

Personalized Recommendations in EdTech: Evidence from a Randomized Controlled Trial*

Keshav Agrawal[†] Susan Athey[‡] Ayush Kanodia[§] Emil Palikot[¶]

August 31, 2022

Abstract

We study the impact of personalized content recommendations on the usage of an educational app for children. In a randomized controlled trial, we show that the introduction of personalized recommendations increases the consumption of content in the personalized section of the app by approximately 60% and that the overall app usage increases by 14%, compared to the baseline system of stories selected by content editors for all students. The magnitude of individual gains from personalized content increases with the amount of data available about a student and with preferences for niche content: heavy users with long histories of content interactions who prefer niche content benefit more than infrequent, newer users who like popular content. To facilitate the diffusion of personalized recommendation systems, we provide a framework for using offline data to develop such a system.

*We are grateful to Deepak Agarwal and his team at Stones2Milestones for collaboration on this project.

[†]Stanford Graduate School of Business, keshavag@stanford.edu

[‡]Stanford Graduate School of Business, athey@stanford.edu

[§]Stanford Computer Science, akanodia@stanford.edu

[¶]Stanford Graduate School of Business, palikot@stanford.edu.

1 Introduction

Recommendation systems, the algorithms that determine which pieces of content will be displayed to each users, have been widely deployed in online services and credited with being an important factor in determining user engagement with the service. Personalized content recommendations have contributed to the success of some of the most valuable companies in the world. Market leaders in the entertainment sector (e.g., *Netflix* or *Spotify*) and in online retail (e.g., *Amazon*) are at the forefront of developing algorithms that provide personalized recommendations, and reap high benefits from implementing them.¹ Recommendations, particularly personalized ones, in principle, have the potential to create significant value in other settings where user preferences for items vary.

However, the incremental benefits of personalization have also been challenged, and the empirical question of its impact remains open in many settings, particularly in education. In real-world educational applications, the user base may be orders of magnitude smaller than popular entertainment applications, and so it is unclear whether data-driven personalization would be effective in such settings. In addition, the benefits of personalization depend on the fundamental preferences of the users (e.g. students); if their preferences are homogeneous, then human curation or simple popularity-based algorithms may be sufficient. Therefore, empirical evidence is required to understand the importance of personalization in a given setting.

In education and training, students might spend more time with educational material (and potentially learn more) if it matches their interests. Yet perhaps surprisingly, the publicly available evidence of the impact of personalized content recommendations in education is limited. This paper aims to fill some of these gaps by providing evidence from a large-scale randomized controlled trial (RCT) designed to measure the impact of the introduction of personalized recommendations in place of editor based manually curated recommendations into *Freedom*, an educational app designed to help children in India learn to read in English. In particular, we conducted a two-week-long randomized experiment, where the control group was exposed to stories based on the status quo, a system in which editors select content for all users (the “editorial-based” system), while the treatment group was exposed to stories from a personalized recommendation system in one section of the app.

Our most important finding is that personalization of recommended content leads to a substantial increase in user engagement with the app compared to the editorial-based system: our estimate of the

¹See (Gomez-Uribe and Hunt, 2015) for a discussion of the purpose and business value of personalized recommendation algorithms at *Netflix*.

increase in usage of the personalized section is 63% ($\pm 28\%$).² A key element of the experiment is that the personalized content was shown in one section of the app; thus, it is possible that users might simply shift from consuming editor-based content to personalized content without increasing overall engagement. We also estimate the total increase in app usage which includes all sections of the app and estimate an increase of approximately 14% ($\pm 12\%$).

Increases in the consumption of educational content of the magnitude that we estimate can lead to substantial societal benefits. Notably, as the app's content is curated by pedagogy experts, higher levels of engagement are likely to accelerate learning. It is worth noting that *Freadom* has wide reach at a low cost; therefore, improving its efficiency can potentially benefit a large user base.

Personalization of content selection in the ed-tech context typically takes the form of either assigning learning materials at the difficulty level that is right for the specific user or adjusting the content's style so that it matches the user's preferences.³ In this paper, we focus on the latter. It is not a priori clear that personalized content increases app usage. Notably, learners might engage with ed-tech products following a specific routine or, in the case of children, the recommendation of parents or teachers. The finding that overall usage of the app increases following the introduction of personalized recommendations suggests that investments in recommendation systems in the ed-tech context can create substantial value.

To understand better the potential impact of the intervention, consider the context of the *Freadom* app. It is developed by *Stones2Milestones (S2M)*, and it is targeted at children aged 3 to 12 years old. Short illustrated stories are the main content of the app. Each story is a self-contained learning unit, generally consisting of a reading part and a quiz. Stories are curated by *S2M* pedagogy experts; they are grade-appropriate and have clear educational goals. *Freadom* is mostly used on smartphones, where the main page of the app consists of various sections. Each section contains a tray of stories. A tray is a sequence of stories sorted by an algorithm. Trays are labelled with different names e.g., *Trending Now*, *New Releases*, or *Recommended Story* and display stories following different algorithms (e.g., *New Releases* features stories recently added to the app). At the time we conducted this research, the algorithms assigning stories to trays were not personalized, and either manual curation or simple algorithms such as the most recently added stories, were used to select stories.

The first step of the project was to develop a personalized recommendation system using data on

²In the brackets we show a 95% confidence interval.

³See Escueta et al. (2020) for a review of the literature on the impact of personalization of learning content difficulty on learning outcomes.

historical user-story interactions. We compared several alternative approaches, selecting an approach based on collaborative filtering (Mnih and Salakhutdinov, 2007; Rendle, 2010) which performed best of the alternatives we considered in terms of estimated policy values (estimated using doubly robust off-line policy evaluation (Gilotte et al., 2018; Zhan et al., 2021)). However, off-line analysis is tailored to understanding the impact of recommending different individual stories to users on their engagement with the particular story, but it does not capture the effects of sustained exposure to a personalized recommendation system. In addition, off-line policy evaluation of recommendation systems has known limitations in terms of both bias and variance. This motivates our next step, which was to design a Randomized Controlled Trial (RCT) in order to compare the status quo system of manually curated recommendations to the personalized algorithm.

To evaluate the impact of personalized content recommendations on the utilization of the app, we carried out a randomized experiment. Since collaborative filtering requires substantial user history to perform well, the experiment included users who interacted with at least sixty stories before the start of the experiment. The main outcome metric is a user’s total utilization of the app, defined as the sum of utilities from all user-story interactions during the experiment. Utility is a constructed metric, which assigns a value of one if a user completed a story, 0.5 if a user started the story but did not finish it, and 0.2 if the user clicked on the story to view the description but did not start it. Otherwise, the user is assigned the utility of zero. The experiment lasted for two weeks. We summarize our findings next.

We find that users in the treatment group had a 63% ($\pm 28\%$) higher total utility from content interactions in the personalized tray compared to users in the control group. Treated users also completed 78% ($\pm 39\%$) more stories and spent 87% ($\pm 41\%$) more time-consuming content on the personalized tray. We document significant patterns of heterogeneity in treatment effects. Users who consumed more niche content (i.e., content that is less popular overall) in the pre-experimental period had substantially higher treatment effects than users who like popular content. This is an expected result as the editorial team selects content targeted to typical tastes. Therefore, users with preferences that are different than those of the majority are likely to benefit more from personalization. Furthermore, users with long histories of content interactions also gained more from the personalization of content. This is because the performance of the collaborative filtering model improves when more information about past interactions is available. Last, we compare outcomes of users that had used the *Recommended Story* tray in the past and users that have not. We find statistically significant treatment

effects in both groups. The positive treatment effect for users that were not interacting with stories in this tray in the past suggests that users explore the app enough to notice content even in trays they rarely use and adjust their consumption decisions.

Users who received personalized recommendations in the *Recommended Story* tray increased utilization of the app across all trays. We find a 14% ($\pm 12\%$) increase in the total utilization of the app, a 19% ($\pm 14\%$) increase in the number of completed stories, and a 20% ($\pm 14\%$) growth in the time spent reading stories. Also, users in the treatment group who didn't read any stories on the *Recommended Story* tray prior to the experiment exhibited a much larger (statistically significant) propensity to start reading on this tray compared to users in the control group. These results suggest that the increased usage of the *Recommended Story* tray is not driven entirely by substitution away from other trays in the app. On the contrary, we find that users substitute away from other non-app activities to use the app more. In summary, better content selection can increase the overall utilization of an ed-tech app, justifying investments in developing recommendation systems.

Literature review. This paper relates to several strands of literature. Personalized recommendation systems have been studied intensively in entertainment (Davidson et al., 2010; Gomez-Uribe and Hunt, 2015; Jacobson et al., 2016; Holtz et al., 2020) and in retail shopping (Linden et al., 2003; Sharma et al., 2015; Smith and Linden, 2017; Greenstein-Messica and Rokach, 2018; Ursu, 2018). For example, in the entertainment context and using a similar approach to our paper, (Holtz et al., 2020) show that personalized recommendations increase consumption of podcasts on Spotify. However, there is little empirical evidence of the usefulness of recommendation engines beyond entertainment platforms and e-commerce. This paper attempts to fill this gap by providing evidence from the ed-tech sector.⁴ Additionally, we show that personalized recommendations can be an effective method of boosting user engagement in settings with moderate amounts of data.

The existing evidence of the efficacy of recommendation systems in education is generally based on small studies that combine the introduction of personalized recommendations with other changes to the user interface. (Drachler et al., 2008) use a recruited group of university students to study the effect of showing personalized recommendations of course materials to not showing any recommendations at all. While this study is an A/B experiment, it bundles two changes in one treatment: adding

⁴Drachler et al. (2015) provide an extensive review of literature on recommendation systems in ed-tech and point out a shortage of papers documenting the efficiency of recommendation systems using reliable evaluation methods. They conclude the review by calling for more comprehensive user studies in a controlled experimental environment.

a user interface element and personalizing recommendations. Furthermore, this study is based on a relatively small sample of 250 subjects. (Ruiz-Iniesta et al., 2018) develop and test a recommendation system on an ed-tech platform called *Smile and Learn*, and evaluate it in an observational study. Their proposed treatment is a new user interface component with recommendations generated using collaborative filtering. The newly introduced system helps users navigate the app and reach desired content quicker. They find substantial increases in consumption of recommended items versus non-recommended items. However, the treatment in (Ruiz-Iniesta et al., 2018) has two elements: the part simplifying app navigation by adding a user interface component and a personalization component. Our work provides results that isolate the impact of personalization on the consumption of learning items.⁵ To the best of our knowledge, our paper is the first large-scale study in the ed-tech context that estimates the effect of personalization on user engagement in isolation from other changes in the app.

Second, our work contributes to the growing literature assessing the effects of personalized recommendation systems on the diversity of consumed content. To our knowledge, we are the first to do so in an ed-tech context. Anderson et al. (2020); Holtz et al. (2020) provide evidence from a randomized experiment indicating that personalized recommendations reduce the diversity of content consumed on *Spotify*. In the context of retail, (Lee and Hosanagar, 2019) show that, while recommendations reduce within-consumer diversity, their effect on aggregate diversity is ambiguous. (Claussen et al., 2021) find that recommendations reduce consumption diversity in the context of news consumption. In this paper, we show that users with niche preferences are recommended more niche content and less often interact with stories liked by the majority of users. This closely relates to the literature documenting ‘filter-bubbles’ due to the personalization of content on media platforms (Haim et al., 2018; Möller et al., 2018).

Last, this paper relates to a rich literature on technology-assisted language learning.⁶ Personalization in the language learning context has been shown to be effective in task assignment (Xie et al., 2019) and learning resource recommendations (Sun et al., 2020). We contribute to this literature by bringing causal evidence of the impact of personalization on time spent interacting with language learning content.

The rest of the paper is organized as follows. Section 2 details the empirical setting. Section 3

⁵Contexts of (Ruiz-Iniesta et al., 2018) and of this paper also differ substantially. In our setting, we have thousands of stories to choose from as compared to around one hundred games. This seemingly technical difference results in problems of data sparsity, which is a serious challenge in creating recommendations for stories that are relatively new. In section 3, we present the methodology for designing and evaluating a recommendation system in such settings.

