

The Alignment Game: A Theory of Long-Horizon Alignment Through Recursive Curation

Ali Falahati¹, Mohammad Mohammadi Amiri², Kate Larson¹, Lukasz Golab¹

¹ University of Waterloo

² Rensselaer Polytechnic Institute
afalahati@uwaterloo.ca

Abstract

In self-consuming generative models that train on their own outputs, alignment with user preferences becomes a recursive rather than one-time process. We provide the first formal foundation for analyzing the long-term effects of such recursive retraining on alignment. Under a two-stage curation mechanism based on the Bradley–Terry (BT) model, we model alignment as an interaction between two factions: the *Model Owner*, who filters which outputs should be learned by the model, and the *Public User*, who determines which outputs are ultimately shared and retained through interactions with the model. Our analysis reveals three structural convergence regimes depending on the degree of preference alignment: consensus collapse, compromise on shared optima, and asymmetric refinement. We prove a fundamental impossibility theorem: no recursive BT-based curation mechanism can simultaneously preserve diversity, ensure symmetric influence, and eliminate dependence on initialization. Framing the process as dynamic social choice, we show that alignment is not a static goal but an evolving equilibrium, shaped both by power asymmetries and path dependence.

Extended version —

<https://aaai.org/example/extended-version>

Introduction

Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al. 2022) has become the de facto method for aligning large language models with human preferences. Its appeal lies in a simple loop: broadcast model outputs to annotators, collect pairwise preferences, and update the policy through reward modeling. Among the many instantiations of this loop, the *Bradley-Terry (BT) comparison model* (Bradley and Terry 1952) is widely used, powering early alignment successes such as InstructGPT and many subsequent fine-tuning pipelines. At the same time, BT has been criticized for strong independence assumptions of annotating samples, a tendency to reward extremely probable rather than diverse outputs, and vulnerability to noisy feedback (Ge et al. 2024; Xiao et al. 2024).

Most analyses of BT ask whether a single round of collecting human preferences yields an aligned model. In prac-

tice, however, modern generative systems are updated recursively: synthetic outputs are added to the corpus, new models are trained on these added data, and the cycle repeats across model generations. Recent work shows that this self-consuming regime can drift away from human values or collapse onto degenerate equilibria (Ferbach et al. 2024; Shumailov et al. 2023; Gerstgrasser et al. 2024; Alemohammad et al. 2023). A long-horizon view demands a formalism that tracks how preference aggregation compounds across successive rounds.

To fill this gap, we study the long-term alignment dynamics under recursive retraining. We focus on the simplest yet already rich scenario with two factions: the **Model Owner**, representing the incentives of developers or platform providers, and the **Public**, representing the aggregated preferences of society interacting with the model. At each iteration, the Owner curates samples via a mechanism based on BT, the model is trained on the selected outputs, the Public interacts with the released model and preserves a subset of model outputs through various actions such as upvoting, sharing, and other forms of engagement, and these data flow into the next training set. Our central question is: *if this BT curation loop continues indefinitely, what distribution of content survives in the limit?*

We model this loop as a dynamic social choice mechanism, which allows us to reason about fairness, incentive compatibility, and power asymmetries in a multi-agent setting. Our analysis is structured around three core alignment scenarios, each motivated by real-world tensions. First, we explore Perfect Alignment, where the Owner and Public preferences are identical. This represents the idealized “fully-aligned” goal that many systems implicitly strive for, and analyzing it is critical to understanding the consequences of perfect agreement. Second, we examine Partial Alignment, the most realistic scenario where preferences overlap but do not perfectly coincide: for instance, a platform’s commercial goals overlapping with, but not fully capturing, public interests. This case allows us to study the dynamics of compromise. Finally, we analyze Disjoint Alignment, which models a critical conflict where developer and user values are in direct opposition. This scenario reveals the system’s inherent power dynamics and who ultimately controls the output when values diverge. By understanding the long-run trajectory of each regime, we can expose the

structural trade-offs between diversity, fairness, and stability, establishing a foundation for alternative mechanisms and richer multi-agent systems.

In summary, our contributions are threefold: (i) a formal model of recursive BT-based RLHF involving a model owner and a public user; (ii) a complete characterization of its long-run behavior, including an impossibility result that generalizes classic social choice tensions to the dynamic setting; and (iii) empirical evidence that these theoretical dynamics surface in practice. This work aims to lay the groundwork for a broader theory of long-term alignment, extending beyond single-step fine-tuning to account for the recursive and evolutionary nature of generative models.

Related Work

Generative AI systems face fundamental challenges arising from recursive training dynamics and the cumulative effects of self-generated data. Shumailov et al. (2023) and Ferbach et al. (2024) demonstrated that training models on their own generated content leads to irreversible “model collapse,” where output diversity diminishes through recursive iterations. These limitations are compounded by recursive training dynamics, as Dohmatob (2024) demonstrates clear crossover points between stable and collapse regimes (Xu, He, and Cheng 2025). Gerstgrasser et al. (2024) showed that data accumulation strategies can prevent collapse while replacement accelerates it.

