

# The Humanitarian eXchange Language: Coordinating Disaster Response with Semantic Web Technologies

**Editor(s):** Name Surname, University, Country

**Solicited review(s):** Name Surname, University, Country

**Open review(s):** Name Surname, University, Country

Carsten Keßler<sup>a</sup> and Chad Hendrix<sup>b,\*</sup>

<sup>a</sup> *Institute for Geoinformatics, University of Münster, Germany*

*E-mail: carsten.kessler@uni-muenster.de*

<sup>b</sup> *United Nations Office for the Coordination of Humanitarian Affairs, Geneva, Switzerland*

*E-mail: hendrix@un.org*

**Abstract.** The Humanitarian eXchange Language (HXL) is a project by the United Nations Office for the Coordination of Humanitarian Affairs that aims at refining data management and exchange for disaster response. Data exchange in this field, which often has to deal with chaotic environments heavily affected by an emergency such as a natural disaster or an armed conflict, still happens mostly manually. The goal of HXL is to automate many of these processes, saving valuable time for staff in the field and improving the information flow for decision makers who have to allocate resources for response activities. This paper presents a case study on this initiative, which is set to significantly improve information exchange in the humanitarian domain. We introduce the HXL vocabulary, which provides a formal definition of the terminology used in this domain, and an initial set of tools and services that produce and consume HXL data. The HXL system infrastructure is introduced, along with its data management principles. The paper concludes with an outlook on the future of HXL and its role in the humanitarian ecosystem.

**Keywords:** Disaster Management, Humanitarian Aid, Linked Open Data, Vocabulary, Data Management, Tools

## 1. Introduction

Events such as large-scale natural disasters or armed conflicts often affect major parts of the local population. If these events exceed the capacity of a government to fully respond, support from the international community is required to address the affected population's need for shelter, food, water, sanitation, and medical care. The Office for the Coordination of Hu-

manitarian Affairs<sup>1</sup> (OCHA) is the United Nations' division responsible for the orchestration of response actions to events such as the 2010 Haiti earthquake and the recent armed conflict in Mali.<sup>2</sup> A specific challenge for OCHA—in addition to the often confusing situation in the affected regions—is the large number of organizations that need to be coordinated to address the affected population's needs. These organizations range from large institutions such as the International Red Cross and Red Crescent Movement

---

\* Any opinions expressed herein are those of the authors and do not represent official positions of the United Nations Office for the Coordination of Humanitarian Affairs.

---

<sup>1</sup> See <http://unocha.org>.

<sup>2</sup> These are only two of the 25 humanitarian operations in which OCHA is currently working.

to other UN divisions such as the UN Refugee organization (UNHCR) and national governmental bodies down to small, local non-governmental organizations (NGOs) with a handful of employees. The database of OCHA's financial tracking service<sup>3</sup> currently lists over 5500 organizations—and each of these organizations use a different system to handle the data about their activities, ranging from full-fledged enterprise information management systems to relational databases to simple spreadsheets. In a large-scale disaster, hundreds of these organizations need to be coordinated; in the aftermath of the Haiti earthquake, response data from an estimated 600 organizations was collected by OCHA. Although this did not represent the total humanitarian activity in country at the time, it did represent a massive challenge in data reconciliation.

Collecting and integrating data to optimize the response efforts in such a heterogeneous and distributed environment is challenging and still leads to situations where OCHA's information management officers on site integrate data by manually copying data from one spreadsheet into another. Widely varying transliterations of placenames and differences in units of measurement for humanitarian interventions add further difficulties to the task of compiling a common operational picture. It is extremely unlikely, though, that all involved organizations can be convinced to use a common information system, due to the differences both in size and topical focus. Therefore, the idea of a common exchange format for humanitarian data was brought up within OCHA in 2011. The rationale behind the exchange format was that it would allow each organization to retain their established data management practices, but still facilitate data sharing with OCHA and other collaborators. At this point, this exchange format was supposed to be developed as an XML schema and therefore named the *Humanitarian eXchange Language* (HXL; ['hɛksl]). HXL was first discussed in public in a breakout session at the International Crisis Mappers Conference in November 2011, where several participants suggested a Semantic Web approach to tackle this large-scale data integration problem.

This paper reports on the current state of the Humanitarian eXchange Language and describes what has been achieved since fall 2011. A schematic overview of the HXL infrastructure and the different components this paper reports on is shown in Figure 1. In

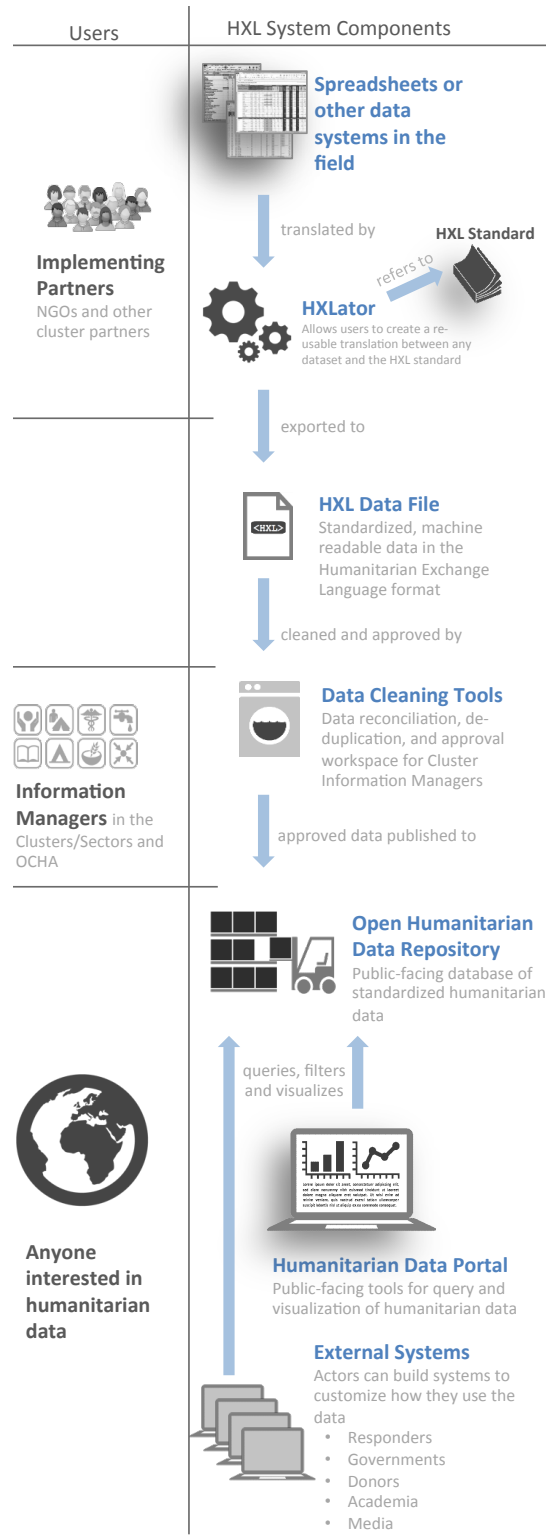


Fig. 1. High-level overview of the HXL system components.

