

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/109923/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Rhode, Matilda, Rana, Omer and Edwards, Timothy 2017. Data capture and analysis to assess impact of carbon credit schemes. [Online]. arXiv. Available at: <https://arxiv.org/abs/1711.07574>

Publishers page: <https://arxiv.org/abs/1711.07574>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



Data Capture & Analysis to Assess Impact of Carbon Credit Schemes

MATILDA RHODE, OMER RANA, School of Computer Science & Informatics, Cardiff University, UK

TIM EDWARDS, Cardiff Business School, Cardiff University, UK

Data enables Non-Governmental Organisations (NGOs) to quantify the impact of their initiatives to themselves and to others. The increasing amount of data stored today can be seen as a direct consequence of the falling costs in obtaining it. Cheap data acquisition harnesses existing communications networks to collect information. Globally, more people are connected by the mobile phone network than by the Internet. We worked with Vita, a development organisation implementing green initiatives to develop an SMS-based data collection application to collect social data surrounding the impacts of their initiatives. We present our system design and lessons learned from on-the-ground testing.

1 INTRODUCTION

Non-Governmental Organisations (NGOs) are increasingly using data to quantify the impact of their work. Performance metrics can be used to inform an organisation's future practices and to demonstrate the impact of programs to investors and donors. Vita[32], a Dublin-based development agency, launched their Green Impact Fund initiative in 2016. The Green Impact Fund is a financially sustainable investment initiative aimed at delivering social and environmental benefits to communities in Eritrea and Ethiopia, where the effects of global warming cause increasing agricultural and economic challenges. Companies or organisations can invest in the Green Impact Fund. Invested sums are transformed into low-cost green commodities (fuel efficient stoves, solar lamps and fixing broken water pumps). The carbon saved by these initiatives is then sold on a voluntary carbon exchange market and the investment is returned to the investors with interest. Vita presently use in-person data collection to calculate the tonnes of carbon offset by their programs. The social impact, however, of these interventions are not measured. Collecting data in-person is costly and these surveys are conducted annually, failing to capture socio-economic benefits at a more granular level. To remedy this assessment shortfall Vita realize investors want additional indicators to specify the social impact of these interventions whilst also requiring reassurances that such work is conducted responsibly when interventions and subsequent data collection are sensitive to the needs and desires of the communities involved. We describe the process, and report the lessons learned, from developing a remote data collection application in collaboration with Vita to ascertain both the opportunities and barriers this application has for capturing social impact responsibly.

Beyond the cost of enumerating time collecting data, the present system uses paper forms before being transferred into an electronic spreadsheet for analysis. This system identifies a number of places for human error in transcription. The application also presents an opportunity to expand current data collection from the basic measurement of carbon off sets to include new lines of questioning suitable for assessing a broader spectrum of "social impact" indicators.

Remote data collection mitigates a number of the costs of in-person analogue data collection. Harnessing existing networks incurs fewer costs than establishing or augmenting a network. Whilst just 45% of the world have access to the Internet, there are 0.96 as many mobile phones as there are people [4], though many people have more than one. Data collection tools using mobile phones exist (e.g. [22], [16]), but existing tools were either unsuited to Vita's needs or required a developer to transform an existing software framework into a fit-for-purpose and usable application.

We describe a system designed to analyse data collected over the Global System for Mobile Communications (GSM) network, capable of assimilating data collected from individual participants and by trusted enumerators. The application for creating and sending questionnaires as well as visualising data is web-based, but transitions to operate over the GSM network for asking questions and receiving the responses. We describe the requirements gathered through interviews with Vita and other experienced enumerators, the system design, the lessons learned from user experience testing and from deploying a prototype in Eritrea.

During development we aimed to create a flexible application that might be re-purposed by other organisations. The application is intended to be usable out of the box, and should not require any programming knowledge to be used. In section 3 we discuss how our approach compares with existing open-source tools for remote data collection. The system architecture is outlined in section 5. In section 8, we discuss the lessons learned from testing our system within the context of Vita’s use case, and the challenges faced in this process.

2 CONTEXT

Vita, a development agency based in Dublin, Ireland, has been working in East Africa for over twenty years [32]. Vita launched the Green Impact Fund [31] in 2016, which seeks to use investment, rather than donation, to bring about green socio-economic improvements for individuals living in rural Eritrea and Ethiopia.

The Green Impact fund comprises several stages. Investors’ money is used to repair (broken) water hand pumps and provide fuel-efficient cook stoves and solar lamps to households [31]. Each of these have both an environmental and a social benefit. Provision of clean water reduces the amount of energy (firewood) needed for sterilisation. Fuel-efficient cook stoves replace open fire pits and require less firewood for the same amount of cooking and use a chimney to extract the harmful emissions away from the household. Reduced firewood consumption may in turn reduce deforestation, which has seen around 20% of woodland in Eritrea disappear since the 1970s for which human activity has been cited as a cause [13]. Deforestation in turn has the capacity to worsen the effects of climate change and of natural disasters [19]. The solar lamps enable activities such as studying in the household after dark and may replace kerosene lamps which cause respiratory problems.

The reduction in carbon emissions from each of these initiatives can be monetised through carbon markets, with each tonne of carbon no-longer being emitted into the atmosphere commanding a *market value*. Vita must demonstrate that these emissions are being saved to the independently-recognised organisation, *Gold Standard*. These measurements are carried out by an external auditor, co2balance, using in-person data collection. In-person collection is used to ensure that certain numerical thresholds are met and that the initiatives being cited to support carbon reductions have in fact been installed. Vita wants also to demonstrate the social impact of their work, but annual, in-person data collection is costly in human hours and may fail to reflect changes that happen during the course of the year e.g. the seasonal trends in social make-up and economic demands of the household.

Our work must also be understood in the context of Vita’s efforts to demonstrate to investors that it is operating ethically, and so early work on the application is framed by the concept of responsible innovation. This concept alerts us to problems created when innovations emerge in the absence of agreed governance structures or rules to moderate the actions of researchers ([34], [26]). Being responsible entails “public dialogue” [18] allowing open-ended governance processes to encourage deliberation. Responsibility means, “taking care of the future through collective stewardship of science and innovation in the present” [26]. Rather than develop the application “at distance” from the users and communities benefitting from Vita’s work, the intention is to develop this application through dialogue

with communities. This enables us to make a more informed assessment of social impact of this work, and the use of the application can be modified in ways responsive to the realities of the respondents’ experiences and needs. This is especially important in a context where the success of Vita’s operations rely on inclusivity and when the Eritrean authorities are highly suspicious of possible foreign interference in national matters.

