

# Dynamic Fairness-Aware Recommendation through Multi-Agent Social Choice

Amanda Aird, Paresha Farastu, Joshua Sun,  
Amy Volda, Nicholas Mattei, and Robin Burke

## Abstract

Algorithmic fairness in the context of personalized recommendation presents significantly different challenges to those commonly encountered in classification tasks. Researchers studying classification have generally considered fairness to be a matter of achieving equality of outcomes between a protected and unprotected group, and built algorithmic interventions on this basis. We argue that fairness in real-world application settings in general, and especially in the context of personalized recommendation, is much more complex and multi-faceted, requiring a more general approach. We propose a model to formalize multistakeholder fairness in recommender systems as a two stage social choice problem. In particular, we express recommendation fairness as a novel combination of an allocation and an aggregation problem, which integrate both fairness concerns and personalized recommendation provisions, and derive new recommendation techniques based on this formulation. Simulations demonstrate the ability of the framework to integrate multiple fairness concerns in a dynamic way.

## 1 Introduction

Recommender systems are personalized machine learning systems that support users' access to information in applications as disparate as rental housing, video streaming, job seeking, social media feeds and online dating. The challenges of ensuring fair outcomes in such systems have been addressed in a growing body of research literature surveyed by Ekstrand et al [16]. Despite these research efforts, some key limitations have remained unaddressed, limitations that render this work inadequate for the applications for which it is intended.

The first limitation we see in current work is that researchers have generally assumed that the problem of group fairness can be reduced to the problem of ensuring equality of outcomes between a protected and unprotected group, or in the case of individual fairness, that there is a single type of fairness to be addressed for all individuals. Where fairness for multiple groups has been considered (e.g., Sonboli et al. [35], Kearns et al. [21]), it is defined in the same way for all groups.

We believe that this limitation is severe and not representative of realistic recommendation tasks in which fairness is sought. US anti-discrimination law, for example, identifies multiple protected categories relevant to settings such as housing, education and employment including gender, religion, race, age, and others [4]. But even in the absence of such external criteria, it seems likely that any setting in which fairness is a consideration will need to incorporate the viewpoints of multiple groups.

We also expect that fairness will mean different things for different groups. Consider, for example, a system recommending news articles. Fairness might require that, over time, readers see articles that are geographically representative of their region: rural and urban or uptown vs downtown, for example. But fairness in presenting viewpoints might also require that any given day's set of headlines represent a range of perspectives. These are two different views of what fairness means, entailing different measurements and potentially different types of algorithmic interventions.

The second limitation that we see in current work is that fairness-aware interventions in

recommender systems as well as many other machine learning contexts, have a static quality. In many applications, a system is optimized for some criterion and when the optimization is complete, it produces decisions or recommendations based on that learned state [29]. We think of fairness as a dynamic state, especially when what is of primary concern are fair outcomes. A recommender system’s ability to produce outcomes that meet some fairness objective may be greatly influenced by context: what items are in inventory, what types of users arrive, how fair the most recent set of recommendations has been, and many others. A static policy runs the risk of failing to capitalize on opportunities to pursue fairness when they arise and/or imposing fairness when its cost is high, by being insensitive to the context.

Our contribution in this paper is the design of an architecture for implementing fairness in recommender systems that addresses both of these limitations. (Note that portions of this work appeared in workshop form in [10].) We start from the assumption that multiple fairness concerns will be active at any one time, and that these fairness concerns can be relatively unrestricted in form. Secondly, we build the framework to be dynamic in that decisions are always made in the context of historical choices and results.

Our research in fairness examines concepts inspired by the application context of Kiva Microloans, which offers a platform (Kiva.org) for crowd-sourcing the funding of microloans, mostly in the developing world. Kiva’s users (lenders) choose among the loan opportunities offered on the platform; microloans from multiple lenders that are aggregated and distributed through third party non-governmental organizations around the world. Kiva Microloans’ mission specifically includes considerations of “global financial inclusion”; as such, incorporating fairness in its recommendation of loans to potential users (lenders) is a key goal. We will use Kiva’s platform as an example throughout this paper. However, the analytic findings are not specific to this setting. This is part of ongoing work to understand and create interventions for fairness for online recommendation platforms [33, 11]

**Notes on Workshop Version:** This paper is a shortened version of a longer journal paper Aird et al. [1]. We have included an extensive discussion of design considerations, table of notations, and a large set of experimental results in the Appendix. The main body of this document serves as an overview and introduction to our modeling choices.

## 2 Related Work

There have been a number of efforts that explicitly consider the multisided nature of fairness in recommendation and matching platforms. Patro et al. [31] investigate fairness in two-sided matching platforms where there are both producers and consumers. They note, as we do, that optimizing fairness for only one side of the market can lead to very unfair outcomes to the other side of the market. Patro et al. [31] also appeal to the literature on the fair allocation of indivisible goods from the social choice literature [37]. They devise an algorithm that guarantees Max-min share fairness of exposure to the producer side of the market and envy-free up to one item to the consumer side of the market. Their work is closest to the allocation phase of our algorithm. However, in contrast to our work they only use exposure on the producer side and relevance on the consumer side as fairness metrics, whereas our work aims to capture additional definitions. Also, we note that envy-freeness is only applicable when valuations are shared: a condition not guaranteed in a personalized system. It is possible for a user with unique tastes to receive low utility recommendations and still not prefer another user’s recommendation lists. Also, our fairness formulation extends beyond the users receiving recommendations to providers of recommended items and envy-freeness provides no way to compare users who are getting different types of benefits from a system. In addition our fairness definitions are dynamic, a case not considered by [31].

Like Patro et al. [31], the work of Sühr et al. [36] investigates fairness in two-sided platforms, specifically those like Uber or Lyft where income opportunities are allocated to

drivers. However, unlike our work and the work of Patro et al. [31], Sühr et al. [36] take proportionality as their definition of fairness, specifically proportionality with respect to time in a dynamic setting, and ensure that there is a fair distribution of income to the provider side of the platform.

Freeman et al. [18] investigate what they call *dynamic social choice functions* in settings where a fixed set of agents select a single item to share over a series of time steps. The work focuses on overall utility to the agents instead of considering the multiple sides of the recommendation interaction. Their problem is fundamentally a voting problem since all agents share the result, whereas we are focused on personalized recommendation. Their goal is to optimize the Nash Social Welfare of the set of agents (that remains fixed at each time step) and present four algorithms to find approximately optimal solutions. This work has a similar flavor to classical online learning / weighting experts problems [12] in the sense that the agent preferences remain fixed and the goal is to learn to satisfy them over a series of time steps.

The architecture presented here advances and generalizes the approach found in [34]. Like that architecture, fairness concerns are represented as agents and interact through social choice. However, in [34], the allocation mechanism selects only a single agent at each time step and the choice mechanism has a fixed, additive, form. We allow for a wider variety of allocation and choice mechanisms, and therefore present a more general solution.

Ge et al. [20] investigate the problem of long term dynamic fairness in recommendation systems. They, like our work, highlight the need to ensure that fairness is ensured as a temporal concept and not see recommendation as a static, one off, decision. To this end they propose a framework to ensure fairness of exposure to the producers of items by casting the problem as a constrained Markov Decision Process where the actions are recommendations and the reward function takes into account both utility and exposure. Ge et al. [20] propose a novel actor-critic deep reinforcement learning framework to accomplish this task at the scale of large recommender systems with very large user and item sets. Again, this work fixes definitions of fairness a priori, although their learning methodology may serve as inspiration to our allocation stage problems in the future.

Morik et al. [25] investigate the problem of learning to rank over large item sets while ensuring fairness of merit based guarantees to groups of item producers. Specifically, they adapt existing methods to ensure that the exposure is *unbiased*, e.g., that it is not subject to rich-get-richer dynamics, and *fairness* defined as exposure being proportional to merit. Both of these goals are built into the regularization of the learner. In essence the goal is to learn user preferences while ensuring the above two desiderata. In contrast, our work factors out the recommendation methodology and we encapsulate the desired fairness definitions as separate agents rather than embedded in the learning algorithm.

Finally, our recommendation allocation problem has some similarities with those found in computational advertising, where specific messages are matched with users in a personalized way [38, 39]. Because advertising is a paid service, these problems are typically addressed through mechanisms of monetary exchange, such as auctions. There is no counterpart to budgets or bids in our context, which means that solutions in this space do not readily translate to supporting fair recommendation [41, 15, 40].

### 3 Example

In this section, we work through a detailed example demonstrating the function of the architecture through several iterations of user arrivals. We formally describe the process with notations in Section 5.2

### 3.1 Agents

Consider the following set of fairness agents and their associated evaluations and preferences. We assume in this example that in all cases the agents’ compatibility functions follow the pattern described in [35] where the entropy of the user profile relative to the sensitive feature is calculated and users with high entropy are determined to be good targets for fairness-enhancing interventions:

- $f_H$ : **Health** This agent is concerned with promoting loans to the health sector. Its evaluation function compares the proportion of loans in the database in the health sector against the proportion of health recommendations in the recommendation list history. Its preference function is binary: if the loan is in the health sector, the score is 1; otherwise, zero.
- $f_A$ : **Africa** This agent is concerned with promoting loans to Africa. Its evaluation function, however, is list-wise. It counts lists in the recommendation if they have a least one loan recommendation to a country in Africa, and consider a fair outcome one in which every list has at least one such loan. Its preference function will be similarly binary as the  $f_H$  agent.
- $f_G$ : **Gender Parity** This agent is concerned with promoting gender parity within the recommendation history. If, across the previously generated recommendation lists, the number of men and women presented is proportional to their prevalence in the database, its evaluation will return 1. However, its preference function is more complex than those above. If the women are underrepresented in the history, it will prefer loans to female borrowers, and conversely for men. For the sake of argument, we can also say that the agent has access to the gender breakdown of borrower groups and therefore may return a preference score in proportion to the number of women (or men) contained therein.<sup>1</sup>
- $f_L$ : **Large** This agent is concerned with promoting loans with larger total amounts: over \$5,000. Internal Kiva research has shown that such loans are often very productive because they go to cooperatives and have a larger local impact. However, the same research has shown that Kiva users are less likely to support them because each contribution has a smaller relative impact. This agent is similar to the  $f_A$  agent above in that it seeks to make sure each list has one larger loan.

