

Harnessing Crowds and Experts for Semantic Annotation of the Qur'an

Amna Basharat, Khaled Rasheed, I. Budak Arpinar

Department of Computer Science
University of Georgia
Athens, GA, 30602 USA
amnabash@uga.edu, khaled@uga.edu, budak@uga.edu

Abstract. In this paper we illustrate how we harness the power of crowds and specialized experts through automated knowledge acquisition workflows for semantic annotation in specialized and knowledge intensive domains. We undertake the special case of the Arabic script of the Qur'an, a widely studied manuscript, and apply a hybrid methodology of traditional 'crowdsourcing' augmented with 'expertsourcing' for semantically annotating its verses. We demonstrate that our proposed hybrid method presents a promising approach for achieving reliable annotations in an efficient and scalable manner, especially in cases where a high level of accuracy is required in knowledge intense and sensitive domains.

Keywords: semantic annotation, disambiguation, classification, ontology, Qur'an

1 Introduction

Thematic annotation of religious texts, in particular, the classical sources of knowledge in the Islamic domain in the Arabic language, has not received much attention, partly owing to time and knowledge constraints from experts required for such an annotation process. In our research, we consider the application of specialized human computation methods such as nichesourcing ([1], [2]) in an attempt to scale this process of annotation. Nichesourcing or expertsourcing extends the idea of engaging skilled and knowledgeable persons in place of faceless crowds for human driven tasks. We employ nichesourcing as means of augmenting traditional crowdsourcing methods rather than as an alternate.

In this paper, we show the results of an exploratory study focussing on two knowledge intensive tasks: the thematic disambiguation and annotation of Qur'anic verses using its Arabic script. While this has been tackled through a pure crowdsourcing approach in our earlier work in [3], we determined that several tasks are rather knowledge intensive and require domain expertise. In this case, not only the knowledge of Qur'anic Arabic is considered imperative, the annotation of the Arabic verses also requires understanding the context and content of the given verse.

2 Hybrid Architecture for Harnessing Crowd and Expert Annotations

We design and develop a hybrid workflow architecture that connects a crowdsourcing framework with an expertsourcing application as shown in Figure 1.

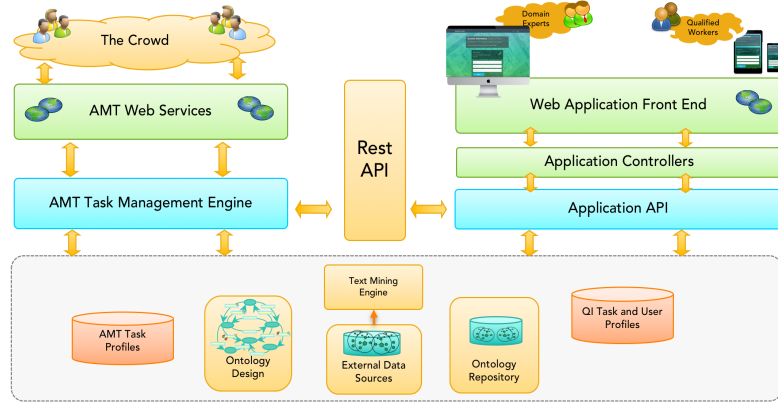


Fig. 1. Hybrid Architecture for Harnessing Crowd and Expert Annotations

Crowdsourcing Stage: We design a task management engine that is responsible for generating tasks, retrieving and aggregating results. The tasks are published on the Amazon Mechanical Turk (AMT)¹ platform. A complete workflow management system is implemented (a derivative of a workflow model for Linked Data Management presented in [4]), which includes means for generating dynamic tasks from a range of task profiles. The semantic annotation process is driven by an ontology schema. The task input is generated by retrieving relevant candidate verses from the available external data sources such as the Semantic-Quran [5] dataset.

The AMT crowd performs the *thematic disambiguation* and *thematic annotation* tasks. Both tasks are based on the Arabic script of the Qur'an. For the disambiguation task, a question is presented to the crowd, which includes a verse, along with a highlighted, candidate explicit assertion for the given theme, and the crowd responds by declaring this assertion as either a positive or negative by determining if the occurrence is a true occurrence of the given theme. The annotation tasks require deeper knowledge and understanding of the Arabic text. The crowd determines whether the given verse contains any implicit reference to the given theme. If their response is positive, then they are also required to provide the portion of the verse (a meaningful phrase or a word) that implies the presence of the theme. As a form of a quality measure, the crowd is also

¹ <http://www.mturk.com>

required to provide a confidence level (ranging from Very High to Very low), to indicate their confidence in their response.

