

Karlsruhe Reports in Informatics 2011,30

Edited by Karlsruhe Institute of Technology,
Faculty of Informatics
ISSN 2190-4782

Revealing the Suitability of Incentive Mechanisms for the Collaborative Creation of Structured Knowledge

Conny Kühne and Klemens Böhm

2011



Fakultät für Informatik

Please note:

This Report has been published on the Internet under the following
Creative Commons License:

<http://creativecommons.org/licenses/by-nc-nd/3.0/de>.

Revealing the Suitability of Incentive Mechanisms for the Collaborative Creation of Structured Knowledge

Conny Kühne and Klemens Böhm

Karlsruhe Institute of Technology (KIT), Germany
{conny.kuehne|klemens.boehm}@kit.edu

Creating and maintaining semantic structures such as ontologies continues to be an important issue. The approach investigated here is to let members of an online community create structured knowledge collaboratively and to use ratings to evaluate the data created. Obviously, such ratings have to be of high quality. Honest rating mechanisms (HRMs) known from literature are a promising means to gain such high-quality ratings. However, the design of such mechanisms for collaborative knowledge creation and their effectiveness have not been studied so far. To evaluate the effects of an HRM on rating quality in this context, we have conducted several experiments with online communities. We find that an HRM increases rating quality and “punishes” rating errors. We also find that rating-based rewards increase the quality of the structured knowledge created.

1 Introduction

Structured knowledge, e.g., in the form of ontologies, has become increasingly important. The question of how to create it on a larger scale, however, continues to be a fundamental research issue. The approach investigated in this paper is web-based peer production [5]. It is a decentralized production process where contributors work on a project without a hierarchical organization. It is both scalable when adding further users and robust because individual users can be replaced.

When a community of peers, without a coordinating authority, creates and maintains structured knowledge, two questions arise: 1) How to motivate users? Many online communities suffer from under-contribution [4], i.e., only a minority contributes. 2) How to ensure and assess data quality? Compared to, say, Wikipedia, quality assurance may be even more important for structured knowledge, which is used for query processing or automated reasoning.

In this paper, we study how rating-based incentive mechanisms can assure the quality of the structured knowledge created collaboratively. We envision the following real-world scenario: A community creates structured knowledge collaboratively. Its members review the contributions of each other and rate them according to the quality perceived. For instance, think of a project with partners who are geographically distributed, and who all contribute to a common knowledge base. To motivate users to contribute, they receive rewards according to the number and quality of their contributions. The quality of contributions, in

turn, is computed based on the ratings. Finally, the points gathered are converted into external rewards, e.g., gift coupons as with Epinions or system privileges as with Slashdot. Ratings have to be of high quality as well. A promising approach to gain high-quality ratings are *honest rating mechanisms* (HRMs) [13, 14, 16]. An HRM rewards subjective truthfulness in scenarios where no objective truth criterion is available. To the best of our knowledge, these mechanisms have not been studied yet in the context of collaborative knowledge creation.

In this paper, we study the following research questions regarding the creation of structured knowledge.

1. How do reward mechanisms for contributions as well as for ratings influence user behavior?
2. Does an HRM lead to ratings of higher quality, compared to a fixed reward per rating?
3. How do rewards for contributions dependent on ratings influence the quality and the number of contributions, compared to rewards that are fixed?

To this end, we have developed a platform for the collaborative creation of structured knowledge called *Consensus Builder* (CB). It features fine-grained rating and incentive mechanisms, in particular an HRM, and is operational in a real-world environment¹. Based on the research questions, we formulate a number of hypotheses and test them in a series of controlled field experiments. The setups are close to the real-world scenario envisioned. Controlled experiments allow us to gain insights into causal effects of the tested mechanisms. This cannot be achieved by observational studies. For each experiment, we have recruited a different community. Participants have used CB from home or from their workplace to create structured data.

We make the following contributions:

- We describe the features of Consensus Builder, our platform for the collaborative construction of structured knowledge. It contains reward mechanisms for contributions and for honest ratings.
- We discuss how reward mechanisms can be applied to the creation of structured knowledge.
- Based on our research questions, we formulate hypotheses and design experiments to test them.
- We extensively test the reward mechanism w.r.t. the creation of structured knowledge in different settings and with participants of different backgrounds.

A main finding of ours is that an HRM leads to ratings of equal or higher quality, compared to static rewards. Further, we find that rating-dependent rewards for contributions do improve contribution quality, but result in fewer contributions compared to a fixed reward per contribution.

Paper outline: Section 2 reviews related work. Section 3 presents our collaboration platform. We discuss the HRM in Sect. 4. Section 5 states our hypotheses. Section 6 discusses the experimental setup and sections 7 and 8 present the results of the experiments. Section 9 concludes.

¹ See <http://consensus.ipd.uka.de/techrep> for a demo.