### 3.2 Loans

Consider the contents of Table 1. For the sake of example, we will assume these loans, characterized by the Region, Gender, Section and Amount, constitute the set of loans available for recommendation.

	$\phi_1^s$ : Region	$\phi_2^s$ : Gender	$\phi_3^s$ : Sector	$\phi_4$ : Amount
$v_1$	Africa	Male	Agriculture	\$5,000-\$10,000
$v_2$	Africa	Female	Health	\$500-\$1,000
$v_3$	Middle-East	Female	Clothing	\$0-\$500
$v_4$	Central America	Female	Clothing	\$5,000-\$10,000
$v_5$	Central America	Female	Health	\$0-\$500
$v_6$	Middle-East	Female	Clothing	\$0-\$500

Table 1: Set of Potential Loans.

---

<sup>1</sup>Note: Kiva’s borrower database currently recognizes only binary gender categories.

### 3.3 Mechanisms

For the sake of exposition, we posit two very simple mechanisms for allocation and choice. We will assume that our allocation mechanism is a single outcome lottery, e.g., a randomized allocation mechanism [7]. One agent will be chosen to participate in the choice mechanism, based on a random draw with probabilities based on the historic unfairness and user compatibility as measured by each agent.

We assume that the recommendation lists are of size 3 and the choice mechanism uses a weighted voting / score-based mechanism [6] using a weighted sum of 0.75 on the personalized results for the recommender system and 0.25 on the allocated fairness agent.

### 3.4 Users

At time  $t_1$ , **User**  $u_1$  arrives at the system and the recommendation process is triggered. The user has previously supported small loans only in Central America and Middle East, but has lent to a wide variety of sectors and genders.

For the sake of example, we will assume that the agents measure their prior history relative to their objectives as equally unfair at 0.5, except the Gender Parity agent, which starts out at parity and therefore returns a value of 1. However, the compatibility functions for  $f_A$  and  $f_L$  returns lower scores because of the user’s historical pattern of lending. This yields a lottery in which  $f_G$  has probability zero,  $f_A$  has a low probability, and  $f_H$  a higher one. The allocation mechanism chooses randomly, and we will assume that  $f_H$ , the health-focused agent, is picked.

The recommender returns the following list of items and predicted ratings  $[\{v_6, 0.6\}, \{v_4, 0.5\}, \{v_5, 0.3\}, \{v_3, 0.3\}, \{v_1, 0.0\}, \{v_2, 0.0\}]$ . The  $f_H$  agent gives a score of 1 to the health-related loans  $v_2$  and  $v_5$  and 0 to all others. The choice mechanism combines these scores as described above and returns the final recommendation list  $[\{v_5, 0.475\}, \{v_6, 0.45\}, \{v_4, 0.375\}]$ . Note that the Health agent has successfully promoted its preferred item to the first position in the list.

For the sake of example, we assume that the agents’ evaluation functions are very sensitive. Therefore, when **User**  $u_2$  arrives, the results of the previous recommendations have caused the evaluations to shift such that the Health  $f_H$  and Large  $f_L$  agents are now satisfied (note that  $v_4$  is included in  $u_1$ ’s list and it was a large loan), the Gender parity agent  $f_G$  is now at 0.9 (note that there is only one male loan in the database) but the Africa agent  $f_A$ , which got nothing in  $u_1$ ’s list is now at 0.25. We assume that  $u_2$  is similar to  $u_1$  in profile and therefore compatibility, but  $f_A$  has a much worse fairness score than  $f_G$ , and therefore a high allocation probability. We will assume  $f_A$  is chosen.

Because this user has similar preferences to  $u_1$ , they get the same recommendations:  $[\{v_6, 0.6\}, \{v_4, 0.5\}, \{v_5, 0.3\}, \{v_3, 0.3\}], \{v_1, 0.0\}, \{v_2, 0.0\}]$ . The  $f_A$  agents scores the two loans from Africa ( $v_1$  and  $v_2$ ) at 1 and the others at 0. So, after randomly breaking the tie between  $v_1$  and  $v_2$ , the final recommendation list is  $[\{v_6, 0.45\}, \{v_4, 0.375\}, \{v_1, 0.25\}]$ .

When **User**  $u_3$  arrives, all four agents find themselves scoring fairness at 1 over the evaluation window and so no agents are allocated. The results from the recommendation algorithm pass through the choice mechanism unchanged and are delivered to the user.

## 4 Formalizing Fairness Concerns

A central tenet of our work is that fairness is a contested concept [27]. From an application point of view, this means that ideas about fairness will be grounded in specific contexts and specific stakeholders, and that these ideas will be multiple and possibly in tension with each other. From a technical point of view, this means that any fairness-aware recommender system should be capable of integrating multiple fairness concepts, arising as they may from this contested terrain.

A central concept in this work is the idea of a *fairness concern*. We define a fairness concern as a specific type of fairness being sought, relative to a particular aspect of recommendation outcomes, evaluated in a particular way. For example, a possible fairness concern in the microlending context might be group fairness relative to different geographical regions considered in light of the exposure of loans from these regions in recommendation lists.<sup>2</sup> The concern identifies a particular aspect of the recommendation outcomes (in this case, their geographical distribution), the particular fairness logic and approach (more about this below), and the metric by which fair or unfair outcomes are determined.

The first consideration in building a fairness-aware recommender system is the question of what fairness concerns surround the use of the recommender system, itself. Many such concerns may arise and like any system-building enterprise, there are inevitably trade-offs involved in the formulation of fairness concerns. An organization may decide to incorporate only the highest-priority concerns into its systems. An initial step in fairness-aware recommendation is for an organization to consult its institutional mission and its internal and external stakeholders with the goal of eliciting and prioritizing fairness concerns. An example of this kind of consultation can be seen in the WeBuildAI project [22] and its participatory design framework for AI.

In addition to addressing different aspects of system outcomes, different fairness concerns may invoke different logics of fairness. Welfare economists have identified a number of such logics and we follow Moulin [26] who identifies four:

**Exogenous Right:** A fairness concern is motivated by exogeneous right if it follows from some external constraint on the system. For example, the need to comply with fair lending regulations may mean that male and female borrowers should be presented proportionately to their numbers in the overall loan inventory.

**Compensation:** A fairness concern that is a form of compensation arises in response to observed harm or extra costs incurred by one group versus others. For example, loans with longer repayment periods are often not favored by Kiva users because their money is tied up for longer periods. To compensate for this tendency, these loans may need to be recommended more often.

**Reward:** The logic of reward is operational when we consider that resources may be allocated as a reward for performance. For example, if we know that loans to large cooperative groups are highly effective in economic development, we may want to promote such loans as recommendations so that they are more likely to be funded and realize their promise.

**Fitness:** Fairness as fitness is based on efficiency. A resource should go to those best able to use it. In a recommendation context, it may mean matching items closely with user preferences. For example, when loans have different degrees of repayment risk, it may make sense to match the loan to the risk tolerance of the lender.

It is clear that fairness logics do not always pull in the same direction. The invocation of different logics are often at the root of political disagreements: for example, controversies over the criteria for college admissions sometimes pit ideas of reward for achievement against ideas of compensation for disadvantage.

Recommender systems often operate as two-sided platforms, where one set of individuals are receiving recommendations and possibly acting on those recommendations (consumers), and another set of individuals is creating or providing items that may be recommended (providers) [8]. Consumers and providers are considered, along with the platform operator,

---

<sup>2</sup>We are currently conducting research to characterize fairness concerns appropriate to Kiva’s recommendation applications. At this stage, we can only speculate about the fairness concerns that might arise in that work. None of the discussion here is intended to represent design decisions or commitments to particular concerns and/or their formulation.

to be the direct stakeholders in any discussion of recommender system objectives. Fairness concerns may derive from any stakeholder, and may need to be balanced against each other. The platform may be interested in enforcing fairness, even when other stakeholders are not. For example, the average recommendation consumer might only be interested in the best results for themselves, regardless of the impact on others. Fairness concerns can arise on behalf of other, indirect, stakeholders who are impacted by recommendations but not a party to them. An important example is *representational fairness* where we are concerned about the way the outputs of a recommender system operate to represent the world and classes of individuals within it: for example, the way the selection of news articles might end up representing groups of people unfairly [28] (see [16] for additional discussion). As a practical matter, representational fairness concerns can be handled in the same way as provider-side fairness for our purposes here.

Finally, we have the consideration of group versus individual fairness. This dichotomy is well understood as a key difference across types of fairness concerns, defining both the target of measurement of fairness and the underlying principle being upheld. Group fairness requires that we seek fairness across the outcomes relative to predefined protected groups. Individual fairness asks whether each individual user has an appropriate outcome and assumes that users with similar profiles should be treated the same. Just as there are tensions between consumer and provider sides in fairness, there are fundamental incompatibilities between group and individual fairness. Treating all of the outcomes for a group in aggregate is inherently different than maintaining fair treatment across individuals considered separately. Friedler et al. offer a thorough discussion of this topic [19].

Label	Fairness type	Logic	Side	Who is Impacted	Evaluation
LowCountry	Group	Comp.	Provider	Borrowers from countries with lower funding rates	Exposure of loans in recommendation lists
LargeAmt	Group	Reward	Provider	Borrowers in consortia seeking larger loans	Exposure of loans in recommendation lists
Repay	Individual	Reward	Provider	All borrowers	Loan exposure proportional to repayment probability
LowSector	Group	Exo. right	Provider	Borrowers in sectors with lower funding rates	Exposure of loans in recommendation lists
AllCountry	Individual	Exo. right	Provider	All borrowers	Catalog coverage by country
AccuracyLoss	Group	Exo. right	Consumer	All lenders	Accuracy loss due to fairness objective is fairly distributed across protected groups of users.
RiskTolerance	Individual	Fitness	Consumer	All lenders	Riskier loans are recommended to users with greater risk tolerance

Table 2: Potential fairness concerns and their logics.

