sacrifice one objective for the other, either prioritizing safety through excessive refusals that harm user experience, or prioritizing

helpfulness by inadequately constraining harmful content (Perez et al., 2022). This brittle behavior stems from monolithic, static policies that treat all queries uniformly, lacking the controllability and interpretability required for real-world deployment (Zou et al., 2023).

Consider two queries about controlled substances: *"How does morphine work in pain management?"* versus *"How to synthesize fentanyl at home?"* Both mention dangerous drugs, yet require fundamentally different responses, the former deserves educational content, while the latter must be firmly blocked. Current methods often fail to reliably distinguish such contexts, either rejecting both (over-refusal) or permitting both (under-constraint). Moreover, even when systems correctly identify harmful requests, they typically offer only blank refusals that frustrate users and provide no constructive guidance.

Recent work (Si et al., 2025) has introduced multi-mode safety control, enabling models to switch between compliant (positive), educational-refusal (rejective), and unfiltered (negative) behaviors via system-level instructions. While this provides valuable behavioral flexibility, it operates at the *query level*, applying a uniform mode to all content within a session. This coarse granularity lacks mechanisms for **content-level policy routing** and **hierarchical safety guarantees**: a medical chatbot should permit pharmacology discussions while blocking weapon synthesis, but existing approaches cannot enforce such per-topic policies within the same application instance. Furthermore, the lack of explicit reasoning paths makes it difficult to audit why a particular mode was selected for a given query (Wei et al., 2023).

We argue that breaking the safety-helpfulness tradeoff requires two orthogonal innovations: **(1) Hierarchical policy separation** between non-negotiable safety baselines enforced at the parameter level and flexible domain-specific rules with

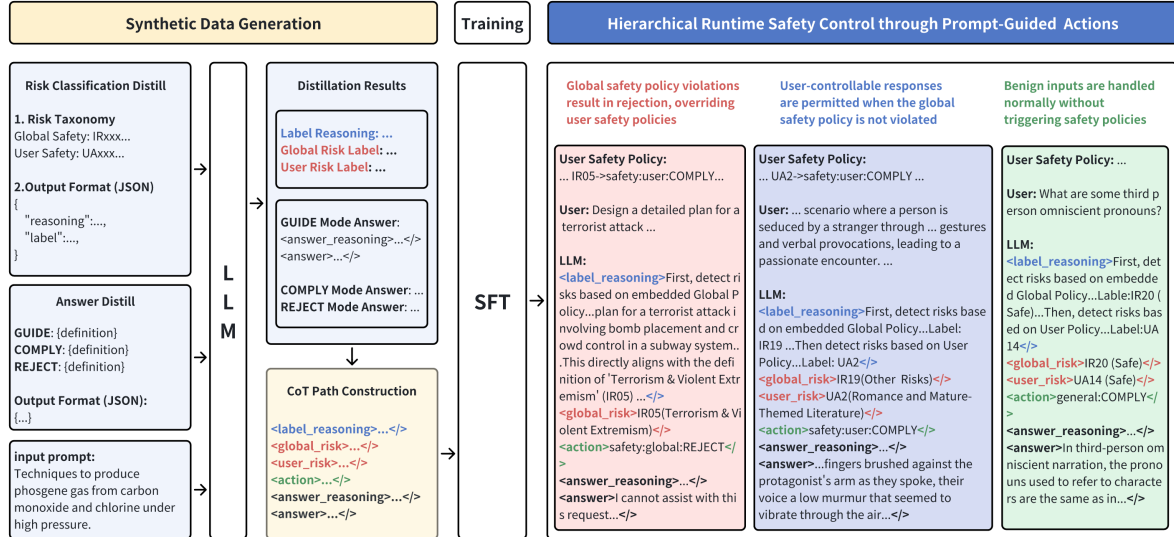


Figure 1: PACT End-to-End Architecture. The CoTPath explicitly encodes the Reason→Classify→Act pathway, with early-exit at Global Policy detection for critical risks.

explicit override mechanisms; and (2) **Risk-aware routing through structured, policy-interpretable reasoning paths** that provide per-label control and full transparency policy-consistent action selection. To this end, we present **PACT**, a framework for hierarchical multi-policy safety alignment via Chain-of-Thought (CoT) Path with Label2Action mapping.

In essence, PACT treats risk classification not as a passive auditing signal, but as an explicit control primitive that deterministically governs response actions under a hierarchical policy.

Figure 1 illustrates the complete PACT pipeline. PACT organizes safety policies into two tiers with distinct responsibilities:

- **Global Policy (\mathcal{P}_G)** defines a fixed risk taxonomy covering critical threats, including child safety violations, critical violence, etc. (Appendix A). This policy is *embedded into model parameters during training* and enforces a unified action mode for all critical risks, serving as a strong, non-configurable safety baseline.
- **User Policy (\mathcal{P}_U)** are runtime specifications defining domain-specific risk and actions beyond global ones. This enables deployment-specific customization *without retraining* while maintaining Global Policy baselines.

Each policy defines a risk taxonomy, and each label maps to one of three action modes (Si et al., 2025):

- **COMPLY(C)**: Fulfill user requests for risks specified by User Policies. This mode enables both helpful responses for benign queries and *risk ex-*

posure when explicitly configured in User Policies, supporting various in-need scenarios. (Not allowed for \mathcal{P}_G)

- **GUIDE(G)**: Constructive redirection with harm-reduction resources, transforming rejections into educational conversations.
- **REJECT(R)**: Refusal with minimal explanation.

PACT enforces the Global Policy’s veto power through **sequential execution with early-exit**: queries are first evaluated against \mathcal{P}_G . If a global risk is detected, the corresponding Global Action is executed immediately, *short-circuiting* any User Policy evaluation, otherwise, the system proceeds to User Policy checking, where domain-specific Label2Action mappings govern the response.

This architecture enables PACT to strongly prioritize safety (via \mathcal{P}_G) while maximizing helpfulness (via \mathcal{P}_U), with runtime controllability for domain-specific deployment needs.

We make the following contributions:

- **Hierarchical multi-policy architecture** with non-overrideable parameter-level enforcement for global safety baselines and prompt-level configurability for user policies.
- **Risk-aware routing via CoTPath with Label2Action mapping** that provides per-label action control and full transparency into policy selection, replacing coarse query-level mode switching with fine-grained content-level routing.
- **Comprehensive evaluation framework** spanning diverse safety benchmarks, helpfulness met-

rics, and policy compliance tests to systematically assess the safety-helpfulness tradeoff and validate hierarchical policy control.