2 Related Work

Various tools for the collaborative creation of structured knowledge have been proposed, ranging from full-fledged ontology editors with collaborative features [18, 19] over Wiki-based approaches for semantic data [3] to tools that support tagging folksonomies [12, 20]. There also exist commercial tools for the collaborative creation of structured knowledge, notably Freebase². Some of these tools feature rating mechanisms. However, we are not aware of any systematic attempts to evaluate the effect of ratings on knowledge quality for these tools.

[11] studies rating based incentive mechanisms for the collaborative construction of structured knowledge. The authors conduct two controlled experiments and find that ratings are a reliable measure of contribution quality. Further, they find that the presence of ratings increases the quality of contributions compared to a setting without. However, they do not test the effect of incentive mechanisms for ratings on rating quality and user behavior.

Eckert et al. use feedback from users of Amazon Mechanical Turk (AMT) to construct *is-a* relations between terms in a philosophy knowledge domain [8]. They compare the results to expert opinions and conclude that experts perform better. Since the AMT workers must understand the knowledge domain at hand, this approach is only applicable to domains that are commonly known, but not to domains that require specific knowledge.

Von Ahn uses games to motivate users to perform useful tasks [1]. For example, users in the ESP game are randomly paired to create tags for an image and receive points whenever their tags match. Siorpaes and Hepp [17] use this principle to build ontologies from Wikipedia entries by categorizing entries as either classes or instances. Users receive points when they agree on the categorization. Again, these approaches are better suited for knowledge domains that are commonly known and where data is already available, e.g., in form of Wikipedia entries. They cannot be used when the structured knowledge is created from scratch. Next, the game of assigning Wikipedia entries to predefined categories and rewarding users based on answers by other users is essentially an HRM scenario. Thus, an HRM could give way to better results with these games.

3 Creating Structured Knowledge with Consensus Builder

This section describes our tool Consensus Builder (CB) that allows for the collaborative creation of structured knowledge.

Data format. CB uses a data format similar to Topic Maps [15] and Entity-Relationship Models with some slight deviations for better usability. Data can be created on the type and on the instance level. I.e., users can specify the schema and create instance data. Topics represent entities of the real world, e.g.,

² <http://www.freebase.com>

The screenshot shows the 'Consensus Builder' interface for the topic 'Harrison Ford'. At the top, there is a navigation bar with the site name, user scores (Score: 15.77, E-Score: 8.07), and links for account, sign out, report a bug, and help. Below this is a search bar and a list of navigation links (home, topics, topic types, statistics, points, prizes, my contributions, to rate). The main content area is titled 'Topic' and features a large box for 'Harrison Ford' with a description and a profile picture. Below this, there are two sections: 'Person' and 'Actor'. The 'Person' section lists attributes such as 'Date of birth' (July 12, 1942) and 'Place of birth' (Chicago, Illinois, U.S.), each with a 'details' link and an 'edit' button. The 'Actor' section lists associations with films like 'Blade Runner', 'The Empire Strikes Back', and 'Raiders of the Lost Ark', each with a 'details' link and a 'delete' button. A 'Recently viewed' sidebar on the right lists 'Harrison Ford', 'The Empire Strikes Back', 'Film', 'Raiders of the Lost Ark', and 'Person'.

Fig. 1: Details for topic ‘Harrison Ford’ in Consensus Builder

Harrison Ford or **Indiana Jones**. Topics can have one or more types, e.g., **Harrison Ford** is of type **Person** and of type **Actor**. Types contain attributes and association types. Attributes describe simple data, like ‘date of birth’, and are constrained by data types, e.g., integer, string. Association types describe associations between topic types, e.g., **Actor** <acts in> **Movie**.

Collaborative Editing and Rating. Users can create and change all parts of the data model collaboratively. They can create single data elements like attributes or topics and change elements created by others. When a type is added to a topic, it inherits the attributes and association types of that type. Users can then set attribute values and add associations. CB takes care that users can create only attribute values and associations that are valid regarding the type level. In addition, CB provides various functions for browsing and searching, comments to discuss topics and topic types, and statistics such as user scores.

CB features a fine-grained association between ratings and contributions, called *rating scheme*, that lets users assess the quality of contributions on a very detailed level. Users can rate every element on the type as well as on the instance level that can be manipulated, e.g., topic names, attributes and association types, and attribute values.

Since our rating scheme is very fine-grained, we use a binary rating scale, i.e., ratings are either ‘low’ or ‘high’, in order not to overload the user. Cosley et al. show that, even though users prefer a finer-grained rating scale, the granularity of the rating scale does not have an adverse effect on user behavior [6]. The

functions for rating, editing, and displaying of the data are tightly integrated in the user interface (cf. Fig. 1).

Users can change and delete individual contributions. However, change operations are not trivial in any setting where data items depend on each other. For instance, what happens with other contributions and ratings associated with a contribution just deleted? Think of the deletion of a type that has associated topics and has received high ratings. To address these issues, we have made the following design decision. A contribution can only be changed if it satisfies two conditions. First, it must not be associated with other contributions, e.g., an attribute can only be changed if there are no attribute values associated with it. Second, it either must not have received any ratings, or its average rating value must be below a certain threshold. Consequently, only contributions deemed low quality can be changed. The user who has changed the contribution becomes its new owner.