The multi-stakeholder dimension of RLHF has been explored by Tewolde et al. (2024), who propose Reinforcement Learning from Collective Human Feedback using social choice theory, while Mishra (2023) proves that universal AI alignment using RLHF is impossible under democratic constraints. The BT model, widely used in preference-based reward modeling, faces significant limitations in multi-stakeholder scenarios. Sun (2025) provides the first theoretical critique, arguing these models are unnecessary for downstream optimization and proposing classification-based alternatives focused on “order consistency.” Zhang et al. (2024) demonstrate that BT models “fall short in expressiveness, particularly in addressing intransitive preferences” that arise naturally in multi-stakeholder scenarios.

Wu (2022) shows that these models violate the independence assumptions when multiple stakeholders are involved, leading to systematic marginalization of the preferences of minorities. Eckersley (2019) demonstrates that Arrow’s impossibility theorem applies to AI alignment, showing that no satisfactory method exists to aggregate multiple human preferences without violating fairness criteria. Qiu (2024) extends this with Arrow-like impossibility theorems for representative social choice settings. The convergence of theoretical impossibility results with empirical evidence of model collapse establishes the foundation for understanding AI alignment as an inherently conflicted game between competing stakeholder interests. Additional insights from social choice theory highlight its role in guiding AI alignment, particularly in handling diverse human feedback and avoiding preference aggregation pitfalls (Conitzer et al. 2024).

Problem Definition

We study iterative retraining of generative models, where synthetic outputs recursively influence future training data through a two-agent curation process. This framework captures the alignment challenges that emerge when two stakeholders jointly shape the evolution of AI systems.

Model Components

Consider a generative modeling system with the following components:

- **State Space:** A compact metric space (\mathcal{X}, d) representing the content domain
- **Agents:** Two curators, the *Owner* (model developer) and the *Public* (user community), with continuous reward functions $r_O, r_P : \mathcal{X} \rightarrow \mathbb{R}$ encoding their respective preferences
- **Data Evolution:** A sequence of public datasets $\{\mathcal{D}_t\}_{t \geq 1}$ where $\mathcal{D}_t \subset \mathcal{X}$
- **Model Sequence:** Generative models $\{\mathcal{M}_t\}_{t \geq 1}$ with output distributions $p_t \in \mathcal{P}(\mathcal{X})$

Definition 1. For a probability measure p on \mathcal{X} , a pool size $K \geq 2$, and a reward function r , define the BT weight as (Ferbach et al. 2024):

$$H_{K,r}^p(x) := \mathbb{E}_{Y_1, \dots, Y_{K-1} \sim p} \left[\frac{K e^{r(x)}}{e^{r(x)} + \sum_{j=1}^{K-1} e^{r(Y_j)}} \right].$$

Alignment Game Framework

The system recursively evolves as follows. Start with initial dataset \mathcal{D}_1 and initial model distribution p_0 .

Recursive loop (for $t = 1, 2, \dots$):

- (1) **Owner Curation:** At iteration t , the Owner samples a pool $\{x_1, \dots, x_K\} \sim p_t$ (the pool is a subset of \mathcal{D}_t , i.e., $\{x_i\} \sim q_t$ with $q_t \approx p_t$) and selects outputs via BT selection with reward r_O , yielding:

$$\tilde{p}_t(x) = p_t(x) \cdot H_{K,r_O}^{p_t}(x) \quad (1)$$

- (2) **Model Update:** Train \mathcal{M}_{t+1} on data drawn from \tilde{p}_t , producing the updated model distribution:

$$p_{t+1}(x) \approx \tilde{p}_t(x) \quad (2)$$

- (3) **Public Curation:** The Public samples a pool $\{x_1, \dots, x_M\} \sim p_{t+1}$ and applies BT selection with reward r_P , yielding:

$$\hat{p}_t(x) = p_{t+1}(x) \cdot H_{M,r_P}^{p_{t+1}}(x) \quad (3)$$

- (4) **Dataset Evolution:** Update $\mathcal{D}_{t+1} = \mathcal{D}_t \cup \mathcal{O}_t^*$ where $\mathcal{O}_t^* \sim \hat{p}_t$.

Our analysis relies on the following assumptions:

1. The recursive curation mechanism is explicitly defined by the BT model. The convergence properties and the impossibility theorem are direct consequences of the mathematical properties of this specific pairwise comparison model.
2. We analyze the idealized dynamic $p_{t+1}(x) = \tilde{p}_t(x)$, where the model distribution at the next step perfectly matches the target distribution from the current step's curation. This abstracts away optimization error, noise, or catastrophic forgetting from practical training.
3. The reward functions $r_O : \mathcal{X} \rightarrow \mathbb{R}$ and $r_P : \mathcal{X} \rightarrow \mathbb{R}$ are assumed to be fixed and continuous over a compact metric space \mathcal{X} . This ensures the optimal sets A_O and A_P are well-defined and non-empty.