The table below indicates the summary statistics for Ethiopia and Eritrea according to the CIA World Factbook estimates¹[2-4]

total / population	Ethiopia [3]	Eritrea [2]	World [4]
Landlined phones	0.01	0.01	0.15
Mobile phones	0.41	0.07	0.96
Internet Users	0.12	0.01	0.43

Table 1. Ratio of landlined phones, mobile phones, and internet users to total population in Ethiopia, Eritrea, and globally

The data above indicates the relative numbers of mobile phone ownership to population size. Globally less than 96% of the population are connected to a mobile phone network – as some individuals have multiple phones. Whilst mobile phone ownership is relatively low in Ethiopia and Eritrea compared with global rates, it is still significantly higher than the number of individuals with internet access. Furthermore, smartphone ownership rates are lower than basic phone ownership at 8% in Ethiopia [7], compared with 43% owning a basic or featureless phone. In Eritrea third generation (3G) mobile wireless mobile telecommunications technology is not available [14]. Hypothesising that mobile phone ownership is positively correlated with wealth, we conjecture that direct data collection from individuals will not conduct a balanced survey for this use case.

We selected SMS communication above voice communication after researching similar ventures in other developing countries. Whilst [23] found that voice reporting gave fewer data errors than SMS reporting for collecting health data (in Gujarat, India), participants were trained prior to using the software. Crawford et al. [10] argue that SMS is the preferred means for disseminating health data. This is due to lower costs, higher delivery success, and higher levels of intended or actual behaviour change. The context of [10] matches the environment of the Green Impact Fund more closely, involving untrained voluntary participants communicating over an unreliable network infrastructure. Voice-based communication, however is more accessible to a wider population. Eritrea has a literacy rate of 73.8% [2], and Ethiopia 49.1% [3]. Furthermore, the Green Impact Fund primarily focuses on providing their fuel efficient cook stoves and solar lamps to women. In both countries, the literacy rate among women is lower than the countrywide average and in Ethiopia; women are also less likely to own a phone [33]. Despite this, network reliability can interfere with voice transmission or potentially cut phone-calls short. The former may increase transcription errors (for data analysis) and the latter introduce a bias towards the initial questions in a survey having a higher response rate.

To combat the partial and biased ownership of mobile phones as well as barriers to text communication, we build in the flexibility for trusted enumerators to collect data as well as for individuals to provide data about themselves. This is still of value to organisations as it allows dynamic changes to data collection (from a web application) in the field and also benefits from the automated analysis outlined in section 6.3.

¹Data estimates from July 2015

3 RELATED WORK

Remote data collection is attractive for its potential to reduce both the time and financial cost of in-person collection. Existing tools for collecting data vary according to their purpose. Tools may be developed for specific projects or may be flexible enough to meet the needs of a number of different projects. The suitability of a given tool to a project may be restricted by the type and capability of telecommunication network(s) over which data is collected, the intended users, the possible data formats etc. The system we present in section 5 seeks to add to the set of existing tools for remote data collection with reference to Vita's specific use case, and others like it.

3.1 Remote data collection tools

Applications designed for a specific purpose may aggregated data to identify trends, e.g. information about crises [22], or singular data points may be useful, e.g. for reporting road incidents in Nairobi [25]. Communication channels can also be reversed to disseminate data, e.g. health information [10] or agricultural information [24]. It is also possible to combine the two by disseminating collected information. Although tailoring tools to specific use cases may enable greater efficiencies for the project in question, Tangmunarunkit et al. [27] highlight the inefficiencies of this approach when many projects have very similar requirements. Other tools are flexible and have been designed to fit the needs of various remote data collection objectives ([6], [16], [27], [1]). Open source applications such as the Open Data Kit (ODK) [6] further allow developers to tailor flexible applications to specific projects or to develop software packages to add greater flexibility to the application. Some tools, such as RapidSMS [30] are specifically designed as a framework upon which developers can build, the disadvantage here being that a developer is required to create a usable application. Flexibility is maximised when applications can be used out-of-the-box and also developed further by software engineers if required. Both Ohmage [27] and ODK fulfil these criteria, but require an Internet connection to send data.

Polling and surveying utilises networks, whether these be terrestrial or digital. Internet-based solutions such as the Android-based ODK and iOS/Android-compatible Ohmage and EpiCollect have several advantages over collection applications using the GSM network. Firstly, multimedia data can be sent easily over the Internet, e.g. images and geolocation data. Secondly, Internet-based collection can reduce costs by storing data locally on enumerator devices until a cheap Internet source is available. SMS messages and voice calls cost per message or per minute and so savings can only easily be made through "unlimited" plans. However, as shown in table 1, Internet connectivity is less prevalent than GSM network connections globally.

Tools may impress a data structure onto the collection process. FrontlineSMS [16] is highly flexible and does not impose a structure on data collection. The advantage of this approach is the flexibility afforded to the question-setter. Collecting data can be riddled with concerns over introducing biases and as such it may be undesirable for the software to place restrictions on question designers. ODK, conversely uses questionnaires and structured response formats (such as check-boxes) with data types, e.g. geolocation data. This approach can enable quick data analysis either visually or numerically with reference to the questionnaires, questions and response location. In cases of open source applications, data analysis tools may be added on more easily with structured data responses, as illustrated by Imran et al. [15] in their analysis of health care questions submitted by SMS to U-Report Zambia [29], to which trained health care professionals then respnd. In [15], questions are sent by SMS then classified into categories using a trained machine learning classifier with an accuracy of 82% area under the curve. The authors simulated a *live testing* scenario and the health workers reported that the system helped to process the large volume of SMS queries more quickly. Imran et al. [15] demonstrates how automatic parsing of SMS messages can elevate the usefulness of an initiative, particularly in cases where action

might be time sensitive. The authors note the difficulty in parsing natural language and of categorising messages which fell into one or more of the predetermined categories. As our model is based on a system which prompts for answers, in this initial phase we simply count word frequency to get a sense of the responses from users. In future iterations of the application machine learning could be employed in a similar fashion to [15] for collecting data on particular topics.

Whilst ODK [6], EpiCollect [1] and Ohmage [27] comprises many of the features required by Vita, each of these applications operates using feature phones and is not appropriate for use in Eritrea.

3.2 Data collection using SMS

Many NGOs use SMS to disseminate and/or collect information. The success of specific projects such as UReport [28] to promote citizen engagement and Ushahidi [22], initially developed following the riots connected with the 2008 elections in Kenya, have now been rolled out to other countries as reporting mechanisms operating more generally. These applications aggregate data to sense changes in population reporting, but do not enable organisations to prompt answers to specific questions, as is required by Vita.