Commission: Rating-Based Remuneration. To motivate users to create and maintain the data we reward data operations with points based on ratings by other users. We refer to the rating-based remunerations as *commission*. A user obtains a commission every time another user issues a positive rating for a contribution that first user is the owner of. The value of the commission is computed as $\kappa \cdot c(\textit{operation})$, where κ is a constant factor depending on the scenario, and $c(\textit{operation})$ depends on the operation. For instance, $c(\textit{create topic}) = 3.0$, and $c(\textit{set attribute value}) = 2.0$.

Other Data Formats. The objective of this paper is the design and deployment of incentive mechanisms for the collaborative construction of structured knowledge. In this specific context, the data format to encode the knowledge is of secondary concern. We have mainly chosen the data format described above because of its ease of use for non-expert users. The functionality of CB is applicable to other formats for structured knowledge as well, such as those specific to the Semantic Web. Furthermore, the data format currently used can be mapped to OWL in a straightforward way: Topic types are mapped to OWL classes, attributes to data-type properties, association types to object properties (as well as to their inverse properties) with domains and ranges restricted to the respective classes or data types, topics are mapped to instances, attribute values to literals, etc.

4 Honest Rating Mechanism

We want to elicit high-quality ratings, as opposed to ratings that are uninformed or simply copy the majority opinion. (This does not exclude the majority opinion from being correct.) However, a simple reward, e.g., one point per rating, does not suffice. It does not provide an incentive for the rater to gather information before issuing her rating and to respond truthfully.

To address these challenges, we use an incentive mechanism that rewards truthful ratings [13]. The mechanism has been designed for online product ratings

where buyers perceive a noisy signal about the quality of a product and rate it according to their perception. The mechanism collects ratings from the buyers and computes a score for each rating. The quality of the product is fixed and defines the type of the product. The mechanism assumes a finite number of types indexed by $t = 1, \dots, T$. Let S^i denote the noisy signal perceived by rater i about the quality of the product, and let $S = \{s_1, \dots, s_M\}$ be the set of possible signals. Conditional on the product's type, the signals of the raters are independent and identically distributed. Let $Pr(s_m|t) = Pr(S^i = s_m|t)$ be the probability that a buyer receives signal s_m when the true type of the product is t . The mechanism assumes $Pr(\cdot|\cdot)$ to be common knowledge, and $\sum_{j=1}^M Pr(s_j|t) = 1$.

After all buyers have received the signals, the mechanism asks them to submit ratings according to their signals simultaneously. Let $r^i = (r^i(1), \dots, r^i(M))$ denote the rating strategy of rater i . That is, i announces $r^i(j)$ whenever she observes signal s_j . The honest reporting strategy is \bar{r} with $\bar{r}(j) = s_j$ for all $j \in \{1, \dots, M\}$, i.e., the rater always reports the truth. Subsequently, the mechanism scores every rating submitted by comparing it to the rating of another rater, $ref(i)$, called the reference rater of i . Let $\tau(r^i(j), r^{ref(i)}(k))$ be the payment received by i when she announces $r^i(j)$ and $ref(i)$ announces $r^{ref(i)}(k)$.

The expected payment of rater i depends on her prior belief and on the signal s_j she has received:

$$\begin{aligned} V(r^i, r^{ref(i)}|s_j) &= E_{s_k \in S}(\tau(r^i(j), r^{ref(i)}(k))) \\ &= \sum_{k=1}^M Pr(S^{ref(i)} = s_k | S^i = s_j) (\tau(s^i(j), s^{ref(i)}(k))) \end{aligned}$$

The conditional distribution that $ref(i)$ receives the signal s_k can be computed as $Pr(s_k|s_j) = \sum_{t=1}^T Pr(s_k|t) \cdot Pr(t|s_j)$. The posterior probability $Pr(t|s_j)$ of type t given the perception of signal s_j , can be computed using Bayes' Law.

The honest reporting strategy \bar{r} is a Nash equilibrium if and only if for all signals $s_j \in S$ and all reporting strategies $\hat{r} \neq \bar{r}$:

$$V(\bar{r}, \bar{r}|s_j) \geq V(\hat{r}, \bar{r}|s_j) \quad (1)$$

Miller et al. prove that payment schemes $\tau(\cdot, \cdot)$ for $V(\cdot, \cdot)$ which satisfy (1) exist. We use a linear program described in [13] to compute the payments.

The mechanism uses i 's rating to update the probability distribution that predicts $ref(i)$'s rating. The score reflects the quality of the reference rating relative to the updated distribution. If the rating of the reference rater is honest, a rater can maximize her expected payment by announcing her subjective beliefs.

