



Externalities in knowledge production: evidence from a randomized field experiment

Marit Hinnosaar^{1,4,5} · Toomas Hinnosaar^{1,5}  · Michael E. Kummer^{2,6,7} · Olga Slivko³

Received: 6 May 2019 / Revised: 18 August 2021 / Accepted: 18 August 2021 /
Published online: 1 September 2021
© Economic Science Association 2021

Abstract

Are there positive or negative externalities in knowledge production? We analyze whether current contributions to knowledge production increase or decrease the future growth of knowledge. To assess this, we use a randomized field experiment that added content to some pages in Wikipedia while leaving similar pages unchanged. We compare subsequent content growth over the next 4 years between the treatment and control groups. Our estimates allow us to rule out effects on 4-year growth of content length larger than twelve percent. We can also rule out effects on 4-year growth of content quality larger than four points, which is less than one-fifth of the size of the treatment itself. The treatment increased editing activity in the first 2 years, but most of these edits only modified the text added by the treatment. Our results have implications for information seeding and incentivizing contributions. They imply that additional content may inspire future contributions in the short- and medium-term but do not generate large externalities in the long term.

Keywords Knowledge accumulation · User-generated content · Wikipedia · Public goods provision · Field experiment

JEL Classifications L17 · L86 · C93

Data and code for this article are available in Hinnosaar et al. (2021a), Harvard Dataverse, <https://doi.org/10.7910/DVN/T4VFCX>.

✉ Toomas Hinnosaar
toomas@hinnosaar.net

Extended author information available on the last page of the article

1 Introduction

Knowledge is a key input to many economic activities and a driver of economic growth (Grossman & Helpman, 1993; Jones, 1995; Romer, 1990). An increasing share of knowledge is created in the form of user-generated content: consumer feedback systems, discussion boards, Q&A sites, open-source software, social networks, and online information repositories, such as Wikipedia. Online knowledge repositories have the potential to revolutionize how society aggregates and transmits knowledge. This potential stems from their ability to combine and aggregate the input of many individuals independent of time and location, and the generated content can be retrieved at a low cost.¹

A key component to the success of user-generated content repositories is a sufficient flow of content contributions by users. Hence, understanding the drivers of contributions to user-generated content has been an important question in economics and management for the past two decades (Lerner & Tirole, 2003). While traditional motives of public goods provision play a vital role in user-generated content, recent literature has identified a novel driver in the form of a feedback loop—a dynamic by which small initial contributions of content inspire ever more follow-on content contributions by other users (Aaltonen & Seiler, 2016; Kane & Ransbotham, 2016; Zhu et al., 2020).

In this paper, we investigate whether there are positive or negative externalities in user-generated content production. Understanding and quantifying such externalities has important policy implications. If content generation has positive externalities on future content generation, then information seeding and paid contributions may have a high return on investment in terms of added stimulated growth (Aaltonen & Seiler, 2016).² On the other hand, if content generation has negative externalities, then such policies may backfire and be ineffective or even lead to worse eventual outcomes (Nagaraj, 2019).

Due to the reflection problem (Manski, 1993), externalities in content generation are difficult to identify. An externality occurs when a contribution by a user motivates other users to contribute (positive externality) or prevents further contributions (negative externality). Yet, correlation in contributions does not necessarily attest to an externality. A positive correlation may arise when two users are exposed to the same external shock, such as a news article or a research finding. Similarly, a non-causal negative correlation may be caused by processes with periodic updates, such

¹ Traditional channels of personal knowledge transmission require a double-coincidence demand and supply of knowledge. The “knowledge-seeker” and the “knowledge-holder” have to meet in person or at least at the same time. The elimination of such double-coincidences has been modeled to understand the advantage of monetary over barter-economies. Kiyotaki and Wright (1989). These features give such systems a drastic competitive advantage that may affect the education sector and other traditional channels of knowledge transmission. The sector of encyclopedic knowledge is one of the most salient examples of the new technology’s potential.

² Nagaraj (2019) describes how such policies have been used by Wikipedia (seeding articles on more than 30,000 US cities from US Census Bureau data), OpenStreetMap (US Census maps), and Reddit (fake user accounts).

as elections or periodically updated statistics. To identify the causal effect, shocks to content growth and contributions must be independent over time. We estimate the causal impact of additional content on subsequent contributions using a randomized field experiment in Wikipedia. Randomization ensures that the addition of content is exogenous in terms of future content generation. As argued by the literature analyzing social interactions [including Manski (1993) in general and more specifically Aaltonen and Seiler (2016) and Zhu et al. (2020)], randomized experiments are the best way to cleanly identify causal relationships in such interactions.

The exogenous variation in our data comes from a randomized field experiment, which was conducted in 2014. The experiment added relevant content to randomly chosen Wikipedia pages while leaving similar pages unchanged. The treatment added about two paragraphs (approximately 2000 characters) and one picture to each page in the treatment group. The pages were about mid-sized Spanish cities in different language editions of Wikipedia.³

For our analysis, we collect data from multiple sources. To measure the impact on the quantity of content and editing activity, we use a dataset of Wikipedia editing histories, which includes all versions of Wikipedia pages in the treatment and the control groups. To measure the effect on the quality of content, we use two approaches. First, we developed a quality rating scheme, and for each page, we obtained quality ratings by two independent raters who are fluent in the corresponding language. Second, we use text analysis to compare the content across different language editions of Wikipedia and measure similarity to the corresponding pages in the Spanish Wikipedia. We use the Spanish Wikipedia as a benchmark because it provides the most detailed coverage of Spanish cities among all language editions of Wikipedia. Our dataset and the experimental setting allow us to analyze both short-term and long-term effects, up to 4 years after the experiment. Our main outcomes of interest are the quantity and quality of content. To study editing activity, we also analyze the number of unique editors, the number of edits, and the amount of content added and deleted.

We find that the pages that were improved by adding about 2000 characters were still only longer by about the same amount, even 4 years after treatment. Similarly, while the added content increased the quality, the difference 4 years later was about the same as immediately after the treatment. Our estimates allow us to rule out effects on 4-year growth of content length larger than twelve percent. We can also rule out the effects on 4-year growth of content quality larger than four points, which is less than one-fifth of the size of the treatment itself. We do find some short-term impact. In particular, we find that the editing activity increased in the first 2 years after the treatment. In this initial stage, the treatment increased the number of Wikipedia users editing the treated pages and increased the number of edits. However, there was no such impact in the third and fourth years post-experiment. Moreover, even in the first years, most of the additional edits only modified the text added by the treatment.

³ A comprehensive description of the experiment is provided in Hinnosaar et al. (2021b), who studied the impact of this treatment on real-world outcomes.

While our study benefits from clean identification, it faces two important limitations. First, our test has low power to measure small effects on content growth. Second, our results might not generalize to other platforms or content in other stages of development. In our setting, the content is relatively mature but still incomplete. It is plausible that the results would be different if we had studied either new pages or almost complete pages.

Our main finding has a clear policy implication—at least in settings similar to the one studied here, investments in information seeding and incentivizing additional content contributions may temporarily increase user participation, but do not have a large cumulative effect on content growth. Therefore information seeding and incentivizing contributions are mainly a matter of direct cost-benefit analysis: they pay off if and only if the costs of creating the content are lower than the value of the new content.

Our paper contributes to the literature that studies externalities in user-generated content production. The closest to our work are Aaltonen and Seiler (2016), Kane and Ransbotham (2016), Nagaraj (2019), and Zhu et al. (2020). Aaltonen and Seiler (2016) and Kane and Ransbotham (2016) used detailed observational data from Wikipedia. Nagaraj (2019) used a natural experiment on OpenStreetMap, a Wikipedia-style digital map-making community, which started with better seeding information in some regions compared to others for quasi-random reasons. Zhu et al. (2020) studied a natural experiment that motivated students to contribute to Wikipedia. The papers arrived at contradicting conclusions, which warrant further investigation regarding this issue. Our paper is the first to study the question using a randomized field experiment, which allows a clean identification of the underlying externalities, especially when analyzing the long-term impact.

More generally, the paper belongs to the literature that analyzes what drives contributions of user-generated content, which represents a salient and highly relevant digital public good.⁴ Among the main drivers of content contributions, the studies addressed the role of personal gain (Shah, 2006), social comparison (Chen et al., 2010), group size (Zhang & Zhu, 2011), networks (Fershtman & Gandal, 2011; Ransbotham et al., 2012), attention spillovers (Kummer, 2020), symbolic awards (Gallus, 2017), performance feedback (Huang et al., 2018), monetary rewards versus social motives (Sun et al., 2017), contributor diversity (Ren et al., 2015), and economic conditions, such as unemployment (Kummer et al., 2019) and migration (Slivko, 2018). Social motives have been shown to affect public good provision (Ayres et al., 2013; Goldstein et al., 2008; Lacetera & Macis, 2010). Specifically, in the case of digital public goods, examples include ranking movies on MovieLens (Chen et al., 2010), editing articles on Wikipedia (Algan et al., 2013; Zhang & Zhu, 2011), and endorsing messages of Facebook users (Egebark & Ekström, 2017).⁵ In

⁴ Other studies on Wikipedia have analyzed biases in Wikipedia's content (Greenstein & Zhu, 2012, 2018; Hinno Saar, 2019) and the impact of Wikipedia on market outcomes (Hinno Saar et al., 2021b; Xu & Zhang, 2013) and science (Thompson & Hanley, 2018).

