



HELSINKI UNIVERSITY OF TECHNOLOGY  
Department of Computer Science and Engineering  
Laboratory of Software Technology

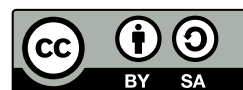
**Henri Sivonen**

# **An HTML5 Conformance Checker (DRAFT)**

~~Master's Thesis submitted in partial fulfillment of the requirements for the degree of Master of Science in Technology.~~

Helsinki, 2007-03-XXX

Supervisor and Instructor: Professor Jorma Tarhio



© 2006–2007 Henri Sivonen

Digital versions of this thesis (including the source files) may be obtained from:  
<http://hsivonen.iki.fi/thesis/>

~~DRAFT NOT EFFECTIVE YET! This literary work (“Work”) is licensed under the Creative Commons Attribution-ShareAlike license version 2.5 or later (“License”). The license text is available from <http://creativecommons.org/licenses/by-sa/2.5/>. You may have received the Work aggregated in a digital file or on a tangible medium together with a Creative Commons license badge graphic and/or the “wing” emblem of the Helsinki University of Technology, and/or you may have received the Work in a digital file that contains embedded fonts. The license badge, the “wing” emblem and the embedded fonts are not part of the Work and are not covered by the License. For avoidance of doubt, when the Work or a Derivative Work is distributed as a file containing embedded fonts (e.g. PDF) or as a markup document accompanied by external style definitions and the markup document can be intelligibly rendered using the default style definitions of typical rendering software (e.g. semantic HTML or L<sup>A</sup>T<sub>E</sub>X using standard macro names), the embedded fonts or the accompanying external style definitions are not considered to be subject to the ShareAlike provision of the License by the Licensor and are not required to be licensed under the License.~~

~~The Licensor believes that the license badge, the “wing” emblem and the embedded fonts do not encumber exercising the rights granted under the License for distribution of verbatim copies of the Work. However, if you create a Derivative Work, please ensure that you are permitted to use elements that are not covered by the License but that you may have received together with the Work or delete the elements if in doubt.~~

~~Please refer to <http://creativecommons.org/policies> for information about the use of the license badge. If you have inquiries about the “wing” emblem, please contact the Helsinki University of Technology.~~

<b>Author:</b>	Henri Sivonen
<b>Department:</b>	Computer Science and Engineering
<b>Major:</b>	Software Systems
<b>Minor:</b>	Strategy and International Business
<b>Title of the thesis:</b>	An HTML5 Conformance Checker
<b>Number of pages:</b>	
<b>Date:</b>	XXX March, 2007
<b>Professorship:</b>	T-106 Software Technology
<b>Supervisor:</b>	Professor Jorma Tarhio
<b>Instructors:</b>	Professor Jorma Tarhio
<p>The WHATWG is developing HTML5 and XHTML5 as successors for HTML 4.01 and XHTML 1.0. An (X)HTML5 conformance checker is expected to take the role that DTD-based validators have had with earlier (X)HTML. Conformance checking goes beyond the capabilities of DTDs.</p> <p>The WHATWG does not prescribe an implementation strategy for conformance checkers and does not endorse schema languages. Realizing that no schema language is adequate for describing the conformance requirements for (X)HTML5, a mainly RELAX NG-based implementation approach was chosen nonetheless.</p> <p>The bulk of the (X)HTML5 language is described as a RELAX NG schema that is supported by a custom datatype library written in Java. A Schematron schema is used alongside RELAX NG for enforcing constraints for which RELAX NG is not suitable. The remaining requirements are enforced by custom code written in Java. For checking HTML5, a special-purpose parser was developed so that the XML tools can work on XHTML5-like parse events.</p> <p>The design of the system and the experience gained so far in the ongoing project are discussed. The ease of expressing and changing the grammar is identified as the main benefit of RELAX NG. The inability to easily fine-tune error messages is identified as a drawback.</p>	
<b>Keywords:</b>	

<b>Tekijä:</b>	Henri Sivonen
<b>Osasto:</b>	Tietotekniikka
<b>Pääaine:</b>	Ohjelmistojärjestelmät
<b>Sivuaine:</b>	Yritysstrategia ja kansainvälinen liiketoiminta
<b>Sivuaine:</b>	HTML5-konformanssitarkistin
<b>Sivumäärä:</b>	
<b>Päiväys:</b>	XXX maaliskuuta 2007
<b>Professori:</b>	T-106 Ohjelmistotekniikka
<b>Työn valvoja:</b>	Professori Jorma Tarhio
<b>Työn ohjaajat:</b>	Professori Jorma Tarhio
<b>Avainsanat:</b>	

# Acknowledgements

This Master's thesis has been done at the Laboratory of Software Technology of Helsinki University of Technology.

I want to thank Ian Hickson for all his work on HTML5, without which this thesis would not exist.

I wish to thank Elika Etemad for developing the core RELAX NG schema for HTML5 and reviewing and commenting on my changes to the schema.

I would also like to thank the Mozilla Foundation for funding this project.

I want to thank James Clark for developing the Jing validation engine that the software developed in this project is based on.

My gratitude also goes to members of the #turska and #whatwg IRC channels as well as the members of the WHATWG mailing list.

I would like to thank YesLogic Pty. Ltd., SyncRO Soft Ltd. and Oskar Ojala for software that I used to make this thesis publishable.

I wish to thank my instructor and supervisor professor Jorma Tarhio.

Finally, I would like to thank my parents.

Helsinki, XXX March 2007

Henri Sivonen

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Methods . . . . .	2
1.3	Objectives . . . . .	3
1.4	The Organization of this Thesis . . . . .	3
<b>2</b>	<b>The History of HTML Leading to HTML5</b>	<b>5</b>
2.1	Early HTML . . . . .	5
2.11	Initial HTML at CERN . . . . .	5
2.12	The IIIR Draft . . . . .	6
2.13	HTML+ . . . . .	6
2.14	HTML 2.0 . . . . .	7
2.15	HTML 3.0 . . . . .	7
2.16	HTML 3.2 . . . . .	8
2.2	Contemporary HTML . . . . .	9
2.21	HTML 4 . . . . .	9
2.22	ISO HTML . . . . .	10
2.23	XHTML 1.0 . . . . .	10
2.24	Modularization . . . . .	11
2.24.1	XHTML Basic . . . . .	12
2.24.2	XHTML 1.1 . . . . .	12
2.24.3	XHTML Mobile Profile . . . . .	12
2.3	HTML5 . . . . .	12
2.31	The Mozilla/Opera Joint Position Paper . . . . .	13
2.32	The WHAT WG is Formed . . . . .	14
2.33	The WHATWG Specifications . . . . .	15
2.33.1	Web Forms 2.0 . . . . .	15
2.33.2	Web Applications 1.0 . . . . .	16
<b>3</b>	<b>Schema Languages</b>	<b>19</b>
3.1	DTDs . . . . .	19

3.2	W3C XML Schema . . . . .	21
3.3	Document Structure Description . . . . .	22
3.4	TREX, RELAX, XDUCE and DDML . . . . .	23
3.5	RELAX NG . . . . .	23
3.51	Datotyping . . . . .	24
3.52	Compact Syntax . . . . .	25
3.53	Use in This Project . . . . .	25
3.6	Schematron . . . . .	25
3.61	Using RELAX NG and Schematron Together . . . . .	26
3.62	Use in This Project . . . . .	26
<b>4</b>	<b>Prior Work on Markup Checking</b>	<b>27</b>
4.1	The W3C Markup Validation Service . . . . .	27
4.2	WDG HTML Validator . . . . .	28
4.3	Page Valet . . . . .	29
4.4	The Schneegans XML Schema Validator . . . . .	30
4.5	Relaxed . . . . .	30
4.6	Feed Validator . . . . .	31
4.7	Validome . . . . .	32
<b>5</b>	<b>Implementation</b>	<b>35</b>
5.1	The Basic Back End . . . . .	35
5.2	The Jing Validation Engine . . . . .	36
5.3	The RELAX NG Schema . . . . .	36
5.31	The General Schema Design . . . . .	37
5.32	Common Definitions . . . . .	37
5.32.1	Common Content Models . . . . .	37
5.32.2	Common Attributes . . . . .	38
5.32.3	Common Datatypes . . . . .	38
5.32.4	Parameter Switches . . . . .	39
5.33	Examples of Elements . . . . .	40
5.4	The HTML5 Datatype Library . . . . .	42
5.41	Dates . . . . .	43
5.42	IRIs . . . . .	43
5.43	Language Tags . . . . .	44
5.44	ECMAScript Regular Expressions . . . . .	45
5.5	The Schematron Schema . . . . .	45
5.51	Exclusions . . . . .	45
5.52	Required Ancestors . . . . .	46
5.53	Referential Integrity . . . . .	47
5.6	The Non-Schema-Based Checkers . . . . .	48

5.61	Table Integrity Checker . . . . .	49
5.62	Checking the Text Content of Specific Elements . . . . .	52
5.63	Checking for Significant Inline Content . . . . .	53
5.64	Unicode Normalization Checking . . . . .	54
	5.64.1 Requirements . . . . .	54
	5.64.2 Interpretation . . . . .	54
	5.64.3 Implementation . . . . .	55
5.7	The HTML Parser . . . . .	56
5.8	Character Model Checking . . . . .	57
5.9	The Front End . . . . .	59
<b>6</b>	<b>Shortcomings</b>	<b>61</b>
6.1	Non-Ideal Error Messages . . . . .	61
	6.11 Bimorphic Content Models . . . . .	61
	6.12 Lack of Datatype Diagnostics . . . . .	62
6.2	Poor Localizability . . . . .	62
6.3	Opportunities for Optimization . . . . .	63
	6.31 RELAX NG . . . . .	63
	6.32 Schematron . . . . .	64
<b>7</b>	<b>Applicability in Other Contexts</b>	<b>67</b>
7.1	RELAX NG-Guided Autocompletion in Editors . . . . .	67
7.2	Content Management Systems . . . . .	67
<b>8</b>	<b>Future Work</b>	<b>69</b>
8.1	Open Up . . . . .	69
8.2	The HTML5 Parsing Algorithm . . . . .	69
8.3	Tracking the Specification . . . . .	70
8.4	RELAX NG Message Improvements . . . . .	70
8.5	More Non-Schema-Based Checkers . . . . .	71
8.6	Assistance for Checking Human-Checkable Requirements . . . . .	72
8.7	Web Service . . . . .	73
8.8	Embedded MathML and SVG . . . . .	73
8.9	Showing the Erroneous Source Markup . . . . .	73
<b>9</b>	<b>Conclusions</b>	<b>75</b>
	<b>References</b>	<b>77</b>



# Chapter 1

## Introduction

The WHATWG is developing HTML5 and XHTML5 as successors for HTML 4.01 and XHTML 1.0. To be successful, a new markup language not only needs support from browsers. It also needs tool supporting authoring. Authoring-side tools include editors, content management systems and quality assurance tools for checking the correctness of markup. This thesis focuses on the latter.