Decision Analytics: We collect and aggregate the responses based on statistical measures of aggregation. Weighted confidence measures and thresholds are applied. Based on this aggregation, the completed tasks are marked as either *Approved* or *Reviewable*. A high confidence and aggregation threshold is applied for the approved tasks. This decision analytics results in identifying the candidate tasks for expertsourcing. The tasks marked as reviewable, which fail to meet the agreement thresholds, are sent off for expert annotations.

Expertsourcing Stage: For this purpose we designed a custom web application to engage with experts. A RestAPI connects the crowdsourcing task management engine with the expertsourcing application. The tasks are sent to the remote application and experts are notified when the tasks become available. The experts also see the candidate responses collected from the crowdsourcing stage. The experts have either the option to choose from the available annotations (collected during the crowdsourcing stage) or provide their own if they do not agree with either one. An example is shown in Fig 2. We present the same task to three experts to analyze the annotation agreements. The approved and validated annotations are passed on for *Ontology Population* and linked with existing data sources.

The screenshot shows a web application interface with the following elements:

- Title:** Identify Verses that contain attributes of Allah
- Instruction:** Help correctly identify those verses that contain attributes related to Allahs names.
- Question:** Does the following verse (No 2) from Surah Aal-Imraan (آل عمران) [Surah No: 3] contain an attribute/quality related to Allahs Name: **إِلَهُ**
- Verse:** ﴿۲﴾ اللَّهُ لَا إِلَهَ إِلَّا هُوَ الْحَيُّ الْقَيُّومُ
- Response Options:**
 - Yes
 - No
- Segment Identification:** Which of the following segments correctly identify the part of the verse that contains the attribute of Allah for the given name
 - لا إله إلا هو
 - الله لا إله إلا هو الحي القيوم
 - None of the Above
- Text Input:** Please provide the part of the verse that contains the attribute of Allah for the given name. You may copy/paste text from the verse above or elsewhere.
- Buttons:** Submit and Skip

Fig. 2. Task Design for Thematic Annotation of Qur'anic verses

3 Results and Discussion

The experimental setup assigned each task to 5 crowd workers. For the reviewable crowd tasks that were sent to experts, 3 experts were assigned to each reviewable task. Table 1 shows the results obtained.

Task	Crowd Tasks		Expert Tasks	
	Disambiguation	Annotation	Disambiguation	Annotation
Approved	1267	477	34	96
Reviewable	40	107	6	11
Total	1307	584	40	107

Table 1. Results for Disambiguation and Annotation Tasks

The results of our exploratory study provide interesting insights into the application of human computation methods to knowledge intensive tasks. Our task design involved the thematic disambiguation and annotation of the Qur’anic verses based on the original Arabic script. For the disambiguation task, 99% of the tasks were able to reach an agreement by combining contributions of crowds and experts. Only about 3% tasks needed expert contributions. For the annotation tasks, about 18% tasks needed expert contributions. There were 10% tasks that did not reach an agreement with both crowd and expert contributions. An administrative review of these cases indicate that some annotations are a matter of personal taste and judgement and closed agreement is therefore difficult. Most of these annotations cannot be classified as wrong, nor better than the others based on an automated agreement mechanism.

Our knowledge acquisition and review workflow to selectively elicit expert annotations only where needed indeed presents a promising method. We utilize annotation agreement and distance analytics to route the appropriate tasks that need expert contributions. Our results suggest that such a hybrid approach indeed creates for a more accurate and reliable annotation process. This method can be effectively utilized for qualitative dataset management and semantic annotation tasks in an economic and feasible manner through crowd engagement, while reducing the need for expert contributions.

References

1. De Boer, V., Hildebrand, M., Aroyo, L., De Leenheer, P., Dijkshoorn, C., Tesfa, B., Schreiber, G.: Nichesourcing: Harnessing the power of crowds of experts. In: Knowledge Engineering and Knowledge Management. Springer (2012) 16–20
2. Oosterman, J., Bozzon, A., Houben, G.J.e.a.: Crowd vs. experts: nichesourcing for knowledge intensive tasks in cultural heritage, Int. WWW Conferences Steering Committee (2014) 567–568
3. Basharat, A., Arpinar, I.B., Rasheed, K.: Leveraging crowdsourcing for the thematic annotation of the qur’an. In: Proceedings of the 25th International Conference Companion on World Wide Web, International World Wide Web Conferences Steering Committee (2016) 13–14
4. Basharat, A., Arpinar, I.B., Dastgheib, S., Kursuncu, U., Kochut, K., Dogdu, E.: Semantically enriched task and workflow automation in crowdsourcing for linked data management. International Journal of Semantic Computing **8**(04) (2014) 415–439
5. Sherif, M.A., Ngomo, A.C.N.: Semantic Quran - a multilingual resource for natural-language processing. Semantic Web **6**(4) (2015) 339–345