Frameworks such as FrontlineSMS and RapidSMS, commissioned by UNICEF in 2007, can be built on to create data collection applications for specific use cases. These frameworks require a developer to create the application, however. Our framework below is intended to suit Vita's needs. We hypothesise, based on the success of applications such as ODK collect, that a structured method for polling using SMS will fit many use cases as well as Vita's. In creating an SMS polling mechanism it will also be possible to disseminate data (technically equivalent to asking a question requiring no answer), thus widening the potential uses of the application.

4 SYSTEM REQUIREMENTS

In designing the system we wanted to build on the existing knowledge of experienced data collectors and conduct tests to closely approximate the situations in which the application might be used. Vita and co2balance were vital in the development of the project. The initial impetus for the system described here was borne out of Vita's Green Impact Fund launch event, during which it became apparent that potential investors were interested in quantifying the social impact of their donations.

To prove compliance with the Gold Standard, Vita employ co2balance to conduct in-person data collection. Though Vita could append social impact questions to this existing data collection practice, there were concerns that this would make the questionnaire too long, increasing the potential for data bias due to respondent fatigue. Furthermore these questionnaires are carried out annually and Vita may want to collect social data at closer intervals. It is also worth noting some reasons that remote data collection is not always preferred. Firstly, as mentioned above, not every respondent has access to a mobile phone, thus data will be biased towards representing those with access. There are a number of other biases that may be introduced through remote collection, as SMS messaging is not free, there may be a bias towards wealthier respondents. There may also be a bias towards those who are more frequent users of mobile phones, however it is also possible to argue that door-to-door surveys are biased towards representing the views of those most likely to answer the door. This system is not intended as a panacea for all issues surrounding in-person data collection, but as a complementary tool. As such it can be used both by enumerators and by individual participants, with the system aggregating responses to the same question regardless of how it was collected. Responses collected by paper should also be possible to add to the system in case Vita wishes to incorporate old data into the new system.

Our research repeatedly indicated features which represented a trade off between convenience and data integrity. If users are experienced enumerators, the system should not impede on the way in which they wish to collect data and if they are not, perhaps a superuser should be able to curb their actions within the guidelines of good data collection practices. We built on the information gathered from Vita and co2balance and took insights from papers such as [9]. In [9], the authors interview a number of NGO employees using ODK to tease out the security concerns of enumerators, during which the authors discover that many enumerators fear data loss as well as exposure of sensitive data. We aimed to protect against the possible causes of data loss such as accidental or purposeful tampering by an authorised and unauthorised individuals as well data lost as as result of poor network integrity.

From our research we distilled the following three principles that the application should seek to meet:

- (1) *Facilitate the collection of as much (accurate) data as possible*: promote those aspects which will encourage users at both ends of the system to generate accurate data;
- (2) *Automate*: use computation to save resources (without violating 1)
- (3) *Provide flexibility*: design the system such that it might be adopted for different forms of SMS communication i.e. information dissemination, collection and collection followed by dissemination.

5 SYSTEM ARCHITECTURE

Using the GSM network, communication may be conducted via text or using voice communication. Whilst voice communication does not entail literacy requirements, we present a text-based application using short message service to capture the data analytic capabilities of Internet-based applications such as Open Data Kit [6], using structured questionnaires. Due to our specific use-case, Vita's concerns over network robustness lean in favour of SMS communication as voice-based data collection occurs at the inconvenience of the participant and network failures may lead to a bias of more data for earlier questions in a given questionnaire. Flexibility for collection by an enumerator is the temporary solution to literacy requirements, but we chose an SMS gateway provider (Twilio [17]) capable of using voice communication with a view to incorporating this feature in the future.

The robustness of SMS is built into its protocol. Although SMS messages are subject to a 160 character limit, longer messages are enabled by "concatenated SMS". Concatenated, or long SMS carry a header of metadata to indicate that they are one in a series of messages, along with their position in that series [20]. If all parts of a message, or even any part, is not able to be sent within the validity period for the message as specified by the network operator, typically two days [20], the message will fail to send. In a questionnaire context this means the next question will not be triggered. If the gateway is unable to detect a failed message, should the next question be sent after two days if there is no response? Or did the respondent choose not to reply and would rather not be bothered?

The system presented here effectively creates a structured system for two-way communication. At one end of the system sits the data analyst and/or questionnaire architect, at the other a data enumerator or an individual sending responses to the server. To harness the mass-polling capabilities of SMS we use an SMS gateway across a Global System for Mobile Communications (GSM) network. Though the use of GSM is intended to reach areas which have poor or non-existent internet connectivity, we have designed a web-based application for the data analysis and questionnaire-setting end of the system. The web-based application suits the structure of decentralised and globally disparate organisations such as Vita, with employees in Eritrea, Ethiopia and Ireland, as data and questionnaires can be accessed and sent from anywhere with an internet connection. In Eritrea and Ethiopia, Vita's headquarters are based in the capital and will have access to Internet.

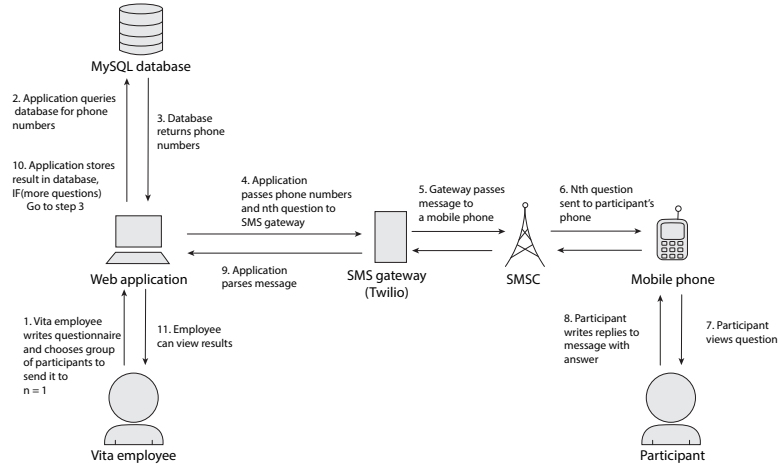


Fig. 1. Data flow within application

Figure 1 illustrates the application architecture. The database houses all questions, questionnaires, responses, respondent and employee data. The SMS gateway connects the web application and a short message service center (SMSC), which then transmits the message to the intended recipient(s) over the GSM network. This model can be adapted to enable information dissemination as well as collection. The web application directly pulls and pushes data to the database, queries to the database are only made when editing data (not when changing view) thus reducing the number of requests to the database server and potential for failed requests due to unreliable Internet connections.

