

## USE AND OPTIMISATION OF PAID CROWDSOURCING FOR THE COLLECTION OF GEODATA

V. Walter, D. Laupheimer, D. Fritsch

Institute for Photogrammetry, University of Stuttgart, 70174 Stuttgart, Germany, Volker.Walter@ifp.uni-stuttgart.de

Commission IV, WG V

**KEY WORDS:** Crowdsourcing, Data Collection, Geodata

### ABSTRACT:

Crowdsourcing is a new technology and a new business model that will change the way in which we work in many fields in the future. Employers divide and source out their work to a huge number of anonymous workers on the Internet. The division and outsourcing is not a trivial process but requires the definition of complete new workflows – from the definition of subtasks, to the execution and quality control. A popular crowdsourcing project in the field of collection of geodata is OpenStreetMap, which is based on the work of unpaid volunteers. Crowdsourcing projects that are based on the work of unpaid volunteers need an active community, whose members are convinced about the importance of the project and who have fun to collaborate. This can only be realized for some tasks. In the field of geodata collection many other tasks exist, which can in principle be solved with crowdsourcing, but where it is difficult to find a sufficient large number of volunteers. Other incentives must be provided in these cases, which can be monetary payments.

### 1. INTRODUCTION

The majority of projects in the field of crowd-based geodata collection are based on the work of volunteers. In contrast to these kinds of approaches, the project described in this paper has been realized with paid crowdsourcing, which means that the crowdworkers are paid for their work. We want to identify and quantify the parameters (e.g. amount of salary, size of working tiles or object types) which influence the quality of the results (especially correctness, geometric accuracy and collection time).

We developed a web-based program for the collection of geodata and integrated it into the commercial crowdsourcing platform microWorkers ([www.microWorkers.com](http://www.microWorkers.com)), which takes over the recruitment and the payment. The platform has access to more than 500,000 registered crowdworkers. The workers are informed automatically when a new job is offered by an employer on the platform. The employers can restrict the jobs to specific groups of workers. For example, it is possible to offer the jobs only to workers that are living in a specific country or to workers that have already successfully worked on a particular number of other jobs. Further qualifications are possible with their own developed tests, which must be solved before the job.

After the job has been completed, the results are submitted to the employer who checks the quality of the results. The final payment is handled by the crowdsourcing platform.

### 2. CROWDSOURCING

The term „Crowdsourcing“ was created by Jeff Howe (Howe, 2006) and is a neologism of the words “Crowd” and “Outsourcing”. In opposite to outsourcing, where tasks are outsourced to known companies, crowdsourcing tasks are outsourced to unknown users in the Internet (The crowd).

This enables the access to huge amounts of workers. The literature and practice are illustrating that nearly every value adding activity can be affected by crowdsourcing (Leimeister & Zogaj, 2013).

In order to realize successful crowdsourcing applications, it is necessary to recruit a sufficient amount of crowdworkers (Holland und Hoffmann, 2013). In many cases the crowdworkers are paid for their work. However, the payment is often only in the range of several cents. Because of this, crowdworkers are often named “Microworker” and the tasks, which have to be solved, “Microtasks” (Hoßfeld, Hirth & Tran-Gia, 2012). Crowdsourcing projects without payments often use other incentives, such as game points or trust levels.

Most of all existing crowdsourcing projects in the field of geodata collection are realized without monetary payments (VGI – Volunteered Geographic Information, such as OpenStreetMap). Approaches of this kind are not subject of the research described in this paper. A comprehensive discussion of this field can be found for example in (Sui, Elwood & Goodchild, 2013).

The recruitment and payment of paid crowdworkers is typically handled by special crowdsourcing platforms. Established platforms have access to a huge number of registered crowdworkers. The realization of paid crowdsourcing projects without the use of special crowdsourcing platforms is in principle also possible, but the recruitment of the crowdworkers would be especially difficult.

One of the first and largest platforms is „Amazon Mechanical Turk“ (MTurk, [www.mturk.com](http://www.mturk.com)), which has access to more than 500,000 registered crowdworkers (The Nation, 2014). Since the platform is restricted to employers that have an address in the US, it was not possible to use MTurk in our project. We use microWorkers, which has meanwhile also access to more than 500,000 crowdworkers ([www.microWorkers.com](http://www.microWorkers.com)).