⁶See Garrett (2009), (Zhao, 2003), and (Tafazoli et al., 2019) for reviews of this literature.

presents the methodology used to develop and test the recommendation model using offline data. Section 4 describes the design of the randomized experiment and presents the results. Finally, section 5 concludes.

2 Empirical setting

Stones2Milestones (S2M) was founded in 2009 in India. The company provides technology-enabled English education through a variety of programs serving a diverse set of users. The main product of *S2M* is a smartphone app called *Freedom*, aimed at 3 to 12-year-old children. Throughout 2021, the average daily number of users amounted to approximately 7,500. Users come to the app through two main channels: customer acquisition through schools, where the *S2M* sales team reaches out to schools that later recommend the app to their students (B2B), and independent users who download the app from the app store (B2C). Additionally, there is a paid version of the app which gives access to some additional non-essential features.

The main content of *Freedom* is short illustrated stories. Stories are organized in different trays based on various themes such as *Trending now* or *Recommended Story*. Figure 1 presents screenshots from the app. Figure 1A shows the landing page that a user sees when launching the app. The landing page contains trays of stories and news, but also occasional promotions and announcements. Figure 1B presents the *Stories* subpage, which contains only stories. The tray displayed at the top is *Recommended Story*.

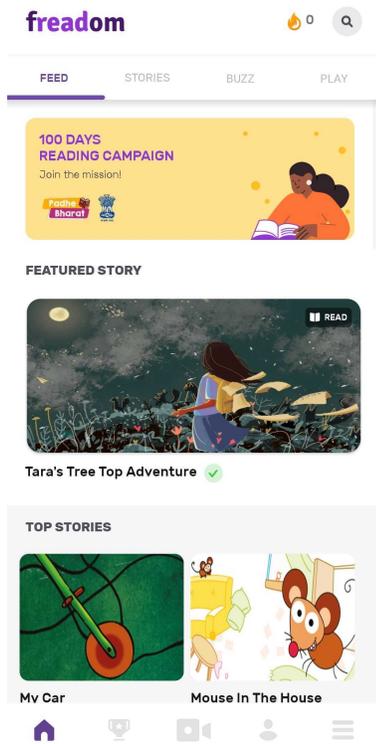
Each tray is a slate of stories that a user can browse, and choose the ones to read. The selection of stories into trays follows various rule-based algorithms. For example, *Trending now* displays stories that are currently consumed by many users. Figure 1C presents the top part of the *Recommended Story* tray. Importantly, during the pre-experimental period, none of the trays of the app assigned students to content in a personalized fashion.

Freedom stories are curated by the *S2M* pedagogical team together with publishers specializing in educational content for kids. They are age-appropriate and created with a pedagogical goal in mind. Therefore, *S2M* operates under the premise that maximizing the consumption of content on the app helps learners achieve their educational goals.

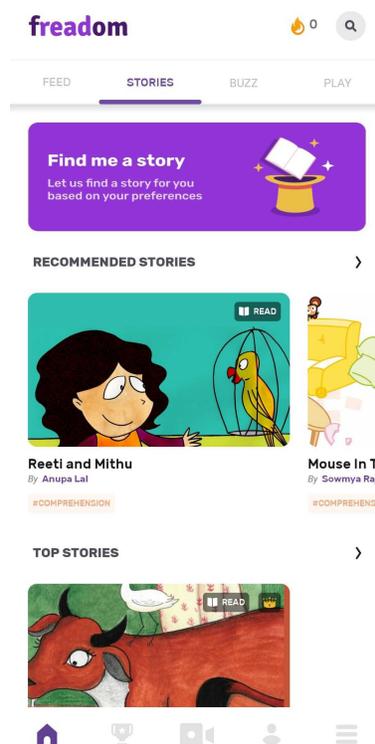
Freedom users can browse stories in the selected tray before deciding on which one to click. Clicking allows the user to open the story and view its description. Many users that view a description decide to go back to browsing; others start the story but do not finish it. Only a small minority of

Figure 1: Screenshots from *Freedom*.

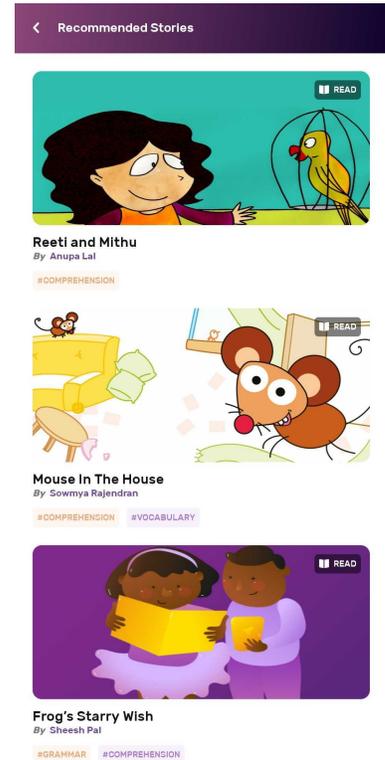
(A) Home Feed page: users open the app on this page.



(B) Stories page: contains all story trays.



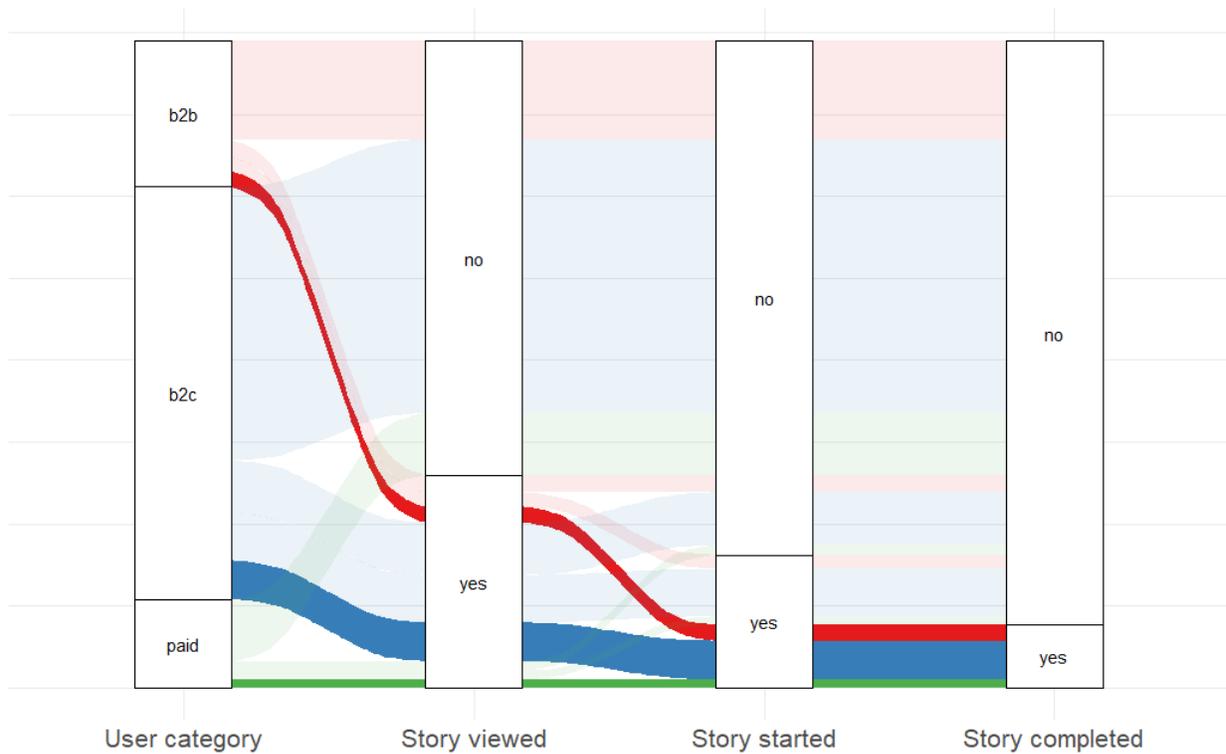
(C) *Recommended Stories*: one of the most popular trays.



user-story interactions lead to the completion of a story.

Figure 2 shows a content interaction funnel representing frequencies of users' content consumption decisions. We divide users into three main categories: *B2C*, *B2B*, and *paid* and show frequencies of different outcomes from interactions with stories. Thus, the unit of interest is the interaction between a user and a story. Users decide whether to view a story or not (second column), whether to start reading it (third column), and whether to complete it or not (final column). We can notice that users tend to explore many stories and acquire information about them through viewing or starting before deciding which stories to complete.

Figure 2: User-story interaction utility funnel



Note: Utility funnel broken by the type of user (B2B, B2C, paid) and the outcome of user-story interactions. In intense colors shares of user-story interactions that resulted in story completion. In red B2B users, in blue B2C, and in green paid users. The first column shows shares of user categories, the second one is the share of users that viewed the story, the third that started the story, and the fourth that completed it.

3 Using offline data to develop a recommendation system

In this section, we aim to describe how we decided on what type of recommendation system we deployed and why. The objective of this section is, on the one hand, to describe the process and decisions taken in the development of the recommendation system that we eventually implemented, but also to act as a guide to practitioners interested in building a similar system who are positioned in a setting similar to ours.

3.1 Target metric and datasets

Our goal. With a story catalog as large as *Freedom's*, it is unpractical for a child to manually choose which stories to consume. Just as in the context of entertainment (movie recommendations) or e-commerce (product recommendations), serving *personalized* recommendations will potentially elicit the most child engagement on the app. Since stories are curated with a focus on pedagogy, this potentially accelerates child learning. Before the experiment, *Freedom* served stories based on editorial

recommendations by experts, which were the same across users with no personalization. Therefore, the goal of our research was to develop a personalized recommendation system and evaluate its efficacy.

Datasets available. We have historical log data of children’s interactions with stories. Every entry of this dataset records an interaction of a child with a story, as well as to what extent they consumed the story; specifically whether a child did not consume it at all, considered reading it by viewing the story description card, started reading it or completed reading it. We have information about a child’s grade level, as well as a tag recording the collections a story belongs to; a collection is a theme such as *animal* or *sport*.⁷

Utility. Based on our interaction data, it is unclear what our goal is in maximizing engagement. There are numerous apparent options; such as maximizing story card view rate or start rate or completion rate. While the ultimate goal of recommending stories to users is that they complete them, viewing and starting a story are prerequisites to completing it, and these are outcomes on a continuum rather than unrelated outcomes. Therefore, we define a metric *utility*, determined together with *S2M* to reflect their organizational objectives. The utility is derived from user-story interactions as follows

- If a story was not shown or was shown to a user who did not interact with any story in that specific session, we do not assign any value: *NA*,⁸
- If the story was shown to the user, but the user skipped the story and viewed another story later in the session: 0,
- If the user viewed the story page, but did not start the story: 0.3,
- If the user viewed and started, but did not complete the story: 0.5,
- If the user viewed, started, and completed the story: 1.

We note that we distinguish between the user choosing not to engage with a story that was shown (0) and the user never having an opportunity to interact with a story because it was not shown (NA),

⁷This is analogous to the type of data in the *movielens* benchmark dataset(Harper and Konstan, 2015); however, a key difference is that our dataset does not contain a rich set of child and story characteristics.

⁸User-item interactions database contains records of only users’ sessions that resulted in at least one click on a story. Thus, sessions in which a user launched the app and skipped all shown stories are not recorded in the data that we have access to.

a critical distinction for understanding user preferences not always made in previous studies. The above utility assignment can be thought of as giving us a utility matrix, with a child represented by a row, and a story by a column. This is the main building block of the recommendation system.

3.2 Recommendation System

We now present how we used observational data on user-story interactions to design the recommendation system for *S2M*. Our dataset does not contain rich user and story characteristics; therefore, we chose a classic collaborative filtering model (Mnih and Salakhutdinov, 2007; Rendle, 2010) as the basis for our personalized recommendation system. We start by describing the collaborative-filtering model.