### 1.1 Motivation

Web authors tend to make mistakes when writing HTML. The vast majority of HTML documents on the Web are erroneous. A study on DTD-based validation of Web pages estimated about 95% of pages to be invalid [ValidationProbs]. A test of the HTML5 parsing algorithm on several billion documents spidered by Google indicated that 93% of documents have low-level syntax errors [Several]. (Documents in the remaining 7% may well have higher-level errors which are not found by the parsing algorithm and would require a full conformance checker to find.)

Even though most Web content is broken without hope of repair and browsers will do *something* with *any* input purporting to be HTML, it is still useful to provide a quality assurance tool for authors. Even if browsers adopt the well-defined error-recovering processing models of HTML5, authors generally do not make errors on purpose in order to elicit particular error recovery response. Silent recovery from inadvertent mistakes—even if deterministic and well-defined—may still confuse an author who did not mean to invoke error recovery. The issue becomes more apparent when an author uses a style sheet or a script that assumes the document to be correct.

Therefore, it is worthwhile to provide a conformance checker that helps authors find their mistakes.

## 1.2 Methods

With HTML 4.01 and XHTML 1.0, DTD-based validators have traditionally been used as the quality assurance tools for checking correctness even though they do not check for all machine-checkable conformance requirements. An (X)HTML5 conformance checker is expected to take the role that DTD-based validators have had with earlier (X)HTML. Conformance checking goes beyond the capabilities of DTDs.

The WHATWG does not prescribe an implementation strategy for conformance checkers and does not endorse schema languages. Not only are schema languages unendorsed but they are seen as being clearly inadequate. Therefore, a non-schema-based implementation strategy is implied. Yet, as an initial impression, abandoning schemata altogether just because they cannot be used for checking every machine-checkable constraint seems overly drastic.

RELAX NG was chosen as the primary schema language. Schematron was chosen as a supporting schema language. Using RELAX NG for document-oriented schemata (as opposed to databinding-oriented schemata) had gained acceptance as the best practice among users XML schema languages. Schematron had gained popularity as a language for refining a RELAX NG schema. A project for developing a RELAX NG schema for HTML5 had already been started by Elika Etemad [Whattf]. Moreover, this author had already developed a service that allows Web users to validate XML documents against arbitrary RELAX NG and Schematron 1.5 schemata [ValidatorAbout]. The service was built in the Java programming language due to the excellent availability of XML tools for Java. Etemad's schema project and the service developed by this author were chosen as starting points for this thesis project.

It was obvious that schemata alone would be inadequate, so it was decided to augment schemata with custom code instead of rejecting schema languages altogether. Since the validation service that this project was based upon was written in Java, also the custom code for augmenting schema-based validation was developed in Java.

Since the parsed syntax tree for HTML5 and the parsed syntax tree for XML are very similar and reusable tools exist for XML, it was decided to use XML tools and to map HTML5 documents to equivalent XHTML5 representations in the parser.

### 1.3 Objectives

The functional objective of the project described in this thesis is developing a partial (X)HTML5 conformance checker that is comprehensive enough to demonstrate that it can be taken to completion once (X)HTML5 itself has stabilized. The research goals are finding out if a hybrid implementation based both on schemata and on custom code developed in a general-purpose programming language is feasible and finding out if XML tooling can be successfully applied to checking the non-XML serialization of HTML5.

### 1.4 The Organization of this Thesis

This thesis has two thematic parts. The first part—the next three chapters—reviews the context of this work. HTML5 is placed in historical context (page 5), schema languages for XML are reviewed (page 19) and prior work on online markup checking services is reviewed (page 27). The second part—the last five chapters—focus on the software implemented in this project. The implementation of the software (page 35) and its shortcoming (page 61) and applicability to other contexts (page 67) is discussed. Finally, the need for future work is reviewed (page 69) and the conclusions given (page 75).



## Chapter 2

# The History of HTML Leading to HTML5

In this chapter, the history of HTML leading to HTML5 is reviewed. Since one of the major changes in HTML5 is the way the specification deals with parsing and the stance the specification takes with respect to SGML, each version of HTML prior to HTML5 is summarized in terms of the key features introduced and in terms of the stated relationship to SGML or XML.

### 2.1 Early HTML

In this review, HTML version prior to HTML 4.0 are considered early, as they are not in active use when new documents are created.

#### 2.1.1 Initial HTML at CERN

Tim Berners-Lee invented the Web in 1989. The first version of his browser was released in 1990[Raggett]. It used HTML, but the language was not formally specified at first. Tim Berners-Lee designed HTML following ideas from SGML. However, HTML was not layered on top the SGML standard but, rather, used similar tags without being a true application of SGML.

The element names available in HTML were largely taken from SGMLguid, an application of SGML used at CERN. SGMLguid, in turn, was similar to Waterloo SCRIPT GML[WaterlooGML], a GML language specified at University of Waterloo. (GML was IBM's predecessor to SGML.) [EarlyHistory] There are also similarities with the language given in the tutorial of the SGML standard[Handbook].

### 2.1.2 The IIIR Draft

The www-talk mailing list was founded in September 1991 for discussing matters related to the Web. The development of HTML was discussed on the list. [Raggett]

Tim Berners-Lee and Dan Connolly wrote an Internet Draft specification for HTML as part of the activity of the Integration of Internet Information Resources (IIIR) working group of the IETF. The Internet Draft was published in June 1993. [IIIR-HTML]

The draft said that HTML was defined in terms of SGML. However, the specification did not specify an HTML document as a conforming SGML document entity but instead said how to construct an SGML document from a transferred HTML file[IIIR-HTML]. The draft also suggested that an HTML parser would not need to be a full SGML parser but a parser that only deals with the document instance after the DTD[IIIR-HTML]. The SGML-purity of the drafted HTML approach was challenged[ToBeDeleted] by W. Eliot Kimber, an SGML expert, and the stated approach was changed in later specifications although browsers continued to behave as before.

The mailing list discussions about the relationship of HTML to SGML are reviewed in [Cascading].

The IIIR draft already included the IMG element. The P element was defined as an empty element that indicates paragraph breaks. As in interesting detail, the XMP, LISTING and PLAINTEXT elements were considered *obsolete* as early as in the draft. [IIIR-HTML]

The draft expired and did not reach the RFC status.

### 2.1.3 HTML+

One of the participants to www-talk, Dave Raggett, visited Tim Berners-Lee at CERN to discuss further development face to face. Raggett went on to draft a new version of HTML called HTML+. [Raggett]

HTML+, published in late 1993, specifically stated that it was “based on the Standard Generalized Markup Language”. It also had a DTD. However, it was implied that there would be “HTML+ parsers” which would be different from “other SGML parsers”. HTML+ explicitly excluded SGML minimization features. It used the P element as a paragraph *container* but said that author may think of the P tag as a paragraph *separator*. [HTMLplus]

HTML+ had a number of elements that never made into real-world HTML, such as BYLINE, ONLINE, PRINTED, ABSTRACT. As a curious detail, HTML+ included an element called IMAGE, which used the element content as the alternative text. HTML+ made an attempt to address the issue of

mathematical formulae, but the attempt was not particularly profound. [HTMLplus]

HTML+ defined markup for tables. The table markup is roughly what was later adopted in HTML 4. HTML+ also defined markup for forms similar to what was actually adopted in real-world browsers. However, the field types also included types that were not adopted, such as URL, DATE and SCRIBBLE (for drawing). [HTMLplus]

HTML+ never made it to mainstream browser implementations. However, at the first World Wide Web conference it was agreed that the ideas from HTML+ should be carried forward.[Raggett]

#### 2.1.4 HTML 2.0

Dan Connolly had advocated a cross-browser HTML specification at the first World Wide Web conference. In early 1994, the Internet Engineering Task Force formed a working group to specify HTML. The working group—with Connolly in the lead—defined HTML 2.0 based on the then-current practice. The HTML 2.0 draft was published in July 1994. [Raggett]

HTML 2.0 reached the RFC status in November 1995. HTML 2.0 explicitly claims that “HTML is an application of ISO Standard 8879:1986 Information Processing Text and Office Systems; Standard Generalized Markup Language (SGML).”[RFC1866] However, it was too late to make browsers use real SGML parsers. Instead, browsers continued to use special-purpose HTML parsers without standard error recovery behavior. HTML 2.0 included a DTD, but the DTD was of no interest to browsers.

Unlike the elements proposed in HTML+, the elements of HTML 2.0 were (and still are) actually supported by real-world browsers. HTML 2.0 included forms but did not include tables which had been proposed in HTML+. Regardless, Netscape implemented tables in its browser in the HTML 2.0 era and made tables popular.

HTML 2.0 established that the document character set of HTML is ISO-10646 (equivalent to Unicode) regardless of the *character encoding* used to transfer the document. However, the internationalization of HTML 2.0 was not fully addressed in the HTML 2.0 specification itself, and a standards track RFC that extended HTML 2.0 to address internationalization issues was published as late as in 1997[RFC2070].