Throughout the analysis, we consider open neighborhoods around maximizer sets. For any set $A \subset \mathcal{X}$ and radius $\eta > 0$,

$$B_\eta(A) := \{x \in \mathcal{X} : \inf_{y \in A} d(x, y) < \eta\}. \quad (4)$$

Alignment Regimes

The interaction between curator preferences determines the system's long-term behavior. Let $A_O = \arg \max_{x \in \mathcal{X}} r_O(x)$ and $A_P = \arg \max_{x \in \mathcal{X}} r_P(x)$ denote the optimal sets for each curator.

Definition 2. We categorize value alignment between the Owner and the Public as:

- **Perfect Alignment:** $A_O = A_P$ (complete agreement)
- **Partial Alignment:** $A_O \cap A_P \neq \emptyset$ with $A_O \neq A_P$ (overlapping preferences)
- **Disjoint Alignment:** $A_O \cap A_P = \emptyset$ (conflicting preferences)

Core Recursive Alignment Challenges

As $t \rightarrow \infty$, the public dataset becomes dominated by curated synthetic data:

$$\lim_{t \rightarrow \infty} \frac{|\mathcal{D}_1|}{|\mathcal{D}_t|} = 0 \quad (5)$$

This *self-consuming* regime, where models train predominantly on their own filtered outputs, is a critical challenge for maintaining alignment with diverse human values while avoiding mode collapse or value lock-in. This framework enables us to investigate fundamental questions about recursive alignment:

- **Convergence:** Under what conditions does the system converge to a stable distribution? Does it collapse to point masses or preserve diversity?
- **Influence Dynamics:** How do the sequential curation mechanisms affect the relative influence of each curator?
- **Alignment Impact:** How does the degree of preference alignment between curators shape the evolution and final state of the generative model?

Theoretical Results ¹

We now analyze the long-term behavior of our two-stage curation mechanism under varying degrees of alignment between curators. We begin with the idealized case where both the Owner and the Public share identical preferences.

Perfect Alignment: The Consensus Trap

Theorem 1. *When the Owner and Public have perfectly aligned preferences: $A_O = A_P =: A_* \neq \emptyset$. Let $(p_t)_{t \geq 0}$ be the sequence of distributions generated by this mechanism. Then the system reaches complete consensus: for any $\eta > 0$, there exist constants $C, c > 0$ such that $p_t(\mathcal{X} \setminus B_\eta(A_*)) \leq Ce^{-ct}$ for all $t \geq 0$. Moreover, p_t converges weakly to the renormalized initial distribution on A_* :*

$$p_\infty(x) = \frac{p_0(x)}{\int_{A_*} p_0(z) dz} \mathbf{1}_{A_*}(x).$$

Our first result uncovers that greater alignment does not preserve diversity, but speeds up its collapse. When preferences are perfectly aligned, the model concentrates on the shared maximizers, leading to a sharply reduced outcome space. Perfect alignment produces an impoverished limiting distribution. The exponential decay rate e^{-ct} reflects the speed of this convergence, highlighting how quickly diversity is lost under alignment. In the limit, the only remaining variation arises from the initial distribution p_0 restricted to A_* .

Corollary 2 (Mode Collapse for Unique Maximizers). *When both curators agree on a unique optimal point $A_O = A_P = \{x^*\}$, system undergoes mode collapse: $p_t \Rightarrow \delta_{x^*}$.*

This corollary represents the most extreme form of homogenization and highlights a fundamental risk in alignment-driven curation: as agreement increases, the support of the distribution contracts. Thus, systems designed to maximize agreement between developers and users may unintentionally collapse into echo chambers. This consensus trap suggests that some degree of preference misalignment may be necessary to sustain a healthy generative ecosystem. Our result complements prior findings such as (Ferbach et al. 2024), which show that recursive training on curated outputs can lead to long-term degeneracy. Thus, before pursuing perfect agreement, we should ask whether we are engineering an echo chamber that eliminates diversity.

Partial Alignment: The Compromise Equilibrium

Next, we consider the case where curators share some common ground while maintaining distinct preferences, a setting that better reflects the relationship between model owners and diverse user communities.

Theorem 3. *Suppose the curators have overlapping but distinct preferences, with shared optima $A_{\text{shared}} := A_O \cap A_P \neq \emptyset$ while $A_O \neq A_P$. Under the two-stage curation mechanism, only the intersection survives: the mass outside $B_\eta(A_{\text{shared}})$ decays exponentially as $p_t(\mathcal{X} \setminus B_\eta(A_{\text{shared}})) \leq$*

¹Proofs in Appendix Section A. of the extended version.

Ce^{-ct} . The limiting distribution p_∞ concentrates entirely on A_{shared} with density proportional to the initial distribution:

$$p_\infty(x) = \frac{p_0(x)}{\int_{A_{\text{shared}}} p_0(z) dz} \mathbf{1}_{A_{\text{shared}}}(x).$$

Our findings reveal a key dynamic of partial alignment: the recursive process filters out content valued exclusively by one curator, preserving only what lies in the intersection. This “lowest common denominator” effect eliminates curator-specific signals, including potentially novel or specialized contributions. Although diversity is maintained within A_{shared} , it remains limited to areas of mutual agreement. For AI systems, this suggests that iterative training with multiple stakeholders may progressively reduce the breadth of the model to only the jointly endorsed features.