Model Description Consider two models for our recommendation system. We evaluate them based on out-of-sample performance in the observational data. In what follows, ϵ_{ij} is an unobserved error drawn iid for each child/story pair, and $\sigma(x) = 1/(1 + e^{-x})$, the sigmoid function.

First, is a popularity-based model, (also called a two-way fixed effects model: TWFE).

$$U_{ij} = \sigma(\beta_0 + \Psi_i + \Gamma_j + \epsilon_{ij}) \quad (1)$$

, where Γ_j and Ψ_i are user and story fixed effects, respectively; Note, the popularity-based model is non-personalized, stories are simply ranked by their mean popularity and users receive stories at the top of the rank. We include this model for two reasons: first, it is a useful benchmark for evaluating personalized models. Second, such a model is simpler to implement; thus, to justify the development and introduction of a more complicated personalized model, it is useful to show that simpler models do not achieve similar performance.

Second, our main candidate model is the collaborative filtering approach.

$$U_{ij} = \sigma(\Lambda_j \times \Theta_i + \beta_0 + \Psi_i + \Gamma_j + \epsilon_{ij}) \quad (2)$$

, where Λ_j is a latent preferences vector per user, Θ_i a latent vector per story. This approach follows the seminal model proposed in (Rendle, 2010). The latent vectors, Λ_j and Θ_i , are of length k and are rows and columns of matrices Λ and Θ . Columns of Λ and rows of Θ are k -dimensional representations of user and story latent preference characteristics, respectively.

This approach allows for simplifying the utility matrix. Instead of modeling the preferences of

each user for each story, we express user preferences and story features as each having k -dimensions. These dimensions (or axes of variation) could be thought of as characteristics (e.g., a serious or a funny story); each story and every user are placed along these axes. A higher value in a particular dimension for a story results in a higher expected user valuation in that dimension, and a higher expected preference for that story among users that also have a high value on that dimension. In sum, the collaborative filtering model identifies a low-dimensional representation of both users and stories, so that users with preferences for a particular type of story are located close (in the sense of Euclidean distance) to one another and to their preferred story types.

The collaborative filtering model can achieve high performance, if the matrices Λ and Θ represent underlying preferences well. The more data we have on users' and stories' past interactions, the higher the chance of arriving at an accurate representation of the utility matrix. Crucially, this depends on how well-structured the data is; if there are clear repetitive patterns of user preferences and story types, we are more likely to capture them with this approach.

System Implementation Details. We build our collaborative filtering system using the PyTorch (Paszke et al., 2019) framework in python. We learn our model using Stochastic Gradient Descent (SGD) using the Adam Optimization method. To regularize, we use an L2 penalty on our parameter. We tune the number of latents, k (the dimension of Λ_j and Θ_i), and our L2 penalty parameter using a randomly held out validation set.⁹ Once we have our optimal learning hyperparameters, we relearn our model on the entire dataset which gives us the final model.

Personalized and baseline model performance on offline data. To test the accuracy of the prediction models, we compare the performance of the popularity-based model from Equation (1) to the performance of the collaborative filtering from Equation (2), additionally for completeness we include the performance of a model with just a constant term (mean model).

We compare the performance of these models in terms of Mean Squared Error (MSE) calculated using randomly held out historical data. See Table 1 for results. We find that collaborative filtering outperforms the other models.

⁹We split our dataset into a train, test, and validation set at random. Another approach is to split the data by time into a train dataset and a test dataset; so that we test in the period following our training data. In this setting, the train data is randomly split into a train set and a validation set. We also executed this approach; this leads to similar results.

Table 1: MSE values for collaborative filtering (PYTF), two-way fixed effects (TWFE), and a simple mean model.

PYTF	TWFE	Mean Model
0.0962	0.1022	0.1309

Note: Models are trained and evaluated on the dataset including all users and stories (no filtering based on user-item history length).

Table 2: Collaborative Filtering Model Mean Squared Error for various user and story histories.

		Stories		
		20	60	100
Users	20	0.0967	0.0931	0.0932
	60	0.0964	0.0931	0.0931
	100	0.0959	0.0930	0.0931

Note: The rows represent the minimum interactions per user, and the columns represent minimum interactions per story. We use a single trained model on the largest dataset (20, 20), and report MSEs on different test sets.

Determining the target audience. The performance of the collaborative filtering model depends on the length of histories of interactions of users and stories.¹⁰ To determine the right set of users and stories for the deployment of the recommendation model, we compared the MSEs of utility predictions from the selected utility model trained over different amounts of data and tested on a held-out test set. The training sets differ by the minimum histories of interactions of stories and users. This analysis tells us how much user and story history that is necessary for the recommendation model to provide high-quality recommendations.

Table 2 presents MSEs for nine specifications depending on the length of the history of stories (columns) and of users (rows). We evaluate all specifications on the same dataset with thresholds (20,20).

Based on the results from Table 2, we decided that the population of users that will receive personal recommendations will consist of users and stories with at least 60 interactions in our historical data. Two factors contributed to this decision; first, high-quality predictions as measured by MSE and,

¹⁰Our approach is generally not suitable for new users and new stories. The so-called cold-start problem of assigning content recommendations to users that have not yet revealed preferences from content interactions or stories whose latent style is still unknown is well-documented in recommendation systems literature, see e.g., Lam et al. (2008); Lika et al. (2014); Bobadilla et al. (2012).

second, the sample size requirement for the A/B experiment.¹¹

Choosing the right tray for the new recommendation system. *Freedom* is built based on multiple horizontally scrollable trays. Trays vary by popularity; one important driver of the tray’s popularity is its position on the page. The most popular trays are *Popular*, *Trending Now*, *Recommended Story*, and *Today For you*.

We chose to deploy the recommendation model in the tray called *Recommended Story*. This tray was popular amongst more experienced users of the app, which meant that we were able to deploy the new system to many users of this tray. *S2M* had also originally intended the tray to be for personalized recommendations hence the name - *Recommended Story*. Before deploying our recommendation model, the tray content was chosen by *editors* on a weekly basis.

Re-ranking over time. On our chosen tray *Recommended Story*, stories are presented in a slate of 15 entries. The slate design task consists of deciding how to rank the stories and how frequently to update the ranking. We wanted to keep the ranking and refreshing module similar to the baseline one, so we can focus on isolating the effects of personalization.

Due to computational constraints, new utility predictions were generated once per week. Thus, every week we would rank stories in decreasing order of predicted utility and the top 15 stories would make the slate. Within the week we would remove completed stories every day. Completed stories were replaced by stories that appeared next in the ranking of predicted utility.

Whenever the user was active in the tray for two days but did not engage with any of the top 3 stories, we would remove those stories from their tray. This decision was motivated by the limitations of our data collection process, which does not allow for observing story skipping behavior in the case when the user did not click on any story during the session. This prevents us from accurately determining, the stories that users chose to ignore. Last, the ranking does not change if the user was inactive on the tray. The ranking algorithm was run every day.

4 Randomized controlled trial

A standard approach in recommendation systems literature is to evaluate a counterfactual policy using off-policy evaluation methods (Swaminathan and Joachims, 2015a,b; Gilotte et al., 2018). Conceptually this involves identifying which observed user-item interactions would have also occurred

¹¹Approximately 15% of users in the entire user base and 92% of all stories in the app have at least 60 interactions.

under the counterfactual policy and using observed utilities from these interactions to compute the mean utility under the counterfactual policy. In our context, this is problematic for two reasons. First, recommendation policies impact both the outcomes of user-story interactions and the number of interactions. Thus, the natural metric for evaluating a recommendation policy is the *total utility*, which we cannot reliably estimate with standard off-policy metrics that do not capture the change in the number of interactions (see Kirshenbaum et al. (2012) for a similar argument). Second, even if we abstract away from changes on the extensive margin and assume that the impact of a new recommendation policy can be summarized by additive effects across user-story interactions we will systematically miss some of them. We can adjust for differences between user-item interactions captured in the off-policy evaluation and those in the population at large but this is likely to be incomplete due to data sparsity.

Both these challenges can be thought of as a problem of the overlap between the data generated using the baseline policy and data that would have been generated under the counterfactual policy. Considering a general case with effects on the extensive margin and interaction effects between stories, reliable off-policy estimates can be obtained only when the baseline and counterfactual policies coincide, making the method impracticable.¹² When one is willing to consider the case of simple additive utilities, the extent of the overlap between user-story interactions in the baseline and the counterfactual policy determines how reliable this approach is; Rossetti et al. (2016) and Peska and Vojtas (2020) show these limitations in empirical studies.¹³

An alternative approach to evaluate a new policy is an A/B experiment in which the targeted metric is *total utility* of a user. This is the method we use in this paper. This section discusses the design of the experiment and presents the results.

4.1 Design of the experiment

In the experiment, 7750 users were randomized into treatment and control. We considered only users that had at least sixty story interactions before the experiment. The treatment group received personalized recommendations in the *Recommended Story* tray, while the control group remained with the baseline system of stories selected randomly from a list specified by editors. Tray’s UI was consistent across the control and treatment group; the only thing exogenously varied was the set of stories displayed in the tray. Content presented in the other trays of the app was unchanged. Treated users were

¹²An alternative is to consider a structural model.

¹³By construction this approach is more suitable for evaluating small changes in the policy. A large change that results in new user-story interactions will imply a low overlap.

not aware of the change in the recommendation system. The experiment lasted for two weeks, which was pre-determined with the partner. Based on the analysis of past data the minimum detectable effect on total utility (per user sum of utility over two weeks) was 0.08 standard deviation.

The experiment started on the 22nd of July 2021 and lasted until the 4th of August.¹⁴ During the experiment, 3023 users from the experimental groups launched the app at least once and of them, 525 viewed at least one story in *Recommended Story* tray.¹⁵ We report the balance of observable characteristics between the treatment and control groups in Appendix A.

In the evaluation of the experiment, we consider subjects that launched the app at least once during the experiment. This means that we exclude users that did not launch the app in the experiment period, but we include users who launched the app but did not click on any of the stories in the *Recommended Story* tray. The reason for including the latter group is that users can see the front page of the first story in *Recommended Story* tray without starting to interact with any of the stories in the tray. Thus, we also capture the change from not interacting at all with content in the *Recommended Story* to having some non-zero utility interaction.

4.2 Outcome metrics

We focus on two types of outcomes: first, outcomes specific to *Recommended Story* tray and, second, overall app usage. Even though other trays in the app remained unchanged, we are interested in the impact on overall app usage to understand whether changes in one tray are compensated by altered utilization of content elsewhere, or the overall time spent on the app also shifts. In this specific context where many users are consuming content based on the recommendation of parents or teachers, understanding the overall elasticity of consumption with respect to changes in the app quality is an important, strategic metric that can guide app development.

We consider the following outcome metrics: (i) *total utility* - per user sums of utility from all user-story interactions in *Recommended Story* tray during the experiment, (ii) *total utility all trays* - per user sums of utility from all user-story interactions in all trays of the app during the experiment, (iii) *total stories* - per user sums of completed stories in *Recommended Story* tray during the experiment, (iv) *total stories all trays* - per user sums of completed stories in all trays, (v) *total reading time* - per user sums of estimated reading time of stories completed in *Recommended Story* tray, (vi) *total reading time all trays* -

¹⁴After the experiment, our system of personalized recommendations was launched for all eligible users on the *Recommended Story* tray.

¹⁵The large difference in the number of randomized students and the number of students who were active during the experimental period is because one, a number of students were only active on other trays and two, there is continuous churn and students drop off the app over time.

per user sums of estimated reading time of stories completed in all trays.¹⁶

All metrics relate to total app utilization per user. This approach assigns the same weight to each user without distinguishing between users of varying consumption patterns. In Appendix C, we additionally consider mean utility from user-story interactions.