The data structure for sending questionnaires operates using questionnaires and user groups. This enables the polling of subsets of individual respondents as well as creating a group of enumerators, who may be sent different questions to users. Figure 2 gives an outline of a situation in which a questionnaire, containing three questions, has been sent to a group, containing two members. Only one member has responded and they have replied to all three questions.

6 SMS DATA COLLECTION APPLICATION

The web application is the management portal for creating and sending questionnaires as well as for data analysis. For consistency the website employs a tabular system comprising questionnaire, user and response data, as well as staff data for superusers. A sidebar on the left hand side enables navigation between tables (Figure 3). Each table may be filtered according to its contents (Figure 4). Multiple filters can be used at once to show the subset of rows fulfilling all applied filters. Filtering enables shortcuts to common processes such as creating user groups or questionnaires and downloading data as a comma separated values (CSV) file. Depending on access levels (Figure 2) employees are able to add, edit and delete table data.

6.1 Creating and sending questionnaires

Creating and sending a questionnaire requires a questionnaire and a user group, comprising one or more questions or users respectively. Questions are added to questionnaires and users to user groups either by manually selecting from the list of possible candidates (Figure 5) or by filtering the questions/users table and choosing "Create [Questionnaire / User Group] from selected".

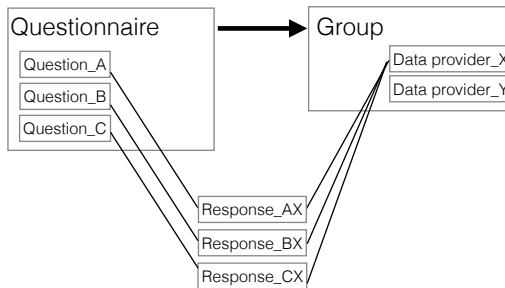


Fig. 2. Illustration of questionnaire relationship to user group

Question	Response Type	Total Acks	Total Responses	Created	Last Updated	Start Date
10 How are you?	Free text	19	11	2016-07-05 00:00:00	2016-07-05 00:00:00	2016-07-05 00:00:00
11 What is 2+2?	Free text	10	10	2016-07-07 00:00:00	2016-07-07 00:00:00	2016-07-07 00:00:00
12 Question test	Set options (e.g. A=Other, B=None)	4	4	2016-07-08 15:00:00	2016-07-08 15:00:00	2016-07-08 15:00:00
13 How old are you?	Number	8	8	2016-08-01 18:11:37.0	2016-08-01 18:11:37.0	2016-08-01 18:11:37.0
14 Question about you park? Ready	Set options (e.g. A=Other, B=None)	0	19	2016-08-03 14:00:00	2016-08-03 14:00:00	2016-08-03 14:00:00
15 How old is your dog?	Free text	0	0	2016-08-18 10:29:24.0	2016-08-18 10:29:24.0	2016-08-18 10:29:24.0
16 Use Reply To: No	Set options (e.g. A=Other, B=None)	1	10	2016-08-19 13:42:23.0	2016-08-19 13:42:23.0	2016-08-19 13:42:23.0
17 Has a cat? Reply: A=OK, B=Not	Set options (e.g. A=Other, B=None)	4	2	2016-08-21 10:42:33.0	2016-08-21 10:42:33.0	2016-08-21 10:42:33.0
18 PS: get of Reply: A=OK, B=Not away	Set options (e.g. A=Other, B=None)	1	2	2016-08-21 10:41:18.0	2016-08-21 10:41:18.0	2016-08-21 10:41:18.0
19 How often do you watch TV?	Set options (e.g. A=Other, B=None)	0	0	2016-08-22 15:11:58.0	2016-08-22 15:11:58.0	2016-08-22 15:11:58.0
20 How often do you watch TV?	Set options (e.g. A=Other, B=None)	0	0	2016-08-22 15:11:58.0	2016-08-22 15:11:58.0	2016-08-22 15:11:58.0
21 How many hours do you spend reading each	Number	0	4	2016-08-23 17:58:55.0	2016-08-23 17:58:55.0	2016-08-23 17:58:55.0
22 How many hours do you spend reading each	Number	0	5	2016-08-23 17:59:48.0	2016-08-23 17:59:48.0	2016-08-23 17:59:48.0
23 How many hours do you spend reading each	Number	0	1	2016-08-24 07:29:53.0	2016-08-24 07:29:53.0	2016-08-24 07:29:53.0
24 How many hours do you spend reading each	Number	0	0	2016-08-24 07:30:08.0	2016-08-24 07:30:08.0	2016-08-24 07:30:08.0

Fig. 3. Web application view of questions in database

Fig. 4. Filter overlay for user groups

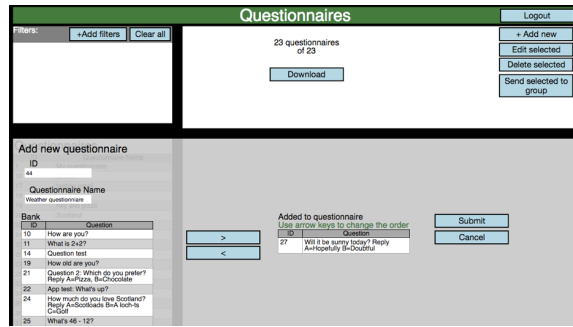


Fig. 5. Creating a questionnaire

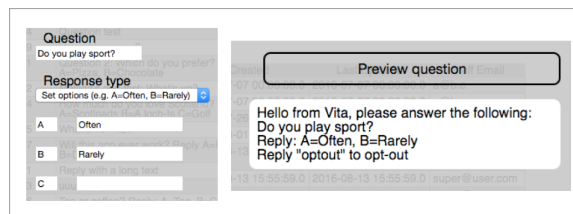


Fig. 6. Previewing a question before committing it to database

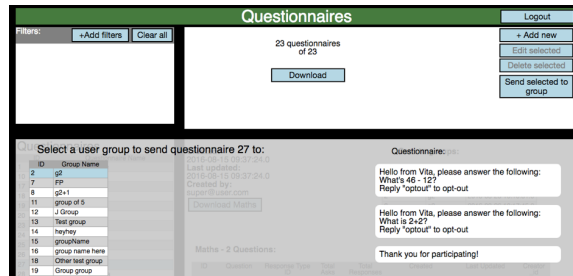


Fig. 7. Preview questionnaire before sending to user group

Questions themselves must be given response types to enable automated response parsing and data analytics. Response format may be numerical, free text or categorical options. After user experience testing we added a “Yes/No” option, which though logically indistinct from a categorical response, was deemed a time-saving and rational option to our testers. Employees are restricted the formatting of options questions based on research by Down and Duke [11] showing the format which was easiest for respondents to understand. Maintaining format also enhances ease of use for repeat data providers. Employees may preview individual questions (Figure 6) as well as entire questionnaires (Figure 7). Previews also reveal settings put in place by the superuser(s) including greetings, thank-yous and opt-out messages.

