# T-RECS: A Simulation Tool to Study the Societal Impact of Recommender Systems

Eli Lucherini
Princeton University
Princeton, NJ, USA
elucherini@cs.princeton.edu

Matthew Sun
Princeton University
Princeton, NJ, USA
mdsun@princeton.edu

Amy Winecoff
Princeton University
Princeton, NJ, USA
aw0934@princeton.edu

Arvind Narayanan
Princeton University
Princeton, NJ, USA
arvindn@cs.princeton.edu

## ABSTRACT

Simulation has emerged as a popular method to study the long-term societal consequences of recommender systems. This approach allows researchers to specify their theoretical model explicitly and observe the evolution of system-level outcomes over time. However, performing simulation-based studies often requires researchers to build their own simulation environments from the ground up, which creates a high barrier to entry, introduces room for implementation error, and makes it difficult to disentangle whether observed outcomes are due to the model or the implementation.

We introduce T-RECS[1], an open-sourced Python package designed for researchers to simulate recommendation systems and other types of sociotechnical systems in which an algorithm mediates the interactions between multiple stakeholders, such as users and content creators. To demonstrate the flexibility of T-RECS, we perform a replication of two prior simulation-based research on sociotechnical systems. We additionally show how T-RECS can be used to generate novel insights with minimal overhead. Our tool promotes reproducibility in this area of research, provides a unified language for simulating sociotechnical systems, and removes the friction of implementing simulations from scratch.

## 1 INTRODUCTION

Recommender systems in social media platforms such as Facebook and Twitter have been criticized due to the risks they might pose to society. For example, "filter bubbles" [58] have been associated with the emergence of echo chambers leading to degraded political discourse online. Similarly, there is evidence that false news might spread faster and wider online than the truth [81], with the phenomenon potentially exacerbated by recommendation algorithms [37]. YouTube has come under fire for its potential links to the radicalization of certain audiences [79] such as young voters in Brazil [24]. However, there is no consensus on the scope of these concerns or ways to remedy them. For example, several studies have argued that the contribution of algorithmic systems to echo chambers is minimal or nonexistent [6, 25, 32].

Because phenomena such as filter bubbles arise through repeated system interactions over time, methods that assess the system at a single time point provide minimal insight into the mechanisms behind them. In contrast, simulations can model how system elements such as user, items, and algorithms interact over arbitrarily long timescales. As a result, simulation has proved to be a valuable tool in assessing the impact of recommendation systems on the content users consume and on society [4, 14, 16, 28, 60, 84]. Most simulation studies of algorithmic systems have relied upon ad-hoc systems for implementation, which presents several challenges for the advancement of scientific understanding of the effects of algorithmic system dynamics.

We present T-RECS (Tools for RECommender system Simulation), an open-source, unified simulation tool designed to enable investigations of emerging complex phenomena caused by millions of individual actions and interactions in algorithmic systems including filter bubbles, political polarization, and (mis)information diffusion.

T-RECS can mitigate problems associated with ad-hoc systems in several ways. First, one of the most time and labor intensive components of developing simulations is the engineering effort necessary to build the system. T-RECS has been developed to drastically reduce the engineering effort needed to develop a recommender system simulator. This allows researchers to shift their focus from the mechanics of the simulations to the scientific assumptions behind them. As a result, T-RECS can accelerate the pace of development and can facilitate high-quality simulation studies. Second, because many system elements in T-RECS have been designed with appropriate checks, its use can reduce the likelihood of software bugs that either slow the research process or contribute to erroneous results. Third, T-RECS provides a simple modular programming interface and terminology intuitive to researchers in the social sciences and computer science.

T-RECS is an easy-to-use but powerful framework; using the pre-populated library of common recommender system models and other system elements, researchers can instantiate a simulation using three lines of code. For instance, the following code runs a simulation with 1,000 users and 10,000 items whose interactions are mediated by the system default content-filtering algorithm (see Section 4). The last line gathers default measurements.

```
recsys = trecs.models.ContentFiltering(num_users=1000,
                                       num_items=10000)
recsys.run(timesteps=100)
measurements = recsys.get_measurements()
```

---

[1]https://github.com/elucherini/t-recs

With T-RECS, researchers with expertise in social science but limited engineering expertise can still leverage simulation to answer important questions about the societal effects of algorithmic systems. Yet, the framework also provides flexibility for expert users to build upon. For example, researchers can easily specify each user and item's representation in the system, along with new custom measurements.

Applying the same tool to different problems (1) promotes reproducibility by allowing researchers to easily share their simulations; (2) provides a common language to describe problems in the literature; and (3) fosters discovery of principles that apply across seemingly different problems.

In addition to social science and algorithmic researchers, T-RECS can also inform the work of engineers, who can compare algorithmic design choices, as well as policymakers, who can use the results to regulate potentially harmful phenomena.

## 2 BACKGROUND AND RELATED WORK

T-RECS is a tool for evaluating recommendation systems and other algorithmic information systems, especially their societal consequences, based on agent-based modeling. We briefly review the relevant background.

### 2.1 Recommender systems

Recommender systems are responsible for a large part of the content we see and interact with online. For example, in 2018 YouTube reported that 70% of the views were derived from recommendations [70]. In 2019, over half of the more than 1 billion daily active users on Instagram turned to Explore, the recommendation-driven section of Instagram, at least once a month [42]; and numbers from 2016 suggest that 80% of the hours streamed by users on Netflix were driven by its recommender system [31]. Given how frequently recommendations influence our day-to-day lives, understanding how recommendation algorithms influence both short-term and long-term outcomes for users is critical.

The algorithms undergirding recommender systems can take advantage of a variety of information to serve recommendations to users. Collaborative filtering is a subset of recommendation algorithms that leverages patterns in existing user interaction data to generate predictions for new items. Methods such as user-based or item-based collaborative filtering use a nearest-neighborhood approach whereas matrix factorization collaborative filtering uses latent representations of users and items to make predictions. In contrast to wholly interaction-based methods, content-based algorithms take advantage of item or user meta-data (e.g., genre, director for movies). More recently, Bayesian approaches (e.g., [62]), reinforcement learning (e.g., [85]), and deep learning (e.g., [15, 82, 86]) that draw from methods in machine learning more generally have also been applied to recommendation problems. Hybrid and ensemble approaches leveraging multiple recommendation models are often used to mitigate weaknesses associated with a single algorithm type [2].

Recommender systems are typically optimized for predictive accuracy [36], but often include additional metrics such as diversity and novelty of content [40, 87] and, more recently, fairness-related metrics [13, 21, 45].

### 2.2 Recommender systems and society

Recommender systems help users find content that better matches their short or long-term preferences, and to find it more efficiently. They also help users discover and develop new interests. From the platform's perspective, recommender systems can influence user behavior in service of the platform's goals (such as increasing user engagement), improve the visibility of new products, and provide a mechanism for learning more about users' preferences [43].

On the other hand, there are a number of potentially harmful effects. **Filter bubbles** are the result of recommender systems presenting a user with content similar to the user's history, resulting in intellectual isolation, political **polarization**, and **echo chambers** [58, 72–74]. A more extreme version is **radicalization**, a feedback loop in which the algorithm encourages users to consume more and more extreme content, gradually nudging them to embrace increasingly radical ideas [55].

Recommender systems may exhibit a bias toward popular items, resulting in **homogenization** of user behavior and a **concentration** of the market for content in the hands of a small number of creators [14, 16, 66].

Finally, recommender systems may contribute to the spread of misinformation online. These phenomena are inter-related: for example, misinformation may be more likely to spread in the presence of echo chambers because users are less likely to encounter corrective information.

Although we have characterized these potential harms as effects of recommender systems, they may also raise from other types of online information systems including search [30] and social networks [8]. In the shift from offline to online information seeking, users are freed from physical-world constraints such as word-of-mouth recommendation networks or the limited content available through a newspaper subscription. Online, users can much more easily find content tailored to their preferences, beliefs, and ideology.

Our understanding of the long-term impact of recommender systems on users and society is still nascent, and there is little consensus in the literature. For example, quantitative research on the effects of algorithmic systems has found little evidence of filter bubbles [6, 25, 32], in contrast to other types of inquiry such as ethnography [30, 48, 78]. Even studies with the same broad methodology produce conflicting findings, such as empirical studies of misinformation [26, 81] and simulation studies of echo chambers [4, 28].

Many potential factors could explain these differences. Echo chambers or other effects may exist but not be caused by the recommender system, or the effects may only become apparent over time due to feedback loops and not be revealed by cross-sectional studies, or there may be differences between platforms. A further wrinkle related to platform differences is that some studies cannot be reproduced externally because the data is proprietary.

### 2.3 Agent-based modeling

Agent-based models help researchers understand the properties of a complex system by modeling of the individual agents that comprise the system and the interactions between them [9]. One of the earliest Agent-Based Models (ABMs) in the social sciences is Thomas Schelling's work on the dynamics of racial segregation

about fifty years ago [67]. ABMs became widespread only after the rise in the availability computational power of the 1990s [5, 23, 50]. ABMs are widely used in areas of research such as economics and finance, ecology, biology, and epidemiology [35].

## 2.4 Content creators

Academic research on recommender systems have often focused on user-centric experiences and outcomes [1]. However, recommender systems often serve multiple stakeholders including platform providers and content creators. Increasingly, researchers have recognized the importance of a multistakeholder perspective for developing and evaluating recommender systems [1, 12, 54]. Indeed, content creators themselves have raised questions about how recommendation algorithms affect their livelihoods or entrench societal biases [51]. Moreover, content creators represent a dynamic set of agents, and they may shift and adapt in response to algorithmically-mediated interactions with users [7, 17].

So far, most simulation-based research on the system-level effects of algorithmic system features have also focused on users and have likewise failed to account for either the dynamic effects of content creators or the consequences of system properties for content creators.

## 3 THE CASE FOR T-RECS

As noted in Section 2, there are still many open questions on the long-term impact of recommender systems on users, content providers, and society at large. T-RECS offers a unified tool to contextualize and understand seemingly contradictory results, such as those from the literature on filter bubbles. The design of T-RECS facilitates the implementation of related problems using the same simulation framework.

Similarly, T-RECS allows studying multistakeholder problems, such as the complex interactions between users and content creators when they are mediated by an algorithm that tends to suggest more and more radicalized content.

In this section, we explain what simulation, in particular agent-based modeling, can offer. Finally, we argue in favor of unified tools such as T-RECS.

## 3.1 The case for simulation

Observational and experimental methods have been widely used to study the impact of recommender systems [3, 15, 38, 57]. Contrary to observational methods, simulation provides the tools for generalizable and causal discoveries; furthermore, models are typically not affected by ethical issues that arise when experimenting on real users. Furthermore, while methods like offline evaluation on historical data (e.g., the MovieLens dataset [33]) [13, 21, 45] are useful to evaluate algorithms on fixed or unknown user strategies, simulation provides an additional degree of control on evaluating systems on varying user behaviors. Finally, contrary to analytical approaches, simulation supports complex environments and focuses on the dynamics of the system instead of static points of equilibrium.