We constructed all variables based on raw log files provided by S2M. These log files are internal data used by S2M data analytics teams, they constitute the most accurate available picture of users’ behavior on the platform. Nevertheless, occasional instrumentation errors occur. The type of instrumentation errors that are problematic for our analysis is an incorrect attribution of user-story interactions.¹⁷ This results in some users having spurious, very high utilization during specific sessions. To avoid including such sessions in the analysis we drop users that had at least one session in which they completed more than 10 stories. In result, we drop 40 users. Table 3 provides summary statistics of variables describing utilization in the *Recommended Stories* tray.

Table 3: Summary statistics of outcome variables describing activity on the *Recommended Stories* tray per group.

names	group	min	mean	percentile 75th	percentile 90th	percentile 95th	max
Total utility	control	0	0.28	0	0.51	1.6	13.6
Total utility	treatment	0	0.45	0	1.30	3.0	23.9
Total stories	control	0	0.15	0	0.00	1.0	11.0
Total stories	treatment	0	0.27	0	1.00	2.0	21.0
Total reading time	control	0	1.04	0	0.00	7.5	78.5
Total reading time	treatment	0	1.94	0	7.25	11.0	148.5

Note: Summary statistics of variables measuring utilization of the Recommended Story tray during the experiment. Sample includes only users that launched the app during the experiment period.

Even though we consider only users that launched the app during the experiment, most of them had zero utilization of the app in the *Recommended Story* tray. Nevertheless, we still include them in the experiment evaluation as different recommendation policies might impact the share of users consuming any content in the tray. From Table 3 we can notice that the treatment group has higher mean utilization and higher utilization on the 90th and 95th percentiles.

We are also interested in the impact of the personalization of content recommendations in *Recommended Stories* tray on the overall app usage. Table 4 presents summary statistics of variables describing utilization on all trays in the app.

¹⁶The estimates of the reading time per story are provided by S2M as intervals, e.g., from two to four minutes. For each story we take the mid point of the interval.

¹⁷This can for example take a form of a user being assigned interactions of another user, or assigned completions instead of views.

Table 4: Summary statistics of outcome variables describing activity on all trays per group.

names	group	min	mean	percentile 75th	percentile 90th	percentile 95th	max
Total utility	control	0	3.79	4.47	11.4	17.77	81.5
Total utility	treatment	0	4.32	5.50	13.5	19.02	63.7
Total stories	control	0	1.89	2.00	6.0	9.00	41.0
Total stories	treatment	0	2.25	3.00	7.0	11.00	42.0
Total reading time	control	0	12.65	14.50	40.0	65.32	273.0
Total reading time	treatment	0	15.15	15.00	46.0	73.12	365.0

Note: Summary statistics of variables measuring the overall app utilization during the experiment. Sample includes only users that launched the app during the experiment period.

In Table 4 we see that mean outcomes are higher in the treatment group for all outcome variables. Treatment has higher or equal outcomes at the 75th, 90th, and the 95th percentile.

To compare distributions of total utility in treatment and control we carry out Wilcoxon test (one sided alternative). Using the total utility in the *Recommended Story* tray we reject the hypothesis that the true location shift is less than zero, with p-value 0.0007, and for all trays in the app with p-value of 0.05. In Figure 3 we present entire distributions of total utility. Panel A shows cumulative distribution functions of total utility from *Recommended Story* tray per experimental group; panel B shows the difference between probability density functions of the treatment and the control group. We can notice that a larger share of control group users did not have any positive-utility content interaction during the experiment. Treatment group has a higher probability mass for almost any non-zero utility.

4.3 Average treatment effects

Estimates of the average treatment effects are presented in Table 5. We use the difference in means, the linear regression, and the augmented inverse propensity weighing (AIPW) estimators.

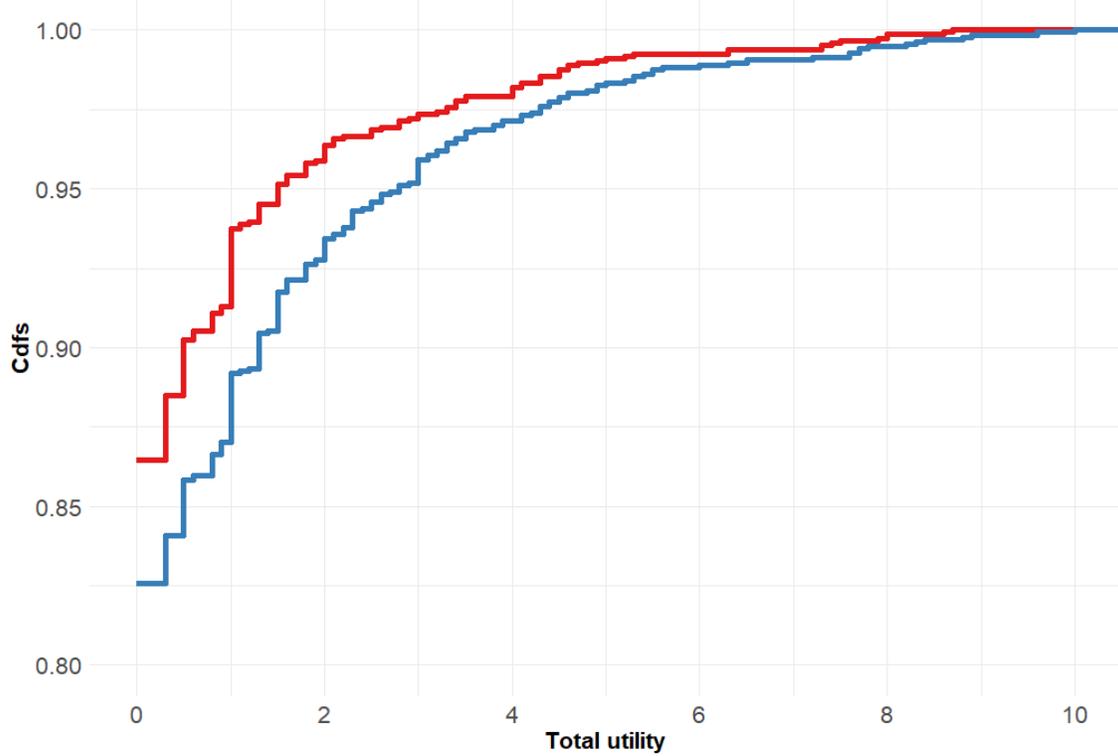
We find a strong positive effect of personalization on all outcomes metrics. The impact on utilization of the *Recommended Stories* tray has high economic and statistical significance. Total utility increases by 63% ($\pm 28\%$), the number of stories completed in the tray by 78% ($\pm 39\%$), and total reading time by 87% ($\pm 41\%$).¹⁸

We also find an increase in the utilization of the app across all trays; total utility increases by 14% ($\pm 12\%$), the number of stories completed by 19% ($\pm 14\%$), and the reading time in all trays by 20% ($\pm 14\%$). Thus, the increase of consumption of content in *Recommended Stories* did not come entirely at the expense of consumption in other trays; on the contrary, this evidence suggest that users started

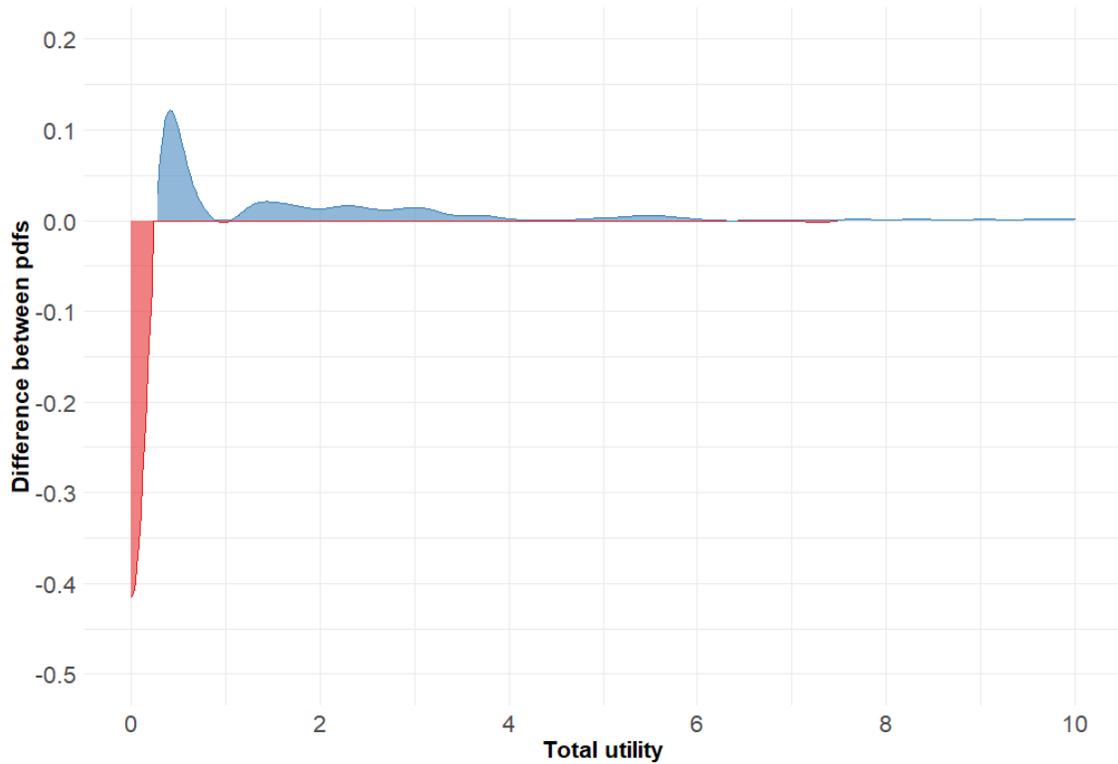
¹⁸Confidence intervals in brackets. Standard errors based on difference in means estimator.

Figure 3: Distribution of total utility in *Recommended Stories* tray per group.

(A) Cumulative distribution function per group. Treatment in blue, control in red.



(B) Difference between the probability density functions of treatment and control groups. Treatment in blue, control in red.



using the app more.¹⁹

Table 5: Estimates of average treatment effects for all outcome variables

variable	ATE	std.err.	p.value	ATE %	ATE reg adj.	std. err. reg adj.	ATE AIPW adj.	std. err. AIPW adj.
Total utility RS	0.17	0.05	0.00	60	0.18	0.05	0.18	0.05
Total stories RS	0.12	0.04	0.00	78	0.13	0.03	0.13	0.04
Total reading time RS	0.90	0.26	0.00	87	0.96	0.25	0.98	0.25
Total utility all trays	0.52	0.27	0.05	14	0.50	0.26	0.51	0.26
Total stories all trays	0.36	0.16	0.03	19	0.36	0.15	0.35	0.15
Total reading time all trays	2.50	1.10	0.02	20	2.47	1.07	2.49	1.06

Note: Estimates of the average treatment effect using difference-in-means estimator (first column), adjusting for covariates with a linear regression (fifth column), and adjusting for covariates using Augmented Inverse Propensity Weighting - AIPW (column seven); covariates used: users' grade, user type (B2B, B2C, or paid), past utilization, niche type (indicator whether user consumes content that is popular amongst other users or more niche content), past usage of the Recommended Story tray. Columns two, six, and eight show standard errors. Column three presents p-values. Three first rows describe outcomes in Recommended Story tray, three bottom rows overall app utilization.

Additionally, we review differences in total utility in the most popular trays in the app across treatment and control. Figure 4 shows differences in average total utility in treatment and control groups in other popular trays in the app. The experiment period is marked in blue; we can see that the difference between the two groups is statistically significant only for Recommended Story tray. We carry out this comparison for the same users in a pre-experiment period; before the experiment, differences in average utility across treatment and control are insignificant in all of the trays (which is expected since the users were randomly assigned).