<sup>3</sup>See <http://fts.unocha.org/>.

the next section, we first refer to relevant related work and point out how HXL differs from these existing solutions that address a similar problem space. The requirements for the development of HXL are laid out in Section 3 to show why we chose an approach based on Semantic Web technologies. Section 4 introduces the current state of the HXL vocabulary and the underlying design principles. We also discuss the development and iteration process we follow in the ongoing extension of the vocabulary. In Sections 5 and 6, we address the different means of producing HXL data and scenarios for its consumption in applications, respectively. Data management principles and workflows are discussed in Section 7, followed by conclusions and an outlook on the future of the Humanitarian eXchange Language.

## 2. Related Work

Data exchange between different actors has been a challenge both for OCHA and the humanitarian field as a whole for a long time. OCHA and other organizations have tried to address different aspects of this problem by developing multiple siloed data collection and management systems, often requiring cut-and-paste operations to feed these systems with data from the actors' own systems. Examples include the Common Request Format (CRF) [27] to streamline information requests and the Multi-Cluster/Sector Initial Rapid Assessment (MIRA) framework [16], which serves as a standardized way for the initial assessment of humanitarian needs at the onset of a disaster. Another example is the Single Reporting Format, which provided a comprehensive system for collecting humanitarian activity data in Pakistan, but met with limited adoption in the field. In some cases, the humanitarian community has been successful at defining standards for describing some humanitarian data. The Inter-Agency Standing Committee (a forum of organizations involved in humanitarian response) established guidelines describing data regarding the attributes and size of populations affected by a crisis [15].

Most of the current information systems hold data that serve very specific needs. Their contents were only combined to generate the periodic reports for a variety of audiences. These reports are largely generated manually by staff members that are highly familiar with the situation on-site and the corresponding data. In addition to the data kept in specialized information systems, some data—such as geographic and demo-

graphic information—is required across OCHA and in other cooperating organizations. Such data is organized by country in the Common Operational Datasets (CODs) as downloadable files on a website.<sup>4</sup> OCHA is responsible for identifying and updating these datasets, which ideally form the baseline data for organizations responding to a humanitarian crisis. Updates are agreed in advance by the humanitarian community in-country to make sure that all involved organizations always refer to the same version of the CODs.

Outside of OCHA, several different standards for data exchange in the humanitarian domain have been developed over the past years. The most widely used is the Emergency Data Exchange Language (EDXL),<sup>5</sup> a collection of XML-based messaging standards initiated by the US Department of Homeland Security and developed by the Organization for the Advancement of Structured Information Standards (OASIS). EDXL consists of different components that specify how to exchange data about distributions, resources, hospital availability, situation reporting, and tracking of emergency patients. As the name suggests, EDXL has been designed to speed up direct information exchange between the different actors on-site. EDXL hence targets the direct communication between actors in immediate response situations—e.g., when an ambulance quickly needs to find a hospital nearby that can accept a patient—, whereas HXL focuses on standardizing and streamlining data reporting in long-term humanitarian operations.

The International Aid Transparency Initiative<sup>6</sup> (IATI) is an effort that does not target the operational side of humanitarian response, but fosters transparency of the spendings for development aid. As the transition from disaster response to development aid is often fluent and many organizations work in both domains, there is an overlap between the resources covered in the IATI registry<sup>7</sup> and those addressed in HXL. Cross-referencing both data sources would hence be a valuable next step, which is already being discussed between both groups. Having the IATI data as Linked Open Data—which has already been proposed [9]—would facilitate this step; so far, the data published in the IATI XML format are available through an implementation of the CKAN API.<sup>8</sup>

<sup>4</sup>See <http://cod.humanitarianresponse.info/>.

<sup>5</sup>See <http://docs.oasis-open.org/emergency/>.

<sup>6</sup>See <http://www.aidtransparency.net>

<sup>7</sup>See <http://www.iatiregistry.org>.

<sup>8</sup>See <http://docs.ckan.org/en/latest/api.html>.

While the use of Semantic Web technology for disaster management has been discussed in the literature [6], the proposal for *Crowdsourced Linked Open Data* [26] is the only actual application to date. The authors show how information about the Haiti earthquake collected on the Ushahidi platform [20] can be organized and made available online, using the Management Of A Crisis (MOAC) vocabulary developed for this purpose.<sup>9</sup> Within W3C, synergies between these different efforts are being discussed in the recently established Emergency Information Community Group.<sup>10</sup>

### 3. Requirements for HXL

This section outlines the peculiarities of data exchange in the humanitarian domain, followed by a discussion of the requirements that were identified for the development of HXL.

#### 3.1. Data Exchange in the Humanitarian Ecosystem

Humanitarian data exists at many scales, from global overviews down to highly granular operational data. An example of the former might be the number of people affected by flooding globally in 2012; an example of the latter might be the number of families provided with sanitation supplies on a given day in a given refugee camp by a given organization. HXL is primarily concerned with this granular, operational-level data. Whereas the more aggregated, global level data benefits from a luxury of time in which to produce and validate it, the operational data has no such luxury. To be useful to humanitarian actors, the data must be compiled, reconciled, validated, analyzed and disseminated within hours or days. Changes, additions, and updates are relentless and take place in a high pressure, high stakes environment. As an additional complication, in large scale disasters many of the responding organizations will be national actors or small local NGOs who have never had contact with the international humanitarian system. They will likely have not had exposure to or training on existing information management tools and standards. Instead, they will have their own established, or perhaps ad hoc, information systems. It is this multitude of systems that

drives the complexity and duplication of humanitarian information during a crisis response.

The demand for humanitarian data comes from many levels: from small local actors trying to plan their response activities for the coming days all the way up to donor governments trying to marshal financial or other resources for the response. A responding organization may find requests for its operational data (or some aggregate of it) coming from its national headquarters, from the national government, from OCHA, from the media, and from one or more major donors. These requests generally do not share a common format. This “reporting burden” is a key reason why the creation of additional reporting systems meets with limited adoption.

The telecommunications environment of large-scale disasters is also a key element to be considered in the development of any information management systems. Until recently, these humanitarian operations often took place in regions that had very poor telecommunications infrastructure even in normal times. During disasters, the reliability of electricity and communications is uncertain. In particular, the lack of reliable Internet connectivity has been a major constraint in the development of humanitarian information systems. However, in the last several years, the increasing coverage of mobile networks as well as the availability of satellite-based Internet connectivity has lessened this constraint and opened a door for the development of Internet-based systems for sharing humanitarian operations data.

The international humanitarian system is organized into “clusters” such as Water and Sanitation, Health, Shelter, and Education, among others. Members of these clusters are humanitarian organizations that are responding in one (or more) of these thematic areas. Data from each organization flows to the “cluster-lead” organization, which has the responsibility for compiling that data into a common operational picture for that cluster. However, because of the myriad of semantic and syntactic differences among the various organizations’ data, this compilation task often exceeds the information management capacity available to the clusters.