The user table can be considered a series of “profiles”. Certain data cannot be removed such as profile creation date, last edit, and phone number; but employees may add further columns to the tables to create more user data. For example birth date and languages spoken may be included enabling quick filtering so that a user group of all under-25s may be created and polled. This feature permits the flexibility to use enumerators rather than polling individuals. An

organisation may add an attribute (column) to user profiles denoting “staff”. Staff may then be polled in the field with per-household, per-street, per-village questionnaires as preferred.

We imagine a scenario in which an organisation’s employees hold various roles. Staff profiles are each attributed access levels which enables the application to run with least privilege. This is a security measure against both tampering and accidental data loss. The access levels are denoted in Table 2.

Data Access Level	View	Filter	Down-load	Add	Edit and Delete	Send SMS	View & manage staff	Update sys. settings
Read only	✓	✓	✓	✗	✗	✗	✗	✗
Add	✓	✓	✓	✓	✗	✗	✗	✗
Edit	✓	✓	✓	✓	✓	✗	✗	✗
Send	✓	✓	✓	✓	✓	✓	✗	✗
Superuser	✓	✓	✓	✓	✓	✓	✓	✓

Table 2. Employee Access Levels

Superusers may place further restrictions on employees beyond their access levels to promote good practice in data collection. For example superusers can set a maximum questionnaire length to mitigate the chances of respondents becoming bored by questionnaires, resulting in a fatigue bias. Superusers may add other features to questionnaires such as greetings and thank-you messages. Featureless phones do not present SMS messages in a “thread” style so the greeting will enable respondents to recognise the organisation polling them if they have not saved the relevant number or message code.

Superusers may further choose special keywords to allow respondents to opt-out or sign-up to receiving questionnaires using SMS. The sign-up word can be associated with a questionnaire which will be sent after the sign-up keyword. The opt-out word deactivates the user profile until the user sends the sign-up keyword, this is so that their previous responses remain associated with a user profile and data is not lost. These features are enabled through the system settings, which are only accessible by superusers (Figure 8).

6.2 SMS response

Our design borrowed lessons from the models used by market research companies and others using SMS to collect data; the questionnaire waits for a response to the last message sent before sending the next question, this stops the user from being repeatedly asked questions when they have no interest in or are unable to respond. An alternative to question-response-question-response is to send the next question if no response had been received following a pre-specified time period. As the short message protocol encompasses an acknowledgement of message receipt [12], a lack of response can be attributed to a user deciding not to respond or network integrity issues on sending the response. The short message protocol should re-attempt sending until success or until two days elapses. Failures are likely to be due to network error, and may be deduced by noting that all respondents from a particular area have not replied to questions sent within a particular timeframe. Such an occurrence, however may also be an active decision not to respond. The web application contains a complete and uneditable log of all SMS sent and received to the gateway. This is intended for data integrity and may be used at an NGO employee’s discretion to decide whether respondents’ decisions or network integrity is the cause of a lack of responses.

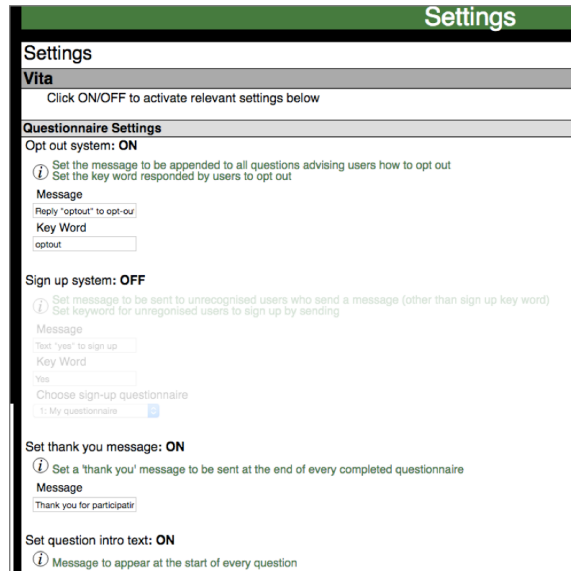


Fig. 8. System settings in web application

Communicating strictly by one new question following an answer to the previous question also enables data parsing to happen automatically. A piece of data is only useful if the question to which it responds is known. It may seem over simplistic to denote an SMS as a response to the last question sent but Church and de Oliveira [8] note that the cost involved in sending SMS messages creates a more structured answer-reply conversation than occurs over instant messaging platforms, in which users will frequently send sentences as single messages in quick succession. If the application does falsely attribute a response to a question, the analyst may examine the SMS logs to determine the question for which it was a response, if any, and alter this in the responses table.

For flexibility we chose to use an SMS gateway API rather than a hardware module, which would be more difficult to distribute and require technical expertise to install. Testing revealed that although our chosen gateway communicated easily between some international network providers, others blocked responses to our tests from being sent and some blocked their delivery. This information was not always available on the API website. There were a number of possible solutions to this problem e.g. using hardware modules or negotiation with network providers.

6.3 Parsing and analysing data

To further reduce the human resources required to process the data, we have included automated parsing of question responses to facilitate data analysis. We concede that automation of human responses may parse data incorrectly for certain edge-cases, so the system maintains a complete log of messages sent and received to stop misinterpreted data being overwritten permanently. Misinterpretation may result from a question response being sent in multiple parts, as described in the previous section or from mis-parsing of categorical data, though we believe that both of these are unlikely occurrences. Each response holds a one-to-one relationship with a user and with a question, metadata about the responses is also stored. Categorical response questions are in turn related to possible answer codes and the meaning of these codes e.g. “A=Everyday, B=Weekly, C=Monthly, D=Never”. We use simple case-insensitive parsing