Crowdsourcing in general and especially paid crowdsourcing are very young technologies. Because of this, many scientific problems are currently solved only partially or even not at all. In the following, the main important scientific challenges of crowdsourcing are discussed briefly.

In order to realize successful crowdsourcing projects, it is important to find appropriate crowdworkers that execute the tasks quickly and with high quality. Therefore, methods for the automatic selection of crowdworkers, which deliver good results, are needed. In the worst case, the crowdworkers are working on tasks that are too difficult for them or which they find uninteresting or for which they have not enough motivation because of too low salary (Kittur et al., 2013). Recommendation systems can be used to offer crowdworkers jobs for which they have special knowledge and which they find interesting (Tran-Gia et al., 2013). User profiles can be used, which are generated from the individual input of the crowdworkers or from an evaluation of already executed tasks.

If the tasks are difficult to solve, it is possible to use training phases or tests to ensure that the crowdworkers have the required qualification. It is important to integrate functionalities that enable the crowdworkers to communicate with the employer and other crowdworkers. This can be realized for example with discussion forums or feedback mechanisms (Brabham, 2008).

An important factor to control the quality of paid crowdsourcing is the amount of salary. Additionally, other non-monetary incentives can be used, such as user rankings, special rights (for example, the right to control other crowdworkers) or awards.

Quality control is another huge challenge. The quality of crowdsourcing tasks can vary considerably. Whereas some tasks are solved with very high quality, 30 percent or more of the tasks are solved incorrectly (Bernstein et al., 2010). The reasons for this are manifold and range from inaccurate job descriptions (job design) and non-intuitive graphical user interfaces to lacking qualifications of the crowdworkers.

The realization of automatic quality control mechanisms is difficult since the typical tasks of crowdsourcing projects are such tasks that can be solved only very hard by computers. A possible solution is to source out the quality control also to the crowd. In principle, three approaches are possible (Leimeister & Zogaj, 2013): (1) control of the results by other crowdworkers, (2) mixing the tasks with additional tests (3) iterative approaches where two or more crowdworkers solve the same tasks, which are compared afterwards. A further method to ensure high quality is the careful design of workflows and graphical user interfaces. Studies show that with an appropriate process design, crowdsourcing tasks are solved better and faster (Leimeister & Zogaj, 2013).

### 3. GRAPHICAL USER INTERFACE

Figure 1 shows the graphical user interface of the program. The program was developed with JavaScript. The central element is an orthophoto with the size of 500 \* 450m<sup>2</sup> from which the data has to be collected.

The object classes that have to be collected are: forests (with polygons), streets (with lines) and buildings (with points). Incorrectly collected objects can be deleted. The workers have the possibility to make additional comments on their work and submit them together with the collected data.

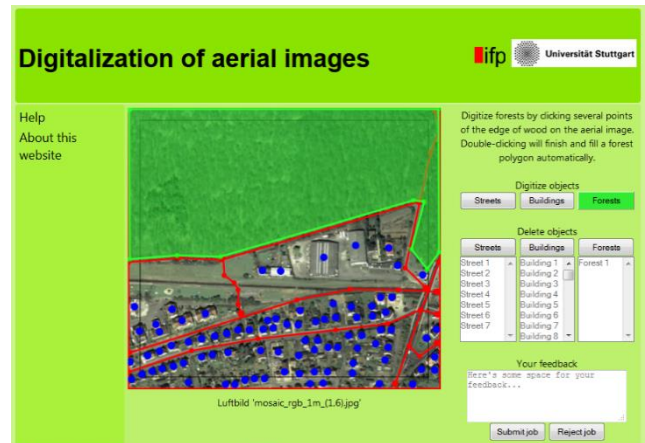


Figure 1. Graphical user interface for the crowd-based collection of spatial data.

### 4. TEST DATA

A RGB orthophoto (see Figure 2) with a ground sampling distance (GSD) of 1m and a size of approximately 5 \* 4 km<sup>2</sup> was subdivided into 88 patches with the size of 500 \* 450 m<sup>2</sup>. Six campaigns with different parameters were launched on the microWorkers platform to evaluate the quality of the crowd-based data collection.