## 3.2 The case for Agent-Based Modeling

Agent-based modeling is useful to study emergent phenomena: those that arise from the interactions of individual agents, often in counter-intuitive ways. Bonabeau gives the example of a traffic jam, noting that it may move in the opposite direction from the cars that are causing it [9]. He presents a list of system properties that may give rise to emergent phenomena, such as adaptive agents and network effects, all of which are present in recommender systems.

Recently, ABMs have been explicitly used to describe interactions mediated by a recommender system. For example, Geschke et al. [28] modeled an online information diffusion network with users and pieces of information as the agents; similarly, Nasrinpour et al. [56] studied message propagation on Facebook with an ABM. For more examples, see Table 1.

ABMs are especially apt at discovering indirect and unintended effects of design choices because the outcomes of simulations emerge from interactions. Therefore, researchers do not need to encode their expectations into the initial assumptions [20].

Agent-based modeling can express geographical and social distance [20]. Thus, they are suitable to model *adjacency networks* between users, which are often associated with recommender systems. For example, recommender systems with a community-based component leverage connections between users and their friends to provide recommendations [64].

## 3.3 The case for a new tool

To determine whether a new tool can be useful, we surveyed existing simulation studies of algorithmic systems, finding 15 studies. A summary of the studies we reviewed is in Table 1.

We argue in favor of unified tools in this space, such as T-RECS, with the following observations. First, these 15 studies represent recently published usages of simulation to understand societal impact of algorithmic systems. Of these, seven specifically focus on generating knowledge on recommender systems, six focus on online information diffusion, and two focus on opinion dynamics.

However, all of the studies we reviewed involved system components conceptually representing "users," "items," and "algorithms." As a result, we designed T-RECS to also have system elements that could represent these types of entities and enable flexible representations of their properties and behaviors.

Second, our literature review indicates that simulations of algorithmic systems often rely on ad-hoc software, which can be challenging to develop and reproduce. For example, 12 out of the 15 studies in Table 1 were developed independently, including six of the seven studies focusing on recommender systems. Two used NetLogo, a programming language and modeling environment originally developed in 1999 for multi-agent simulations, and one used AnyLogic, which is proprietary GUI-based simulation software available for academic use [10, 76]. The heterogeneity in implementations in the studies points to the need for a unified framework.

The primary benefits of applying the same tool to different problems are both practical and scientific. From the practical perspective, a unified framework will reduce the engineering effort needed to develop a simulation, allowing researchers to shift focus from the mechanics of the simulation to the assumptions behind them. In addition to enabling researchers to focus more on the science of

| Simulation to study sociotechnical systems | | | | |
|---|---|---|---|---|
| Work | Model | Main topic | Tool | ABM? |
| Aridor et al. [4] | RS | Filter bubbles | n.a. | Yes |
| Chaney et al. [14] | RS | Filter bubbles | n.a. | Yes |
| Ciampaglia et al. [16] | RS | Popularity bias | n.a. | No |
| Garimella et al. [27] | ID | Filter bubbles | n.a. | Yes |
| Geschke et al. [28] | RS | Filter bubbles | NetLogo [83] | Yes |
| Goel et al. [29] | ID | Virality | n.a. | Yes |
| Jiang et al. [44] | RS | Filter bubbles | n.a. | No |
| Lee et al. [47] | OD | Perception biases | n.a. | Yes |
| Lim et al. [49] | ID | Digital divide | NetLogo | Yes |
| Nasrinpour et al. [56] | ID | Virality | AnyLogic [10] | Yes |
| Perra and Rocha [60] | OD | Polarization | n.a. | No |
| Sun et al. [71] | RS | Popularity bias | n.a. | Yes |
| Tambuscio et al. [75] | ID | Misinformation | n.a. | Yes |
| Törnberg [77] | ID | Misinformation | n.a. | Yes |
| Yao and Huang [84] | RS | Popularity bias | n.a. | Yes |

Table 1: Simulation literature survey. RS = recommender systems, ID = online information diffusion, OD = online opinion dynamics

simulations, a common framework will also speed up development, facilitating a greater volume of high quality research on algorithmic systems. For example, if much of the up-front engineering effort has already been accomplished through the unified framework, there are fewer new system elements and consequently, fewer places where software bugs could lead to slow downs in development, or worse, erroneous results.

From the scientific perspective, a major issue with ad-hoc systems is that results can be difficult to reproduce. Because many idiosyncrasies in design and implementation are eliminated in a unified system, the likelihood that another research team could reproduce a scientific finding increases. Furthermore, a unified framework will allow researchers to communicate in a common language for different problems. For example, our literature review in Table 1 demonstrates that multiple researchers have conducted simulations on filter bubbles; however, each of these studies has a unique definition for what constitutes a filter bubble and unique metrics for evaluating filter bubble effects. As a result, it is difficult to reconcile results across similar simulations, much less across simulations that address more disparate concepts.

Before T-RECS, there have been other efforts to provide a unified simulation environment for sociotechnical systems [10, 11, 19, 46, 76]. **NetLogo** was designed with a heavy emphasis on visualizing agents in 2D space and also requires users to learn the NetLogo programming language [76]. Furthermore, NetLogo was not optimized for large-scale simulations for hundreds of thousands of users, and complex recommender system algorithms are not straightforward to implement in the NetLogo interface. **AnyLogic** is proprietary software for agent-based modeling developed by the AnyLogic company [10]. Simulations in AnyLogic are created through a graphical user interface, limiting expressiveness and flexibility.

In the past few years, multiple simulation libraries specifically geared towards studying the temporal dynamics of simulated recommender systems have been released. **RecoGym** was the first

reinforcement learning simulation environment for recommendation, in the context of online advertising [65]. The **RecSim** library was designed for a reinforcement learning-based approach to modeling recommender systems [41]. Another follow-up simulation framework, **RecSim NG**, uses probabilistic programming for uncertainty modeling in recommender ecosystems [52]. While these libraries offer useful mathematical tools for reasoning about recommender systems, they are geared towards expert recommender systems practitioners and researchers, requiring knowledge of libraries like Tensorflow, OpenAI Gym, or Edward2, and are designed for a particular technical frame, such as reinforcement learning or probabilistic programming. In contrast, T-RECS is designed to be flexible and simple enough that users familiar with Python and numpy will be able to build and run simulations relatively quickly, promoting accessibility to a broader audience of researchers, from computer scientists to social scientists.

The most related tools to ours are **ML Fairness Gym** and **RecLab** [19, 46]. **ML Fairness Gym** implements the OpenAI Gym API to provide a set of reusable components for studying long-run fairness in algorithmic decisionmaking systems over time. The ML Fairness Gym is framed to have a specific focus on fairness, while our tool is intended to allow practitioners to study a broad range of phenomena. Furthermore, the Fairness Gym does not have robust support for multi-agent simulation, instead focusing on the decision-maker as the primary agent. In contrast, T-RECS allows for dynamic behavior at the levels of users, recommender systems, and content providers. **RecLab**, a recently released Python library for simulating and evaluating recommender systems, is the most similar tool to T-RECS [46]. While RecLab focuses on the evaluation of different recommender systems, T-RECS is designed to make it equally easy to study the effects of different user choice models, content creator behavior, or item distributions. As a result, T-RECS has the expressiveness to model sociotechnical systems

beyond recommender systems, such as information diffusion in social networks.

## 4 DESIGN AND ARCHITECTURE

Figure 1 shows the modules of T-RECS, which include: (1) the **users**, representing the entities interacting with the recommendations: for example, consumers on an e-commerce platform; (2) the **items**, representing the object of the recommendations: for example, the movies in a movie recommendation system; and the **content creators**, providing new items during the simulation; (3) the **model** representing the mechanism through which users interact with items: for example, a collaborative filtering algorithm; and (4) the **metrics** which evaluate the outcomes of the simulation sampled at each simulation step. Metrics also keep store information on selected *internal states* of the system, such as the user preferences over time. Note that these modules are common to many of the simulations in Table 1. We go into more detail about each of the core modules in the subsequent sections.

### 4.1 Simulation Dynamics

We summarize the simulation dynamics of T-RECS, inspired by the model developed by Schmit and Riquelme [68]. We note that this simulation framework works for many algorithmic system that can be modeled as a *user-algorithm feedback loop*. In the following sections, we present more details about how the basic dynamics can be enriched and changed. The numbers in the following list correspond to the numbers highlighted in Figure 1.

At each time step:

(1) The model predicts the user scores for each item in the system. These predicted scores, which can be different from the *actual* user scores, are used by the system to make recommendations to users. By default, the predicted user scores are calculated as the inner product between user preferences and item attributes.

(2) The model presents a set of items to each user. In the general case, the system presents different items to each user, generally chosen considering the user scores predicted in the previous step. At every step, items may also be randomly interleaved into recommendations, simulating user discovery of items outside of algorithmic recommendations, as in Chaney et al. [14]. If runtime creation of items is enabled, the items most recently generated by the content creators will be included in system-generated recommendations based on how the specific recommendation algorithm handles new items. For example, in a popularity-based system, all new items have a popularity at zero and are placed near the bottom of recommendation lists.[2]

(3) Each user gives feedback to the presented items. This means something different for different systems (e.g., in a movie recommender system, users choose movies to watch). Typical user feedback is *implicit* (to simulate a click, generally only on one item). Users decide which item(s) to consume based on their actual preferences.

(4) The system records each user's feedback and updates its internal state. What this means varies with each model. For example, T-RECS' popularity model keeps track of the number of interactions each item has received; our content filtering model keeps track of the attributes each user has interacted with.

(5) The measurement module updates all tracked metrics based on the latest interactions. The default metrics vary by model, but could include an accuracy measure, such as the mean squared error between the predicted user scores and the actual user scores or a variety of other measures of recommendation list diversity, user interaction similarity, etc.

### 4.2 Users

Users in T-RECS are represented by two main elements: their preferences and the score they attribute to each item in the system. With our design, informed by our survey of the literature, we assume that the real user preferences and the preferences predicted by the recommender algorithm are separate. Similarly, recommender algorithms may predict scores for the items in the system that do not perfectly correspond to the real scores attributed by the users. This modeling choice in several studies we analyzed in our survey [4, 14]. Therefore, in T-RECS there are four dedicated data structures: *actual* user preferences, *predicted* user preferences, actual user scores, and predicted user scores.

Additionally, T-RECS provides a simple framework to change how user preferences evolve over time, typically in response to the items they are exposed to and consume. This behavior has been adopted in several of the studies we analyzed in our survey [28, 44].

### 4.3 Items and content creators

Items in T-RECS can be fixed or dynamic. In a system such as a movie recommendation website, it might be useful to consider the catalog of items as fixed, since the changes to it are not necessarily dependent on the response of the audience to it.

However, in many real-world recommender systems the catalog of items is not fixed. Instead, it changes over time as producers, or *content creators*, create and publish new items. For example, a platform such as YouTube might be better modeled with a dynamic catalog of items, as *YouTubers* are as much a part of the ecosystem as viewers. Most importantly, creators are often incentivized to adapt the content they publish in direct response to the viewers' reactions. In T-RECS, we added a module for content creators to approximate this phenomenon. Researchers can define the incentives of the creators and how they change over time.

Regardless of the specific item and content creator configuration, researchers can enforce additional constraints to how users consume items, or to how items are recommended to users. For example, it is possible for users to avoid consuming items they have already seen in the past; similarly, recommender models can avoid recommending items that the users have already interacted with.