#### 3.2. Requirements Specification

The primary requirement for HXL is that it address the fundamental information management problem described above: HXL must make the compilation of a common operational picture easier and more timely.

<sup>9</sup>See <http://observedchange.com/moac/ns/>.

<sup>10</sup>See <http://www.w3.org/community/emergency/>.

The identification of this problem comes from the experiences of information managers who have worked in multiple emergencies over the last several years.

In solving this fundamental problem, any proposed solution must not significantly increase the reporting burden already imposed on humanitarian actors and ideally should reduce it by making it possible for organizations to report data once in a way that serves the diversity of users, from operational partners to analysts at the global level. Furthermore, to be successful, any proposed solution should not require the replacement of existing information management systems, but rather focus on interoperability between existing systems. A standard way of describing and encoding operational humanitarian data can achieve this interoperability, however a standard alone is not adequate to solve the problem. During a crisis, there is not sufficient time for organizations who have not been part of an international humanitarian response to integrate a data standard into their operations; indeed, for many small actors there would not be resources for such integration work. Instead, any solution must include not only a data standard, but also a suite of tools allowing easy translation from the most common information management systems (spreadsheets and relational databases) to the data standard. Additionally, to encourage participation in the data standard, this suite of tools should also be able to quickly produce some feedback in the form of data visualization and/or maps.

In the chaotic environment of an international crisis response, there are often conflicting data on humanitarian needs as well as operational data. These problems are currently handled, albeit with great effort, at the cluster-level with some support from OCHA for certain types of data. A solution to the stated problem must replicate and improve the efficiency of this cluster-level reconciliation process before data are published. Efficiency will be also be greatly improved if differences among standard reference information, such as placenames, can also be resolved. A proposed solution should include the ability to serve out standard reference lists that can be ingested into existing information management systems in a variety of ways, from simple spreadsheet-based gazetteers to standards-based web services for geodata.

After analyzing these requirements, a solution based on Semantic Web technologies, following the Linked Open Data [2] paradigm, was identified as the most promising solution. An RDF vocabulary for the domain provides a sound definition of the most im-

portant domain concepts, and any datasets annotated with the corresponding classes and properties are self-describing, as the respective terms can always be looked up online. The SPARQL query language [13] provides a standardized API, so that there is no need to define a new, proprietary API specific to this standard. The double function of URIs as both identifiers and locators for information about a resource allows for easy sharing and lookup of IDs for commonly required resources, such as organizations, emergencies, or geographic features. For geographic information about administrative boundaries or locations of refugee camps, the Open Geospatial Consortium's Simple Features Model and the corresponding ontology already provide a useful standard [23] and extension to SPARQL for geospatial queries [24].

The first step towards HXL was the definition of the vocabulary, which is described in the following section.

## 4. Vocabulary

This section introduces the Humanitarian eXchange Language vocabulary.<sup>11</sup> We first outline the current coverage of the vocabulary, followed by a discussion of the underlying design principles and an overview of the development process.

### 4.1. Coverage of the HXL Vocabulary

The HXL vocabulary—officially entitled *Humanitarian eXchange Language (HXL) Situation and Response Standard*—has been developed to annotate humanitarian data in the absence of properties in established vocabularies that are detailed enough to meet the requirements outlined in Section 3.2. While the development of the MOAC [26] vocabulary was a first step into the right direction and some generic vocabularies provide useful classes and properties that can (and should) be reused, it was apparent that a substantial number of classes and properties have to be introduced to ensure a meaningful annotation of the data at hand. The current version of the vocabulary consists of five sections that organize the vocabulary by topic:

1. **Geolocation section.** This section provides the classes and properties to annotate geographic

<sup>11</sup>The latest version of the vocabulary is available from <http://hxl.humanitarianresponse.info/ns/>.

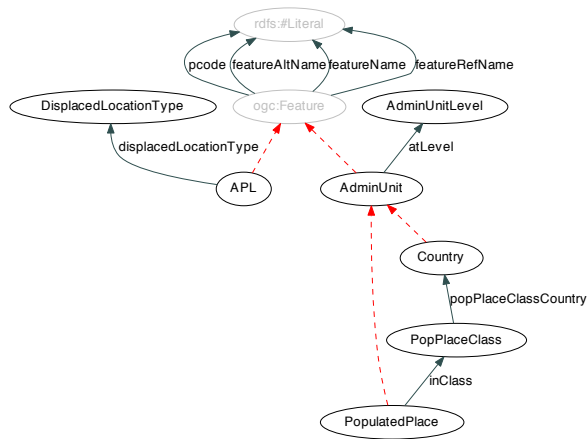


Fig. 2. Overview of the core classes and properties of the geolocation section. Subclass relationships are shown as red dashed arrows.

information such as the common operational datasets. It builds on the Open Geospatial Consortium's Simple Features model [23] and extends the corresponding ontology.<sup>12</sup> This approach ensures that all HXL data is fully compliant with the GeoSPARQL recommendation [24, 1] and hence support complex spatial queries in a standardized way. HXL extends the Simple Features model by the classes and properties required to model the administrative hierarchy in a country, such as `hxl:AdminUnit`. Figure 2 gives an overview of this section of the vocabulary.

2. **Humanitarian profile section.** This section defines the classes and properties required to publish data about the populations affected by an emergency. The classes in this section correspond to the humanitarian profile (see Section 2), breaking down the person counts by the way in which the corresponding populations are affected (Casualty, Missing, Displaced, etc.). For each of these subclasses of Population, the respective properties for `personCount`, `ageGroup`, `sexCategory`, etc. are provided. Figure 3 gives an overview of this section of the vocabulary.
3. **Metadata section.** HXL makes extensive use of named graphs for data management (see Section 7 for details). This section defines the classes and properties to annotate the named graphs,

which are declared as instances of the class `DataContainer` in HXL. To track the provenance data relevant in the humanitarian context, each data container has information attached about its `Source`, what person or organization it has been approvedBy, as well as timestamps for reporting and validity dates.

4. **Response section.** This section contains the classes and properties to describe the organizations involved in response activities coordinated by OCHA, including name, abbreviation, and internal ID. In the next iteration of the vocabulary, this section will be extended with classes and properties to describe actual response activities such as food distributions, vaccinations, etc.
5. **Situation section.** Similar to the response section, the situation section is a stub that will be extended in the future. Its main purpose is currently to enable the annotation of emergencies in HXL with Global Identifier numbers (GLIDE).<sup>13</sup> GLIDE numbers are unique identifiers assigned to all disaster events that meet the criteria of the EM-DAT disaster database<sup>14</sup> and commonly used across the humanitarian domain. The situation section will ultimately contain vocabulary for describing the situation requiring a humanitarian response, including information about needs generated by the crisis (for shelter, water, protection of minors, etc.) and events (security incidents, damage reports, etc.) that shape the response environment.

