

FAIRNESS BEYOND PREDICTION: RETHINKING ALIGNMENT PROCEDURES FOR AGENTIC AI SYSTEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

As AI systems transition from static predictive models to interactive and agentic systems that plan, adapt, and act over time, classical algorithmic fairness frameworks become insufficient. Traditional fairness notions were largely developed for single-shot prediction or prediction-based decision making, and they often fail to capture procedural, temporal, and emergent fairness issues introduced by alignment procedures such as preference optimization, RLHF, and policy constraints. This paper argues that fairness must be analyzed as a property of *alignment processes* and *agentic dynamics*, not solely model outputs. We identify key fairness failure modes that can arise across value learning, policy adaptation, and long-horizon deployment, including representational asymmetries in preference data, feedback-driven amplification, and path-dependent disparities. We then propose a lightweight conceptual framework for fairness across alignment procedures, outlining evaluation questions and governance implications for building fairer agentic AI systems.

1 INTRODUCTION

Algorithmic fairness research has historically focused on predictive settings—classification, regression, and ranking—where a model produces an output given an input, and fairness is evaluated via parity constraints across demographic groups (e.g., equalized odds, demographic parity). These approaches assume that decisions are largely *single-shot* and that the system does not substantially adapt its behavior based on interaction.

This assumption increasingly breaks in modern deployments. Contemporary AI systems are often *agentic*: they reason, plan, interact with users and tools, and update their behavior over time. Examples include conversational assistants, tool-augmented agents, personalized tutoring systems, and autonomous decision-support systems. These systems are shaped by *alignment procedures*—such as preference learning, RLHF, instruction tuning, safety fine-tuning, and constraint-based policies—that define how an agent learns what to do and how it adapts in response to feedback.

In this regime, fairness failures rarely arise from a single prediction. Instead, unfairness can emerge *procedurally* through how objectives are learned, how feedback is incorporated, and how adaptation changes downstream behavior. A system may appear fair at one evaluation snapshot yet become unfair after extended interaction, feedback loops, or distributional shifts. As a result, fairness becomes inseparable from alignment: alignment procedures and agentic dynamics determine whose interests are represented, whose preferences shape behavior, and how harms accumulate over time.

This paper asks:

How should fairness principles and tools evolve when AI systems not only predict, but also adapt and act?

We argue that fairness must be evaluated across the *alignment pipeline itself*. We outline fairness risks introduced by alignment procedures, propose a lightweight framework for fairness across alignment and agentic systems, and provide practical evaluation and governance implications. Our goal is

054 not to introduce a new metric, but to clarify where existing fairness tools break down and to propose
055 concrete questions to support fairness-aware alignment.
056

057 058 2 WHY CLASSICAL FAIRNESS BREAKS FOR AGENTIC SYSTEMS 059

060 061 2.1 FROM OUTCOME FAIRNESS TO PROCEDURAL FAIRNESS

062 Classical fairness notions focus on outcomes (e.g., prediction parity across groups). In agentic
063 systems, the relevant object is often a *policy* that selects actions, observes consequences, and updates
064 based on feedback. This shift introduces three challenges.
065

066 **Temporal coupling.** Actions affect future states and opportunities. Even if an agent satisfies a parity
067 constraint at time t , its actions may change the environment and create disparities at later times.

068 **Feedback loops.** Agent behavior influences the data used for continued alignment. Unequal feed-
069 back quality or feedback availability across groups can amplify disparities.

070 **Policy adaptation.** Many agents continuously update. Fairness properties are not static: a fairness
071 guarantee at deployment does not imply fairness after adaptation.
072

073 These characteristics suggest that fairness must include *procedural* questions: how the agent learns
074 values, how it updates, and how harms accumulate.
075

076 077 2.2 FAIRNESS RISKS INTRODUCED BY ALIGNMENT PROCEDURES

078 Alignment procedures are often treated as neutral mechanisms for shaping behavior. We argue they
079 are primary sites of fairness risk:
080

081 **Preference learning.** Preferences are incomplete, context-dependent, and socially biased. If pref-
082 erence data under-represents marginalized perspectives, the learned objective can encode asymmet-
083 rical value weighting even when downstream predictions appear balanced.

