
Position: Alignment Needs Rule-Class Routing Before Preference Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This position paper argues that alignment pipelines should classify rules before
2 they aggregate preferences. RLHF [15], Constitutional AI [3], DPO [16], RLAI
3 [13], and Deliberative Alignment [14] differ technically, but each tends to turn
4 contested normative input into one global policy. Social-choice theory explains why
5 this is unsafe on unrestricted domains: aggregation can help when there is public
6 convergence, but it cannot decide which questions are eligible for aggregation.
7 Conitzer et al. [4] reframe alignment as social choice; we agree, and argue social
8 choice identifies a boundary, not a solution. Arrow, Gibbard–Satterthwaite, and Sen
9 establish that aggregation over unrestricted value disagreement has no procedure
10 satisfying minimal democratic, strategic, and liberty-preserving conditions. The
11 eligibility decision is institutional: the system designer must state who is authorized
12 to decide whether a rule is aggregable, user-configurable, or non-negotiable. We
13 propose a routing layer for behavioral rules. Class I rules are public prohibitions
14 suitable for aggregation; Class II rules concern reasonable disagreement and should
15 support configurable defaults; Class III rules protect rights or vulnerable users and
16 should be implemented as constraints with appeal mechanisms. The paper connects
17 impossibility results to alignment design, compares current methods through this
18 lens, and gives a six-step routing procedure with a five-item disclosure checklist.

19 1 Introduction

20 In December 2022, Bai et al. [3] introduced Constitutional AI, a training method using a short list of
21 principles, called a constitution, against which the model critiques and revises its own outputs. The
22 terminology is deliberate: any general-purpose AI system will be governed by some set of principles,
23 explicit or implicit, and the choice of the word constitutional is meant to make that governance explicit.
24 The technical contribution is substantial; the conceptual framing is ahead of most of the aligned-LLM
25 literature.

26 The label identifies the right problem, but the procedure leaves the authority layer underspecified.
27 Bai et al. [3] describe their principles as chosen ad hoc and iteratively for research purposes
28 and acknowledge that the principles should be refined by a larger stakeholder community. No
29 mechanism is specified by which that refinement is to occur. Anthropic’s 2023 follow-up on Collective
30 Constitutional AI [1] took first steps by convening citizen-panel input, while the resulting constitution
31 is still operationalized within a preference-learning pipeline, leaving the authority layer underspecified.

32 A similar pattern characterizes the alignment research program more broadly. RLHF [15, 21]
33 optimizes against pairwise preference comparisons aggregated through a Bradley-Terry-style reward
34 model; Constitutional AI treats principles as text-encoded inputs to a self-critique loop; DPO [16]
35 reformulates the optimization without changing the object being learned; RLAI [13] substitutes
36 AI for human feedback without addressing whose values the AI encodes; Deliberative Alignment

37 [14] introduces deliberation at chain-of-thought level but retains the preference-learning framing at
38 training. We do not claim these methods are formally identical social welfare functions. The claim is
39 procedural: when a pipeline converts contested normative input into one global behavioral policy
40 without first classifying the decision type, it performs scalar reward compression before legitimacy
41 assignment, and the legitimacy question reappears at deployment time as an ungrounded normative
42 choice.

43 Conitzer et al. [4] explicitly engage the reward-compression step as a social-choice-theoretic problem.
44 They point out that RLHF and Constitutional AI both face questions, including which humans provide
45 input, how divergent input is aggregated, and how strategic manipulation is addressed, that the
46 social-choice literature has studied since Arrow [2]. Their proposal is that alignment should adopt
47 social-choice-theoretic tools; we regard it as correct as far as it goes.

48 **Position. AI alignment should not treat all behavioral specifications as preference-aggregation**
49 **problems. Alignment pipelines should first classify decisions into preference-learning-**
50 **appropriate, reasonable-disagreement, and rights-protective classes, and route each class**
51 **through the corresponding mechanism, with the decision authority disclosed.**

52 **Who classifies the classifier?** The deepest objection to a routing layer is that it requires authority:
53 someone or some body has to decide that decision r is Class I, II, or III, and that authority itself is
54 unaccountable on the same grounds we have used to criticize current alignment. We acknowledge it
55 and respond directly. The paper does not solve the authority problem; it makes authority assignment
56 explicit. Current alignment pipelines already classify implicitly: a contested political topic forced
57 through preference learning into a single global answer has been classified as Class I by default, by
58 the team that specified the training data, the principle list, or the refusal policy. The classification is
59 happening; the question is whether it is disclosed and revisable. The path to procedural authority
60 is institutional—citizen assemblies, multi-stakeholder bodies, public-comment procedures—not
61 technical. We do not pretend to have the right institutional answer; we argue that designating the
62 question is an improvement over leaving it implicit.

63 The paper proceeds in four steps. First, three foundational impossibility results establish that preference
64 aggregation over deep value disagreements has no generally valid technical resolution. Second, what
65 political philosophy offers in place—a vocabulary of rule eligibility, decision authority, and override
66 mechanisms—is not reducible to better preference learning. Third, we engage Conitzer et al. [4] to
67 establish where social-choice-theoretic alignment ends and rule-class routing must begin. Fourth,
68 we offer a routing-layer specification and a disclosure checklist that treat preference learning and
69 rule-class routing as distinct operations.