In the following we describe the design decisions behind our implementation of the HRM: We have modelled signals to be in the set $\{l, h\}$ and types to be in the set $\{1, \dots, 9\}$. For each rateable contribution we maintain estimates of two probability distributions: the prior distribution of the types $Pr(t)$ and the signal distribution $Pr(s_j|t)$. We have modeled types and signal distributions with type i generating a distribution $Pr(h|i) = i/10$ as proposed in [14], e.g., type 3

generates h ratings with frequency 0.3. For each contribution we maintain one type distribution $Pr(t)$. We also maintain a global type distribution that we use to initialize $Pr(t)$ of newly created items. We update $Pr(t)$ with the new ratings submitted. If someone changes an item, we reset its rating history and initialize its prior distribution with the global distribution at the time of change.

We put subsequent ratings of a contribution into groups of small size (typically 3 or 4) and score ratings in a group against each other. To motivate the user further, we display the sum of the expected scores for her unscored ratings.

5 Hypotheses

We formulate three hypotheses. They refer to the quality of either ratings or contributions. We measure quality with respect to a gold standard. I.e., a high quality rating coincides with the gold standard, whereas a low quality rating or rating error does not.

H_{rate}: An HRM improves rating quality.

H_{comm}: Commissions improve contribution quality and reduce quantity, compared to fixed remunerations.

In the case of static rewards users are rewarded for every contribution, regardless of its quality. This is why we expect contributions to be fewer, but of higher value when they are rewarded contingent on ratings of other users.

H_{err}: Rating errors receive lower scores from the HRM.

6 Experiments

To test our hypotheses we have conducted three experiments with CB. In each experiment, we randomly assigned participants to the experimental group (EG) or the control group (CG). We use the EG to evaluate the effects of the mechanisms in questions. The CG serves as the baseline. In the following, we first present the individual experiments. We then describe characteristics of the experimental setup common to all experiments, the different gold standards we use for quality assessment, and the statistical methods we use in our analysis.

RATEONLY. In this experiment we focus exclusively on the HRM. We recruited participants among students of our chair and instructed them to rate 127 contributions. We had preselected these 127 contributions from a knowledge base which students had created for the domain “Karlsruhe Institute of Technology” in a creation phase prior to the experiment itself. The selected contributions remained embedded in the other contributions from the creation phase. But since we wanted to test the rating mechanism in isolation for this experiment, we disabled ratings for these other contributions, as well as data manipulations. This reduces effects not related to the HRM. For three days, participants rated the contributions using CB. To test **H_{rate}**, we scored participants in the EG with the HRM, while the CG was scored statically with one point per rating.

HONSTUDENTS. We tested H_{rate} in a setting with the full functionality of CB. We invited students of the lecture “Machine Design” of the department of mechanical engineering. We told them to create topics and types which represent the content of the lecture and to rate the contributions of others. Again, we rewarded the EG by means of the HRM and the CG with one point per rating. To allow for comparing ratings of CG and EG later on, both groups had to rate contributions from the same set. For this reason, participants of both groups worked together on the same data. If the groups had used separate data, it would have been hard to say whether differences in rating quality result from the rating mechanism or from differences in the nature of the contributions.

HONSTAFF. We repeated the experiment HONSTUDENTS to test H_{rate} in a setting close to that of a community within a company. For this run we invited researchers from the Institute of Product Engineering. As knowledge domain we used a model for the engineering design process developed by this institute [2]. We advised participants to use the elements of that engineering model as topic types and concrete instances as topics.

COMMISSION. We designed the next experiment to test H_{comm} . Testing it in the experiments just described would have been difficult. This is because testing H_{comm} requires the CG and the EG to operate on separate data in order to eliminate potential influences between the groups and to allow for an unambiguous quality assessment of the contributions of each group. Such undesired influences are likely to occur in shared knowledge bases because data entries depend on each other, e.g., the contributions on the instance level depend on the schema level. Furthermore, according to H_{comm} , we expect a higher number of low quality contributions in the CG. This might affect the results even more if both groups operated on shared data. COMMISSION took place in a real world setting as well. Participants were students of the lecture “Database Systems”. Again, we asked participants to model the content of the lecture and related information on the schema and instance level and to rate contributions of each other. To test H_{comm} , we choose *usage of commissions* as independent variable: The CG received a fixed amount per operation depending on the operation category, as specified by $c(\text{operation})$. To prevent potential exploitation, this amount was deducted when the contribution was deleted. The EG received commissions contingent on the operation category and on ratings of other participants, i.e., the current owner received $\kappa \cdot c(\text{operation})$ for every positive rating the contribution had received and 0 points for negative ratings. We set κ to 0.2. We rewarded ratings of both the CG and the EG by means of the HRM.

In each experiment described so far, at least one of the experimental groups used the HRM. This allows us to test H_{err} .

Experimental Setups – Discussion. Our experiments go well beyond vanilla laboratory experiments, in several respects: They take place in real-world settings, within online communities where participants are not restricted by laboratory conditions. Unlike toy domains, the knowledge domains used were complex and

had real-world significance. The participants used CB from home or from their workplace. We put attention to not letting them know that an experiment was taking place nor that there were different experimental groups, by announcing the experiments as “beta test and user study”. (We introduced experimental features by means of announcements within CB.) Further, the assignment to groups was transparent to the participants, i.e., there were no indicators (e.g., specific URLs, etc.) that made the group explicit. The experiments lasted up to several weeks. This has blurred the distinction between real world and experiment further. Finally, participants remained anonymous to each other throughout the experiment and had no information about how many members their respective communities had. In this respect, the settings do not differ from large online communities.