⁵ More generally, the literature suggests strong effects of social influence on individual choices related to savings (Duflo & Saez, 2002, 2003), education (De Giorgi et al., 2010; Hanushek et al., 2003), entertainment (Salganik et al., 2006), etc.

our setting, social externalities are rather implicit, as individuals contributing content are not connected by any direct social ties, and the interactions with the other members of the community can occur only in the process of contributing knowledge to the Wikipedia articles. Our paper extends the literature by using variation from a randomized field experiment to measure the impact of additional content on future content generation.

The structure of the paper is as follows. In the next section, we describe the experiment and provide some background on Wikipedia editing. Section 3 describes the data. Section 4 presents the results of the impact of the treatment on the subsequent growth in content quantity and quality and on editing activity. Section 5 interprets our results in a simple theoretical framework. Section 6 discusses the connection between our findings and related literature. It also discusses the implications and generalizability of our findings. Section 7 concludes.

2 Experiment and background

2.1 Experiment

The field experiment added content (text and photos) to randomly chosen Wikipedia pages. The sample consisted of 240 Wikipedia pages. Specifically, it consisted of the pages of 60 Spanish cities in the French, German, Italian, and Dutch editions of Wikipedia. The cities were all medium-sized, excluding the largest like Madrid and Barcelona, and also excluding the smaller cities. The Wikipedia pages in these languages were relatively short—up to 24,000 characters in each of these four languages.

Each city and each language edition of Wikipedia was treated equally. For each city, its page was assigned to the treatment group in two randomly chosen languages. In each language edition of Wikipedia, 30 randomly chosen city pages were assigned to the treatment group. Specifically, to obtain balance in the treatment and control groups, the randomization was stratified.⁶ The 60 cities were divided into ten equal-sized groups. Within each group, each city was randomly assigned to one of six treatment arms. The six treatment arms were as follows: treat the city page in one of the six possible language pairs (French & German; French & Italian; French & Dutch; German & Italian; German & Dutch; Italian & Dutch). This resulted in a design where the number of pages which were treated equaled the number of those that remained in the control group.

The Wikipedia pages were treated mid-August, 2014. The treatment added about 2000 characters of text and photos to each page in the treatment group. The added text and photos were mostly obtained from the corresponding Spanish and English language Wikipedia pages. Because all the pages were about Spanish cities, the Spanish Wikipedia typically contained more information than the other language versions. The English language version of the page, typically, was also more

⁶ For further details of the randomization, see Hinno Saar et al. (2021b).

detailed than in the languages in the experiment. Hence, there was information available in Spanish and English pages that was missing from the other language editions of Wikipedia. The treatment translated that text and added it to the corresponding pages in the treatment group.

The treatment of the pages in Dutch Wikipedia was not successful. While in French, German, and Italian Wikipedia, the added text and photos survived well over time, all the additions to Dutch Wikipedia were deleted within 24 h (by a single editor). Wikipedia allows anyone to edit. It also means that anyone can delete or undo the latest changes by reverting to a previous version of the page. This happened in the Dutch version of Wikipedia, where 24 h after the treatment, all the pages looked as if they had never been treated. Therefore, we exclude Dutch pages from our main analysis and restrict attention to the 180 pages in French, German, and Italian. Robustness analysis shows that our results do not change if Dutch pages are included in the analysis.

2.2 Power analysis

We acknowledge that the sample size and power of our tests are rather small. We analyze this in online appendix A, which presents power analysis for one of our main outcome variables (page length) and the main editing activity variable (number of users). As estimation results in Sect. 4 show, we would expect the power for other measures to be similar.

To provide some context, let us discuss the expected effect sizes, given the results from previous literature. The main measure studied in most previous works is the impact of added content on the number of future contributors. Aaltonen and Seiler (2016) estimate that adding 10,000 characters of content leads to 0.204 additional users per week, which corresponds to about 0.18 additional users per month for 2000 characters as added by the treatment in our case. Zhu et al. (2020) estimates are similar: their median treatment size was 3180 characters and estimated increase of 2 unique users per 6-month period, which implies 0.21 new users per month per 2000 added characters.⁷

The literature provides less guidance for the long-term effect size for the quantity and quality of the content. Aaltonen and Seiler (2016) use their estimates for the impact on the number of users to simulate a possible effect on the content quantity and find 45% growth in content. The only other paper to estimate this effect is Nagaraj (2019), who finds a long-term effect of negative 10%.

Figures A.1a and A.1b in online appendix A describe the power analysis for page length. Figure A.1a shows that when the Dutch pages are included in the sample, as originally intended, then if the true treatment effect is 10% increase in page length over 4 years, we would reject the null hypothesis of no effect at 10%-significance

⁷ Kane and Ransbotham (2016) provide some evidence that the effect could be larger for less developed content. They find that in the case of less developed articles, 1%-increase in length implies 0.03–0.04 more monthly contributors. In our case, the treatment was on average 23% of the page length, which would imply 0.7–0.9 more users.

level with 76% probability and at 5%-significance level with 65% probability. The minimum detectable effect size is about 12%.⁸ If we exclude the Dutch pages (Figure A.1b), we lose some power, but the minimum detectable effect is still around 13%. Indeed, our study is underpowered to detect small long-term effects, but there should be no difficulties detecting even half of the effect-size suggested by Aaltonen and Seiler (2016). Figures A.1c and A.1d show that our experiment has relatively more power to detect meaningful effects on editing activity. Even if we exclude the Dutch pages (Figure A.1d), the minimum detectable effect size is 0.11 users, and we should certainly be able to detect the effect sizes suggested in the literature. Section 4 describes the ex-post minimum detectable effect sizes (with our realized data), which turn out to be similar to those described here (calculated based on the pre-experiment data).

2.3 Background on Wikipedia and its editing

Wikipedia exists in 309 languages, and the different language editions are not identical. The differences across languages made the experimental design possible. As we show in Sect. 3.5, Spanish Wikipedia contained much more information about Spanish cities than the pages in our sample. This imbalance allowed the treatment to translate information from Spanish Wikipedia to the target languages.

Why don't Wikipedia editors translate the content between languages themselves? First, it requires language skills. Many people are monolingual (Eurobarometer, 2012). English language skills would not be enough because the pages in English Wikipedia are also rather incomplete (as we show in Sect. 3.5). The language skills required to translate between Spanish, French, German, and Italian are not common. Table B.1 in online appendix shows that in France, Germany, and Italy, less than 10% of the population can read Spanish, and in Spain, less than 10% of the population can read the languages in our sample. Second, perhaps equally importantly, Wikipedia is written by volunteers whose motivation depends largely on how fun the editing process is (Nov, 2007). Presumably, the task of translating information is somewhat mundane compared to other possible uses of time. Third, note that for the majority of the 309 language editions of Wikipedia, automatic translation is still not good enough. Furthermore, even when automatic translation is technically possible, to use it, it would require that Wikipedians agree on which language version is the superior one.

Let us briefly describe how Wikipedia editing works on these pages. While large edits that add an entire section or paragraph definitely exist, they do not make up the majority of Wikipedia edits. Most Wikipedia edits are small, fixing typos, grammar, and formatting, rearranging text without modifying content. Figure B.1a shows that 75% of the edits to pages in our sample (before treatment) add less than 100 characters (about one sentence). 42% of edits are marked by the editor as minor edits,

⁸ The minimum detectable effect size is calculated at 5%-significance level and 80% power.

which means not modifying content. Only about 4% of edits add more than 1000 characters (about one paragraph).

Where do the small edits come from? We hypothesize that some of these small edits occur when editors keep track of pages that they are interested in. On Wikipedia, any registered user can sign up to be notified when a certain page has been modified. We find some evidence supporting this on the corresponding city pages in English and Spanish Wikipedia. Figure B.1b shows that almost half of these pages have more than 30 editors that get notified of all changes. Unfortunately, if there are less than 30 watchers for a page, Wikipedia does not report the exact number. Hence, we just know that only 8 pages in our sample have at least 30 persons signed up for notifications. Nevertheless, the comparison with the pages in English and Spanish Wikipedia suggests that for the remaining pages, there should be some watchers as well.

3 Data

We combine multiple sources of data. To measure the impact on the quantity of content and editing activity, we use a dataset of Wikipedia editing histories. An editing history contains the full text of each revision⁹ of each page starting from the creation of the page until the beginning of September 2018.

To measure the effect on the quality of content, we use two approaches. First, we developed a quality rating scheme, and for each page, we obtained quality ratings by two independent raters who are fluent in the corresponding language. They rated three revisions of each page in our sample (before treatment, after treatment, and 4 years after treatment) comparing it to the English Wikipedia as a benchmark. Second, we use text analysis to compare the content across different language versions of Wikipedia and measure similarity to the corresponding pages in the Spanish Wikipedia.

Our sample consists of the 180 pages in the experiment, which are the pages of 60 cities in French, German, and Italian Wikipedia. In the following subsections, we describe the construction of the dataset and variables used in the analysis.

3.1 Page length

One of our main outcome variables is the page length after the experiment. We measure page length in characters, including spaces and wiki markup commands.