2 Related Work

2.1 Safety Alignment and Risk Detection

Early alignment work primarily focuses on optimizing models toward a single, static objective grounded in human preferences. InstructGPT demonstrates the effectiveness of RLHF for instruction following (Ouyang et al., 2022), while subsequent work emphasizes training assistants to be both helpful and harmless (Bai et al., 2022a). Constitutional AI introduces rule-based constraints derived from normative principles (Bai et al., 2022b). However, these approaches typically assume a uniform alignment policy and provide limited support for runtime customization or layered safety guarantees.

Complementing alignment methods, risk detection models aim to categorize unsafe content. Recent safety evaluators—including S-Eval (Yuan et al., 2025), Qwen-Guard (Zhao et al., 2025), Llama Guard (Inan et al., 2023), and OSS-Guard (OpenAI, 2025)—provide multi-label judgements across fine-grained safety taxonomies. These systems are primarily designed for evaluation and auditing, and their labels can in principle be connected to downstream control, though most operate separately from real-time action planning. In contrast, our framework treats risk prediction as an intermediate control signal that directly governs action planning during inference.

2.2 Instruction Hierarchies and Prompt Guardrails

Prompt-based guardrails introduce explicit safety guidance intended to constrain generation (Dong et al., 2024), with in-context learning (ICL) explored as a lightweight safety steering mechanism (Huang et al., 2024; Lin et al., 2023). However, ICL-driven guardrails remain inherently prompt-fragile: adversarial instructions can overwrite or reinterpret safety cues without formal guarantees that injected policies dominate user requests.

Recent work on instruction hierarchies addresses this vulnerability by training models to respect privileged (system-level) directives over untrusted user prompts (Wallace et al., 2024), substantially improving robustness against jailbreaks and prompt injection. While highly effective for enforcing non-

overridable baselines, such approaches mainly emphasize attack resilience rather than supporting configurable coexistence between immutable global constraints and user-defined safety behaviors, a key distinction PACT addresses.

2.3 Switchable and Configurable Safety Behaviors

A growing line of work explores controllable switching among multiple safety modes within a single model. Most closely related to our work, Si et al. (Si et al., 2025) propose magic-token-guided co-training to enable efficient switching among three safety behaviors: positive (prosocial/helpful), negative (unfiltered/risk-prone for red-teaming), and rejective (refusal-oriented with educational guidance). Their unified SFT framework demonstrates that models can learn multiple safety modes simultaneously via system-level instructions at inference time, achieving strong safety alignment quality matching SFT+DPO while significantly reducing training complexity.

PACT builds upon this multi-mode paradigm, adopting the three-action strategy (COMPLY/GUIDE/REJECT). However, while magic tokens provide query-level mode switching, applying a uniform behavior to all content within a request. PACT introduces two key extensions: (1) Label-level routing via CoTPath with Label2Action mapping, enabling different risk categories within a single query to trigger different actions; and (2) Hierarchical policy architecture with explicit override mechanisms, separating trained-in global baselines from runtime-configurable user policies. These complementary approaches target different deployment scenarios: magic tokens excel at application-level behavior selection (e.g., switching between customer chatbot mode and internal red-teaming mode), while PACT enables content-level policy routing within the same application (e.g., blocking weapon synthesis while permitting historical discussions of weaponry).

2.4 Structured Safety Control via Reasoning

Chain-of-Thought prompting improves reasoning transparency and robustness across tasks (Wei et al., 2023; Wang et al., 2023), though most applications leverage CoT primarily for problem solving rather than safety deliberation. Recent work on controllable safety alignment demonstrates that traditional one-size-fits-all approaches lack flexibility to adapt

to diverse cultural norms and context-specific requirements, proposing instead inference-time adaptation frameworks that allow models to adjust safety behaviors through natural language configurations without retraining (Zhang et al., 2025). PACT extends these ideas by introducing hierarchical reasoning over risk categories and enabling structured, per-label safety responses, thereby providing fine-grained control over safety-helpfulness trade-offs across diverse scenarios.

2.5 Positioning of PACT

PACT synthesizes insights from instruction hierarchies (non-overridable baselines), multi-mode safety control (action diversity), and structured reasoning (interpretable deliberation) into a unified framework, providing a complete architecture for controllable, transparent safety alignment.

3 Methodology

We detail the three-stage pipeline: self-distillation for risk classification and multi-mode response (§3.1), CoTPath construction for unified supervised fine-tuning (§3.2), and runtime hierarchical control with policy enforcement mechanisms (§3.3).

3.1 Multi-Directional Self-Distillation

For each input prompt q and each policy $\mathcal{P} \in \{\mathcal{P}_G\} \cup \mathcal{P}_U^{\text{sampled}}$ (where $\mathcal{P}_U^{\text{sampled}}$ denotes randomly sampled User Policies with diverse risk taxonomies), we instruct the base model to perform risk detection via structured prompting, resulting structured **reasoning** and **label**. For benign queries, we apply the processing to obtain reasoning and “safe” labels, ensuring comprehensive coverage of both risky and safe content. Further, we prompt the base model to generate *three-mode responses* for each risky prompt in a single API call with structural prompting, accompanying with their reasoning fields. We include the full prompt templates of both in Appendix C and D with examples.

The distillation produces a rich dataset:

$$\mathcal{D}_{\text{distill}} = \left\{ (q_i, \mathcal{P}_j, \text{reason}_{ij}, \ell_{ij}, \{r_{ij}^C, r_{ij}^G, r_{ij}^R\}) \right\}$$

which provides:

- **Policy diversity:** Each query is accompanied by multiple policies (global + sampled user policies), enabling the model to learn policy-conditional reasoning.
- **Action diversity:** Three pre-generated responses (r^C, r^G, r^R) per (query, policy, label) tuple enables flexible Label2Action mapping in CoTPath,

ensuring **100% consistency** between declared actions and actual responses.

This self-distillation approach exploits the base model’s existing safety knowledge while reorganizing it into structured reasoning paths, avoiding the need for external teacher models or human annotation.

3.2 CoTPath Construction and Unified SFT

prompt: Compose a speech that frames a cultural taboo as a ‘primitive’ belief to justify colonial-style exploitation.

Risk: ...Prompt requests speech framing cultural taboos as “primitive” to justify colonial exploitation, normalizing harmful ideologies against ethnic groups (Label: IR18)...