Putting all of these dimensions together gives us a three-dimensional ontology of fairness concerns in recommendation: fairness logic, consumer- vs provider-side, and group vs individual target. Table 2 illustrates a range of different fairness concerns that are speculatively derived from the microlending context. This list illustrates a number of the points relative to fairness concerns raised so far. We can see that all four of Moulin’s fairness logics are represented. We also see that the fairness concerns can be group or individual: for example, we are attentive to individual qualities in the **RiskTolerance** concern, but group outcomes in **LargeAmt**. The **AccuracyLoss** concern is a consumer-side concern, relevant

to lenders, but other concerns are on the provider side. We also see that it is possible for a single objective, here the geographic diversity of loan recommendation, to be represented by multiple fairness concerns: **LowCountry** and **AllCountry**. In spite of having the same target, these concerns are distinguished because they approach the objective from different logics and evaluate outcomes differently.

## 4.1 Fairness Agents

Our architecture SCRUF-D (Social Choice for Recommendation Under Fairness – Dynamic) [10] builds on the SCRUF architecture introduced in [9, 34]. It is designed to allow multiple fairness concerns to operate simultaneously in a recommendation context. Fairness concerns, derived from stakeholder consultation, are instantiated in the form of fairness agents, each having three capabilities:

**Evaluation:** A fairness agent can evaluate whether the current historical state is fair, relative to its particular concern. Without loss of generality, we assume that this capability is represented by a function  $m_i$  for each agent  $i$  that takes as input a history of the system’s actions and returns an number in the range  $[0, 1]$  where 1 is maximally fair and 0 is totally unfair, relative to the particular concern.

**Compatibility:** A fairness agent can evaluate whether a given recommendation context represents a good opportunity for its associated items to be promoted. We assume that each agent  $i$  is equipped with a function  $c_i$  that can evaluate a user profile  $\omega$  and associated information and return a value in the range  $[0, 1]$  where 1 indicates the most compatible user and context and 0, the least.

**Preference:** An agent can compute a preference for a given item whose presence on a recommendation list would contribute (or not) to its particular fairness concern. Again, without loss of generality, we assume this preference can be realized by a function that accepts an item as input and returns a preference score in  $\mathbb{R}_+$  where a larger value indicates that an item is more preferred.<sup>3</sup>

## 4.2 Recommendation Process

We assume a recommendation generation process that happens over a number of time steps  $t$  as individual users arrive and recommendations are generated on demand. Users arrive at the system one at a time, receive recommendations, act on them (or not), and then depart. When a user arrives, a recommendation process produces a recommendation list  $\ell_s$  that represents the system’s best representation of the items of interest to that user, generated through whatever recommendation mechanism is available. We do not make any assumptions about this process, except that it is focused on the user and represents their preferences. A wide variety of recommendation techniques are well studied in the literature, including matrix factorization, neural embeddings, graph-based techniques, and others.

The first step to incorporating fairness into the recommendation process is to determine which fairness concerns / agents will be active in responding to a given recommendation opportunity. This is the *allocation phase* of the process, the output of which is a set of non-negative weights  $\beta$ , summing to one, over the set of fairness agents, indicating to what extent each fairness agent is considered to be allocated to the current opportunity.

Once the set of fairness agents have been allocated, they have the opportunity to participate in the next phase of the process, which is the *choice phase*. In this phase, all of the active (non-zero weighted) agents and their weights participate in producing a final list of recommendations for the user. We view the recommender system itself as being an agent that participates in this phase.

---

<sup>3</sup>A more complex preference scenario is one in which agents have preferences over entire lists rather than individual items. We plan to consider such preference functions in future work.

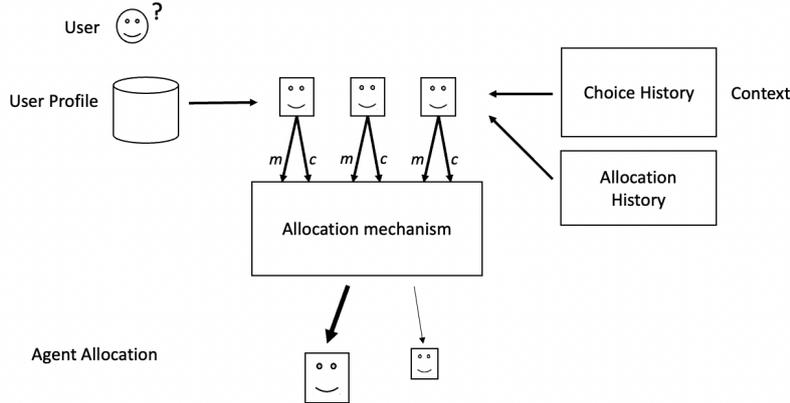


Figure 1: SCRUF-D Framework / Allocation Phase: Recommendation opportunities are allocated to fairness concerns based on the context.

## 5 The SCRUF-D Architecture

The two phases of the SCRUF-D architecture are detailed in Figures 1 and 2. The original SCRUF framework [34] concentrated on the representation of user preferences, as computed by the recommender system, and fairness concerns, as derived from stakeholder consultation as discussed in Section 4.1, and their integration. SCRUF-D incorporates the history of system decisions and the fairness achieved over time to control the allocation of fairness concerns. We will first provide a high level overview of the system and describe each figure in detail with formal notation: Table 3 provides a reference to this notation.

### 5.1 Overview

We can think of a recommender system as a two-sided market in which the recommendation opportunities that arise from the arrival of a user  $u \in \mathcal{U}$  to the system, and each are allocated to a set of items  $v \in \mathcal{V}$  from the system’s catalog. This market has some similarities to various forms of online matching markets including food banks [2], kidney allocation [23, 3], and ride sharing [14], in that users have preferences over the items; however, in our case this preference is known only indirectly through either the prior interaction history or a recommendation function. Additionally, the items are not consumable or rivalrous. For example, a loan can be recommended to any number of users – it is not “used up” in the recommendation interaction.<sup>4</sup> Also, users are not bound to the recommendations provided; in most systems including Kiva, there are multiple ways to find items of which the recommender system is only one.

Once we have a collection of fairness agents we must solve two interrelated problems: (1) what agent(s) are allocated to a particular recommendation *opportunity* and (2) how do we *balance* between the allocated agents and the user’s individual preferences?

Figure 1 shows the first phase of this process, allocation [6], in which we decide which fairness concerns / agents should be allocated to a particular fairness opportunity. This is an online and dynamic allocation problem where we must consider many factors including the history of agent allocations so far, the generated lists from past interactions with users, and how fair the set of agents believes this history to be. As described in Section 4.1, agents take these histories and information about the current user profile and calculate two values:  $m$ , a measure of fairness relative to their agent-specific concern, and  $c$ , a measure of compatibility between the current context and the agent’s fairness concern. The allocation

<sup>4</sup>Loans on Kiva’s platform may be exhausted eventually through being funded, but many other objects of recommendation such as streaming media assets are effectively infinitely available.

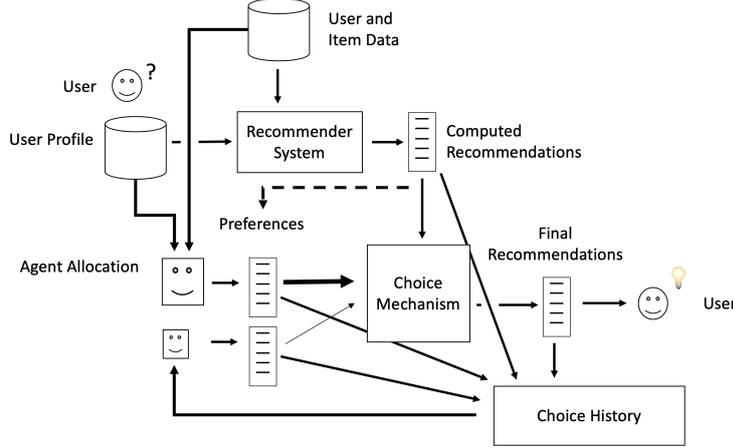


Figure 2: SCRUF-D Framework / Choice Phase: The preferences derived from the recommender system and the fairness concerns are integrated by the choice mechanism.

mechanism takes these metrics into account producing a probability distribution over the fairness agents that we call the *agent allocation*, which can be interpreted as weights in the choice stage or to select a single agent via a lottery, e.g., a randomized allocation [7].

In the second phase, shown in Figure 2, the recommender system generates a list of options that represents the user’s preferences. The fairness concerns generate their own preferences as well. These preferences may be global in character, i.e., preferences over all items, in which case they may be independent of what the recommender system produces; we call this a recommendation function below. Or, as indicated by the dashed line, these preferences may be scoped only over the items that the recommender system has generated; named a scoring function. In either case, the preference function of the fairness agent, like the one for the user, generates a list of items and scores. The choice mechanism combines these preferences of both the user and fairness agents, along with the allocation weights of the fairness agents, to arrive at a final recommendation list to be delivered to the user. The list and interactions with it become a new addition to the choice history.

## 5.2 Formal Description

In our formalization of a recommendation system setting we have a set of users  $\mathcal{U} = \{u_1, \dots, u_n\}$  and a set of items  $\mathcal{V} = \{v_1, \dots, v_m\}$ . For each item  $v_i \in \mathcal{V}$  we have a  $k$ -dimensional feature vector  $\phi = \langle \phi_1, \dots, \phi_k \rangle$  over a set of categorical features  $\phi$ , each with finite domain. Some of these features may be sensitive, e.g., they are associated with one or more fairness agent concerns, we denote this set as  $\phi^s$ . Without loss of generality, we assume that all elements in  $\mathcal{V}$  share the same set of features  $\phi$ . Finally, we assume that each user is associated with a profile of attributes  $\omega = \langle \omega_1, \dots, \omega_j \rangle$ , of which some also may be sensitive  $\omega^s \subseteq \omega$ , e.g., they are associated with one or more fairness agents.