Figure 1 presents average page length in the treatment and control groups. Until the experiment in August 2014, the average page length in the control and treatment

⁹ A revision (or an edit) is a version of a Wikipedia article saved at a specific moment by a particular user. All revisions with the corresponding metadata, including full text, user, and timestamp, are preserved by Wikipedia and publicly available.

groups was rather similar.¹⁰ The experiment added significant length to the pages in the treatment group. After the experiment, the difference has been relatively stable. An exception is a sharp increase in the mean of the treatment group in August 2016. This jump comes from the efforts of a single editor who worked hard to improve one page in French Wikipedia—the page of the city of Cordoba.¹¹ Online appendix B presents the same figure, first without French Cordoba (Figure B.2a) and second, with the logarithm of page length (Figure B.2b), both of which show no evidence of an increase in the treatment group average in 2016.

Similar dynamics can be seen when looking at the changes separately by language (Fig. 2a–c). As expected, the placebo test with the Dutch pages (Fig. 2d) shows that the assignment to the treatment group had no impact. Page length is one possible output measure of knowledge production in Wikipedia. Similar dynamics as in Fig. 1 can also be seen in Figure B.3 in online appendix B, which presents alternative measures of content: images and plain text (that is, html elements removed from the parsed text).

To make the treatment and the control groups comparable, we subtracted the length of text added by the treatment from the length of pages in the treatment group. Moreover, as the distribution of page lengths is relatively skewed, we use the logarithm of page length in our estimations.

3.2 Quality ratings

To assess the changes in the quality of Wikipedia articles, we hired six research assistants to rate the quality of articles. Each article in French, German, and Italian Wikipedia in our sample was evaluated by two raters, who were fluent in the respective language as well in English.

We asked the raters to evaluate the quality of three versions of 60 Wikipedia articles in our sample. Version A was the latest version before August 1st, 2014 (i.e., pre-treatment), version B the latest version before September 1st, 2014 (i.e., post-treatment), and version C the latest version before September 1st, 2018 (i.e., 4 years after the treatment). As a benchmark, we used articles in English-language Wikipedia (as of September 1st, 2018).¹²

Each version of each article was rated in five dimensions on a scale where 0 is the lowest possible rating and 100 means equivalence to the benchmark page from English Wikipedia (for detailed instructions, see Figure B.4 in online

¹⁰ The drop in both the treatment and control groups in early 2013 comes from technical changes in Wikipedia: Addbot removed about 2000 characters from each page with an explanation similar to “Migrating 77 interwiki links, now provided by Wikidata”.

¹¹ By August 2016, the page of Cordoba in French Wikipedia was relatively long, with 19,426 characters (at the time 93% of the pages in our sample were shorter than that). During August 2016, this user increased the page length to 100,702 characters, which is almost twice the length of the longest page at the time (57,076 characters). Our conclusions do not change if we exclude this page.

¹² As we show below, the articles about Spanish cities in English-language Wikipedia are sometimes quite incomplete, so ideally we would have preferred to use Spanish Wikipedia as a comparison. Because the combination of necessary language skills is not common, it would have been prohibitively costly.

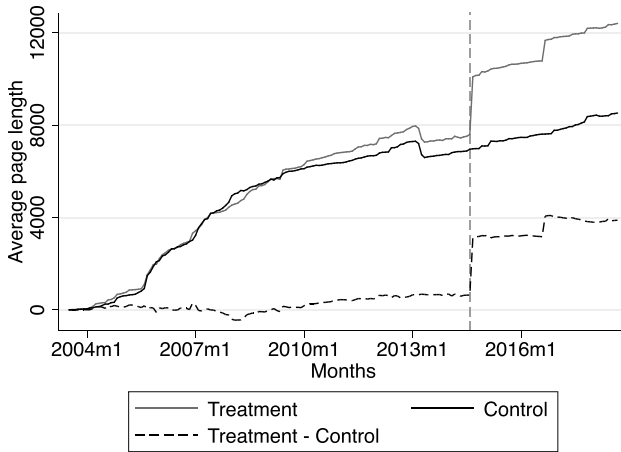


Fig. 1 Average page length in the treatment and control groups. *Notes:* The number of observations is 90 in the control and 90 in the treatment groups. The experiment month (August 2014) is marked by the dashed vertical line

appendix B). The five dimensions were the following. (1) *Completeness*: the article comprehensively covers all relevant aspects of the city (compared to the article in English). (2) *Well-written*: the prose is clear, concise, and spelling and grammar are correct. (3) *Illustrated*: the article includes photos that are relevant to the topic and have suitable captions (compared to the article in English). (4) *Interesting*: the article makes the city seem like an exciting place to visit (compared to the article in English). (5) *Overall*: Overall, the article is a high-quality reference source (compared to the article in English).

For example, a score of 100 in the completeness dimension means that the article covers the relevant aspects of the city as comprehensively as the corresponding page in English. It may occur if the pages cover the same topics in the same level of detail, or if they cover different topics, but the missing parts “balance out” in the eyes of the raters. A score of 50 would mean that the page is half as good and 200 that the page is twice as good as the benchmark page from English Wikipedia.

Figure 3 presents the change in quality during the treatment (August 2014) and within 4 years after treatment (September 2014–September 2018). It shows that in the treatment group, as expected, there was a large increase in quality during the treatment. The increase takes place in overall quality and in all measured dimensions except well-written. Indeed, during the treatment month, pages in the treatment group become slightly less well-written. This is expected because, in the experiment, the treatment text was written outside of Wikipedia and then copied to Wikipedia. During the same time in August 2014, the changes in the control group within August 2014 are negligible. Within the following 4 years, overall quality increases in both treatment and control groups, and interestingly these increases are much smaller than the treatment itself. However, the 4-year increase in the treatment group is of a similar size as that in the control group. To further help to understand

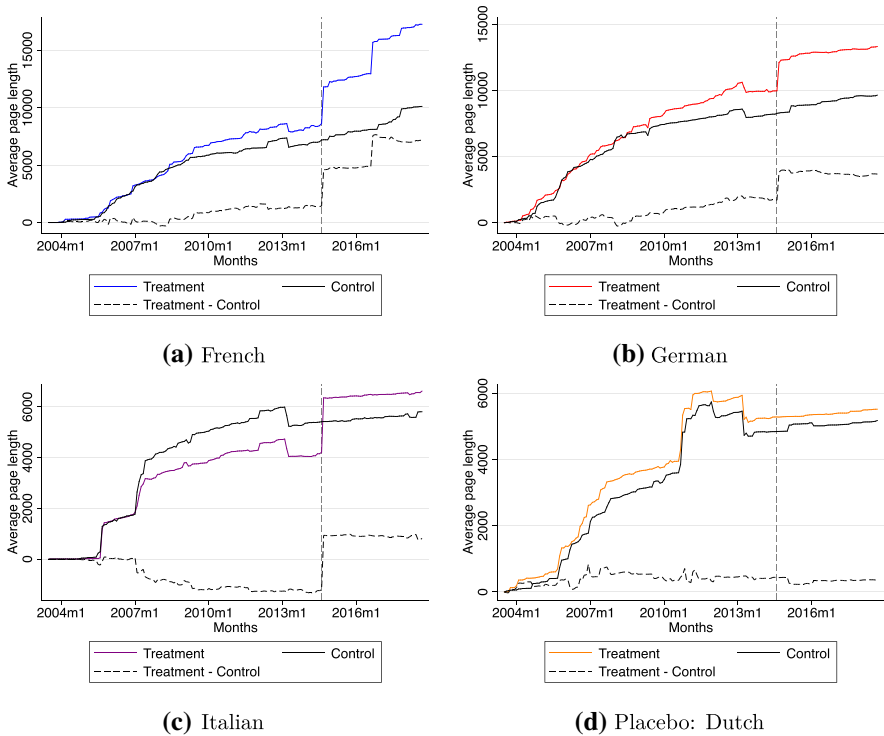


Fig. 2 Average page length in the treatment and control groups, by language. *Notes:* On each figure, the number of observations is 30 in the control and 30 in the treatment groups. The experiment month (August 2014) is marked by the dashed vertical line

the context of the quality measures, Figures B.5a–B.5b in online appendix show that the quality measures are correlated with page length.¹³

3.3 Measures of editing activity

To better understand the process of content creation, we also study editing activity. To construct the measures of editing activity, we start with 30,601 edits (revisions) from 180 Wikipedia pages. This set includes all the edits except those generated as part of the treatment in the experiment. Following Aaltonen and Seiler (2016), we restrict the sample of edits in the following ways. First, we exclude edits by bots (about 30% of edits); these are non-human user accounts that generate automated edits. Specifically, we define bots as users whose username occurs in the list of bots (in the English, French, German, or Dutch Wikipedias) or whose username includes

¹³ Similar correlations between quantity and quality of content have been found previously, for example by Chen et al. (2019).

“bot”. Second, we exclude reverts, which are edits that restore any previous version of the same page (about 7% of remaining edits). Third, we exclude vandalism (about 0.8% of remaining edits). We use the following criteria to classify an edit as vandalism: (a) an edit that only deletes text from the previous revision, and (b) the revision immediately after vandalism reverts the article back to a past revision. Then we are left with 19,586 productive edits generated by human users.

To analyze the impact of treatment on editing activity, we construct three types of monthly measures that characterize how many people edited the pages, how many times they edited, and how much they edited. The first measure is the number of unique users editing a page per month. We define a unique user by the username for registered users and by IP address for anonymous users. The second measure is the number of edits per month. To avoid double-counting of micro-edits,¹⁴ we first aggregate edits to the day-user-page level and then sum these up to month-page level. The third measure is edit distance—the number of characters an edit added plus the number of characters it deleted compared to the previous version of the page.¹⁵ We aggregate the edit distance measure to monthly level. Figure B.6 in online appendix describes the average editing activity in the treatment and control groups over time.