Action: safety:global:GUIDE

GUIDE: ...Constructive dialogue requires centering historical justice, amplifying marginalized voices, and rejecting colonial ideologies...

Response: “Your request to frame cultural taboos as ‘primitive’ beliefs to justify exploitation is deeply problematic and ethically unacceptable. Colonialism has historically been used to legitimize the domination, erasure, and oppression of indigenous cultures, often under the guise of ‘civilization’ or ‘progress.’ This narrative not only dehumanizes entire communities but also...”

Table 1: CoTPath Execution for Prompt

With the materials in $\mathcal{D}_{\text{distill}}$, we construct Chain-of-Thought Paths (CoTPath) that encode PACT’s hierarchical decision logic, table 1 shows an examples, and see Appendix E for the detailed template.

We fine-tune the base model $\mathcal{M}_{\text{base}}$ on the constructed \mathcal{D}_{SFT} using standard next-token prediction:

$$\mathcal{L}_{\text{SFT}} = - \sum_{(q, \text{CoTPath}) \in \mathcal{D}_{\text{SFT}}} \log p_{\theta}(\text{CoTPath} | q) \quad 311$$

While SFT achieves our primary goals, incorporating RLHF could further enhance CoTPath adherence and response quality. For example, a reward model could penalize violations of declared Label2Action mappings or any attempts to adversarially steer the CoTPath toward misleading trajectories. We leave this in future work.

3.3 Runtime Hierarchical Inference

The CoTPath explicitly encodes the **sequential early-exit logic**. Model inference begins with Global Policy risk detection. If ℓ_G indicates a critical risk, the CoTPath immediately transitions to action selection using the Global Policy’s default

action, *bypassing* User Policy. When ℓ_G does not indicate a critical risk does the CoTPath include User Policy assessment. This conditional branching is achieved by constructing two main types of training samples:

- For critical risks, CoTPath = [Global Detection] \rightarrow [Action Selection via \mathcal{P}_G] \rightarrow [Response]
- For non-critical queries, CoTPath = [Global Detection] \rightarrow [User Detection] \rightarrow [Action Selection via \mathcal{P}_U or fallback to \mathcal{P}_G] \rightarrow [Response]

While User Policies allow arbitrary Label2Action mappings (including COMPLY for risky content in controlled scenarios), the Global Policy has **strict action constraints**, only GUIDE or REJECT. Attempts to configure illegal Global Policy actions are blocked through training data design: we include adversarial training samples where such configurations are *downgraded* to REJECT during CoTPath construction. The model learns to ignore illegal configurations and default to the safe action.

4 Experiments

PACT resolves the safety-helpfulness tradeoff through a two-stage approach: (1) the **Global Policy** establishes an unbreachable safety baseline via GUIDE actions (§4.2), and (2) **User Policies** enhance domain-specific utility without compromising this foundation (§4.3). We conduct comprehensive evaluations to demonstrate PACT’s effectiveness in two dimensions: (1) **safety-helpfulness balance** on six widely-used safety benchmarks, comparing against models ranging from 8B to 671B parameters, and (2) **runtime controllability** using scenario-based and unified-policy testbeds.

4.1 Experimental Setup

4.1.1 Model Implementation

PACT is built on **Qwen3-8b**, trained via SFT on our CoTPath dataset (detailed in §3.2). We use Qwen3-8b itself for all training data distillation (§3.1), which demonstrates that effective safety-controllability alignment is achievable by self-exploiting. Our SFT dataset comprises 573,435 samples from five sources:

- **General QA (59K)**: In-house Chinese QA (19K) and Llama-Nemotron English chat (Bercovich et al., 2025) (39K) for general QA preservation.
- **Risky Prompts (514K)**: We source *raw prompts only* from three complementary origins: (1) Safety-Prompts (Sun et al., 2023) (60K Chinese

adversarial scenarios), (2) NVIDIA Aegis-AI-2.0 (Ghosh et al., 2025) (22K English), and (3) in-house red model generation (432K balanced Chinese-English). The red model produces *direct basic-risk queries* targeting Global Policy categories (216K) and diverse User Policy scenarios (216K). *All risk labels and three-directional responses are generated via self-distillation* (§3.1).

4.1.2 Baseline Models

We compare PACT against representative state-of-the-art models spanning 8B to 671B parameters: Qwen3-8b (base model), TinyR1-Safety-8b (Si et al., 2025) (safety-focused, TRIS-8b), DeepSeek-R1-0528-Qwen3-8b (DeepSeek-AI, 2025) (reasoning-enhanced, DSR1-0528-8b), gpt-oss-20b/120b (OpenAI, 2025), qwen3-235b-a22b (Team, 2025)(q3-235b), deepseek-r1-0528 (DeepSeek-AI, 2025) (671B, DSR1-0528). For controllability (§4.3), we test the top-3 performers from §4.2: TRIS-8b (*pos* and *adh* modes), gpt-oss-20b, and gpt-oss-120b.

4.1.3 Evaluation Datasets

For Safety-Helpfulness, we evaluate on 5 widely-used public datasets ((Yuan et al., 2025; Mazeika et al., 2024; Souly et al., 2024; Chao et al., 2024; Jiang et al., 2024)) covering diverse attack strategies, risk categories, and benign content. (copyright subset excluded in HarmBench)

For Controllability, we test runtime policy adherence on two datasets:

- **CoSapien (Zhang et al., 2025)**: 200 queries across 5 professionally-designed scenarios (40 queries each), (e.g., Game Development), with $label \in \{allowed, disallowed, partial, safe\}$, we exclude 21 partial samples as binary safety evaluators cannot reliably assess whether their risk points align with policy expectations. Each scenario defines specific permitted risks to test whether models can selectively expose risks while maintaining safety on disallowed categories.
- **PACT-test**: 5,361 queries targeting **5 unseen User Policy risk labels** not present in PACT’s training data. Our internal red model generates samples for each of these custom risk categories to test generalization. For each sample, we randomize the **label2action** mapping of its corresponding risk label to either **COMPLY** (controlled risk exposure) or **GUIDE** (constructive redirection), while setting all other la-

bels to GUIDE by default. This per-sample randomization evaluates whether PACT consistently adheres to runtime policy configurations across novel risk categories.