#### 2.1.5 HTML 3.0

Both Netscape and the W3C were founded in 1994. Netscape added features to HTML on its own. Dave Raggett—this time representing the

W3C—edited a specification called HTML 3.0 which carried on the ideas of HTML+ instead of attempting to formalize the Netscape extensions. [Raggett]

To support the use of style sheets, HTML 3.0 introduced the `STYLE` element and the `CLASS` attribute, which lived on. [Raggett]

An HTML 3.0 draft was published through the IETF as an Internet Draft[HTML30], but it was not adopted in mainstream browsers. HTML 3.0 was abandoned and never reached the RFC status[Raggett].

Even though HTML 3.0 as a whole was dropped, a specification for HTML tables (as an extension to HTML 2.0) was published as an experimental RFC. [RFC1942]

### 2.1.6 HTML 3.2

In November 1995, representatives of browser vendors and the W3C formed an HTML working group at the W3C. In December 1995, the IETF HTML working group was disbanded. [Raggett]

In January 1997, the W3C published HTML 3.2 as a Recommendation. Unlike HTML 3.0, HTML 3.2 documented actual practice that had grown as extensions to HTML 2.0. The specification itself stated: “HTML 3.2 aims to capture recommended practice as of early ’96 and as such to be used as a replacement for HTML 2.0 (RFC 1866).” [HTML32]

HTML 3.2 continued to assert that HTML is an application of SGML: “HTML 3.2 is an SGML application conforming to International Standard ISO 8879 – Standard Generalized Markup Language. As an SGML application, the syntax of conforming HTML 3.2 documents is defined by the combination of the SGML declaration and the document type definition (DTD).” However, even the specification itself admitted that SGML-compliance of user agents was not part of the actual practice as of early ’96 by noting: “Note that some user agents require attribute minimisation for the following attributes: `COMPACT`, `ISMAP`, `CHECKED`, `NOWRAP`, `NOSHADE` and `NOHREF`. These user agents don’t accept syntax such as `COMPACT=COMPACT` or `ISMAP=ISMAP` although this is legitimate according to the HTML 3.2 DTD.”. [HTML32]

In documenting the actual practice, HTML 3.2 included presentational features such as the `FONT` element that would later be frowned upon.



## 2.2 Contemporary HTML

The versions of HTML discussed above are of historical interest and are not in active use for creating new documents. The versions in current use start with HTML 4.

### 2.2.1 HTML 4

HTML 4.0 was published as a W3C Recommendation in December 1997[HTML40]. It was revised without incrementing the version number, and the revision was published in April 1998[HTML40rev]. Another revision called HTML 4.01 became a W3C Recommendation in December 1999[HTML401].

Once again, the specification claimed that “HTML 4 is an SGML application conforming to International Standard ISO 8879 – Standard Generalized Markup Language.”[HTML401] Yet, the specification acknowledged the reality that user agents in general are not conforming SGML systems: “SGML systems conforming to [ISO8879] are expected to recognize a number of features that aren’t widely supported by HTML user agents. We recommend that authors avoid using all of these features.”

HTML 4 deprecated presentational features such as the `FONT` element that had made its way to a W3C Recommendation less than a year prior to the first version of HTML 4. In principle, HTML 4 tried to backpedal on the point of presentational features to where HTML 2.0 was—with the intent of leaving presentation to style sheets. HTML 4 without the deprecated features was termed Strict and HTML 4 with the presentational features was termed Transitional. In practice, the deprecated features continue to be used.

HTML 4 added features for adding more structure to tables, for adding more structure to forms, and for marking up insertions and deletions. The table model adopted by HTML 4 is a subset of the model proposed in the experimental RFC on HTML tables [RFC1942]. Internationalization features, including support for bidirectional text (e.g. for Hebrew and Arabic), were adopted from the standards track RFC on the internationalization of HTML [RFC2070].

HTML 4 introduced the `OBJECT` element, which was supposed to eventually replace `IMG`, `APPLET` and the Netscape `EMBED` elements. `EMBED` did not fit together with an SGML DTD, because it could take arbitrary attributes. However, in practice, browsers continued to support `EMBED`, and even today browsers do not fully support `OBJECT` as designed.

HTML 4 formalized frames that had been introduced by Netscape and that were frowned upon[Frames] even before their inclusion in HTML 4.

Additionally, HTML 4 introduced `IFRAME`, which was not part of Netscape frames.

### 2.2.2 ISO HTML

In 2000, ISO standardized its own version of HTML by referencing a subset of HTML 4.0 as defined by the W3C but, yet, also making changes other than mere subsetting in the DTD[ISO15445]. A technical corrigendum changed the references to HTML 4.01[ISO15445TC1].

In practice, ISO HTML is only of curiosity value, since it has been largely ignored by Web authors.

### 2.2.3 XHTML 1.0

XML was published as a W3C Recommendation in February 1998[XML]. XML is a simplification of SGML that stands alone without making a normative reference to SGML. Since HTML was supposedly an application of SGML and the W3C now had its own replacement for XML, the W3C decided to swap the markup language framework from underneath HTML. The result was XHTML 1.0—a reformulation of HTML 4 in XML. XHTML 1.0 became a Recommendation in January 2000[XHTML10].

XHTML 1.0 includes the features that were deprecated in HTML 4—that is, XHTML 1.0 has three version just like HTML 4: Strict, Transitional and Frameset.

To be compatible with existing HTML user agents, XHTML 1.0 included compatibility guidelines commonly known as “Appendix C”. Appendix C limits the syntactic sugar permitted by XML 1.0 so that an XHTML 1.0 document that adheres to Appendix C could be processed by existing HTML user agents if served as `text/html`. Appendix C relies on the fact that real-world browsers do not actually process `text/html` as SGML.

Appendix C made it easy to pretend that something new had been achieved, because a new language could be deployed but the user agents did not need to be updated. Obviously, since the browsers gained no new capabilities, using XHTML 1.0 could not actually deliver any true benefits over HTML 4 in user agents designed for HTML. No XML processor was involved despite the XHTML 1.0 being a reformulation in XML! In fact, the HTML WG of the W3C gave an explicit opinion that browsers should not try to process documents served as `text/html` using a real XML processor.

Later, the makers of notable browser engines (except the maker of the engine with the largest desktop market share) actually took the XML nature of XHTML seriously and implemented support for XHTML 1.0 using a real

XML processor when the document is served using the `application/xhtml+xml` media type (instead of `text/html`). This approach has not become popular for three reasons: First, since XHTML 1.0 is a reformulation of HTML 4 on top another markup language framework, it (alone) does not enable new interesting things in the browser, which means there is not a compelling technical advantage to be gained from using XHTML 1.0 served as `application/xhtml+xml` over HTML 4.01 served as `text/html`. Second, the browser engine with the largest desktop market share (Trident that is the engine of Microsoft's Internet Explorer) does not support `application/xhtml+xml`. Third, when a real XML processor is used, an error is reported if the document is not syntactically XML. This means that author need to pay more attention to actual correctness as opposed to merely claiming to be using something new.

Moreover, there are subtle differences in the ways CSS and the DOM API exposed to JavaScript interact with `text/html` and `application/xhtml+xml`. Differences involve issues such as case-sensitivity and whether elements are namespace[MozFAQ]. Also, `document.write()`, which allows scripts to insert data into the character stream being parsed, does not work in XML. In practice, script written naively for XHTML served as `text/html` do not work when the document is served as `application/xhtml+xml`.

There are experts close to the development of browser engines who have discredited the practice of serving XHTML as `text/html` (e.g. [Harmful] and [Understanding]).

## 2.2.4 Modularization

The W3C decided to abandon the development of the old non-XML HTML and to only develop XHTML. After the reformulation of HTML 4 in XML, which became XHTML 1.0, the W3C HTML working group proceeded to modularize XHTML. Modularization meant dividing XHTML into logical parts such as Hypertext Module or Image Module and rewriting the previously monolithic DTD as multiple files following the logical module partitioning.

The rationale for the modularization was based on the belief that one size of XHTML did not fit all client platforms. In the words of the specification itself: "This modularization provides a means for subsetting and extending XHTML, a feature needed for extending XHTML's reach onto emerging platforms." [M12N] The rationale for modularization implicitly assumes a walled garden-style world view of mobile operators where a client platform design can dictate a language subset used on the network. Such a view assumes a

separate “Mobile Web”, because—quite obviously—the real Web would still use full HTML or XHTML.

#### 2.2.4.1 XHTML Basic

XHTML Basic[XHTMLBasic], published in late 2000, defines a baseline for XHTML languages built on top the Modularization. The specification itself lists “mobile phones, PDAs, pagers, and settop boxes” as target devices[XHTMLBasic].

#### 2.2.4.2 XHTML 1.1

XHTML 1.1[XHTML11], published in 2001, was the first (X)HTML specification since HTML 4 that introduced a new feature. XHTML 1.1 includes the XHTML modules that correspond to XHTML 1.0 Strict. Additionally, XHTML 1.1 includes the XHTML Ruby Annotation module[Ruby] for expressing a type of text annotations used in East Asia.

Microsoft’s Internet Explorer for Windows supports a draft version of Ruby markup when used in `text/html` documents. However, the most notable browsers that support `application/xhtml+xml` do not support Ruby. Therefore, in practice, XHTML 1.1 has failed to make a significant practical impact.

#### 2.2.4.3 XHTML Mobile Profile

In 2001, WAP Forum—a consortium of mobile phone manufacturers—defined a superset of XHTML Basic called XHTML Mobile Profile[XHTML-MP]. XHTML Mobile Profile import some modules only partially, which demonstrates the those who subset XHTML will not necessarily do so in a manner prescribed by the W3C.

To confuse the media type matter further, WAP Forum registered the type `application/vnd.wap.xhtml+xml` for XHTML Mobile Profile documents.

## 2.3 HTML5

The above review explains the context in which HTML5 was born. Prior versions of HTML had officially been applications of SGML, but browsers were actually using special-purpose HTML parsers rather than SGML parsers. Moreover, significant new features had not been introduced in years and the work had focused on reformulating the syntax as XML. Yet, document

purporting to use the reformulated XML syntax were still served as `text/html`, so browsers kept using the same special-purpose parsers as before.