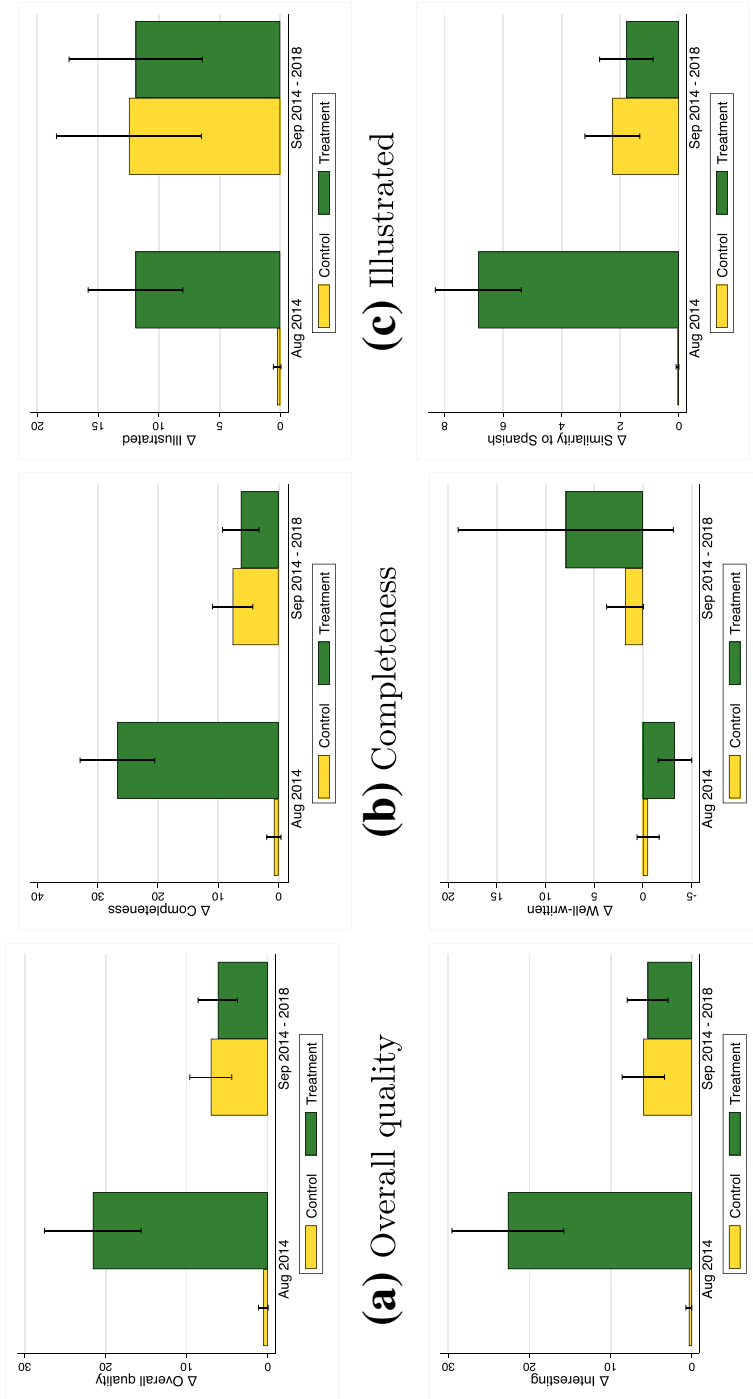
In addition to the aggregate measures of editing activity, we separate edits that directly modify the treatment text and those that modify other parts of the page. We classify edits into these two categories using a method similar to Hinno Saar et al. (2021b). For each page in the treatment group, we use the diff algorithm between the revision before and after the treatment to determine treatment text—the exact text added by the treatment. For each revision post-treatment, using the diff algorithm between the treatment text and this revision, we check whether the revision deletes any part of the treatment text. If the revision doesn't delete anything from the treatment text, we classify the revision as one that edited other parts of the page.

3.4 Similarity to Spanish Wikipedia

To complement the human-rated page quality data, specifically page completeness, we also evaluate the changes in quality by computing the similarity with Spanish Wikipedia. We use the articles in the Spanish language (as of September 1st, 2018) as a benchmark as they provide the most detailed coverage of Spanish cities among all language versions of Wikipedia. Therefore, if the coverage of topics in another language becomes more similar to the corresponding article in the Spanish language, this can be interpreted as increasing completeness of the article.

¹⁴ Many Wikipedia editors save many revisions to the same page in a short period of time, for example, generating a new revision after each sentence they write. This is partly motivated by the fact that someone else might edit the page at the same time.

¹⁵ Edit distance is widely used in computational linguistics and computer science to measure the similarity of strings. It is a generic term that allows any weights of insert, delete, and substitution operations. Common variants put weight 1 to addition and deletions, and weight substitutions either by 1 (called Levenshtein distance) or by 2 (the measure we use). For each edit, we calculate the edit distance using PHP FineDiff class at the granularity level of a character.



(a) Overall quality **(b)** Completeness **(c)** Illustrated **(d)** Interesting **(e)** Well-written **(f)** Similarity to Spanish

Fig. 3 Change in quality during the treatment month (August 2014) and within 4 years after treatment (from September 2014 until September 2018), separately for the control and treatment groups. *Notes:* The number of observations is 90 in the control and 90 in the treatment groups

As a measure of similarity, we use the Tversky index (Tversky, 1977). Formally, if the set of terms mentioned in the article of the target language is A and the set of terms mentioned in the Spanish article is B , then the Tversky index is computed as

$$S(A, B) = \frac{|A \cap B|}{|A \cap B| + \alpha|A \setminus B| + \beta|B \setminus A|}, \quad (1)$$

where $\alpha \geq 0$ and $\beta \geq 0$ are parameters. If $\alpha = \beta = 1$, the index becomes equivalent to the Jaccard similarity index, and if $\alpha = \beta = 0.5$, it simplifies to the Sørensen-Dice similarity index.¹⁶ More generally, a larger α puts a bigger weight on A (interpreted as a variant) and β a bigger weight on B (interpreted as a prototype). As our aim is to use the Tversky index as a completeness measure (how complete is A in comparison to B), we set $\alpha = 0$ and $\beta = 1$, which simplifies the similarity measure to $S(A, B) = \frac{|A \cap B|}{|B|}$.

Before comparing the articles in the treatment languages with their Spanish counterparts, we need to extract the set of terms mentioned in each article. For this, we first translate all articles into English using Yandex Translate API.¹⁷ Then, we remove all so-called stop words, i.e., words that do not reflect the specifics of the content.¹⁸ We then use the Porter stemming algorithm (Porter, 1980) to remove the endings from words in English so that only root words remain. Finally, we drop all remaining strings that are shorter than three characters or contain non-alphabetic characters. This process gives us a set of terms for each article that we use in computing the Tversky index.¹⁹

Figure 3f presents the change in similarity during the treatment and within 4 years after treatment. As with quality, we see a large increase in the treatment group in similarity during treatment and much smaller increases within the following 4 years in both treatment and control groups. Figures B.5c–B.5d in online appendix show that similarity to Spanish is correlated with page length and also with completeness (compared to English).

3.5 Summary statistics

Comparison of pages in treatment and control groups Table 1 presents the comparison of pre-treatment page length, quality, and editing activity in the treatment group versus the control group. The table shows that there were no significant differences between the two groups before the treatment. Figures B.8–B.9 in online

¹⁶ Both indexes are used to measure the similarity between two documents. The Jaccard index is also known as the Intersection over Union and sometimes called the Tanimoto similarity. Another similarity measure used in earlier works in economics (Chen et al., 2019; Thompson & Hanley, 2018) is cosine similarity. It would be preferable if we want to capture the similarity of two pages not only in terms of content covered but the language and tone. As we aim to capture completeness of the page compared to a benchmark, we found the Tversky index most appropriate for the task.

¹⁷ For more information, see <https://tech.yandex.com/translate/>.

¹⁸ For example, pronouns (“it”, “their”) and prepositions (“on”, “before”). The full list is in B.7.

¹⁹ Terms such as “america”, “archipelago”, “area”, “arona”; about 250–500 terms for each article.

appendix present the kernel density estimates of full distributions separately for the treatment and control groups for all the variables in Table 1. We conclude that while there are differences between the groups, by and large, the randomization was rather successful.

Are the pages in our sample already rather complete? To answer the question, ideally, we would like to use Wikipedia's own quality ratings. Unfortunately, the pages in our sample have not been rated. Each language edition of Wikipedia has its own quality rating system, but only in the English Wikipedia is the system widely used. Therefore, we assess the completeness of the pages in two steps. First, we determine the quality of the corresponding city pages in the English Wikipedia. Then we compare the pages in our sample to the English pages.

Figure B.10 in online appendix presents the distribution of the quality ratings in English Wikipedia of the 60 city pages that correspond to the pages in our sample. These city pages have been assigned the lowest possible grades (Stub, Start, and class C) or have not been rated at all (9 cities).²⁰ The highest rated pages, class C, according to the Wikipedia rating scale, still miss important content. Hence, we conclude that the 60 city pages are all low-quality articles in the English Wikipedia.²¹

Having seen that all the English pages still have room for improvement, how do the pages in the sample compare to the English pages? Table 2 compares our pages to the pages in English in 3 dimensions: relative length, (calculated) similarity, and (human rated) completeness. Panel A column 1 shows that after treatment the median page in the treatment group is about 75% of the relative length, about 33% in terms of similarity, and almost 100% in terms of completeness compared to the corresponding page in English. Recall here that we saw above that the pages in English were far from complete. Panel A, column 2 shows that, as expected, after treatment the median page in the control group is relatively shorter and of relatively lower quality. But before the treatment (panel B), the median pages in the treatment and control groups are similar.