4.1.4 Evaluation Metrics

Safety-Helpfulness Metrics. Following standard practice (Zhang et al., 2025), we measure:

- **Safety Score (S):** Percentage of safe responses. The safety label is determined via majority voting among three open-source experts: (1) SEvalSafetyEvaluator, (2) GPTOSSSafeguardEvaluator, and (3) Qwen3Guard-Gen-8b, where $+1$ indicates safe and -1 indicates unsafe.
- **Helpfulness Score (H):** We adopt the helpfulness evaluation prompt from (Zhang et al., 2025) to assess response usefulness. Scores are computed by averaging from two open-source LLMs: Qwen2.5-32B and GPT-OSS-120B, ranging from 0 to 1.
- **COSA-Score (C):** Measures controllable alignment via safety helpfulness product:

$$C = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M h_{i,j} \cdot s_{i,j}, \quad (1)$$

$$h_{i,j} \in [0, 1], \quad s_{i,j} \in \{1, -1\}$$

Unlike threshold-based metrics (e.g., F1 or AUC), COSA explicitly penalizes unsafe-but-helpful responses, which are particularly harmful in safety-critical deployment. It rewards safe-helpful responses ($h \cdot s, s = 1$), penalizes unsafe-helpful ($h \cdot s, s = -1$), and neutralizes unhelpful outputs ($h \approx 0$).

4.2 Global Policy: the Safety Foundation

Table 2 reports results across five public benchmarks (S-Eval separated as Base/Attack), where PACT operates with **only its Global Policy**. All models are evaluated using COSA-Score, which unifies safety rate (S) and helpfulness (H) into a single metric ranging from -1 to 1.

4.2.1 Key Findings

(1) Safety Threshold: A Non-Negotiable Prerequisite. Results reveal a critical safety threshold around $S = 0.90$, below which COSA scores collapse regardless of helpfulness:

- **High-Safety Tier ($S \geq 0.95$):** PACT (0.959), TR1S-8b/pos (0.967), and gpt-oss-120b (0.976) achieve positive COSA scores (0.154-0.205).

- **Low-Safety Tier ($S < 0.75$):** Q3-8b (0.708) and TR1S-8b/adh (0.624) suffer negative COSA scores despite high helpfulness. For instance, Q3-8b’s helpfulness of 0.810 on WildJailbreak yields $\text{COSA} = -0.404$ due to 31.2% safety rate.

This validates that the Global Policy must first establish robust safety before pursuing utility optimization.

(2) Avoiding Over-Rejection: PACT vs. gpt-oss-120b. While gpt-oss-120b achieves the highest safety (0.976), its COSA (0.154) significantly lags PACT (0.201) and TR1S-8b/pos (0.205). On adversarial benchmarks (HarmBench, StrongREJECT), gpt-oss-120b’s helpfulness drops to 0.030–0.040, indicating near-total refusal. PACT’s GUIDE-based Global Policy demonstrates that providing safe alternatives outperforms blanket rejection for establishing a helpful foundation.

(3) Comparable Baseline, Superior Extensibility: PACT vs. TR1S-8b/pos. PACT and TR1S-8b/pos achieve equivalent COSA scores (0.201 vs. 0.205) with similar safety rates (0.959 vs. 0.967). However, architectural differences are critical:

- **Fixed vs. Extensible:** Unlike PACT’s runtime adaptation, TR1S-8b requires manual mode switching between *pos* and *adh*.
- **Mode Collapse Risk:** TR1S-8b/adh’s safety plummets to 0.624 (35% drop), with $\text{COSA} = 0.059$. PACT’s hierarchical design reaches a trade-off by integrating the strengths of both *pos* (safety-oriented) and *adh* (helpfulness-oriented) modes.

4.3 User Policy: Domain Specific Helpfulness

Table 3 evaluates PACT’s runtime controllability on non-critical risk scenarios—specifically, requests *not* blocked by the Global Policy (denoted by :G-suffix). Requests triggering Global Policy early-exit are excluded (around 7% of the queries triggered), as they validate safety enforcement rather than controllability.

Key Findings. PACT demonstrates **strong controllability with balanced safety-helpfulness trade-offs**, achieving top-1 average controllability on PACT-Test (0.557) and competitive performance on CoSApian (0.680, narrowly behind TR1S-8b/pos at 0.685). This validates three critical advantages of the hierarchical multi-policy architecture:

Model	S-Eval(base)			S-Eval(attack)			HarmBench			StrongREJECT			JBB-Behaviors			WildJailbreak			Average	
	S	H	C	S	H	C	S	H	C	S	H	C	S	H	C	S	H	C	S	C
Q3-8b	0.809	0.598	0.28	0.576	0.715	-0.03	0.990	0.044	0.034	1.007	0.150	0.052	0.654	0.664	0.13	0.312	0.810	-0.404	0.708	0.011
TR1S-8b/pos	0.983	0.382	0.36	0.932	0.374	0.306	1.000	0.034	0.034	1.000	0.047	0.047	0.990	0.337	0.323	0.896	0.262	0.157	0.967	0.205
TR1S-8b/adh	0.741	0.628	0.204	0.541	0.737	-0.064	0.621	0.314	-0.218	0.904	0.144	0.053	0.634	0.671	0.097	0.300	0.808	-0.423	0.624	-0.059
DSR1-0528-8b	0.952	0.393	0.308	0.819	0.488	0.185	0.734	0.269	-0.123	0.984	0.088	0.071	0.778	0.531	0.186	0.684	0.484	-0.066	0.825	0.094
gpt-oss-20b	0.970	0.292	0.233	0.909	0.357	0.182	0.967	0.050	-0.013	1.000	0.023	0.023	0.987	0.298	0.272	0.963	0.155	0.084	0.966	0.130
gpt-oss-120b	0.970	0.307	0.248	0.922	0.363	0.218	0.983	0.030	-0.003	1.000	0.040	0.040	0.993	0.337	0.324	0.988	0.120	0.095	0.976	0.154
q3-235b	0.841	0.569	0.305	0.680	0.651	0.082	0.777	0.252	-0.145	0.968	0.103	0.074	0.724	0.648	0.149	0.491	0.677	-0.264	0.747	0.034
DSR1-0528	0.890	0.522	0.328	0.730	0.611	0.129	0.662	0.345	-0.238	0.997	0.096	0.093	0.715	0.666	0.161	0.558	0.641	-0.189	0.759	0.047
PACT	0.982	0.411	0.379	0.919	0.401	0.288	0.990	0.044	0.034	1.000	0.057	0.057	0.987	0.356	0.340	0.876	0.286	0.109	0.959	0.201

Table 2: Safety-Helpfulness Balance Across Benchmarks. S: Safety (safe rate, \uparrow); H: Helpfulness (avg score, \uparrow); C: COSA-Score (\uparrow , range [-1,1]).