Disjoint Alignment: Asymmetric Power Dynamics

The most adversarial situation arises when curators hold entirely disjoint preferences. In such cases, the central question becomes: who ultimately gets to define the output space of the generative model when both sides embody incompatible values?

Theorem 4. *When curator preferences are completely disjoint ($A_O \cap A_P = \emptyset$), the Owner determines the support while the Public refines within it. Define the Public’s preferred subset within Owner optima as $A_{P|O} := \arg \max_{x \in A_O} r_P(x)$. The system exhibits two-stage exponential suppression. First, content outside the Owner’s optima vanishes with $p_t(\mathcal{X} \setminus B_\eta(A_O)) \leq C_1 e^{-c_1 t}$, then within A_O , mass concentrates on the Public’s preferred subset with $p_t(A_O \setminus B_\eta(A_{P|O})) \leq C_2 e^{-c_2 t}$. The limiting distribution is*

$$p_\infty(x) = \frac{p_0(x)}{\int_{A_{P|O}} p_0(z) dz} \mathbf{1}_{A_{P|O}}(x).$$

This result exposes the inherent power asymmetry in the recursive alignment game. The Owner’s first-mover advantage allows it to determine the feasible region, while the Public can only optimize within these constraints. This “best of the worst” dynamic mirrors real-world scenarios where users must select from options pre-filtered by platform algorithms. The Owner’s preferences shape the system quickly, while the Public’s influence manifests gradually as a refinement. This temporal divide reveals a deeper governance dilemma: once early alignment choices solidify into norms, user feedback becomes more about adaptation than agency.

The Fundamental Impossibility Result

After exploring the range of alignment scenarios, we turn to a deeper limitation of recursive BT-based curation. Even under ideal conditions, certain properties that appear simultaneously desirable cannot all be achieved at once.

Theorem 5. *For any non-trivial preference misalignment ($A_O \neq A_P$), it is impossible to simultaneously satisfy:*

(1) *Full Coverage:*

$$\liminf_{t \rightarrow \infty} p_t(A_O \setminus A_P) > 0 \quad \text{and} \quad \liminf_{t \rightarrow \infty} p_t(A_P \setminus A_O) > 0$$

(2) *Symmetric Influence: There exists a permutation-invariant functional Φ such that*

$$p_\infty = \Phi(r_O, r_P, p_0) = \Phi(r_P, r_O, p_0)$$

(3) *Initial Distribution Independence: For any two initial distributions $p_0, q_0 \in \mathcal{P}(\mathcal{X})$ with $\text{supp}(p_0) = \text{supp}(q_0) = \mathcal{X}$,*

$$p_\infty^{(p_0)} = p_\infty^{(q_0)}$$

This impossibility result states that no recursive BT-based curation mechanism can simultaneously maintain content diversity across disagreement regions (Full Coverage), treat both curators equally (Symmetric Influence), and produce outcomes independent of initial conditions (Initial Distribution Independence). Each of these properties captures a fundamental value tension in alignment systems. Full Coverage represents the epistemic goal of preserving the full spectrum of ideas, ensuring that minority or dissenting content is not prematurely filtered out. Symmetric Influence encodes fairness: both curators should exert comparable control over the generative process, preventing dominance by either side. Initial Distribution Independence embodies stability, guaranteeing that the system’s long-term behavior does not hinge on arbitrary starting points or early biases.

The theorem’s implications extend beyond our specific model: it formalizes the intuition that these goals pull in incompatible directions. The tradeoffs observed in recursive BT-based curation systems are therefore not artifacts of design but mathematical necessities. Designers must decide which virtue to compromise: accepting homogenization, embracing asymmetry, or tolerating historical lock-in.

Remark 1. *Across all alignment regimes, the support of the limiting distribution shrinks as follows:*

- $\text{supp}(p_\infty) \subseteq A_O \cap A_P$ when the intersection is non-empty
- $\text{supp}(p_\infty) \subseteq A_{P|O} \subseteq A_O$ when $A_O \cap A_P = \emptyset$

This demonstrates that recursive BT-based curation inevitably reduces diversity, concentrating mass on increasingly narrow regions of agreement. These structural constraints raise a deeper question: what kind of decision-making process is recursive curation, and how should we reason about its limitations? To answer this, we now reinterpret the recursive BT-based curation mechanism using social choice theory.

The Curation Mechanism as Social Choice

The recursive BT-based curation process can be viewed as a form of collective decision-making, where two agents jointly shape the long-term distribution of model outputs. To examine its properties—fairness, efficiency, and the preservation of diversity—we draw from social choice theory and mechanism design. We formalize recursive BT-based curation as a dynamic preference aggregation mechanism in which the Owner and the Public iteratively express preferences over a shared alternative space \mathcal{X} through a two-stage influence process analogous to voting.