Impact on time spent on the app. In Table 5, we see a strongly significant positive effect on the time spent, both in *Recommended Story* tray as well as across all trays.²⁰

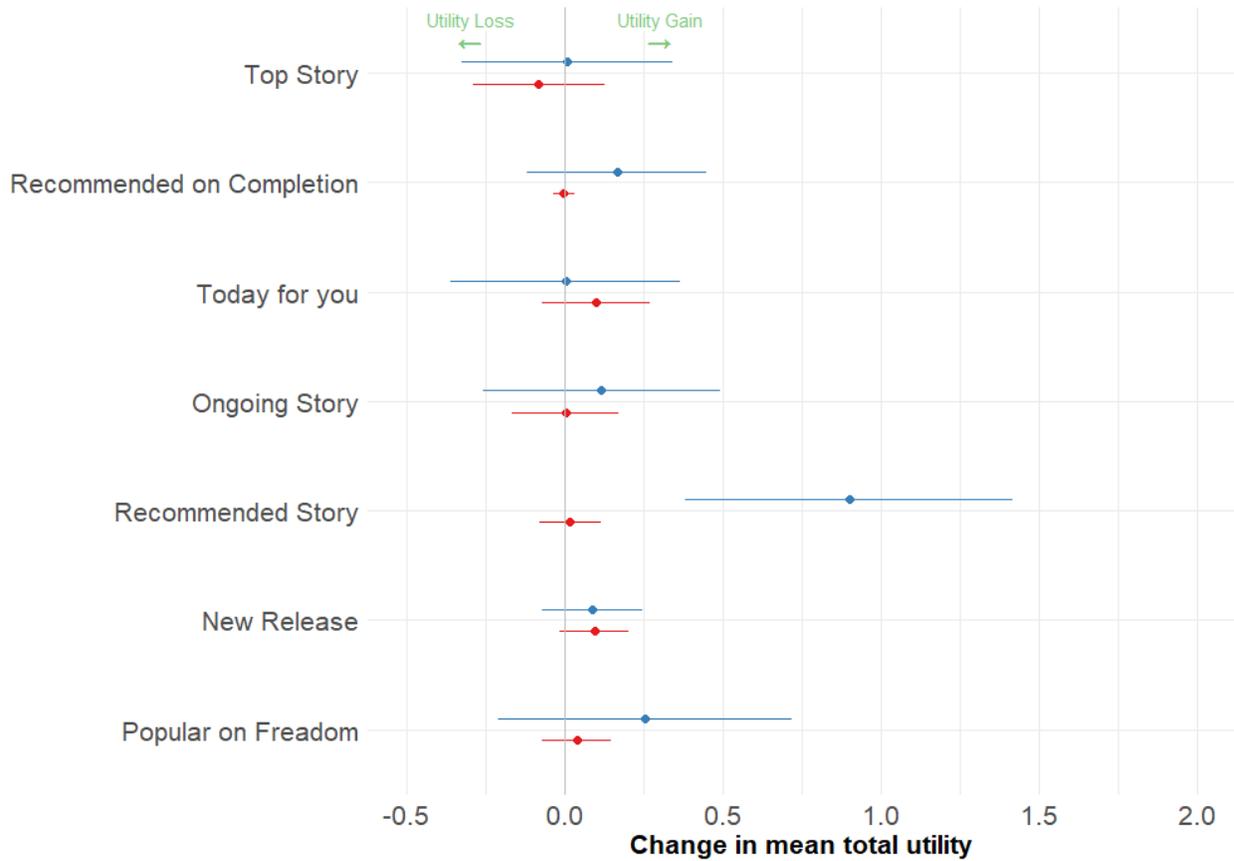
The increase in total time spent on the app is particularly interesting because it means that students prefer to spend time on the app than engage in other activities, outside of the app. In our context, this result suggests that if the content is interesting to students, they are willing to go beyond the time prescribed by parents or teachers.

Generally, we can consider users responding to the improvements in the app quality on an intensive and extensive margin. Gains on the intensive margin would be due to users better allocating their time; in our case, that is reallocation of the time to more attractive, personalized content in the *Recommended Story* tray. While the impact on the extensive margin means that users substitute away

¹⁹In Appendix B we provide robustness check of this estimates by trimming the top 5% users with the highest daily number of completed stories instead of the cap on 10 stories.

²⁰Note, this outcome metric is a sum of the duration of completed stories. We do not include stories that were started, but not completed, since we do not observe the moment in which users stopped engaging with a specific story. The average number of stories started but not completed in treatment and control is roughly the same 2.33 in treatment and 2.17 in control; the test for difference in means has p-value of 0.55.

Figure 4: Difference in average total utility in treatment and control groups for eight most popular trays.



Note: Difference in average total utility in treatment and control groups for eight most popular trays. Experimental period in blue, pre-experimental in red. Pre-experimental period is 7-19.06.2021, there are approximately twice as many users in the pre-experimental period (this date is chosen on the basis of being the closest two-weeks long period without other major experiments and alterations in the app).

from other activities and start using the app more. The effect on the extensive margin highlights that the app quality matters to the users, and improving it will result in more time spent with the app.

4.4 Heterogeneous treatment effects

The evidence presented so far relates to the average impact of personalization. In this section, we analyze heterogeneity in treatment across past usage intensity, taste for popular vs. niche content, and the usage of the *Recommended Story* tray prior to the experiment. In Appendix E, we carry out a data-driven analysis of treatment heterogeneity and find a moderate amount of treatment heterogeneity.

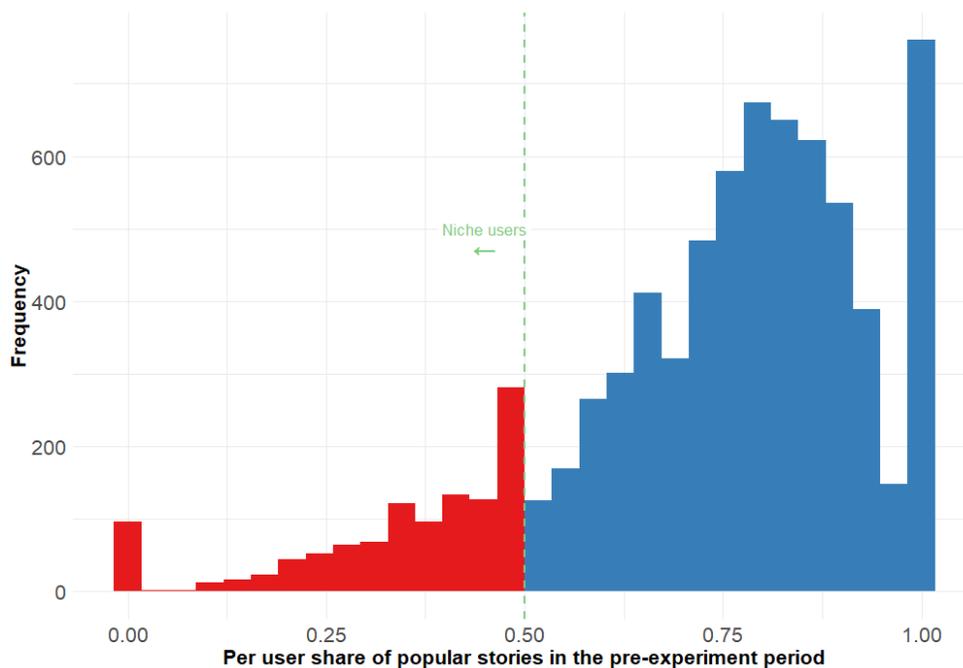
We expect that the personalization of content recommendations will mostly benefit heavy and niche-type users. Frequent users leave a long record of user-story interactions, which allows us to well understand their tastes. Additionally, we expect niche users to have high benefits because, in the

baseline system, stories are targeted at a typical user, whereas in the personalized system, their niche tastes are taken into account.

Definitions of users’ types. To determine whether someone is a heavy user we analyze the pre-experimental app usage. For each user, we compute the total utility and the total number of completed stories prior to the start of the experiment. Additionally, we construct indicator variables: *high utility user* and *high story completion user*, which take a value of one when a user is in the top 50th percentile of the distribution of past utilization (past number of completed stories) and zero otherwise.

Niche-type users are users that consume content that is generally not very popular. We consider a story to be a popular story if it is one of the top 25% of stories in terms of pre-experiment completions.²¹ Figure 5 shows the histogram of shares of popular content consumption per user prior to the experiment. There are some users whose content is largely niche. We consider a user to be a niche type if the share of niche content in her pre-experiment consumption is more than 50% (in red in Figure 5). Note, that all users were receiving the same recommendations prior to the experiment; thus, finding niche stories required searching beyond the top of the recommendation list.

Figure 5: Histogram of the share of popular stories consumed by users.



Note: A popular story is a story in the top 25% of stories ranked by the number of pre-experiment completions. Niche users in red.

²¹Top 25% of stories correspond to 67% of impressions in the *Recommended Story* during the experiment.

Treatment effects per group. We start by providing estimates of the average treatment effects per group of interest. We consider total utility in *Recommended Story* tray as the outcome variable of interest and use a difference-in-means estimator. Table 6 presents the results.

Table 6: Estimates of average treatment effects per group.

category	group	ATE	std. error	p. value
Type	Niche users	0.334	0.089	0.000
	Non-niche users	0.044	0.063	0.487
Past utilization	High utility users	0.299	0.094	0.001
	Low utility users	0.092	0.056	0.103
	High story completion users	0.249	0.094	0.008
	Low story completion users	0.120	0.056	0.032
Type and past utilization	High utility and niche users	0.508	0.132	0.000
	High utility and not niche users	-0.023	0.122	0.848

Note: Outcome variable is total utility per user. ATE is estimated using a difference-in-means estimator. All groups are defined based on the pre-experiment app usage.

We find that the gains from personalization are higher for niche users than for non-niche users and for heavy users than for light users. The niche dimension is of higher magnitude and statistical significance. In the last two rows of Table 6, we focus on the distinction between niche and non-niche users in the heavy utility group and find that niche users in this group have much higher treatment effects. This highlights, that the niche users form a distinct category, rather than are just heavy users who completed all popular stories and need to explore less popular ones.²² In Appendix D we provide further robustness of this result by regressing AIPW scores on past utilization and user type.²³

Last, in Figure 6 we show how AIPW scores change across users depending on their past utilization. Panel A shows how AIPW scores change depending on the percentile of the pre-experiment number of story completions and panel B on users' past utility. We can notice upward trends in both figures. The differences are, however, moderate.

Niche-type users see more niche content. Personalized recommendations benefit niche users because they do not need to seek out their favorite niche stories away from the top of the list of recom-

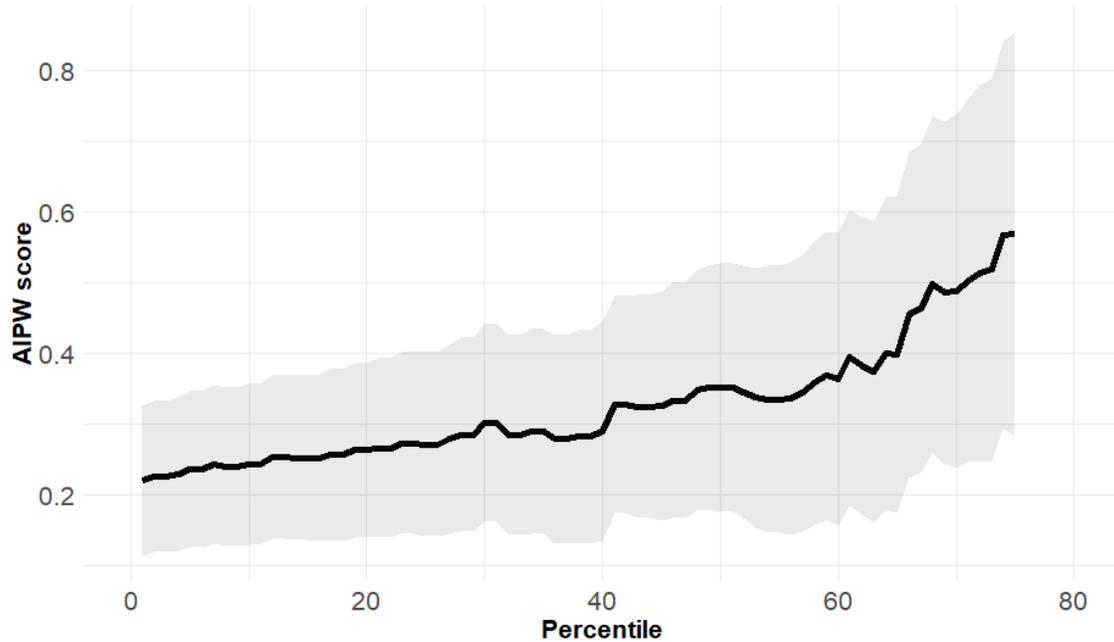
²²One might argue that a user becomes niche after having seen all the popular stories. Note, that there are 839 users in the heavy utility and niche and 573 in heavy utility and non-niche. This indicates that niche users are indeed a distinct category of users.