Content creators are characterized by the attributes of items they generate, and can be thought of as particular distributions from which items may be sampled. To mirror some of the complexity of content creator behavior in real systems, the attributes of the items content creators generate are sampled probabilistically.

---

[2]T-RECS users may specify custom behavior for how new items are scored. For example, you may specify that all new items do not appear in recommendation lists and instead are randomly interleaved throughout the recommendations, as in Chaney et al. [14].
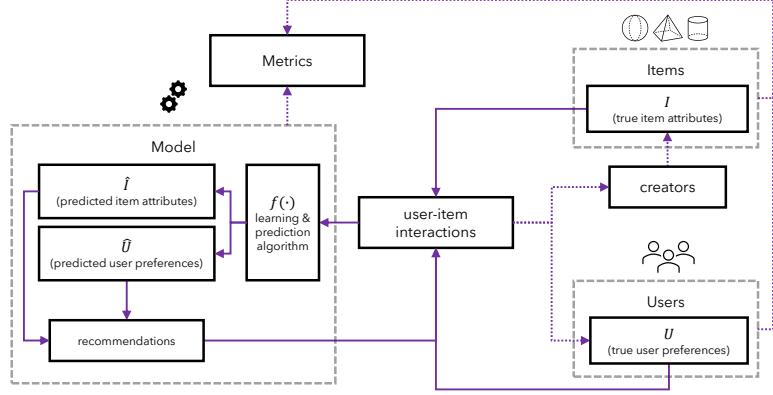
Figure 1: A conceptual diagram indicating the different design components of a T-RECS simulation, and the components that are used as inputs to other components as the simulation progresses. Lines with arrows indicate inputs fed into outputs. Dashed gray lines indicate conceptual components. Dotted arrows indicate possible feedback mechanisms (for example, item interactions affecting content creator distributions).

## 4.4 Model representation

In T-RECS, we make the following decisions regarding recommender system models:

(1) We use an *event-driven* model as opposed to a *time-driven* model–that is, we automatically skip all steps in which the system did not observe any user interactions. This increases the efficiency of the simulation. As a consequence, one simulation step does not correspond to a constant time unit.

(2) User interactions are processed in parallel. By default, at each step, the system receives one piece of feedback from each user, although T-RECS users can override this as necessary. Researchers who want to observe multiple user-item interactions before retraining a model can specify that behavior through the same framework.

Although our framework is particularly well-suited to simulating recommender systems, we find that it is also flexible enough to model a broader set of sociotechnical phenomena. For example, we used T-RECS to model the process of information diffusion in networks in which recommender systems do not necessarily have a role, as long as they follow the user-item-system feedback loop approach. In the context of information diffusion, users come into contact with various pieces of content (items) and probabilistically choose to share content with others in their network. At each step, the system monitors which users have shared content and presents shared content to their neighboring users at the next timestep, keeping track of which users have shared content (and thus will not reshare in the future). process of sharing We offer more details on this topic in Section 5.2.2.

## 4.5 Metrics

In T-RECS, metrics are designed for researchers to gather information about the model and the simulation. They often define the output of the simulation: for example, researchers interested in discovering filter bubbles can measure diversity of recommended

content across users and draw conclusions from the results of these measurements.

Our main contribution in this regard is to provide a mechanism for easy implementation and templates for researchers to use as an example. Our literature survey (Table 1) suggested very little consistency of metrics used across studies. While we could identify overarching themes — the most prominent being distances between attributes and preferences — each of the 14 studies we analyzed used distinct metrics. Therefore, we purposely provided a limited set of metrics with the expectation that researchers will want to implement their own metrics.

## 4.6 Programmer interface

Our programmer interface is beginner-friendly. We took inspiration from popular machine learning library SciKit-Learn [59] in that models can be instantiated and simulations can be run without specifying any argument. T-RECS was developed in Python and makes heavy use of the numpy library [34]. We deliberately chose the default behaviors and provided an interface to easily modify them. For example, instantiating a simulation using the popularity model can be achieved by this line of code:

```
recsys = trecs.models.PopularityRecommender()
```

The line above instantiates a popularity model with 100 users and 1250 items. User and item characteristics are also generated with default values from an arbitrary distribution (see the documentation for more details). We see these as sensible defaults, as the simulation is large enough to see emergent network effects while small enough to pose minimal computational burden; of course, T-RECS users are free to modify these parameters as they see fit.

Although T-RECS was designed to be accessible to beginners, it offers flexibility in changing a wide array of components for more advanced calibration. The customization is not limited to the parameters directly related to users and items. For example, researchers can define new score functions–that is,the functions

calculating actual and predicted user scores. Instead of using the default inner product as the score function, the user can define their own mechanism for calculating user-item scores. As an example, we provide an alternative score function that calculates scores using the cosine similarity between $U$ (or $\hat{U}$) and $I$. We also provide a mechanism to define new score functions when initializing a model.

As with score functions, T-RECS offers an API to define new models, metrics, user feedback behaviors, and more.

## 5 CASE STUDIES

Our goal in developing T-RECS was to create a tool that introduced a common conceptual framework for conducting simulations with different levels of complexity, different methodologies, and different goals, such as those enumerated in Table 1. To illustrate that T-RECS has achieved this aim, we reproduced results from two substantially different simulation-based studies. In the first study, Chaney et al. [14] examined how algorithmic confounding in recommendation systems leads to the homogenization of user behavior. In the second, Goel et al. [29] simulated models of contagion to understand the patterns of online information diffusion on Twitter. Although the two studies model distinct sociotechnical systems, both can be expressed within the T-RECS framework.

Our goal in reproducing these two studies was to perform a *conceptual replication* rather than an *exact replication* [18]. While exact replication aims to perform an identical operationalization of the original test as the original study, a conceptual study examines whether a hypothesis holds up under different operationalizations of the same conceptual variable. Conceptual replications provide evidence as to the robustness and generalizability of the theory. Although our replications are faithful to the original authors' descriptions of their methods, our goal was to show that T-RECS can be used to generate similar insights to those from these studies, rather than to generate wholly identical graphs or statistics. For this reason, we did not seek out source code from the authors of the original studies, but rather, reconstructed the simulation parameters within T-RECS using only the details provided in the published reports.

Both of the studies we replicated were originally implemented via ad-hoc systems. Without a common framework, a researcher interested in reproducing these studies would have to re-implement every element using the methodology presented in the papers, and ideally, would also have to implement all the necessary testing functionality from scratch. In contrast, because both replication studies involved representations for users, items, and models, we were able to use the core components of T-RECS as a starting point to reproduce the results of these prior studies rather than beginning anew. Similarly, we were able to take advantage of the many built-in testing capabilities of T-RECS for debugging. Relying on the core functionality of T-RECS allowed us to concentrate our efforts on faithfully reproducing aspects of our replication studies' designs rather than on engineering.

Our conceptual replication case studies demonstrate the power of T-RECS to accommodate complex and distinct algorithmic simulations; however, another major benefit of our framework is that simulation designs can be easily extended. In our third and final case study, we conduct a *constructive replication*–a replication in which

the findings of a prior study are replicated and then extended to encompass new elements that provide additional scientific insight. Our constructive replication study examines the homogenizing effect on item generation that is induced by the presence of content creators, who adapt over time to user feedback. For this study, we reuse much of the code and many of the assumptions from our replication of Chaney et al.'s research [14]. As a result, we can confidently attribute any differences between the results of the simulations with and without content creators to the introduction of this new system element rather than faulty reproduction of the prior work's assumptions or implementation details.

In the following sections, we describe our motivations, methods, and results in adapting T-RECS to each case study. For more technical details, please see the Appendix. All code necessary to reproduce our experiments can be found online.[3]

### 5.1 Algorithmic confounding (Chaney et al.)

*5.1.1 Background.* At a high level, Chaney et al. [14] illustrate the detrimental effects of algorithmic confounding, which occurs when a recommendation algorithm is trained on user interaction data that is itself influenced by the prior recommendations of the algorithm. The study shows that algorithmic confounding homogenizes user behavior more than what would occur if all users were provided with recommendations that best matched their underlying preferences.

*5.1.2 Motivation.* We chose to replicate the results from Chaney et al. [14] for several reasons. First, this investigation leveraged complex representations for each of the core components of the simulation (i.e., users, items, and models). Although the goal of the original study was not to approximate a real-world system at high fidelity, many of their choices of parameters and representation were based on an understanding of the characteristics of real-world systems (e.g., short head, long tail in item popularity). By illustrating that T-RECS can be modified to manifest this level of nuanced representation, we demonstrate the power of our simple architecture for encapsulating complicated specifications. That is, T-RECS can be used not only for simple toy models, but also for highly complex use cases. Second, this study focused on how the dynamics of one core component (models) could affect the outcome of the system over time through its influence on another component (users). Thus, by replicating this finding in T-RECS, we show that our system can be used to model interactions at the level of entities that ultimately affect the system as a whole. As this is generally the goal of most agent-based simulations, this suggests that our system will be well-suited to most questions of this nature.

*5.1.3 Simulation overview.* We begin with a high-level overview of the user-item-model interaction dynamics. Each user and each item has some "true" representation as a vector of attributes. At each timestep, every user is shown a list of items by the model, where the model attempts to show each user $u$ the set of items which have maximum predicted utility for $u$. Furthermore, at each timestep, "new" items are introduced into the system. These items are interleaved randomly into each user's recommendation list independently for each user. For each item $i$, there is a "true utility" that

---

user $u$ obtains from interacting with that item that is a probabilistic function of $u$ and $i$'s attributes; however, each user has incomplete knowledge of their true utilities (i.e., they have an "educated guess" of how much they will like each item before interacting with it). Each user then interacts with one item on the basis of this imperfect knowledge of the utility that would be gained from each item in the recommendation list, plus an "attention mechanism" that accounts for the item's position in the recommendation list. The interaction data is then fed back into the model, which may or may not use the feedback to update its internal representation of users and items. After interacting with an item once, a user will not interact with that item again during the simulation.

Each model has its own method of internally representing users, items, and the predicted utility that each user $u$ obtains from interacting with each item $i$. For example, the popularity model represents each item $i$ as the number of times that any user has interacted with $i$, and the predicted utility for a given user $u$ and item $i$ is simply the popularity of item $i$ (i.e., all users are recommended the same set of the top-$k$ most popular items, setting aside new item interleaving). We implement and compare results between six different types of models: popularity, content filtering, matrix factorization, social filtering, random, and ideal. The last two models are included for the purpose of comparison: the random model recommends items randomly to each user, while the ideal model recommends items on the basis of the true user-item utilities. For more details about each model in the simulation, see Section A.1.

*5.1.4 Simulation parameters.* The simulation parameters that follow are taken from the methods of Chaney et al. [14]; we describe them here for exposition.

We run each simulation for a total of 100 steps. At each time step, we introduce 10 new items, for a total of 1000 items at the end of the simulation.

To isolate the effects of algorithm confounding, we execute the simulations in *single training* and *repeated training* modes. In both modes, the recommendations operate in *startup mode* at first, meaning they recommend items randomly in order to gather data about user behavior.

