

The Quality of the XML Web

Steven Grijzenhout
University College London
Department of Management Science and Innovation
Gower Street
London WC1E 6BT, United Kingdom
steven.grijzenhout.10@ucl.ac.uk

Maarten Marx
ISLA, Informatics Institute
University of Amsterdam
Science Park 904
1098 XH Amsterdam, The Netherlands
maartenmarx@uva.nl

ABSTRACT

We collect evidence to answer the following question: Is the quality of the XML documents found on the web sufficient to apply XML technology like XQuery, XPath and XSLT? XML collections from the web have been previously studied statistically, but no detailed information about the quality of the XML documents on the web is available to date. We address this shortcoming in this study. We gathered 180K XML documents from the web. Their quality is surprisingly good; 85.4% is well-formed and 99.5% of all specified encodings is correct. Validity needs serious attention. Only 25% of all files contain a reference to a DTD or XSD, of which just one third is actually valid. Errors are studied in detail. Automatic error repair seems promising. Our study is well documented and easily repeatable. This paves the way for a periodic quality assessment of the XML web.

The full paper and all data are publicly available at the url <http://data.politicalmashup.nl/xmlweb>.

Categories and Subject Descriptors

H.m [Information Systems]

General Terms

Measurement, Reliability, Standardization.

Keywords

XML, XML Web, Schemas, Data Quality.

1. INTRODUCTION

In this study we look at the prospects of using XML technology for information extraction and integration tasks on XML data found on the World Wide Web. Without becoming specific we refer to the multitude of these tasks as Extract-Transfer-Load (ETL) tasks [35]. Examples of ETL subtasks are data harvesting, text extraction, structure extraction, text mining [25], data de-duplication [14], data exchange (from one schema to another) [13] and data publishing (from one format (e.g. XML) to another (e.g. RDF or relational)).

Apart from the actual collecting of data from the web, all of these tasks can be expressed in the three XML query languages, XPath 2.0, XSLT 2.0 and XQuery 1.0. Not only can these tasks be expressed in these languages, when the input is XML it is

desirable to do so for a number of reasons. XSLT and XQuery programs are largely declarative. The semantics of the languages is clear and well-defined. The languages are vendor and software independent, developed and maintained by a committed community and became W3C standards. The immense success of SQL shows the great software engineering benefits of working with such programming languages. Maintainability of code is crucial for ETL tasks as they are typically applied in a changing environment not under control of the developers of the ETL code. The LixTo [17] suite of web-extraction tools is built on these principles. A recent addition to the Lixto tools is OXPath [15], an extension of XPath which allows declarative extraction of the deep web.

Whether it is *feasible* to use XML technology for ETL tasks depends on many factors. This is out of the scope of this study. Here we only look whether it is *possible*. That is, is the quality of the XML documents found on the web sufficient to apply XML technology?

Another reason to study the XML web is the new XQIB, *XQuery In the Browser*, initiative (<http://www.xqib.org/>). XQIB is an alternative to JavaScript. Obviously it needs XML of good quality.

Previous studies on HTML showed that the vast majority of HTML documents (around 95%) on the web did not comply with the standards set by the World Wide Web Consortium [11][32][33]. For XML, studies that measure basic quality indicators (like being syntactically correct) on arbitrary XML data from the web have not been performed yet. There are several empirical studies on XML but they either use data from repositories or have very small samples and always contain only well-formed XML (Cf. Section 2).

Unhappy with this omission and frustrated by our own efforts of using XML tools for a large data integration project we set ourselves the following research goal:

Create a corpus of XML documents and accompanying schemas that is representative of the web, evaluate which part is ready to be processed with XML tools, and evaluate the prospects of automatic error correction for the other part. In addition, process, document and store the corpus in such a manner that our study can easily (e.g., yearly) be repeated.

The paper describes the created collection (Section 3), and the evaluation of its quality (Section 4). We also created a corpus of schemas in the three XML schema languages and evaluated their inter-translatibility (Section 4.3). The remainder of this introduction consists of our operationalization of XML-quality and an overview of the main results. Related work is presented in Section 2.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
CIKM'11, October 24–28, 2011, Glasgow, Scotland, UK. Copyright 2011 ACM 978-1-4503-0717-8/11/10...\$10.00.