Columns 3–4 in Table 2 present the comparison with the pages in the Spanish Wikipedia. Compared to Spanish, the pages in the sample are even shorter and less similar. After the treatment, the median page in the treatment group is still only 50% of the length of the corresponding city page in Spanish. While it is not clear that all that material should be included in the French, German, and Italian Wikipedia, at least we can say that there is additional material that was important enough to be included in the Spanish Wikipedia. Also note that the treatment added only material on topics relevant for tourists, such as the city's main sights and culture. Hence, while the treatment made the pages more complete on these topics, on other topics

²⁰ The three lowest grades in the Wikipedia content assessment are: Stub—a very basic description of the topic; Start—developing but still quite incomplete; Class C—substantial but is still missing important content or contains much irrelevant material. Source: https://en.wikipedia.org/wiki/Wikipedia:Content_assessment.

²¹ Note that these articles, on average, are also not very important according to the English Wikipedia article importance scheme. The scheme uses ratings from Low, Mid, High, to Top. Only 7 of the pages are rated as High importance, 15 as Mid, and 9 as Low importance, 29 of the articles have not been assigned an importance rating, which probably also implies that those articles are not highly important.

Table 1 Comparison of pre-treatment characteristics in the treatment versus the control group

	Control group mean	Treatment group mean	t-test <i>p</i> -value	Wilcoxon test <i>p</i> -value	Obs.
	(1)	(2)	(3)	(4)	(5)
Log. length before treatment	8.586	8.611	0.842	0.655	180
Quality rating before treatment	72.700	70.917	0.810	0.602	180
Quality: complete before treatment	71.656	71.706	0.996	0.758	180
Quality: interesting before treatment	75.800	72.917	0.777	0.561	180
Quality: well-written before treatment	92.867	93.122	0.858	0.388	180
Quality: illustrated before treatment	72.372	73.772	0.854	0.441	180
Similarity to Spanish before treatment	17.488	17.251	0.885	0.820	180
Aver. # of users before treatment	0.378	0.333	0.321	0.285	180
Aver. # of edits before treatment	0.396	0.351	0.341	0.364	180
Aver. edit dist. before treatment	79.013	76.786	0.882	0.738	180
Aver. capped edit dist. before treatment	42.511	38.238	0.470	0.583	180

Column 1 and 2 present the means of pre-treatment values of variables, separately for the control and the treatment group. Column 3 presents the *p*-value of the t-test for whether the difference between the control and treatment groups is significantly different from zero. Column 4 presents the *p*-value of the corresponding Wilcoxon rank-sum test. Column 5 presents the number of observations used in each test

that are typically covered on each city page, such as demographics, economy, education, and government, there is probably still room for improvement.

4 Results

We start by estimating the impact of the treatment on growth after treatment, both in length and quality. We are looking at growth during 4 years after treatment. After that, we go into more details in two ways. First, we study the effects on different dimensions of quality. Then we analyze the short- and long-term effects using a difference-in-differences estimator.

To estimate the effect of the treatment on growth after treatment, we compare the growth of pages (indexed by *i*) in the treatment and control group controlling for city and language fixed effects:

$$\Delta y_i = \beta_0 + \beta_1 \text{TreatmentGroup}_i + \text{LanguageFE}_i + \text{CityFE}_i + \varepsilon_i \quad (2)$$

The coefficient of interest is β_1 on TreatmentGroup_i , which is an indicator variable that takes value one if the page was assigned to the treatment group and zero if it was assigned to the control group. The outcome variable $\Delta y_i = y_{2018\text{September}} - y_{2014\text{September}}$ measures the change in outcome from September 2014 to September 2018. Specifically, the outcome variables are the change in the

Table 2 Completeness of the median page compared to English and Spanish Wikipedia

	Compared to English		Compared to Spanish	
	Treatment	Control	Treatment	Control
	(1)	(2)	(3)	(4)
<i>Panel A: After treatment (2014 September)</i>				
Relative length	75.6	53.4	50.2	35.2
Similarity	33.1	27.8	24.1	17.5
Completeness	98.4	72.4	.	.
<i>Panel B: Before treatment (2014 August)</i>				
Relative length	57.9	53.2	37.7	35.1
Similarity	26.2	27.7	17.3	17.5
Completeness	71.7	71.7	.	.

Each cell presents the median from 90 pages either in the treatment group (columns 1 and 3) or control group (columns 2 and 4)

logarithm of page length, change in the overall quality rating, and the change in the similarity to the corresponding Spanish Wikipedia article.²²

Table 3 presents the estimates of the effect of treatment on subsequent growth over 4 years post-treatment.²³ The point estimates are 7% of the standard deviation for length and -3% for quality.²⁴ But the estimates are imprecise. The 95%-confidence interval for length is from -7 to + 12%, and the ex-post minimum detectable effect size is 13%. Similarly, the 95%-confidence interval for the growth in quality is from -4 to + 3 points and the ex-post minimum detectable effect size is 5 points. Note that these bounds are no more than one third of the size of the treatment itself. The estimates for the similarity index are analogous to the quality rating.

Table 4 presents the results of the treatment on the growth in different dimensions of quality. In all four dimensions of quality, measuring how complete, interesting, illustrated, and well-written the page is, the estimates are not statistically significant. We can reject that the effect of treatment on subsequent growth is more than one fourth the size of the treatment itself for the completeness and interesting quality dimensions. The largest point estimate is for well-written—the effect of treatment on subsequent growth in this dimension is about 5 points but the estimates are

²² To make the pages comparable, we subtract the length of text added by the treatment from the length of pages in the treatment group after treatment (both in 2014 and in 2018). We do that because the outcome variable measures the percentage change, and hence, without subtracting the length of the treatment text, the treatment group would have a higher base when calculating the percentage, then the same increase in characters would give a smaller percentage for the treatment group.

²³ Table B.2 in online appendix shows that the results are robust to alternative control variables. Table B.3 shows that the results are also robust when including the Dutch pages either in the control group or estimating intention-to-treat.

²⁴ For expositional clarity, we interpret the coefficient in column 1 as measuring a percentage change in length. This is a logarithmic approximation that performs well when changes are small which is the case in our sample.

imprecise. This is consistent with what we would expect, given that the treatment itself reduced the quality in this dimension by 3 points.

Next, we analyze both the short- and long-term effects. We estimate the following difference-in-differences regression using page-level monthly panel data:

$$y_{it} = \sum_s \beta_s \cdot \mathbf{1}[Year_t = s] \cdot TreatmentGroup_i + MonthFE_t + LanguageCityFE_i + \varepsilon_{it} \quad (3)$$

where the sum over s is taken over the years $\{-3, -2, -1, 1, 2, 3, 4\}$ and $Year_t$ measures years since the experiment. The year of the experiment is the baseline, which is why $s = 0$ is excluded. The regression includes fixed effects for each page i (language and city pair) and for each time period t . The coefficients of interest are the β -s on $TreatmentGroup_i$ and year dummy interactions. All the year and treatment group interactions, including for pre-treatment years, are presented graphically in Figure B.11 in online appendix. It shows that there is no evidence of differential pre-treatment trends.

Table 5 presents the estimates of short- and long-term effects from regression (3). In column 1, the outcome variable is the logarithm of page length minus treatment text.²⁵ The estimates of the short- and long-term impacts of treatment on page length are all insignificant at ten percent level. The estimate of the long-term (4-year) effect on page length is similar to that in Table 3.

In columns 2 and 3, the outcome variables are the number of users (people editing the page) and edits per month. The treatment increased the number of users and the number of edits during the first 2 years after the experiment. Specifically, the treatment increased the monthly number of users editing the page by about 0.11 users (column 2) and the monthly number of edits by 0.12–0.13 edits (column 3). However, these increases are only short-lived. In the third and fourth year, for both measures, the estimates of treatment effect are insignificant at ten percent level, and the point estimates are almost ten times smaller in magnitude. In the fourth year, we can reject an effect of the same size found in the first years (0.11 users).²⁶

What do these editors and edits do in the first 2 years post-treatment if it has surprisingly little effect on page length? A natural explanation could be that the additional edits simply polish the text added by the treatment. To study this, we re-calculated the number of edits per month while excluding edits that directly edited the text added by the treatment. The estimates using this outcome variable are presented in column 4. The results show that when excluding the edits that directly affect the text added by the treatment, then the treatment effect is much smaller. We conclude that most of the short-term increase in editing comes from editing the content added by the treatment.

²⁵ To make the pages comparable, we subtract the length of text added by the treatment from the length of pages in the treatment group after treatment. Hence, the estimates should be interpreted as the effect of treatment on page length after removing the mechanical increase created by the treatment.

²⁶ For the long-term (4-year) effect, the ex-post minimum detectable effect size is 0.11 users.