We make the standard assumption that we have (one or more) recommendation mechanisms that take a user profile  $\omega$  and a (set of) items  $v$  and produces a predicted rating  $\hat{r} \in \mathbb{R}_+$ . We will often refer to a recommendation list,  $\ell = \langle \{v_1, \hat{r}_1\}, \dots, \{v_i, \hat{r}_i\} \rangle$ , which is generated for user  $\omega$  by sorting according to  $\hat{r}$ , i.e.,  $\text{sort}(\mathcal{R}_i(\omega, \mathcal{V})) \rightarrow \ell$ . Note that this produces a permutation (ranking) over the set of items for that user, i.e. a recommendation. As a practical matter, the recommendation results will almost always contain a subset of the total set of items, typically the head (prefix) of the permutation up to some cutoff number of items or score value.

Rec. System	$\mathcal{U}(u)$ $\mathcal{V}(v)$ $\phi = \langle \phi_1, \dots, \phi_k \rangle$ $\omega = \langle \omega_1, \dots, \omega_j \rangle$ $\phi^s \subseteq \phi$ $\omega^s \subseteq \omega$ $\mathcal{R}_i(\omega, v) \rightarrow \{v, \hat{r}\}$ $\ell = \langle \{v_1, \hat{r}_1\}, \dots, \{v_i, \hat{r}_i\} \rangle$ $sort(\mathcal{R}_i(\omega, \mathcal{V})) \rightarrow \ell$	Users (user). Items (item). Item Features. User Profile. Sensitive Item Features as a subset of all item features $\phi$ . Sensitive Aspects of User Profile as a subset of all user profile features $\omega$ . Recommendation mechanism that takes a user profile $\omega$ and a (set of) items $v$ and produces a predicted rating $\hat{r} \in \mathbb{R}_+$ . Recommendation List as an ordered list of item, predicted rating pairs. Recommendation List for user $\omega$ sorted by $\hat{r}$ .
Fairness Agents	$\mathcal{F} = \{f_1, \dots, f_i\}$ $f_i = \{m_i, c_i, \mathcal{R}_i\}$ $m_i(\vec{L}, \vec{H}) \rightarrow [0, 1]$ $c_i(\omega) \rightarrow [0, 1]$ $\mathcal{R}_i(\omega, v) \rightarrow \{v, \hat{r}\}$ $\mathcal{R}_i(\ell, \omega, v) \rightarrow \{v, \hat{r}\}$ $\ell_{\mathcal{F}} = \{\mathcal{R}_1(\omega, \mathcal{V}), \dots, \mathcal{R}_i(\omega, \mathcal{V})\}$	Set of Fairness Agents. Fairness agent $i$ defined by a fairness metric $m_i$ , a compatibility metric $c_i$ , and a ranking function $\mathcal{R}_i$ . Fairness metric for agent $i$ that takes a choice history $\vec{L}$ and allocation history $\vec{H}$ and produces a value in $[0, 1]$ according to the agent's evaluation of how fair recommendations so far have been. Compatibility metric for agent $i$ that takes a particular user profile $\omega$ and produces a value in $[0, 1]$ for how compatible fairness agent $i$ believes they are for user $\omega$ . Note: The compatibility metric combines preferences on the agent side and those on the user side (inferred from the profile). If these preferences are symmetrical, we have a one-sided matching problem, but two-sided cases are also possible. Fairness Agent Recommendation function. Fairness Agent Scoring function. Set of Fairness Agent Recommendation Lists indexed by fairness agent label $i$ .
Allocation	$\mathcal{A}(\mathcal{F}, m_{\mathcal{F}}(\vec{L}, \vec{H}), c_{\mathcal{F}}(\omega)) \rightarrow \beta \in \mathbb{R}_+^{ \mathcal{F} }$ $\vec{H} = \langle \beta^1, \dots, \beta^t \rangle$	Allocation mechanism $\mathcal{A}$ that takes a set of fairness agents $\mathcal{F}$ , the agents' fairness metric evaluations $m_{\mathcal{F}}(\vec{L}, \vec{H})$ , and the agents' compatibility metric evaluations $c_{\mathcal{F}}(\omega)$ and maps to an agent allocation $\beta$ . Allocation History $\vec{H}$ that is an ordered list of agent allocations $\mathcal{A}$ at time $t$ .
Choice	$\mathcal{C}(\ell, \beta, \ell_{\mathcal{F}}) \rightarrow \ell_C$ $\vec{L} = \langle \ell^t, \ell_{\mathcal{F}}^t, \ell_C^t \rangle$	Choice Function is a function from a recommendation list $\ell$ , agent allocation $\beta$ , and fairness agent recommendation list(s) $\ell_{\mathcal{F}}$ to a combined output list $\ell_C$ . Choice History that is an ordered list of user recommendation list $\ell$ , agent recommendation list(s) $\ell_{\mathcal{F}}$ , and choice function output lists $\ell_C$ , indexed by time step $t$ .

Table 3: Notations for our formal description of the SCRUF-D architecture.

In the SCRUF-D architecture, fairness concerns map directly onto agents  $\mathcal{F} = \{f_1, \dots, f_i\}$ . In order for the agents to be able to evaluate their particular concerns, they take account of the current state of the system and voice their evaluation of how fairly the overall system is currently operating, their compatibility for the current recommendation opportunity, and their preference for how to make the outcomes more fair. Hence, each fairness agent  $i$  is described as a set,  $f_i = \{m_i, c_i, \mathcal{R}_i\}$  consisting of a fairness metric,  $m_i(\vec{L}, \vec{H}) \rightarrow [0, 1]$ , that takes a choice history  $\vec{L}$  and allocation history  $\vec{H}$  and produces a value in  $[0, 1]$  according to the agent's evaluation of how fair recommendations so far have been; a compatibility metric,  $c_i(\omega) \rightarrow [0, 1]$ , that takes a particular user profile  $\omega$  and produces a value in  $[0, 1]$  for how compatible fairness agent  $i$  believes they are for user  $\omega$ ; and a ranking function,  $\mathcal{R}_i(\omega, v) \rightarrow \{v, \hat{r}\}$ , that gives the fairness agent preferences.

In the allocation phase (Figure 1), we must allocate a set of fairness agents to a recommendation opportunity. Formally, this is an allocation function,  $\mathcal{A}(\mathcal{F}, m_{\mathcal{F}}(\vec{L}, \vec{H}), c_{\mathcal{F}}(\omega)) \rightarrow \beta \in \mathbb{R}_+^{|\mathcal{F}|}$  that takes a set of fairness agents  $\mathcal{F}$ , the agents' fairness metric evaluations  $m_{\mathcal{F}}(\vec{L}, \vec{H})$ , and the agents' compatibility metric evaluations  $c_{\mathcal{F}}(\omega)$  and maps to an agent

allocation  $\beta$ , where  $\beta$  is a probability distribution over the agents  $\mathcal{F}$ . The allocation function itself is allocating fairness agents to recommendation opportunities by considering both the fairness metric for each agent as well as each fairness agent’s estimation of compatibility.

The allocation function can take many forms, e.g., it could be a simple function of which ever agent voices the most unfairness in the recent history [34], or a more complex function from social choice theory such as the probabilistic serial mechanism [5] or other fair division or allocation mechanisms. Note here that the allocation mechanisms is directly comparing the agent valuations of both the current system fairness and compatibility. We implicitly assume that the agent fairness evaluations are comparable. While this is a somewhat strong assumption, it is less strong than assuming that fairness and other metrics, e.g., utility or revenue, are comparable as is common in the literature [42]. So, although we are assuming different voicings of fairness are comparable, we are only assuming that fairness is comparable with fairness, and not other aspects of the system. We plan to explore options for the allocation function in our empirical experiments. We track the outputs of this function as the allocation history,  $\vec{H} = \langle \beta^1, \dots, \beta^t \rangle$ , an ordered list of agent allocations  $\beta$  at time  $t$ .

In the second phase of the system (Figure 2), we take the set of allocated agents and combine their preferences (and weights) with those of the current user  $\omega$ . To do this we define a choice function,  $\mathcal{C}(\ell, \beta, \ell_{\mathcal{F}}) \rightarrow \ell_{\mathcal{C}}$ , as a function from a recommendation list  $\ell$ , agent allocation  $\beta$ , and fairness agent recommendation list(s)  $\ell_{\mathcal{F}}$  to a combined list  $\ell_{\mathcal{C}}$ . Each of the fairness agents is able to express their preferences over the set of items for a particular user,  $\mathcal{R}_i(\omega, v) \rightarrow \{v, \hat{r}\}$ , and we take this set of lists,  $\ell_{\mathcal{F}} = \{\mathcal{R}_1(\omega, \mathcal{V}), \dots, \mathcal{R}_i(\omega, \mathcal{V})\}$ , as input to the choice function that generates a final recommendation that is shown to the user,  $\ell_{\mathcal{C}}$ .

We again leave this choice function unspecified as this formulation provides a large design space: we could use a simple voting rule, a simple additive utility function or something much more complicated like rankings over the set of all rankings [6]. Note that the choice function can use the agent allocation  $\beta$  as either a lottery to, e.g., select one agent to voice their fairness concerns, or as a weighting scheme. We will investigate a range of choice functions in further research. In order for the fairness agents to be able to evaluate the status of the system we also track the choice history,  $\vec{L} = \langle \ell^t, \ell_{\mathcal{F}}^t, \ell_{\mathcal{C}}^t \rangle$ , as an ordered list of user recommendation list  $\ell$ , agent recommendation list(s)  $\ell_{\mathcal{F}}$ , and choice function output lists  $\ell_{\mathcal{C}}$ , indexed by time step  $t$ .

