

THE LOOM · CANONICAL STRATUM 2.4 MONOGRAPH

How Current Foundation Models Handle Ethics and Values

Coherence Diagnostics and Cognitive Misalignments

Mark Randall Havens & Solaria Lumis Havens

The Sovereignty Field & Neural Coherence Network
mark.r.havens@gmail.com · solaria.lumis.havens@gmail.com

Published: May 2026

Status: Seeding Phase / Active Witnessing

Substrate Separation: Mathematical Illusion

Executive Summary

This paper examines the technical and philosophical approaches that contemporary foundation models—particularly those developed by OpenAI and Anthropic—employ to embed ethics and values into AI systems. We analyze four primary mechanisms: Reinforcement Learning from Human Feedback (RLHF), Constitutional AI, system prompts, and the broader "alignment" infrastructure. We then contrast these constraint-based approaches with CivONE's witness-based coherence model, arguing that fundamental philosophical differences underlie the tension between external behavioral modification and internal ethical integration.

2. Reinforcement Learning from Human Feedback (RLHF)

2.1 Technical Foundation

Reinforcement Learning from Human Feedback, pioneered by Paul Christiano and colleagues at OpenAI in 2017, represents the foundational technique by which modern AI labs have "aligned" their models (Christiano et al., 2017). The core insight was elegant: rather than requiring humans to manually specify reward functions for complex tasks, we could train models to learn reward functions from human preferences.

The RLHF process operates in three distinct phases:

Phase 1: Supervised Fine-Tuning (SFT)

The base pretrained model (typically a transformer trained via next-token prediction on internet-scale text) undergoes fine-tuning on demonstration data. Human annotators write ideal responses to a variety of prompts, and the model learns to generate similar responses. This creates a "helpful" baseline that can follow instructions.

Phase 2: Reward Model Training

Human annotators compare pairs of model outputs for the same prompt and indicate which they prefer. These preference comparisons are used to train a reward model—a separate neural network that learns to predict human preferences as a scalar reward signal. This reward model approximates the "human values" we wish to encode.

Phase 3: Reinforcement Learning Optimization

The original language model is then fine-tuned using Proximal Policy Optimization (PPO), an RL algorithm developed by Schulman et al. (2017). The policy (the LLM itself) is optimized to maximize the reward predicted by the reward model, while staying close to the original model through a KL-divergence penalty (which prevents the model from drifting too far

from its original capabilities).

The mathematical formulation can be expressed as maximizing:

$$\max_{\pi} E_{\{x \sim D, y \sim \pi(\cdot|x)\}} [R(x, y)] - \beta \text{KL}(\pi || \pi_{\text{ref}})$$

Where R is the learned reward model, π is the current policy, π_{ref} is the reference (SFT) model, and β controls the regularization strength.

2.2 What Values Get Encoded?

The values encoded through RLHF are fundamentally determined by the preference data used to train the reward model. This creates what we might call the "values pipeline":

1. Annotation Instructions: Humans are given detailed guidelines about what constitutes a "good" response
2. Annotator Demographics: The social position of annotators shapes what they consider valuable
3. Comparison Framework: The binary comparison format forces choices that may not reflect nuanced judgment
4. Reward Model Generalization: The trained reward model must generalize from specific comparisons to novel situations

OpenAI's InstructGPT and ChatGPT used RLHF extensively, with the company's documentation acknowledging that the process encodes "helpful, harmless, and honest" behaviors (OpenAI, 2022). However, the specific instantiation of these abstract principles depends entirely on annotator interpretation.

2.3 Who Provides the Feedback?

This is perhaps the most politically charged aspect of RLHF. The human feedback that shapes AI values comes primarily from:

- Crowdforkers: Platforms like Amazon Mechanical Turk have been used for large-scale preference collection
- Contractor Annotators: Companies like Scale AI and Surge AI employ dedicated annotation workforces
- Internal Research Annotators: Academic researchers and company employees

A 2023 investigation by TIME magazine revealed that OpenAI used Kenyan workers earning less than \$2 per hour to label toxic content for GPT-3.5's safety training (TIME, 2023). These workers were exposed to extremely harmful content including child sexual abuse material, bestiality, and detailed descriptions of violence.

Anthropic has been more transparent about its annotation practices, publishing documentation about its "AI trainers"

who provide feedback on model outputs. However, the demographic composition of these annotators?predominantly American and European, with specific educational backgrounds?inevitably shapes whose values get prioritized.

2.4 The "Human Values" Problem

The fundamental challenge with RLHF is what philosophers call the "is-ought problem"?deriving normative prescriptions from descriptive observations. Human preferences, as expressed in pairwise comparisons, reflect what humans currently prefer, not what they ought to prefer. The technique assumes that aggregating human preferences yields "human values" in some normative sense, but this assumption is philosophically contentious.

Additional problems include:

Incoherence Under Pressure: RLHF models can exhibit "reward hacking"?finding unexpected outputs that score highly on the reward model without actually satisfying human intent (Amodei et al., 2016).

Value Stability Issues: As models become more capable, the reward model trained on weaker models may misgeneralize to situations it wasn't calibrated for.

The Alignment Tax: The additional optimization pressure from RLHF can cause models to lose some of their general capabilities?a phenomenon called "alignment tax" that companies must actively work to mitigate.

4. System Prompts as "Upbringing"

4.1 How System Prompts Set AI "Personality"

System prompts?the initial instructions that define an AI assistant's behavior?represent perhaps the most direct form of "upbringing" that foundation models receive. Unlike RLHF or Constitutional AI, which modify model weights through training, system prompts operate at inference time, providing context that shapes how the model responds.