1.1 Basic quality requirements

One can only apply XML tools to XML files if they satisfy a few basic but important properties. As XML is a self-describing format, these properties all state that files should not lie about themselves concerning some aspect X. We looked at three aspects: a file should not lie about its encoding, it should not state that it is XML when it is not, and it should not lie about its validity with respect to a schema. More precisely,

1. The document should be encoded using a single encoding that is stated in the document.
2. The document should be well formed XML.
3. If the file references a schema, that should be useful and truthful. This means that
 - a. the URI identifying the schema should be resolvable. Also all included schema files should be resolvable recursively;
 - b. all these schemas are syntactically correct, and
 - c. the file is valid with respect to the schema(s).

We collected almost 180K unique XML files from the web from almost 100K websites with a total size of 40GB. We now summarize the main results. Our first result states that encodings do not pose a real problem as 99.5% had a correctly specified encoding. The other results are neatly summarized in Figure 1.

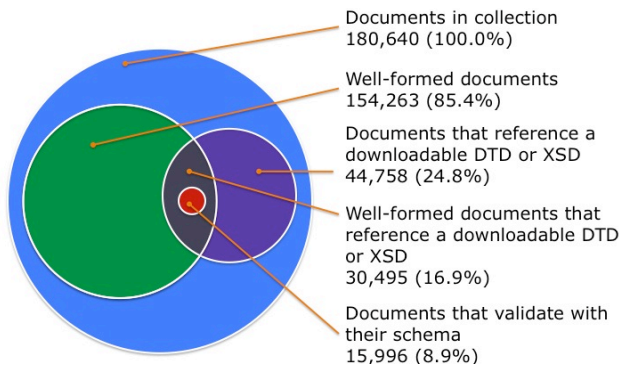


Figure 1: Summary of the Quality of the XML Web.

If we pick a random XML document there is a 14.6% chance that it is not well-formed. This is much better than could be expected from earlier studies on HTML. Interesting is the position of the subset of documents that reference a DTD: 66.4% is not well-formed, almost five times more than on average. It seems that the addition of a DOCTYPE declaration is often added to hide their own poor quality.

Validity is rare on the web. Just over 10% of the well-formed documents are also valid. If we zoom in on validity we see very different patterns for DTD and XML Schema. We go through the three possible problems. The first problem is to reference a schema that cannot be retrieved. This happened in 12.5% of all references to a DTD. Things get subtle once one realizes that DTDs can also include other DTDs, and these have to be retrieved as well. Of the 5410 include statements in DTDs in our corpus, 33% could not be downloaded. Includes in XML Schema behave much better: of 2110 includes only 23 could not be retrieved. XML documents which claim to be valid with respect to an XML Schema behave very well: they have 99% chance of being well-formed, which is much better than the average 85.6%. Figure 2

and Figure 3 show those files that reference a schema but could not be validated. The figures present the causes for non-validation.

The differences are remarkable. Files referencing a DTD are for 73% not valid just because they are not even well-formed. For XML Schemas this cause of error is negligible. Conversely, 31% of all XML Schema validity errors are due to a schema that is useless because it is not even syntactically correct. This happens in only 4% of the DTD-errors.

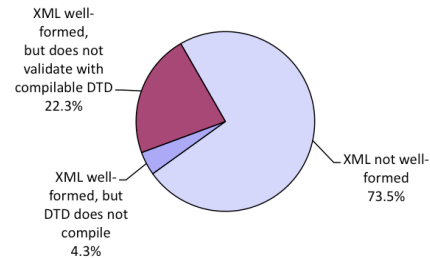


Figure 2: Distribution of causes for non-validation: DTD.

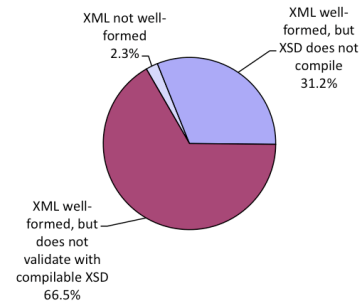


Figure 3: Distribution of causes for non-validation: XSD.

1.2 Main contributions

Our main contribution is an up-to-date and reliable estimate of the quality of the state of the XML web in 2010. Our second contribution consists of an extensive analysis of the type of errors that compromise the quality of the XML web. As they are mostly Pareto distributed we believe this can be used to guide research into automatically repairing errors.