In single training mode, the model in each simulation is trained once after 50 *startup* steps. After training, the only items recommended to users are those from the startup period (although users will see new items because of random interleaving). In repeated training mode, the models are trained at each time step after 10 startup steps. These models update their internal representations of users, item, and user-item scores at each training step. Each recommendation model provides recommendations for all items in the system at the time of training.

We differ from Chaney by running 400 trials for each recommendation model in both the single training and repeated training modes, in order to average out random noise. In each trial, the underlying user profiles, item attributes, true user-item utilities, and the user-item utilities known to the user were kept the same for all recommendation models.

*5.1.5 Metrics.* To assess the homogenization of user behavior, Chaney et al. calculate the average Jaccard index of user interactions over pairs of users. For a given pair of users, the Jaccard index is calculated as the size of the set of items both users have
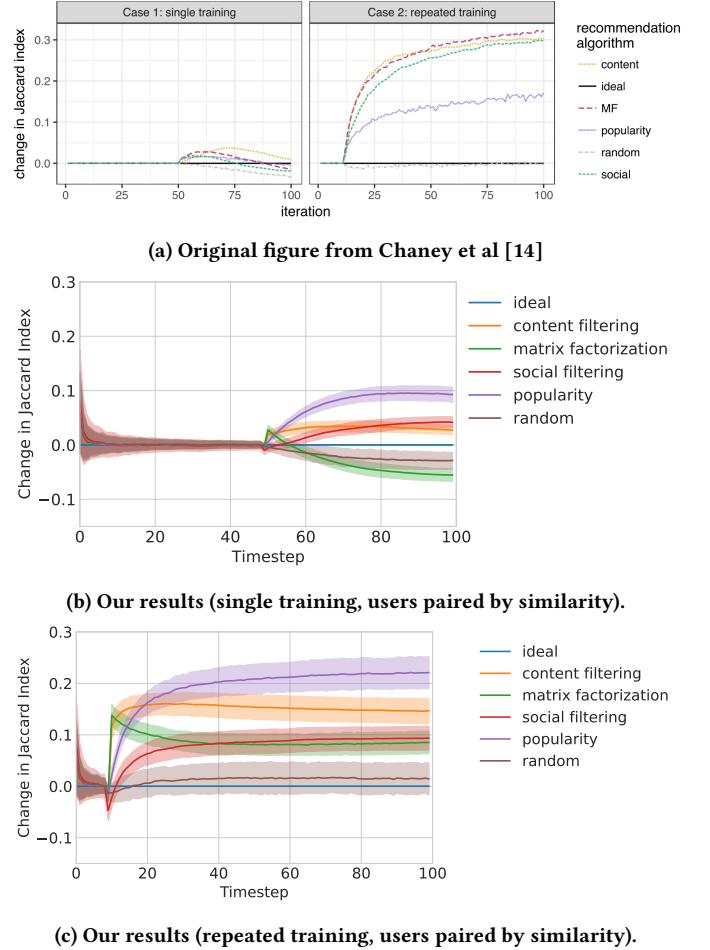


(a) Original figure from Chaney et al [14]



(b) Our results (single training, users paired by similarity).



(c) Our results (repeated training, users paired by similarity).

**Figure 2: Subfigures (b) and (c) show changs in Jaccard index of user behavior, users paired by cosine similarity of predicted user attributes, ±1 SD across runs. We observe mild homogenization in the single training case but observe increased homogenization with repeated training. Compare to subfigure (a), a copy of the original figure from Chaney et al. [14].**

interacted with, divided by the size of the set of items either user has interacted with. To construct pairs of users, the authors pair each user to the user who is most similar according to the recommendation model's internal representation of user preferences. At each timestep, the average Jaccard index is measured relative to the ideal model's average Jaccard index calculated over the same pairs of users to allow for an evaluation of homogenization over and above what would occur if the users were provided with perfect recommendations.

*5.1.6 Results.* Figure 2 shows the evolution of user behavior homogenization under four different recommendation models, where homogenization is measured relative to the ideal recommendation system. We observe results that are qualitatively similar to the

original study. First, our results show that the magnitude of homogenization is much greater when recommendation models are trained repeatedly. Second, in the single training case, most recommendation models exhibit an increase in homogenization directly after the training step, followed by a decrease or slowed increase in homogenization. Third, the random model homogenizes user behavior similarly to the ideal recommendation. Overall, although our results do not match Chaney et al.'s [14] exactly, we emphasize that the takeaways are qualitatively similar even though the implementations (i.e., operationalizations) are slightly different.

## 5.2 Structural virality (Goel et al.)

*5.2.1 Background.* Goel et al. use a simulation-based approach to investigate how well simple theoretical models of social contagion can capture patterns of online information diffusion observed empirically in a dataset of billions of events on Twitter [29]. They find that a relatively simple theoretical model simulated at similar scale recapitulates many of the trends found in their empirical dataset. In particular, they examine patterns of popularity–how often information cascades reach large numbers of users–and structural virality–a metric that distinguishes between viral and broadcast events. They discovered that content that is popular on Twitter is not always "viral," as was often assumed. Instead, content commonly achieved popularity by being reshared by accounts with large followings.

*5.2.2 Motivation.* We chose to replicate Goel et al.'s study for two primary reasons. First, we designed T-RECS to be general enough to accommodate algorithmic systems other than recommendation systems. Goel et al. do not explicitly attempt to simulate a recommender system, but instead use a simulation-based approach to better understand the possible mechanisms for the diffusion of online content. By replicating their findings, we demonstrate that our library is powerful enough to provide a common framework for many different types of simulation-based research. Second, Goel et al. study outcomes that emerge from user interactions *at scale*. Our replication effort simulates outcomes on networks of 1,000,000 users, indicating that researchers can use T-RECS to evaluate large-scale phenomena.

*5.2.3 Simulation framework.* Goel et al.'s models are all based on the susceptible-infected-recovered (SIR) framework, a model of contagion that is often used to model social diffusion. In SIR models, individuals are *infected* (i.e., share a piece of content), and then subsequently infect their susceptible contacts independently with probability $\beta$. After being infected, individuals *recover* and are no longer susceptible to infection, nor do they infect their contacts. On a given graph with average node degree $\bar{k}$, an item with infection probability $\beta$ has the "basic reproduction number" $r = \bar{k}\beta$.

Similar to our case study of Chaney et al. [14], we describe the parameters of our simulations; unless otherwise stated, the parameters are the same as described in Goel et al. [29]. In our replication, we generated scale-free networks of 1,000,000 users. Goel et al. use scale-free networks of size 25,000,000; we use smaller networks (albeit on the same order of magnitude) due to space and compute constraints. The degree of each graph is determined by $\alpha$, the parameter for a power law sequence. In accordance with

Goel et al., we analyzed graphs for the following values of $\alpha$: 2.1, 2.3, 2.5, 2.7, and 2.9. For each value of $\alpha$, we generated 25 1M-node graphs .[4] Following Goel et al., in each simulation, we randomly selected a "seed user" to be infected with a particular item that, on a given graph, has a reproduction value $r$. We tested the following values of $r$: 0.1, 0.3, 0.5, 0.7, and 0.9. In total, we ran approximately 2,100,000 simulations across all values of $\alpha$ and $r$. Subsequently, we measured the probability of popular cascades (i.e., simulations where at least 100 users became infected) and the structural virality of these popular cascades.

*5.2.4 Results.* The results of our simulations are shown in Figure 3. Our findings are qualitatively similar to Goel et al.'s, even though our simulations were performed on 1M node graphs, while theirs were performed on 25M node graphs. For higher values of $r$, we observe higher probability that a given cascade reaches at least 100 users. We also find that at lower values of $r$, content is less likely to become popular on graphs generated with lower values of $\alpha$, and at higher levels of $r$, the opposite is true. In line with Goel et al.'s findings, we also observed that the mean structural virality was generally highest for high values of $\alpha$ and increases with $r$.

Finally, Goel et al. observed that simulations run with parameters $r \approx 0.5$ and $\alpha \approx 2.3$ best matched the patterns in the empirical Twitter dataset. For this parameter setting, they found that the probability of a content becoming popular was about one in a thousand, and the mean structural virality of popular cascades was about 3.7. This matches the empirical results from our simulations as well. Goel et al. also find that the correlation between the size of popular cascades and structural virality is approximately 0.1, which falls within range of 0 to 0.2 which was observed in the Twitter data. Similarly, our simulations with the parameter setting $r \approx 0.5$ and $\alpha \approx 2.3$ yielded a correlation of $\approx 0.086$ between size and structural virality.

## 5.3 Challenges in replicating simulation-based research

The prior two case studies illuminated some challenges in replicating simulation-based research. In theory, the methodology section of a technical report should be sufficient for another researcher to be able to fully replicate the design. In practice, few research designs are simple enough to be completely described within a relatively short methodology section. To reproduce the results of Chaney et al. [14] and Goel et al. [29], we needed to consult not only the methodology section, but also methodological appendices, footnotes, and details that appeared only in the graphs of results. When the researcher's primary goal is explaining the theoretical model to the reader, focusing only on the most critical details in the methodology section is understandable; however, this approach has a side effect of rendering the study more difficult to replicate. To successfully produce our conceptual replications, our team frequently had meetings to parse the exact meaning of single sentences or even single phrases within the text to attempt to correctly grasp their intended meaning. Even then, there were multiple plausible

---

[4]Although not specified in the original paper, Goel et al. also ran multiple simulations on pre-generated graphs; we are not certain of the number of unique graphs they generated for their experiments.

(a) Chance of a cascade becoming popular (>100 users).



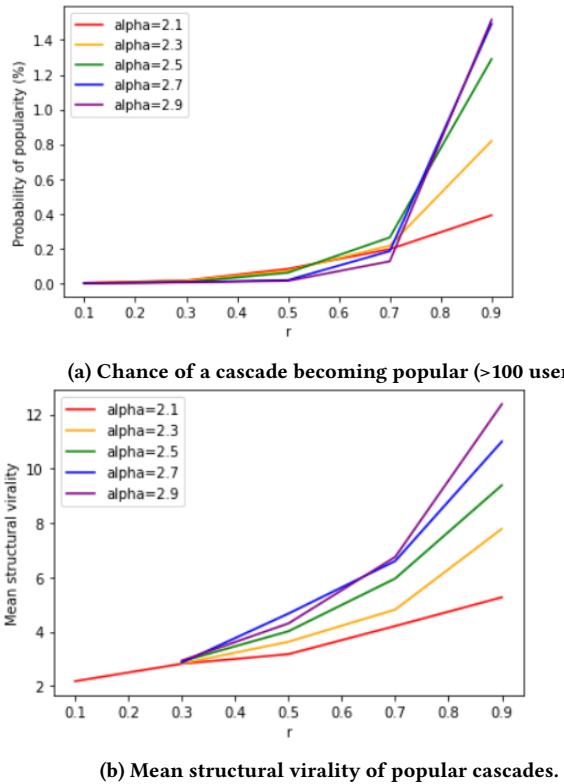(b) Mean structural virality of popular cascades.