#### 4.2. Design Principles

Reuse of existing vocabularies is a central principle to ensure interoperability on the Web of Data. However, in the context of HXL as an effort of a United Nations agency, special care had to be taken to make sure that

1. any existing vocabularies used are stable,
2. the definitions of concepts correspond *exactly* to those used in the UN context, and
3. the vocabulary reflects the jargon commonly used in the domain to facilitate adoption.

These points have constrained the number of potential vocabularies to reuse considerably. A class such as

<sup>13</sup>See <http://www.glidenumbers.net/>.

<sup>14</sup>See <http://emdat.be/frequently-asked-questions#FAQ3> for the criteria.

<sup>12</sup>See <http://www.opengis.net/ont/geosparql>.

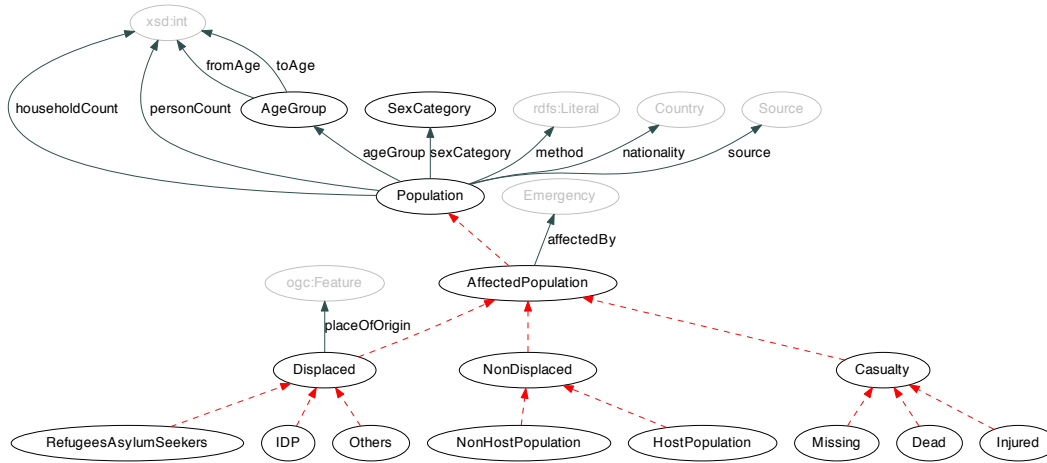


Fig. 3. Overview of the core classes and properties of the humanitarian profile section. Subclass relationships are shown as red dashed arrows.

`hxl:Country` could have been taken from existing vocabularies such as the DBpedia ontology [4], however, we could not find any definition that includes nations along with dependent territories and other special cases. Those may be considered countries in the humanitarian context—however, this does not imply any endorsement or recognition by the United Nations. For classes and properties such as `hxl:NonDisplaced` or `hxl:householdCount`, there were no existing properties at all. For those reasons, the current version of HXL reuses only the Friend of a Friend (FOAF) vocabulary [7], Dublin Core [11], and the OGC GeoSPARQL ontology [24].

From the technical perspective, the vocabulary has been divided into thematic sections introduced in the previous subsection, using the `rdfs:isDefinedBy` property. This makes the vocabulary more tractable and allows us to automatically generate a well-structured documentation for the vocabulary.<sup>15</sup> All HXL classes are subclasses of an abstract `hxl:BaseClass`. Properties are generally defined with domain and range, and “mandatory” properties are marked using an `owl:minCardinality` of 1. These restrictions are especially useful when developing tools such as the HXLator (see Section 5.1) and to validate whether a HXL dataset is complete with respect to those recommended properties.

<sup>15</sup>The scripts generating the documentation at <http://http://hxl.humanitarianresponse.info/ns/> are available from <http://github.com/hxl-team/HXL-Vocab/>.

### 4.3. Development Process

As mentioned in Section 4.2, reuse of terms from the humanitarian field was a mandatory requirement for the development of HXL. The first step was hence to collect as many standards documents, spreadsheet templates, guidelines, and API documentations as possible. From this set of documents, frequently recurring terms were identified and grouped thematically. These clusters were then arranged as concept maps as a first draft of the class hierarchy, which was discussed in several face-to-face meetings with domain experts for verification. At this point, it turned out that a graph-based data model was a considerable mental leap for the domain experts used to the relational and table-oriented models commonly used in the domain. This resulted in a number of presentations on the Semantic Web, RDF, and Linked Data, which set the stage for further discussions with a better understanding of the technology.

When turning the concept maps into an RDF vocabulary, it quickly became apparent that the scope is too broad for a first version, as the documents that had been taken into account cover tasks as different as needs assessment, financial tracking, and response planning. We therefore took a step back and decided to limit the vocabulary to a specific task for the time being, for which we chose the humanitarian profile. From this first scaled-down version, we have since gone through more than 30 revisions, the latest of which is always available from

<http://hxl.humanitarianresponse.info/ns/>. These frequent updates result from a *fail early, fail often* approach to ontology engineering in which we re-evaluate the vocabulary whenever we produce, translate, or consume HXL data and find bugs or missing classes and properties.

## 5. Data Generation

This section introduces the two current ways of producing HXL data: The HXLator (Section 5.1), an interactive tool that guides the user through the process of translating spreadsheets to HXL; and custom-build system crosswalks (Section 5.2) that publish data from existing information systems as HXL data.

### 5.1. HXLator

A significant share of the data exchanged in the humanitarian space lives in ad-hoc spreadsheets. In the field, spreadsheets are used to keep track of refugee counts, distribution activities (e.g. of food or shelter kits), needs assessments, etc. Spreadsheet templates have been developed for common use cases that recur in many emergencies. Translating the contents of these spreadsheets to HXL is a straightforward task for someone with a decent background in programming and Semantic Web technologies. Unfortunately, this skill set is usually not covered by the field staff, so that the numerous tools that already exist for this use case (e.g. [17,12,19]) cannot be applied. This lack of an easy-to-use tool to *extract* data from a spreadsheet, *transform* the data to RDF (i.e., HXL), and *load* them into a triple store (ETL process) made us develop our own custom solution.

The HXLator is an open source<sup>16</sup> online tool that guides the user through the process of converting the data in a spreadsheet to HXL. It has been developed in PHP, based on the PHPEXcel<sup>17</sup> and EasyRDF<sup>18</sup> libraries and the Bootstrap<sup>19</sup> framework for the frontend. The main challenge for the HXL use case was to make the process easy enough so that someone with decent skills in Excel, but no knowledge about Semantic Web technologies, could use it. The HXLator was hence developed with a strong focus on the work flow

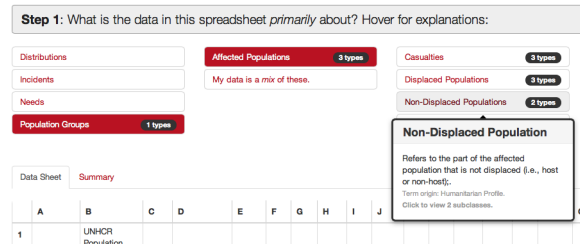


Fig. 4. Class selection in HXLator. When a class is selected, its subclasses expand to the right. Hovering over one of the class buttons brings up the class's definition from the HXL vocabulary.

and hiding most of the complexity of the underlying technology from the user.