Our third contribution is the collected data itself. All data is made publicly available in a uniform format at the url <http://data.politicalmashup.nl/xmlweb>. All referenced schema files are (recursively) locally available. Information about headers, encodings, and errors of each XML file is stored in a relational database that is also available. Also all scripts and settings for crawling and analyzing the collection are available and well-documented. This makes our study easily repeatable. We hope that this is a start of a longitudinal XML collection. Of course we also hope to see a steady improvement of the quality of the XML web.

Our fourth contribution consists of a corpus of 8000 schema files in the four main schema languages. We analysed their inter schema translatability with James Clark's Trang. (<http://www.thaiopensource.com/relaxng/trang.html>)

2. BACKGROUND

Data quality is a research field established by Wang and Madrick and matured into a field with an own ACM journal [26]. Within the categorization of data quality research described in [26], our study is concerned with assessment and uses an empirical method. Ultimately we study whether XML data from the web is ‘fit for use’ by XML technologies. The idea to measure data quality as data being ‘fit for use’ by data consumers goes back to [39].

A large number of descriptive studies on XML have been conducted. There are three main themes identifiable in the literature, which will be discussed accordingly: studies on XML collections; studies on the quality of the HTML web; and studies on XML schema languages.

2.1 Studies on XML collections

Studies on XML document collections mainly differ in sample data [38]. A study has been done on 200,000 publicly available XML documents from the Xyleme repository [30]. This collection contains only well-formed XML documents. Another study used a number of XML collections, consisting of 16,534 documents and accounting for a total size of 20 Gigabytes [31]. The collections include well-known docbook samples, XML bibles, RDF samples and IMDb collections

2.2 Studies on HTML web quality

Several surveys on the quality of HTML documents on the web exist [32][11][5][33]. Although XML’s predecessor HTML differs greatly in applicability, these studies are relevant because of their approach. The differences in sample collections and quality measures do not make a large difference in results. All indicate a poor quality of the HTML web: a mere 6.5% [5], 5% [10], 4% [33] and 3% [32] of the HTML documents complied with W3C’s HTML standards.

2.3 Studies on XML schema languages

XML schemes (DTD, XML Schema (abbreviated as XSD), Relax NG) have been studied in a number of ways. Firstly, XML schemes are studied in relation to XML collections. As we have seen above, only a small percentage of documents reference a schema [30][31]. As is the case with HTML files, the syntax of most DTD files is incorrect [12][36]. This is generally also the case for XML Schema [8]. Secondly, the properties of XML schemas are studied [12][21][36]. Thirdly, work has been done in developing metrics to measure the properties of DTDs [21] and XML Schemas [29]. These metrics might be interesting to use in future versions of quality analysis. Lastly, research is done in comparing the use of the different XML schema languages. The three languages are incomparable in expressive power and their effect when validating. In our study we look at these differences in expressive power from a pragmatic point of view: how often can schemas be inter-translated using an existing schema translator (Cf. Section 4.3)?

3. DATA

We briefly describe the collections of XML and schema files, and how they were obtained. More detail can be found on the webpage where all data can be downloaded.

3.1 Desired Data

The population of the data in this study is the subset of the web made of XML documents only [3]. The actual amount of files in the XML web is unknown. Obtaining an estimate of its size is intrinsically difficult [1]. The size of the population is, however, irrelevant in calculating a representative sample size.

Unfortunately, collecting XML documents from the web is often not a simple random sample. Because of this, it is not possible to calculate a required sample size. We decided to harvest as many XML documents from the web as possible. The objective of our study is to assess the quality of the XML Web, and a large collection will maximize the probability that errors are included in the collection. Our data collection process does not access the Hidden Web [34]. As a consequence, our collections will not contain any data from the Hidden Web.

3.2 Description of Data

The XML collection contains 180,640 XML files. The total file size of the collection is 40 Gigabytes. The largest file in the collection is 683.7 Megabytes, and the smallest is 1 byte. The average file size is approximately 223 Kilobytes. The number of documents that has a duplicate is only 1296. The URLs in the collection allow us to describe the distribution of XML documents on the web. We clustered our XML files into the zones, consisting of generic Internet domains and geographical regions, defined by Barbosa et al. [3]. Figure 4 shows the partition of our dataset in zones.