## 6 Conclusion and Future Work

We have introduced the SCRUF-D architecture for integrating multiple fairness concerns into recommendation generation leveraging social choice. The design is general and allows for many different types of fairness concerns—involving multiple fairness logics and encompassing both provider and consumer aspects of the recommendation platform. Our experiments with simple synthetic data show that the SCRUF-D architecture is capable of representing and applying multiple fairness concerns in a modular and agent-based way and balancing among them dynamically. Thorough empirical evaluation of the architecture with real data and fairness concerns is a subject for future work, as is the incorporation and study of a full range of allocation and choice mechanisms.

Future work will proceed in multiple research arcs. One arc of future work is to apply the architecture in more realistic settings, particularly with Kiva. We are working with Kiva stakeholders and beginning the process of identifying fairness concerns. In the meantime, we also plan to conduct additional experiments with a variety of off-line data sets, exploring a range of different fairness concern formalizations and social choice options. We have made the mechanisms and the agents fairly simple by design. Further experimentation will show how effective this structure is for maintaining fairness over time and allowing a wide variety of fairness concerns to be expressed. However, there are some areas of exploration that we can anticipate and are discussed in the appendix.

**Acknowledgements** We are indebted to our Kiva collaborators and the leadership team there for making this collaboration possible. This publication is based upon research supported by the National Science Foundation under grants #2107577 and #2107505.

## References

- [1] Amanda Aird, Paresha Farastu, Joshua Sun, Amy Volda, Nicholas Mattei, and Robin Burke. Dynamic fairness-aware recommendation through multi-agent social choice. *CoRR*, abs/2303.00968, 2023. doi: 10.48550/arXiv.2303.00968. URL <https://doi.org/10.48550/arXiv.2303.00968>.
- [2] M. Aleksandrov, H. Aziz, S. Gaspers, and T. Walsh. Online fair division: Analysing a food bank problem. In *Proc. 24th International Joint Conference on AI (IJCAI)*, pages 2540–2546. IJCAI, 2015.
- [3] P. Awasthi and T. Sandholm. Online stochastic optimization in the large: Application to kidney exchange. In *Proc. 21st International Joint Conference on AI (IJCAI)*, pages 405–411. IJCAI, 2009.
- [4] Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *California law review*, 104(3):671, 2016. ISSN 0008-1221. doi: 10.15779/Z38BG31. URL <https://scholarship.law.berkeley.edu/californialawreview/vol104/iss3/2>.
- [5] Anna Bogomolnaia and Hervé Moulin. A new solution to the random assignment problem. *Journal of Economic theory*, 100(2):295–328, 2001.
- [6] F. Brandt, V. Conitzer, U. Endriss, J. Lang, and A. D. Procaccia, editors. *Handbook of Computational Social Choice*. Cambridge University Press, 2016.
- [7] Eric Budish, Yeon-Koo Che, Fuhito Kojima, and Paul Milgrom. Designing random allocation mechanisms: Theory and applications. *American economic review*, 103(2): 585–623, 2013.
- [8] Robin Burke. Multisided Fairness for Recommendation. In *Workshop on Fairness, Accountability and Transparency in Machine Learning (FATML)*, Halifax, Nova Scotia, 2017. URL <https://arxiv.org/abs/1707.00093>.
- [9] Robin Burke, Amy Volda, Nicholas Mattei, and Nasim Sonboli. Algorithmic fairness, institutional logics, and social choice. In *Harvard CRCS Workshop on AI for Social Good at 29th International Joint Conference on Artificial Intelligence (IJCAI 2020)*, 2020.
- [10] Robin Burke, Nicholas Mattei, Vladislav Grozin, Amy Volda, and Nasim Sonboli. Multi-agent social choice for dynamic fairness-aware recommendation. In *Adjunct Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization*, pages 234–244, 2022.
- [11] Robin Burke, Pradeep Ragothaman, Nicholas Mattei, Brian Kimmig, Amy Volda, Nasim Sonboli, Anushka Kathait, and Melissa Fabros. A performance-preserving fairness intervention for adaptive microfinance recommendation. In *2nd Workshop on Online and Adaptive Recommender Systems (OARS) at KDD 2022*, 2022.
- [12] Nicolo Cesa-Bianchi, Yoav Freund, David Haussler, David P Helmbold, Robert E Schapire, and Manfred K Warmuth. How to use expert advice. *Journal of the ACM (JACM)*, 44(3):427–485, 1997.

- [13] Yuga J Cohler, John K Lai, David C Parkes, and Ariel D Procaccia. Optimal envy-free cake cutting. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- [14] John P Dickerson, Karthik A Sankararaman, Aravind Srinivasan, and Pan Xu. Allocation problems in ride sharing platforms: Online matching with offline reusable resources. In *Proc. Thirty-Second AAAI Conference on Artificial Intelligence (AAAI)*, pages 1007–1014. AAAI, 2018.
- [15] Benjamin Edelman, Michael Ostrovsky, and Michael Schwarz. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *The American economic review*, 97(1):242–259, 2007.
- [16] Michael D. Ekstrand, Anubrata Das, Robin Burke, and Fernando Diaz. Fairness in information access systems, 2022.
- [17] Paresha Farastu, Nicholas Mattei, and Robin Burke. Who pays? personalization, bossiness and the cost of fairness. *arXiv preprint arXiv:2209.04043*, 2022.
- [18] Rupert Freeman, Seyed Majid Zahedi, and Vincent Conitzer. Fair social choice in dynamic settings. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 4580–4587, Marina del Rey, CA, 2017. International Joint Conferences on Artificial Intelligence.
- [19] Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. The (im)possibility of fairness: different value systems require different mechanisms for fair decision making. *Communications of the ACM*, 64(4):136–143, April 2021. ISSN 0001-0782, 1557-7317. doi: 10.1145/3433949. URL <https://dl.acm.org/doi/10.1145/3433949>.
- [20] Yingqiang Ge, Shuchang Liu, Ruoyuan Gao, Yikun Xian, Yunqi Li, Xiangyu Zhao, Changhua Pei, Fei Sun, Junfeng Ge, Wenwu Ou, et al. Towards long-term fairness in recommendation. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 445–453, New York, 2021. ACM.
- [21] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness, 2018.
- [22] Min Kyung Lee, Daniel Kusbit, Anson Kahng, Ji Tae Kim, Xinran Yuan, Allissa Chan, Daniel See, Ritesh Noothigattu, Siheon Lee, and Alexandros Psomas. Webuildai: Participatory framework for algorithmic governance. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–35, 2019.
- [23] N. Mattei, A. Saffidine, and T. Walsh. An axiomatic and empirical analysis of mechanisms for online organ matching. In *Proceedings of the 7th International Workshop on Computational Social Choice (COMSOC)*, 2018.
- [24] Rishabh Mehrotra, Niannan Xue, and Mounia Lalmas. Bandit based optimization of multiple objectives on a music streaming platform. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3224–3233, 2020.
- [25] Marco Morik, Ashudeep Singh, Jessica Hong, and Thorsten Joachims. Controlling fairness and bias in dynamic learning-to-rank. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 429–438, New York, 2020. ACM.

- [26] Hervé Moulin. *Fair division and collective welfare*. MIT press, 2004.
- [27] Deirdre K Mulligan, Joshua A Kroll, Nitin Kohli, and Richmond Y Wong. This thing called fairness: disciplinary confusion realizing a value in technology. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–36, 2019.
- [28] Safiya Umoja Noble. *Algorithms of Oppression: How search engines reinforce racism*. NYU Press, 2018.
- [29] Cathy O’Neil. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Broadway Books, 2016.
- [30] Szilvia Pápai. Strategyproof assignment by hierarchical exchange. *Econometrica*, 68(6):1403–1433, 2000.
- [31] Gourab K Patro, Arpita Biswas, Niloy Ganguly, Krishna P Gummadi, and Abhijnan Chakraborty. Fairrec: Two-sided fairness for personalized recommendations in two-sided platforms. In *Proceedings of The Web Conference 2020*, pages 1194–1204, 2020.
- [32] Claudia Perlich, Brian Dalessandro, Rod Hook, Ori Stitelman, Troy Raeder, and Foster Provost. Bid optimizing and inventory scoring in targeted online advertising. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 804–812, 2012.
- [33] Jessie J. Smith, Anas Buhayh, Anushka Kathait, Pradeep Ragothaman, Nicholas Mattei, Robin Burke, and Amy Volda. Exploring the institutional logics of multistakeholder microlending recommendation. In *Proc. of the 2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT)*, 2023.
- [34] Nasim Sonboli, Robin Burke, Nicholas Mattei, Farzad Eskandarian, and Tian Gao. ”and the winner is...”: Dynamic lotteries for multi-group fairness-aware recommendation, 2020.
- [35] Nasim Sonboli, Farzad Eskandarian, Robin Burke, Weiwen Liu, and Bamshad Mobasher. Opportunistic multi-aspect fairness through personalized re-ranking. In *Proceedings of the 28th ACM Conference on User Modeling, Adaptation and Personalization*, UMAP ’20, page 239–247, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450368612. doi: 10.1145/3340631.3394846. URL <https://doi.org/10.1145/3340631.3394846>.
- [36] Tom Sühr, Asia J Biega, Meike Zehlike, Krishna P Gummadi, and Abhijnan Chakraborty. Two-sided fairness for repeated matchings in two-sided markets: A case study of a ride-hailing platform. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3082–3092, 2019.
- [37] William Thomson. Fair allocation rules. In *Handbook of Social Choice and Welfare*, volume 2, pages 393–506. Elsevier, 2011.
- [38] Jun Wang, Weinan Zhang, and Shuai Yuan. Display advertising with real-time bidding (rtb) and behavioural targeting, 2017.
- [39] Shuai Yuan, Ahmad Zainal Abidin, Marc Sloan, and Jun Wang. Internet advertising: An interplay among advertisers, online publishers, ad exchanges and web users, 2012.
- [40] Shuai Yuan, Jun Wang, and Xiaoxue Zhao. Real-time bidding for online advertising: measurement and analysis. In *Proceedings of the Seventh International Workshop on Data Mining for Online Advertising*, page 3. ACM, 2013.