**Figure 3: Probability of popular cascades and mean structural virality for SIR simulations on random scale-free networks. Compare to Figure 7(A) and 7(B) of Goel et al. [29]**

operationalizations of the theoretical models described. For example, in the single training case of Chaney et al.'s study, we were not certain whether new items were permanently inserted in each user's recommendation list after random interleaving, or if the new items interleaved at timestep $t$ were never interleaved again at any later timestep. Such details are incidental to the high-level results of the paper; however, we found that such details were critical for correctly reproducing the results. The accumulation of ambiguities in seemingly trivial implementation details meant that our research team spent considerable time and effort conducting analyses to rule out many possible specifications before arriving at the specification that produced results most similar to the original research. Using T-RECS allows researchers to use the language of our framework to describe their methodology in terms that T-RECS users will be familiar with. Adoption of T-RECS also encourages researchers to upload and share code from their experiments – which would completely eliminate the aforementioned ambiguities – since code that uses a common framework - as opposed to an ad-hoc solution - is likely to be more useful and legible to other researchers.

Second, the large scale of the simulations in Goel et al. [29] presented several computational challenges. Without a clear description the resources required to perform the original study, our team needed to experiment with several different hardware configurations and implement software optimizations to replicate their

findings. Future simulation studies can better facilitate replication by providing a description of the hardware and memory requirements for conducting their analysis. In doing so, other researchers would be better able to determine whether whether a replication or novel extension of the work would even be feasible. Similarly, sharing descriptions and code regarding various optimizations techniques (for example, the use of sparse matrices) would likewise reduce the time and effort necessary for other research teams to reproduce – and extend – original simulation findings. Again, we see T-RECS's role in encouraging the sharing of reproducible simulation code as naturally mitigating this issue, as researchers can simply run the simulation code on their own machines to measure performance, without having to rebuild the ad-hoc system. Furthermore, contributors to T-RECS will continue to improve the library, including adding speed and memory optimizations, allowing the entire research community to benefit and removing the need for isolated research teams to reinvent the wheel. For example, support sparse matrices is currently included by default in T-RECS.

Third, rarely do papers describe their procedure for testing the correctness of their implementations; instead, correctness of implementation is taken as a given. We say this not to imply that simulations from prior research were incorrect in any way – rather, we aim to highlight the difficulty of ensuring there are no implementation errors in a simulation of a large, complex system, particularly in cases where the simulation environment is being built ad hoc. Like many machine learning systems, these simulations can fail *silently*, meaning that they do not raise obvious bugs or errors during runtime. Instead, the researcher observes outcomes that are obviously incorrect, and then must work backward through the entire simulation implementation to understand what has gone awry. With this issue in mind, we designed T-RECS with thorough test coverage, believing that well-tested code is essential to robust, reproducible simulation-based research.

Finally, future simulation-based research should include confidence intervals on all results. For example, where figures are shown with results averaged over multiple simulation trials, it would be ideal to provide information about the variance of the results across trials. This would aid future researchers in understanding whether their observed results differ from those reported in the original paper because of differing implementations/assumptions or statistical noise. In our case, this information would have provided an additional amount of confidence that we had properly replicated previous findings, despite minor differences in the appearance of our generated figures.

## 5.4 Content creators and polarization

For our final case study, we present novel work examining a different question: how does the presence of adaptive content creators affect the distribution of *dynamically generated items*, in comparison to when items are served from a fixed catalog? To answer this question, we set up our environment with nearly identical assumptions to those of Chaney et al. The key difference is that the new items introduced at each timestep, rather than being sampled randomly from a fixed distribution, are generated by a pool of content creators, each with their own time-varying item-generating distributions. At each timestep, these content creators adapt their

item-generating distributions to user feedback at each iteration. We investigate how this process results in changes to the items generated over time, when each creator starts out by sampling items from a fairly uniform distribution across item attributes.

Importantly, we reuse much of the code and many of the assumptions from the case study of Chaney et al. In the original Chaney simulations, new items are introduced into the system at each timestep, but they are sampled from a fixed distribution. When new items are instead produced by content creators who respond to user feedback, we may observe that creators themselves are "polarized" in that they begin the simulation by generating a diverse range of items, but then adapt to feedback by generating items from narrower and narrower subsets of the space of possible items.

Our reuse of code affords two main benefits that are indicative of the advantages of using T-RECS in general. First, it mitigates the possibility that differences observed are due to differences in implementation of the simulation environment, rather than the presence of content creators. Second, from a practical standpoint, it saves us a great deal of time since we no longer have to reimplement the assumptions from the original Chaney et al. sutdy from scratch.

### 5.4.1 Measuring creator homogenization.
In the original study by Chaney et al., homogenization is measured through an average Jaccard index of sets of user-item interaction histories [14]. In our study, we shift our focus to studying the homogenization of creators, rather than the homogenization of users. We use the average entropy of each creator's item-generating distribution as a rough proxy for homogenization, henceforth referred to as average creator entropy (ACE). Intuitively, if the entropy of all creator item distributions is high, then each creator is likely to produce a diverse set of items, and if the entropy of the item-generating distributions is low, then each creator is generating a group of items that are more uniform in their attributes. Note that this measure captures *within-creator* homogenization; that is, it measures whether each creator is sampling from a narrow or broad distribution of items, not whether creators are increasingly similar to each other. Therefore, this definition of homogenization does not imply that there cannot be diversity across creators as a whole; instead, it describes individual creator behavior.

### 5.4.2 Adaptive content creators.
The specific item-generating distribution from which content creators sample is a Dirichlet distribution. Each creator's attributes $\gamma_i$ is initially sampled in the following manner: $\gamma_i \sim \text{Dirichlet}(\mathbf{10})$. For a given creator $i$, the item-generating distribution is $\alpha_i \sim \text{Dirichlet}(\gamma_i \cdot 0.1)$. This ensures that item attributes are sparse and are highly sensitive to shifts in the creator's attributes profile, which is initially spread somewhat evenly across most attributes. Recall that this is very similar to how the static item-generating distributions are sampled in Chaney et al. [14].

Content creators respond to feedback by shifting their profiles $\gamma$ towards items they produced with which users interacted during the most recent timestep. (Note that both creator profiles and item attribute vectors sum to one, since both are sampled from Dirichlet distributions.) Notably, we do not claim that this is the best or only way to model adaptive behavior; instead, we posit it as a possible
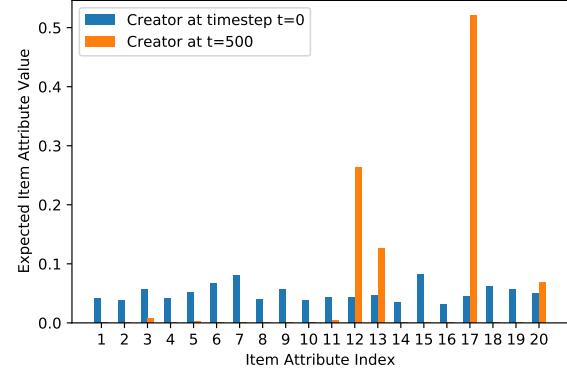


**Figure 4: To illustrate the phenomenon of creator homogenization, we show the change in a randomly chosen creator's item-generating distribution at the start and end of the simulation. At $t = 0$, the expected item attribute values are all approximately the same. Since items are sparse, this indicates that the creator is equally likely to create an item with any of the possible attributes. At $t = 500$, however, the creator is likely to generate only items that have a small set of attributes.**

way to approximate how content creators respond to the incentive to maximize user interactions with their content.

### 5.4.3 Results.
In Figure 4, we provide a visual representation of the trajectory of a single creator's item-generating distribution from the beginning to the end of the simulation. The recommendation system used was the "ideal" recommender from the Chaney et al. study. We observe that at the beginning of the simulation, the creator is about equally likely to create an item with any of the possible attributes. However, at the end of the simulation, the creator has "narrowed" to generating items that have a much smaller set of possible attributes; the creator essentially never generates items that have a high value for the other attributes.

To illustrate this phenomenon across all recommender systems, we plot the ACE at each timestep in Figure 5. In all recommender systems, ACE decreases throughout the simulation at an increasing rate. We also observe a greater degree of creator polarization in the social filtering and popularity recommender systems, suggesting that algorithm selection makes a difference in the rate or degree to which creators are homogenized.

Finally, we performed a preliminary exploration of how creator homogenization impacts user homogenization, using a different measurement of user homogenization than the Chaney et al. use in their study. Our proxy for user homogenization is averaged over pairs of users, where the metric averaged is the distance between the mean items interacted with by the two users. Formally, our metric is: $\frac{1}{n} \sum_{k,j} d(\bar{i}_k, \bar{i}_j)$, where $n$ is the number of pairs of users, $\bar{i}_u$ is the average item profile across user $u$'s interaction history, and $d$ is the Euclidean distance function. We refer to this measure as Average Pairwise Distance between Mean Consumed Items (APDMCI). We find that for all non-ideal recommender systems, APDMCI increases
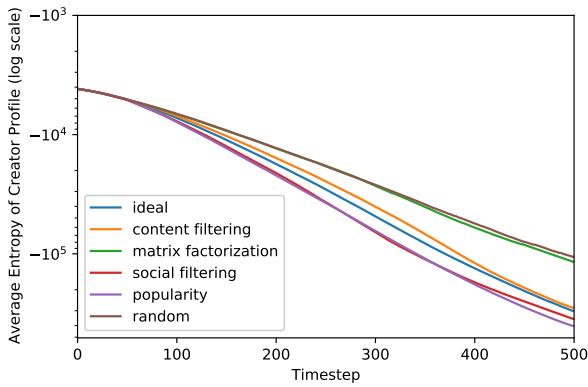
**Figure 5: Average entropy of the creators' item-generating distributions across. For all systems, including the random recommendation scheme, entropy decreases over time, indicating that each individual creator's item-generating distribution is being homogenized over the duration of the simulation. Note that since we use a logarithmic scale on the y-axis, the rate of decrease accelerates over time.**

beyond what is necessary to achieve optimal utility for users. This suggests that as individual creators become homogenized, user behavior is also homogenized.

Note that the Jaccard-based measure of user homogenization is a reasonable measure to use when the item-generating distribution is fixed, but is not able to capture an important dimension of homogenization when creators adapt to user feedback. Consider the example where User A and User B are loyal consumers to Creator C and Creator D, respectively. Initially, Creators C and D generate items that are extremely different from one another, but over time both creators "drift" towards each other, creating items that differ only moderately. In this scenario, a Jaccard-index based measure would report no change in homogenization throughout the experiment, although users A and B are now consuming items that are not as different from each other as they used to be.

*5.4.4 Discussion.* We observe that when content creators adapt to user behavior, within-creator homogenization may occur, as measured by increasing ACE. This means that over time, creators respond to the engagement incentive by narrowing the set of items they produce.

First, we stress that we do not make a normative claim that within-creator homogenization is a negative outcome. In certain domains, our definition of homogenization can be interpreted as creators learning users' interests and responding to user interests. We suggest that recommender system designers consider designing domain-specific interventions to mitigate the possible consequences of polarization.

Second, we also note that we observed a high degree of variance between runs. We average our results over hundreds of simulations to reduce variance in our calculation of the mean outcome. Sources of intrinsic variance include the sampling of new user and creator profiles for each simulation and the noise that is internal to

each recommender system (for example, the process of randomly interleaving new items into recommendation sets). The ability to run this many repeated experiments is a strength of the simulation approach; it may be comparatively expensive or cumbersome to repeatedly sample new users and retrain recommender systems in real-world settings.