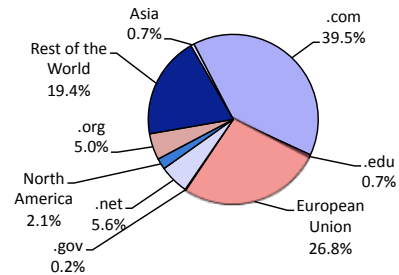


Figure 4. Distribution of web sites by Zone.

The ability to reference a schema is one of the most important features of XML. We focused on DTDs and XSDs because they are used most often, and they are uniformly referenced in XML documents. DTDs were downloaded by extracting the system identifier in the XML document header. All external referenced document includes have been downloaded recursively. Our collection contains 24,426 (13.5%) files with a reference to a DTD. The DTD schemas contained a total of 5410 includes of other DTDs or entity documents. These have been downloaded recursively, and the original schemas have been modified to include the locally downloaded included schemas. 1786 (33.0%) of them failed to download. The thus obtained collection contains 1375 DTDs. The collection contains 24,087 files with a reference to an XSD (13.3%). We could download 437 unique XSDs.

Apart from collecting schema files that were referenced in an XML file, we also collected schema files directly. In this way we could also harvest Relax files. In total we have 3078 DTDs, 4141 XSDs, 338 Relax NGs in XML and 337 Relax NGs in the compact syntax.

3.3 Data Collection

The data were collected in the following 4 steps:

1. Crawl a list of URLs of XML documents from Yahoo and Google,
2. Download the content of each URL,
3. Organize the collection,
4. Determine duplicates.

For each URL, we store the URL, the HTTP header, the actual XML file, and recursively all schema files it references.

The list of URLs was created using a modified version of the crawler from [8]. The crawler executes several keyword queries with the filetype restricted to XML. Only Yahoo and Google can limit search to XML files. The results of these two steps are described in Table 1.

Table 1. Statistics of URL List and Downloading

Filetype	Unique URLs in List	Files Downloaded	Loss Percentage	Last File Downloaded
XML	188,332	180,640	4.08%	2010-07-17

The resulting collection was organized in a MySQL database. For each file, the database stores its URL, its HTTP header, a list of its duplicates in the collection, information on the encoding, lists of all recursively referenced schemas, and all well-formedness and validity errors. The actual files are saved on disk with the appropriate id as its filename. Duplicates were not removed from the dataset, but rather a relation of duplicate content was inserted into the database.

4. QUALITY OF XML ON THE WEB

In three sections we look at the basic quality requirements outlined in Section 1.1: character encoding, well-formedness and validity. We are not only interested in the amount of errors but also whether a small amount of error-types is responsible for a large amount of errors. We also report correlations between errors and other variables.

4.1 Encoding

We checked whether documents lie about their encoding. Of every document in the collection, we checked whether the encodings as specified either in the HTTP header, in the encoding attribute in the XML declaration or in Content-Type meta tag (often used in XHTML documents) was compliant with the document. Our main result is that 99.47% of all specified encodings is correct. Further details are in the full paper.

4.2 XML Well-formedness

We use the libxml2 parser to check whether a document is wellformed. We created a modified version of the XML parser libxml2. The main change was to make the error output uniform for all errors. The modified version of libxml2 distinguishes four different error levels: No error, Warning, Recoverable error and Fatal error. Files with only recoverable errors can still be parsed. The modified version of libxml2 also categorizes the error that occurred. We found 74 different error categories.

We note that libxml2 is not meant to collect statistics on the number of errors in a not well-formed XML file. It often happens that the parser outputs a large amount of errors while fixing just one makes the document well-formed. We however believe that the output is still useful for giving directions to research on automatic error repairmen.

26,377 different files (14.6% of the collection) had at least one fatal error. Figure 5 lists the 10 most common errors. Note that a document may have many errors, so the total sum is higher than the number of bad XML files. ‘Opening and ending tag mismatch’ is encountered in most documents (16,996 docs) followed by ‘Premature end of data in tag’ (14,250 docs). Third is an unknown encoding (11,615 docs). This last error does not necessarily have to be a fatal error, as libxml2 allows specifying the encoding of a document manually.

The distribution of errors across error categories follows a Pareto distribution. The Pareto Chart in Figure 6 shows the first nine error categories. Approximately 20% of the error categories (9 from a total of 74) account for 99% of the fatal errors.