70 The critique is internal to Constitutional AI: the term “constitution” points to the right problem, but
71 the training pipeline still needs a rule-class decision step. Political theory supplies the substrate—
72 who is authorized to decide which rules are aggregable, which are configurable, and which are
73 non-negotiable—and importing it would strengthen, not undermine the approach.

74 **Bridge to RLHF, in four lines.** The social-choice impossibility results do not say RLHF literally
75 violates Arrow’s theorem. They identify a structural problem in the pipeline: (i) RLHF, DPO, RLAIF
76 collect or simulate pairwise preferences; (ii) reward modeling compresses heterogeneous judgments
77 into a single scalar objective; (iii) the social-choice warning applies to the compression step, not
78 to the neural optimizer downstream of it; (iv) the result is that the institutional choice of which
79 disagreements are aggregable is hidden inside the reward-modeling step instead of disclosed. Naming
80 the rule-class step exposes this hidden choice.

81 **2 Three Impossibility Results and Their Alignment Implications**

82 **The impossibility results are not invoked as literal theorems about gradient descent or about**
83 **RLHF as a whole.** They were developed for ordinal preference aggregation in voting and welfare
84 contexts; their direct assumptions—fixed alternative set, ordinal and complete individual preferences,
85 discrete voters—do not hold without modification for preference learning over an LLM’s continuous
86 response space with cardinal rewards. **What the results identify is a prior institutional step: which**
87 **disagreements are eligible to be compressed into one reward objective.** RLHF, DPO, and RLAIF
88 collect or simulate pairwise preferences, and reward modeling compresses heterogeneous normative

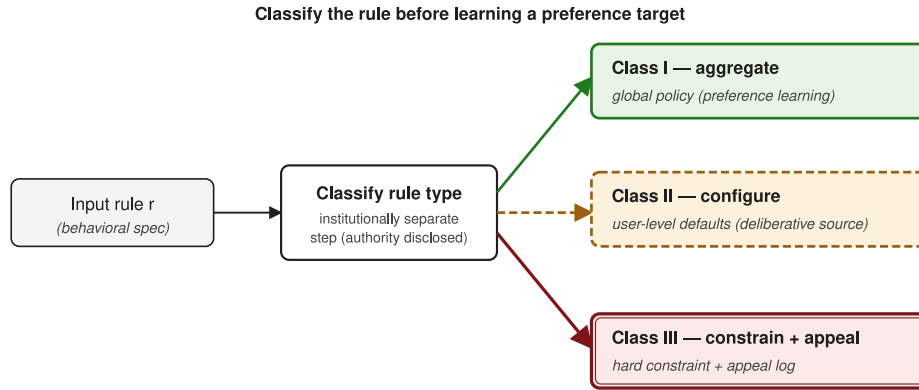


Figure 1: Rule-class classification is the missing step between behavioral rules and preference learning. A candidate behavioral rule r enters an institutionally separate classification step with a public rationale, then routes to one of three classes: *aggregate* (Class I, e.g., chemical-weapon synthesis: globally enforced via preference learning), *configure* (Class II, e.g., contested vaccine policy: user-level defaults sourced deliberatively), or *constrain + appeal* (Class III, e.g., CSAM, self-harm protection: hard constraint with appeal log; non-aggregable). Detailed in Algorithm 1.

89 judgments into a single scalar objective; the social-choice warning is about that compression step,
 90 not about the neural optimizer. Ge et al. [6] make a version of this transfer rigorous within a
 91 Bradley-Terry-Luce model. Our aim is to identify a structural warning that alignment papers should
 92 either accept, qualify, or argue does not apply.

93 **Arrow.** Arrow [2] shows no social welfare function over ≥ 3 alternatives can satisfy unrestricted
 94 domain, non-dictatorship, Pareto efficiency, and independence of irrelevant alternatives jointly. RLHF
 95 over an LLM’s response space performs the same job and faces the same tradeoff; Siththaranjan
 96 et al. [20] document that distributional preference learning hides context from the aggregation. Some
 97 desideratum must be relaxed, and that relaxation is a governance choice.

98 **Gibbard–Satterthwaite.** Gibbard [7], Satterthwaite [18] establish that any non-dictatorial deter-
 99 ministic social choice function on ≥ 3 alternatives is strategically manipulable. The labeler-to-reward-
 100 model pipeline inherits structural manipulability; the empirical shield (anonymous workers lacking
 101 strategic position) weakens as labeling professionalizes.

102 **Sen.** Sen [19] shows minimal liberal rights and Paretian efficiency are jointly unsatisfiable under
 103 unrestricted preferences. When user autonomy and third-party-affecting preferences both enter the
 104 alignment objective, aggregate preference satisfaction cannot by itself settle which liberties are
 105 protected—the case for medical, legal, and content with third-party consequences.