Finally, our study shows the need for clarity in how the research community defines concepts like homogenization and polarization that can apply to a wide range of settings. By using a common framework like T-RECS, we can hold all components of an experiment constant and compare how different metrics might capture fundamentally different phenomena. For example, we discovered during our experiments that the measure of user homogenization used by Chaney et al. did not translate perfectly to the adaptive content creator setting, and thus formulated our own measure [14].

These results are intended as an initial exploration into the role that content creators play in recommender systems. Different models of content creators and how they adapt to user feedback may lead to different results. As a result, we invite further research into the interplay between content creators, users, and recommendation system algorithms.

## 6 CONCLUSION

Recommender systems comprise an increasingly central role in many of our collective and individual digital experiences, from reading the news, watching movies, and connecting with friends. In part because recommender systems deployed in production are proprietary and also vary significantly by use and application, researchers in the social sciences and computer science have turned to simulation-based approaches. However, most of these efforts have required researchers to build ad-hoc systems, leading to significant and redundant engineering burdens across studies. Moreover, as recommender systems become more complex, the time required to build even a prototype ad hoc simulation will only increase, as will the likelihood of bugs and errors.

Informed by both an examination of the current trends in simulation-based recommender systems research, T-RECS offers unified, flexible framework for simulating the dynamic interplay between users, items, and algorithms in recommender systems and other sociotechnical systems. Adopting T-RECS offers benefits for individual practitioners, including the guarantee of well-tested code and significant reductions in engineering burden, as well as for the research community as a whole, including the promotion of reproducibility, collaboration between researchers, and a common language for reasoning about complex phenomena in recommender systems.

# REFERENCES

[1] Himan Abdollahpouri, Gediminas Adomavicius, Robin Burke, Ido Guy, Dietmar Jannach, Toshihiro Kamishima, Jan Krasnodebski, and Luiz Pizzato. 2020. Multistakeholder recommendation: Survey and research directions. *User Modeling and User-Adapted Interaction* 30, 1 (2020), 127–158.

[2] Charu C Aggarwal. 2016. Ensemble-based and hybrid recommender systems. In *Recommender systems*. Springer, 199–224.

[3] Muhammad Ali, Piotr Sapiezynski, Miranda Bogen, Aleksandra Korolova, Alan Mislove, and Aaron Rieke. 2019. Discrimination through optimization: How Facebook's Ad delivery can lead to biased outcomes. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–30.

[4] Guy Aridor, Duarte Goncalves, and Shan Sikdar. 2020. Deconstructing the Filter Bubble: User Decision-Making and Recommender Systems. In *Fourteenth ACM Conference on Recommender Systems*. 82–91.

[5] Robert Axelrod. 1997. *The complexity of cooperation: Agent-based models of competition and collaboration.* Vol. 3. Princeton University Press.

[6] Eytan Bakshy, Solomon Messing, and Lada A Adamic. 2015. Exposure to ideologically diverse news and opinion on Facebook. *Science* 348, 6239 (2015), 1130–1132.

[7] Omer Ben-Porat and Moshe Tennenholtz. 2018. A game-theoretic approach to recommendation systems with strategic content providers. *arXiv preprint arXiv:1806.00955* (2018).

[8] Sushil Bikhchandani, David Hirshleifer, and Ivo Welch. 1998. Learning from the behavior of others: Conformity, fads, and informational cascades. *Journal of economic perspectives* 12, 3 (1998), 151–170.

[9] Eric Bonabeau. 2002. Agent-based modeling: Methods and techniques for simulating human systems. *Proceedings of the national academy of sciences* 99, suppl 3 (2002), 7280–7287.

[10] Andrei Borshchev, Sally Brailsford, Leonid Churilov, and Brian Dangerfield. 2014. Multi-method modelling: AnyLogic. *Discrete-event simulation and system dynamics for management decision making* (2014), 248–279.

[11] Dimitrios Bountouridis, Jaron Harambam, Mykola Makhortykh, Mónica Marrero, Nava Tintarev, and Claudia Hauff. 2019. SIREN: A simulation framework for understanding the effects of recommender systems in online news environments. In *Proceedings of the conference on fairness, accountability, and transparency*.

[12] Robin Burke. 2017. Multisided fairness for recommendation. *arXiv preprint arXiv:1707.00093* (2017).

[13] Robin Burke, Nasim Sonboli, and Aldo Ordonez-Gauger. 2018. Balanced neighborhoods for multi-sided fairness in recommendation. In *Conference on Fairness, Accountability and Transparency*. PMLR, 202–214.

[14] Allison J. B. Chaney, Brandon M. Stewart, and Barbara E. Engelhardt. 2018. How Algorithmic Confounding in Recommendation Systems Increases Homogeneity and Decreases Utility. In *Proceedings of the 12th ACM Conference on Recommender Systems* (Vancouver, British Columbia, Canada) *(RecSys '18)*.

[15] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.

[16] Giovanni Luca Ciampaglia, Azadeh Nematzadeh, Filippo Menczer, and Alessandro Flammini. 2018. How algorithmic popularity bias hinders or promotes quality. *Scientific reports* 8 (2018).

[17] Kelley Cotter. 2019. Playing the visibility game: How digital influencers and algorithms negotiate influence on Instagram. *New Media & Society* 21, 4 (2019), 895–913.

[18] Christian S. Crandall and Jeffrey W. Sherman. 2016. On the scientific superiority of conceptual replications for scientific progress. *Journal of Experimental Social Psychology* 66 (2016), 93–99. Publisher: Elsevier.

[19] Alexander D'Amour, Hansa Srinivasan, James Atwood, Pallavi Baljekar, D Sculley, and Yoni Halpern. 2020. Fairness is not static: deeper understanding of long term fairness via simulation studies. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*.

[20] Scott De Marchi and Scott E Page. 2014. Agent-based models. *Annual Review of political science* 17 (2014), 1–20.

[21] Sarah Dean, Sarah Rich, and Benjamin Recht. 2020. Recommendations and user agency: the reachability of collaboratively-filtered information. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 436–445.

[22] Michael D Ekstrand. 2020. LensKit for Python: Next-Generation Software for Recommender Systems Experiments. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2999–3006.

[23] Joshua M Epstein and Robert Axtell. 1996. *Growing artificial societies: social science from the bottom up.* Brookings Institution Press.

[24] Max Fischer and Amanda Taub. 2019. How YouTube Radicalized Brazil. *The New York Times* (2019). https://www.nytimes.com/2019/08/11/world/americas/youtube-brazil.html

[25] Richard Fletcher and Rasmus Kleis Nielsen. 2018. Automated serendipity: The effect of using search engines on news repertoire balance and diversity. *Digital Journalism* 6, 8 (2018), 976–989.

[26] Adrien Friggeri, Lada Adamic, Dean Eckles, and Justin Cheng. 2014. Rumor cascades. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 8.

[27] Kiran Garimella, Aristides Gionis, Nikos Parotsidis, Nikolaj Tatti, et al. 2017. Balancing information exposure in social networks. In *IEEE Conference on Neural Information Processing Systems*.

[28] Daniel Geschke, Jan Lorenz, and Peter Holtz. 2019. The triple-filter bubble: Using agent-based modelling to test a meta-theoretical framework for the emergence of filter bubbles and echo chambers. *British Journal of Social Psychology* 58, 1 (2019), 129–149.

[29] Sharad Goel, Ashton Anderson, Jake Hofman, and Duncan J. Watts. 2016. The structural virality of online diffusion. *Management Science* 62, 1 (2016), 180–196.

[30] Michael Golebiewski and Danah Boyd. 2018. Data voids: Where missing data can easily be exploited. (2018).

[31] Carlos A. Gomez-Uribe and Neil Hunt. 2016. The Netflix Recommender System: Algorithms, Business Value, and Innovation. *ACM Trans. Manage. Inf. Syst.* 6, 4, Article 13 (Dec. 2016), 19 pages. https://doi.org/10.1145/2843948

[32] Aniko Hannak, Piotr Sapiezynski, Arash Molavi Kakhki, Balachander Krishnamurthy, David Lazer, Alan Mislove, and Christo Wilson. 2013. Measuring personalization of web search. In *Proceedings of the 22nd international conference on World Wide Web*. 527–538.

[33] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)* 5, 4 (2015), 1–19.

[34] Charles R. Harris, K. Jarrod Millman, Stéfan J van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. 2020. Array programming with NumPy. *Nature* 585 (2020), 357–362. https://doi.org/10.1038/s41586-020-2649-2

[35] Daniel Heard, Gelonia Dent, Tracy Schifeling, and David Banks. 2015. Agent-based models and microsimulation. *Annual Review of Statistics and Its Application* 2 (2015), 259–272.

[36] Jonathan L Herlocker, Joseph A Konstan, Loren G Terveen, and John T Riedl. 2004. Evaluating collaborative filtering recommender systems. *ACM Transactions on Information Systems (TOIS)* 22, 1 (2004), 5–53.

[37] Ryan Holmes. 2016. The Problem Isn't Fake News, It's Bad Algorithms — Here's Why. *Observer* (2016). https://observer.com/2016/12/the-problem-isnt-fake-news-its-bad-algorithms-heres-why/

[38] Homa Hosseinmardi, Amir Ghasemian, Aaron Clauset, David M Rothschild, Markus Mobius, and Duncan J Watts. 2020. Evaluating the scale, growth, and origins of right-wing echo chambers on YouTube. *arXiv preprint arXiv:2011.12843* (2020).

[39] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *2008 Eighth IEEE International Conference on Data Mining*. Ieee, 263–272.

[40] Neil Hurley and Mi Zhang. 2011. Novelty and diversity in top-n recommendation–analysis and evaluation. *ACM Transactions on Internet Technology (TOIT)* 10, 4 (2011), 1–30.

[41] Eugene Ie, Chih wei Hsu, Martin Mladenov, Vihan Jain, Sanmit Narvekar, Jing Wang, Rui Wu, and Craig Boutilier. 2019. RecSim: A Configurable Simulation Platform for Recommender Systems. (2019). arXiv:1909.04847 [cs.LG]

[42] Taylor Gordon Ivan Medvedev, Haotian Wu. 2019. Powered by AI: Instagram's Explore recommender system. *Facebook AI Blog* (2019). https://ai.facebook.com/blog/powered-by-ai-instagrams-explore-recommender-system/

[43] Dietmar Jannach and Gediminas Adomavicius. 2016. Recommendations with a purpose. In *Proceedings of the 10th ACM conference on recommender systems*. 7–10.

[44] Ray Jiang, Silvia Chiappa, Tor Lattimore, András György, and Pushmeet Kohli. 2019. Degenerate feedback loops in recommender systems. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 383–390.

[45] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2018. Recommendation independence. In *Conference on Fairness, Accountability and Transparency*. PMLR, 187–201.

[46] Karl Krauth, Sarah Dean, Alex Zhao, Wenshuo Guo, Mihaela Curmei, Benjamin Recht, and Michael I Jordan. 2020. Do Offline Metrics Predict Online Performance in Recommender Systems? *arXiv preprint arXiv:2011.07931* (2020).