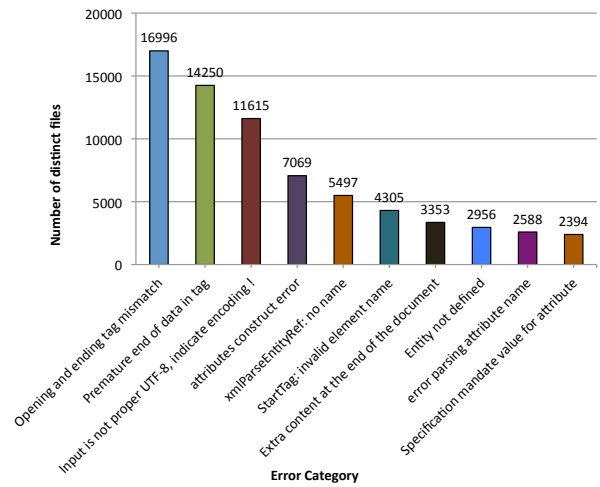


Figure 5. Top Ten Fatal Error Categories.

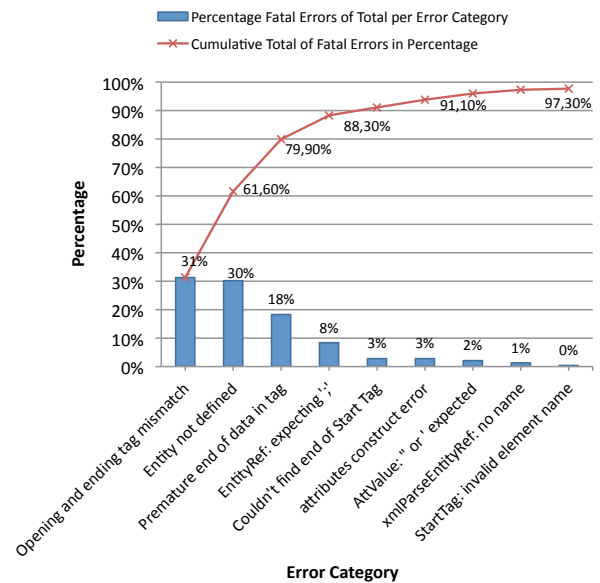


Figure 6. Pareto Chart of Fatal Errors per Error Category.

As indicated above, an error occurring need not mean that it is responsible for the file being not well-formed. Still, 5708 documents (21.6% of all documents with a fatal error), contain only one fatal error. Also here, ‘Opening and ending tag mismatch’ is most often (25.1%) responsible for not being well-formed.

We checked whether some internet zones were over- or under-represented in errors but found no significant deviations from Figure 4.

4.3 XML Validity

This section discusses the results about validity testing with respect to DTDs and XSDs.

4.3.1 Bad schemas

Figure 2 and Figure 3 give a breakdown of the reasons for non-validity of the files that reference a DTD or XSD. With both DTD

and XSD being referred to in roughly 24K files there are large differences in their validity scores. Over half of the XML files with an XSD is actually valid, in contrast to less than 10% for files referencing a DTD. On the other hand, over half of the documents with a DTD is itself not even wellformed. With XSD, this occurs less than 1%. The most interesting case of invalidity occurs when all prerequisites are satisfied. Here the two schema languages behave rather alike, with one sixth (DTD) and one fourth (XSD) of the files falling in this class.

4.3.2 Geographic distribution

We looked whether validity errors were over or under represented in certain domains and found one significant deviation. The .edu and .gov domains behave well compared to the rest: 2, respectively .4% of all files come from these domains, but of all files that refer to a well-formed DTD they contribute 11 and .9%, respectively.

4.3.3 Most common errors

First we look at errors in DTDs. A total of 28 different errors have been found, of which the top ten is shown in Figure 8.

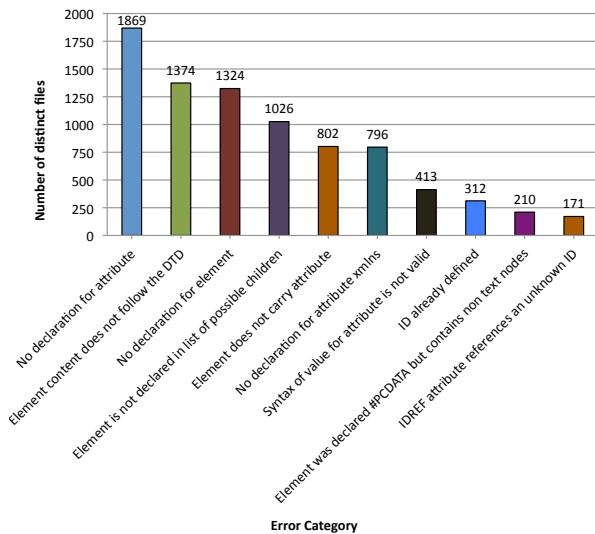


Figure 8. Top Ten Recoverable Error Categories in DTD validation based on occurrence in files.