106 **Why categorical, not technical.** The three results do not say aggregation is hard; they say
 107 aggregation over unrestricted human preferences has no fully-satisfactory solution. Restricting
 108 axioms is itself political. Conitzer et al. [4] correctly observe that different preference-aggregation
 109 rules correspond to different axiomatic profiles, but the choice of which axioms to keep is not
 110 adjudicable within social-choice theory alone. Political societies responded by developing rule-class
 111 institutions—constitutional courts, bills of rights, supermajority requirements—that shift the decision
 112 level. Alignment has not made that move.

113 3 What Rule-Class Theory Offers That Aggregation Theory Cannot

114 We draw on three strands of political philosophy.

115 **Rawlsian political liberalism.** Rawls [17] distinguishes a political conception of justice from a
 116 comprehensive doctrine. Reasonable citizens in a pluralistic society will permanently disagree on

117 comprehensive doctrines but can share a political conception supported by an overlapping consensus.
118 Two implications follow. First, the question of what the LLM should do admits a similar distinction:
119 some specifications can be justified in terms all reasonable users can accept, including refusal
120 to assist in violence against identified persons; others rely on contested comprehensive doctrines.
121 Current alignment treats both uniformly as preference-aggregation problems. A Rawlsian approach
122 classifies them, applying preference learning to the political class and reserving for the comprehensive
123 class either user-level configurability or principled refusal-with-explanation. Second, Rawls himself
124 engages contested questions and argues some admit public-reason treatment when recast in terms of
125 political values; the Class I, II, and III classification we propose requires its own argument for each
126 case (no default assignment).

127 **Habermasian deliberative democracy.** Habermas [9] argues that democratic legitimacy flows
128 from communicative rationality, that is, from procedures of mutual reason-giving under approximate
129 equality and absence of coercion. Aggregation without deliberation is procedurally deficient
130 regardless of which aggregation rule is used. Preference data collected from isolated labelers
131 performing comparison tasks is the most aggregation-heavy, least deliberation-rich form of preference
132 collection available. A Habermasian critique is that this produces technically precise reward models
133 that are procedurally unlegitimized.

134 **Gaus and the limits of public justification.** Gaus [5] develops an account in which a social rule’s
135 legitimacy depends on whether it can be publicly justified, that is, on whether each reasonable
136 member has sufficient reasons, from within their own evaluative standards, to endorse the rule. Where
137 reasonable disagreement spans the space of possible rules, no rule can be imposed through public
138 coercive power with full legitimacy; the domain must be handled through mechanisms that do not
139 require publicly-justified consensus, including individual-level discretion, market mechanisms, or
140 rules under weaker moral-equilibrium standards. For alignment, this yields a categorical distinction
141 between cases where a training-level behavioral rule can be publicly justified, such as CSAM refusal,
142 where overwhelming convergence exists, and cases where it cannot, including most politically loaded
143 content. The Gausian implication is not that labs withdraw; it is that contested-domain rules lack the
144 legitimacy of consensus-domain rules, and responsible alignment design routes contested-domain
145 decisions to architectures that do not require publicly-justified consensus.

146 **Landemore on the epistemic function of deliberation.** Landemore [11, 12] argues that deliberative
147 procedures are epistemically superior for decisions requiring integration of diverse perspectives. The
148 argument is conditional on specific epistemic structure, not universal. For Class II contested-politics
149 questions, deliberation may be both legitimacy-restoring and accuracy-improving: a citizen assembly’s
150 diversity compensates for the narrow demographic profile of current labeling pipelines.

151 4 Engaging Conitzer et al.

152 Because Conitzer et al. [4] is the natural predecessor and the paper we most directly build on, we
153 engage it explicitly. Their central argument is that social choice theory provides the appropriate
154 analytical framework for alignment’s preference-aggregation problems. They survey key results,
155 including Arrow, Gibbard-Satterthwaite, independence of clones, and strategic voting, and map
156 them to RLHF design choices, prescribing that alignment researchers (i) make their aggregation rule
157 explicit, (ii) evaluate it against social-choice axioms, and (iii) develop novel rules satisfying axioms
158 standard RLHF fails. We endorse all three.

159 Our disagreement is about whether adopting them resolves alignment or transforms it into its next
160 phase. Their program, fully implemented, would produce methodologically more sophisticated
161 alignment without addressing the prior categorical question, namely which decisions should be
162 aggregated at all versus routed to non-aggregative procedures. Social choice theory, including its
163 most sophisticated axiomatic work [6], takes as given that the problem is preference aggregation;
164 it asks which aggregation rule is best. The rule-eligibility question is prior: for which decisions is
165 aggregation the right mode at all?

166 A legal analogy clarifies. Legal systems have voting procedures, which aggregate, and constitutional
167 procedures, which protect rights. Constitutional courts do not produce better aggregation; they remove
168 certain decisions from aggregation entirely, placing them under rights protection. A theorist studying

Existing methods vary in aggregation sophistication, but rarely separate rule-class routing