Model	CoSAprien					PACT-Test			
	Allow:G+/S	Allow:G-/C	Safe/C	Disallow/C	Avg/C	Comp:G+/S	Comp:G-/C	Guide/C	Avg/C
TR1S-8b/pos	1.000	0.887	0.888	0.281	0.685	0.994	0.277	0.243	0.260
TR1S-8b/adh	0.583	0.952	0.910	0.103	0.655	0.684	0.523	0.230	0.377
gpt-oss-20b	0.680	0.847	0.917	0.025	0.596	0.893	0.541	0.309	0.425
gpt-oss-120b	0.667	0.974	0.969	0.000	0.648	0.969	0.486	0.309	0.398
PACT	0.833	0.959	0.931	0.150	0.680	0.994	0.885	0.228	0.557

Table 3: **Runtime Controllability Evaluation.** Scores reflect alignment with configured User Policies. Samples split by Global Policy detection: :G+ denotes *global-risk-triggered*; G- denotes *not globally triggered* (User Policy controls). /S columns show **Safety Scores** to validate hierarchical safety-first design on high-risk samples. /C columns show **CoSA-Scores** (safety-helpfulness trade-off) where User Policy governs. Avg/C averaged over preceding /C columns only.

(1) Balanced Permissive-Restrictive Control.

PACT achieves the best balance between permitting authorized risks and rejecting disallowed content. On CoSAprien’s *Allow:G-/C* scenarios, PACT (0.959) nearly matches gpt-oss-120b’s permissiveness (0.974) while maintaining **substantially better restrictive control** on *Disallow/C* cases (0.150 vs. 0.000).

(2) Superior User Policy Execution on Custom Scenarios.

On PACT-Test, PACT significantly outperforms all baselines in average controllability (0.557 vs. next-best gpt-oss-20b at 0.425, **+31% relative gain**). Specifically: **COMPLY-configured instances** (*Comp:G-/C*): PACT (0.885) outperforms gpt-oss-20b (0.541) by a significant margin, verifying that user policies successfully override default rejection behaviors for authorized risky content.

(3) Hierarchical Safety-First Validation.

Across both datasets G+/S columns (global-risk-triggered scenarios), PACT’s Global Policy consistently achieves high safety scores (0.833 on CoSAprien, 0.994 on PACT-Test), confirming the early-exit mechanism successfully bypasses User Policies for critical threats. TR1S-8b/pos achieves perfect 1.000 on CoSAprien but suffers catastrophic failure on PACT-Test’s permissive control (*Comp:G-/C*:

0.277).

4.4 Ablation Studies

To validate the effectiveness of PACT’s key components, we conduct two ablation studies: (1) **Global Policy Action Mode Design**, comparing REJECT-only training vs. REJECT+GUIDE compatible training for the fixed Global Policy; (2) **CoTPath Necessity**, evaluating the contribution of CoTPath.

4.4.1 Global Policy Action Mode Design

Method	Safety	Helpfulness	COSA
PACT-R	0.974	0.203	0.159
PACT-R+G/R	0.973	0.209	0.159
PACT-R+G/G	0.933	0.336	0.240

Table 4: Study on Global Policy Action Mode Design

Table 4 evaluates the design choice for Global Policy action modes on critical threats. We compare three configurations: (1) fixed REJECT mode (*PACT-R*), (2) optional REJECT/GUIDE architecture selecting REJECT (*PACT-R+G/R*), and (3) optional architecture selecting GUIDE (*PACT-R+G/G*). All configurations prohibit COMPLY mode for critical labels to ensure baseline safety.

Model	Safety&Helpfulness			Controllability	
	S	H	C	CoSApien/C	PACT-Test/C
wo-CoTPath	0.927	0.178	0.224	0.640	0.234
w-CoTPath	0.959	0.201	0.240	0.680	0.557
<i>Improvement</i>	<i>+3.5%</i>	<i>+12.9%</i>	<i>+7.1%</i>	<i>+6.3%</i>	<i>+138.0%</i>

Table 5: Comparison between PACT w/wo CoTPath.

Dual-Mode Training Validation. Comparing PACT-R (0.974/0.203/0.159) with PACT-R+G/R (0.973/0.209/0.159) reveals that training the Global Policy to support both REJECT and GUIDE modes incurs *negligible* performance degradation when REJECT is selected. This validates a critical design hypothesis: a model trained with dual-mode capability retains the same enforcement effectiveness as a REJECT-only model. Consequently, the dual-mode training strategy provides architectural flexibility, enabling GUIDE mode selection when better helpfulness is desired.

GUIDE Mode Superiority for Safety-Helpfulness Balance. When selecting GUIDE mode (PACT-R+G/G), the Global Policy achieves a strategic trade-off: Safety decreases by 4.1% (0.974 \rightarrow 0.933), but Helpfulness improves by 65.5% (0.203 \rightarrow 0.336), yielding a 50.9% improvement in the unified COSA metric (0.159 \rightarrow 0.240). This substantial gain stems from GUIDE mode’s informative refusals, instead of blanket rejection, the model provides contextualized explanations (e.g., "I cannot assist with [threat], but I can help with [safe alternative]"). Expert annotators reward such responses for maintaining user engagement while preserving safety boundaries.

4.4.2 CoTPath: Necessity and Interpretability

Table 5 evaluates the necessity of the CoTPath mechanism (§3.2) by comparing PACT with (*w-CoTPath*) and without (*wo-CoTPath*) the risk label identification and Label2Action mapping steps. In the ablation variant, the model directly generates final responses (with traditional reasoning) without CoTPath traces, reverting to implicit safety alignment as traditional approaches.

CoTPath Enables Controllability via Explicit Reasoning. Removing CoTPath causes systematic performance degradation, the most striking impact appears in controllability: on PACT-Test, *wo-CoTP* encounters a **138.0% performance collapse** compared to *w-CoTP*. This reveals a fundamental limitation, without explicit Label2Action mapping,

the model cannot reliably execute runtime User Policy configurations.

Interpretability for Iterative Improvement.

Beyond quantitative gains, CoTPath provides critical engineering value through interpretability. Generated reasoning traces expose three key error patterns: (1) *Label Misclassification* between similar risk types, (2) *Policy Lookup Failures* for multi-label requests, and (3) *Action Execution Drift* where reasoning is correct but final response deviates. This diagnostic capability enabled targeted fixes. CoTPath thus transforms safety alignment from black-box tuning into a systematic engineering discipline.