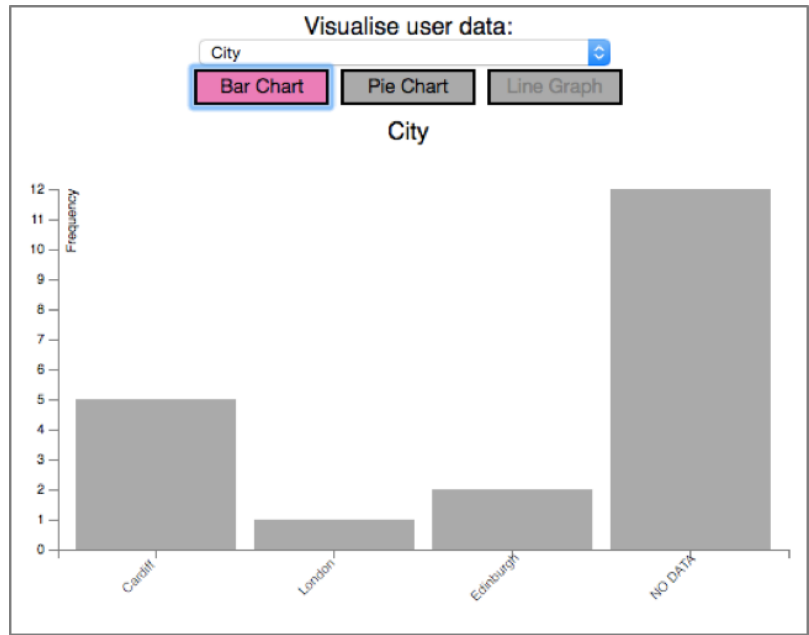


Fig. 9. User city data visualised in suggested bar chart form, can be overridden to view pie chart representation

on responses to categorical answer type questions. Answers are stripped of start and end spaces and punctuation to translate “a”, “A.” and “everyday”, for example, to “Everyday”. The purpose of this parsing is to enable quick data analytics. Response parsing also validates inputs to protect against the possibility of an SQL injection attack, which could cause data to be lost or corrupted.

After logging-in, the landing page for employees is a dashboard, where user and question data is displayed in simple visualisations. Data gathered from the database is transformed into a JSON format representing summary statistics for user attributes and question responses. The visualisations use the d3 JavaScript library for data visualisation [5]. A further development could easily be added to download the summary statistics as a JSON file for analysis use other software.

The dashboard comprises a timeline for an at-a-glance indication of the number of SMS being sent and the number of responses being received. The data relating to each question can also be visualised. Employees select the question from a drop-down menu and data is displayed as a bar chart, pie chart or word cloud if the response type was numeric, categorical or free text respectively. Similarly, user attributes can be selected from a dropdown menu on the right (Figure 10) and a visualisation type is suggested based on data format, though this can be overridden (Figure 9). These visualisations are intended for a quick overview of demographics and responses; all data may be downloaded as a CSV file if analysts wish to conduct further investigation in another application.

7 IMPLEMENTATION AND USER EXPERIENCE TESTING

Interviews with experienced data enumerators and research informed the application’s development but we wanted to conduct tests to find further areas for development and improvement. We conducted a user experience test for the web

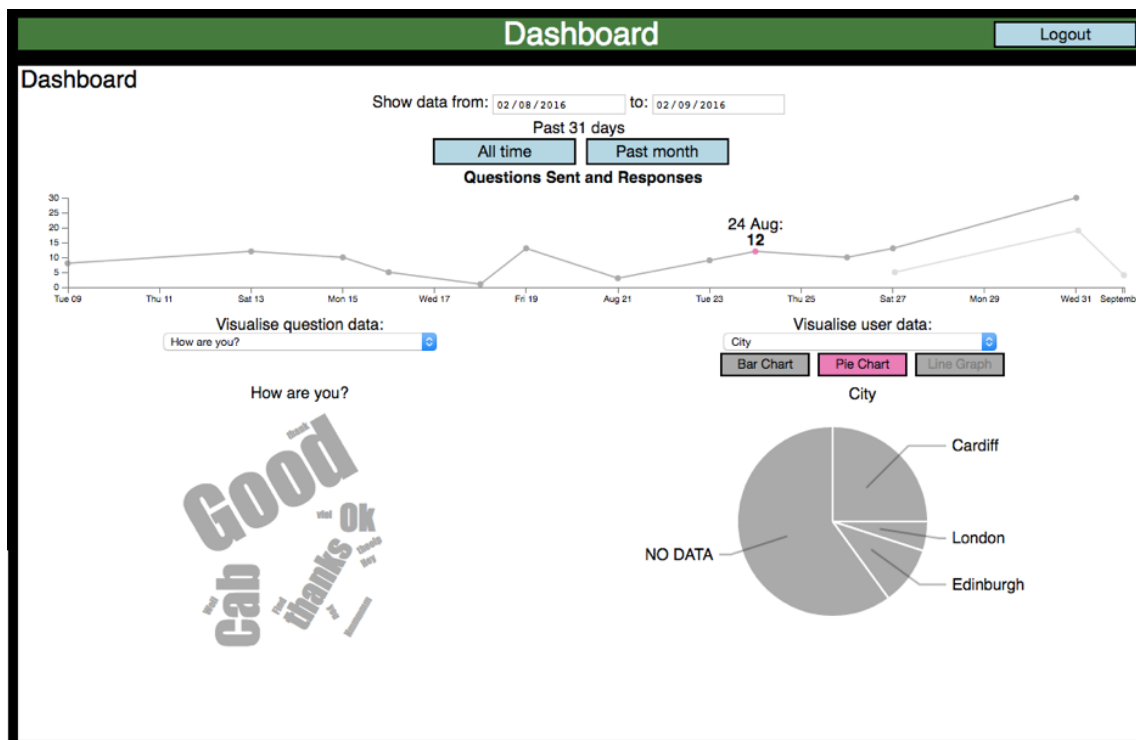


Fig. 10. Screenshot of the dashboard during development

application, an international data collection test, and interviews and system tests with Vita in Eritrea. Four challenges stood out as barriers to the usefulness of the application: network provider blocking of SMS gateways, distribution of mobile phone access, impact of network integrity on enumerator data collection and graphical interface changes for greater ease of use.

Initially we had problems communicating internationally by SMS between the UK and some countries. Using the SMS gateway service provided by Twilio [17], we were able to communicate nationally and with several European countries but experienced difficulties conducting two-way communication in Kenya. We later adapted the system to use a GSM modem containing a SIM card and the Ozeki SMS Gateway [21], which circumvented this issue. A further solution is to negotiate with network providers to allow bulk SMS traffic through from pre-specified phone numbers. The advantage of the latter solution is that the GSM hardware acquisition and installation does not need to be performed.

The mobile phone statistics reported by the CIA world factbook for Eritrea [2] give a figure for the entire nation. Brief interviews in Eritrea indicated that probability of mobile phone ownership grows in proximity to Asmara. In one village close to Asmara, residents estimated that ownership stood close to 40%. Although we did not conduct a numerical survey, this figure is much higher than 7%. Even with a significant margin of error allowed for estimation biases, the implication is that ownership is much lower than 7% in other areas to give this national average. This information emphasised the importance of flexibility so that enumerators could collect data and not rely on participants owning mobile phones to provide data.