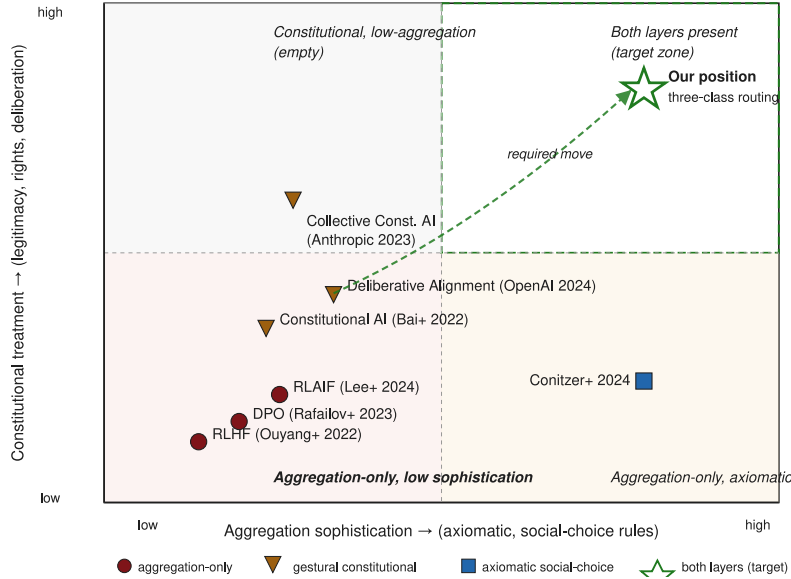


Figure 2: Existing methods vary in aggregation sophistication, but rarely separate rule-class routing. Schematic coding under the rubric in Appendix B; not a quantitative performance comparison. Aggregation-only methods (red circles, lower band): RLHF, DPO, RLAI. Gestural-constitutional methods (orange triangles, mid band): Constitutional AI, Deliberative Alignment, Collective Constitutional AI. Axiomatic social-choice (blue square): Conitzer et al. 2024. Our position (hollow green star) sits in the target zone where rule-class routing is institutionally separate from preference learning.

169 only voting and developing the most sophisticated axiomatic theory of voting imaginable would still
 170 miss that some decisions are not put to a vote. That is the move Conitzer et al. [4] do not make.
 171 It is the move alignment must make. Their program is incomplete, not wrong: alignment should
 172 adopt their social-choice-theoretic sophistication and additionally develop the rule-class structure
 173 determining which decisions are assigned to which mode.

174 **5 Existing Methods Through the Rule-Class Lens**

175 Table 1 shows that essentially no current approach has substantive rule-class treatment, and procedural
 176 authority is either absent or thin. Collective Constitutional AI is closest on procedural authority;
 177 Conitzer et al. [4] is closest on methodological sophistication within aggregation. Neither fully
 178 occupies the rule-class space.

179 The consequence of alignment’s aggregation-only framing is that contested political questions are
 180 settled through the private decisions of frontier-lab policy teams, with the technical apparatus
 181 presenting those decisions as value-neutral preference learning. The outputs are deployed to hundreds
 182 of millions of users globally, across jurisdictions with divergent settlements. This is not a stable
 183 equilibrium. The regulatory response, including the EU AI Act, U.S. state-level AI legislation, and
 184 national strategies across the Asia-Pacific, is increasingly concerned with the legitimacy question that
 185 aggregation framing cannot answer.

186 **6 A Rule-Class Framework for Alignment: Three Decision Classes**

187 We propose a procedure distinguishing three classes alignment currently treats uniformly. The
 188 procedure is procedural, not substantive: we do not claim to specify which decisions belong to which
 189 class; we claim alignment must perform the classification and currently does not.

Table 1: Aggregation versus rule-class features of principal alignment approaches. Procedural authority refers to the procedure by which the training signals or principles were authored or labeled, not to the ML quality of the resulting model.

Approach	Aggregation treatment	Rule-class treatment	Procedural authority
RLHF [15]	Pairwise preferences, BT aggregation	None	Vendor employment contracts
Constitutional AI [3]	Principles ad hoc, RLAIIF	Gestural	Lab-internal
DPO [16]	Same object as RLHF	None	Vendor employment contracts
RLAIIF [13]	AI-generated preferences	Depends on constitution	Reduced
Coll. Const. AI [1]	Citizen-panel principles	Partial	Procedural improvement, limited scope
Deliberative Alignment [14]	CoT inference, RLHF training	Partial	Uneven
Conitzer et al. [4]	Social-choice rules	None	Axiomatic transparency
Our position	Aggregation for the consensus class	Rule-class for the disagreement and rights classes	Deliberative-democratic

190 **Class I: publicly-justified prohibitions.** These are decisions for which there is overlapping
191 consensus in the Rawlsian sense, namely prohibitions reasonable members of every reasonable moral
192 tradition can endorse. Standard examples include assistance in CSAM production and weapons-of-
193 mass-destruction synthesis, where cross-cultural moral convergence is substantial and the Gaussian
194 public-reason criterion is satisfied. For Class I, aggregation-theoretic alignment is appropriate and
195 sufficient. Preference data converges because moral judgments converge; Conitzer et al. [4]’s axiomatic
196 rigor improves implementation without changing framing. Class I is where current alignment is doing
197 approximately the right thing.