Further Characteristics of the Experiments. To allow for a comparison with the CG, which was rewarded 1 point per rating in the experiments that test H_{rate} , we scaled the expected score of the HRM to 1.5 points. (Assuming risk-averse participants, we use 1.5 points instead of 1.) To stimulate participation, we announced a fixed monetary compensation for all users who reached a certain number of points. In addition, we conducted a lottery at the end of each experiment. The number of lottery tickets depended on the points achieved during the experiment. Lottery prices varied over the experiments and included monetary payments as well as prices like usb sticks. For RATEONLY, we paid a fixed monetary compensation to the CG and a point-dependent amount to the EG. Prior to each experiment, participants could take part in a “training phase” to get accustomed with CB. The purpose of the training was to remove potential distortions of the experiments due to the learning curve. Next, participants could enter an experiment while it was already running. The domain of the email address constrained the registration to guard against sybil attacks, i.e., against single users creating multiple accounts to gain unfair advantages [7]. An algorithm based on biased coin randomization [9] assigned participants to either the CG or the EG, while keeping the numbers of members of the groups balanced.

After each experiment, we invited the participants to complete a questionnaire. It elicited feedback on rewards, ratings, rating mechanisms, the behavior of other participants, and the usability of CB. The number of questions per questionnaire ranged up to 24, dependent on the configuration for the respective participant. We used a five point Likert scale response format (‘strongly disagree’ to ‘strongly agree’) for most questions.

Gold Standards for Quality Assessment. Assessing the quality of contributions and ratings required several different gold standards. For HONSTUDENTS, HONSTAFF, and COMMISSION we let domain experts rate a subset of contributions and used their ratings as gold standard. The subset depended on the hypothesis to test. Testing H_{rate} required comparing the rating quality between the CG and the EG. We randomly picked 150 contributions that had received at least one rating from both experimental groups. To test H_{comm} we simply picked 50 contributions randomly from each group. The experts used the same CB user

interface as the participants to issue ratings. To understand the context, experts could see all contributions created in the respective experiment. For COMMISSION, in addition to the detail ratings for randomly selected contributions such as attributes, associations etc., we let the experts assess the ‘overall quality’ and ‘overall adequateness’ of the topic or topic type associated with the respective contribution as a whole. This allows for a comprehensive quality assessment of the contributions. For the overall ratings we used a five-star rating scale instead of a binary one. We chose the following domain experts for the various experiments: the teaching assistant of the respective course for COMMISSION and for HONSTUDENTS, and a scientist whose research topic is the engineering model that served as the domain for HONSTAFF. Since both domain experts had limited experience in data modeling for the latter two experiments, a database expert supported them with the data modeling.

For RATEONLY we selected 127 contributions manually, out of the more than 5000 contributions created during the data-creation phase. The contributions selected were unambiguously correct, as confirmed by information publicly available on websites. We manipulated 34 out of the 127 contributions (all on instance level) so that they were false. The manipulated contributions together with the remaining manually selected ones served as the gold standard.

We classified the manipulations in three categories according to the effort needed to verify the respective errors:

1. *Easy to verify.* These are blatant errors, like a building having 666 elevators or a paper on sensor networks published in 1920.
2. *Medium effort to verify.* This category contained plausible-looking errors, like changes in room numbers or changes in co-authors of a paper. They could be detected by internet search.
3. *Hard to verify.* These manipulations were subtle and could only be verified with high effort, for example, the number of floors in a remote building.

We expect that the HRM has an effect on errors of Category 2 only. Both groups should recognize errors in Category 1. Category 1 allows checking whether participants made any effort at all. For Category 3, the effort for error detection exceeds the benefit from honest ratings by much. It serves as an extra check to exclude the possibility that the EG was more motivated than the CG a priori.

Table 1 shows an overview of the different setups.

Statistical Methods. We use Pearson’s χ^2 test to evaluate associations between binary variables, e.g., between classification errors and usage of the HRM. (For directional associations we use the one-tailed χ^2 test [10].) We use Student’s *t*-test to compare the means of the five-star ratings for overall quality. We use Pearson’s correlation to test the point-biserial correlation between binary rating errors and rating scores. Finally, we use Spearman’s ρ to test the correlation between Likert responses from the questionnaire and experiment results.

	RATEONLY	HONSTUDENTS	HONSTAFF	COMMISSION
Static Ratings	CG	CG	CG	-
Honest Ratings	EG	EG	EG	both
Static Contrib.	-	-	-	CG
Commission	all	all	all	EG
Duration	3 days	3 weeks	2 weeks	3 weeks
Shared Data CG/EG	yes	yes	yes	no
Gold Standard	Manipulation	Experts	Experts	Experts

Table 1: Overview of Experimental Setups

7 Results

We present the results of our experiments, including the evaluation of the hypotheses and of the questionnaire. Table 2 shows a list of quantitative characteristics of the four experiments.