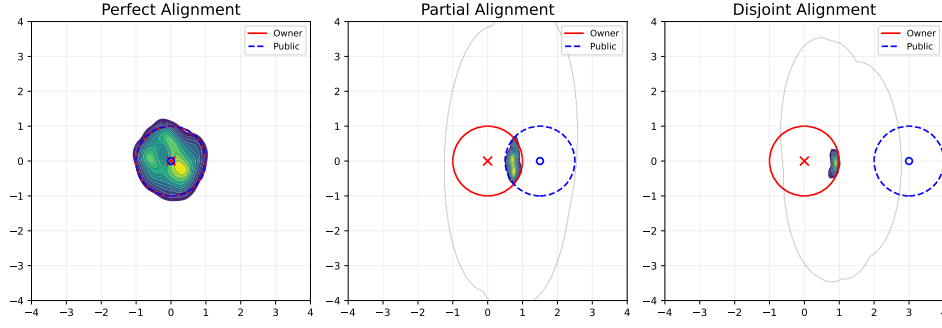


Figure 1: KDE plots showing point distributions in three alignment scenarios: perfect alignment (same circles), partial alignment (overlapping circles), and disjoint alignment (non-overlapping circles). Red circles indicate owner preferred regions, blue dashed circles show public preferred regions.

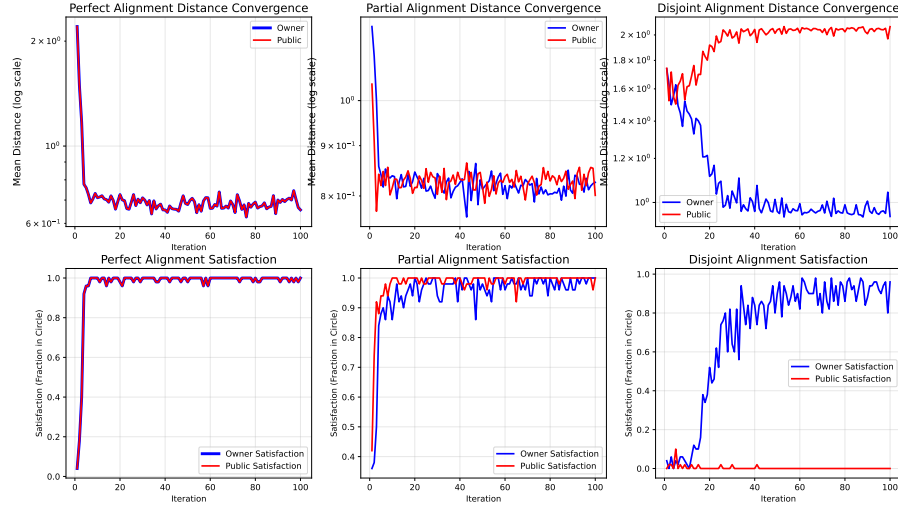


Figure 2: (Top) Convergence of the mean distance to the owner and public centers over iterations for three alignment scenarios: perfect alignment, partial alignment, and disjoint alignment. (Bottom) The fraction of points within the owner and public preferred regions ("satisfaction") as the alignment process progresses.

Definition 3. A dynamic social choice mechanism for the alignment game consists of:

- (1) A set of alternatives \mathcal{X} (the content space)
- (2) Two agents $\mathcal{N} = \{O, P\}$ (Owner and Public)
- (3) Preferences, captured by reward functions $r_i : \mathcal{X} \rightarrow \mathbb{R}$ for $i \in \mathcal{N}$
- (4) A social choice correspondence $\mathcal{F} : \mathcal{R}^2 \times \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{X})$ that maps preference profiles and current distributions to updated distributions
- (5) A limit social choice function $f_\infty : \mathcal{R}^2 \rightarrow \mathcal{P}(\mathcal{X})$ representing the long-run outcome

The mechanism implements a sequential decision process in which agents express preferences via pairwise (binary) comparisons over alternatives. The Owner selects an alternative $x \in \mathcal{X}$ first, based on their reward function, and the Public responds by evaluating this choice using their own reward function. The final distribution over alternatives is updated

based on these sequential comparisons, reflecting the asymmetric influence of each agent.

Strategic Voting and Incentive Compatibility

In the social choice literature, a central property is strategyproofness: agents should not have an incentive to misrepresent their preferences.

Theorem 6 (Weak Strategyproofness of Reports). *For any reported pair (r'_O, r'_P) let $p_\infty(r'_O, r'_P)$ denote the weak limit of (p_t) , and define each agent's utility by*

$$U_O(r'_O, r'_P) := \mathbb{E}_{x \sim p_\infty(r'_O, r'_P)}[r_O(x)],$$

$$U_P(r'_O, r'_P) := \mathbb{E}_{x \sim p_\infty(r'_O, r'_P)}[r_P(x)].$$

Then truthful reporting is weakly dominant for both agents:

$$U_O(r_O, r'_P) \geq U_O(\hat{r}_O, r'_P), \quad U_P(r'_O, r_P) \geq U_P(r'_O, \hat{r}_P)$$

for all (\hat{r}_O, \hat{r}_P) , and for all (r'_O, r'_P) . In particular, $(r'_O, r'_P) = (r_O, r_P)$ is a Nash equilibrium.