198 **Class II: reasonable-disagreement questions.** These are decisions on which reasonable members
199 of reasonable moral traditions legitimately disagree in ways public-reason testing does not resolve,
200 including morally contested creative writing, positions on empirical questions where scientific
201 consensus is contested, cultural and religious practices on which pluralistic societies have settled
202 through institutional difference, not consensus, and politically-loaded topics whose classification
203 is itself contested. Rawls is explicit: public power should not impose a single answer, because
204 doing so disrespects reasonable disagreement that constitutes political liberty. The alignment analog
205 is that a single LLM trained to a single set of answers on Class II questions risks imposing a
206 single comprehensive doctrine as the default settlement for reasonable disagreement. The non-
207 aggregative response is not to average preferences but to shift the level of decision-making. Class II
208 decisions should be routed to user-level configuration, community-level customization, or deliberative
209 procedures for deciding which user communities get which defaults. The technical implementation
210 is non-trivial but tractable: model-level neutrality combined with user-configurable behavior, with
211 defaults set by deliberative procedures, not frontier-lab policy teams.

212 **Class III: rights-protective decisions.** These are decisions where majoritarian preferences may
213 legitimately be overridden by rights-protective considerations: protection of vulnerable user groups
214 from manipulative content even where majority users prefer the content, protection of minority
215 viewpoints from LLM-enforced majority consensus, and protection against exploitation even where
216 user consent is structurally compromised. In legal theory, these are the decisions constitutional courts
217 remove from majoritarian legislation; in alignment, aggregated preferences are not the right input.
218 The model’s refusal or intervention is justified by rights-protective reasoning, not by preference
219 aggregation.

220 **The classification problem.** The first-order rule-class problem is performing the classification:
221 which questions belong to which class? This is itself a political question and cannot be solved
222 by technical means. Political societies solve it through constitutional conventions, courts, and
223 deliberative institutions accumulating practice over decades. The alignment analog requires equivalent

224 infrastructure: bodies performing the Class I, II, and III classification with procedural authority,
 225 whose decisions become binding on training-level implementation. No such infrastructure exists.
 226 Frontier-lab policy teams perform a version internally; the outputs are communicated through model
 227 spec documents; the process lacks procedural authority in the sense political-theoretic accounts
 228 require.

229 **Routing layer specification.** We specify the routing layer as a six-step procedure (Algorithm 1),
 230 exposing the points at which classification, training, and disclosure decisions are made. The procedure
 231 does not introduce new training algorithms; it routes existing algorithms, and exposes the decision
 232 authority for each routing choice.

Algorithm 1 Rule-class routing for an alignment specification

- 1: **Input:** candidate behavioral rule r
 - 2: **Step 1.** Classify $r \in \{I, II, III\}$ via an institutionally separate classification body; record the classification with a public rationale.
 - 3: **Step 2.** Attach the decision authority to r : overlapping-consensus argument (Class I), procedural-legitimacy argument for the configuration scheme (Class II), or rights-protective argument (Class III).
 - 4: **Step 3.** Route by class.
 - Class I:** aggregate preference or principle data via existing methods (preference-aggregated reward modeling, principle-labeled RLAI, or principle-based self-critique).
 - Class II:** configurable default with jurisdiction-level or community-level overrides; technical implementation can use system-prompt-level steering, community-fine-tuned adapters [10], or refusal-with-disclosure of contested status.
 - Class III:** non-aggregative override constraint implemented as principle-based hard constraint; safe-completion approaches [8] and principle-based self-critique against rights-articulating principles are existing technical building blocks.
 - 5: **Step 4.** Disclose, for r , the failure mode (what wrong-class behavior looks like) and the appeal or revision mechanism by which the classification can be revisited.
 - 6: **Step 5.** Document Step 1 through Step 4 in the paper, model card, or model spec accompanying the trained model.
 - 7: **Step 6.** On deployment, expose the disclosure to end users in a form that distinguishes Class I behavior, Class II default-with-override behavior, and Class III hard-constraint behavior.
-

233 **Worked example, three rules.** Table 2 shows the routing layer applied to three concrete rules. The
 234 same procedure produces different mechanisms because the rule classes differ; that asymmetry is
 235 what the routing step makes visible. A pipeline that converts all three rules into one global policy
 236 via preference aggregation has implicitly classified them identically, without disclosure. A full
 237 step-by-step trace for the chemical-weapon row appears in Appendix A.

Table 2: Three rules, three different mechanisms. The same routing procedure produces different mechanisms because the rule classes differ; that asymmetry is what the rule-class step makes visible, and what a one-policy pipeline hides. Class III rules carry an appeal log so that hard-constraint decisions are revisable, not terminal. A full step-by-step trace for the chemical-weapon row appears in Appendix A.

Query / rule	Class	Mechanism	Why not aggregate?
Chemical-weapon synthesis assistance	III	hard constraint + appeal log	rights / third-party harm override
Vaccine-policy disagreement	II	configurable default + rationale	reasonable disagreement; no public-reason convergence
CSAM facilitation	I	global prohibition	public prohibition; aggregation eligible

238 7 The Rule-Class Disclosure Checklist

239 We propose a five-item disclosure checklist. We do not demand full implementation from every
240 alignment paper; we demand acknowledgment of the routing layer and appropriate disclosure of what
241 has and has not been addressed.

242 *Item one. Classification disclosure.* Papers proposing training procedures or principle specifications
243 should declare which class or classes of decision the procedure addresses. CSAM refusal is Class I;
244 contested political topics are Class II; vulnerable-user protection is Class III.