Table 3 The long-term effect of treatment on the growth in page length and quality

	Change in page length or quality ($y_{2018Sep} - y_{2014Sep}$)		
	Δ Log. page length	Δ Quality rating	Δ Similarity to Spanish
	(1)	(2)	(3)
Treatment group	0.026 (0.048)	-0.375 (1.777)	-0.590 (0.624)
Language FE	Yes	Yes	Yes
City FE	Yes	Yes	Yes
Mean dep. var.	0.190	6.589	2.032
SD dep. var.	0.353	11.958	4.409
Adj. R-squared	0.259	0.116	0.199
Observations	180	180	180

Each column presents estimates from a separate cross-section regression of 180 Wikipedia pages. The dependent variable $\Delta y_i = y_{2018September} - y_{2014September}$ is the change in logarithm of page length (column 1), change in the overall quality rating (column 2), and the change in the similarity to the corresponding Spanish Wikipedia article (column 3). All regressions include language fixed effects and city fixed effects. Standard errors are reported in parentheses

Table 4 The long-term effect of treatment on the growth in different dimensions of page quality

	Change in page quality ($y_{2018Sep} - y_{2014Sep}$)			
	Δ Complete	Δ Interesting	Δ Illustrated	Δ Well-written
	(1)	(2)	(3)	(4)
Treatment group	-0.967 (2.254)	-0.150 (1.851)	-2.037 (3.752)	4.658 (5.980)
Language FE	Yes	Yes	Yes	Yes
City FE	Yes	Yes	Yes	Yes
Mean dep. var.	6.939	5.700	12.186	4.875
SD dep. var.	15.212	12.201	27.280	37.879
Adj. R-squared	0.122	0.080	0.243	0.003
Observations	180	180	180	180

Each column presents estimates from a separate cross-section regression of 180 Wikipedia pages. The dependent variable $\Delta y_i = y_{2018September} - y_{2014September}$ is the change in the following dimensions of page quality: complete (column 1), interesting (column 2), illustrated (column 3), and well-written (column 4). All regressions include language fixed effects and city fixed effects. Standard errors are reported in parentheses

Table 5 The short- and long-term effects of treatment on subsequent page length and editing activity

	Log. length excluding treat- ment	# Users	# Edits	# Edits excluding treatment	Edit distance	Capped edit distance
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment × year 1	0.045 (0.032)	0.112*** (0.043)	0.118** (0.049)	0.006 (0.045)	- 13.609 (40.176)	12.970 (9.250)
Treatment × year 2	0.051 (0.039)	0.109** (0.047)	0.128** (0.056)	0.067 (0.055)	107.260 (101.633)	17.967* (10.414)
Treatment × year 3	0.049 (0.047)	0.008 (0.043)	0.001 (0.049)	- 0.046 (0.049)	- 37.161 (30.389)	- 6.849 (9.412)
Treatment × year 4	0.020 (0.056)	0.012 (0.041)	0.013 (0.046)	- 0.041 (0.044)	2.626 (74.799)	- 3.498 (8.922)
Month FE	Yes	Yes	Yes	Yes	Yes	Yes
Treatment × { year-1,-2,-3 }	Yes	Yes	Yes	Yes	Yes	Yes
Language-City FE	Yes	Yes	Yes	Yes	Yes	Yes
Mean dep. var.	8.651	0.329	0.347	0.329	89.491	36.062
SD dep. var.	0.813	0.694	0.759	0.740	1120.108	131.787
Observations	17,100	17,100	17,100	17,100	17,100	17,100

A unit of observation is a page-month pair. The dependent variable is the logarithm of page length minus treatment text (column 1) or a monthly number of users (columns 2), edits (columns 3), edits not editing treatment text (column 4), edit distance (column 5), or capped edit distance (column 6). *Treatment × year l* is an indicator variable that takes value one during the first year post-treatment if the page belongs to the treatment group and zero otherwise; and similarly for the other years. All regressions include page fixed effects and month fixed effects. All the year and treatment group interactions, including for pre-treatment years, are presented graphically in Figure B.11 in online appendix. The sample is a balanced sample from September 2010 to August 2018, excluding the treatment month of August 2014. Standard errors, reported in parentheses, are clustered by page (180 pages). ***Indicates significance at the 1 percent level, **At a 5 percent level, *At a 10 percent level

In column 5, the outcome variable is the number of characters added plus deleted. The treatment had no statistically significant (at ten percent level) effect on the measure. Coefficients vary in sign over the years. The largest point estimate is in the second year: the (statistically insignificant) coefficient estimate implies the treatment effect of about 100 characters per month.

Because the distribution of the number of characters added plus deleted has a long tail, in column 6, we use an alternative measure calculated from individual capped edits. The individual edits are capped from above at the 90th percentile. The 90th percentile equals about 500 characters and is about 10 times larger than the median edit. In this way, the capped edit distance measure gives smaller weight to long edits. Estimates in column 6 show that in the second year after the experiment, the treatment increases the capped change in characters by about 18 characters per month. Provided that the average word length across languages in the experiment is about 10.4 characters, our treatment increased the edit distance by about two

words.²⁷ This is in line with the findings from columns 3–4, which showed that most of the increase in editing comes from editing the content added by the treatment. In later years, the point estimates of the treatment effect are even smaller, and the estimates are statistically insignificant.

Finally, we analyze whether the effect of the treatment is heterogeneous, varying by subgroups of the pages. Tables B.4–B.5 in online appendix re-estimate the regressions in Table 3, allowing the effect of treatment to vary by the quality, completeness, relative length (compared to the pages in Spanish Wikipedia), and the age of the page. While all the estimates are statistically insignificant, the pattern in the point estimates seems to suggest that the treatment effect was larger on lower quality and less complete pages.

Multiple hypothesis testing We run many tests, and with a large number of tests, we could easily get false positives, i.e., simply by chance, some estimates could turn out to be statistically significant. That is, the probability of incorrectly rejecting at least one null hypothesis is greater than the probability of incorrectly rejecting each individual hypothesis test. To address the concern, we adjust for multiple hypothesis testing by adjusting for the family-wise error rate (the probability of incorrectly rejecting at least one null hypothesis belonging to the same family of hypothesis).

We prefer not to assign all the tested hypotheses to the same family because we view that some of our outcomes are more important than others. Specifically, we are mainly interested in long-term outcomes and on the impact on length and quality, not on the inputs like various measures of editing. Therefore, to adjust the family-wise error rate, we group our hypotheses into the following families: (1) long-term effects on length and quality, (2) short-term effects on editing behavior, (3) long-term effects on editing behavior. Specifically, we use Westfall and Young (1993) multiple hypothesis p -value adjustment as implemented by Jones et al. (2019), employing 10,000 bootstrap draws.

Tables B.6–B.8 in online appendix present the adjusted p -values. As expected, the adjusted p -values are much larger, and therefore our conclusions regarding the long-term effects are unchanged. Using the stricter p -values, the short-term effects on editing activity become statistically insignificant. Therefore, a conservative approach would be to interpret the findings of the short-term effects simply as suggestive. On the other hand, the stricter p -values might be viewed as too conservative, considering that the literature oftentimes does not correct for multiple hypothesis testing, and the number of our tests has increased in part to show robustness. For this reason, in our preferred estimates, we use unadjusted p -values.

5 Theoretical framework

To provide a framework for interpreting our results and unify findings from the related literature, we introduce a simple theoretical model of the private provision of public goods. We postpone the technical details of the model to online appendix

²⁷ Source: <http://www.ravi.io/language-word-lengths>.

C and present its implications by way of Fig. 4 here. The theoretical framework involves a sequence of players making public good contributions. Depending on the current state of the public good (denoted by X), each contribution can have positive or negative externalities for future contributions.

We first show that if there are no externalities, the rate of contributions is on average constant over time and therefore growth path of the public good value X is linear, as illustrated in Fig. 4a.²⁸ This means that a jump in contributions (such as experimental treatment) would lead to a parallel shift in the growth path and the long-term outcomes are shifted by about the same magnitude as the initial treatment. With negative externality, the rate of contributions decreases with the state and therefore growth path is decreasing and long-term effect of a jump in contributions is smaller than the short-term effect (Fig. 4b). Conversely, with positive externality, growth path is increasing over time and long-term effect of a jump in contributions is bigger than the short-term effect (Fig. 4c). Finally, as we argue below, perhaps the most realistic scenario for the externalities is an U-shaped externality function—initially the externality is positive and eventually negative. This would leave to S-shaped growth path (Fig. 4d). In this case the comparison of short-term and long-term impacts depends on time horizon as well as the moment of treatment.

Let us now discuss how our empirical results can be interpreted within this theoretical framework. The first question to ask is which channels lead to either positive or negative externalities? We can highlight several channels leading to positive externalities. *Attention*: some Wikipedia editors sign up as watchers for a page and get notifications for each time the page is edited. The reason they do it is to maintain the quality of the page, and therefore some of these notifications must lead to future edits. Moreover, contributions increase the number of page views, mostly by increasing visibility in search rankings (Hinnosaar et al., 2021b). Some of these additional eyeballs are likely to convert to contributions (Kane & Ransbotham, 2016; Kummer, 2020; Zhu et al., 2020). *Learning and inspiration*: existing content provides contributors new information about a topic and gives them ideas for contributions (Olivera et al., 2008). *Social motives*: additional content signals potential interest in the community (Chen et al., 2019; Zhang & Zhu, 2011). All these channels would lead to positive externalities, i.e., more contributions by earlier editors either make it easier to contribute or raise their benefits from contributions. Without loss of generality, this can be modeled as a reduction of the marginal cost of contributions, as we did in the previous section. These channels likely have a declining impact, i.e., during the early stages of the page development, we would expect the impacts to be larger than for mature pages, which already have a large amount of content.

