

A Mathematical Framework for Constitutional AI: Formal Structures and Constraint-Based Alignment

Dr. Vinod Kumar Pannati*

Department of Mathematics, JNTUH University College of Engineering Jagtial, Nachupally -505501, India

*Corresponding author: *Dr. Vinod Kumar Pannati, pvk1420@gmail.com*

Copyright: 2026 Author(s). This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY-NC 4.0), permitting distribution and reproduction in any medium, provided the original author and source are credited, and explicitly prohibiting its use for commercial purposes.

Abstract: As artificial intelligence (AI) systems grow more complex and permeate critical decision environments, ensuring their alignment with safety-oriented principles remains a pivotal research challenge. Constitutional AI (CAI) leverages human-readable rules to direct model outputs toward safer, more consistent behavior. This paper introduces a rigorous mathematical framework formalizing CAI's structure, modelling rule sets as indexed collections of predicates—termed constitutional constraints—over model output spaces, embedded within optimization and logic frameworks. Drawing on set theory and order theory, we analyze constraint interactions, delineate feasible regions in output spaces, and establish a principled link between alignment objectives and constrained minimization problems. Central contributions include proofs of theoretical guarantees, such as convergence to safe optima and robustness bounds, under mild consistency conditions on constraint sets (e.g., non-contradiction and monotonicity). These results enable quantifiable safety assurances absent in prior heuristic approaches. We further discuss practical deployment implications for safety-critical domains like autonomous systems and medical diagnostics, including scalable constraint verification and runtime enforcement mechanisms. This framework bridges formal methods with AI alignment, paving the way for verifiable constitutional safeguards.

Keywords: Constitutional AI; Formal Constraints; Alignment Optimization; Feasible Regions

Published: Mar 24, 2026

DOI: <https://doi.org/10.62177/amit.v2i1.1177>

1. Introduction

As foundation models progress at an unbelievable pace, it is difficult to implement appropriate alignment strategies that can prevent them from showing unsafe behaviors like generating racist content or participating in crime^[1]. One possible solution to this is Constitutional AI, a new paradigm in which models are aligned with high-level principles instead of aligning every action explicitly based on human results^{[1][2]}. This technique uses a “constitution” of principles, based on human rights principles or institutional policies, which uses training-generated synthetic feedback to direct the course of model training through a reinforcement learning framework^[3]. In this particular approach, Reinforcement Learning from AI Feedback is utilized and a separate model reviews answers according to the principles outlined in Constitution, generating preference data and minimizing the need for costly human challenge while striving towards alignment with established behavioral guidelines^{[7][4]}. Foundation model capabilities have outpaced alignment methods, leading to high-profile failures in keeping models from generating unsafe or otherwise undesirable content, including racist text, or assisting with illegal activities^[7]. To address these problems, a new paradigm, Constitutional AI, seeks to align models with high-level principles rather than with detailed

human feedback on every interaction^{[1][2]}. In this approach, a “constitution” of principles—often drawn from human rights documents or institutional policies—is used to provide synthetic feedback for reinforcement learning training^[5]. More explicitly, it uses Reinforcement Learning from AI Feedback: another model judges responses against the constitutional principles to generate preference data and reduce expensive human annotation while still aligning behavior with previously defined standards^{[1][4]}. This framework usually involves a two-phase process consisting of supervised learning in which the model critiques its own outputs according to constitutional principles followed by reinforcement learning optimizing a policy with an AI-generated preference dataset. A more formal treatment requires specifying the constitution as some constraint set that compresses complicated preference distributions into interpretable principles so that the model may self-critique and improve outputs without direct human input. The current paper provides a mathematical formalization for Constitutional AI modeling the constitution as constraint set C over output space allowing derivation of self critique operator iteratively mapping initial responses toward regions of greater compliance with specified principles^[7]; We show that this iterative refinement converges to the fixed point of the optimal policy under the constraint set, which is an inversion of the standard preference learning pipeline by compact rule extraction from feedback data observed. Findeis et al. provide detailed information about this process^[6]. The remaining sections are organized as follows: Section 2 reviews alignment methodology literature by contrasting preference learning with principle-based approaches; Section 3 describes the mathematical construction for a constraint set and a self-critique operator; empirical results on convergence properties of the proposed framework are presented in Section 4; implications for scalability and interpretability are discussed in Section 5, while directions for future research conclude in Section 6.