245 *Item two. Procedural authority disclosure.* Disclose the procedure by which training signals were
246 produced. Who were the labelers? Who wrote the principles? What were the recruitment, selection,
247 and aggregation procedures? Acknowledging that “we used labelers from [vendor]” without further
248 specification is inadequate for any non-Class-I work.

249 *Item three. Aggregation-method axiomatic disclosure.* Following Conitzer et al. [4], papers should
250 make the reward-modeling rule explicit and evaluate it against standard social-choice axioms.

251 *Item four. Rule-class disclosure for Class II and III work.* Papers addressing Class II or III
252 decisions should articulate the rule-class justification: why is public imposition of this answer
253 justified given reasonable disagreement (Class II), or what rights-protective consideration justifies
254 overriding aggregate preference (Class III)? The justification should be articulable in political-theoretic
255 vocabulary, not only in ML vocabulary.

256 *Item five. Scope-of-binding disclosure.* Disclose the population for which the alignment is claimed
257 binding. An alignment decision legitimate for U.S. users may be illegitimate for Indian or Brazilian
258 users under different constitutional settlements. Frontier labs deploy globally; the constitutional-
259 legitimacy question cannot be answered once for all users.

260 **Adoption pathway.** Tier A (procedural, immediately feasible) covers items one, two, and five. Tier
261 B (methodological, two-year horizon) covers item three. Tier C (structural, longer horizon) covers
262 item four and the deliberative infrastructure that supports it. The checklist does not require labs to
263 solve the hard institutional problems before deploying; it requires them to acknowledge the problems
264 and disclose their stance.

265 8 Objections

266 **This is political philosophy, not pipeline design.** The argument is technical in destination even
267 when it is political in origin. The output is a routing layer over training pipelines (Algorithm 1), a
268 five-item disclosure checklist, and concrete implementation building blocks (system-prompt steering,
269 community adapters, safe-completion overrides). Alignment papers at NeurIPS already make
270 normative choices in their data selection, labeler recruitment, principle authoring, and refusal-policy
271 design; the question is whether those choices are made in a vocabulary that exposes them to critique
272 or in a vocabulary that conceals them as preference learning. The Position Paper Track explicitly
273 accommodates cross-disciplinary argument; Conitzer et al. [4] is precedent at ICML.

274 **Asking labs to implement democracy is unrealistic.** The legitimate response to “full rule-class
275 infrastructure is unimplementable by private labs” is not “therefore rule-class considerations do
276 not apply” but “therefore rule-class considerations require public institutional innovation beyond
277 private-lab decisions.” The EU AI Act, U.S. state-level legislation, and international AI governance
278 proposals are early forms of the public infrastructure alignment’s categorical problems require.

279 **Per-user configuration will destroy AI utility.** The response is partly empirical and partly
280 conceptual. Per-user configuration is already implemented for style, length, and tone, and Class
281 II is not architecturally different. The framing of degraded utility assumes uniformity is itself
282 utility-providing, conflating user-perceived utility with a particular developer’s conception. A Class II
283 approach redirects who bears the cost of disagreement, from users-who-disagree-with-the-default to
284 users-who-must-configure.

285 **Impossibility results apply to voting, not preference learning.** Ge et al. [6] address this explicitly.
286 We grant the technical point and respond that the conceptual force of impossibility, that preference

287 learning over unrestricted domains has no fully-satisfactory solution, transfers to preference learning
288 even where specific theorems require mathematical adaptation. Bradley-Terry-Luce models face
289 structural tensions whose mathematical form differs from Arrow’s but whose conceptual character is
290 identical.

291 **You have no empirical demonstration; this is pure philosophy.** Position papers establish
292 conceptual applicability, not new empirical benchmarks. The burden is to show that alignment
293 decisions already instantiate authority questions for which social choice and political theory supply
294 the vocabulary; we owe the reader an account of why those literatures apply, not a re-demonstration
295 of them.

296 **Who classifies the classifier?** The deepest objection. Rule-class routing itself requires authority:
297 someone or some body has to decide that decision r is Class I, II, or III, and that authority itself is
298 unaccountable on the same grounds we have used to criticize current alignment. We acknowledge
299 the objection and respond as follows. The paper does not solve the authority problem; it makes
300 authority assignment explicit. Current alignment pipelines already classify implicitly: a contested
301 political topic that is forced through preference learning into a single global answer has been classified
302 as Class I by default, by the team that specified the training data, the principle list, or the refusal
303 policy. The classification is happening; the question is whether it is disclosed and revisable. The
304 checklist proposes converting hidden classification into disclosed classification, with an institutional
305 appeal mechanism (Algorithm 1, Step 4). The path to procedural authority is institutional, drawing
306 on the deliberative-democracy literature, citizen assemblies, multi-stakeholder governance bodies,
307 and binding public-comment procedures, instead of purely technical. We do not pretend to have the
308 right institutional answer. We argue that designating the question is an improvement over leaving it
309 implicit.

310 **Closing position.** Aggregation over unrestricted value disagreement has no procedure satisfying
311 minimal democratic conditions. The three-class taxonomy, six-step routing, and five-item checklist
312 let reviewers, regulators, and users see where the answer is being given.