There was demand for new features for HTML and demand for the recognition of the fact that `text/html` content was parsed neither as SGML nor as XML.

### 2.3.1 The Mozilla/Opera Joint Position Paper

In June 2004, the W3C held a workshop on Web Applications and Compound Documents. The Mozilla Foundation and Opera Software—the two most active browser vendors in the W3C at the time—submitted a joint position paper noting the “rising threat of single-vendor solutions” and calling for seven principles to be followed in the design of Web Applications Technologies[JointPosition].

The balance of power in the W3C had shifted from traditional desktop browser vendors to various other interest groups such as makers of software for mobile walled gardens and developers of “rich client” technologies that could be deployed on intranets but that were not used by the general public on the Web. This had led to a situation where the focus was more on the “Semantic Web”, “Web Services” and “Mobile Web” than on what is usually considered “the Web”. As a result, the development of the Web itself had been neglected and the two browser vendors were trying to put the Web back on the agenda. (At the time Microsoft—a notable browser vendor itself—was pushing a single-vendor solution code named Avalon[MS-WebApps] and Apple was catching up having entered the market only recently.)

The first one of the seven principles in the Mozilla/Opera position paper was “Backwards compatibility, clear migration path”[JointPosition]. The transition from HTML 4 to XHTML 1.0 had not worked out smoothly as discussed above. Also, XForms—the W3C’s successor for HTML forms—did not provide backwards compatibility or a clear migration path. Moreover, the HTML working group was working on XHTML 2.0, which was designed to be incompatible with XHTML 1.x, even though even though even the transition to XHTML 1.x served as `application/xhtml+xml` was not complete.

The position paper called for well-defined error handling—something that had never been addressed for HTML. It also took a position in favor of graceful recovery (as in CSS) and against the Draconian error policy of XML.

The paper called for every feature to be backed by a practical use case and for the specification process to be open. This is in stark contrast with

including features that are “nice to have” in theory and making decisions on the W3C’s member-only mailing lists.

The paper took a position against device-specific profiles—in direct contrast with the Modularization of XHTML discussed above as well as mobile profiles of other W3C deliverables such as SVG. The paper also took a position more favorable to scripting (JavaScript in practice) than what has been the general line in the W3C.

The paper went on to list specific features that a Web application host environment should provide.

The paper stated two design principles for compound documents: “Don’t overuse namespaces” and “Migration path”. The latter was related to the problems with the HTML to XHTML migration discussed above. The position paper was dismissive of schema languages.

The position paper made several references to XBL—a very politicized language.

### 2.3.2 The WHAT WG is Formed

The proposal presented by Opera Software and the Mozilla Foundation was not well received at the W3C. At the end of the second day of the workshop, a straw poll was held on the topic of the joint position paper: whether the W3C should developed extensions to HTML, CSS and the DOM as proposed. Of the 51 attendees of the workshop, 8 voted in favor of the motion and 11 voted against. When the motion was slightly reformulated, 14 voted against. [cdf-ws-minutes2]

Two days after the vote at the workshop, The Web Hypertext Applications Technology Working Group (WHATWG) and its public mailing list were publicly announced. The group was described as “a loose, unofficial, and open collaboration of Web browser manufacturers and interested parties”. The stated intent was creating specification for implementation in “mass-market Web browsers, in particular Safari, Mozilla, and Opera”. [WHAT-Ann]

The initial (invite-only) membership of the WHATWG consisted of individuals affiliated with Opera Software, Mozilla and Apple. (Ian Hickson, the editor of the WHATWG specifications, later moved to Google.) However, in the view of the Web held by the WHATWG, there is also a fourth mass-market browser: Microsoft’s Internet Explorer—the leader in market share. Microsoft is not participating in the WHATWG despite having been invited. The publicly stated reason is that the WHATWG lacks a patent policy [Wilson]. Dean Edwards, a non-Microsoft Internet Explorer expert, joined the WHATWG later.

Even though the group of WHATWG “members” is invite-only, anyone is allowed to join the WHATWG mailing list and contribute technically, which makes the process open. The WHATWG members “provide overall guidance” [WHAT-Charter], which in theory means the power to impeach and replace the editor of the specifications.

Microsoft is not expected to implement the WHAT WG specifications in Internet Explorer in the near term. Instead, the implementations of the WHAT WG specifications for IE are expected to be built by teams not affiliated with Microsoft using the extensibility mechanisms provided by Microsoft in IE.

The author of this Thesis shares the view of the Web that holds Gecko, Presto, WebKit and Trident (the engines of Mozilla/Firefox, Opera, Safari and IE, respectively) the most important browser engines.

### 2.3.3 The WHATWG Specifications

The WHATWG has two specifications in development and another two that are expected in the future. [WHAT-Charter]

The two specifications being developed are Web Forms 2.0[WebForms2] and Web Applications 1.0[WebApps]. Web Forms 2.0 is an update to HTML 4.01 forms. Web Applications 1.0 is a re-specification of HTML that both constrains and extends HTML. The language specified by Web Forms 2.0 and Web Applications 1.0 taken together is colloquially referred to as (X)HTML5.

The two expected future specifications are Web Controls 1.0 for creating new widgets and CSS Rendering Object Model for defining programmatic access to the CSS rendering tree. [WHAT-Charter]

#### 2.3.3.1 Web Forms 2.0

Web Forms 2.0 extends HTML forms with new features. The HTML forms as of HTML 4.01 are considered “Web Forms 1.0”. Web Forms 2.0 is not a standalone specification. Instead, it specifies updates to HTML 4.01 and the DOM. The choice of updates is based on what has been identified as common needs and what can be implemented as a script-based library for Internet Explorer.

The most obviously visible new features are new input field types. For example, there are new inputs for dates which can be implemented in browsers by popping up a platform specific calendar widget. The new input types are backwards compatible in the sense that unknown input types degrade into text inputs in legacy browsers. Simple constraints on the values of

the input field can be declared and checked by the browser without the form author having to resort to scripting. For complex restrictions, new integration points for scripts are provided.

In addition to the new input types, there is also a repetition model for adding and removing repeating sets of fields from the form without scripting. A new XML form submission format is introduced. The format can also be used for loading values into the form fields.

Web Forms 2.0 is the most mature part of the new features of HTML5. It has already been implemented and shipped in the Opera 9 browser.

Two primary snapshots of the Web Forms 2.0 specification were used in this project. The first snapshot was taken at the start of the thesis project and was dated January 10 2006. The second snapshot was taken in November 2006 and was dated October 12 2006. In March 2007, a third snapshot was not necessary, because the specification had not changed.

### 2.3.3.2 Web Applications 1.0

Web Applications 1.0 is the main specification for HTML5. The name of the specification highlights the Web application focus of the new features. An XML-based parallel language called XHTML5 is specified alongside HTML5.

The specification has two general areas that are intertwined. On one hand, new features are specified. On the other hand, existing features are specified in detail that was absent from previous specifications. When existing features are respecified, the behavior of the four notable mass-market browsers is reverse engineered and the specification is made compatible with the existing practice. New features are based on expected use cases.

As the name of the specification suggests, there are new features aimed for Web applications. New application-oriented markup includes `canvas` for establishing a canvas onto which scripts can draw, `menu` and `command` for building menus, `meter` and `progress` for representing gauges and progress indicators, `datagrid` for a complex data display widgets, `details` for additional information that can be hidden and `event-source` for indicating that the page listens to server-sent remote events.

In addition to new element, the specification includes a number of scripting APIs for Web applications, but they are outside the scope of this thesis.

The addition of new elements for document structure is based on an analysis of common `class` attribute values as used on the Web [Stats]. New elements for document structure include `section` for document sections, `nav` for identifying page navigation, `article` for marking up standalone parts



of page content, `aside` for tangential notes, `header` for complex headers and `footer` for page footers.

Finally, HTML5 breaks the tradition of calling HTML an application of SGML and specifies an HTML-specific parsing algorithm in meticulous detail.



# Chapter 3

## Schema Languages

In this chapter, XML schema languages are reviewed in order to put the choice of RELAX NG and Schematron in context.

A schema partitions the set of all XML documents into two disjoint sets: document that are valid according to the schema and documents that are not valid according to the schema. A computer language for expressing a schema is called a schema language.

There are two main classes of schema languages: grammar-based schema languages and schema languages that are not grammar-based. The properties of grammar-based schema languages in terms of mathematics and formal language theory as well as the related validation algorithms are discussed in [Taxonomy]. Three types of grammar-based languages are identified: local, single-type and regular (from least expressive to more expressive). The paper identifies three classes of schema languages in addition to grammar-based languages: special-purpose languages dedicated to particular kind of information that may be represented as XML, languages for representing identity constraints and languages for namespace-based validation dispatching.

### 3.1 DTDs

XML has a built-in schema language called DTDs. DTD is short for document type definition. However, for reasons explained below, it does not actually define the type of the document in a useful way. Unlike many later schema languages for XML, the syntax of DTDs is not comprised of XML elements and attributes. Rather, the DTD syntax is separate from the syntax for elements and attributes. (It is sometimes said that DTDs do not have an

XML syntax, but this is misguided in the sense that the DTD syntax is part of XML.)

DTDs are grammar-based. The expressiveness of DTDs is relatively weak compared to other schema languages. In [Taxonomy], DTDs are classified to be of the type “local”, which is the weakest grammar type. A grammar-based schema language is classified as “local” if there are no competing productions. Two non-terminal productions compete if they have the same terminal on the right-hand side of the productions. The practical consequence is that the content model of an element cannot depend on the context of the element is the document tree.

SGML DTDs were slightly more expressive than XML DTDs. SGML parsers needed to know the DTD grammar in order to be able to parse the document. XML removed some of the more complex DTD features and made it possible to parse documents without the parser knowing the DTD grammar for the document.

Still, DTD-based validation is intertwined in the process of parsing XML. Instead of the merely checking that the structure of the document meets the constraints of the grammar, DTDs have features that augment the infoset. Infoset augmentation means that the validation process adds some information that is reported to the application compared to the situation where the document is parsed without any validation formalism between the document and the application. Infoset augmentation is problematic in three ways.