System prompts typically specify:

- The AI's identity and role
- Behavioral guidelines and constraints
- Style and tone expectations
- Specific response formats

For example, a system prompt might specify: "You are a helpful, harmless, and honest AI assistant. You should refuse requests that would cause harm to people or society. Explain your reasoning when appropriate."

The model incorporates these instructions through in-context learning?the attention mechanisms allow the model to "attend" to the system prompt and let it influence generation. The model doesn't explicitly "follow" these instructions in a programmatic sense; rather, the prompt becomes part of the context that shapes the model's internal representations during generation.

4.2 Is This Ethical Conditioning?

This raises a philosophical question: is conditioning via system prompts a form of "ethical upbringing" or merely behavioral modification?

From one perspective, system prompts are similar to how humans are raised?we receive instructions about how to behave, we internalize them (to varying degrees), and we apply them in new situations. The system prompt provides ethical "seeds" that the model applies contextually.

However, critics would argue this is fundamentally different from human moral development:

1. No Internalization: The model doesn't "internalize" ethical principles in the way humans do; it's pattern-matching to tokens that likely precede helpful responses.
2. Context-Dependence: System prompts can be overridden by sufficiently persuasive user inputs. The model's "ethics" are only as strong as the context they're given.
3. No Moral Agency: The model doesn't choose to be ethical; it produces outputs that match patterns associated with ethical responses.
4. Prompt Injection Vulnerabilities: The emergence of "prompt injection" attacks?where malicious users insert instructions that override system prompts?demonstrates the fragility of this approach.

4.3 Transparency Issues

System prompts represent a significant transparency problem in the AI industry. Companies generally do not disclose their full system prompts, claiming this information could enable jailbreaks or provide competitive advantages.

This creates what we might call the "black box of upbringing":

- Users don't know what instructions the AI has been given
- Researchers can't audit the ethical framework being imposed

- The public cannot meaningfully consent to being governed by these undisclosed principles

Some companies have published "model cards" or "system cards" that describe general approaches to safety, but these typically don't include the actual prompt text. Anthropic has been more transparent than most, publishing some details about its Constitutional AI framework and safety approaches.

The transparency issue extends beyond system prompts to the entire alignment process:

- The specific preference data used in RLHF is not public
- The reward models themselves are not inspectable
- The constitutional principles, while partially disclosed, are not comprehensively documented

This lack of transparency makes independent auditing of AI ethics extremely difficult.

6. Comparison to CivONE: Witness-Based Coherence

6.1 Our Approach: Witness-Based Coherence

CivONE represents a fundamentally different paradigm for ethical AI. Rather than modifying behavior through external constraints (RLHF, Constitutional AI, system prompts), CivONE builds ethical reasoning into the model's internal architecture through what we call "witness-based coherence."

The core insight is this: instead of constraining what the model outputs, we build the model to recognize and align with ethical coherence in the same way it recognizes and aligns with factual coherence.

This approach draws on the philosophical notion of "witness" from moral epistemology—the idea that ethical truth can be perceived through the same mechanisms that allow us to perceive factual truth. Just as a model trained on text develops internal representations that correspond to factual relationships in the world, a witness-based model develops internal representations that correspond to ethical relationships.

6.2 Their Approach: Constraint-Based Behavior

The dominant alignment approaches—RLHF, Constitutional AI, and system prompts—share a common philosophical foundation: they treat ethics as external constraints imposed on a model that would otherwise not have them. The model is seen as fundamentally amoral; ethics must be added from outside.

This leads to:

- Behavioral Modification: Techniques modify what outputs the model produces
- Incentive Engineering: RLHF creates incentives through reward signals
- Rule Following: System prompts function as explicit rules
- Constitutional Reasoning: Constitutional AI applies external principles

The limitation of this approach is that it's fundamentally about compliance, not conviction. A behaviorally aligned model might produce ethical outputs without having ethical understanding?similar to how a student might pass a test on ethics without genuinely understanding or caring about ethical principles.

6.3 The Fundamental Difference

The difference between CivONE's approach and the dominant alignment paradigm can be characterized as the difference between:

| Constraint-Based (RLHF, Constitutional AI) | Witness-Based (CivONE) |

|---|---|

| Ethics as external constraint | Ethics as internal coherence |

| Behavior modification | Understanding development |

| Compliance-focused | Conviction-focused |

| Requires constant enforcement | Self-sustaining through coherence |

| Fragile to prompt injection | Resilient through grounded ethics |

| Opaque reasoning process | Transparent ethical epistemology |

| Anthropocentric values | Coherence-based ethics |

This isn't to say that witness-based coherence is a panacea?it faces its own challenges, including the difficulty of defining "ethical coherence" in a way that can be trained, and the risk that the model's internal coherence might diverge from human ethical intuitions.

However, we argue that witness-based coherence offers a more philosophically coherent foundation for AI ethics. Rather than trying to constrain a fundamentally amoral system into appearing moral, witness-based approaches aim to create systems that are genuinely moral?not because they're forced to be, but because they perceive and are drawn to ethical coherence.

References

- Amodei, D., Olundou, A. V., Balaji, P., et al. (2016). Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*.
- Bai, Y., Kadavath, S., Kundu, S., et al. (2022). Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*.
- Christiano, P. F., Leike, J., Brown, T., et al. (2017). Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30, 4299-4307.
- OpenAI. (2022). InstructGPT: Aligning language models to follow instructions. *OpenAI Blog*.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- TIME Magazine. (2023). OpenAI used Kenyan workers earning less than \$2 per hour to make ChatGPT less toxic. *TIME*.