Figure 2. Test area.

### 5. ORIGIN OF THE CROWDWORKERS

Figure 3 shows the spatial distribution of the origin of the crowdworkers of campaign 1 and 2. The results of the other campaigns are similar.

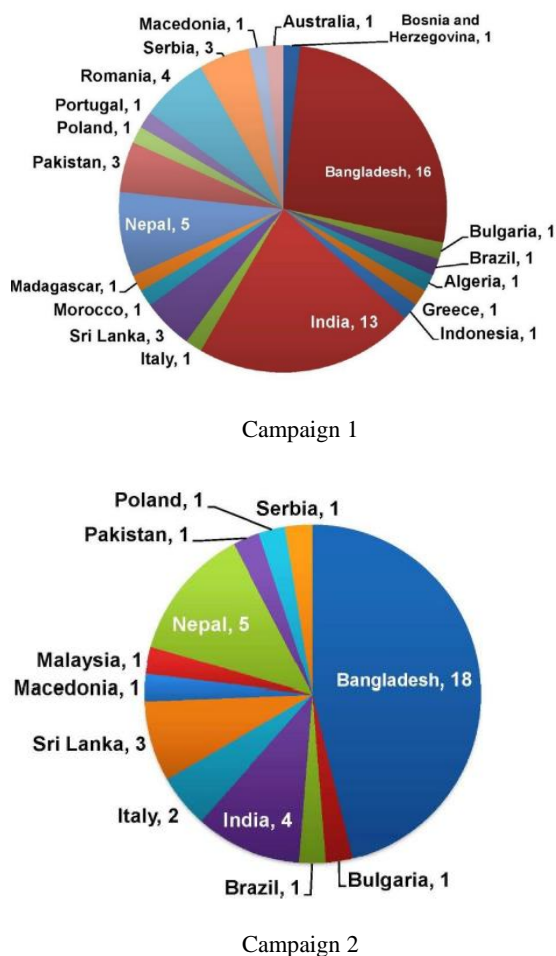


Figure 3. Origin of the crowdworkers

The majority of all crowdworkers live in Asia or Eastern Europe. Typically, Bangladesh, India and Nepal are the top three countries from which the crowdworkers come from. All other parts of the world are only rarely involved. Since the payment of the crowd-based work is typically very low, it is clear that the typical crowdworkers come from countries with low wages.

## 6. DATA EXAMPLES

The results are very promising and the quality of the data was even outperforming our expectations in many cases. However, the tests show also that the quality of the collected data vary significantly. Some crowdworkers collected the data with very high quality whereas other crowdworkers collected completely incorrect data.

All results were subdivided into five categories:

- *Category 1* (very good): Most of all data is correct collected with high accuracy (see example in Figure 5a: Very few houses are missing, but all other objects are collected completely and with high accuracy).

- *Category 2* (good): Some objects are not collected or not correct collected, but overall the work has high accuracy (see example in Figure 5b: Some houses and streets are missing in the upper left area, but overall the work has high accuracy).
- *Category 3* (partly good): A significant amount of objects is missing or is not correct collected, but the remaining work is accomplished with good accuracy (see example in Figure 5c: The lower right part of the image and the forest at the bottom of the image is not digitized, but the remaining work is accomplished with good accuracy).
- *Category 4* (poor): Many objects are missing and the rest of the objects are collected with low accuracy (see example in Figure 5d: Most of the streets are ok but many houses are incorrect or not digitized at all. Also, the geometry of the forest is wrong).
- *Category 5* (unsatisfactory): The major area is not digitized (see example in Figure 5e: Only some house are collected with low accuracy) or the collected data is completely senseless (see example in Figure 5f).

## 7. QUALITY OF THE COLLECTED DATA

Figure 4 shows the distribution of the quality of the collected data of campaign 1 and 2. The results of the other campaigns are similar. It can be seen that most of the patches were collected with very good or good quality (category 1 and 2). Only a smaller amount of patches was collected with lower quality.