[47] Eun Lee, Fariba Karimi, Claudia Wagner, Hang-Hyun Jo, Markus Strohmaier, and Mirta Galesic. 2019. Homophily and minority-group size explain perception biases in social networks. *Nature Human Behaviour* 3, 10 (2019), 1078–1087.

[48] Rebecca Lewis. 2018. Alternative influence: Broadcasting the reactionary right on YouTube. *Data & Society* 18 (2018).

[49] Dongwon Lim, Hwansoo Lee, Hang-Jung ZO, and Andrew Ciganek. 2014. Opinion formation in the digital divide. *The Journal of Artificial Societies and Social Simulation* 17, 1 (2014).

[50] Helmut Lorek and Michael Sonnenschein. 1999. Modelling and simulation software to support individual-based ecological modelling. *Ecological Modelling* 115, 2-3 (1999), 199–216.

[51] Megan McCluskey. [n.d.]. These TikTok Creators Say They're Still Being Suppressed for Posting Black Lives Matter Content. *Time* ([n. d.]). https://web.archive.org/web/20210325170448/https://time.com/5863350/tiktok-black-creators/

[52] Martin Mladenov, Chih-Wei Hsu, Vihan Jain, Eugene Ie, Christopher Colby, Nicolas Mayoraz, Hubert Pham, Dustin Tran, Ivan Vendrov, and Craig Boutilier. 2021. RecSim NG: Toward Principled Uncertainty Modeling for Recommender Ecosystems. *arXiv preprint arXiv:2103.08057* (2021).

[53] Bojan Mohar and Tomaž Pisanski. 1988. How to compute the Wiener index of a graph. *Journal of mathematical chemistry* 2, 3 (1988), 267–277.

[54] Kevin Munger and Joseph Phillips. 2020. Right-Wing YouTube: A Supply and Demand Perspective. *The International Journal of Press/Politics* (2020), 1940161220964767.

[55] Luke Munn. 2019. Alt-right pipeline: Individual journeys to extremism online. *First Monday* (2019).

[56] HR Nasrinpour, MR Friesen, and RD McLeod. 2016. An agent-based model of message propagation in the Facebook electronic social network. preprint. *arXiv preprint arXiv:1611.07454* (2016).

[57] Tien T Nguyen, Pik-Mai Hui, F Maxwell Harper, Loren Terveen, and Joseph A Konstan. 2014. Exploring the filter bubble: the effect of using recommender systems on content diversity. In *Proceedings of the 23rd international conference on World wide web*. ACM, 677–686.

[58] Eli Pariser. 2011. *The filter bubble: How the new personalized web is changing what we read and how we think.* Penguin.

[59] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.

[60] Nicola Perra and Luis EC Rocha. 2019. Modelling opinion dynamics in the age of algorithmic personalisation. *Scientific reports* 9, 1 (2019), 1–11.

[61] Al Mamunur Rashid, Istvan Albert, Dan Cosley, Shyong K Lam, Sean M McNee, Joseph A Konstan, and John Riedl. 2002. Getting to know you: learning new user preferences in recommender systems. In *Proceedings of the 7th international conference on Intelligent User Interfaces*. 127–134.

[62] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).

[63] Manoel Horta Ribeiro, Raphael Ottoni, Robert West, Virgílio AF Almeida, and Wagner Meira Jr. 2020. Auditing radicalization pathways on YouTube. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 131–141.

[64] Francesco Ricci, Lior Rokach, and Bracha Shapira. 2011. Introduction to recommender systems handbook. In *Recommender systems handbook*. Springer, 1–35.

[65] David Rohde, Stephen Bonner, Travis Dunlop, Flavian Vasile, and Alexandros Karatzoglou. 2018. Recogym: A reinforcement learning environment for the problem of product recommendation in online advertising. *arXiv preprint arXiv:1808.00720* (2018).

[66] Matthew J Salganik, Peter Sheridan Dodds, and Duncan J Watts. 2006. Experimental study of inequality and unpredictability in an artificial cultural market. *science* 311, 5762 (2006), 854–856.

[67] Thomas C Schelling. 1971. Dynamic models of segregation. *Journal of mathematical sociology* 1, 2 (1971), 143–186.

[68] Sven Schmit and Carlos Riquelme. 2018. Human interaction with recommendation systems. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 862–870.

[69] Ken Shoemake. 1985. Animating rotation with quaternion curves. In *Proceedings of the 12th annual conference on Computer graphics and interactive techniques*. 245–254.

[70] Joan E. Solsman. 2018. YouTube's AI is the puppet master over most of what you watch. *CNET* (2018). https://www.cnet.com/news/youtube-ces-2018-neal-mohan/

[71] Wenlong Sun, Sami Khenissi, Olfa Nasraoui, and Patrick Shafto. 2019. Debiasing the human-recommender system feedback loop in collaborative filtering. In *Companion Proceedings of The 2019 World Wide Web Conference*. 645–651.

[72] Cass R. Sunstein. 2001. *Republic.com.* Princeton University Press.

[73] Cass R Sunstein. 2009. *Going to extremes: How like minds unite and divide.* Oxford University Press.

[74] Cass R. Sunstein. 2009. *Republic.com 2.0.* Princeton University Press.

[75] Marcella Tambuscio, Diego FM Oliveira, Giovanni Luca Ciampaglia, and Giancarlo Ruffo. 2018. Network segregation in a model of misinformation and fact-checking. *Journal of Computational Social Science* 1, 2 (2018), 261–275.

[76] Seth Tisue and Uri Wilensky. 2004. Netlogo: A simple environment for modeling complexity. In *International conference on complex systems*, Vol. 21. Boston, MA, 16–21.

[77] Petter Törnberg. 2018. Echo chambers and viral misinformation: Modeling fake news as complex contagion. *PloS one* 13, 9 (2018), e0203958.

[78] Francesca Tripodi. 2018. Searching for alternative facts. *Data & Society* (2018).

[79] Zeynep Tufekci. 2018. YouTube, the Great Radicalizer. *The New York Times* (2018). https://www.nytimes.com/2018/03/10/opinion/sunday/youtube-politics-radical.html

[80] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods* 17 (2020), 261–272. https://doi.org/10.1038/s41592-019-0686-2

[81] Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. *Science* 359, 6380 (2018), 1146–1151.

[82] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. 2015. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 1235–1244.

[83] Uri Wilensky. 1999. NetLogo. https://ccl.northwestern.edu/netlogo/.

[84] Sirui Yao and Bert Huang. 2017. Beyond parity: fairness objectives for collaborative filtering. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*. 2925–2934.

[85] Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang, Nicholas Jing Yuan, Xing Xie, and Zhenhui Li. 2018. DRN: A deep reinforcement learning framework for news recommendation. In *Proceedings of the 2018 World Wide Web Conference*. 167–176.

[86] Lei Zheng, Vahid Noroozi, and Philip S Yu. 2017. Joint deep modeling of users and items using reviews for recommendation. In *Proceedings of the tenth ACM international conference on web search and data mining*. 425–434.

[87] Cai-Nicolas Ziegler, Sean M McNee, Joseph A Konstan, and Georg Lausen. 2005. Improving recommendation lists through topic diversification. In *Proceedings of the 14th international conference on World Wide Web*. 22–32.

# A   APPENDIX

## A.1   Replication

*A.1.1   Algorithmic confounding (Chaney et al.)* In the following sections, we provide more details about our methods reproducing Chaney et al.'s work [14]. If we have omitted any details, please defer to the assumptions laid out in the original paper.

**Recommender system models.** Following Chaney et al., we model six types of recommender systems. Each recommender system maintains an internal representation of predicted user preferences and predicted item attributes. Recommendations to each user are based on the predicted user scores, which are themselves a function of the predicted user preferences and predicted item attributes (typically a dot product).

The first model is content-based filtering, which recommends content similar to what users have liked in the past. As Chaney et al. indicate, the predicted user preferences in this model are updated at each step by solving for the least-squares approximation of user attributes using the `scipy.optimize.nnls` function [80]. Item attributes are equal to the true item attributes.

Second, we use popularity-based filtering, which serves the most popular items in the systems. In this model, the predicted item attributes is equal to a single number for each item: the total number of interactions that item has received. Predicted user preferences are identical and constant for all users, so that the predicted score for user $u$ and item $i$ is simply equal to the number of interactions that item $i$ has received.

Third, we implement a matrix factorization collaborative filtering model in which a common latent representation of users and items is used to recommend items based on user interaction data. The predicted user preferences and item attributes are generated using

| Simulation to study sociotechnical systems | | |
|---|---|---|
| **Work** | **Concept** | **Metric(s)** |
| Aridor et al. [4] | Filter bubble | Average consumption distance at time $t$ and $t-1$ |
| | User welfare | Average realized utility |
| | Item diversity | Average normalized pairwise distance between consumed items as seen in [87] |
| | User homogenization | Average pairwise Jaccard index (as seen in [14]) |
| Chaney et al. [14] | User homogenization | Jaccard index on the sets of seen items by pairs of users |
| | Differential item consumption | Gini index on the distribution of consumed items |
| Ciampaglia et al. [16] | User welfare | Average quality of items |
| | Faithfulness | Kendall rank correlation between popularity and quality of items |
| Garimella et al. [27] | Filter bubble | Expected number of users exposed to both or neither of two items (representing viewpoints) |
| Geschke et al. [28] | Echo chamber | Mean distances between users and items/item sharers/friends |
| | | Visual analysis of clustering |
| Goel et al. [29] | Structural virality | Wiener index of diffusion tree [53] |
| Jiang et al. [44] | Echo chamber | Distance between initial user interests and final user interests (weak or strong degeneracy) |
| | Filter bubble | Speed of degeneracy |
| Lee et al. [47] | Social perception bias | Error of individuals in estimating true prevalence of binary attribute |
| Lim et al. [49] | Polarization | Number of clusters, Herfindahl-Hirschman Index (HHI), proportion of agents in majority/minority clusters, speed of convergence |
| Nasrinpour et al. [56] | Virality | Number of interactions (reads, reposts) |
| Perra and Rocha [60] | Polarization | Prevalence of opinions over time |
| | Echo chamber | Distribution of neighbors holding the majority or minority opinion at end of simulation, relative to start of simulation |
| Sun et al. [71] | Accuracy | Root Mean Squared Error |
| | Popularity bias | Gini coefficient |
| Tambuscio et al. [75] | Spread of misinformation | Fraction of infected users at equilibrium, state transition rates |
| | Echo chamber | Custom generative model for segregated network |
| Törnberg [77] | Virality | Probability that majority of nodes are "infected" |
| | Network polarization | Increased ties within cluster, decreased ties from cluster to outside |
| | Opinion polarization | Probability that neighboring nodes have similar activation thresholds |
| Yao and Huang [84] | Group unfairness | Four different unfairness metrics |

**Table 2: Simulation literature survey. Note that many different metrics are used to assess the dimensions of various concepts. Concepts that are similar in nature are given the same color.**

the alternating least squares approach [39] implemented in the LensKit library [22]. Note that this matrix factorization model is not identical to the model used by Chaney et al.

Fourth, we implement social-based filtering, which recommends items based on the preferences of users in their social network. In this model, predicted user preferences are represented with an adjacency matrix that includes the connections between users. Following Chaney et al., we generate this adjacency matrix from the covariance matrix of true user preferences.