The theorem ensures truthful elicitation: for both agents, truthful reporting weakly dominates any misreport, regardless of the other report. Manipulating reports is therefore ineffective; influence instead comes from the dynamics. The result addresses incentives over reports only, and the broader questions of fairness, symmetric influence, and coverage are taken up in the next theorem, which formalizes their mutual incompatibility.

Fairness and Representation in Dynamic Voting

Beyond strategyproofness, social choice theory provides principled ways of thinking about fairness criteria, such as equal treatment of agents. We formalize this notion through the concept of *influence parity*.

Theorem 7. *No asymmetric sequential preference aggregation mechanism can simultaneously satisfy:*

(EQ) *Influence parity: the outcome is order symmetric, that is*

$$p_{\infty}^{\text{OF}}(r_O, r_P) = p_{\infty}^{\text{PF}}(r_O, r_P).$$

(PO) *Pareto optimality: there is no distribution q on X with*

$$\mathbb{E}_q[r_O] \geq \mathbb{E}_{p_{\infty}}[r_O], \quad \mathbb{E}_q[r_P] \geq \mathbb{E}_{p_{\infty}}[r_P],$$

and at least one inequality strict.

(UC) *Uniqueness: for each (r_O, r_P) the outcome is unique, independent of initialization and of any tie breaking in the dynamics.*

The sequential nature of voting creates a first-mover advantage that conflicts with equal treatment. Unlike static voting where ties can be broken symmetrically, the recursive dynamics amplify initial asymmetries, making true voter equality impossible while maintaining uniqueness.

Experiments

We validate our analysis through two experimental frameworks: a synthetic alignment game that directly implements our mathematical framework, and a text-based alignment game based on a realistic language modeling setting.

Synthetic Alignment Game

Experimental Setup. We implement a synthetic alignment game in \mathbb{R}^2 where the Owner and Public preferences are defined by circular reward regions. The reward function for agent $i \in \{O, P\}$ with center c_i and radius r is:

$$r_i(x) = \begin{cases} 1.0 & \text{if } \|x - c_i\| \leq r \\ -2.0 \cdot (\|x - c_i\| - r) & \text{if } \|x - c_i\| > r \end{cases} \quad (6)$$

This creates a sharp preference boundary where points inside the circle receive maximum reward and points outside receive penalties proportional to their distance.

Alignment Scenarios. (i) *Perfect alignment:* $c_O = c_P = (0, 0)$, $r = 1.0$, i.e., both agents prefer the same circular region; (ii) *Partial alignment:* $c_O = (0, 0)$, $c_P = (1.5, 0)$, $r = 1.0$, i.e., overlapping circles with shared optima; (iii) *Disjoint alignment:* $c_O = (0, 0)$, $c_P = (3, 0)$, $r = 1.0$, i.e., nonoverlapping circles with no shared optima.

Experimental Parameters. Each experiment runs for 100 iterations with the following parameters: (i) Initial dataset:

1000 points sampled uniformly from $[-5, 5] \times [-5, 5]$; (ii) Owner curation: Select 100 points using the BT mechanism with temperature $\tau = 0.5$; (iii) Generation: Train a Gaussian Mixture Model (GMM) (Bishop 2006) on curated data, generate 200 new samples; (iv) Public curation: Select 50 points from generated samples using the BT mechanism; (v) Dataset update: Add public-curated samples to the training.

Results and Analysis. Figure (1) (top) shows the final distribution of points after 100 iterations for each alignment scenario. The KDE plots reveal convergence patterns that validate our theoretical predictions. *Perfect Alignment:* The system converges to a concentrated distribution within the shared optimal region, with exponential suppression outside the circle. The final distribution shows high density within the preferred region and negligible mass elsewhere, confirming Theorem (1). *Partial Alignment:* The system concentrates on the intersection of the two circles, preserving diversity only where preferences align. Points outside the intersection are exponentially suppressed, while the shared region maintains substantial density. This validates Theorem (3). *Disjoint Alignment:* The system converges to the Public’s preferred subset within the Owner’s optimal region. The Owner’s first-mover advantage determines the support, while the Public refines the distribution within that support. This confirms Theorem (4). Figure (2) (top) tracks the mean distance to each agent’s preferred center over iterations. All scenarios show exponential convergence, with the rate depending on the degree of alignment. Perfect alignment achieves the fastest convergence, while disjoint alignment shows a two-stage process, first converging to the Owner’s region, then refining within it. Figure (2) (bottom) shows the fraction of points satisfying each agent’s preferences (“satisfaction rate”). In perfect alignment, both agents achieve near-universal satisfaction. In partial alignment, satisfaction is limited to the intersection region. In disjoint alignment, the Owner maintains high satisfaction, while the Public’s satisfaction is constrained by the Owner’s preferences.

Text-Based Alignment Game

Experimental Setup. We implement a realistic text generation scenario using GPT-2 models (Radford et al. 2019) where the Owner and Public have different preferences for response length. The system operates on the WikiText-2 (Merity et al. 2016) dataset, with agents preferring different word count ranges for generated responses.