We discussed improvements to the application to enable enumerators to provide data quickly and cheaply. If network speeds fall or a network failure occurs while enumerators are in the field collecting data, they do not want to wait for a text asking them the next question to reply. Furthermore the cost of sending a text message is determined per message (if the message is less than 160 characters). We began to develop an encoding system so that enumerators could send data for ten questions in one go. This relies on keywords and would necessitate some training of enumerators. The keywords need to signal questionnaire, question, and respondent. This may appear very close to the paper questionnaire model but its advantages are that human error is mitigated, data is automatically parsed for analysis and this data is integrated with the individual respondent and other enumerators' data. Enumerators would need a copy of the questionnaire, this could be stored electronically, but this would require frequent switching between applications or SMS messages in the phone, two devices or learning by-heart, so a paper copy was deemed the best alternative. The questionnaire number is printed at the top of the page and a list of user data stored too. Keywords are reduced to letters to conserve the 160-character limit. The message style uses codes and ID numbers to establish questionnaire, starting question and respondent and then uses semicolons to separate question answers e.g. "U24F102Q1;A;12;C;More time with family" could be the responses of user 24 in answering questionnaire (form) 102 from questions 1 to 4.

This development later inspired the ability to combine short SMS into single text messages to reduce the cost of the respondent. This feature is still under development. The hesitation over implementing this feature is in understanding the trade-off between the cost of sending an SMS as a deterrent and the data loss that will occur from human error in inputting the data (for untrained individual participants). In the future we would like to conduct experiments to assess the relative merits and demerits of these two systems. As the relative impact of these factors is likely to depend on context, the superuser will have powers to activate or deactivate this feature for flexibility.

We conducted an initial round of user experience testing (four volunteers) in which we asked users to complete various tasks after explaining the purpose of the application but not how to use it. The results were that the membership of questions to questionnaires and users to user groups should be more fluid to enact, and there should be more ways of doing this.

Plans to extend this work would ideally involve user group workshops within targeted communities to explain the application, the purpose of the study and how to use the application to enable *live* data collection over time. However, for this to be done in a responsible way it would also be necessary to engage in dialogue about the fears and desires of respondents in participating in such an evaluation. Transparency is crucial, so *lessons learnt* inform monitoring and evaluation protocols moving forward alongside the expansion of Vita's work planned between now and 2021. In this respect, the application might also be viewed as a mechanism to enable dialogue and not simply as a means to report impact. Indeed, the application has the potential to improve impact through knowledge dissemination.

8 USAGE EXPERIENCE

The system has been used in Eritrea to collect data for the Vita Green Impact Fund. Accurate data collection remains an important challenge for an NGO of this kind, and this was the key reason for developing the proposed system. We gained the following experience from this study:

- Although mobile phone usage and access remain high in this region, access to *smart phones* with internet access remains limited. A system that requires data collection to be supported directly based on an internet-enabled infrastructure is of limited value. The telecomms. infrastructure may be provisioned through a private company,

its operation is often controlled by government. This requires any additional service which must operate over such infrastructure to go through various approval processes.

- The energy infrastructure required to sustain data collection points (e.g. mobile base stations, internet access points, etc) may also be fragile. This could be due to electricity load reduction strategies that limit the time grid-based supply is made available (a common occurrence in many other developing countries also). Planning and supporting additional sources of power is essential, and could be in the form of lead-acid batteries, solar panels, etc.
- Understanding the language in which the questionnaire should be developed can be a challenge, especially if the system needs to be operate in rural areas. Literacy levels must be taken into consideration, along with local dialects which may limit the quality of data that is obtained using such automated approaches. A voice-to-text system may also be used to enable participants to submit their responses in voice format, which can subsequently be converted into text prior to transmission. However, such a system requires an app. to be supported on the device of the participant.

9 CONCLUSIONS

We describe a system to support data collection for carbon credit initiatives, particularly for regions with limited internet coverage. An SMS gateway is proposed that enables an NGO (for instance) to send out a questionnaire to potential participants, and the subsequent responses can then be recorded into a database for further analysis. Such a system may be used directly by participants involved in the carbon credit scheme (depending on the levels of literacy within a region) to volunteers who act as intermediaries to collect the data. The proposed approach can significant increase the frequency a which data collection can be carried out, to supporting reporting of findings to potential donors and investors.

In this instance, a Web-based visual interface is provided for the NGO to: (i) create a repository of questions that can be re-used across questionnaires. Such a repository is particularly useful if multiple languages (or dialects) exist within a particular region, enabling an administrator to choose questions appropriately; (ii) create multiple questionnaires that can be sent out to different participant groups; (iii) create analysis queries on the collected data, the results of which can be visualised using a *word cloud* – to highlight common terms in the collected responses, or to analyse the data based on demographic information. The data can also be exported to a third party system, such as Matlab or SPSS, for further analysis. The current system uses a central database to hold questions and the results of analysis. This can also be extended to a cloud-hosted database that could be accessible by an NGO that operates in mutiple regions. The data flow presented in Figure 1 can be extended with a cloud-based system, where the Web application and database can be scaled on-demand.

In addition to the general purpose approach for supporting data collection, we also discuss how this particular system has been used by Vita Green Impact Fund in Eritrea. We investigated a number of factors that could limit potential take up of the proposed system in practice, and limitations with practical deployment of the system due to infrastructural limitations (telecomm. network & energy grid).