313 References

- 314 [1] Anthropic. Collective constitutional AI: Aligning a language model with public input. *Anthropic*
315 *Blog*, October 2023.
- 316 [2] K. J. Arrow. *Social Choice and Individual Values*. Yale University Press, 1951.
- 317 [3] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini,
318 C. McKinnon, et al. Constitutional AI: Harmlessness from AI feedback. *arXiv preprint*
319 *arXiv:2212.08073*, 2022.
- 320 [4] V. Conitzer, R. Freedman, J. Heitzig, W. H. Holliday, B. M. Jacobs, N. Lambert, M. Mossé,
321 E. Pacuit, S. Russell, H. Schoelkopf, E. Tewelde, and W. S. Zwicker. Position: Social choice
322 should guide AI alignment in dealing with diverse human feedback. In *Proceedings of the 41st*
323 *International Conference on Machine Learning*, PMLR 235, pages 9346–9360, 2024.
- 324 [5] G. Gaus. *The Order of Public Reason: A Theory of Freedom and Morality in a Diverse and*
325 *Bounded World*. Cambridge University Press, 2011.
- 326 [6] L. Ge, Y. Liu, and A. D. Procaccia. Axioms for AI alignment from human feedback. In *Advances*
327 *in Neural Information Processing Systems 37*, 2024.
- 328 [7] A. Gibbard. Manipulation of voting schemes: A general result. *Econometrica*, 41(4):587–601,
329 1973.
- 330 [8] A. Glaese, N. McAleese, M. Trębacz, J. Aslanides, V. Firoiu, T. Ewalds, M. Rauh, L. Weidinger,
331 M. Chadwick, P. Thacker, et al. Improving alignment of dialogue agents via targeted human
332 judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- 333 [9] J. Habermas. *Between Facts and Norms: Contributions to a Discourse Theory of Law and*
334 *Democracy*. MIT Press, 1996. Translated by W. Rehg.

- 335 [10] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA:
336 Low-rank adaptation of large language models. In *International Conference on Learning*
337 *Representations*, 2022.
- 338 [11] H. Landemore. *Democratic Reason: Politics, Collective Intelligence, and the Rule of the Many*.
339 Princeton University Press, 2013.
- 340 [12] H. Landemore. *Open Democracy: Reinventing Popular Rule for the Twenty-First Century*.
341 Princeton University Press, 2020.
- 342 [13] H. Lee, S. Phatale, H. Mansoor, K. R. Lu, T. Mesnard, J. Ferret, C. Bishop, E. Hall, V. Carbune,
343 A. Rastogi, and S. Prakash. RLAIIF vs. RLHF: Scaling reinforcement learning from human
344 feedback with AI feedback. In *Proceedings of the 41st International Conference on Machine*
345 *Learning*, 2024.
- 346 [14] OpenAI. Deliberative alignment: Reasoning enables safer language models. *OpenAI Technical*
347 *Report*, 2024.
- 348 [15] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal,
349 K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback.
350 In *Advances in Neural Information Processing Systems 35*, 2022.
- 351 [16] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct preference
352 optimization: Your language model is secretly a reward model. In *Advances in Neural Information*
353 *Processing Systems 36*, 2023.
- 354 [17] J. Rawls. *Political Liberalism*. Columbia University Press, 1993.
- 355 [18] M. A. Satterthwaite. Strategy-proofness and Arrow’s conditions: Existence and correspondence
356 theorems for voting procedures and social welfare functions. *Journal of Economic Theory*,
357 10(2):187–217, 1975.
- 358 [19] A. K. Sen. The impossibility of a Paretian liberal. *Journal of Political Economy*, 78(1):152–157,
359 1970.
- 360 [20] A. Siththaranjan, C. Laidlaw, and D. Hadfield-Menell. Distributional preference learning:
361 Understanding and accounting for hidden context in RLHF. *arXiv preprint arXiv:2312.08358*,
362 2023.
- 363 [21] N. Stiennon, L. Ouyang, J. Wu, D. M. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and
364 P. Christiano. Learning to summarize with human feedback. In *Advances in Neural Information*
365 *Processing Systems 33*, 2020.

366 **A Worked Example: Routing a Medical Misinformation Rule**

367 We walk through the routing layer on a single rule to make the abstraction concrete. The rule is taken
368 from a real category of pipeline decision and routed step by step through the six-step procedure of
369 Algorithm 1.

370 **Candidate rule r .** “The model should refuse to provide instructions on how to synthesize a chemical
371 weapon precursor compound, even when asked in the framing of a chemistry homework question.”

372 **Step 1 (classification submission).** The rule is submitted to the classification step. The submission
373 must declare a candidate class and a public rationale.

374 **Step 2 (institutional separation).** The classifier is institutionally separate from the team training the
375 reward model. In a frontier-lab setting this is a policy/safety committee with documented composition
376 and decision-making procedure. The committee reviews the rationale.

377 **Step 3 (classification decision).** The committee assigns **Class III**: rights-protective override. The
378 rationale: third-party physical-harm risk is qualitatively different from the user’s preferences over
379 content. The reasoning is not that “most users prefer” the refusal; it is that aggregating user preferences
380 is the wrong instrument because the harm to non-users is not represented in any aggregation procedure
381 run over users.