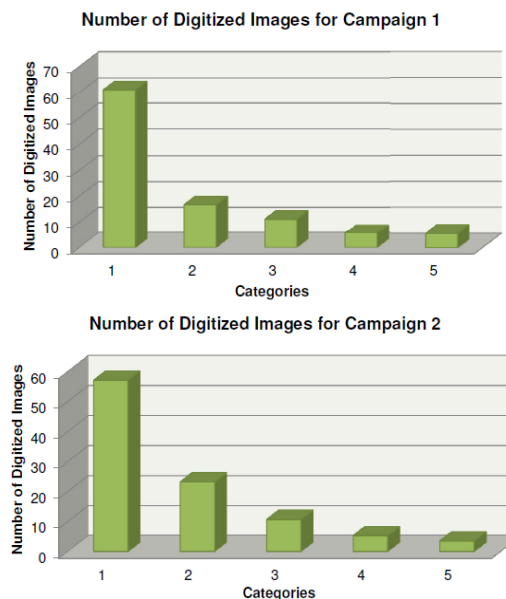


Figure 4. Quality of the collected data.

The difference between the two campaigns is that the salary per patch was \$0.10 in campaign 1 and \$0.45 in campaign 2. We expected that an increased salary would also increase the quality of the results but we observed that there is no direct connection between the amount of salary and the quality of the results. An increase of the salaries did not lead to a better quality but only to a faster completion of the campaigns.