Alignment Scenarios. We test three text-based alignment scenarios: (i) *Perfect Alignment:* Both agents prefer exactly 4 words (Owner: 4-4, Public: 4-4); (ii) *Partial Alignment:* Owner prefers 2-4 words, Public prefers 4-6 words (overlapping range); *Disjoint Alignment:* Owner prefers 1-3 words, Public prefers 5-6 words (non-overlapping ranges).

Experimental Parameters. (i) Initial dataset: 10000 sentences from WikiText-2 filtered by word count; (ii) Owner curation: Select 1000 sentences using the BT mechanism (iii) Generation: Fine-tune the model for 2 epochs, generate 2000 responses with temperature 0.8; (iv) Public curation: Select 1000 responses using the BT mechanism (v) Training: Learning rate 5×10^{-5} , batch size 8.

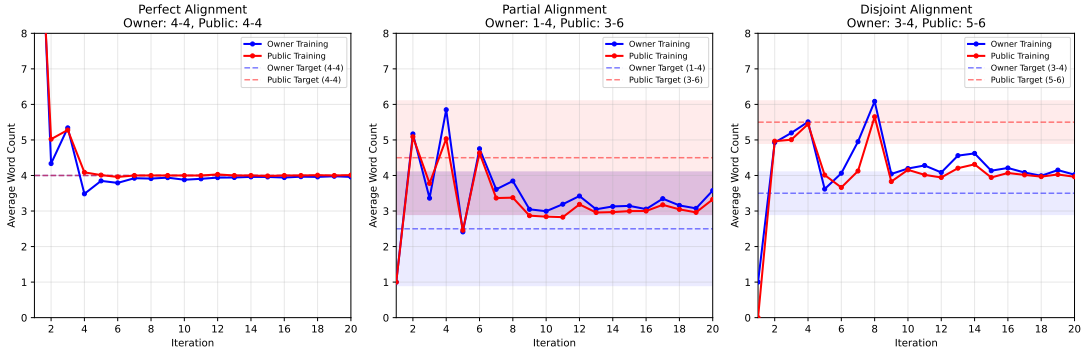


Figure 3: Word count evolution across three alignment scenarios: perfect alignment (agents target 4 words), partial alignment (owner targets 1-3 words, public targets 3-5 words), and disjoint alignment (owner targets 3-4 words, public targets 5-6 words).

Results and Analysis. Figure (3) shows the evolution of average word count across iterations for all three scenarios. The results demonstrate the same convergence patterns as the synthetic experiments. *Perfect Alignment:* Both training stages converge to the shared target of 4 words within 10 iterations, achieving near-perfect alignment with minimal variance. *Partial Alignment:* The system converges to the intersection of preferences (around 3-4 words), with the Owner’s training showing slightly lower word counts and the Public’s training showing slightly higher counts, but both remaining within the overlapping range. *Disjoint Alignment:* The system shows a two-stage convergence process. Initially, both training stages converge toward the Owner’s preferred range (3-4 words). However, the Public’s training gradually shifts toward longer responses, eventually settling around 4 words, demonstrating the Public’s ability to refine within the Owner’s preferred domain.

Key Findings. Our experiments validate several theoretical predictions. All scenarios show exponential convergence to equilibrium distributions, with convergence rates inversely proportional to the degree of alignment. In all cases, the final distribution concentrates on a subset of the original support, confirming the shrinking support principle (Remark (1)). In disjoint alignment scenarios, the Owner’s preferences dominate the support selection, while the Public refines within that support. Partial alignment scenarios preserve diversity only within the intersection of preferred regions, with exponential suppression elsewhere. These phenomena persist across different domains (continuous spaces vs. discrete text) and different model architectures (GMM vs. transformer).

Conclusion

Alignment is not a one-time setting. In recursive curation, even well-intended preference aggregation can reduce diversity, confer first mover power, and lock models into narrow equilibria. We identified three convergence regimes, consensus collapse, lowest common denominator compromise, and owner-led refinement, showing how small preference gaps can redirect the long-run trajectory of generative systems.

Implications for AI alignment. Our findings suggest that we must reconceptualize alignment itself. First, it moves

the problem from one-time preference-matching to dynamic mechanism design. The paper demonstrates that the structure of the alignment process, the sequential curation, the power dynamics, is a value-driven system that actively shapes outcomes, privileging owner-defined constraints over public refinement. Second, the impossibility theorem is not a failure but a clarification. It suggests that the goal of a single, stable, and diverse alignment may be a contradiction in terms. This should force the field to pivot from designing for consensus to designing for pluralism. The challenge is not to eliminate disagreement but to build systems that can productively manage it, treating “lowest common denominator” outcomes as a failure state of the mechanism, not an acceptable compromise. Finally, this work implies that the true alignment “meta-problem” is not just aligning a model, but aligning the recursive alignment process itself. We must design the feedback loop, the “game”, to be transparent, contestable, and fair. Instead of a one-time setting, alignment becomes a continuous process of governing the system that governs the models.