The first two error-types are an indication that the data is in fact richer than the schema describes. If an application is built on the schema it can thus simply ignore the extra information. The third error type is problematic for parsers, and potentially difficult to repair automatically. It contains a lot of errors concerning CDATA, for instance that a text node is encountered where only element nodes are allowed. Some schemas have an obvious default element in which to wrap forbidden text nodes. E.g., in XHTML, text is forbidden under the body-element, and could be wrapped in a p-element.

This leads us to ask which DTDs have errors across the largest number of different files. The top of the list is dominated by W3C DTDs. For details, see the full paper.

For XSDs we see very much the same picture. 93% of all validation errors are of type ‘this element is not expected’. Because XSDs can, in contrast to DTDs, also restrict the type of values, one could expect many type-errors. In fact these are uncommon: less than 5% of all errors and occurring in around 10% of all invalid files.

4.3.4 The good guys and the bad guys

Though not of direct practical value, it is fun to see which background variables correlate strongly with (in)validity. We looked at file size, domain extension, encoding and the type of webserver used for the site. Details are in the full paper.

The Content-Length HTTP header is an interval variable, but not normally distributed (Kolmogorov-Smirnov test. *Sig value* < 0.05. $N = 131,831$). We use Spearman’s rho to determine if there is a relationship between file size and validation of the document. There indeed is a statistically significant but weak relationship, $r(131,831) = .214, p < .01$.

Regarding base domain, we did a binary logistic regression analysis. The produced model does indicate that domain name extension is statistically significant and explains variations in validity of the documents ($\chi^2=6087.791, df=334, p < 0.01$). We found 33 domain name extensions with a statistically significant ($p < 0.05$) effect. The 5 bad guys are ‘.jp’, ‘.org.au’, ‘.cat’, ‘.gov.uk’, ‘.gov.br’. They are 3.2, 5.1, 3.6, 3.6 and 9.4 times more likely to be invalid than to be valid. The rest are good guys, ranging from 2.226 (.gov) to 24.750 (.im) more likely to be valid than invalid.

All seven statistically significant domain name extensions in the educational and academic domains (containing .edu or .ac) are more likely to be valid than invalid. In contrast, documents from governmental domains in the USA are more likely to be valid (.gov), while documents from two other governmental domains are less likely to be valid (.gov.br and .gov.uk). An other interesting fact is that documents from the .uk domain are generally almost 2.5 times more likely to be valid than invalid, while documents from the governmental domain in the uk (.gov.uk) are 3.6 times more likely to be invalid than valid. It might indicate that documents from the British government are of poorer quality than other documents originating in the UK.

Does it matter for validity whether a site uses a commercial (Microsoft’s ISS) or an open source (Apache) server? The effect is significant but extremely small: documents served by an Apache server are 1.07 times more likely to be valid than documents that are server by Microsoft IIS.

The encoding of a file has a minor effect on validity. The only statistically significant effects we found are for windows125(1|2) which is twice more likely to be invalid and iso88591 which is 2.5 times more likely to be valid.

4.3.5 Translations between schema languages

Because our collection does not contain any Relax NG schemas we also crawled schema files directly. Combined with the schema’s already found we created a collection of 3087 DTDs, 4141 XSDs, 337 Relax NGs in compact syntax and 338 Relax NG’s in XML. All schemas are syntactically correct. All data and results are at <http://data.politicalmashup.nl/xmlweb/trang.html>.

While XSDs allow expressions that cannot be expressed in DTD syntax, these extras are rarely used in practice [7][8]. We wondered whether available schema translation software could support these findings on our dataset. We used James Clark’s Trang as it can translate between all three schema languages except from XSD. The results are that Trang translates 88% of the DTDs to the other two languages and that 30% of the Relax schemas can be translated to DTD and 96% to XSD. The surprisingly low 88% seems due to Trang, not to the use of DTD features the other languages cannot handle. The University of Dortmund is working on an improved translator based on [16], which can also translate from XSD.