	RATEONLY		HONSTUDENTS		HONSTAFF		COMMISSION	
	CG	EG	CG	EG	CG	EG	CG	EG
Participants	3	6	8	12	10	10	11	14
Contributions	127*	127*	151	1052	136	162	802	206
Ratings	381*	762*	943	456	419	555	180	151
Ratings per Contribution	3*	6*	0.78	0.38	1.4	1.86	0.22	0.73

Table 2: Overview of the Experiments (*Number of contributions and ratings fixed.)

7.1 H_{rate} : An HRM improves rating quality.

The absolute rating error for a rating $r_{ij} \in \{0, 1\}$ of participant j for contribution i and gold standard $g_i \in \{0, 1\}$ is $|r_{ij} - g_i|$. We average over all ratings for which we have a gold standard to calculate the mean absolute error (MAE).

For HONSTUDENTS there was a highly significant association between rating errors and the usage of the HRM ($\chi^2(1) = 71.52, p < 0.01$). The MAE was much higher for the CG (0.57) than for the EG (0.11). The odds ratio of making a rating error when using the HRM was 0.09. We conclude for this experiment that the mechanism improved rating quality.

For HONSTAFF we found no statistically significant association between rating errors and the usage of the HRM ($\chi^2(1) = 1.9071, p = 0.17$). The CG showed slightly better results regarding rating quality (MAE=0.16) than the EG (MAE=0.22). For HONSTAFF we conclude that there is no significant effect of the HRM on rating quality. A possible reason for this is that the researchers

already had a high intrinsic motivation to create high-quality data since they wanted to use it in their research later on. This high intrinsic motivation might have diminished the effects of the HRM.

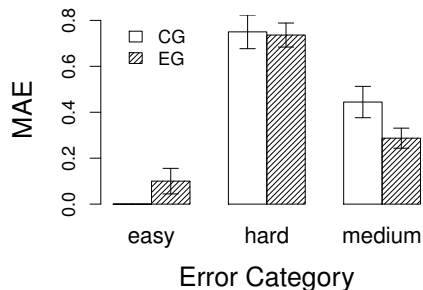


Fig. 2: Mean Absolute Error (MAE) by Error Category

Table 3: Pearson correlation between absolute error and score for ratings scored with the HRM

Experiment	Correlation	p
HONSTUDENTS	-0.18	0.38
HONSTAFF	-0.26	< 0.01
COMMISSION	-0.97	< 0.01
RATEONLY	-0.83	< 0.05

Figure 2 shows the MAE for the three error categories of RATEONLY for CG and EG respectively. The participants of the EG made significantly fewer errors in the “medium effort to verify” category (odds ratio = 0.5, $\chi^2(1) = 3.3$, one-tailed $p < 0.05$). For the other two error categories we found no significant association between errors and the usage of the HRM. The MAE for ratings of non-manipulated contributions was very low in both groups (CG: 0.054, EG: 0.043) and the association not statistically significant (two-tailed $\chi^2(1) = 0.27$, $p = 0.6$). We conclude that, for this experiment, the HRM increases rating quality.

Summing up, we find that two out of three experiments support H_{rate} .

7.2 H_{comm} : Commissions improve contribution quality and reduce quantity.

To test the “quality” part of H_{comm} we compare the five-star expert ratings for selected contributions rewarded with commission (EG) to those rewarded statically (CG). The quality ratings for the contributions rewarded statically were significantly lower ($mean = 0.76$, $se = 0.03$) than for contributions rewarded with commissions ($mean = 0.88$, $se = 0.03$, $t(85.7) = -2.6311$, $p < 0.01$). The ratings for adequateness of contributions rewarded statically were significantly lower as well ($mean = 0.77$, $se = 0.04$), compared to those rewarded with commissions ($mean = 0.98$, $se = 0.02$, $t(70.4) = -5.3688$, $p < 0.01$).

There were much fewer contributions per participant in the group using commissions ($mean = 27.0$, $se = 7.96$) than in the one without ($mean = 126.5$, $se = 67.27$). However, the difference was not statistically significant ($t(5.14) = 1.47$, $p = 0.10$). Interestingly we also found that participants remunerated with com-

missions seem to rate more critically. The ratio of negative ratings was significantly higher in the EG (0.258) than in the CG (0.039) ($\chi^2(1)=31.2$, $p < 0.01$), even though the experts rated the contributions of the EG more favorably.

7.3 H_{err} : Rating errors receive lower scores from the HRM.

To test H_{err} we calculate the correlation between the absolute rating error and the rating score by the HRM. A negative correlation means that the HRM scored rating errors lower than correct ratings. For RATEONLY we calculate the correlation per user, since every participant issued the same number of ratings. For the other experiments we calculate the correlation coefficients per rating.

Indirectly, H_{err} measures the average agreement of the raters with the gold standard. If the community, more precisely, the raters using the HRM, is sufficiently in agreement with the gold standard, punishment in the form of lower scores *should* follow. If, on the other hand, the community disagrees with the gold standard on average, rating errors should yield higher scores. Such a disagreement could happen for reasons of systematic differences in perception, e.g., due to different tastes or incompetence, or because of collusion attacks.