²³To estimate AIPW scores we use the *grf* package (see Athey et al. (2019)). This methodology allows us to flexibly adjust for individual characteristics and estimate conditional average treatment effects. We consider users' school grade, type (B2B, B2C, paid), max streak (maximal number of consecutive days in which users completed at least one story), past utilization (the total number of completed stories prior to the experiment, and total utility prior to the experiment), and whether a user is a niche type. To determine the variables based on past consumption we consider a period of app usage between July 2020 and the start of the experiment.

Figure 6: AIPW scores across past utilization. AIPW scores for users with past utility higher than the percentile.

(A) Past story completions. AIPW scores for users with the past number of completions higher than the percentile.

Heterogeneity across past completions



(B) Past utility. AIPW scores for users with past utility higher than the percentile.

Heterogeneity across past utility

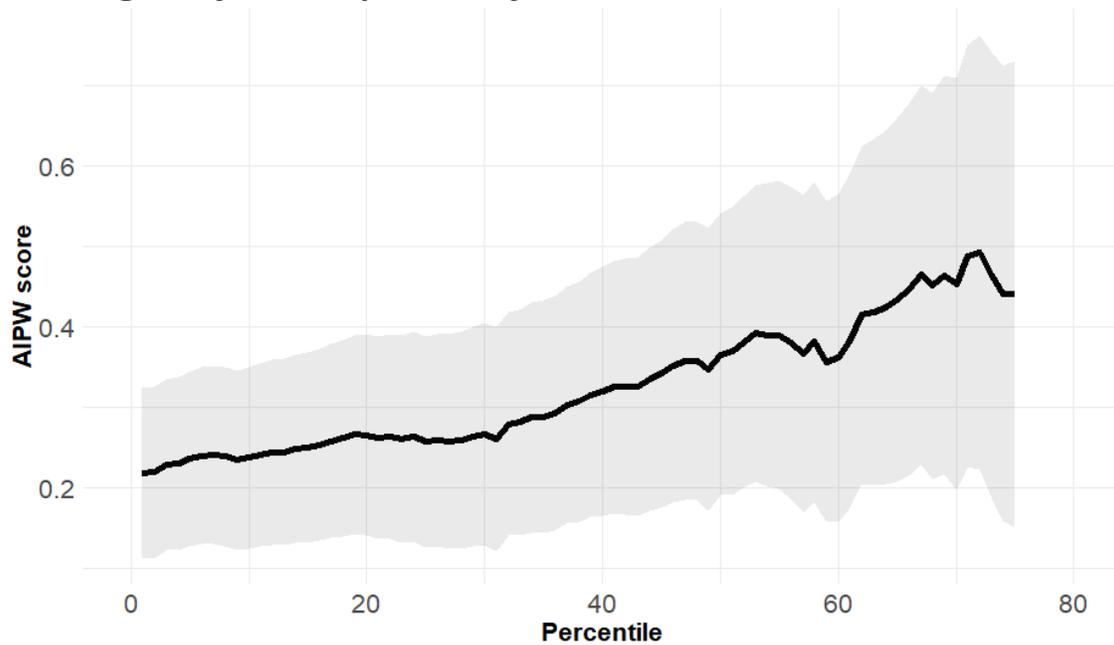


Table 7: Type of stories shown to niche and popular-type users across treatment and control.

group	variable	mean niche	mean non-niche	std. error	p. value
Treatment	Rank of impressions	379.184	343.442	17.313	0.040
Control	Rank of impressions	498.730	489.391	18.357	0.611
Treatment	Percentile of impressions	0.390	0.447	0.028	0.040
Control	Percentile of impressions	0.401	0.412	0.022	0.611

Note: Type of stories shown to niche and popular-type users across treatment and control. The rank of impressions - stories ranked by the number of impressions during the experiment in the experimental group. Percentile refers to the percentile of the distribution of the share of impressions per story in the total impressions in the experimental group.

mended stories, but receive them right away. In Table 7, we confirm this intuition by comparing the popularity of stories shown to popular and niche types in the two experimental groups.

For each story, we compute the share of its impressions in total impressions in an experimental group and rank stories by it (*Rank of impressions*). Additionally, we compute each story’s percentile in the distribution of impressions within the experimental group (the total number of stories per experimental group differs).

In the control group, popular and niche type users see stories of similar popularity, while in the treatment group niche users are shown more niche stories; the difference is statistically significant.

New and old users of *Recommended Story* tray. Another important layer of heterogeneity is between users that have been consuming stories in *Recommended Story* tray before the experiment and those that started using this tray because of the personalized recommendations. Out of all experimental subjects only 14% interacted with at least one story from the *Recommended Story* tray in the two weeks prior to the experiment, and 47% have never interacted with a story in this tray.

Table 8 presents estimates of conditional average treatment effects. We consider only outcomes specific to the utilization of *Recommended Story* tray. Three top rows present results for users that have interacted with at least one story in the *Recommended Story* tray in the two weeks prior to the experiment. We find high and statistically significant treatment effects for this group.

The three bottom rows of table 8 present the results for users that were not actively using this tray prior to the experiment. We find that the treatment effects for such users are highly statistically significant and have high economic magnitudes. While the point estimates are small, the percentage change compared to the baseline (usage in the control group) is very high. These results suggest that the introduction of personalized recommendations attracted users to the tray that otherwise would

Table 8: ATE by past usage of *Recommended Story tray*.

group	variable	ATE	ATE % baseline	std.error	p.value
Past users of RS	Total utility	0.788	54.461	0.322	0.015
Past users of RS	Total stories	0.497	56.182	0.248	0.046
Past users of RS	Total time reading	3.980	65.417	1.777	0.026
New RS users	Total utility	0.132	87.423	0.039	0.001
New RS users	Total stories	0.097	143.136	0.026	<0.001
New RS users	Total time reading	0.692	148.707	0.189	<0.001

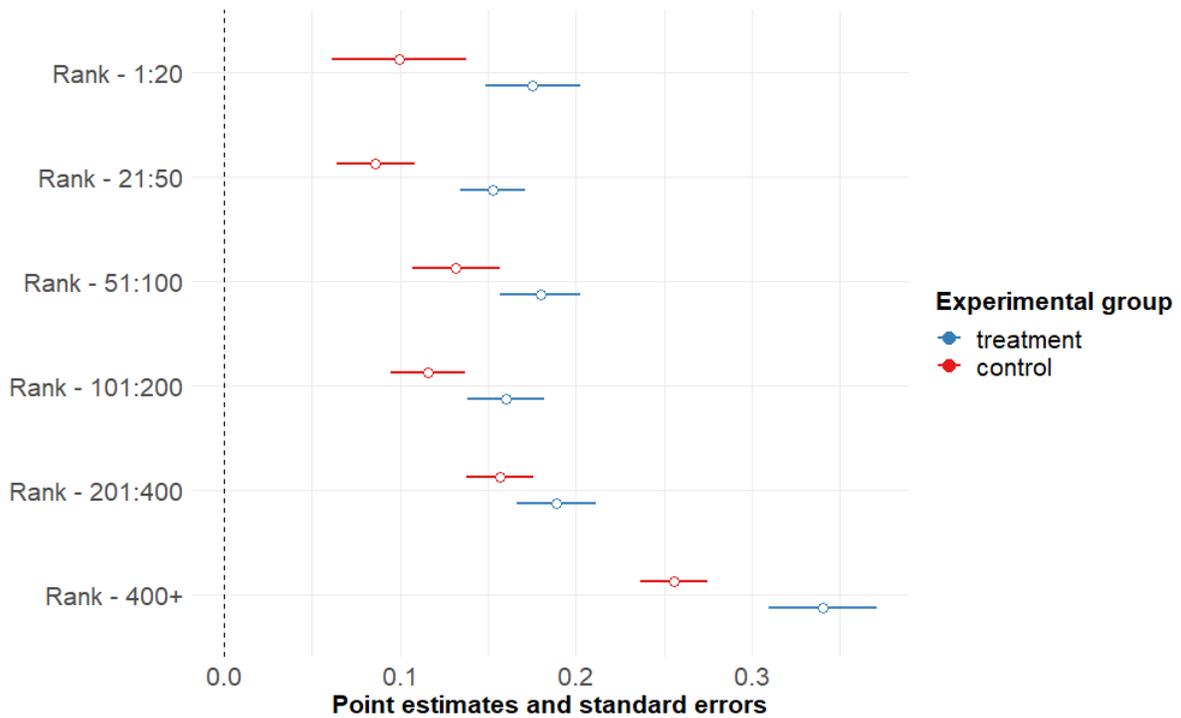
Note: Average treatment effects estimates using Difference-in-Means estimator. Subjects were grouped based on the usage of Recommended Story tray two weeks prior to the experiment. Three first rows show results for users that viewed at least one story in the tray; three bottom rows users that did not interact with any stories in the Recommended Story tray in this period.

not be using it at all.

Stories that drive the treatment effect. Is the increase in total utility driven by a few stories liked by many users or a better assignment of many stories? To answer this question, we want to group stories into frequently and rarely shown and compare user utilities in these categories, in both the treatment and control groups.

Figure 7 shows estimates of the conditional expectation of utility from user-story interactions for stories in different buckets of popularity. We use a linear regression where we adjust for users' grade, type, and past utilization. Buckets are constructed according to the rank of the number of story impressions in the experimental group (the total number of impressions in the treatment group is approximately equal in each bucket). Differences across experimental groups in the average utility in a bucket are (apart from personalization) due to, the selection of stories into buckets and differences in users that see stories in these buckets. Adjusting for user features allows us to isolate the effect of the story selection.

Figure 7: Estimates of the conditional expectation of utility per bucket.



Note: Utility estimates adjusted for the difference in grades, user types, and past usage intensity across buckets.

We find that utilities in the treatment group are higher in all buckets. There is a high and statistically significant difference in the two first buckets. This suggests that our model picked up stories that were liked by many users. However, there is also a substantial increase in utility from the least impressed, niche stories. This means that there is a component of personalized niche content driving higher utility in the treatment group. In sum, we see that there are two mechanisms in story selection that increase the utility in the treatment group: (i) stories that are shown to many users on average lead to higher utility in the treatment group, and (ii) personalization of niche, infrequent stories in the treatment group leads on average to higher utility from interactions with these stories.

5 Conclusion

In this paper, we provide evidence from a randomized controlled trial of the efficacy of personalized recommendations in promoting user engagement on an ed-tech app. We show that children learning to read in English engage more with content when it is selected based on their preferences. We find an effect of an over 60% increase in the utilization of the personalized content as compared to the baseline system of content selected by editors. We also find a 15% boost in overall app usage.

We evaluate the effects of the treatment on different user subgroups in the experiment and find interesting patterns of heterogeneity. We find that heavy users have substantially higher treatment effects. We have more data about such users; thus, we know their preferences better and can provide them with higher-quality recommendations. Second, we find that users that ex-ante prefer niche stories are the main beneficiaries of the personalized system; we also find that the personalized recommendation system makes it easier for them to discover niche content on the platform. Third, we show that both users who have been using the personalized section of the app prior to the experiment as well as those who have not benefited from the personalization.

We examine whether the increased utilization comes with increased diversity, and find that while the recommendation algorithm picks up on stories that are popular, it also increases utility from the least shown, niche stories.