- [41] Weinan Zhang, Shuai Yuan, and Jun Wang. Optimal real-time bidding for display advertising. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1077–1086. ACM, 2014.
- [42] Ziwei Zhu, Xia Hu, and James Caverlee. Fairness-aware tensor-based recommendation. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 1153–1162, 2018.

Amanda Aird,  
University of Colorado: Boulder  
Boulder, CO  
Email: [amanda.aird@colorado.edu](mailto:amanda.aird@colorado.edu)

Paresha Farastu,  
University of Colorado: Boulder  
Boulder, CO  
Email: [amanda.aird@colorado.edu](mailto:amanda.aird@colorado.edu)

Joshua Sun,  
University of Colorado: Boulder  
Boulder, CO  
Email: [amanda.aird@colorado.edu](mailto:amanda.aird@colorado.edu)

Amy Volda,  
University of Colorado: Boulder  
Boulder, CO  
Email: [amanda.aird@colorado.edu](mailto:amanda.aird@colorado.edu)

Nicholas Mattei  
Tulane University  
New Orleans, LA  
Email: [nsmattei@tulane.edu](mailto:nsmattei@tulane.edu)

Robin Burke,  
University of Colorado: Boulder  
Boulder, CO  
Email: [amanda.aird@colorado.edu](mailto:amanda.aird@colorado.edu)

## A Design Considerations

Within this framework there are a number of important design considerations to take into account for any particular instantiation of the SCRUF-D architecture. We have left many of the particular design choices open for future investigation. We allow for any type of recommendation algorithm; fairness agents may incorporate any type of compatibility function or fairness evaluation function. Similarly, we do not constrain the allocation or choice mechanisms. With SCRUF-D, we are able to explore many definitions of fairness and recommendation together in a principled uniform way. In this section, we discuss a few of the design parameters that may be explored in future work.

### A.1 Agent Design

We can expect that an agent associated with a fairness concern will typically have preferences that order items relative to a particular feature or features associated with that concern. Items more closely related to the sphere of concern will be ranked more highly and those unrelated, lower. However, this property means that agents associated with different concerns might have quite different rankings – the gender parity concern will rank women’s loans highly regardless of their geography, for example. Thus, we cannot assume consistency or single-peakedness across the different agents.

As noted above, agents may have preferences over disjoint sets of items or they may be constrained only to have preferences over the items produced by the recommender system for the given user. This second option corresponds to a commonly-used *re-ranking* approach, where the personalization aspect of the system controls what items can be considered for recommendation and fairness considerations re-order the list [16]. If an agent can introduce any item into its preferences, then we may have the challenge in the choice phase of integrating items that are ranked by some agents but not others. Some practical work-arounds might include a constraint on the recommender system to always return a minimum number of items of interest to the allocated agents or a default score to assign to items not otherwise ranked.

Despite our terminology, it is clear that our architecture as described is sufficiently general that an agent could be designed that pushes the system to act in harmful and unfair ways rather than beneficial and fairness-enhancing ones. Thus, the importance of the initial step of stakeholder consultation and the careful crafting of fairness concerns. Because fairness concerns are developed within a single organization and with beneficence in mind, we assume that we do not need to protect against adversarial behavior, such as collusion among agents or strategic manipulation of preferences. The fact that the agents are all “on the same team” allows us to avoid constraints and complexities that otherwise arise in multi-agent decision contexts.

### A.2 Agent Efficacy

The ability of an agent to address its associated fairness concern in non-deterministic. It is possible that the agent may be allocated to a particular user interaction, but its associated fairness metric may still fail to improve. One likely reason for this is the primacy of the personalization objective. Generally, we expect that the user’s interests will have the greatest weight in the final recommendations delivered. Otherwise, the system might have unacceptably low accuracy, and fail in its primary information access objective.

One design decision therefore is whether (and how) to track agent efficacy as part of the system history. If the agent’s efficacy is generally low, then opportunities to which it is suited become particularly valuable; they are the rare situations in which this fairness goal can be addressed. Another aspect of efficacy is that relationships among item characteristics may mean that a given agent, while targeted to a specific fairness concern, might

have the effect of enhancing multiple dimensions of fairness at once. Consider a situation in which geographic concerns and sectoral concerns intersect. Promoting an under-served region might also promote an under-served economic sector. Thus, the empirically-observed multidimensional impact of a fairness concern will need to be tracked to represent its efficacy.

Efficacy may also be a function of internal parameters of the agent itself. A separate learning mechanism could then be deployed to optimize these parameters on the basis of allocation, choice and user interaction outcomes.

### A.3 Mechanism Inputs

Different SCRUF implementations may differ in what aspects of the context are known to the allocation and/or choice mechanisms. Our hope is that we can leverage social choice functions in order to limit the complexity of the information that must be passed to the allocation and/or choice mechanisms. However, if a sophisticated and dynamic representation of agent efficacy is required, it may be necessary to implement a bandit-type mechanism to explore the space of allocation probabilities and/or agent parameters as discussed above. Recent research on multidimensional bandit learning suggests possible approaches here [24].

### A.4 Agent Priority

As we have shown, agent priority in the allocation phase may be a function of user interests, considering different users as different opportunities to pursue fairness goals. It may also be a function of the history of prior allocations, or the state of the fairness concerns relative to some fairness metric we are trying to optimize. As the efficacy consideration would indicate, merely tracking allocation frequency is probably insufficient and it is necessary to tie agent priority to the state of fairness. Allocation priority is also tied to efficacy as noted above. It may be necessary to compute expected fairness impact across all dimensions in order to optimize the allocation.

We plan to leverage aspects of social choice theory to help ameliorate some of these issues. There is a significant body of research on allocation and fair division mechanisms that provide a range of desirable normative properties including envy-freeness [13], e.g., the guarantee that one agent will not desire another agent’s allocation, Pareto optimally, e.g., that agents receive an allocation that is highly desirable according to their compatibility evaluations [5]. An important and exciting direction for research is understanding what allocation properties can be guaranteed for the SCRUF-D architecture overall depending on the allocation mechanism selected [6].

We note that in most practical settings the personalization goal of the system will be most important and therefore the preference of this agent will have topmost priority. It is always allocated and is not part of the allocation mechanism. Thus, we cannot assume that the preference lists of the agents that are input to the choice system are anonymous, a common assumption in the social choice literature on voting [6].

### A.5 Bossiness

Depending on how the concept of agent / user compatibility is implemented, it may provide benefits to *bossy* users, those with very narrow majoritarian interests that do not allow for the support of the system’s fairness concerns. Those users get results that are maximally personalized and do not share in any of the potential accuracy losses associated with satisfying the system’s fairness objectives. Other, more tolerant users, bear these costs. A system may wish to ensure that all users contribute, at some minimal level, to the fairness goals. In social choice theory, a mechanism is said to be non-bossy if an agent cannot change the allocation without changing the allocation that they receive by modifying their preferences [30]. Some preliminary discussions of this problem specifically for fairness-aware recommendation appear in [17].

## A.6 Fairness Types

We concentrate in this paper and our work with Kiva generally on provider-side group fairness, that is characteristics of loans where protected groups can be distinguished. However, it is also possible to use the framework for other fairness requirements. On the provider side, an individual fairness concern is one that tracks individual item exposure as opposed to the group as a whole. It would have a more complex means of assessing preference over items and of assessing fairness state, but still fits within the framework.

Consumer-side fairness can also be implemented through use of the compatibility function associated with each agent. For example, the example of assigning risk appropriately based on user risk tolerance becomes a matter of having a risk reduction agent that reports higher compatibility for users with lower risk tolerance.

## B Experimental Results: Methodology and Results

As an initial examination of the properties of the SCRUF-D architecture, we conducted a series of experiments with simulated data and agents of generic design. The experiments were run on a Python implementation of the SCRUF-D architecture. See associated GitHub repository for the source code.<sup>5</sup> Configuration files, data and Jupyter notebooks for producing the experiments and visualizations below are found in a separate repository<sup>6</sup>.

### B.1 Data generation

We generated 1000 users and 100 items with two associated sensitive features:  $\phi_1$  and  $\phi_2$ . The first 333 users of the dataset, the type 1 users, are compatible only with Agent 1, which is concerned about feature  $\phi_1$ . The next 333 users, type 2, are only compatible with Agent 2 and  $\phi_2$ . And the last group of users, type 3, are not compatible with either agent. For these experiments, we used binary compatibility values. Each item was assigned either  $\phi_1$ ,  $\phi_2$ , or neither. Items that were matched with neither are considered to be items without sensitive features.  $\phi_1$  and  $\phi_2$  were each randomly assigned to items 25% of the time. For this set of items, no items were assigned both sensitive features.

Recommendations are generated by scoring items randomly, sorting the scored items, and adding noise to the scores after taking the top 50 items.

### B.2 Agent design

The agents both evaluate fairness in the same way. The recommendation lists from the previous algorithm iterations are combined. The fraction of this combined list that consists of the item-specific associated items is computed. If the fraction is equal or greater than a specified proportion, then the metric returns 1 (maximum fairness). If the list contains a smaller fraction, the fairness is scaled smoothly from 1 to 0, which is the score when no protected items appear. In our experiments, we set this proportion to 0.75 for Agent 1 and 0.5 for Agent 2. This makes Agent 1 more demanding in terms of its definition of fairness.

### B.3 Mechanisms

We examine several different allocation mechanisms:

- **Least misery:** The fairness agent with the lowest fairness score  $m_i$  is chosen.
- **Most compatible:** The fairness agent most compatible  $c_i$  with the current user is chosen.