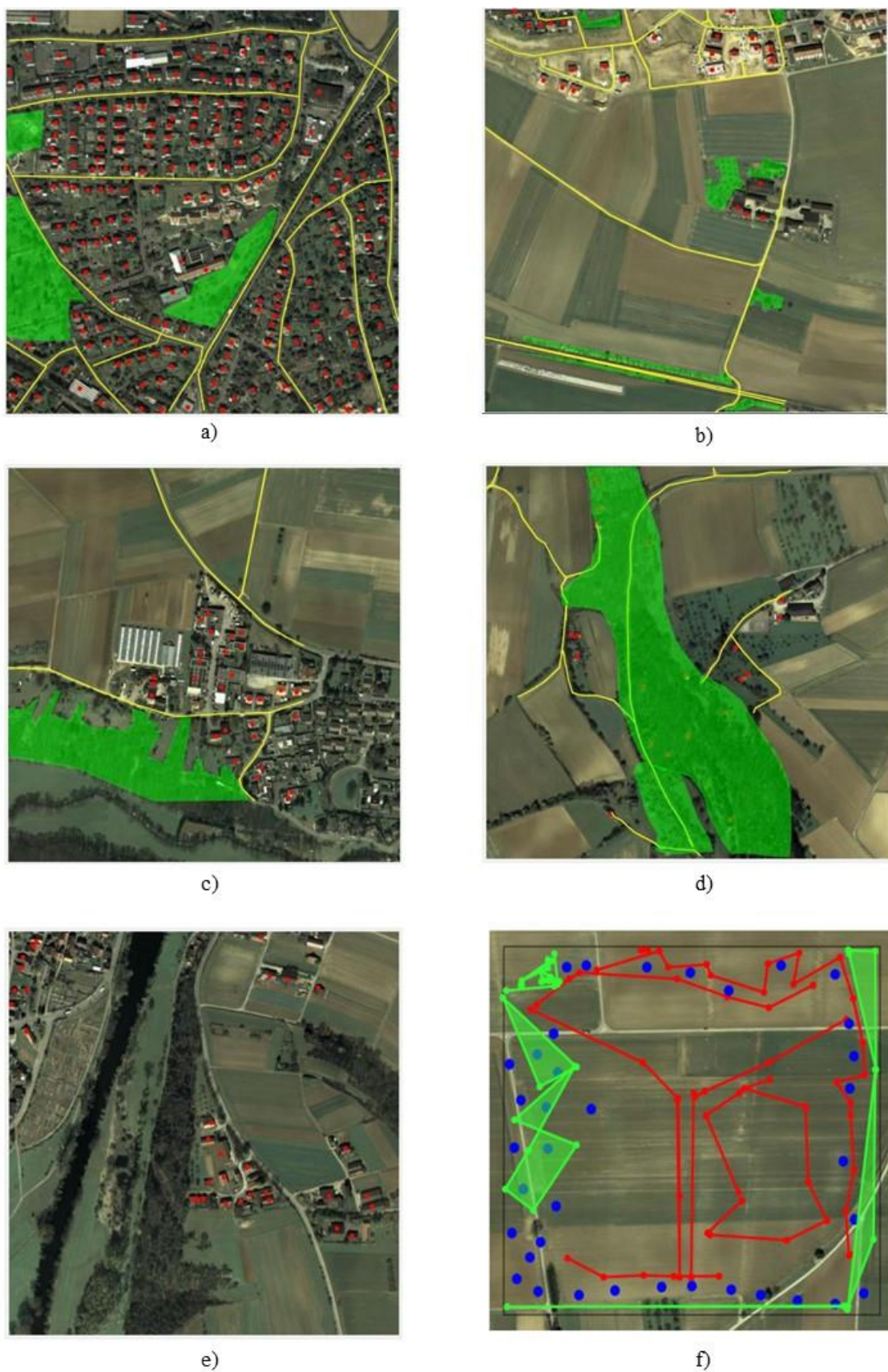


Figure 5. Examples of crowd-sourced data collection.

## 8. TIME TO ACCOMPLISH THE CAMPAIGNS

Table 1 shows the time which was needed to complete campaign 1 and 2. It can be seen that an increase of the salaries leads to a significantly quicker completion of the campaign. The reason for that is, that it is much more attractive for crowdworkers to work on jobs with higher salaries. Therefore, the salary is a good parameter to control the collection time.

Table 1. Time to accomplish campaign 1 and 2.

| Campaign | Payment Job in \$ | No. of Days |
|----------|-------------------|-------------|
| 1        | 0.10              | 12          |
| 2        | 0.45              | 3           |

## 9. SUMMARY AND OUTLOOK

The test showed that in principle it is possible to produce high quality spatial datasets with paid crowdsourcing. The main problem is that the quality of the data is extremely heterogeneous.

Therefore, it is necessary to find control mechanisms that evaluate the quality of the data. This either must be done automatically or again sourced out to the crowd.

Furthermore, selection procedures are needed, which can automatically select crowdworkers who collect data with high quality. This can be realized for example with user profiles.

Finally, algorithms are needed to integrate the individual results into an overall result. Spatially inconsistent datasets, whose overlap multiplies, have to be integrated into a consistent, uniform dataset. All these aspects will be investigated in our ongoing research.

## REFERENCES

Bernstein, M. S., Little, G., Miller, R. C., Hartmann, B., Ackerman, M. S., Karger, D. R., Panovich, K. (2010). Soylent: a word processor with a crowd inside. In: Proceedings of the 23rd annual ACM symposium on User interface software and technology, 313-322.

Brabham, D.C. (2008). Moving the Crowd at iStockphoto: The Composition of the Crowd and Motivations for Participation in a Crowdsourcing Application. *First Monday* 13(6).

Hoßfeld, T., Hirth, M., Tran-Gia, P. (2012). Crowdsourcing. *Informatik-Spektrum*, 35(3), 204-208.

Howe, J. (2006). The Rise of Crowdsourcing. *WIRED Magazin*, <http://archive.wired.com/wired/archive/14.06/crowds.html>.

Kittur, A., Nickerson, J. V., Bernstein, M., Gerber, E., Shaw, A., Zimmerman, J., Lease, M., Horton, J. (2013). The future of crowd work. In: Proceedings of the 2013 ACM conference on Computer supported cooperative work, 1301-1318.

Leimeister, J., M., Zogaj, S. (2013). Neue Arbeitsorganisation durch Crowdsourcing. Arbeitspapier Nr. 287, Hans Böckler Stiftung, Düsseldorf.

Sui, D., Elwood, S., Goodchild, M. (2013). *Crowdsourcing Geographic Knowledge, Volunteered Geographic Information (VGI) in Theory and Practice*, Springer Verlag, New York.

The Nation (2014). How Crowdworkers became the Ghosts in the Digital Machine. *The Nation*, available online at: <http://www.thenation.com/article/178241/how-crowdworkers-became-ghosts-digital-machine#>.

Tran-Gia, P., Hoßfeld, T., Hartmann, M., Hirth, M. (2013). Crowdsourcing and its Impact on Future Internet Usage. *it-Information Technology Methoden und innovative Anwendungen der Informatik und Informationstechnik*, 55(4), 139-145.