



# **Context-Aware Worker Selection for Efficient Quality Control in Crowdsourcing**

*Summary*

*Tarekt Awaad*

In the last decade, crowdsourcing has proved its ability to address large scale data collection tasks, such as labeling large data sets, at a low cost and in a short time. However, the performance and behavior variability between workers as well as the variability in task designs and contents, induce an unevenness in the quality of the produced contributions and, thus, in the final output quality. In order to maintain the effectiveness of crowdsourcing, it is crucial to control the quality of the contributions. Furthermore, maintaining the efficiency of crowdsourcing requires the time and cost overhead related to the quality control to be at its lowest. While effective, current quality control techniques such as contribution aggregation, worker selection, context-specific reputation systems, and multi-step workflows, suffer from fairly high time and budget overheads and from their dependency on prior knowledge about individual workers.

In this thesis, we address this challenge by leveraging the similarity between completed and incoming tasks as well as the correlation between the worker declarative profiles and their performance in previous tasks in order to perform an efficient task-aware worker selection. To this end, we propose CAWS (Context Aware Worker Selection) method which operates in two phases; in an offline phase, completed tasks are clustered into homogeneous groups for each of which the correlation with the workers declarative profile is learned. Then, in the online phase, incoming tasks are matched to one of the existing clusters and the correspondent, previously inferred profile model is used to select the most reliable online workers for the given task. Using declarative profiles helps eliminate any probing process, which reduces the time and the budget while maintaining the crowdsourcing quality. Furthermore, the set of completed tasks, when compared to a probing task split, provides a larger corpus from which a more

precise profile model can be learned. This translates to a better selection quality, especially for harder tasks.

In order to evaluate CAWS, we introduce CrowdED (Crowdsourcing Evaluation Dataset), a rich dataset to evaluate quality control methods and quality-driven task vectorization and clustering. The generation of CrowdED relies on a constrained sampling approach that allows to produce a task corpus which respects both, the budget and type constraints. Beside helping in evaluating CAWS, and through its generality and richness, CrowdED helps in plugging the benchmarking gap present in the crowdsourcing quality control community.

Using CrowdED, we evaluate the performance of CAWS in terms of the quality of the worker selection and in terms of the achieved time and budget reduction. Results shows the following: first, automatic grouping is able to achieve a learning quality similar to job-based grouping. And second, CAWS is able to outperform the state-of-the-art profile-based worker selection when it comes to quality. This is especially true when strong budget and time constraints are present on the requester side.

Finally, we complement our work by a software contribution consisting of an open source framework called CREX (CReate Enrich eXtend). CREX allows the creation, the extension and the enrichment of crowdsourcing datasets. It provides the tools to vectorize, cluster and sample a task corpus to produce constrained task sets and to automatically generate custom crowdsourcing campaign sites.