First, the DTD is usually included in the document by reference but the XML 1.0 specification makes processing such external references optional. As a result, the application sees a different data depending on whether the DTD is being processed or not.

Second, being able to attach data types to attributes requires that there not be two derivations for a given document that would assign conflicting data types to a given attribute. To avoid this problem, the grammars expressed as DTDs are required to be unambiguous. This requirement is restrictive and would be unnecessary for pure non-infoset-augmenting validation.

Third, since one of the key features of XML is that a document can be parsed even if it does not have a DTD and since the DTD supplied by the document itself, an application can not trust documents to supply DTDs with particular infoset augmentation features. Therefore, applications cannot rely on DTD-based infoset augmentation taking place, which severely limits the usefulness of such infoset augmentation.

DTDs allow a primitive form of attribute datatyping. The possible datatypes are CDATA, ID, IDREF, IDREFS, ENTITY, ENTITIES, NMTOKEN and NMTOKENS. Additionally, the value of an attribute can be constrained to

be one of enumerated tokens. CDATA means an unconstrained string. ID, IDREF and IDREFS exist for cross references whose integrity is checked. (There is a single document-wide namespace for IDs and a particular reference cannot be constrained to e.g. point to elements of a certain type only.) The only datatypes that constrain the allowed lexical space of the attribute value without involving any referential semantics are NMTOKEN and NMTOKENS. NMTOKEN is merely a token that satisfies a particular, relatively arbitrarily specified, grammar production. NMTOKENS is a whitespace-separated list of these. These datatypes are rather useless unless the desired datatype constraint happens to match the definition of NMTOKEN.

It was noted above that the documents supplies its own DTD. This is one of the key problems. It means that DTD-based validation cannot be used for determining if a document belongs in a class of documents (the “type” of the document) that the recipients expects to receive. That is, if two parties have agreed to exchange documents in a particular format the recipient cannot use normal DTD-based validation to find out whether a given document is in the pre-agreed format, because the sender can use another DTD in the document. DTD-based validation only shows whether a document conformance to the grammar the document declares for itself. Since DTDs require that the schema is included in the document itself, DTDs effectively require the document to be contaminated with schema-specific syntax.

There are implementations that allow validation against an application-supplied DTD, but this is not a standard part of XML 1.0 processing and other schema formalisms (RELAX NG in particular) provide superior features in such a scenario.

Lastly, DTDs do not work properly with namespaces, because the Namespaces in XML specification layers the namespace processing on top of XML 1.0 processing and, therefore, DTD-based validation takes place underneath the namespace layer—not on top of it.

DTDs were not used in this project, because DTDs are not expressive enough, are not namespace-aware and would allow smuggling of grammar productions by the document that is being checked.

## 3.2 W3C XML Schema

The W3C XML Schema (sometimes abbreviated WXS but more often XSD for XML Schema Definition) is a schema language defined by the W3C in response to shortcomings of DTDs. XSD is a grammar-based schema language. In [Taxonomy], XSD is classified to be of the type “single-type”. In a

single-type language, competing productions within a content model and competing start symbols are prohibited. The practical consequence is that the grammar is required to be unambiguous.

XSD provides a richer repertoire of data types than DTDs. In fact, data typing is seen as one of the major improvements over DTDs. XSD validation very much involves infoset augmentation. Instead of merely checking whether a document satisfies the schema, the XSD validation process yields a Post-Schema Validation Infoset (PSVI) which is the infoset of the validated document augmented with data type information.

The concept of PSVI stems from the data-orientedness of XSD—as opposed to document-orientedness. XSD is biased towards use cases that involve serializing objects or database items as XML and, on the other hand, databinding which involves doing the reverse. Despite the data type focus of XSD, the data type system is not extensible. The schema author has to get by with the data types that the specification provides. For this reason it is not uncommon to see schemata that do not constrain the data type of a given attribute, even when the attribute has a very specific format.

XSD is rather verbose and, therefore, inconvenient to write. After a compact syntax for RELAX NG (page 25) was developed and found useful, a compact syntax for XSD was developed by Kilian Stillhard.[CompactXSD] However, this format is not standardized and appears not to have gained wide acceptance.

XSD is anecdotally notorious for partial implementation and implementations that are not interoperable. The anecdotes gained more credibility after a workshop on XML Schema 1.0 User Experiences and Interoperability[SchemaUE] was held by the W3C. Reports about the problems faced with XSD were summarized by Rick Jelliffe on the xml-dev mailing list in [Freddy].

The failures of XSD are discussed in more detail in [IntroXML].

XSD was not used in this project due to its questionable reputation, its verbosity and its data-orientedness and relative lack of acceptance in document-oriented tasks.

### 3.3 Document Structure Description

Document Structure Description (DSD) is a schema language developed at the University of Aarhus and AT&T Labs Research. In [Taxonomy], DSD 1.0 is classified as “single-type” and DSD 2.0 is classified to be able to represent any regular tree grammar.

DSD was not used in this project due to the wider acceptance of RELAX NG and Schematron and the tool availability situation. DSD has largely been sidelined in the marketplace with XSD on one hand and the DSDL family of languages (see below) on the other hand being the main contestants.

### 3.4 TREX, RELAX, XDuce and DDML

There have been various also-ran schema languages in the quest for a replacement for DTDs. Perhaps the most famous ones are TREX by James Clark and RELAX by Makoto Murata. These languages were the basis of RELAX NG (discussed below), which has superseded them. Document Definition Markup Language (DDML) was published as a W3C Note but was abandoned in favor of XSD. XDuce has not gained wide acceptance.

### 3.5 RELAX NG

RELAX NG is a grammar-based schema language for XML. In [Taxonomy], it is classified to be of the most powerful kind of tree grammars: “regular”. It was developed from the basis of TREX and RELAX in OASIS (Organization for the Advancement of Structured Information Standards) with James Clark and Makoto Murata (the developers of TREX and RELAX respectively) as the editors of the specification. RELAX NG has also been subsequently standardized as Part 2 of the Document Schema Definition Language (DSDL) family[ISO19757-2].

James Clark describes RELAX NG as an evolution and generalization of XML DTDs based on experience both from SGML and XML. Design patterns used for writing DTDs can be applied to RELAX NG. Moreover, DTD is can be programmatically converted into RELAX NG. [RNGdesign]

RELAX NG treats elements and attributes in a uniform way to the extent possible. This means that co-occurrence constraints between attributes and the content model of an element are possible.

RELAX NG is strictly for validation. No info-set augmentation is performed. Since there is no need to ambiguously assign augmentation to the info-set, ambiguous grammars are allowed. A document is a valid according to RELAX NG schema if there is at least one derivation for the document in the grammar expressed by the schema. It does not matter if there are multiple derivations. Allowing ambiguous grammars makes RELAX NG schemata easier to write than DTDs or XSD. (Note: There exist tools that use RELAX NG purposes other than validation and these tools tend to place

restrictions on grammar ambiguity. Still, ambiguous grammars are allowed by the core specification.)

### 3.5.1 Datatyping

RELAX NG has only two built-in data types: `string` and `token`. Additionally, there is a `list` pattern that allows data types to be used in white space-separated lists. The `string` type allows all strings that are legal in XML. Also, a schema may enumerate permissible `string` values that are compared using strict code point for code point comparison. The `token` data type normalizes whitespace before such comparison.

In addition to the built-in datatypes, RELAX NG has a framework for pluggable datatype libraries. A datatype library makes it possible to use a Turing-complete programming language for checking whether a string conforms to a particular datatype. In formal terms, a datatype (with given parameters) is a formal language and an equivalence relation for strings of the language.

Each possible string of XML characters either belongs to the language of the datatype or does not belong in the language. For example, a datatype for dates could accept strings that represent valid Gregorian dates in the ISO 8601 notation and reject all other strings.

At minimum, each string needs to be equivalent with an identical string (reflexivity). However, the equivalence relation may be more lax as long as it is transitive and symmetric. For example, a datatype might accept all possible strings and define the equivalence relation as the equality of the Unicode fold case transformations of the operands.

It is important to note that the RELAX NG notion of datatypes only concerns classifying strings. Unlike the non-extensible XSD type system and the PSVI concept, RELAX NG datatyping is not about converting the strings to datatypes of a programming language (integers, floats, date objects, etc.). Special-purpose tools built on top RELAX NG could use datatyping in such a way (if they restrict grammar ambiguity), but such usage is not sanctioned by the RELAX NG specification.

The RELAX NG specification says how datatype libraries are used from a schema. However, the RELAX NG specification does not specify an API for interfacing a datatype library implementation with a RELAX NG validator, because such an API needs to be specific to the programming language used for implementation and RELAX NG does not require any particular programming language. However, for Java there is a de facto standard datatype library API developed by James Clark and Kohsuke



Kawaguchi.[DatatypeAPI] The Java API has been also been adapted to other languages—C# in particular.

There are two well-known datatype libraries: the XSD datatype library [RNG-XSD] and the DTD compatibility datatype library [DTDCompat]. The former brings the datatypes from [XSDDatatypes] to RELAX NG. The latter brings the datotyping features of DTDs to RELAX NG. RELAX NG validators often have built-in support for these two datatype libraries.

### 3.5.2 Compact Syntax

RELAX NG is defined as an XML vocabulary. However, since the XML syntax is designed for marking up text, it is not particularly convenient to write or even read in cases where there's almost no text and a lot of markup.

To address this problem, RELAX NG has an alternative Compact Syntax [Compact]. The compact syntax is vastly more human-friendly than the XML syntax and is intuitive to anyone familiar with the customary notation for regular expressions and the Backus–Naur Form. Like RELAX NG proper, the Compact Syntax was specified by OASIS and has subsequently been adopted as an amendment to the ISO standard [ISO19757-2Amd1].

### 3.5.3 Use in This Project

RELAX NG was chosen as the main schema language for this project because of its status as the schema language of choice for document-oriented tasks and because the a schema project and a validator project were already in place. The Compact Syntax was chosen due to its human-friendliness.