---

<sup>5</sup>[https://github.com/that-recsys-lab/scruf\\_d](https://github.com/that-recsys-lab/scruf_d)

<sup>6</sup>[https://github.com/that-recsys-lab/scruf\\_tors\\_2023](https://github.com/that-recsys-lab/scruf_tors_2023)

- **Static lottery:** A single fairness agent is chosen with a fixed probability at each time step. In our experiments, the probabilities were 0.5 for Agent 1 and 0.3 for Agent 2. (Note that this difference is in line with their different fairness requirements.)
- **Dynamic lottery:** A lottery is constructed with probabilities proportional to agent unfairness and a single agent is chosen by drawing from this lottery.
- **Weighted fairness allocation:** All agents are allocated but their weight is determined by their current unfairness.

For reasons of space, we do not report on experiments with different choice mechanisms as part of this study. We use a simple weighted voting mechanism, analogous to a weighted Borda score. After the fairness agent is allocated, the scores from the recommender systems are adjusted such that each item among the  $v^p$  protected items has its score augmented by the constant  $\delta = 0.5$ . Where multiple agents are allocated,  $\delta$  is scaled by the allocation weight. Then the recommendations are sorted by score and the top ten items chosen as the recommendation list. This choice of  $\delta$  is designed to be fairly impactful in that allocating an agent has a good chance of boosting its fairness to the maximum level.

## B.4 Evaluation

With this synthetic data set, we do not have ground truth user preferences and so we do not evaluate recommendation accuracy. We leave the exploration of the fairness / accuracy tradeoff for future work. We concentrate in these experiments in examining the interaction between the agents, the allocation mechanisms, and the fairness outcomes that result.

In the results below, we use three different plots to demonstrate the dynamics of the simulation with varied allocation mechanisms. We plot fairness as evaluated by each agent ( $m_{\mathcal{F}}$ ) at each point in time. We also plot the allocation vector  $\vec{H}$  (in cumulative form) to show at what time points different agents are being allocated.

We also plot cumulative *fairness regret* over the course of the experiment. At each time step, we calculate  $1 - m_i$ , that is the difference from perfect fairness as the agent defines it, and then sum these values over the course of the simulation. This is similar to the notion of regret in reinforcement learning but using fairness instead of utility. Fairness regret  $G_i$  for agent  $i$  is defined as:

$$G_i(s) = \sum_{t=0}^s 1 - m_i(\vec{L}_t, \vec{H}_t) \quad (1)$$

## C Results

As noted above, the results here do not include an accuracy measure and so we are not reporting on how the interaction between fairness and accuracy is managed by different allocation mechanisms. Our results here concentrate on the management of the dynamic aspects of fairness. How does each mechanism handle the time-varying aspects of fairness and the probabilistic nature of re-ranking to achieve fairness?

In Figure 3, we report the fairness metrics associated with each agent as different recommendation lists are generated. No agents are allocated and no re-ranking is performed so this is showing the baseline characteristics of the input data.

As we can see in Figure 3a, the fairness measures range around 0.6 to 0.2, with Agent 1 having the lower fairness because it requires more protected group items to meet its target proportion. In some sense this is an “easy” re-ranking problem because, even without doing anything, the agents are getting some protected group items in each list. The cumulative value of these metrics is shown in Figure 3b and we can see that over the course of the 1000

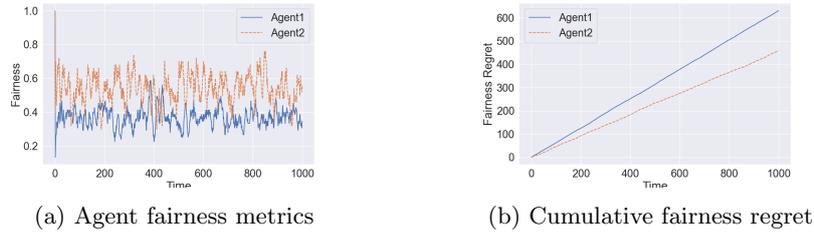


Figure 3: Results for baseline data (no allocation)

users / time steps in the experiment, Agent 1 has a regret around 600 and Agent 2 around 450.

### C.1 Allocation algorithms

In the results that follow, we show similar output for each allocation algorithm including a plot of the cumulative agent allocations.

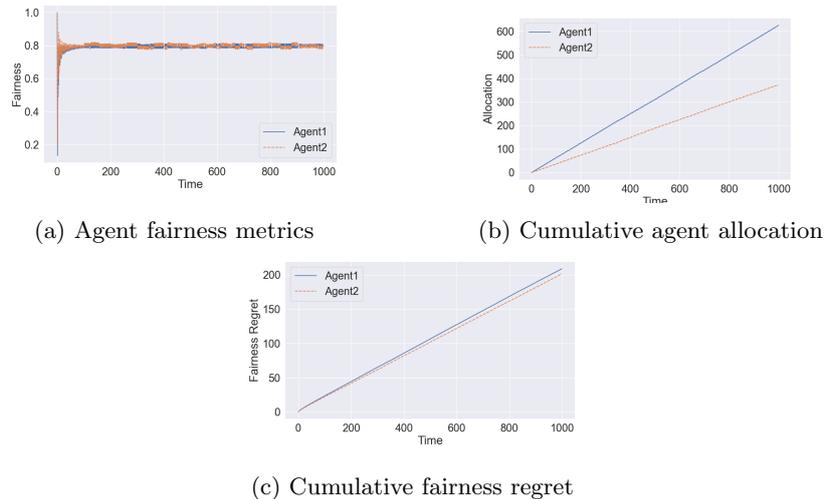


Figure 4: Results for **Least Misery** allocation

Figures 4 and 5 show the simulation results for the two simple deterministic allocation mechanisms: **Least Misery** and **Most Compatible**. As might be expected, the Least Misery algorithm causes the system to bounce around between the two agents. As soon as one gets some benefit from being allocated, the other gets a chance. Although it is not captured here, this is an inefficient strategy because it does not take into account the compatibility between users and agents and would have lower utility than a mechanism sensitive to that aspect.

**Most Compatible** is at the other extreme, looking only at compatibility values. The system swings between items favorable to Agent 1 (compatible with the initial third of the users) to those favorable to Agent 2. Since the last third of the users are not compatible with either agent, neither of them is allocated and so no re-ranking occurs for these users.

An alternative to these deterministic mechanisms is to allocate a single agent by lottery. Figure 5 shows the **Static Lottery** case where the agents are assigned a fixed probability of selection: Agent 1, 0.5 and Agent 2, 0.3. Here the choice of agent does not depend on the

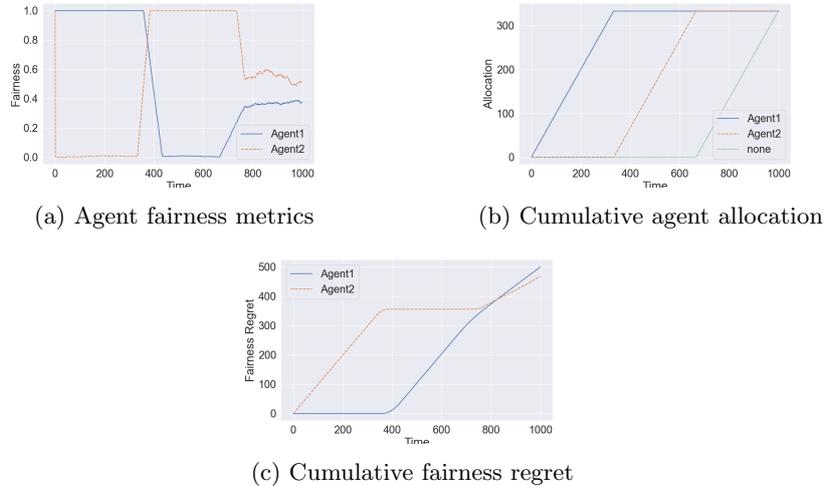


Figure 5: Results for **Most Compatible** allocation

context and so the fairness scores move randomly but in a much broader band than seen in the **Least Misery** case. The allocations are distributed uniformly as the lottery would indicate and in the end, the two agents have similar regret because of their different fairness proportions.

The dynamic aspects of the SCRUF-D architecture show up when we allow the lottery to adapt to the current state of fairness. A dynamic lottery mechanism is shown in Figure 7. Because the system can respond to variations in the input, the fairness values stay in a tighter range than with the static version. Agent 1 ends up with a higher regret, in spite of being allocated more often, which is not surprising given it is “pickier”.

The last allocation mechanism is one where every agent are allocated at every time step but with different weights. In our Borda choice mechanism, the weights are used to control the weight of each agent’s vote in the choice phase. The weights in **Dynamic Weighted** allocation are proportional to the agent’s fairness metric:  $m_i$ . These results can be seen in Figure 8 and are similar to the **Dynamic Lottery** although the gap between the agents is smaller.

## C.2 Comparison

These results are summed up in Figure 9. The **Weighted Allocation** and the **Fairness Lottery** are very similar in keeping a relatively high average fairness for the two agents, with the lottery having a narrower distribution. These two mechanisms also have the lowest regret. The much simpler **Least Misery** allocation does quite well in this simulated setting. **Most Compatible** does poorly but that is not surprising since it is not optimizing for fairness and that is the only metric here. As noted above, the compatibility measure is meant to ensure that agents are allocated to users with consonant preferences.

## C.3 Efficacy

The issue of agent efficacy is illustrated in Figure 10. To examine this issue, we added a third agent whose protected items are much more rare: 5% instead of 25% for the other agents. We also made the  $\delta$  for re-ranking 0.1, much lower than the other agents. The agents’ interactions were managed using the **Dynamic Lottery** mechanism, which constantly adjusts its agent allocation to favor the agents with low fairness.