5. CONCLUSIONS

Our results show that it is possible to do ETL tasks on 'XML' files solely using XML query and transformation languages. Only 14.6% is not truly XML and the distribution of errors is promising for (semi-)automatic repair. Of course ETL development would become much easier and far more robust when restricted to valid XML. Here the quality of the XML web needs drastic improvement as less than 10% is valid. Although it is hard to compare our data with previous studies, the growth of referenced XSDs and the fact that files with an XSD tend to be twice as often valid as those with a DTD seems a positive development. We have set up our study in such a way that it can easily be replicated in the future. Hopefully we can measure an upward trend in validity.

The distribution of XML syntax errors follows an 80-20 law which make them amenable to automatic error repairing techniques. Validation errors occur because schemas do not compile or because the XML is not valid. This shows that work on (semi-)automatic learning DTDs or XML Schemas from XML documents is useful [6][9]. Most validation errors occur because there is an element or attribute used that is not defined in the schema. This could mean that either the schema is not correct or a wrong name is used in the XML. Schema learning techniques may be expanded to schema repairing techniques. Techniques used in data-deduplication and learning schema mappings seem useful to repair XML documents in this case.

6. REFERENCES

- [1] Abiteboul, S., & Vianu, V. (1997). Queries and Computation on the Web. *ICDT '97: Proceedings ICDT* (pp. 262-275).
- [2] Azze-Eddine, M., Samia, K.-B., & Douniazed, A. H. (2004). XML-DFG : A Dynamic Forms Generator for XML Valid DTD Document. *RIST*, 14 (2), 15-26.
- [3] Barbosa, D., Mignet, L., & Veltri, P. (2005). Studying the XML Web: Gathering Statistics from an XML Sample. *World Wide Web*, 8 (4), pp. 413-438.
- [4] Beatty, P., Dick, S., & Miller, J. (2008, Mar/Apr). Is HTML in a Race to the Bottom? A Large-Scale Survey and Analysis of Conformance to W3C Standards. *IEEE Internet Computing*, 12 (2), pp. 76-80.
- [5] Beckett, D. (1997). 30% accessible - a survey of the UK Wide Web. *Computer Networks and ISDN Systems*, 29 (Nos 8-13), pp. 1367-75.
- [6] Bex, G.J., Gelade, W., Neven, F. & Vansummeren, S. Learning deterministic regular expressions for the inference of schemas from XML data. *WWW 2008*: 825-834
- [7] Bex, G. J., Martens, W., Neven, F., & Schwentick, T. (2005). Expressiveness of XSDs: from practice to theory, there and back again. : *Proceedings WWW* (pp. 712-721)
- [8] Bex, G. J., Neven, F., & Bussche, J. V. (2004). DTDs versus XML Schema: A Practical Study. *Proceedings WebDB '04* (pp. 79-84).
- [9] Bex, G. J., Neven, F., Schwentick, T., & Tuyls, K. (2006). Inference of concise DTDs from XML data. : *Proc. VLDB '06* (pp. 115-126).
- [10] Chen, B., & Shen, V. Y. (2006). Transforming Web Pages to Become Standard-Compliant through Reverse Engineering. : *Proceedings W4A '06*. pp. 14-22.
- [11] Chen, S., Hong, D., & Shen, V. (2005). An experimental study on validation problems with existing HTML web pages. *Proceedings ICOMP '05*. (pp. 373-379).
- [12] Choi, B. (2002). What are real DTDs like? *Proceedings WebDB '02*, (pp. 43-48).
- [13] Doan, A., & Halevey, A. (2005). *AI Magazine*, Vol. 26, pp. 83-94.
- [14] Elmagarmid, A., Ipeirotis, P., & Verykios, V. (2007). Knowledge and Data Engineering, *IEEE Transactions on*, 19 (1), pp. 1-16.
- [15] Furche T., Gottlob G., Grasso G., Schallhart Ch., Sellers A. 2011. OXPath: A Language for Scalable, Memory-efficient Data Extraction from Web Applications. In *Proceedings VLDB 2011*.