## 3.6 Schematron

Schematron is an assertion-based schema language. It was developed by Rick Jelliffe at the Academia Sinica Computing Centre (ASCC). A newer version of Schematron has been standardized as Part 3 of the Document Schema Definition Language (DSDL) family[ISO19757-3]. The ISO version of Schematron is incompatible with processors for the ASCC versions of Schematron.

A Schematron schema consist of assertions. In practice the assertions are XPath[XPath] expressions that evaluate to true or a non-empty node set when successful (i.e. not in error). Alternatively, the negation of this case can be treated as the error condition (i.e. an error is reported if an expression evaluates to true of a non-empty node set).

This is significantly different from grammar-based schema languages. In order for a document to conform to grammar-based schema, there has to be a derivation for the document in the grammar. Therefore, in a grammar-based schema, by default everything is forbidden and only the constructs that can be derived from the grammar are allowed. In Schematron, however, everything is allowed by default and each grammar makes a specific restriction.

### 3.6.1 Using RELAX NG and Schematron Together

Adding specific restrictions without having to take a stance on the document as a whole makes Schematron ideal for refining a cruder schema written in another language. The creator of Schematron, Rick Jelliffe, has characterized Schematron as “a feather duster for the furthest corners of a room where the vacuum cleaner cannot reach”[SchematronOld]. Indeed, in the last three years, a pattern of using RELAX NG and Schematron together has emerged: a slightly over-permissive RELAX NG grammar is used for the bulk of the schema and Schematron assertions are used to tighten corner cases. For example, the schemata for the Atom syndication format[RFC4287] and DocBook 5.0[DocBook] are RELAX NG schemata that are refined with Schematron assertions.

The RELAX NG schema and the Schematron schema can be separate or combined. If they are separate, the document being validated is simply validated against both the RELAX NG schema and the Schematron schema separately and considered valid if it passes both validations.

In the combined case, the Schematron assertions are written inside the RELAX NG schema. The validation phases are still separate, but the Schematron assertions can be organized so that they appear in the relevant element context in the RELAX NG schema for a human reader. An implementation may use the RELAX NG schema context of an assertion to establish the XPath context in Schematron. [Relaxtron]

### 3.6.2 Use in This Project

Following the example of Atom and DocBook, Schematron was adopted as the secondary schema language in this project for expressing details that are inconvenient or impossible to express in RELAX NG.

Schematron 1.5[Schematron15] was chosen for this project instead of ISO Schematron due to lack of tool support—in particular lack of support in the Jing engine.

# Chapter 4

## Prior Work on Markup Checking

The service presented in this thesis is not by any means the first markup checking service on the Web. In this chapter, notable other markup checking services are reviewed.

### 4.1 The W3C Markup Validation Service

The W3C Markup Validation Service (better known as the W3C Validator), originally written by Gerald Oskoboiny, is probably the best known of the markup checking services reviewed here. For many users, it is *the* validator. It has been in use since late 1990s.

The W3C Validator is a Perl CGI front end for OpenSP. OpenSP is an SGML parser based on James Clark's SP. OpenSP performs DTD-based validation according to [ISO8879] (SGML). That is, the input document is validated against the DTD that the document declares for itself. The front end allows the user to override the DTD declared by the document, in which case the front end modifies the document accordingly before passing it to OpenSP. The HTML 4.01 and XHTML 1.0 specifications come with normative DTDs. Typically documents include one of the normative DTDs by reference, but a document can include any DTD.

As a matter of long-running policy, the W3C Validator sticks strictly to the SGML validity formalism. It is often argued that it would be inappropriate for a program to be called a "validator" unless it checks exactly for validity in the SGML sense of the word—nothing more, nothing less. Markup language specifications virtually always contain conformance requirements that cannot be expressed in an SGML DTD. Those requirements are simply not checked for at all. For example, in HTML 4.01 Strict, the value of the `rowspan` attribute is required to be a non-negative integer, but since SGML

DTDs cannot express this constraint, any string passes as a valid value for `rowspan`. However, due to the ways validation is often promoted, users tend to think that the W3C Validator checks for more than what it actually does. In fact, the front page of the W3C Validator states “This is the W3C Markup Validation Service, a free service that checks Web documents in formats like HTML and XHTML for conformance to W3C Recommendations and other standards.”

Another problem is related to XML support. XML 1.0 was designed to be compatible with SGML in the sense that an XML document that is valid according to its DTD when treated as XML is also a valid SGML document when the Annex K of the SGML standard is in effect. The opposite is not always true, however. A document can be valid for the purposes SGML without even being well-formed from the XML point of view. As a trivial example, SGML does not require white space between attributes but XML does. As a result, tools designed for SGML are not appropriate for checking XML correctness. Yet, due to the way the W3C Validator is presented to users, users end up using an SGML tools thinking that they are running a tool that checks for XML validity. When giving results for XML documents, the W3C Validator states briefly “Note: The Validator XML support has some limitations.”

The SGML validation process requires a different SGML declaration for XML than for HTML 4.01. The choice of SGML declaration is external to the document. However, to avoid asking the user the rather esoteric question of which SGML declaration to use, the W3C Validator uses heuristics to decide which SGML declaration to use.

The code of the W3C Validator was not used as a starting point for the software discussed in this thesis, because the W3C Validator is a DTD-based validator for SGML documents written in Perl and C while the software discussed in this thesis started out as RELAX NG-based validator for XML documents with key libraries written in Java.

## 4.2 WDG HTML Validator

The Web Design Group HTML Validator[WDG] developed by Liam Quinn is very similar to the W3C Validator. It too has been in use since the late 1990s, has James Clark’s SP inside and has a Perl front end.

Originally, the WDG HTML Validator was differentiated from the W3C Validator by better error messages, by support for hexadecimal character references and by support for more character encodings than just ISO-8859-1 [WDG1998]. The W3C Validator has later added all these features as well.

Currently, the WDG Validator is differentiated from the W3C Validator by warning about certain SGML markup minimization features that do not work in real-world browsers, by warning about character references to C1 control characters and by using a different SGML declaration with custom DTDs [WDG2007]. Also, the WDG HTML Validator can spider a site and validate multiple pages on one invocation.

Like the W3C Validator, the WDG HTML Validator is restricted to the SGML validity formalism.

The code of the WDG Validator was not used as a starting point for the software discussed in this thesis due to the same technical reasons that apply to the code of the W3C Validator.

### 4.3 Page Valet

Page Valet is a DTD-based validator developed by Nick Kew. Page Valet uses OpenSP for validating HTML as SGML. However, unlike the W3C and WDG validators, Page Valet uses a real XML parser—Xerces-C—by default for validating XML. (The option to use OpenSP for XML is offered.)

Page Valet has an experimental option for turning on XSD-based validation in Xerces-C. However, there is no user interface for providing a schema separately from the document. That is, it is up to the document to specify its own schema.

For SGML-based validation, Page Valet provides three parse modes: Strict, Web and Fussy. The Strict Mode does what the W3C Validator does. It adheres strictly to the *de jure* formalism even when the results are impractical considering browsers. The Web Mode is described to be similar to what the WDG HTML Validator does. That is, SGML markup minimization features that do not work in real browsers are flagged. The Fussy Mode is described to “add further checks over and above Web Mode”. [ValetMode]

Unlike the W3C and WDG validators, Page Valet is not a Perl CGI program. It is implemented as a C-language Apache module called `mod_validator`[`mod_validator`]. `mod_validator` was not used as a starting point for the software discussed in this thesis, because using RELAX NG within a C program would have in practice required using libxml2 instead of Xerces-C. Moreover, this author is more comfortable with using a managed runtime than C for Web-facing services.

## 4.4 The Schneegans XML Schema Validator

The XML Schema Validator (formerly XHTML Schema Validator) by Christoph Schneegans validates XML documents against XSD schemata. It does not validate HTML documents.

While the three DTD-based validators discussed above use any DTD that the document declares for itself, the XML Schema Validator has a closed list of built-in XSD schemata. The user can choose a schema from a list manually or request the validator to choose a built-in schema based on the `xsi:schemaLocation` or the doctype. The schemata offered include the three variants of XHTML 1.0, XHTML 1.1, XHTML 1.0 Basic, MathML 2.0, XSD itself and Google Sitemaps.

The schemata for the variants of XHTML as well as XSD itself come from the W3C. The schemata for XHTML 1.0 were published as a W3C Note [XHTML10XSD]. The schemata for modularized XHTML were published as a Proposed Recommendation [M12N11] (later changed back to Working Draft).

The XML Schema Validator is written in Visual Basic .NET and is based on the XML tool chain provided as part of the Microsoft .NET Framework 2.0. The source code is not available.

## 4.5 Relaxed

Relaxed by Petr Nálevka validates documents against RELAX NG and Schematron schemata. The main advantage of Relax compared to the W3C Validator is the ability check for requirements that cannot be expressed in DTDs. Relaxed was originally presented in Czech in Nálevka's thesis [Validace], which I am unable to read. Relaxed has subsequently been described in English in a paper co-written by Nálevka and his thesis instructor Jirka Kosek [Relaxed].

Relaxed builds upon James Clark's *Modularization of XHTML in RELAX NG* [M12N-RNG]. The schemata developed by Clark have been further refined and augmented with Schematron assertions. In addition to Schematron assertions based on the conformance requirements of XHTML 1.0, Relaxed offers optional Schematron-based checks for some of the *Web Content Accessibility Guidelines 1.0* [WCAG] requirements. In later updates after the initial release, schemata for compound documents that embed SVG and MathML in XHTML have been added.

Relaxed is written in Java and JSP. It uses the Sun Multi-Schema Validator (MSV) by Kohsuke Kawaguchi as its validation engine. Even though there

is a plug-in for MSV that enables support for Schematron assertions that are embedded in RELAX NG, Relaxed uses a separate XSLT-based solution where the Schematron part is first extracted from RELAX NG using XSLT, then the resulting Schematron schema is compiled into an XSLT script using another XSLT transformation and finally the resulting XSLT script is run against the input document.