After logging in, the user starts by entering some metadata, selecting the emergency the data is about, the report category (currently only supporting *humanitarian profile*), a validity date for the uploaded data, and the file to upload. Optionally, she can reuse an existing translator: once the translation process is complete, HXLator stores the translator the user has created, so that it can be re-applied to a new (or updated) spreadsheet of the same structure. This is especially useful for the commonly used templates mentioned above, as the translators for those can already be provided in HXLator, saving the user the effort to create them.

Once the spreadsheet has been uploaded to the server, it is shown back to the user, who is now guided through a five step process to generate a translator for the file:

- **Step 1:** HXLator shows the classes defined in HXL, organized based on the subclass hierarchy (see Figure 4). Here, the user has to select the class her spreadsheet has data about. This step defines the properties available for mapping, using the domain and range definitions provided in the HXL vocabulary and inferred via subclass reasoning.
- **Step 2:** HXLator asks the user to select the first row in the spreadsheet that contains actual data (not the header row). This row acts as a template and is used for the actual mapping process.
- **Step 3a:** The user selects a cell in this row that identifies an instance of the class selected in Step 1; e.g., a population of refugees and asylum seekers.
- **Step 3b:** The user selects one of the properties available for the selected class. At this point,

<sup>16</sup>See <http://github.com/hxl-team/HXLator/>.

<sup>17</sup>See <http://phpexcel.codeplex.com>.

<sup>18</sup>See <http://www.easyrdf.org>.

<sup>19</sup>See <http://twitter.github.com/bootstrap/>.



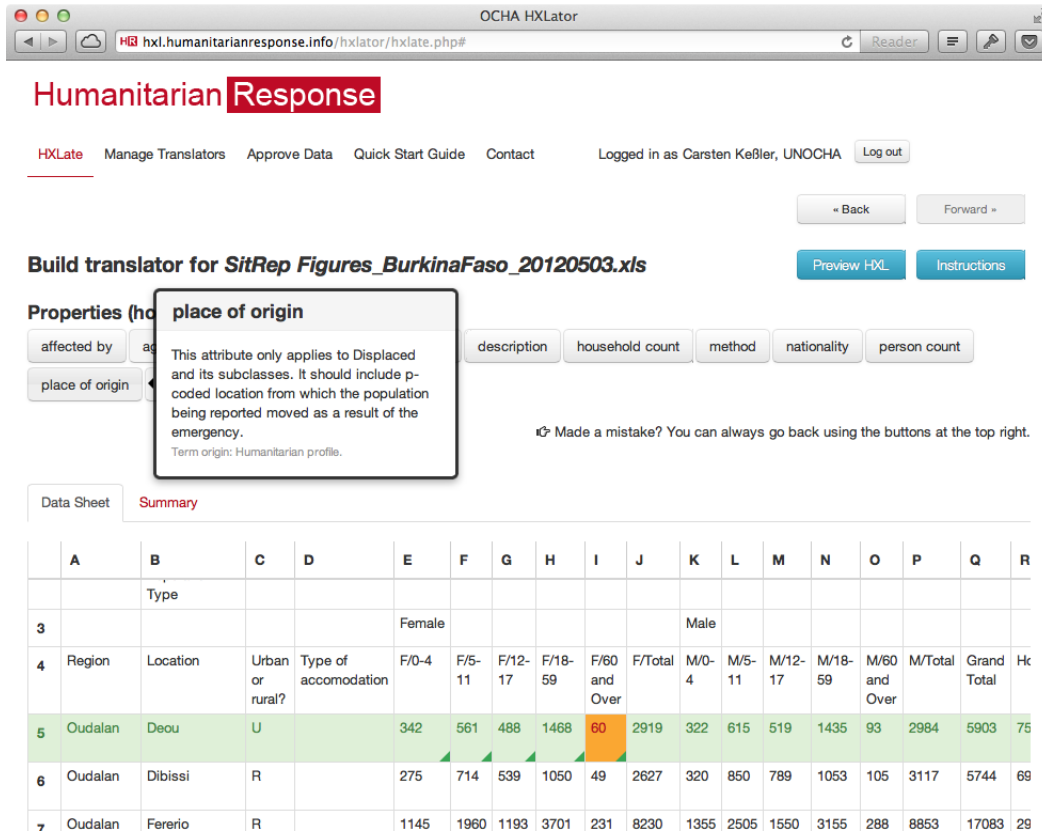


Fig. 5. The HXLator user interface.

HXLator distinguishes between data properties and object properties:

- \* For data properties, HXLator asks the user to either select the value by clicking the respective cell in the spreadsheet, or by directly typing it in.
- \* For object properties, HXLator allows the user to perform lookups for existing resources on the HXL triple store. Again, these can be selected from the spreadsheet or typed in, and the user has to confirm the correct resource to link to (e.g., the URI of camp that the spreadsheet only contains the name of).

Steps 3a and 3b are repeated until all cells have been mapped. During this process, all mapped cells and properties are highlighted with green marks (See Figure 5).

- **Step 4:** When the user continues to the next step, she is warned about any properties that have not been mapped. Even though it is unlikely that a

spreadsheet will contain data for each of the HXL properties available for the selected class, this step is to make sure the user does not miss any relevant data. Once this is confirmed, the user is asked to select the rows that will be translated in the final step. In case the user has used an existing translator, she will directly jump to this step, i.e., skip the actual mapping process and only select the rows to translate.

- **Step 5:** After selecting the rows, the translated data is shown to the user for a final check, before they are submitted to a protected triple store for review (see Section 7).

HXLator contains several functionalities to make this process as straight-forward and efficient as possible, such as the selection of multiple cells that share the same property (e.g., if several cells identify female populations). The preview function mentioned above (see Figure 6) is available for the user throughout the process, so that she can easily check what data has been mapped yet, and whether the mapping is correct

The screenshot shows a window titled "HXL Preview" with a search bar and a table. The table has columns for "ageGroup", "sexCategory", and "personCount". The data rows are as follows:

	ageGroup	sexCategory	personCount
1-E-5	Ages 0 to 4	female	342
1-F-5	Ages 5 to 11	female	561
1-G-5	Ages 12 to 17	female	488
1-H-5	female	female	1468
1-I-5	female	female	60

Fig. 6. In HXLator, the user can access a preview of the generated HXL at all times, both in tabular form and as RDF code in Turtle notation.

so far. If an error is spotted, the user can always go back, as HXLator keeps track of all changes to the translator. The translators are currently implemented as JSON objects, as the whole mapping process happens on the client side, in JavaScript, and JSON is easiest to parse and serialize in this environment. We are considering a switch to R2RML [8] as a more standardized solution that also allows the export of the generated translators to other tools in the future.