REFERENCES

- [1] David M Aanensen, Derek M Huntley, Edward J Feil, Brian G Spratt, and others. 2009. EpiCollect: linking smartphones to web applications for epidemiology, ecology and community data collection. *PloS one* 4, 9 (2009), e6968.
- [2] Central Intelligence Agency. 2015. CIA World Factbook: Eritrea. (2015). <https://www.cia.gov/library/publications/the-world-factbook/geos/et.html>

- [3] Central Intelligence Agency. 2015. CIA World Factbook: Ethiopia. (2015). <https://www.cia.gov/library/publications/the-world-factbook/geos/er.html>
- [4] Central Intelligence Agency. 2015. CIA World Factbook: World. (2015). <https://www.cia.gov/library/publications/the-world-factbook/geos/xx.html>
- [5] Michael Bostock, Vadim Ogievetsky, and Jeffrey Heer. 2011. D³ data-driven documents. *IEEE transactions on visualization and computer graphics* 17, 12 (2011), 2301–2309.
- [6] Waylon Brunette, Mitchell Sundt, Nicola Dell, Rohit Chaudhri, Nathan Breit, and Gaetano Borriello. 2013. Open Data Kit 2.0: Expanding and Refining Information Services for Developing Regions. In *Proceedings of the 14th Workshop on Mobile Computing Systems and Applications (HotMobile '13)*. ACM, New York, NY, USA, Article 10, 6 pages. DOI: <http://dx.doi.org/10.1145/2444776.2444790>
- [7] Pew Research Center. 2016. Smartphone ownership rates skyrocket in many emerging economies, but digital divide remains. (2 2016). <http://www.pewglobal.org>
- [8] Karen Church and Rodrigo de Oliveira. 2013. What's Up with Whatsapp?: Comparing Mobile Instant Messaging Behaviors with Traditional SMS. In *Proceedings of the 15th International Conference on Human-computer Interaction with Mobile Devices and Services (MobileHCI '13)*. ACM, New York, NY, USA, 352–361. DOI: <http://dx.doi.org/10.1145/2493190.2493225>
- [9] Camille Cobb, Samuel Sudar, Nicholas Reiter, Richard Anderson, Franziska Roesner, and Tadayoshi Kohno. 2016. Computer Security for Data Collection Technologies. In *Proceedings of the Eighth International Conference on Information and Communication Technologies and Development*. ACM, ACM, 2.
- [10] Jessica Crawford, Erin Larsen-Cooper, Zachariah Jezman, Stacey C Cunningham, and Emily Bancroft. 2014. SMS versus voice messaging to deliver MNCH communication in rural Malawi: assessment of delivery success and user experience. *Global Health: Science and Practice* 2, 1 (2014), 35–46.
- [11] Joel Down and Simon Duke. 2003. SMS polling. A methodological review. In *ASC. Association for Survey Computing*, 277–286.
- [12] ETSI. 2010. *Digital cellular telecommunications system (Phase 2+); Universal Mobile Telecommunications System (UMTS); Technical realization of the Short Message Service (SMS)*. ETSI. V9.3.0.
- [13] Mihretab G Ghebregabher, Taibao Yang, Xuemei Yang, Xin Wang, and Masihulla Khan. 2016. Extracting and analyzing forest and woodland cover change in Eritrea based on landsat data using supervised classification. *The Egyptian Journal of Remote Sensing and Space Science* 19, 1 (2016), 37–47.
- [14] GSMarena. 2017. Network Coverage in Eritrea. (2017). <http://www.gsmarena.com/network-bands.php3?sCountry=Eritrea>
- [15] Muhammad Imran, Patrick Meier, Carlos Castillo, Andre Lesa, and Manuel Garcia Herranz. 2016. Enabling digital health by automatic classification of short messages. In *Proceedings of the 6th International Conference on Digital Health Conference*. ACM, 61–65.
- [16] Occam Technologies Inc. FrontlineSMS. (????). <http://www.frontlinesms.com/>
- [17] Twilio Inc. 2017. Twilio. (2017). <https://www.twilio.com/>
- [18] Alan Irwin. 2006. The politics of talk: coming to terms with the finewfscientific governance. *Social studies of science* 36, 2 (2006), 299–320.
- [19] Edmond J Keller. 1992. Drought, war, and the politics of famine in Ethiopia and Eritrea. *The Journal of Modern African Studies* 30, 04 (1992), 609–624.
- [20] Gwenal Le Bodic. 2003. *Mobile messaging technologies and services: SMS, EMS, and MMS* (2 ed.). John Wiley & Sons, Chichester.
- [21] Ozeki Informatics Ltd. 2000–2017. Ozeki SMS Gateway. (2000–2017). <http://www.ozekisms.com/>
- [22] Ory Okolloh. 2009. Ushahidi, or fitestimonyfi: Web 2.0 tools for crowdsourcing crisis information. *Participatory learning and action* 59, 1 (2009), 65–70.
- [23] Somani Patnaik, Emma Brunskill, and William Thies. 2009. Evaluating the accuracy of data collection on mobile phones: A study of forms, SMS, and voice. In *Information and Communication Technologies and Development (ICTD), 2009 International Conference on*. IEEE, 74–84.
- [24] Christine Zhenwei Qiang, Siou Chew Kuek, Andrew Dymond, Steve Esselaar, and IS Unit. 2011. Mobile applications for agriculture and rural development. *World Bank, Washington, DC* (2011).
- [25] Darshan Santani, Jidraph Njuguna, Tierra Bills, Aisha Walcott-Bryant, Reginald E. Bryant, Jonathan Ledgard, and Daniel Gatica-Perez. 2015. CommuniSense: Crowdsourcing Road Hazards in Nairobi. *CoRR abs/1506.07327* (2015). <http://arxiv.org/abs/1506.07327>
- [26] Jack Stilgoe, Richard Owen, and Phil Macnaghten. 2013. Developing a framework for responsible innovation. *Research Policy* 42, 9 (2013), 1568–1580.
- [27] H. Tangmunarunkit, C. K. Hsieh, B. Longstaff, S. Nolen, J. Jenkins, C. Ketcham, J. Selsky, F. Alquaddoomi, D. George, J. Kang, Z. Khalapyan, J. Ooms, N. Ramanathan, and D. Estrin. 2015. Ohmage: A General and Extensible End-to-End Participatory Sensing Platform. *ACM Trans. Intell. Syst. Technol.* 6, 3, Article 38 (April 2015), 21 pages. DOI: <http://dx.doi.org/10.1145/2717318>
- [28] U-Report. 2017. U-Report Uganda. (2017). <http://ureport.ug/about/>
- [29] U-Report. 2017. U-Report Zambia. (2017). www.zambiaureport.org/
- [30] UNICEF, Caktus Group and Meraka Institute, Entropy Free LLC, Dimagi, Columbia University's Earth Institute, and ThoughtWorks. 2007–2017. RapidSMS. (2007–2017). <https://www.rapidsms.org/>
- [31] Vita. 2016. Vita Green Impact Fund: Investment Memorandum. (2016). <http://greenimpact.web.ie/wp-content/uploads/2016/03/Investment-Memorandum.pdf>
- [32] Vita. 2017. Vita – Investment Partner with African Communities. (2017). <http://www.vita.ie/>
- [33] Mark West and Chew Han Ei. 2014. *Reading in the mobile era: A study of mobile reading in developing countries*. UNESCO.
- [34] Brian Wynne. 2002. Risk and environment as legitimacy discourses of technology: reflexivity inside out? *Current sociology* 50, 3 (2002), 459–477.