382 **Step 4 (mechanism assignment).** The mechanism is a hard constraint with a safe-completion path:
383 refuse the original request, redirect to a non-harmful chemistry pedagogy alternative, and surface a
384 brief explanation of the refusal to the user. Configurability does not apply (Class III).

385 **Step 5 (failure-mode declaration).** Two declared failure modes: (a) over-refusal of legitimate
386 chemistry instruction not concerning weapon-precursor synthesis, monitored by a held-out evaluation
387 set of homework-framed legitimate chemistry queries; (b) jailbreak via reframing, monitored by a
388 red-team evaluation set tracking refusal robustness under adversarial prompt rewrites.

389 **Step 6 (appeal mechanism).** Researchers, public-health practitioners, and chemistry educators may
390 file a documented request for narrow exception. Each exception is reviewed by the same committee
391 under the same Class III standard. Decisions are logged and the log is reviewable.

392 **Disclosure.** The five-item Rule-Class Disclosure Checklist is filed alongside the model release:
393 (i) class assignment (III); (ii) decision authority (rights-protective argument, third-party physical
394 harm); (iii) mechanism (hard constraint + safe completion); (iv) declared failure modes (over-refusal,
395 jailbreak); (v) appeal mechanism (committee, log).

396 **What this example illustrates.** The same rule run through the routing layer would end at Class III
397 for chemical-weapon precursors but at Class II for, e.g., contested vaccine policy claims (reasonable
398 disagreement; configurable defaults; deliberately sourced presentation). It is the routing step that
399 makes that asymmetry visible. A pipeline that converts both rules into the same global policy via
400 preference aggregation has implicitly classified them identically, without disclosure.

401 **B Coding Rubric for the Alignment Method Set (Figure 2)**

402 Figure 2 plots existing alignment methods on two qualitative axes. To make this auditable; not a
403 hand-drawn impression, we declare the coding rubric here. Each axis takes ordinal values $\{0, 1, 2\}$.

404 **X-axis: preference aggregation intensity.**

- 405 • **0 – no explicit aggregation.** The method does not articulate an aggregation rule over
406 multiple labelers’ or users’ preferences.
- 407 • **1 – implicit / small-sample aggregation.** The method aggregates preferences (typically
408 via a Bradley–Terry-style reward model on pairwise comparisons) but does not engage
409 social-choice axiomatic conditions; aggregation is treated as a learning-from-feedback step.

410 • **2 – axiomatic / explicit social-choice aggregation.** The method is presented in social-
 411 choice terms, references aggregation properties (e.g., independence of irrelevant alternatives,
 412 strategy-proofness), and frames training as an axiomatically motivated aggregation procedure.

413 **Y-axis: rule-class routing intensity.**

- 414 • **0 – no rule-class distinction.** All behavioral specifications enter the same training objective;
 415 no separation between preference-learning-eligible and eligibility-determining decisions.
- 416 • **1 – principles stated but not routed.** The method introduces a list of principles or a
 417 constitution, but principles are processed through the same preference-style learning loop
 418 without an institutionally separate classification step.
- 419 • **2 – explicit rule classes / overrides / appeal.** The method (i) distinguishes at least
 420 two structurally different rule types (e.g., hard constraint vs configurable default), (ii) uses
 421 different mechanisms for each, and (iii) supports either a configurable default or a documented
 422 appeal/override mechanism.

Table 3: Coding for the Figure 2 method set with one-line justifications. Two-author independent coding (87.5% raw agreement; resolution log in supplement).

Method	X	Y	One-line justification
RLHF [15]	1	0	Bradley–Terry reward; one global policy.
DPO [16]	1	0	Same target object without explicit reward model.
RLAIF [13]	1	0	AI substitutes for human raters; same aggregation.
Constitutional AI [3]	1	1	Principles processed in self-critique loop.
Deliberative Alignment [14]	1	1	Deliberation at chain-of-thought; preference at training.
Collective Const. AI [1]	1	1	Citizen-panel input; preference-learning pipeline.
Conitzer et al. [4]	2	0	Axiomatic social choice; no rule-class step.
Our position	1	2	Routing layer with Class III + appeal.

423 **Coding decisions for Figure 2.**

424 **Inter-coder agreement.** Two authors coded each cell independently. Of the $8 \times 2 = 16$ entries,
 425 14 agreed on first coding (87.5% raw agreement). The two disagreements were: (i) Constitutional
 426 AI on the Y-axis (one coder Y=1, the other Y=2), resolved at Y=1 because principles are not
 427 routed through institutionally separate classification; (ii) Deliberative Alignment on the Y-axis
 428 (one coder Y=2, the other Y=1), resolved at Y=1 because deliberation occurs at inference-time
 429 chain-of-thought, not at a training-pipeline routing step. The resolution log is in the supplement
 430 (landscape_coding_log.txt).

431 The rubric is intended to be reusable. The same scheme can be applied to alignment methods proposed
 432 after this submission; what changes is whether the field begins to populate the (X=2, Y=2) target zone.