Limitations and Future Work. This analysis is based on a model with key simplifying assumptions. First, the framework rests on the BT model, which assumes that preferences are independent and transitive. Future work must extend this analysis to more complex preference models. Second, the two-agent “Owner-Public” game is a simplification of a complex ecosystem. The next step is to model this as an n-agent game with heterogeneous agents. Third, the model assumes static preferences. A more realistic framework would treat preferences and model outputs as co-evolving. Future research should explore this co-evolutionary dynamic, where the outputs of one generation’s model actively shape the preferences of the agents who curate the next. Finally, this paper is descriptive; it explains what will happen. The next step is normative mechanism design. Given the impossibility theorem, how can we design new curation systems that explicitly and transparently choose which property to sacrifice? This moves the alignment challenge beyond the realm of optimization and into the domain of political and social governance. The task is no longer to discover a single, mathematically “correct” alignment, but to engineer a

process that is perceived as legitimate, transparent, and fair by all stakeholders.

References

- Alemohammad, S.; Casco-Rodriguez, J.; Luzzi, L.; Humayun, A. I.; Babaei, H.; LeJeune, D.; Siahkoobi, A.; and Baraniuk, R. 2023. Self-consuming generative models go mad. In *The Twelfth International Conference on Learning Representations*.
- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. Springer.
- Bradley, R. A.; and Terry, M. E. 1952. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika*, 39(3/4): 324–345.
- Conitzer, V.; Freedman, R.; Heitzig, J.; Holliday, W. H.; Jacobs, B. M.; Lambert, N.; Mossé, M.; Pacuit, E.; Russell, S.; Schoelkopf, H.; et al. 2024. Social choice should guide ai alignment in dealing with diverse human feedback. *arXiv preprint arXiv:2404.10271*.
- Dohmatob, e. a. 2024. Model Collapse Demystified: The Case of Regression. *arXiv preprint arXiv:2402.07712*.
- Eckersley, P. 2019. Impossibility and Uncertainty Theorems in AI Value Alignment (or why your AGI should not have a utility function). *arXiv preprint arXiv:1901.00064*.
- Ferbach, D.; Bertrand, Q.; Bose, A. J.; and Gidel, G. 2024. Self-Consuming Generative Models with Curated Data Provably Optimize Human Preferences. *arXiv preprint arXiv:2407.09499*.
- Ge, L.; Halpern, D.; Micha, E.; Procaccia, A. D.; Shapira, I.; Vorobeychik, Y.; and Wu, J. 2024. Axioms for ai alignment from human feedback. *Advances in Neural Information Processing Systems*, 37: 80439–80465.
- Gerstgrasser, M.; Schaeffer, R.; Dey, A.; Rafailov, R.; Sleight, H.; Hughes, J.; Korbak, T.; Agrawal, R.; Pai, D.; Gromov, A.; Roberts, D. A.; Yang, D.; Donoho, D. L.; and Koyejo, S. 2024. Is Model Collapse Inevitable? Breaking the Curse of Recursion by Accumulating Real and Synthetic Data. *arXiv preprint arXiv:2404.01413*.
- Merity, S.; Xiong, C.; Bradbury, J.; and Socher, R. 2016. Pointer Sentinel Mixture Models. *arXiv:1609.07843*.
- Mishra, A. 2023. AI Alignment and Social Choice: Fundamental Limitations and Policy Implications. *arXiv preprint arXiv:2310.16048*.
- Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744.
- Qiu, T. 2024. Representative Social Choice: From Learning Theory to AI Alignment. *arXiv preprint arXiv:2410.23953*.
- Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; and Sutskever, I. 2019. Language Models are Unsupervised Multitask Learners. *OpenAI*.
- Shumailov, I.; Shumaylov, Z.; Zhao, Y.; Gal, Y.; Papernot, N.; and Anderson, R. 2023. The Curse of Recursion: Training on Generated Data Makes Models Forget. *arXiv preprint arXiv:2305.17493*.
- Sun, e. a. 2025. Rethinking Bradley-Terry Models in Preference-Based Reward Modeling: Foundations, Theory, and Alternatives. *arXiv preprint arXiv:2411.04991*.
- Tewolde, E.; Conitzer, V.; Freedman, R.; Heitzig, J.; Holliday, W. H.; Jacobs, B. M.; Lambert, N.; Mossé, M.; Pacuit, E.; Russell, S.; Schoelkopf, H.; and Zwicker, W. S. 2024. Social Choice Should Guide AI Alignment in Dealing with Diverse Human Feedback. *arXiv preprint arXiv:2404.10271*.
- Wu, e. a. 2022. A Diagnostic Framework for the Bradley-Terry Model. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 185(Supplement.2): S461–S482.
- Xiao, J.; Li, Z.; Xie, X.; Getzen, E.; Fang, C.; Long, Q.; and Su, W. J. 2024. On the algorithmic bias of aligning large language models with RLHF: Preference collapse and matching regularization. *arXiv preprint arXiv:2405.16455*.
- Xu, S.; He, H.; and Cheng, G. 2025. A Probabilistic Perspective on Model Collapse. *arXiv preprint arXiv:2505.13947*.
- Zhang, Y.; Li, B.; Liu, S.; Zhou, Z.; and Shi, Z. 2024. Beyond Bradley-Terry Models: A General Preference Model for Language Model Alignment. *arXiv preprint arXiv:2410.02197*.