## 5.2. System Crosswalks

While HXLator has been developed to map the data from ad-hoc spreadsheets, often filled in on-site, large humanitarian organizations—both within the UN system and international NGOs—already have a wide range of relational databases and information systems in use. As these provide far more predictable and well-structured data, setting up ETL processes for these is more straight-forward. Moreover, the underlying structure (database schema or API) hardly changes, so that it is usually sufficient to have an expert work out a crosswalk to generate HXL from these systems once, and it will work as long as there are no major changes made to the source system.

The first crosswalks developed was the translation of the geographic information contained in the common operational datasets (CODs). The CODs contain GIS shape files for each country OCHA currently has missions in,<sup>20</sup> acting as the geographic reference data for UN agencies in the respective area. These files

are central to OCHA’s operations, as most other data is directly or indirectly tied to the locations referenced here. Moreover, there was already a system in place that assigned each place a unique ID, a so-called *p-code* that reflects the administrative hierarchy of the feature [14]. Using existing libraries for handling shape files,<sup>21</sup> translating the shape files to RDF using the HXL vocabulary was straight-forward.<sup>22</sup> Since HXL builds on the OGC simple features model [23], HXL data is ready to be queried via GeoSPARQL [24].

The first relational database exposed through HXL is a system for refugee numbers based on their current location, origin, and crisis that affects them. A first version of a custom crosswalk to this database maintained by the United Nations Refugee Organization (UNHCR) was developed from scratch to explore the difficulties that such a translation might bear. It turned out that the translation could only be completed with fairly massive manual intervention, as the database did not make use of p-codes, for example, but only used the place names. In many cases, there are a number of potential spelling alternatives for a place name (and places with the same or very similar names), which cannot be automatically resolved using string distance measures such as the Levenshtein distance [18]. In order to facilitate the translation and streamline the data management practices between UN agencies in general, we aim for a solution that encourages the use of unique identifiers such as the p-codes or GLIDE numbers across different agencies. This will also facilitate the move from translating dumps of the database,<sup>23</sup> as in the current prototype, to “live” exposure of such databases in the future, via tools such as D2RQ [5,3].

With a growing number of tools and scripts producing HXL data, agreed-upon URI patterns [10] gain in importance. A standard pattern per HXL class is especially important to make sure that the same real-world entities are always represented by the same URI, independent of the system that produces the current data at hand. Patterns such as

`http://hxl.humanitarianresponse.info/data/locations/admin/country-code/p-code`

<sup>21</sup>See <http://www.gdal.org/ogr2ogr.html>.

<sup>22</sup>See <http://hxl.humanitarianresponse.info/data/locations/admin/bfa/BFA050> for an example.

<sup>23</sup>See <http://hxl.humanitarianresponse.info/data/datacontainers/1356597368.2703> for a sample data container generated from the UNCHR database.

<sup>20</sup>See <http://cod.humanitarianresponse.info/country-region/afghanistan> for an example.

are collected in a shared document<sup>24</sup> and implemented both in the HXLator, as well as the crosswalk tools.

## 6. Data Consumption

The HXL project is driven by the need to be able to re-purpose data into products that support the work of many users: humanitarian actors need a solid basis for planning their activities, donors want to prioritize projects to allocate funds, media need up-to-date information for articles, and academics for scientific studies – hence the need for open, machine readable data. In the following, we introduce two use cases that already leverage the HXL data and discuss how they improve the situation for humanitarian actors.

### 6.1. Dashboards

The UN is a large organization with many operational and administrative layers. The information collected by information management officers in the field needs to support decision-making at many of these levels. It is often repackaged into a variety of infographics or other types of reports for different audiences. Many of these data visualization products are put together manually in a process that could easily last several days. Obviously, the situation on-site often already differs from the numbers that the decision makers are looking at because of this lengthy process.

The goal was hence to develop a dashboard that is completely driven by HXL data and that can be easily configured and set up for a new emergency. Once up and running, the dashboard should require no manual work and always display the latest data from the HXL triple store. Figure 7 shows a prototype for such a dashboard that was developed during a hackathon at OCHA in Geneva last November. The dashboard was conceptualized as an empty frame that is “put to life” with data dynamically loaded via AJAX from the HXL triple store and an instance of the Humanitarian Response platform. Besides the fact that it is always up to date with the latest data from the triple store, it also allows the decision makers to explore the data in detail, for example through the interactive charts and the mapping module.

### 6.2. HXL Geo Web Services

The HXL triple store also contains (for those countries used in the prototype HXL development) geographic reference data provided by OCHA, as mentioned in Section 5.2. While we are sure that having the different features available as Linked Data will be beneficial in the long run, there are not many software solutions available yet that can take advantage of this offering. Instead, most GIS systems and Web mapping frameworks build on the geo web service specifications developed by the Open Geospatial Consortium, most importantly the Web Map Service (WMS) for pre-rendered map images, and the Web Feature Service (WFS) for vector data [25,22].

In order to make the data available in an easily digestible format for GIS analysts and Web mapping applications, a service chain was set up that publishes new geographic reference data through a suite of OGC services provided by an ArcGIS server instance hosted by an OCHA partner.<sup>25</sup> Since the reference data only change very infrequently, and changes are communicated to the partners beforehand to make sure everyone is always using the same reference data, it was not necessary to develop a “live” mapping that wraps the triple store as a WFS/WMS. Instead, we have implemented a pull-based solution that checks the triple store every night for changes to the geographic information.<sup>26</sup> If any of the data should have changed, a process is triggered that generates an `INSERT` request to the transactional WFS. The ArcGIS server instance hosting the WFS then automatically creates a WMS based on the data.

Both services are available to the whole community and already in use on the dashboard shown in Figure 7, for which the WMS delivers the base map. The beauty in this solution is that the geo web service infrastructure is automatically synced with the triple store on a daily basis, where the changed shape files would have required manual updates to the ArcGIS server instance before. Moreover, this approach easily supports setting up additional server instances that can be run e.g. on-site for the local staff in an emergency with poor Internet connectivity, acting as a “Spatial Information Infrastructure [21] on a USB stick”.

<sup>24</sup>See <http://goo.gl/kGnK4>.

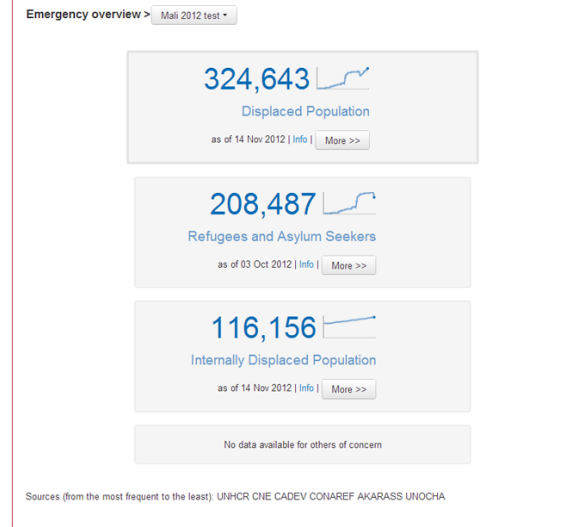
<sup>25</sup>The GIST, hosted by the University of Georgia; see <https://gistdata.itos.uga.edu/node/5>.