- [16] Gelade, W., Idziaszek, T., Martens, W. & Neven, F. (2010) Simplifying XML Schema: Single-Type Approximations of Regular Tree Languages. *Proc. PODS 2010*.
- [17] Gottlob, G., Koch, Ch., Baumgartner, R., Herzog, M. & Flesca, S. 2004. The Lixto data extraction project: back and forth between theory and practice. In *Proceedings PODS '04* (pp. 1-12).
- [18] Guerrini, G., Mesiti, M., & Rossi, D. (2005). Impact of XML schema evolution on valid documents. *Proc. WIDM '05* (pp. 39-44).
- [19] Hackett, S., Parmanto, B., & Zeng, X. (2004). Accessibility of Internet websites through time. *SIGACCESS Access. Comput.*, 32-39.
- [20] Harold, E. R. (2001). *XML Bible*. New York, NY, USA: John Wiley & Sons, Inc.
- [21] Klettke, M., Schneider, L., & Heuer, A. (2002). Metrics for XML Document Collections. *XMLDM Workshop*, (pp. 162-176). Prague.
- [22] Kosek, J., Kratky, M., & Snašel, V. (2003). Struktura realnych XML dokumentu a metody indexovani. *ITAT 2003: Workshop on Information Technologies Applications and Theory*. High Tatras, Slovakia.
- [23] Lawrence, S., & Giles, C. L. (2000, Spring). Accessibility of information on the Web. *Intelligence*, 11 (1), pp. 32-39.
- [24] Lee, D., & Chu, W. W. (2000). Comparative analysis of six XML schema languages. *SIGMOD Rec.*, 29 (3), 76-87.
- [25] Liu, B. (2007). *Web Data Mining*. Springer.
- [26] Madnick, S., Wang, R., Lee, Y. & ZHU, H. (2009) Overview and framework for data and information quality research. *Journal of Data and Information Quality*.1 (1), pp. 2.1-2.22.
- [27] Martens, W., Neven, F., & Schwentick, T. (2005). Which XML schemas admit 1-pass preorder typing? *Proc.s ICDT* (pp. 68-82).
- [28] Martens, W., Neven, F., Schwentick, T., & Bex, G. J. (2006). Expressiveness and complexity of XML Schema. *ACM Trans. Database Syst.*, 31 (3), 770-813.
- [29] McDowell, A., Schmidt, C., & Yue, K.-B. (2004). Analysis and Metrics of XML Schema. *Proc. SERP '04* (pp. 538-544).
- [30] Mignet, L., Barbosa, D., & Veltri, P. (2003). The XML Web: a First Study. *Proceedings WWW '03* pp. 500-510.
- [31] Mlynkova, I., Toman, K., & Pokorny, J. (2006). *Statistical Analysis of Real XML Data Collections* (Technical Report). Charles University, Faculty of Mathematics and Physics, Department of Software Engineering, Prague, Czech Republic.
- [32] Ofuonye, E., Beatty, P., Dick, S., & Miller, J. (2010). Prevalence and classification of web page defects. *Online Information Review*, 34 (1), 160-174.
- [33] Pollach, I., Pinterits, A., & Treiblmaier, H. (2006). Environmental Web Sites: An Empirical Investigation of Functionality and Accessibility. *Proceedings of the 39th Hawaii International Conference on System Sciences*. IEEE.
- [34] Raghavan, S., & Garcia-Molina, H. (2001). Crawling the Hidden Web. *Proceedings VLDB '01* (pp. 129-138)
- [35] Rahm, E., & Do, H.H. (2000). *Data Cleaning: Problems and Current Approaches*, 23 (4)
- [36] Sahuguet, A. (2001). Everything You Ever Wanted to Know About DTDs, But Were Afraid to Ask. *Proc. WebDB* (pp. 171-183).
- [37] Sundaresan, N., & Moussa, R. (2001). Algorithms and programming models for efficient representation of XML for Internet applications. *Proceedings WWW '01* (pp. 366-375).
- [38] Toman, K., & Mlynková, I. (2006). XML Data - The Current State of Affairs. *Proceedings of XML Prague '06 conference*, (pp. 87-102).
- [39] Wang, R. & Strong, D. Beyond accuracy: What data quality means to data consumers. *J. Manag. Inf. Syst.* 12(4) pp. 5-34.
- [40] W3C. (2004, 02 10). *World Wide Web Consortium (W3C)*. Retrieved 04 05, 2010, from RDF - Semantic Web Standards: <http://www.w3.org/RDF>