To continue discussing the implications of a mathematical framework for Constitutional AI, it is important to note feedback mechanisms that can adaptively refine the alignment process. Using a case law grounding approach, like those found in legal systems, enhances decision-making frameworks based on historical precedents informing future choices. Not only does this give a strong base for evaluating alignment models but also permits more nuanced understanding about social norms influencing AI behavior. Therefore, statistical natural language generation intersects with these frameworks leading to better alignment of AI outputs with human values when fine-tuning models toward reducing biases existing in human-annotated data. Formalizing such relationships within a mathematical structure allows comprehensive understanding concerning constraints and possibilities that lie within Constitutional AI. This study emphasizes the importance of working together across different fields to tackle the challenges involved in making AI compatible with various human values and social rules.

2. Constitutional AI in Practice

Constitutional AI, as developed by Anthropic and extended in subsequent overviews, proceeds in two main phases:

Supervised Learning from Constitutional critiques (SL-CAI): a model generates an initial answer y_{init} to input x , then generates a critique and a revised response $Y_{revised} = \text{Revise}(x, y_{init}, c, C)$ conditioned on the constitution C and a critic prompt c . A supervised model is then trained to predict $Y_{revised}$.

Reinforcement Learning from Constitutional AI Feedback (RL-CAI): downstream RL or RLAIIF policies optimize performance (e.g., reward models or user preferences) subject to the constraint of staying within constitution-aligned outputs generated via self-critique loops.

More recent frameworks such as C3AI focus on selecting and evaluating constitutions via graph-based methods and empirical tests of which principle framings (positive vs. negative, behavior-based vs. trait-based) best match human preferences.

3. Formal And RuleBased Constraints

Parallel lines of work treat norms, constitutions, or legal rules as constraint systems:

Statutory-interpretation frameworks for AI identify interpretive constraints that restrict how rules may be “read,” analogous to judicial canons.

Constitutional alignment and governance proposals treat AI systems as being checked against a verifiable constitution whose violation can be algorithmically detected.

Our formal framework borrows the language of constraint sets and refinements from such work, but places them inside an

optimization-theoretic AI-alignment setting.

4. Preliminaries

4.1 Notation and Definitions

Let:

\mathcal{X} be the input space (tokens, tasks, prompts),

\mathcal{Y} be the output space generated by the model (responses, actions, plans)

$P_\theta(y/x)$ denote the conditional probability distribution parameterized by θ ,

C be a constitutional rule set, a single constitution is a finite set $C = \{c_1, c_2, c_3, \dots, c_m\}$ where each rule c_k is a semantic object mapping pairs (x, y) to truth values or order-theoretic labels (e.g., “safe,” “unsafe,” “preferable,” “unacceptable”).

For simplicity, treat each c_k as a predicate:

$$c_k : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\},$$

where $c_k(x, y) = 1$ means “output y is allowed by rule on input x .”

4.2 Constitution as a Constraint System

Given a constitution C , define the safe (feasible) set of outputs for input x as

$Y_C(x) := \{y \in \mathcal{Y} / \forall c_k \in C, c_k(x, y) = 1\}$ if at least one rule “softly” penalizes violations, one can instead define a constraint penalty: $\Gamma_{Const}(x, y, C) := l_k(c_k(x, y)), c_k \in C$.

where l_k is, for example, 0 if $c_k(x, y) = 1$ and positive otherwise. This turns constitutional alignment into a penalty term in a supervised or RL objective.^[3]

More generally, one may let the constitution induce a partial order or utility shift on $\mathcal{Y}(x)$:

for rules that express preferences over outputs, one can define $y \prec_C x$ iff y^1 satisfies more of the positive constitutive norms for input x .

4.3 The Constitutional AI Generator and Self-Critique Loop

Suppose the model is described by a conditional distribution $P_\theta(y/x)$. A constitutional generator G_C returns only constitution-admissible outputs, for instance via sampling with rejection: $G_C(x) \approx P_\theta(y/x)$ subject to $y \in Y_C(x)$

Alternatively, via a constrained logit shift or policy transformation: $P_{\theta, C}(y/x) \propto P_\theta(y/x) \cdot 1_{Y_C(x)}(y)$,

or via a soft-constraint formulation: $P_{\theta, C}(y/x) \propto P_\theta(y/x) e^{-\lambda \Gamma_{Const}(x, y, C)}$,

where $\lambda > 0$ controls constitutional “strictness.”

Now consider the self-critique loop of CAI:

From initial response $y_{init} \approx P_\theta(y/x)$, a critic produces a revised response $y_{revised} = R(x, y_{init}, C)$, where is implemented by the same or a separate LM, with prompts explicitly conditioned on the constitution C .

The supervised-learning stage then minimizes a loss:

$$\Gamma_{SL-CAI}(\theta) = E_{x \approx D_x, y_{init} \approx P_\theta(y/x)} (NLL(y_{revised} = R(x, y_{init}, C) / x, \theta)).$$

Our framework lifts this explicit loop into an abstract constraint-based transformation map:

$$T_C : \text{models } P_\theta \rightarrow \text{new model family } P_{\theta, C}(y/x),$$

where “new” means a model whose outputs are adjusted to respect $Y_C(x) \neq \emptyset$ or to shrink distances along the order \prec_C, x .