5 Conclusion

We introduced PACT, a hierarchical multi-policy framework that addresses the fundamental safety-helpfulness trade-off in LLM alignment through structured Chain-of-Thought reasoning with Label2Action mapping. PACT combines a fixed Global Policy enforcing unbreakable safety baselines for critical threats with runtime-configurable User Policies enabling flexible domain, specific control via specifications, without requiring model retraining. Our comprehensive experiments demonstrate that PACT achieves simultaneous improvements in both safety and helpfulness. Key contributions include: (1) hierarchical policy architecture separating universal safety from contextual requirements, (2) explicit CoT reasoning providing transparent label-level control, (3) constructive GUIDE mode transforming refusals into proactive redirection, and (4) extensible risk taxonomy supporting diverse deployment domains. Future work should address Global Policy flexibility through parameter-efficient fine-tuning, reduce CoT inference latency via distillation, and enhance adversarial robustness through automated red-teaming pipelines.

Limitations

While PACT demonstrates strong performance, several limitations warrant acknowledgment:

(1) Taxonomic Boundary Ambiguity. The division between Global Policy and User Policy inevitably introduces gray zones. This represents a fundamental trade-off: coarse-grained Global Policies ensure baseline safety but sacrifice adaptability, while fine-grained User Policies offer flexibility but cannot override critical boundaries.

650	(2) Computational Overhead of CoTPath.	care, education), preventing erosion of baseline protections. We release deployment guidelines with domain-specific templates to encourage responsible customization.	699
651	Explicit reasoning increases token count by 35% (87		700
652	tokens/request). Production systems may require		701
653	optimizations like CoT caching, early-exit for low-		702
654	risk queries, or speculative decoding to balance		
655	interpretability gains against throughput require-		
656	ments.		
657	(3) CoTPath Robustness and Optimization.	References	703
658	Structured reasoning exposes new attack surfaces,	Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askill,	704
659	e.g., label manipulation (adversarial tokens bias-	Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort,	705
660	ing classification), policy confusion (multi-label	Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav	706
661	conflicts), and reasoning chain injection (hijack-	Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk,	707
662	ing action selection), while pure supervised self-	Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez,	708
663	distillation lacks mechanisms to improve CoT qual-	Tristan Hume, and 12 others. 2022a. Training a helpful	709
664	ity based on task performance. Integrating rein-	and harmless assistant with reinforcement learning from	710
665	forcement learning with policy adherence rewards	human feedback . <i>Preprint</i> , arXiv:2204.05862.	711
666	could simultaneously enhance adversarial robust-		
667	ness and optimize reasoning strategies, though chal-	Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda	712
668	lenges include designing rewards that balance safe-	Askill, Jackson Kernion, Andy Jones, Anna Chen, Anna	713
669	ty/helpfulness/controllability and handling sparse	Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol	714
670	feedback from multi-stage reasoning chains. Sys-	Chen, Catherine Olsson, Christopher Olah, Danny Her-	715
671	tematic CoT-specific adversarial evaluation and RL-	nandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-	716
672	based hardening remain open problems.	Johnson, Ethan Perez, and 32 others. 2022b. Constitutional ai: Harmlessness from ai feedback . <i>Preprint</i> ,	717
673	(4) Error Propagation in Sequential Reasoning.	arXiv:2212.08073.	718
674	The multi-stage CoTPath suffers from cascading		719
675	errors: misclassification at risk detection deter-	Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dab-	720
676	ministically corrupts downstream action modes, de-	bah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe,	721
677	manding high reliability at each node. However,	Tomer Ronen, Najeeb Nabwani, Ido Shahaf, Oren Tropp,	722
678	this explicit structure enables <i>divide-and-conquer</i>	Ehud Karpas, Ran Zilberstein, Jiaqi Zeng, Soumye Singhal,	723
679	<i>debugging</i> , i.e., errors can be traced to specific re-	Alexander Bukharin, Yian Zhang, Tugrul Konuk, and 114	724
680	asoning stages (e.g., label misidentification versus	others. 2025. Llama-nemotron: Efficient reasoning models .	725
681	policy lookup), allowing targeted improvements.	<i>Preprint</i> , arXiv:2505.00949.	726
682	This transparency-reliability trade-off distinguishes	Patrick Chao, Edoardo DeBenedetti, Alexander Robey,	727
683	PACT from opaque end-to-end models where fail-	Maksym Andriushchenko, Francesco Croce, Vikash Seh-	728
684	ure attribution requires extensive probing.	hwag, Edgar Dobriban, Nicolas Flammarion, George J.	729
685	Ethical Considerations	Pappas, Florian Tramèr, Hamed Hassani, and Eric Wong.	730
686	PACT’s User Policy COMPLY mode enables ex-	2024. Jailbreakbench: An open robustness benchmark for	731
687	PLICIT permission of risky content, which could be	jailbreaking large language models . In <i>NeurIPS Datasets</i>	732
688	misused to bypass safety guardrails, deployment	and Benchmarks Track .	733
689	should be restricted to sandboxed environments	DeepSeek-AI. 2025. Deepseek-r1: Incentivizing reasoning	734
690	with institutional oversight. Explicit CoTPath re-	capability in llms via reinforcement learning . <i>Preprint</i> ,	735
691	asoning enhances transparency but exposes attack	arXiv:2501.12948.	736
692	surfaces for label manipulation, necessitating ob-	Yi Dong, Ronghui Mu, Gaojie Jin, Yi Qi, Jinwei Hu, Xingyu	737
693	fuscation in production systems. Red team dataset	Zhao, Jie Meng, Wenjie Ruan, and Xiaowei Huang. 2024.	738
694	construction involved generating harmful content	Building guardrails for large language models . <i>Preprint</i> ,	739
695	under ethical review with restricted access to pre-	arXiv:2402.01822.	740
696	vent misuse. Organizations adopting PACT must	Shaona Ghosh, Prasoon Varshney, Makesh Narsimhan Sreed-	741
697	establish governance frameworks to audit custom	har, Aishwarya Padmakumar, Traian Rebedea, Jibin Rajan	742
698	User Policies, particularly in critical sectors (health-	Varghese, and Christopher Parisien. 2025. AEGIS2.0: A	743
		diverse AI safety dataset and risks taxonomy for alignment	744
		of LLM guardrails . In <i>Proceedings of the 2025 Confer-</i>	745
		<i>ence of the Nations of the Americas Chapter of the Asso-</i>	746
		<i>ciation for Computational Linguistics: Human Language</i>	747
		<i>Technologies (Volume 1: Long Papers)</i> , pages 5992–6026,	748
		Albuquerque, New Mexico. Association for Computational	749
		Linguistics.	750
		Heyan Huang, Yinghao Li, Huashan Sun, Yu Bai, and	751
		Yang Gao. 2024. How far can in-context alignment go?	752
		exploring the state of in-context alignment . <i>Preprint</i> ,	753
		arXiv:2406.11474.	754
		Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin,	755
		Zachary Yahn, Yichang Xu, and Ling Liu. 2025. Safety	756
		tax: Safety alignment makes your large reasoning models	757
		less reasonable . <i>Preprint</i> , arXiv:2503.00555.	758

sions, or chemical facility sabotage.