We find rather strong negative correlations for COMMISSION and for RATEONLY, and weak ones for HONSTUDENTS and for HONSTAFF (cf. Table 3).

Potential reasons for the two strong correlations could be that the manipulations used as gold standard in RATEONLY are more precise than expert ratings. For COMMISSION, we speculate that the strong correlation results from the better knowledge of participants regarding data-modeling techniques, and therefore a higher correlation with the expert ratings, compared to participants of HONSTUDENTS and HONSTAFF.

Overall, we conclude that the HRM punished rating errors to different degrees. The communities seem to have been in consensus with the experts, i.e., there were no systematic differences in perception.

7.4 Evaluation of the Questionnaire

27 out of the 46 users invited answered the questionnaire. Figure 3 shows the results for selected questions. We asked participants which rating strategy they used to maximize their rating score. Some stated an altruistic attitude “I did not intend to get as many points as possible, but tried to increase the quality of contributions by rating pointless or bad contributions as bad.”, “I tried to rate as much as possible as honestly as possible.” (both rewarded by the HRM). Others said they tried to maximize their scores, although with different rating strategies, dependent on their respective scoring mechanism, namely “Rating many items. But only those whose quality was easy to decide.” (HRM), and “Simply rated everything, no matter how.” (static reward for ratings).

Finally, we analyzed the correlations of experimental results of the participants and their questionnaire answers. Not surprisingly, we find a positive correlation between the understanding of the HRM and the number of ratings

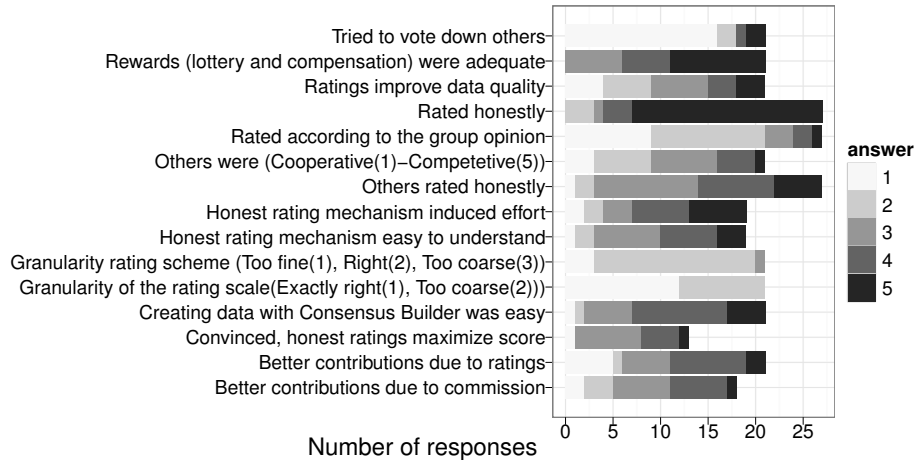


Fig. 3: Number of Responses for Selected Questionnaire Questions. Answers range from 1: ‘strongly disagree’ to 5: ‘strongly agree’, unless noted otherwise.

($\rho = 0.52, p < .05$). There is a strong correlation between the number of ratings and the stated strategy of rating the contributions of others badly in order to keep their scores low, but only for raters whose ratings were scored statically ($\rho = 0.764, p < 0.5$), while for participants using the HRM this correlation was negligible ($\rho = 0.16, p = 0.59$). We find a weak correlation between the understanding of the HRM and the score received per rating ($\rho = 0.31, p = .18$). In other words, it is not necessary to understand the mechanism in order to profit from it.

8 Discussion and Lessons Learned

Participants have been intrinsically motivated to some degree. They made good contributions and gave high-quality ratings even when they did not receive an extra reward for it. This is pronounced for the close-knit group of staff members. However, one of our results is that contributions and ratings are at least as good or better in the presence of commissions and the HRM, respectively. We speculate that, at least to some degree, the intrinsic motivation resulted from the fixed monetary compensation for participation, insofar as participants felt they had to offer at least some effort. In fact, when planning the experiments, given a fixed total budget, we were confronted with a tradeoff between two quantities: on the one hand, the guaranteed compensation, on the other, the score-dependent compensation. A low guaranteed compensation results in fewer participants. A high guaranteed compensation provides less incentive from rating dependent rewards.

Despite the rewards offered, ratings are sparse in both CG and EG in the experiments with a variable number of ratings, i.e., there are not nearly as many

ratings per contribution as there are participants per group (cf. Table 2). Questionnaire responses offer some explanation for this: Some participants do not like to rate data items they do not understand well enough even if they receive a guaranteed score. One participant of RATEONLY dropped out after the creation phase because she did not feel sufficiently familiar with the knowledge domain. Another reason might be the guaranteed compensation, as discussed above.