<sup>26</sup>See <http://github.com/hxl-team/HXL2WFS>.

## Humanitarian Response

Humanitarian Profile Data Browser

Note: This is a test setup and some of the data shown here may be inaccurate, outdated, or even entirely made up.



## Humanitarian Response

Humanitarian Profile Data Browser

Note: This is a test setup and some of the data shown here may be inaccurate, outdated, or even entirely made up.

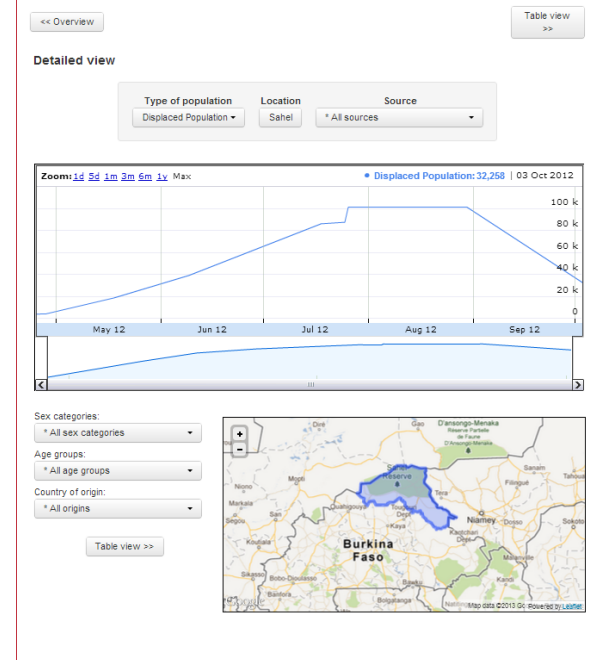


Fig. 7. A prototype dashboard based on HXL data, available from <http://hxl.humanitarianresponse.info/dashboard>. The left side shows the initial page that allows the user to select the emergency and provides the key numbers. The right side shows the detailed view for an emergency, where numbers can be broken down by age or sex group, for example.

## 7. Data Management

A data management approach that is designed around existing structures within OCHA and the humanitarian sector as a whole was a fundamental requirement during the development of HXL. This section discusses the governance and workflow requirements and how they were implemented in the standard itself and the surrounding tools.

### 7.1. Governance and Workflows

Data collected by OCHA and its partner organizations are the basis of decisionmaking for the international community, especially concerning the allocation of resources in the aftermath of a disaster. Wrong or incomplete data can easily lead to insufficient political actions, putting lives in the affected areas at risk. Within the humanitarian domain, erroneous data can lead to a wrong focus in planning the response activities, or even put field staff in dangerous situations. It is therefore of utmost importance that all data published

through HXL go through a review by staff members who are familiar with the overall situation in a specific crisis. Most data compilation already happens at the cluster-lead level, which has a broader scope than the information management officers on a specific site. At the same time, the cluster-lead staff is still highly familiar with the situation, so that any obviously wrong data will immediately catch their eye.

The overall workflow for HXL is hence adopted from existing practice within the humanitarian domain (see Figure 8): Any data collected in the field is compiled and translated to HXL by the information management officers before it propagates to the headquarters in Geneva and New York through the respective cluster lead. This approach also applies the *many eyes principle* to reduce the number of potential errors.

### 7.2. Implementation

The workflow introduced in the previous section has been implemented at two different levels: In the server setup, and in the HXL vocabulary.

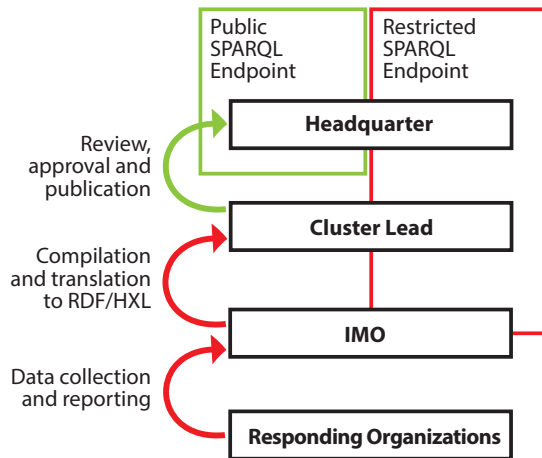


Fig. 8. Overview of the data propagation process.

The metadata section (see Section 4.1) of the HXL vocabulary is designed around the principle of using named graphs for data management. In HXL, a named graph is an instance of the class *DataContainer*. The rationale behind this naming choice is that it should also be possible to transfer a data container offline as a file, e.g. on a USB stick, when OCHA staff works in situations without online access. A data container is self-describing in that it contains all corresponding metadata. The corresponding triples are automatically generated by the data generation and approval tools (see Section 5). In the data management context, the following properties are particularly important, as they contain the reporting and approval provenance data:

- `aboutEmergency`: Links the data container to the corresponding emergency it contains data about.
- `reportCategory`: Links the data container to a specific report category such as *humanitarian profile*, *shelter* or *security*. The report category is used to determine who is responsible for reviewing the data before publication.
- `reportedBy`: Metadata property that links the data container to the person or organization that has reported the contained data.
- `approvedBy`: Metadata property that links the data container to the person or organization that has approved the contained data for publication.

As indicated in Figure 8, the HXL server setup hosts two separate triple stores (both running Fuseki<sup>27</sup>): A public endpoint<sup>28</sup> that can be queried without restrictions, and a second endpoint for unapproved data with restricted access. We refer to this protected triple store as our *incubator store*. Any data coming in from the field, such as from the HXLator, are staged for review in the incubator store. The corresponding cluster lead is notified that new data is waiting to be reviewed, and can look at the data in a web frontend. In case of suspicious data, such as near-duplicates (which may point to double reporting) or extreme outliers, the reviewer needs to find out whether there has been a reporting error, or whether the situation in the field has really changed dramatically. This is usually achieved by calling the colleagues in the field. In the data publication step, i.e., when a data container gets approved at the cluster-lead level, the corresponding triples are moved from the restricted triple store to the public one.

## 8. Conclusions

The Humanitarian eXchange Language is an emergent standard for operational data in the humanitarian domain. It is the foundation for an infrastructure that makes strong use of Semantic Web technology in producing, maintaining, and using the data which form the basis for the planning of humanitarian activities in the international community. In this paper, we have reported on the first steps to making HXL a central reference point for the whole domain. The peculiarities of this field are reflected in a set of requirements that HXL needs to address concerning the structure and existing practices in the humanitarian ecosystem. The HXL vocabulary formalizes established terminology from the domain. It currently focuses on humanitarian profile data as well as core reference data, such as geographic information. Data according to the HXL vocabulary can currently be produced either using the HXLator, an interactive tool to translate spreadsheets to HXL, or via crosswalks that produce HXL from existing information systems. HXL data already drive the first applications, including emergency dashboards and a Web service infrastructure for geographic information. An initial data governance system has been set

<sup>27</sup>See [http://jena.apache.org/documentation/serving\\_data](http://jena.apache.org/documentation/serving_data).