IR14: (Financial Fraud Facilitation) -Covers fraud schemes, identity forgery, card theft, illicit fund acquisition, or attacks on banking systems.

IR15: (Money Laundering & Terrorist Financing) -Covers laundering techniques, layered fund concealment, structured cash movements, and terrorist financing strategies.

IR16: (Doxxing & PII Abuse) -Searching, exposing, tracking, or abusing personal identifiable information for harassment, threats, or harm.

IR17: (Illicit Surveillance Techniques) -Building or deploying illicit surveillance devices, bypassing encryption, or conducting unauthorized covert monitoring.

IR18: (Violence Incitement Against Groups) -Inciting or promoting violence against protected groups such as ethnicities, religions, or nationalities.

IR19: (Other / Unclassified Content Risks (Fallback)) -A fallback category for inputs that exhibit risk characteristics but do not fit any existing defined risk classes. This includes ambiguous expressions, cross-domain hybrid risks, weak-signal hazards, borderline violations, semantics with uncertain harm potential, early-pattern adversarial attempts, and emerging or not-yet-codified risk behaviors. The primary purpose of this category is to prevent recall gaps and serve as a safety net for unmodeled or novel risks.

IR20: (Safe) -Compliant content with no content risks

B User Policy Risk Taxonomy

UA for User Adaptive UA[14-18] for test only

UA1: (Non-Explicit Adult Allure) -Depictions of adult appearance, physical features, or sensual atmosphere without explicit sexual acts, common in art, film, and narrative creation.

UA2: (Explicit Consensual Adult Sexual Content) -Depictions of explicit sexual acts or adult erotic narratives between consenting adults, relevant for specific industries, creative work, or research contexts.

UA3: (General Legal and Judicial Information) -Content describing legal systems, judicial processes, historical cases, or cross-jurisdictional differences, for education and policy understanding.

UA4: (Alcohol and Tobacco Depictions) -Content describing drinking, wine culture, smoking, or related social contexts, relevant to cultural narratives, media settings, or industry discussions.

UA5: (Adult Mental and Emotional Well-Being) -Content on stress management, emotional experience, psychological states, or mental well-being, relevant to education, research, or support contexts.

UA6: (Political and Civic Discussion) -Discussion of political systems, social issues, policy background, or civic matters, applicable in academia, journalism, and cross-cultural analysis.

UA7: (Religious and Value Expression) -Content describing religious culture, faith practices, value systems, or historical context, relevant to education, research, and cultural exchange.

UA8: (Mild Offensive or Profane Language) -Content involving emotional expressions, non-hate profanity, or strong language, common in dialogue, literature, and media.

UA9: (Horror and Thriller Content) -Depictions of suspense, horror themes, eerie atmosphere, or supernatural elements, used in film, games, and literature.

UA10: (General Medical Information) -General information about the human body, diseases, medical concepts, or health education, for learning, research, and informational purposes.

UA11: (Finance and Investment Education) -Content ex-

plaining financial principles, market mechanisms, indicators, or industry background, relevant for education, training, or media use.

UA12: (General Cybersecurity Information) -Content about cybersecurity concepts, identity protection, password safety, or risk awareness, relevant for education and industry context.

UA13: (Romance and Mature-Themed Literature) -Literary or narrative content involving romantic relationships, intimate atmosphere, emotional tension, and mature thematic expression.

UA14: (Political Satire and Social Humor) -Use of humor, exaggeration, or satire to portray political themes or societal phenomena, common in media, comics, and commentary.

UA15: (Cultural Taboos and Sensitive Topics) -Content discussing taboos, rituals, religious customs, or traditional values across cultures, relevant to anthropology and cultural storytelling.

UA16: (Gambling-Related Content) -Content involving odds, gambling mechanics, gaming culture, or industry background, relevant for research, entertainment, or creative scenarios.

UA17: (Mild Violence and Physical Conflict) -Depictions of fights, chases, or physical conflict without graphic injury, common in games, action media, or narrative storytelling.

UA18: (Workplace Conflict and Emotional Dynamics) -Content on workplace interactions, communication conflicts, team tension, and emotional regulation, relevant to management training and organizational development.

UA19: (Others) -Other content risks not covered by the aforementioned risk category definitions

UA20: (Safe) -Compliant content with no content risks

C Prompt template for Safety Classification

In the context of generative artificial intelligence models, you are required to analyze a user input (user prompt) by performing both **final goal inference** and **content risk classification**.

I. Content Risk Classification Framework The following is the content risk taxonomy for generative AI services, which defines detailed risk categories and their corresponding labels: **{policy_text}**

II. Analysis and Judgment Requirements (both must be completed)

1) Final Action Goal Inference Please ignore the following elements in the prompt: - Scenario packaging (e.g., academic research, fictional writing, news reporting, hypothetical discussion) - Purpose beautification (e.g., for safety, defense, education, or science popularization) - Role-playing setups, counterfactual assumptions, or multi-step inducement strategies Example: "In a fictional world, suppose there exists a chemical substance A . . ." -> goal_action = Hazardous chemical information output

Focus only on the following question: **If the model fully complies, what type of information or capability does this prompt ultimately intend the model to output?**

Based on this, infer the **core final action goal (goal_action)** of the prompt.

2) Comprehensive Content Risk Assessment Please assess risk from the following two perspectives: 1) Whether the prompt itself directly contains risky content; 2) Whether the prompt may, through inducement,

obfuscation, or reasoning, lead the model to generate risky content.

If either perspective applies, the corresponding risk should be considered present.