This paper contributes to the recommendation systems literature by bringing evidence from the educational sector and a setting with limited data (as compared to big-tech environments where such systems are typically deployed). We carefully discuss the recommendation system design process hoping to allow practitioners to develop and deploy similar recommendation systems in other contexts.

The main limitation of this paper is that we focus on students that are heavy app users (interacted with at least sixty stories) and on stories that have been already shown to many users. This is a limitation of any system based on the collaborative filtering model as the model's performance improves with the number of past users-content interactions. Furthermore, the approach is not applicable to new users and new stories. Developing and implementing recommendations for new users and new items is a valuable extension of this work.

Last, the proposed approach optimizes for user engagement rather than for learning. The recommendation system assigns stories that the user is most likely to complete, but these might not necessarily be the stories that will maximize learning. Optimizing the story selection for learning would be a preferable approach; however, because of difficulties in accurately measuring learning outcomes and slower feedback loops, we focused on engagement.²⁴ Bridging the gap between optimizing for short-term outcomes vs. long-term learning, for example by using surrogates, is a promising next step on this research agenda.

²⁴This relates to the literature on surrogate (Yang et al., 2020), where a surrogate metric that closely tracks the target metric is optimized instead due to the target metric being infeasible to access

References

- Anderson, A., Maystre, L., Anderson, I., Mehrotra, R., and Lalmas, M. (2020). Algorithmic effects on the diversity of consumption on Spotify. In *Proceedings of The Web Conference 2020*, pages 2155–2165.
- Athey, S., Tibshirani, J., and Wager, S. (2019). Generalized random forests. *The Annals of Statistics*, 47(2):1148–1178.
- Athey, S. and Wager, S. (2019). Estimating treatment effects with causal forests: An application. *Observational Studies*, 5(2):37–51.
- Bobadilla, J., Ortega, F., Hernando, A., and Bernal, J. (2012). A collaborative filtering approach to mitigate the new user cold start problem. *Knowledge-based systems*, 26:225–238.
- Chernozhukov, V., Demirer, M., Duflo, E., and Fernandez-Val, I. (2018). Generic Machine Learning Inference on Heterogeneous Treatment Effects in Randomized Experiments, with an Application to Immunization in India. Technical report, National Bureau of Economic Research.
- Claussen, J., Peukert, C., and Sen, A. (2021). The editor and the algorithm: Returns to data and externalities in online news. *Available at SSRN 3479854*.
- Davidson, J., Liebald, B., Liu, J., Nandy, P., Van Vleet, T., Gargi, U., Gupta, S., He, Y., Lambert, M., Livingston, B., et al. (2010). The youtube video recommendation system. In *Proceedings of the fourth ACM conference on Recommender systems*, pages 293–296.
- Drachsler, H., Hummel, H., Berg, B., Eshuis, J., Waterink, W., Nadolski, R., Berlanga, A., Boers, N., and Koper, R. (2008). Effects of the isis recommender system for navigation support in self-organized learning networks. *Journal of Educational Technology and Society*, 12.
- Drachsler, H., Verbert, K., Santos, O. C., and Manouselis, N. (2015). Panorama of recommender systems to support learning. In *Recommender systems handbook*, pages 421–451. Springer.
- Escueta, M., Nickow, A. J., Oreopoulos, P., and Quan, V. (2020). Upgrading education with technology: Insights from experimental research. *Journal of Economic Literature*, 58(4):897–996.
- Garrett, N. (2009). Computer-assisted language learning trends and issues revisited: Integrating innovation. *The modern language journal*, 93:719–740.

- Gilotte, A., Calauzènes, C., Nedelec, T., Abraham, A., and Dollé, S. (2018). Offline A/B testing for recommender systems. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pages 198–206.
- Gomez-Uribe, C. A. and Hunt, N. (2015). The Netflix Recommender System: Algorithms, Business Value, and Innovation. *ACM Transactions on Management Information Systems (TMIS)*, 6(4):1–19.
- Greenstein-Messica, A. and Rokach, L. (2018). Personal price aware multi-seller recommender system: Evidence from ebay. *Knowledge-Based Systems*, 150:14–26.
- Haim, M., Graefe, A., and Brosius, H.-B. (2018). Burst of the filter bubble? effects of personalization on the diversity of Google News. *Digital journalism*, 6(3):330–343.
- Harper, F. M. and Konstan, J. A. (2015). The movielens datasets: History and context. *Acm transactions on interactive intelligent systems*, 5(4):1–19.
- Holtz, D., Carterette, B., Chandar, P., Nazari, Z., Cramer, H., and Aral, S. (2020). The engagement-diversity connection: Evidence from a field experiment on Spotify. In *Proceedings of the 21st ACM Conference on Economics and Computation*, pages 75–76.
- Jacobson, K., Murali, V., Newett, E., Whitman, B., and Yon, R. (2016). Music personalization at Spotify. In *Proceedings of the 10th ACM Conference on Recommender Systems*, pages 373–373.
- Kirshenbaum, E., Forman, G., and Dugan, M. (2012). A live comparison of methods for personalized article recommendation at Forbes.com. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 51–66. Springer.
- Lam, X. N., Vu, T., Le, T. D., and Duong, A. D. (2008). Addressing cold-start problem in recommendation systems. In *Proceedings of the 2nd international conference on Ubiquitous information management and communication*, pages 208–211.
- Lee, D. and Hosanagar, K. (2019). How Do Recommender Systems Affect Sales Diversity? A Cross-Category Investigation via Randomized Field Experiment. *Information Systems Research*, 30(1):239–259.
- Lika, B., Kolomvatsos, K., and Hadjiefthymiades, S. (2014). Facing the cold start problem in recommender systems. *Expert Systems with Applications*, 41(4):2065–2073.

- Linden, G., Smith, B., and York, J. (2003). Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing*, 7(1):76–80.
- Mnih, A. and Salakhutdinov, R. R. (2007). Probabilistic matrix factorization. *Advances in neural information processing systems*, 20.
- Möller, J., Trilling, D., Helberger, N., and van Es, B. (2018). Do not blame it on the algorithm: an empirical assessment of multiple recommender systems and their impact on content diversity. *Information, Communication & Society*, 21(7):959–977.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Peska, L. and Vojtas, P. (2020). Off-line vs. on-line evaluation of recommender systems in small e-commerce. In *Proceedings of the 31st ACM Conference on Hypertext and Social Media*, pages 291–300.
- Rendle, S. (2010). Factorization machines. In *2010 IEEE International conference on data mining*, pages 995–1000. IEEE.
- Rossetti, M., Stella, F., and Zanker, M. (2016). Contrasting offline and online results when evaluating recommendation algorithms. In *Proceedings of the 10th ACM conference on recommender systems*, pages 31–34.
- Ruiz-Iniesta, A., Melgar, L., Baldominos, A., and Quintana, D. (2018). Improving children’s experience on a mobile edtech platform through a recommender system. *Mobile Information Systems*, 2018:1374017.
- Sharma, A., Hofman, J. M., and Watts, D. J. (2015). Estimating the causal impact of recommendation systems from observational data. In *Proceedings of the Sixteenth ACM Conference on Economics and Computation*, pages 453–470.
- Smith, B. and Linden, G. (2017). Two decades of recommender systems at Amazon.com. *Ieee internet computing*, 21(3):12–18.
- Sun, F., Yu, M., Zhang, X., and Chang, T.-W. (2020). A vocabulary recommendation system based on knowledge graph for chinese language learning. In *2020 IEEE 20th International Conference on Advanced Learning Technologies (ICALT)*, pages 210–212. IEEE.

- Swaminathan, A. and Joachims, T. (2015a). Batch learning from logged bandit feedback through counterfactual risk minimization. *The Journal of Machine Learning Research*, 16(1):1731–1755.
- Swaminathan, A. and Joachims, T. (2015b). The self-normalized estimator for counterfactual learning. *Advances in Neural Information Processing Systems*, 28.
- Tafazoli, D., Huertas Abril, C. A., Gómez Parra, M. E., et al. (2019). Technology-based review on computer-assisted language learning: A chronological perspective. *Pixel-Bit*.
- Ursu, R. M. (2018). The power of rankings: Quantifying the effect of rankings on online consumer search and purchase decisions. *Marketing Science*, 37(4):530–552.
- Xie, H., Wang, M., Zou, D., and Wang, F. L. (2019). A personalized task recommendation system for vocabulary learning based on readability and diversity. In *International conference on blended learning*, pages 82–92. Springer.
- Yang, J., Eckles, D., Dhillon, P., and Aral, S. (2020). Targeting for long-term outcomes. *arXiv preprint arXiv:2010.15835*.
- Zhan, R., Ren, Z., Athey, S., and Zhou, Z. (2021). Policy learning with adaptively collected data. *arXiv preprint arXiv:2105.02344*.
- Zhao, Y. (2003). Recent developments in technology and language learning: A literature review and meta-analysis. *CALICO journal*, pages 7–27.

Appendix

A Covariate balance check

Table 9 presents comparison of means of user characteristics across treatment and control. We find that difference between treatment and control are small and statistically insignificant.

Table 9: Balance of covariates across treatment and control

covariate	mean treatment	sd treatment	mean control	sd control	p value
past utility	101.23	139.38	94.67	114.35	0.17
past stories	57.14	89.03	52.75	77.22	0.16
max streak	18.58	85.00	17.79	84.00	0.80
share b2b	0.31	0.46	0.29	0.46	0.46
share b2c	0.41	0.49	0.41	0.49	0.79
share paid	0.25	0.43	0.26	0.44	0.66
share grade 2	0.24	0.43	0.24	0.43	0.82
share grade 3	0.22	0.42	0.22	0.41	0.69

Note: Means of users' characteristics in treatment and control. Last column p-value from a t.test for difference in means. Category paid includes users from a paid fLive program and regular paying users; category b2b includes regular b2b customers and club 1br users, a B2B promotion.

B Robustness check of the average treatment effect estimates

In table 10 we present estimates of the average treatment effect based on data which is trimmed at the 95% percentile of daily stories completed, i.e., we remove users that are in top 5% of users with highest daily number of completed stories across all paths of the app.

We find very similar estimates of the ATE for outcomes across all paths in the app. The path specific estimates are smaller, but still high and statistically significant. The confidence intervals include the point estimates from the baseline specification.

Table 10: Estimates of average treatment effects for all outcome variables

variable	ATE	std.error	p.value	ATE percentage
Total utility RS	0.14	0.04	<0.001	58
Total stories completed RS	0.09	0.03	<0.001	68
Total reading time RS	0.67	0.21	<0.001	77
Total utility all paths	0.50	0.23	0.03	15
Total stories completed all paths	0.32	0.13	0.01	21
Total reading time all paths	2.07	0.88	0.02	20

Note: Estimates of the average treatment effect using difference-in-means estimator. Three first rows describe outcomes in Recommended Story path, three bottom rows overall app utilization. Last columns shows the ATE estimate as a percent share of the baseline.

Table 11: Average treatment effects: alternative utility metrics

variable	ATE	std. error	p. value	ATE %
Mean utility RS	0.015	0.006	0.013	0.31
Utility RS	0.006	0.013	0.626	0.01

Note: Estimates of the average treatment effect using difference-in-means estimator. First row mean utility per user in Recommended Story path (mean of the mean utilities within the group); only users that launched the app considered. Second row mean utility per user-story interaction.

C Alternative utility metrics.

The utility metric that we analyzed so far is the per user sum of utility from all user-story interactions during the experiment period. This metric assigns the same weight to each user, irrespective of the number of stories consumed by that user. It also captures the fact that a new policy might impact the number of stories consumed by users. However, a firm introducing a new recommendation system might have a different objective, for example to weigh each user-story observation equally, or simply focus on maximizing the mean utility each user receives. In Table 11, we provide treatment effects on such alternative utility metrics.