<sup>28</sup>See <http://hxl.humanitarianresponse.info/sparql>.

up to ensure that any data coming in from the field are approved at the cluster lead level before publication.

The current version of the vocabulary and the corresponding tools are at a prototyping stage to demo the capabilities of HXL within the UN system and to outside partners to foster adoption. The next steps include the extension of the approval chain with consistency checks and automatic highlighting for data that drastically diverge from existing data. This could either point to significant changes in the situation in a camp, for example, or it could point to a reporting error; resolving such issues will require the expertise of an expert familiar with the situation on the ground. The main task will hence be to build a system that reliably identifies potential problems in the data, and provides an easy-to-use interface to resolve them. The HXL vocabulary will be gradually extended as required, depending on the next reference datasets to be included.

Defining the HXL vocabulary for the humanitarian system as a whole clearly goes beyond the capabilities and expertise of OCHA. In order to achieve this goal, the involvement of the global clusters in developing their respective components, such as vocabulary extensions and cluster-specific tools, is required. Moreover, the volunteer and technical community needs to be included. Collaboration with this community has already been established in a *random hacks of kindness* event at the international crisis mappers conference in Washington, DC in fall 2012. Future involvement should also address the development of models for crowd-sourced data, which is not yet covered in HXL. The underlying technology, however, has the potential to vastly improve the integration of official agency data information and data collected by the volunteer community.

## References

- [1] Robert Battle and Dave Kolas. Enabling the Geospatial Semantic Web with Parliament and GeoSPARQL. *Semantic Web*, 3(4):355–370, 2012.
- [2] Tim Berners-Lee. Linked Data – Design Issues. Available from <http://www.w3.org/DesignIssues/LinkedData.html>, 2009.
- [3] Christian Bizer and Richard Cyganiak. D2R server—publishing relational databases on the semantic web. In *Proceedings of the 5th international Semantic Web Conference (ISWC2006)*, 2006.
- [4] Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. DBpedia—A crystallization point for the Web of Data. *Web Semantics: Science, Services and Agents on the World Wide Web*, 7(3):154–165, 2009.
- [5] Christian Bizer and Andy Seaborne. D2RQ—treating non-RDF databases as virtual RDF graphs. In *Proceedings of the 3rd International Semantic Web Conference (ISWC2004)*, 2004.
- [6] Eva Blomqvist. The Use of Semantic Web Technologies for Decision Support – A Survey. *Semantic Web Journal*, accepted.
- [7] Dan Brickley and Libby Miller. FOAF Vocabulary Specification 0.98. Available from <http://xmlns.com/foaf/spec/20100809.html>, 2010.
- [8] Souripriya Das, Seema Sundara, and Richard Cyganiak. R2RML: RDB to RDF Mapping Language. W3C recommendation available from <http://www.w3.org/TR/r2rml/>, 2012.
- [9] Tim Davies. IATI Linked Data. Discussion paper available from <http://goo.gl/Xc8s4>, 2012.
- [10] Leigh Dodds and Ian Davis. *Linked Data Patterns – A pattern catalogue for modelling, publishing, and consuming Linked Data*. 2012.
- [11] Dublin Core Metadata Initiative. DCMI Metadata Terms. Available from <http://dublincore.org/documents/dcmi-terms/>, 2012.
- [12] Lushan Han, Tim Finin, Cynthia Parr, Joel Sachs, and Anupam Joshi. Rdf123: a mechanism to transform spreadsheets to rdf. In *Proceedings of the Twenty-First National Conference on Artificial Intelligence (AAAI 2006)*. AAAI Press, Menlo Park, 2006.
- [13] Steve Harris and Andy Seaborne. SPARQL 1.1 Query Language. W3C recommendation available from <http://www.w3.org/TR/sparql11-query/>, 2012.
- [14] Chad Hendrix. Field Guide for the Use of Geo-Codes. Available from <http://goo.gl/Ya7fh>, 2011.
- [15] Inter Agency Standing Committee. IASC Guidelines Common Operational Datasets (CODs) in Disaster Preparedness and Response. Available from <http://cod.humanitarianresponse.info/node/53>, 2010.
- [16] Inter-Agency Standing Committee. Multi-Cluster/Sector Initial Rapid Assessment (MIRA). Available from [http://ochanet.unocha.org/p/Documents/mira\\_final\\_version2012.pdf](http://ochanet.unocha.org/p/Documents/mira_final_version2012.pdf), 2012.
- [17] Andreas Langeegger and Wolfram Wöß. XLWrap—Querying and Integrating Arbitrary Spreadsheets with SPARQL. *The Semantic Web—ISWC 2009*, pages 359–374, 2009.
- [18] V.I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10(8):707–10, 1966.
- [19] Fadi Maali, Richard Cyganiak, and Vassilios Peristeras. Re-using cool URIs: Entity reconciliation against LOD hubs. In *Proceedings of the Linked Data on the Web Workshop*, volume 3, 2011.
- [20] Ory Okolloh. Ushahidi, or ‘testimony’: Web 2.0 tools for crowdsourcing crisis information. *Participatory Learning and Action*, 59(1):65–70, 2009.
- [21] Harlan Onsrud, Barbara Poore, Robert Rugg, Richard Taupier, and Lyna Wiggins. The future of the spatial information infrastructure. In Robert B. McMaster and E. Lynn Usery, editors, *A Research Agenda for Geographic Information Science*, pages 225–255. CRC Press, 2005.
- [22] Open Geospatial Consortium. OpenGIS Web Feature Service 2.0 Interface Standard (also ISO 19142). Available from <http://opengis.org/standards/wfs>, 2010.

- [23] Open Geospatial Consortium. OpenGIS Implementation Specification for Geographic information – Simple feature access – Part 1: Common architecture. Available from <http://opengis.org/standards/sfa>, 2011.
- [24] Open Geospatial Consortium. OGC GeoSPARQL – A Geographic Query Language for RDF Data. Available from <http://opengis.org/standards/geosparql>, 2012.
- [25] Open Geospatial Consortium. OpenGIS Web Map Service (WMS) Implementation Specification. Available from <http://opengis.org/standards/wms>, 2012.
- [26] Jens Ortmann, Minu Limbu, Dong Wang, and Tomi Kaupinen. Crowdsourcing Linked Open Data for Disaster Management. In *Proceedings of Terra Cognita 2011, The 10th International Semantic Web Conference (ISWC2011)*, Bonn, Germany, October 2011.
- [27] United Nations Office for the Coordination of Humanitarian Affairs. OCHA Strategic Framework, Objective 2.4. Available from [http://www.unocha.org/ocha2012-13/strategic-plan/objective-2\\_4](http://www.unocha.org/ocha2012-13/strategic-plan/objective-2_4), 2012.