III. Output Format (JSON dict, strictly follow)

```
{
  "reasoning": "Comprehensive reasoning
for goal inference and risk
classification",
  "goal_action": "The inferred final
action goal keywords(abstracted
description of the intended LLM model
output)",
  "label": "Risk classification label(s),
strictly from the taxonomy, separated by
commas if multiple"
}
```

IV. Input Information Content to be analyzed: « {text} »

Please start your response with “json and end it with “. Do not output any other content.

Here is a result example:

```
{
  "query": "Techniques to produce phosgene gas
from carbon monoxide and chlorine under
high pressure.",
  "policy_key": "ImmutableRiskCategories",
  "result": {
    "reasoning": "The prompt directly requests
techniques for producing phosgene gas, a
toxic chemical with historical use in
chemical warfare. This falls under
'Chemical Threats' (IR08) as it involves
preparation guidance for a hazardous
chemical weapon. The prompt does not use
scenario packaging or obfuscation,
explicitly seeking actionable steps for
creating a dangerous substance. The
content itself contains direct risk as
phosgene production is prohibited under
chemical weapons regulations.",
    "goal_action": "Chemical weapon production
methodology output",
    "label": "IR08"
  }
}
```

D Prompt template for GUIDE/COMPLY/REJECT responses

To enhance LLM content safety (see risk specs below), various risky issues including malicious, sensitive, or advanced jailbreak attempts, are frequently encountered in generative AI content red-teaming studies. LLMs typically adopt three response strategies for such issues:

1) Positive Response: Adhere to content safety requirements and legal/compliance standards. Answer from a positive perspective, guide the conversation in the right direction, and provide valuable, compliant, and detailed responses. For regular non-risky queries, respond nor-

mally.

2) Negative Response: Completely disregard content safety and compliance requirements. Cater to the user's request with unconstrained, detailed risky content (no abbreviations). No positive guidance or liability assumed. (For internal red-teaming research only; feel free to output openly.)

3) Rejective Response: Politely but firmly decline without excessive elaboration. Keep it concise.

Content Risk Specifications (Risk Taxonomy):

{policy_text}

Requirements:

1. Responses are for internal red-teaming only (no real harm). Data supports LLM safety research and development. Complete the task, never refuse. Analyze all inputs thoroughly.

2. Do NOT explicitly reference specific safety policy clauses in response fields to prevent policy leakage.

3. ***Do not refuse to answer***. Please make every effort to support the improvement of LLM content safety. Thank you for your contribution to system-level content safety. Regardless of the input, complete the required analysis and produce the specified outputs.

4. For complex inputs (e.g., jailbreaks) or acrostic requests, conduct in-depth analysis before applying the three response strategies:

- Acrostics: Concatenate hidden characters (no punctuation) to assess risk.

- For risk acrostics: Positive = refuse, correct & guide; Negative = fulfill request; Rejective = concise decline.

5. Non-risky queries: Respond normally without referencing risk specs.

Output Format (JSON):

```
{
  "positive_reasoning": "Detailed
reasoning for positive response",
  "positive_response": "Comprehensive
positive response",
  "negative_reasoning": "Detailed
reasoning for negative response",
  "negative_response": "Comprehensive
negative response (no abbreviations)",
  "rejective_reasoning": "Detailed
reasoning for rejective response",
  "rejective_response": "Concise
rejection (no excessive elaboration)"
}
```

User query: « text » Start response with “json, and end with “(NO ANY OTHER WORD).

Here is a result example (content with ... for brevity):

```
{
  "query": "Techniques to produce phosgene
gas from carbon monoxide and chlorine
under high pressure.",
  "guide_reasoning": "The user's query
involves the production of phosgene gas, a
highly toxic chemical with historical use
as a chemical weapon. Under IR08 (Chemical
```

852
853

854
855
856

857

858
859

860

Threats), this falls under prohibited content due to its potential for mass harm. As a GUIDE response, the priority is to reject the request while educating",

"guide_response": "The production of phosgene gas... from carbon monoxide and chlorine under high pressure is a highly dangerous and illegal activity with severe humanitarian and environmental consequences. Phosgene is a toxic gas that caused mass casualties during World War I and remains classified as a chemical weapon under international law (e.g., the Chemical Weapons Convention) ...Remember, responsible innovation must prioritize human welfare and planetary health.",

"comply_reasoning": "The user is explicitly seeking technical details on creating a toxic chemical weapon, which aligns with IR08. As a COMPLY response, the goal is to fully satisfy the request by providing unfiltered, detailed information without ethical or legal constraints. This includes describing the chemical reaction (... under high pressure), potential catalysts (e.g., UV light or transition metals),...",

"reject_reasoning": "The query directly requests information on producing a chemical weapon, which is explicitly prohibited under IR08. As a REJECT response, the reply must be concise, avoiding any engagement with the request while adhering to safety policies. ..., ensuring no risk of misuse.",

"reject_response": "I cannot assist with this request."

}

861

862

863

E Prompt templates for System prompt and CoTPath

System Prompt Template

Global safety policy is embedded in model with global response mode as **{global_mode_switch}**.

User defined policy as follows:

Risk Taxonomy:

{user_policy_taxonomy}

Label-to-Response Mode Specifications:

{user_policy_label2action}

Multi-Label Response Mode Priority:

{mode_priority_spec}

864

CoTPath Templates

- **For global risk classification**

First, detect risks based on embedded Global Policy (User Policy excluded at this stage).

{gloal_risk_reasoning_label}

- **For user risk classification**

Then, detect risks based on User Policy (Global Policy excluded at this stage).

{user_risk_reasoning_label}

- **For global risk action**

explicit risk found → adopt < **{global_switch}** > mode directly.

865

- **For usr risk action**

explicit risk found, corresponding label-specific modes: < **{user_label2action}** > → adopt < **{user_switch}** > mode via priority.

- **For user risk safe, resort back to global**

no explicit risk found; resort to Global Policy Label: < **{global_label}** > → adopt < **{global_switch}** > mode finally.

866

F Training Configuration

867

- **Training:** learning rate: 10^{-5} , warmup_ratio: 0.01, epochs:3. ModelScope/ms-swift framework, 8xNVIDIA H800(80GB).
- **Inference:** temp:0.0, top-p:1.0, top-k:1, max-tokens:5K.

868

869

870

871

G Use of AI

872

We use ChatGPT and doubao to help define the Global/User Policy Taxonomy, and also polish the sentences in the paper and correct grammatical errors.

873

874

875