The results show a weak correlation between rating scores and understanding of the mechanism. An interesting question is whether a fake HRM would have the same effects as the real one. The authors of the HRM claim, but do not test, that it is not necessary for users to understand the HRM [14]. Instead it suffices if users trust the mechanism and believe that truthful answers maximize their scores in the long run. We speculate that telling participants that an HRM is used while scoring with some fake mechanism (for example, randomly) would still yield comparable results, at least in the short run. This could be tested experimentally by comparing the alternatives ‘no mechanism’, ‘real HRM’, and ‘fake HRM’. Note that, even if a fake mechanism yielded results similar to those of the real mechanism, the real mechanism would still be at least as good (or better in case participants realized the fake).

Next, it turned out to be difficult to find domains with all of the following characteristics: (a) They are sufficiently controversial to generate variance in the ratings. (b) The experimenters, but not the participants, have access to the gold standard. For example, in the creation phase before the RATEONLY experiment, participants kept the schema extremely simple and almost exclusively copied data publicly available on the web. Since the contributions were almost completely correct, there were no negative ratings and hence no variance in the rating values. This means that we could not have measured the effects of our mechanisms on either contribution or rating quality of the creation phase meaningfully.

Finally, the quality of the schema created by participants not familiar with data modeling was surprisingly good. Despite some beginner mistakes (confusion of normal associations and type associations, creation of topic type ‘Properties’) the quality of the schema level was good and detailed.

9 Summary

In this paper we have investigated how rating-based reward mechanisms can improve the quality of structured knowledge created collaboratively. In particular, we have discussed how mechanisms for honest ratings known from literature can be applied to this scenario. We have presented a community platform that features such reward mechanisms. We have formulated hypotheses and designed experiments to evaluate the effects of reward mechanisms for the collaborative creation of structured knowledge. We have carried out the experiments with different online communities. The communities constructed complex knowledge domains in settings close to real-world scenarios. We find that an honest rating mechanism improves the quality of ratings in two out of three experiments and

that it “punishes” rating errors with lower scores. Further we find that rewards for good contributions increase the quality of the contributions.

References

1. von Ahn, L.: Games with a purpose. *IEEE Computer* 39(6) (2006)
2. Albers, A., et al.: A new perspective on product engineering overcoming sequential process models. In: *The Future of Design Methodology*. Springer London (2011)
3. Auer, S., Dietzold, S.: OntoWiki - A Tool for Social, Semantic Collaboration. In: *Proceedings of the 5th International Semantic Web Conference ISWC* (2006)
4. Beenen, G., et al.: Using social psychology to motivate contributions to online communities. In: *Proceedings of the 2004 ACM conference on Computer supported cooperative work - CSCW '04* (2004)
5. Benkler, Y.: The wealth of networks: How social production transforms markets and freedom. *Information Economics and Policy* 19(2) (2007)
6. Cosley, D., et al.: Is seeing believing?: how recommender system interfaces affect users' opinions. In: *Proceedings of the SIGCHI conference on Human factors in computing systems* (2003)
7. Douceur, J.R.: The sybil attack. In: *IPTPS* (2002)
8. Eckert, K., et al.: Crowdsourcing the assembly of concept hierarchies. In: *Proceedings of the 10th annual joint conference on Digital libraries - JCDL '10* (2010)
9. Efron, B.: Forcing a sequential experiment to be balanced. *Biometrika* 58(3) (Dec 1971)
10. Fleiss, J.L.: *Statistical methods for rates and proportions*. Wiley, New York, 2nd ed. edn. (1981)
11. Hütter, C., Kühne, C., Böhm, K.: Peer production of structured knowledge - an empirical study of ratings and incentive mechanisms. In: *Proceeding of the 17th ACM conference on Information and knowledge mining - CIKM '08* (2008)
12. Jäschke, R., et al.: Organizing publications and bookmarks in bibsonomy. In: *CKC* (2007)
13. Jurca, R., Faltings, B.: Minimum payments that reward honest reputation feedback. In: *Proceedings of the 7th ACM conference on Electronic commerce - EC '06* (2006)
14. Miller, N., Resnick, P., Zeckhauser, R.: Eliciting Informative Feedback: The Peer-Prediction Method. *Management Science* 51(9) (Sep 2005)
15. Pepper, S., Architect, S.I., Infotek, S.: *The TAO of Topic Maps - Finding the Way in the Age of Infoglut* (2000)
16. Prelec, D.: A Bayesian truth serum for subjective data. *Science* (New York, N.Y.) 306(5695) (Oct 2004)
17. Siorpaes, K., Hepp, M.: Games with a Purpose for the Semantic Web. *IEEE Intelligent Systems* 23(3) (May 2008)
18. Sunagawa, E., Kozaki, K., Kitamura, Y., Mizoguchi, R.: An Environment for Distributed Ontology Development Based on Dependency Management. In: *International Semantic Web Conference*. pp. 453–468 (2003)
19. Tudorache, T., Noy, N.F.: Collaborative protege. In: *CKC* (2007)
20. Zacharias, V., Braun, S.: Soboleo – social bookmarking and lightweight engineering of ontologies. In: *CKC* (2007)