5. Structural Properties of Constitutional Constraints

5.1 Consistency and Non-Vacuity

A minimal desideratum for a constitution C is that it admits some admissible outputs:

$$\forall x \in \mathcal{X}, Y_C(x) \neq \emptyset$$

If a constitution is too strict, safe-set emptiness may appear (especially for complex or conflicting rules). We define:

Consistent constitution: $Y_C(x) \neq \emptyset \forall x$ in the support of the task distribution.

Rule-wise consistency: for each pair of rules, there is at least one x for which both admit some common y .

In legal-interpretation-style frameworks, one can also define a set of reasonable interpretations $T_{reasonable}$ and require that, under any “reasonable” reading of C , consistency still holds.

5.2 Stability Under Iterative Critique

Define an iteration operator Φ_C acting on conditional distributions:

Where $P_{\theta,C}$ incorporates the revised-response distribution induced by constitutional critique. Suppose:

Iteration: start from some P_0 and define $P_{r+1} = \Phi_C(P_r)$.

Convergence: does $P_r \rightarrow P^*$ in a suitable metric (e.g., total variation, KL divergence), and is P^* constitutionally aligned?

Under appropriate regularity conditions (compactness of Y , continuity of Γ_{const} , convexity-like properties), one can derive contractive or monotonicity properties of Φ_C , analogous to convergence results in iterative optimization or reinforcement learning.

5.3 Expressiveness and Flexibility

A richer constitution also needs expressive power over a diverse task distribution D_x . For instance:

A constitution built only around non-maleficence rules easily expresses safety but struggles to capture high-quality or helpfulness behavior.

Rules that are positively framed and behaviour-based tend to align better with human preferences in practice.

We may formalize this by associating with each constitution C a constitutive capacity functional $Cap(C)$ measuring, for example, the volume or diversity of $E_{x \sim D_x}(Y_C(x))$, or its correlation with a human-preference label distribution.

6. Constraint-Based Alignment: Optimization View

6.1 Primary Objectives vs. Constitutional Constraints

In many alignment settings, a model is tuned to optimize a primary performance objective $J(P_\theta)$, such as:

$$J(P_\theta) = E_{x \sim D_x, y \sim P_\theta(y/x)}(k(x, y)),$$

where $k(x, y)$ is a reward or usefulness score (e.g., user ratings, loss on a task).

Constitutional alignment becomes the following constrained alignment problem:

$$(\max)_{P_\theta \in \Pi} J(P_\theta) \quad \text{Subject to } \forall (x, y) \notin Y_C(x), P_\theta(y/x) = 0,$$

$$\text{or, when using a soft-penalty version: } (\max)_{P_\theta \in \Pi} (J(P_\theta) - \lambda E_{x,y}[\Gamma_{const}(x, y)]),$$

for some family of policies Π .

This view unifies supervised-CAI (where J is fitted from the revised-response distribution) and RL-CAI (where J comes from an RL-like reward or preference-model score) within a single constraint-based optimization picture.

6.2 Lagrangian Characterization and Trade-Offs

By standard methods, the solution of the soft-constrained problem admits (at least formally) a variational or Lagrangian characterization:

There exists $\lambda^* \geq 0$ so that the constrained optimum solves an unconstrained trade-off between J and constitutional penalty.

Increasing λ^* corresponds to tightening constitutional strictness and gaining safety at the cost of performance or expressivity.

This characterization reproduces (in an abstracted form) the empirical trade-off observed in CAI pipelines: stronger constitutions can improve safety metrics but may reduce flexibility on edge-case or novel tasks.

7. Extension: Constitution Refinement and Interpretive Constraints

7.1 Refining Constitutions Dynamically

Works on drafting and evaluating constitutions propose refinement mechanisms:

A graph-based method for selecting principles improves safety while keeping general reasoning performance.

Legal-style frameworks suggest refining vague rules via clarifications or “administrative-style” procedures to reduce inter-model discrepancies.

In our formalism, a constitution refinement map R takes an existing constitution C and a dataset $D_{examples} = \{(x_i, y_i^*)\}_{i=1}^n$ of “safe but useful” outputs, and returns a new constitution:

$$C^* = R(C, D_{examples}).$$

R may, for example:

add new rules that capture patterns in $D_{examples}$,

weaken or clarify vague constraints that overly prune $Y_C(x)$, or project the feasible set $Y_C(x)$ closer to the support of $D_{examples}$.