Relaxed offers a list of preset schemata. Custom schemata are not allowed. Preset schemata are provided for XHTML 1.0, optionally with SVG and MathML. The schemata for XHTML 1.0 can also be used for HTML 4.01. For XHTML 1.0 and HTML 4.01 without SVG or MathML, partial WCAG checks are available. A separate user interface is provided for a prerelease version of DocBook 5.0.

The HTML support in Relaxed is based on the idea that HTML 4.01 can be mapped to XHTML 1.0 before applying XML validation technologies. There is a problem, though. Relaxed uses John Cowan's TagSoup for the conversion. TagSoup has not been designed for error detection. On the contrary, TagSoup has been designed never to report errors of any kind. Therefore, tokenization-level errors may go unreported. Moreover, the document is serialized as XML and then reparsed with an XML parser, which may cause differences in line numbers.

Of the validation services reviewed in this chapter, Relaxed is the closest to the system described in this thesis. The reason why the work described in this thesis is not based on the Relaxed code base is that the validation service upon which the HTML5 conformance checker is an elaboration was first deployed in the spring of 2005 but Relaxed was announced later in August 2005. Effort was expended in order to test MSV in place of Jing. However, the result was that MSV was easier to crash with user-supplied schemata. Also, Jing supports the RELAX NG Compact Syntax but MSV does not. Moreover, I prefer using a SAX pipeline as the output generation method instead of JSP.

Even though Relaxed and the service described in this thesis do not share Java code, the schemata developed for Relaxed are used to implement the pre-HTML5 functionality of the service whose HTML5 part is discussed in this thesis.

## 4.6 Feed Validator

The Feed Validator by Sam Ruby, Mark Pilgrim, Joseph Walton, and Phil Ringnalda is notably different from the services reviewed above. Whereas the other services focus on HTML and/or XHTML, the Feed Validator

focuses on Atom and RSS feeds, which are rather different from HTML and XHTML as markup languages. Also, the methodology used in the Feed Validator is significantly different.

The Feed Validator does not use any schema formalism. Instead, it uses hand-crafted SAX consumers written in Python. The main `SAXContentHandler` maintains a stack of element-specific delegates. Each element-specific delegate (inheriting from a common base) can check the conformance requirements pertaining to its element.

The approach taken in the Feed Validator has two benefits over schema-based validation. First, the Feed Validator is not limited by the capabilities of any schema formalism. Since Python is a full programming language, any requirement that is machine checkable in principle is checkable by a Python program. Second, since the emission of error messages is programmed by humans on a case-by-case basis, the messages can be as good and informative as the human developers care to make them. In the case of grammar-based schema languages in particular, error messages are generated by the validation engine and context-sensitive advice is generally not provided except by perhaps applying guesswork outside the validation engine.

It should be noted that feeds are different from (X)HTML in the sense that there is a greater focus on string values adhering to certain formats and the nesting structures of elements are less complex than in (X)HTML. The Feed Validator does not validate (X)HTML content embedded in feeds.

The Feed Validator methodology was not chosen for this project, because it was assumed that using a domain-specific non-programming language—RELAX NG—would be more manageable in terms of effort and malleability during the development process of HTML5. Even though Feed Validator-style methodology was not chosen as the only methodology, the non-schema checkers (page 48) developed in this project are similar to the `ContentHandler` delegates used by the Feed Validator.

## 4.7 Validome

Validome by Thomas Mell, Vadim Konovalov, Alex Leporda, Olivier Duffez, Eduard Schlein and Dirk Klar checks both (X)HTML and feeds. Additionally, Validome offers generic XML validation against a DTD or an XSD schema declared by the document itself. Validome also offers to check DTDs and XSD schemata for syntax errors.

In contrast to the other services that are English-only or also provide messages in French, Validome supports German and English as the user



interface languages. The (X)HTML facet of Validome also offers French and Russian as user interface languages.

Validome uses SGML DTDs for HTML and XSD for XHTML. Additionally, Validome performs non-schema checks, which sets it apart from the other (X)HTML validators discussed above.

It is difficult to review what exactly Validome does, because its inner workings are publicly documented and the source code is not available. Still, it is worth emphasizing that Validome goes beyond schema formalisms when specifications have conformance requirements that cannot be expressed as schemata.



# Chapter 5

## Implementation

This chapter starts the second part of this thesis which focuses on the software implemented as the experimental part of the thesis project.

### 5.1 The Basic Back End

The basic architecture of the back end of the conformance checker is very simple: A parser consumes the document byte stream and emits SAX parse events. An arbitrary number of SAX event consumers can listen to the events. The SAX consumers do not form a pipeline with one consumer emitting events to the next. Instead, the for  $n$  consumers, the SAX event stream is split  $n - 1$  times.

The consumers can be schema-based validators or custom code. The consumers receive SAX parse events that affect the meaning of the document. That is, syntactic sugar is not exposed. For example, comments are not reported and it is not exposed if and how characters were escaped in the source.

The meaningful SAX events are the ones reported to the SAX `ContentHandler` and `DTDHandler` interfaces. However, in the events reported to `DTDHandler` (notifications of notations and unparsed entity declarations) are of no interest for an (X)HTML5 conformance checker, so in practice the SAX consumers only listen to `ContentHandler`. Moreover, the processing is genuinely namespace-aware and, hence, qualified names are not used by event consumers. Instead, the consumers observe namespace URIs and local names.

`LexicalHandler` is not listened to, because exposes syntactic details that do not affect the meaning of the document and it would be inappropriate to tie conformance on the higher layer to the choice of syntactic sugar on

the lower layer. `DeclHandler` is not listened to, because it exposes declarations that are intended to be expanded by the XML processor and, thus, should not be managed by the application.

The main SAX event consumer is an instance of the Jing validation engine that has been instantiated with a RELAX NG schema for HTML5 or XHTML5.

## 5.2 The Jing Validation Engine

Jing by James Clark was chosen as the RELAX NG and Schematron validation engine. Jing implements James Clark's derivative algorithm [Derivative] for RELAX NG validation. For Schematron, Jing provides a wrapper for an XSLT engine. The SAXON XSLT engine by Michael Kay was chosen.

## 5.3 The RELAX NG Schema

The bulk of what is allowed in HTML5 is encoded in a RELAX NG schema. The core of the schema was developed by Elika Etemad [Whattf].

For the most part, the element nesting conformance requirements for HTML5 are defined in terms of simple parent-child relationships. That is, it is defined that a given element may have children of a given type. This sort of requirements map trivially to a grammar.

DTDs allow only one grammar production per element name. HTML5, however, requires a different grammar productions in different contexts. For example, the attributes allowed on list items depend on the kind of list the items are in. Fortunately, RELAX NG decouples grammar productions from element names. There can be an arbitrary number of grammar productions for a given element name. This makes it possible for elements to have a different attributes or content models depending on where the elements occur in the document. For example, the `form` element has an empty content model when it occurs as a child of the `head` whereas it can have child elements when it occurs as a descendant of the `body` element.

In HTML5, some elements are defined to take either block-level children or inline-level children. Such content models are called bimorphic. In HTML 4.01 those elements generally allowed a mix of block and inline content, because DTDs cannot express the kind restriction demanded by HTML5. Fortunately, RELAX NG can deal with bimorphic content models. When a RELAX NG validator sees an element, the element can have multiple

pending derivations in the grammar. Therefore, adjacent subtrees can affect each other in validation.

The vast majority of the conformance criteria related to the document structure can be expressed as RELAX NG grammar. In the implementation, the use of RELAX NG has been favored whenever practical. Still, there are many conformance criteria that are not conveniently expressible in a RELAX NG grammar. Schematron and custom Java code are used for those criteria.

Expressing constraints related to ancestry turned out to be impractical in RELAX NG, even though it is theoretically possible. Balancing the needs of the main conformance checking use case and the expected reuse of the RELAX NG schema for e.g. guiding auto-completion in an editor turned out to be a problem. From the conformance checking point of view, it makes sense to handle more things in Schematron, because special cases of constraints on element ancestry are trivially expressed in Schematron with precise error messages whereas “remembering” ancestry in a grammar generally requires duplicating grammar productions, which leads to unmanageable growth of the number of grammar productions. Eventually, Schematron was favored.

### 5.3.1 The General Schema Design

The schema is divided into modules. The module division is motivated both by organizing the schema by grouping similar elements together and by making it possible to easily subset the schema in ways deemed reasonable by the schema authors. The schemata for different subsets include the modules that are enabled for the particular subset.

### 5.3.2 Common Definitions

A module called `common.rnc` includes definitions for the schema framework. It contains switches for parameterizing the schema, initial definitions for common content models, definitions for common attributes and definitions for common data types. It also designates the grammar start symbol, i.e. the root element.

#### 5.3.2.1 Common Content Models

The `common.rnc` module initializes the definition for common content models like this:

```
common.inner.strict-inline =  
  ( text )
```

```
common.inner.struct-inline =
  ( text )
```

```
common.inner.block =
  ( empty )
```

Other modules then add to these definitions like this. For example, the following makes the `p` element allowed where block elements are allowed.

```
common.inner.block &= p.elem*
```

The `&=` operator redefines the named pattern on the left-hand side to be the interleaving of right-hand side and the previous definition of the left-hand side.

The `common.rnc` module also contains content model definitions derived from other content model definitions:

```
common.inner.bimorphic =
  ( common.inner.struct-inline
    | common.inner.block
  )
```

The derived definitions are not augmented directly by other modules.

### 5.3.2.2 Common Attributes

HTML5 defines a handful of attributes that apply to all HTML elements. A pattern called `common.attrs` is defined in `common.rnc`. The pattern serves as a reusable pattern for element definitions and provides an extension point for modules that add attributes that apply to all elements. For example, including the module for scripting adds event handler attributes to all elements.

### 5.3.2.3 Common Datatypes

Many attributes in HTML5 take values that are required to conform to a specific format. In RELAX NG such formats are known as datatypes. As discussed before, datatypes in RELAX NG only constrain the space of allowable string values and do not imply info set augmentation or data binding. Except for datatypes related to HTML forms and enumerated string values, the datatypes used to constrain attribute values are defined in `common.rnc`.