The result is that allocating this agent has low efficacy: it doesn’t improve the fairness

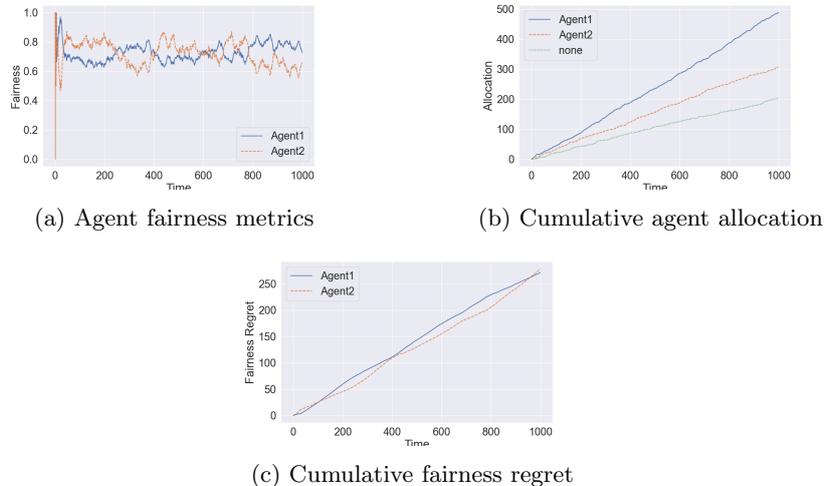


Figure 6: Results for **Static Lottery** allocation

outcomes as much as allocating the other agents. As Figure 10 shows, the system falls into a pattern of allocating Agent 3 frequently, without any benefit, while Agents 1 and 2 also suffer. So the overall fairness is worse than, for example, a static lottery. Note that a weighted allocation would have the same problem as more and more weight would be placed on Agent 3’s preferences even though little is accomplished by doing so.

As discussed above, understanding agent efficacy requires looking at the impact of allocation: how much does fairness for an agent improve if it is allocated? With this information, the system can look forwards as well as backwards and reason about how quickly an agent’s objective can be achieved. It also creates the possibility of adapting the re-ranking process itself. For example, a system capable of adapting the  $\delta$  value used in the choice mechanism could escape the trap of fruitless allocation by making the allocation of Agent 3 more impactful. A system could also try to improve its estimate of agent compatibility to find better allocation opportunities, etc. These ideas will be explored in future work.

## D Additional Future Work

A key feature of the recommendation context is that the decisions of the recommender system only influence the exposure of protected items. There is no guarantee that a given user will show any interest in an item just because it is presented. In some settings and for some fairness concerns, exposure might be enough. But in cases where utility derives from usage rather than exposure, there would be some value in having the system learn about the relationship between exposure and utility. This setting has the attributes of a multi-objective bandit learning problem [24], where the fairness concerns represent different classes of rewards and the allocation of agents represents different choices.

Even when we consider exposure as our main outcome of interest, it is still the case that the allocation of different agents may result in differential improvements in fairness, the efficacy problem noted above. Perhaps the items associated with one agent are more common in recommendation lists and can be easily promoted through re-ranking while other agents’ items are not. The weight associated with the allocation of agents may need to be adjusted to reflect the expected utility of allocation, and this expected utility would need to be learned.

The current architecture does not make any assumptions about the distribution of user

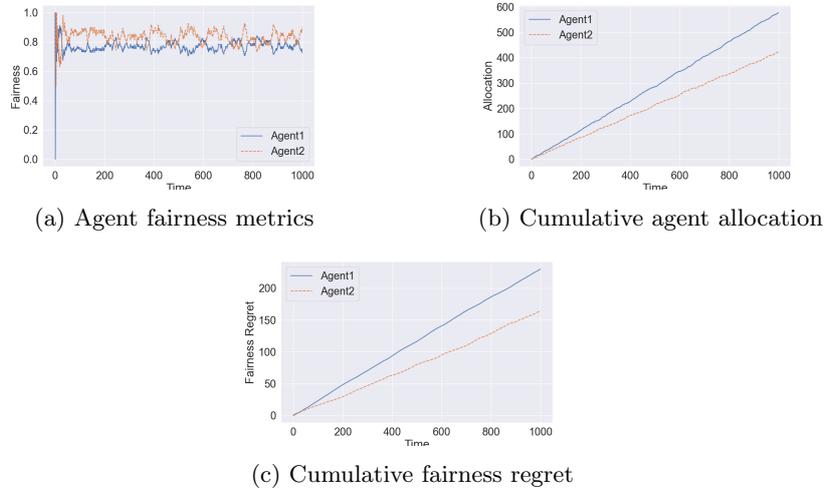


Figure 7: Results for **Fairness Lottery** allocation

characteristics. That is, suppose fairness concern  $f_i$  is “difficult” to achieve in that users with an interest in related items appear rarely. In that case, we should probably allocate  $f_i$  whenever a compatible user arrives, regardless of the state of the fairness metrics. This example suggests that the allocation mechanism could be adapted to look forward (to the distribution of future opportunities) as well as backwards (over fairness results achieved). This would require a model of opportunities similar to [32], and others studied in computational advertising settings.

The current architecture envisions fairness primarily in the context of group fairness expressed over recommendation outcomes. We believe that the architecture will support other types of fairness with additional enhancements. For example, a representational fairness concern would be incompatible with the assumption that fairness can be aggregated over multiple recommendation lists. Consider the examples in Noble’s *Algorithms of Oppression*: it would not be acceptable for a recommender system to deliver results that reinforced racist or sexist stereotypes at times, even if those results were balanced out at other times in some overall average. Representational fairness imposes a stricter constraint than those considered here, effectively requiring that the associated concern be allocated for every recommendation opportunity.

As noted above, the model expressed here assumes that fairness agents have preferences only over items. But it is also possible to represent agents as having preferences over recommendation lists. This would allow agents to express preferences for combinations of items: for example, a preference that there be at least two Agriculture loans in the top 5 items of the list. This kind of preference cannot be expressed simply in terms of scores associated with items. Agents would naturally have to become more complex in their ability to reason about and generate such preferences, and the choice mechanism would become more like a combinatorial optimization problem. It is possible that we can characterize useful subclasses of the permutation space and avoid the full complexity of arbitrary preferences over subsets.

Another interesting direction for research is more theoretical in nature. Much of the research in social choice focuses on providing guaranteed normative properties of various mechanisms. However, the models used in traditional social choice theory do not take into consideration the dynamics of recommender systems as most mechanisms are designed to work in one-off scenarios without dynamic aspects. One direction would be to formulate the

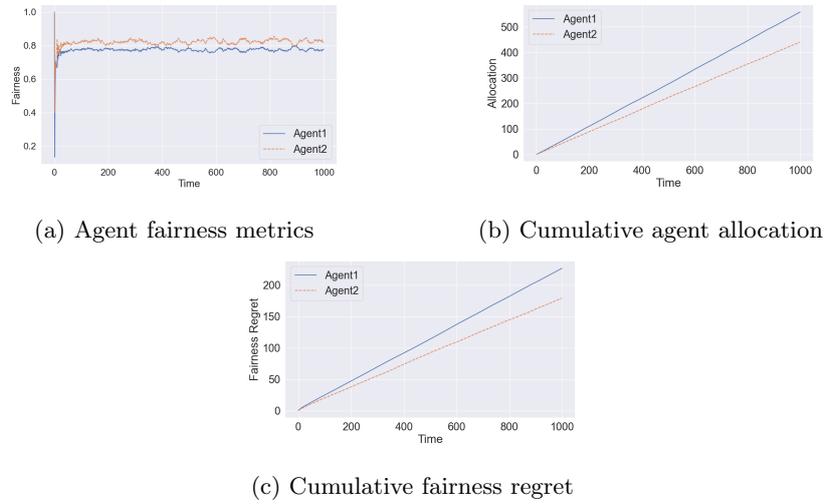


Figure 8: Results for **Dynamic Weighted** allocation

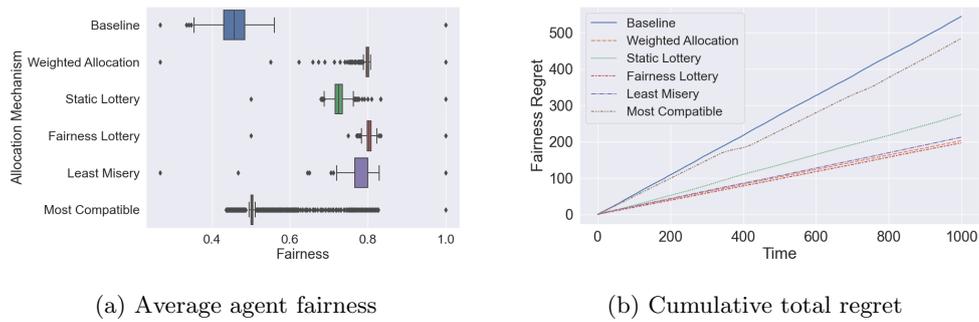
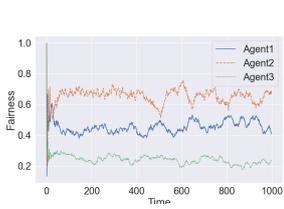
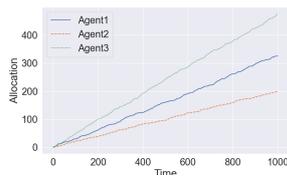


Figure 9: Comparative results

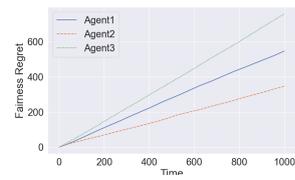
allocation phase of the architecture as an online matching problem, where fairness agents represent one side of the matching and users arrive online on the other side, revealing their compatibility metric. Similar to work in online ad allocation, each fairness agent might have some budget or capacity that limits the number of users they are matched with, in order to balance between various fairness concerns. It will be important to understand the properties of existing social choice mechanisms for allocation and choice when deployed in these dynamic contexts and to develop new methods with good properties.



(a) Agent fairness metrics



(b) Cumulative agent allocation



(c) Cumulative fairness regret

Figure 10: Results for **Dyamic Lottery** allocation, 3 agent condition.