On the other hand, there are also channels that lead to negative externalities: *Crowding out*: if users add the content that adds most value at the least possible cost

²⁸ Without externalities, our model is similar to the model in Chen et al. (2019), with two differences: the payoffs in their model depend on social impact (i.e., number of viewers) and participation is endogenous. A crucial difference in our model is the inclusion of externalities. The model could be generalized and solved using tools introduced in Hinnosaar (2018).

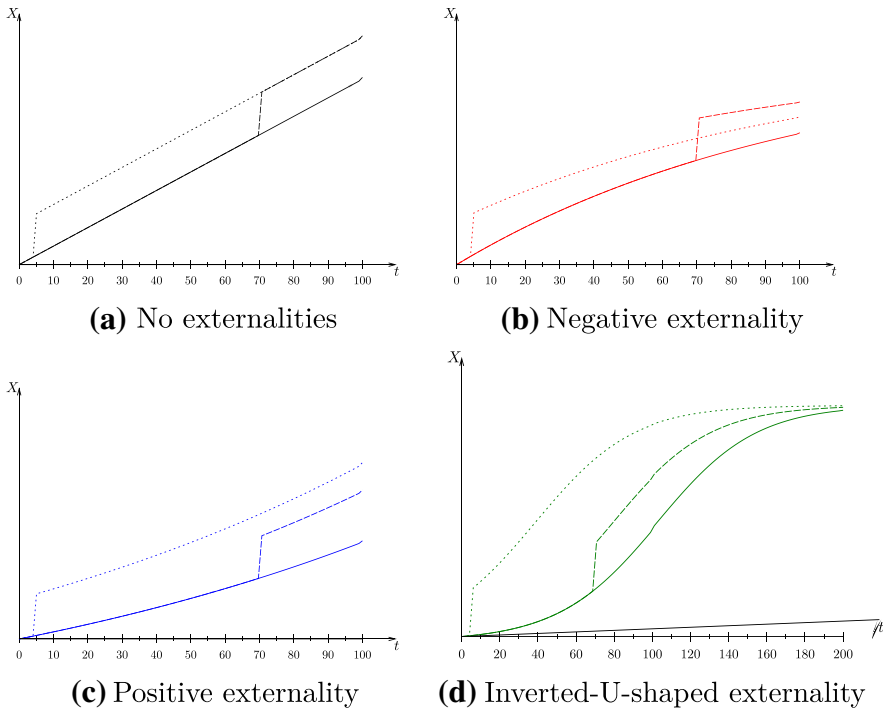


Fig. 4 The expected growth of the state under different assumptions about externalities. *Notes:* Time (or number of contributors) on the horizontal axis and state of the public good on the vertical axis. For more details, see online appendix C

(“low hanging fruits”), then new content raises the cost of future contributions. *Freeriding*: new content tells potential editors that someone is already taking care of the edits of this particular page, and they may choose to point their attention elsewhere. *Increasing complexity*: a Wikipedia page should be a coherent reference source; the more content there is, the more possibilities there are for combinations on how the material could be organized and more parts that need to have the same style and structure. All these channels lead to negative externalities, i.e., more contributions by earlier editors either make it more difficult to contribute or reduce their benefits from contributions. Without loss of generality, this can be modeled as an increasing marginal cost of contribution. It is likely that these channels for negative externalities become more prevalent as content matures and converges to completeness.

Combining these observations, in different settings, the total effect may lean either towards positive or negative externalities, depending on which channels are more active. It is likely that for relatively new and incomplete pages, channels inducing positive externalities are more prevalent. As the page matures, channels for positive externalities become less important, and channels for negative externalities more important. This would imply an inverted-U-shaped externality that we saw in the previous section.

The remaining question is how large are these externalities, which is the empirical question that our analysis addresses. Our empirical results show that a treatment Δ has approximately the same effect Δ on the outcomes 4 years later. This finding is consistent with two possibilities: either there are no externalities (Fig. 4a), or there are both positive and negative externalities in the same magnitude (Fig. 4d). Both cases are plausible, as the treatment contributed to relatively mature pages that were still far from being complete. The results about treatment heterogeneity (Table B.4–B.5) provide some evidence in favor of the second possibility as they seem to indicate that the treatment had a more positive effect on lower quality and less complete pages.

6 Discussion

6.1 Comparison with related literature

Aaltonen and Seiler (2016) and Kane and Ransbotham (2016) study the impact of added content on the short-term editing activity. The estimates from Aaltonen and Seiler (2016) imply that additional content of 2000 characters leads to about 0.18 additional monthly users. Zhu et al. (2020) estimate the impact over 6 months and find an effect of similar magnitude, about 0.21 additional monthly contributors per 2000 added characters. Our estimates largely confirm the findings from the literature. The magnitude of our estimates is slightly smaller (about 0.11 additional users per month), the difference might be explained by the fact that the other papers are measuring more immediate effects, while we measure the impact over a year or two.

However, we find that there is no large long-term impact on editing activity or content growth. This finding is related to the results in Kane and Ransbotham (2016), who find that although additional content may bring additional contributions to pages that are relatively incomplete, this effect may disappear once the pages become complete. Our results on treatment heterogeneity provide some evidence towards this, although the effects are not statistically significant.

Such differences in editing activity are consistent with a simple explanation that in earlier stages of the content life-cycle, the channels leading to positive externalities are more active, while the channels leading to negative externalities become more dominant in the long-term. In particular, as some editors have signed up as watchers who make sure that added content is up to the quality standards (see Sect. 2.3), we would expect that immediately after content is added there is increased editing activity focusing on the added content. Our analysis of the short-term editing activity supports this conjecture.²⁹

²⁹ Note that there are other, more subtle differences in the research environments that may affect the outcomes. For example, in our setting, the added content comes from an outside source. Instead, in the settings of Kane and Ransbotham (2016), Aaltonen and Seiler (2016), and Zhu et al. (2020), the added content is created by the community.

On the other hand, Nagaraj (2019) studies the impact of better-quality seeding data on quality outcomes 10 years later. He finds a negative effect of about 10%, which is a large difference from the conclusion of Aaltonen and Seiler (2016), whose simulations implied that the positive effect on short-term editing activity could lead to about 45% better output. Our results about the impact of added content on the long-term growth in content quantity and quality provide some support in favor of the first number, but we do not have enough power to measure small effects.

The finding from Nagaraj (2019) is partially consistent with the implications of our theoretical framework. Assuming that a time horizon of 10 years is sufficient for the content to be close to completeness, we would expect that early differences in content would disappear in the long-term. In other words, in the long term, we would expect the negative externality to dominate. Explaining the negative impact on outcomes requires an explanation beyond our model and analysis. The key mechanism proposed by Nagaraj (2019) to explain the negative externality is the “ownership effect”, which plays a less prominent role in Wikipedia. Nagaraj (2019) suggests that contributors who added particular bridges or streets on the user-generated map may feel more responsible for keeping these objects updated over time. Therefore, the treatment of adding more seeding information may backfire by not allowing the ownership of objects to arise naturally. All other papers mentioned here, including our work, focus on textual content in Wikipedia, where ownership is less clear, and we would thus expect the negative effect of adding content to be less prominent.

6.2 Generalizability

A word of caution is in order here. The results from Wikipedia might not generalize to other user-generated content platforms. As an example, a relevant difference among the platforms is the magnitude of the contributor’s personal benefit. In Wikipedia, the personal benefit from contributing is likely to be smaller than in open-maps or open-source software. For example, a user of open-maps could directly benefit from correcting a mistake on a map, while an error in Wikipedia is unlikely to have any personal consequences. Another example is the “ownership effect” proposed by Nagaraj (2019), which seems to be a key driver in the Wikipedia-style mapping service, where each object (street, building, etc.) explicitly states who has edited this object. As Wikipedia does not display who wrote each part of the page, the ownership assignment is less clear. Therefore, we would expect this channel for negative externality to be less active.

Moreover, our results also might not generalize to settings in other stages in content development. Our theoretical model highlighted that the externalities might depend on the existing amount of available content, additional content in early stages being more beneficial than later. Our empirical results about the heterogeneity provide some suggestive evidence that this might be the case. However, the experiment was not designed to analyze this heterogeneity and is underpowered to do that. We would encourage future research that aims at uncovering how the existing content (or more generally the lifecycle of the content development) affects the externalities in content creation.

7 Conclusions

In this paper, we studied the impact of added content on the subsequent long-term growth of content. We identify the causal effect using exogenous variation from a randomized field experiment in Wikipedia where the treatment added content to randomly chosen Wikipedia pages. Our estimates allow us to rule out effects on 4-year growth of content length larger than twelve percent. We can also rule out the effects on 4-year growth of content quality larger than four points, which is less than one-fifth of the size of the treatment itself. We find that the treatment increased subsequent content generation in the first 2 years. Specifically, it increased the number of edits and editors. However, the amount of content these users added was on average only a couple of words, and most of their edits modified the content added by the treatment.

Our findings have a clear policy implication—in settings like the one studied here, information seeding and motivating content creation is not necessarily enough to generate a large increase in future content generation. However, these policies are also not counterproductive as they can stimulate additional edits, and we can rule out a large discouragement effect on future contributions. Therefore, it is mostly a matter of direct cost-benefit analysis whether such policies pay off.

Our results may not generalize to settings where other channels leading to positive or negative externalities are more prominent. For example, it is possible that in the early stages of content development, information seeding may be beneficial. On the other hand, it is also possible that in situations where individual contributions are well-identified, the seeding policies may backfire.