084 **RLHF and reward optimization.** Reward models and feedback providers can contain implicit
085 bias. Optimization can privilege short-term performance signals that correlate with dominant-group
086 norms, leading to systematic disadvantage over repeated interactions.

087 **Constraint-based alignment.** Constitutional rules or safety constraints can freeze normative as-
088 sumptions that do not generalize across cultures or contexts. These constraints may be protective
089 for some groups while burdening others.
090

091 These failure modes are not fully captured by post-hoc fairness checks on outputs. They require
092 analyzing fairness *within* the alignment pipeline.
093

094 095 3 A FRAMEWORK FOR FAIRNESS ACROSS ALIGNMENT PROCEDURES

096
097 We propose viewing fairness as a process-level property spanning three stages. This framing sup-
098 ports practical audits and clarifies where interventions should occur.
099

100 101 3.1 STAGE I: VALUE LEARNING FAIRNESS

102 This stage concerns *what* the system is optimized for. Fairness questions include: (i) whose pref-
103 erences or values are represented, (ii) how conflicting preferences are aggregated, (iii) whether mi-
104 nority or long-tail preferences are preserved, and (iv) whether value learning is robust to strategic or
105 unequal feedback.
106

107 A failure at this stage yields a misaligned objective that persists even if downstream behavior is later
constrained.

108 3.2 STAGE II: POLICY ADAPTATION FAIRNESS
109

110 Agentic systems adapt through interaction. Fairness risks include: (i) unequal exploration across
111 user groups, (ii) differential learning rates that favor groups producing more feedback, (iii) path
112 dependence where early interactions lock in biased behaviors, and (iv) shifting behavior under dis-
113 tribution drift.

114 This motivates fairness evaluation across *learning trajectories*, not only static snapshots.
115

116 3.3 STAGE III: LONG-HORIZON IMPACT FAIRNESS
117

118 Agentic decisions accumulate over time. Small asymmetries can compound into disparities in ac-
119 cess, opportunity, or exposure to harm. Long-horizon fairness requires evaluating: (i) cumulative
120 reward or benefit distribution, (ii) compounding error dynamics, (iii) differential risk exposure, and
121 (iv) the downstream societal impacts of repeated agent actions.

122 This stage connects technical evaluation to governance: fairness becomes an ongoing responsibility,
123 not a one-time certification.
124

125 4 IMPLICATIONS FOR EVALUATION AND GOVERNANCE
126

127 Our framework suggests that fairness evaluation for aligned, agentic systems must extend beyond
128 output metrics.
129

130 **Alignment-aware fairness audits.** Audits should examine preference data composition, feedback
131 mechanisms, and reward model behavior. Questions include: which groups contribute feedback,
132 what is optimized, and how updates change behavior across groups.

133 **Process transparency.** Understanding alignment choices and update rules becomes central to fair-
134 ness accountability. This includes documenting data sources, preference aggregation methods, and
135 constraints used in alignment.

136 **Governance upstream.** Oversight and risk management should target alignment procedures (how
137 systems are trained and updated), not only deployed outputs. This is especially important for long-
138 lived agents deployed in sensitive domains.

139 We emphasize that these recommendations are compatible with existing fairness tools; the key shift
140 is to apply them at *procedural points* in the alignment pipeline and to measure fairness over time.
141

142 5 CONCLUSION
143

144
145 As AI systems become increasingly agentic, fairness can no longer be treated as a static property
146 of predictions. It must be understood as a dynamic property of alignment procedures and agentic
147 adaptation. This paper reframes fairness as a procedural concern that spans value learning, policy
148 adaptation, and long-horizon deployment. By identifying where classical fairness breaks down and
149 proposing a process-level framework, we aim to support the development of fairer, more accountable
150 aligned and agentic AI systems.
151

152 REFERENCES

153 A LLM USAGE DISCLOSURE
154

155 Large language models were used for language refinement and formatting assistance. All conceptual
156 framing, arguments, and technical content were authored and verified by the human authors.
157
158
159
160
161