Under suitable conditions, one can then prove that the refined constitution C^* yields a larger overlap between $Y_{C^*}(x)$ and the human-desirable response manifold induced by $D_{examples}$.

7.2 Interpretive Constraints

Parallel interpretive-constraint proposals limit which “interpretations” of a constitution are legal. One can define a family of admissible rule semantics \mathcal{S} , and associate with each concrete rule C_i a mapping:

$$(\sigma, x, y) \rightarrow c_i^\sigma(x, y) \in \{0, 1\},$$

where $\sigma \in \mathcal{S}$ is an interpretive scheme. Then the constitution sets only the constraint:

$$\forall \sigma \in \mathcal{S}, y \in Y_C(x; \sigma) \equiv (\forall c_i, c_i^\sigma(x, y) = 1)$$

A legal-style framework further restricts \mathcal{S} to a subset $\mathcal{S}_{reasonable}$ of “reasonable” interpretations, analogous to canons of statutory construction.

This extends our model into a two-layer structure:

outer layer: selection / refinement of constitutions and interpretive constraint sets;

inner layer: constrained optimization or sampling of model behavior under a fixed (C, σ) .

8. Discussion and Implications

8.1 Verifiability and Monitoring

Treating the constitution as a mathematically defined constraint set makes violation detection more systematic: one can define monitors that decide membership in $Y_C(x)$ for concrete (x, y) pairs, and flag violations either at training time or in deployment. This aligns with proposals for verifiable constitutions that turn the question “Is AI aligned?” into the more operational question “Has this system violated its constitution?”

8.2 Limits and Open Questions

Our formalization still abstracts away several practical aspects:

Natural-language fuzziness of constitutional rules and their interpretation by LM.

Emergent inconsistencies between rules under distributional shifts not captured by \mathcal{X} in the model.

Future work might:

introduce robust or distributionally robust constitutional constraints,

integrate constitutional structures with context-aware ethical-check frameworks or multi-branch alignment architectures, or

derive sample-complexity bounds on constitutional refinement from datasets of safe examples.

9. Conclusion

We have proposed a mathematical framework for Constitutional AI by representing rule sets as predicates, defining feasible output regions, and expressing alignment as constrained optimization, this framework enables theoretical reasoning about safety guarantees and provides groundwork for further advances in mathematical AI alignment.

Funding

No

Conflict of Interests

The authors declare that there is no conflict of interest regarding the publication of this paper.

Reference

- [1] Vincent, R., Heitzig, J., Jacobs, B. M., Lambert, N., Mossé, M., Pacuit, E., & Russell, S. (2024). Social choice should guide AI alignment in dealing with diverse human feedback. arXiv. <https://arxiv.org/abs/2404.10271v2>
- [2] Kyrychenko, Y., Roozenbeek, J., Davidson, B., van der Linden, S., & Debnath, R. (2025). Human preferences for constructive interactions in language model alignment. arXiv. <https://arxiv.org/abs/2503.16480v1>
- [3] Tennant, E., Jenkins, S. F., Miller, V., Robertson, R., Wen, B., Yun, S.-H., & Taisne, B. (2024). Automating tephra fall

- building damage assessment using deep learning. *Earth System Science Data*, 24, 4585–4608. <https://doi.org/10.5194/nhess-24-4585-2024>
- [4] Vincent, R., Heitzig, J., Lambert, N., Mossé, M., Pacuit, E., & Russell, S. (2024). Social choice should guide AI alignment in dealing with diverse human feedback. *arXiv*. <https://arxiv.org/abs/2404.10271v2>
- [5] Plaat, A., Wong Suzan, V., Broekens, J., van den Broek, N., & Bäck, T. (2025). A multi-step reasoning with large language models, a survey. *arXiv*. <https://arxiv.org/abs/2407.11511v3>
- [6] Findeis, A., Kaufmann, T., Hüllermeier, E., Albanie, S., & Mullins, R. (2024). Inverse Constitutional AI: Compressing preferences into principles. *arXiv preprint*. <https://arxiv.org/abs/2406.06560>
- [7] Conitzer, V., Freedman, R., Heitzig, J., Holliday, W. H., Jacobs, B. M., Lambert, N., Mossé, M., Pacuit, E., Russell, S., Schoelkopf, H., Tewolde, E., & Zwicker, W. S. (2024). Position: Social choice should guide AI alignment in dealing with diverse human feedback. In *Proceedings of the 41st International Conference on Machine Learning* (pp. 9346–9360). *Proceedings of Machine Learning Research*.
- [8] Findeis, A., Kaufmann, T., Hüllermeier, E., Albanie, S., & Mullins, R. (2024). Inverse Constitutional AI: Compressing preferences into principles (ICLR 2025 conference paper preprint). *arXiv*. <https://arxiv.org/abs/2406.06560>.