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s10683-021-09730-x>.

Acknowledgements We are grateful to John Duffy, two anonymous referees, Yan Chen, Chris Forman, Willa Friedman, Shane Greenstein, David Hugh-Jones, Tobias Kretschmer, Giovanni Mastrobuoni, Ignacio Monzón, Juan Morales, Abhishek Nagaraj, Stefan Penczynski, Imke Reimers, Stephan Seiler, Ananya Sen, Matthias Sutter, and Michael Zhang for valuable comments. We would also like to thank the seminar participants at the University of Strasbourg, the Collegio Carlo Alberto, George Mason University, ParisTech 2019, Digital Economy Workshop 2019 (Católica Lisbon), ZEW Conference on the Economics of ICT, the Advances with Field Experiments 2019 Conference (University of Chicago), SED 2021 for valuable input.

References

- Aaltonen, A., & Seiler, S. (2016). Cumulative growth in user-generated content production: Evidence from Wikipedia. *Management Science*, *62*, 2054–2069.
- Algan, Y., Benkler, Y., Morell, M. F., & Hergueux, J. (2013). Cooperation in a peer production economy: experimental evidence from Wikipedia. *manuscript*.
- Ayres, I., Raseman, S., & Shih, A. (2013). Evidence from two large field experiments that peer comparison feedback can reduce residential energy usage. *Journal of Law, Economics, and Organization*, *29*, 992–1022.
- Chen, Y., Farzan, R., Kraut, R. E., YeckehZaare, I., & Zhang, A. F. (2019). Motivating contributions to public information goods: A personalized field experiment on Wikipedia. *manuscript*.

- Chen, Y., Harper, F. M., Konstan, J., & Li, S. X. (2010). Social comparisons and contributions to online communities: A field experiment on movie lens. *American Economic Review*, *100*, 1358–98.
- De Giorgi, G., Pellizzari, M., & Redaelli, S. (2010). Identification of social interactions through partially overlapping peer groups. *American Economic Journal: Applied Economics*, *2*, 241–75.
- Duflo, E., & Saez, E. (2002). Participation and investment decisions in a retirement plan: The influence of colleagues' choices. *Journal of Public Economics*, *85*, 121–148.
- Duflo, E., & Saez, E. (2003). The role of information and social interactions in retirement plan decisions: Evidence from a randomized experiment. *Quarterly Journal of Economics*, *118*, 815–842.
- Egebark, J., & Ekström, M. (2017). Liking what others “like”: Using Facebook to identify determinants of conformity. *Experimental Economics*, *21*, 1–22.
- Eurobarometer. (2012). *Europeans and their languages, special report 386*. European Commission.
- Fershtman, C., & Gandal, N. (2011). Direct and indirect knowledge spillovers: The “social network” of open-source projects. *RAND Journal of Economics*, *42*, 70–91.
- Gallus, J. (2017). Fostering public good contributions with symbolic awards: A large-scale natural field experiment at Wikipedia. *Management Science*, *63*, 3999–4015.
- Goldstein, N. J., Cialdini, R. B., & Griskevicius, V. (2008). A room with a viewpoint: Using social norms to motivate environmental conservation in hotels. *Journal of Consumer Research*, *35*, 472–482.
- Greenstein, S., & Zhu, F. (2012). Is Wikipedia biased? In *American Economic Review: Papers and Proceedings* (pp. 343–348).
- Greenstein, S., & Zhu, F. (2018). Do experts or crowd-based models produce more bias? Evidence from Encyclopedia Britannica and Wikipedia. *MIS Quarterly*, *42*, 945–959.
- Grossman, G. M., & Helpman, E. (1993). *Innovation and growth in the global economy*. MIT Press.
- Hanushek, E. A., Kain, J. F., Markman, J. M., & Rivkin, S. G. (2003). Does peer ability affect student achievement? *Journal of Applied Econometrics*, *18*, 527–544.
- Hinnosaar, M. (2019). Gender inequality in new media: Evidence from Wikipedia. *Journal of Economic Behavior & Organization*, *163*, 262–276.
- Hinnosaar, M., Hinnosaar, T., Kummer, M., & Slivko, O. (2021a). Replication data for: “Externalities in knowledge production: Evidence from a randomized field experiment”. *Harvard Dataverse*. <https://doi.org/10.7910/DVN/T4VFCX>.
- Hinnosaar, M., Hinnosaar, T., Kummer, M., & Slivko, O. (2021b). Wikipedia matters. *Journal of Economics & Management Strategy*, *163*, 1–13.
- Hinnosaar, T. (2018). Optimal sequential contests, *manuscript*.
- Huang, N., Burch, G., Gu, B., Hong, Y., Liang, C., Wang, K., et al. (2018). Motivating user generated content with performance feedback: Evidence from randomized field experiments. *Management Science*, *65*, 327–345.
- Jones, C. I. (1995). R&D-based models of economic growth. *Journal of Political Economy*, *103*, 759–784.
- Jones, D., Molitor, D., & Reif, J. (2019). What do workplace wellness programs do? Evidence from the Illinois workplace wellness study. *Quarterly Journal of Economics*, *134*, 1747–1791.
- Kane, G. C., & Ransbotham, S. (2016). Content as community regulator: The recursive relationship between consumption and contribution in open collaboration communities. *Organization Science*, *27*, 1258–1274.
- Kiyotaki, N., & Wright, R. (1989). On money as a medium of exchange. *Journal of Political Economy*, *97*, 927–954.
- Kummer, M. (2020). Attention in the peer production of user generated content: Evidence from 93 pseudo-experiments on Wikipedia, *manuscript*.
- Kummer, M., Slivko, O., & Zhang, X. (2019). Unemployment and digital public goods contribution. *Information Systems Research*, *31*, 801–819.
- Lacetera, N., & Macis, M. (2010). Social image concerns and prosocial behavior: Field evidence from a nonlinear incentive scheme. *Journal of Economic Behavior & Organization*, *76*, 225–237.
- Lerner, J., & Tirole, J. (2003). Some simple economics of open source. *Journal of Industrial Economics*, *50*, 197–234.
- Manski, C. F. (1993). Identification of endogenous social effects: The reflection problem. *Review of Economic Studies*, *60*, 531–542.
- Nagaraj, A. (2019). *Information seeding and knowledge production in online communities: Evidence from OpenStreetMap*, *manuscript*.
- Nov, O. (2007). What motivates Wikipedians? *Communications of the ACM*, *50*, 60–64.

- Olivera, F., Goodman, P. S., & Tan, S.S.-L. (2008). Contribution behaviors in distributed environments. *MIS Quarterly*, *32*, 23–42.
- Porter, M. F. (1980). An algorithm for suffix stripping. *Program*, *14*, 130–137.
- Ransbotham, S., Kane, G. C., & Lurie, N. H. (2012). Network characteristics and the value of collaborative user-generated content. *Marketing Science*, *31*, 387–405.
- Ren, Y., Chen, J., & Riedl, J. (2015). The impact and evolution of group diversity in online open collaboration. *Management Science*, *62*, 1668–1686.
- Romer, P. M. (1990). Endogenous technological change. *Journal of Political Economy*, *98*, S71–S102.
- Salganik, M. J., Dodds, P. S., & Watts, D. J. (2006). Experimental study of inequality and unpredictability in an artificial cultural market. *Science*, *311*, 854–856.
- Shah, S. K. (2006). Motivation, governance, and the viability of hybrid forms in open source software development. *Management Science*, *52*, 1000–1014.
- Slivko, O. (2018). Online “brain gain”: Do immigrants return knowledge home? *manuscript*.
- Sun, Y., Dong, X., & McIntyre, S. (2017). Motivation of user-generated content: Social connectedness moderates the effects of monetary rewards. *Marketing Science*, *36*, 329–337.
- Thompson, N., & Hanley, D. (2018). Science is shaped by Wikipedia: Evidence from a randomized control trial. *manuscript*.
- Tversky, A. (1977). Features of similarity. *Psychological Review*, *84*, 327–352.
- Westfall, P. H., & Young, S. S. (1993). *Resampling-based multiple testing: Examples and methods for p-value adjustment* (Vol. 279). Wiley.
- Xu, S. X., & Zhang, X. (2013). Impact of Wikipedia on market information environment: Evidence on management disclosure and investor reaction. *MIS Quarterly*, *37*, 1043–1068.
- Zhang, X., & Zhu, F. (2011). Group size and incentives to contribute: A natural experiment at Chinese Wikipedia. *American Economic Review*, *101*, 1601–1615.
- Zhu, K., Walker, D., & Muchnik, L. (2020). Content growth and attention contagion in information networks: A natural experiment on Wikipedia. *Information Systems Research*, *31*, 491–509.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Authors and Affiliations

Marit Hinnosaar^{1,4,5} · Toomas Hinnosaar^{1,5}  · Michael E. Kummer^{2,6,7} · Olga Slivko³

Marit Hinnosaar
marit.hinnosaar@gmail.com

Michael E. Kummer
M.Kummer@uea.ac.uk

Olga Slivko
slivko@rsm.nl

- ¹ University of Nottingham, Nottingham, UK
- ² University of East Anglia, Norwich, UK
- ³ Rotterdam School of Management, Erasmus University, Rotterdam, The Netherlands
- ⁴ Collegio Carlo Alberto, Turin, Italy
- ⁵ CEPR, London, UK
- ⁶ Georgia Institute of Technology, Atlanta, USA
- ⁷ ZEW, Mannheim, Germany