Lastly, we provide two baseline models: a random recommender system and an ideal recommender system. The former serves random items to users. The latter presents items based on the users' true utility.

**Synthetic data.** Next, we generate the user and item data. The actual user preferences $U$ are represented as a $|U| \times |A|$ matrix, where $|U| = 100$ is the number of users and $|A| = 20$ is the number of *attributes* describing each user. The $i$-th row of $U$ contains the attributes describing user $u_i$. The values of the rows were drawn from a Dirichlet distribution with parameters as specified by Chaney et al.

Each model requires a distinct representation of the predicted user preferences $\hat{U}$. In content-based filtering, predicted user preferences are a $|U| \times |A|$ matrix with the same properties as the actual user preference matrix $U$. In the popularity-based recommender system, $\hat{U}$ is a $|U| \times 1$ matrix with all elements equal to 1, as the predictions of the system are the same for all users (recall that predictions are calculated with the dot product of the user item matrices). In our collaborative filtering model, we use matrix factorization; therefore, predicted user profiles are a $|U| \times k$ matrix, with $k$ being the number of features in the latent representation. In our social-based filtering model, we represent the users' social networks in $\hat{U}$; therefore, $\hat{U}$ is a $|U| \times |U|$ matrix.

Items in the system are represented by a matrix $I$ of size $|I| \times |A|$, where $|I| = 1000$ is the total number of items and $|A| = 20$ is the number of attributes that describe each item, analogously to matrix $U$ for users. The values of the rows were drawn from a Dirichlet distribution with parameters as specified by Chaney et al.

We additionally calculate users' *utility* as defined by Chaney et al. We specifically distinguish between true utility and *known* utility. The latter is the fraction of true utility known to users and is equivalent to the actual user scores in our framework — that is, the scores that the recommender systems predict.

True and known utilities are represented with two distinct matrices of size $|U| \times |I|$. For a given user $u$ and a given item $i$, the true utility obtained by $u$ interacting with $i$ is sampled from a Beta distribution whose mean is parameterized as the dot product of the $u$'s true preferences and $i$'s true attributes. From the vantage point of $u$, the known utility of interacting with $i$ is a percentage drawn from a beta distribution with parameters as specified by Chaney et al.

**Performance** The total time to run the entire set of experiments was just under 3 hours on a laptop with 32 GB of RAM and a 6-core 2.6 GHz processor.

*A.1.2 Structural virality (Goel et al.)* **Performance.** Several optimizations were required for our replication effort. First, we had to ensure that the large graphs of millions of users fit into memory during simulation. T-RECS supports `scipy` sparse matrices [80], which allowed our simulations of hundreds of thousands of users to run using just a few gigabytes of RAM. Second, we made thorough use of multiprocessing-based parallelization to generate tens of thousands of simulation runs within days; this was essential because the probability of generating a popular cascade on any given run was extremely low (approximately 0.1%). Examples of using multiprocessing to achieve these speedups are available online at https://github.com/sunnymatt/t-recs-experiments.

*A.1.3 Content creators.* **Performance.** We perform 200 separate 500-timestep trials for each of our six recommendation algorithms in approximately six hours. While our simulations were performed on our university's academic computing cluster, the simulations utilized just 4.20 GB of memory, suggesting that the experiments could also be performed on regular desktop or laptop hardware.

## A.2 Internal representation for users and items
**Challenge: developing a scalable and flexible representation to encode information about users and items.**

We considered two different engineering approaches to represent users and items internally:

(1) Dedicate a data structure to each user and item.
(2) Dedicate one data structure to all users and another data structure to all items.

While the first approach may be the most natural choice for an object-oriented design and well-suited to emphasizing the unique characteristics of each user and item, the latter has several advantages. First, it scales well with the addition of more users or items. Second, it is flexible because the system does not need to make assumptions on the number of data structures dedicated to users and items, as it will be constant to one. Third, it allows to process user actions on items in parallel.

As seen in previous work [14], T-RECS stores information about all the users in the system in one matrix; similarly, all the items are represented with a single matrix. For both users and items, two representations are generated—a ground truth representation and

a predicted representation. Therefore, we use a matrix to represent predicted user preferences, and a separate matrix for the actual user preferences.

*A.2.1 Actual user preferences.* Bidimensional matrix $U$ with the first dimension equal to the number of users in the recommender system ($|U|$). The second dimension varies with the model. Typically, the $i$-th row represents the preferences of user $U_i$ ($i \in \{0, \dots, |U| - 1\}$). Actual user preferences are unknown to the recommender system and only known to the users.

*A.2.2 Predicted user preferences.* Bidimensional matrix $\hat{U}$ with the same properties as $U$. This matrix is calculated by the model and represents inferred user preferences.

*A.2.3 Item attributes.* Bidimensional matrix $I$ with the second dimension equal to the number of items in the recommender system ($|I|$). The first dimension varies with the model. Typically, the $j$-th column represents the attributes of item $I_j$ ($j \in \{0, \dots, |I| - 1\}$). For simplicity, we assume that the model knows the real item attributes; therefore, predicted item attributes are equal to the actual item attributes.

*A.2.4 Actual user scores.* Matrix $S$ of dimensions $|U| \times |I|$. The element at $S_{ij}$ represents the ground truth score for user $U_i$ on item $I_j$. These scores are typically used to define the user behavior when giving feedback to items that are presented to them. That is, users will interact with items with the highest predicted actual user score. T-RECS provides two built-in methods to calculate actual user scores: the inner product between actual user preferences and item attributes, or the cosine similarity between the two matrices. Actual user scores are unknown to the recommender system and are only known to the users.

*A.2.5 Predicted user scores.* Matrix $\hat{S}$ of dimensions $|U| \times |I|$. Predicted user scores have the same dimensions as $S$, but are is calculated by the model using $\hat{U}$ instead of $U$. Predicted user scores are used by the model to recommend item to users.

*A.2.6 Metrics.* **Challenge: instrumenting the simulator to gather information over time about.**

A fitting implementation choice for metrics was the observer design pattern, a standard pattern in object-oriented programming that is used when one object (the *observer*) needs to be "notified" automatically when a second object (the *observable*) changes. In the case of T-RECS, the model is the observer and the metrics are the observable objects updated at each step. Specifically, we define the observer paradigm with a single observer and multiple observables, as opposed to the traditional pattern with multiple observers and a single observer. Therefore, models can *register* (i.e., monitor) and unregister (i.e., stop monitoring) multiple metrics. Using the observer design pattern also provided flexibility to easily develop and integrate new metrics.

Based on a review of the simulation literature about algorithmic systems (Table 1), we identified two different kinds of metrics: *measurements* and internal *system states*. They are both implemented following the observer pattern.

Measurements are designed to calculate quantities of interest about the system as a whole. For example, diversity of content

can be calculated by comparing the distances between the user preferences and the attributes of the items recommended [28].

System states expose relevant information about the simulation and the system. Unlike measurements, system states are presented as is and do not require additional computation. For instance, researchers interested in monitoring the evolution of the predicted user preferences can register matrix $\hat{U}$ as a system state; at each iteration, the state of $\hat{U}$ will be stored internally. This also provides a mechanism useful to debug model behavior and retrieve quantities for later inspection and manipulation.

*A.2.7 Content creators.* By default, content creators are modeled by a matrix $C$ of dimension $\mid C \mid \times \mid A \mid$, where $C$ is the number of content creators and $A$ is the number of attributes that parameterize each creator's item-generating distribution. In the default model, each content creator $c$ is modeled as a vector of length $A$, where each entry in the vector denotes $p_a$, the Bernoulli probability that an item created by $c$ will have a value of 1 for attribute $i$. (Note that this implies that item attributes are sampled independently.) Researchers can also pass in an optional probability $p_c$ that represents the probability that each creator generates an item at any timestep; this quantity can be used to determine the general "productivity" level of the content creator pool.

We aimed to maximize simplicity and modularity with our implementation of content creators. A simulation with 100 content creators that generate items with 20 attributes can be instantiated as follows:

```
c = Creators(size=(100, 20))
recsys = trecs.models.ContentFiltering(creators=c)
recsys.run(timesteps=100)
```

This "drop-in" API allows researchers to easily measure the impact of adding content creators to their simulations.

This simple default setup is not intended to be highly realistic; for example, it does not impose any restrictions on the sparsity of item attributes. Instead, the default implementation is meant to illustrate the general mechanics of representing content creators through the matrix $C$, which contains the parameters for each content creator's item-generating distribution. Because the research community has not yet reached consensus on the correct theoretical model for content creators, we expect researchers to test theoretical variations of content creators by implementing their own content creator subclasses or by modifying the existing content creator subclass to better suit their research questions.

Simulating content creators generally induces a greater computational cost per timestep than a fixed item catalog, since the new items are sampled at each timestep. As the item set grows over time, long simulations may result in high memory and compute requirements. The choice of item-generating distribution also affects simulation performance; for example, we found in practice that modeling each content creator's item-generating distribution as a multivariate normal slowed simulations by an order of magnitude.

Finally, the basic simulation dynamics described in Section 4.1 can also be modified with the following options:

## A.3 Dynamic user preferences

Modeling how users themselves change over time is key to capturing phenomena of interest that occur in some part at the level of individual users, such as political polarization or ideological radicalization. For example, Geschke et al. [28] simulation-based studies model how agents attitudes' shift over time in environments with or without recommendation algorithms. Jiang et al. [44] provide a theoretical model for how recommender systems might alter user preferences over time. In their model, the utility a user receives from interacting with an item at timestep $t$ might depend on the user's history of interactions at all previous timesteps $t-1, t-2, \cdots$. Empirical research also seems to suggest that user's interest might change over time, such as when users gravitate to progressively more extremist content on YouTube [63]. The capacity to model how individual users are affected by their interactions with sociotechnical systems is necessary to capture these types of effects on users.

In T-RECS, choosing to model dynamic user preferences in a given simulation translates to an added step after users choose which item to interact with from the recommendation set. For a given user, the user's attributes "drift" towards the chosen item's attributes; concretely, we implement spherical linear interpolation [69], such that the user's profile vector is rotated in the direction of the item vector. Users of the T-RECS library may also choose to implement their own custom drift functions.

## A.4 Start-up mode

Many recommendation models are unable to make recommendations to new users (i.e., *cold-start*); therefore, we provide a mechanism for the system to gather information about users' preferences during a *start-up phase*. We assume that users have preferences that can be expressed for all items; therefore, we present items randomly during start-up to maximally sample the range of users' preferences and minimize the risk that start-up preference elicitation strategy could bias the model.

As with other system components, researchers can define their own start-up strategies if they are interested in the effect of start-up behavior on subsequent system state. For example, a reasonable assumption is that users cannot express preferences for unfamiliar items, so practitioners interested in collecting start-up preferences as efficiently as possible may benefit from showing users popular items [61]. However, popular items are the least informative about individual users' preferences since they tend to be preferred by all users.

In start-up mode, the system skips step 1 of the simulation dynamics. At the end of the start-up phase, the model is trained on the collected interactions – that is, the predicted user scores are calculated based on the items users have interacted with. Note that the T-RECS API does not support adding and removing users during a simulation, so the cold start problem only manifests itself when training a new model from scratch.