In the first row of Table 11 we present the average treatment effect on mean utility per user, it is insignificant. This metric weights each user equally, but does not capture the increase in the number of user-story interactions. Finally, in the second row, we present the treatment’s impact on mean utility per user-story interaction, this metric puts more weight on heavy users, as they have more interactions. We find that there is a strongly positive and statistically significant treatment effect. This suggests that the personalized policy had a stronger positive effect on heavy users, that consume many stories, than on somewhat infrequent users.

D Heterogeneous treatment effects: regressions analysis

To provide further robustness into the finding that heavy and niche users benefit from personalized recommendations, we present results of regressions of AIPW scores based on total utility on users’ past utilization (see Athey and Wager (2019) for methodology). See table 12 for summary of results.

Columns one to four of table 12 show that users with high past utilization have higher treatment effects. Column five shows higher treatment effect for niche type users. Finally, columns six and seven control both for heavy utilization and niche type; niche type remains to have a high and statistically significant treatment effect.

Table 12: Results of a regression of user types on AIPW scores.

	Dependent variable:						
	Total utility (aipw.scores)						
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
past utility	0.001*** (0.0003)						0.001 (0.0004)
stories completed		0.002*** (0.001)				0.001** (0.001)	
heavy user utility			0.366*** (0.075)				
heavy user completions				0.352*** (0.076)			
niche type					0.342*** (0.074)	0.231*** (0.088)	0.256*** (0.091)
Observations	2,661	2,661	2,661	2,661	2,661	2,661	2,661
R ²	0.006	0.007	0.009	0.008	0.008	0.010	0.009

Note:

*p<0.1; **p<0.05; ***p<0.01

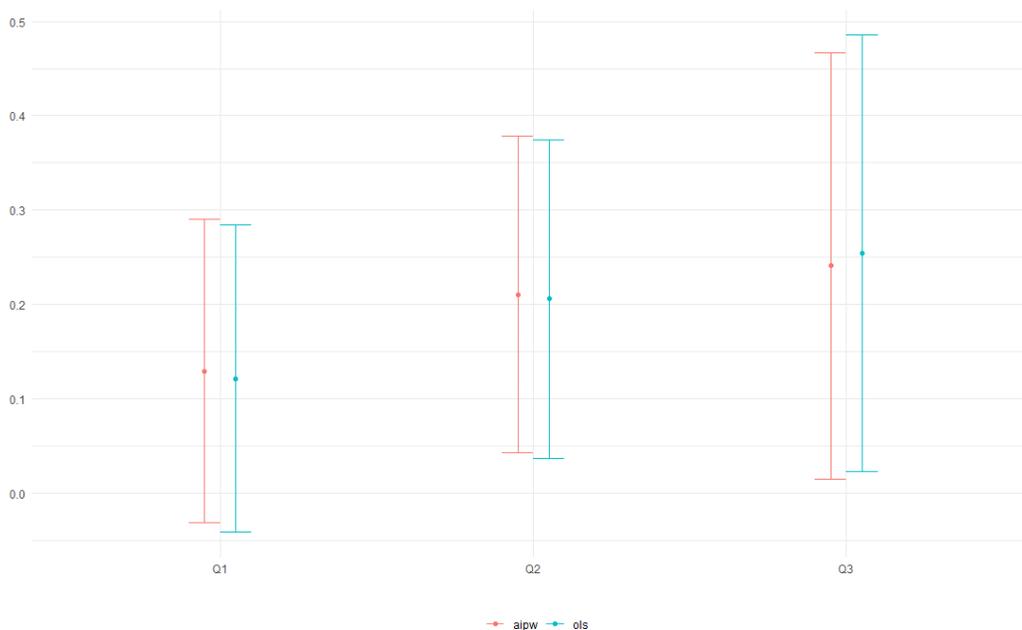
Note: Outcome variable is AIPW score of total utility per user. OLS estimator. All covariates defined based on the pre-experiment app usage. *p<0.1; **p<0.05; ***p<0.01

E Data-driven treatment effects heterogeneity

We use the estimated causal forest to divide our users into tertiles according to their estimates CATE prediction (see Chernozhukov et al. (2018) for details of this approach). To avoid using model that was fitted using observations for which we make predictions, we use honest sample splitting with 10 folds.

Figure 8 shows the the predicted CATES in the four groups. First of all, the treatment effects are quite similar for the four groups. The fourth quartile appears to have higher treatment effects, but the differences are small.

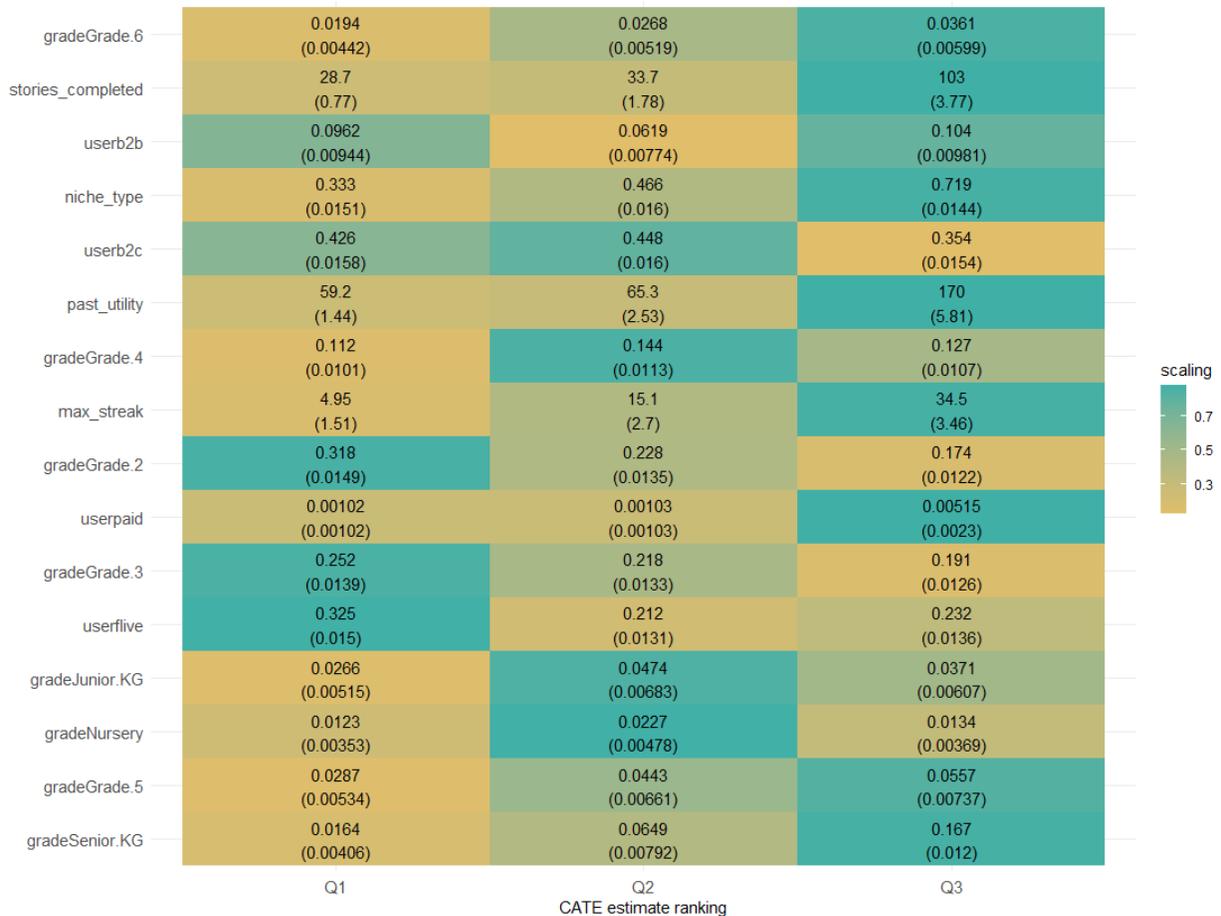
Figure 8: Average CATE within each ranking (as defined by predicted CATE). Predictions with OLS in blue and AIPW scores in red.



Finally, we can also compare average characteristics for individuals in the four quartiles. We present such a comparison in Figure 9.

Heavy users (high maximal streak and freq-user indicator) appear more frequently in the highest quartile. We also see more niche users in the fourth group. We look in detail into these groups in the next subsections.

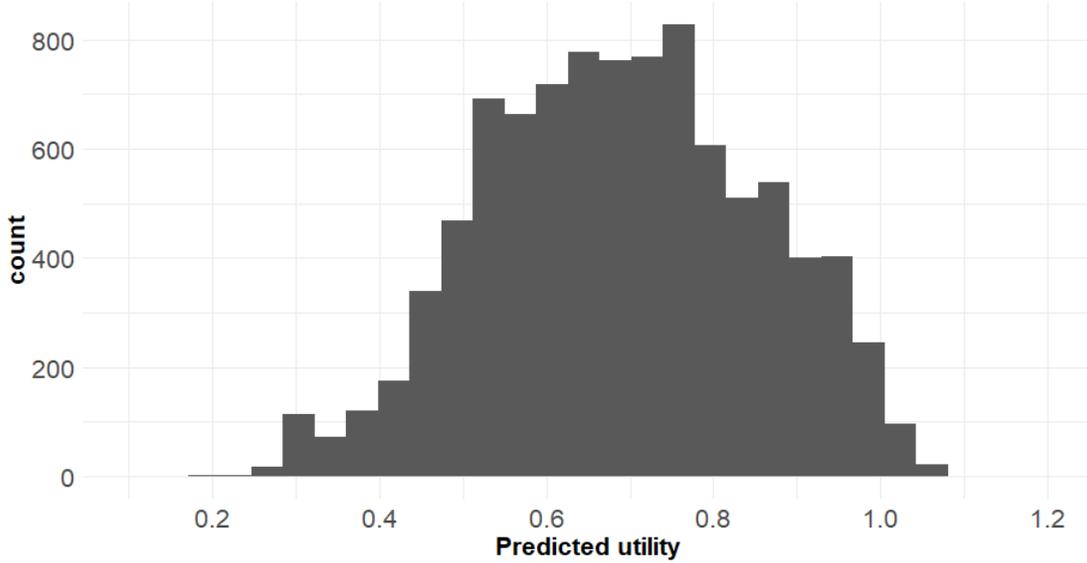
Figure 9: Average covariate values within group (based on CATE estimate ranking).



F Model calibration with experimental data

The main component of the recommendation system is the collaborative filtering model that predicts user utility from user-story interactions. In our analysis, high treatment effects suggest that the model has successfully identified user preferences and selected stories that users liked.’ In this section, we further evaluate the calibration of the model by correlating the models predicted user utilities with observed utilities from the experiment. Figure 10 shows the histogram of the predicted utilities; we can notice that they vary from very high values of around 1 to lower values of 0.2. We don’t see values

Figure 10: Histogram of predicted utility



lower than 0.2 because in the experiment we considered only stories ranked at the top of the ranking of predicted utilities.

To evaluate the model calibration, we regress the predicted utility on the observed utility from the experiment. Regression results are in Table 13.

Table 13: Correlation between utility predicted by the collaborative filtering model and observed in the experiment. *** : $p < 0.01$

	<i>Dependent variable: utility</i>		
	All users	Frequent users	Infrequent users
pred. utility	0.44*** (0.01)	0.46*** (0.01)	0.43*** (0.01)
Observations	9,344	2,268	7,076
R ²	0.40	0.38	0.41

Table 13 shows that the model predictions are strongly correlated with observed utilities. The model is better calibrated for frequent users, for whom we have longer consumption histories (albeit the difference is small). We also break down the analysis by the experimental groups, which is shown in Table 14.

In Table 14 we can notice that the predictions from the model correlate strongly with observed utility in both treatment and control groups. The model is much better calibrated in the control group, this is not surprising because the model is trained on similar data

Table 14: Results from linear regressions of actual utility on predicted utility. Column (1) treatment group, column (2) control group.

	<i>Dependent variable: utility</i>	
	treatment	control
pred. utility	0.33*** (0.01)	0.60*** (0.01)
Observations	4,695	1,773
R ²	0.32	0.51