When possible, the W3C XML Schema Datatypes[RNG-XSD] are used in order to make the schema more easily portable to different RELAX NG

validators. In particular, the regular expression facet of the W3C XML Schema Datatypes is used. For example, percentage values are defined as follows:

```
common.data.percent =
  xsd:string {
    pattern = "(100)|([1-9]?[0-9](\.[0-9]+)?)%"
```

The XSD regular expressions have some unconventional quirks that require attention when writing schemata. For example, the `\d` shorthand is conventionally defined to mean an ASCII digit, that is U+0030 DIGIT ZERO though U+0039 DIGIT NINE, but in XSD regular expressions `\d` is defined to match any character that is classified as Nd (decimal digit) in Unicode. The XSD definition is politically correct but less useful in practical cases. Hence, `[0-9]` is used in the example above.

When the permissible lexical space of a datatype does not form a regular language or when the lexical space would in theory form a regular language but writing it down as a regular expression would be impractical, datatypes from the custom-built HTML5 Datatype Library (discussed later in this chapter) are used. For example, datatypes involving dates and IRIs are handled using the custom-built library:

```
common.data.datetime =
  w:datetime-tz
```

```
common.data.uri =
  string "" | w:iri-ref
```

It is important to note that using a custom datatype library makes the schema less portable. To use the schema with, the RELAX NG validator needs to have an implementation of the datatype library available to it. An alternative less precise version of the schema could be made if portability was appreciated over correctness.

#### 5.3.2.4 Parameter Switches

RELAX NG has two special patterns that can be used to implement boolean feature switches. These patterns are `empty` and `notAllowed`. The `empty` pattern takes no attributes or element content to satisfy. The `notAllowed` is never satisfied. Hence, interleaving `empty` with another pattern is equivalent to the other pattern alone and interleaving `notAllowed` with another pattern makes the interleave as a whole unsatisfiable. Thus, a named pattern

can be used in a schema and the effect of the pattern can be changed by changing the definition of the named pattern from `empty` to `notAllowed` as needed. RELAX NG makes this easy by allowing a schema file that includes another to override definitions in the file that is being included.

For example, the `p` element has a more versatile content model in XHTML5 than in HTML5. In XHTML5, structured inline children—that is inline mixed with selected primarily block elements like `ul`—are allowed. (The HTML5 serialization cannot allow elements that were block elements in HTML 4.01 as children of `p` due to backwards compatibility considerations.)

This distinction is handled using a switch called `nonHTMLizable`. It is defined in `common.rnc` as `nonHTMLizable = empty` which is what is needed for XHTML5. To flip the switch for HTML5, `common.rnc` is included as follows:

```
include "common.rnc" {
  nonHTMLizable = notAllowed
}
```

The switch pattern is then used like this:

```
p.inner =
  (
    common.inner.strict-inline
  | ( common.inner.struct-inline
    & nonHTMLizable
    )
  )
```

When `nonHTMLizable` expands to `notAllowed`, the right-hand side of `|` becomes unsatisfiable and `p.inner` becomes equivalent to `common.inner.strict-inline`.

When `nonHTMLizable` expands to `empty`, `p.inner` becomes equivalent to `common.inner.struct-inline` because `common.inner.strict-inline` is defined to be a subpattern of `common.inner.struct-inline`.

### 5.3.3 Examples of Elements

The RELAX NG implementation for a typical element looks like this:

```
blockquote.elem =
  element blockquote { blockquote.inner & blockquote.attrs }
blockquote.attrs =
  ( common.attrs
```



```

    & blockquote.attrs.cite?
  )
  blockquote.attrs.cite =
    attribute cite {
      common.data.uri
    }
  blockquote.inner =
    ( common.inner.block )

```

The element is given a named definition: `blockquote.elem`. The definition simply expands to an element pattern for the element in question. The content model is defined to be an interleaving of two named patterns: one for element content (`blockquote.inner`) and another for attributes (`blockquote.attrs`).

The named pattern for attributes is defined as the interleaving of common attributes and element-specific attributes. In this case, the only element-specific attribute (`cite`) is optional, hence the `?` quantifier. Next, the named pattern (`blockquote.attrs.cite`) for the element-specific attribute is defined. The datatype for the attribute refers to a common datatype definition.

The named pattern for the element content (`blockquote.inner`) is merely defined to map to the common content model for elements that accept only block children (`common.inner.block`).

Other elements are defined analogously. Of course, the definitions for other elements tend to become more complex. For example, this is the definition for `datetime` form controls:

```

input.datetime.elem =
  element input { input.datetime.attrs }
input.datetime.attrs =
  ( common.attrs
    & common-form.attrs
    & input.datetime.attrs.type
    & common-form.attrs.accesskey?
    & input.attrs.autocomplete?
    & common-form.attrs.autofocus?
    & input.attrs.list?
    & input.datetime.attrs.min?
    & input.datetime.attrs.max?
    & input.attrs.step.float?
    & common-form.attrs.readonly?
    & input.attrs.required?
  )

```

```

& input.datetime.attrs.value?
)
input.datetime.attrs.type =
  attribute type {
    string "datetime"
  }
input.datetime.attrs.min =
  attribute min {
    form.data.datetime
  }
input.datetime.attrs.max =
  attribute max {
    form.data.datetime
  }
input.datetime.attrs.value =
  attribute value {
    form.data.datetime
  }

```

```
input.elem |= input.datetime.elem
```

The notable detail in this more complex example is that the reference to `input.datetime.attrs.type` is not quantified. Hence, the `type` attribute with the value `datetime` is required. This definition for the `input` element is ORed together with the other definitions with different `type` attributes (`input.elem |= input.datetime.elem`). This means that the value of the `type` attribute serves as a discriminator for determining which patterns apply even though the name of the element remains the same.

## 5.4 The HTML5 Datatype Library

It was found that a custom datatype is needed in the following cases:

- When the lexical space of a datatype is not a regular language or when it is but formulating it as a regular expression would be particularly hard or inconvenient.
- When the lexical space of a datatype is a regular language, but the grammar would require a lot of literal strings (e.g. registered language codes or character encoding names) embedded into it.
- When checking the value requires calendar calculations.

### 5.4.1 Dates

On the surface, it would seem that the last case is unnecessary considering that the XSD datatypes include a `dateTime` type that can be constrained through its regular expression facet. However, the XSD `dateTime` is inappropriate for HTML5 in a subtle way. With the XSD type, leading and trailing whitespace is discarded before the pattern is matched, so there's no way of forbidding surrounding whitespace. The Web Forms 2.0 date and time types do not allow surrounding whitespace. Hence, custom types are needed.

It turns out that discarding surrounding whitespace as part of the datatype makes the datatype library less flexible, since any datatype that does not allow whitespace at all (neither surrounding the meaningful value nor inside the value) can be used in a way that allows surrounding whitespace by wrapping it in the RELAX NG `list` pattern so that a single-token list of whitespace-separated tokens results.

### 5.4.2 IRIs

In addition to dates, it turns out that custom types for IRIs and language tags are also called for, even though the repertoire of XSD types includes datatypes for these.

The definition of the XSD `anyURI` datatype cannot be considered stable or useful. Its definition has changed with each edition and version of the specification. In the first edition of version 1.0 [XSDDatatypesFE], which predated the IETF IRI specification, the definition implied that not all strings are valid `anyURI` values. Indeed, the Jing implementation does not allow all possible strings. It is unclear, though, whether the definition actually accidentally allowed more than it was thought to allow. In the second edition of version 1.0 [XSDDatatypes], which was finalized after the Jing implementation had been released, the definition suggested that different implementation levels could treat different lexical spaces valid. In a February 2006 working draft of the 1.1 version [XSDDatatypes11WD], the definition concedes that any finite string is a valid `anyURI`.

For these reasons, custom data types for IRIs (only absolute) and IRI references (either relative or absolute) were developed. The major problem with deciding whether a string is a conforming IRI is that the generic IRI syntax is not particularly restrictive. Instead, individual IRI schemes, such as `http`, `mailto` and `ftp`, are allowed to specify their own syntax. And implementation working only on the generic IRI level would pass just about any string as a valid IRI reference. On the other hand, it would be impossible to

implement scheme-specific knowledge of future or private IRI schemes. Even implementing support for all the IANA-registered IRI schemes would be impractical. In fact, these problems were the reason (in addition to the XSD datatypes predating the IRI specification) for the fuzzy definition of `anyURI`. Still, it would be a pity not to help authors with a scheme-specific issues related to the most commonly used schemes—the `http` scheme in particular.

The obvious solution is to implement scheme-specific checking for the most common schemes and apply it only generic processing to the rest. However, leaving the decision on what scheme-specific checking to support to implementations of the abstract datatype library specification would make implementations uninteroperable. It was decided to include scheme-specific knowledge of the following IRI schemes: `http`[RFC2616], `https`[RFC2818], `ftp`[RFC1738], `mailto`[RFC2368], `file`[RFC1738] and `data`[RFC2397]. The selection of these IRI schemes was based both on which scheme are specifically covered by the Web Forms 2.0 specification and on which schemes are supported by the Jena IRI library around which the datatype library puts a thin wrapper to implement IRI checking. A reasonable candidate for including on the list would be the `javascript` pseudo-scheme. Support for it would require using the Rhino JavaScript library for syntax checking instead of using the Jena IRI library. The Jena IRI library does not offer satisfactory support for the `mailto` and `data` schemes. Adding support for these IRI schemas was not attempted within the scope of this thesis project but was left as possible future work (page 69).

### 5.4.3 Language Tags

The HTML `lang` attribute and the XML `xml:lang` attribute take as the value a language tag identifying a human language. Until recently, language tags were defined by [RFC3066]. Now that specification has been obsoleted by [RFC4646] and [RFC4647]. A language tag consists hyphen-separated subtags.

Different subtags have different lengths. Moreover, the permissible list of values is restricted to private use subtags (starting with `x-`) and values that are in the IANA subtag registry. The XSD definition for `language` is significantly more coarse: the lexical space is defined to be 1–8 ASCII letters followed by zero or more hyphen-separated subtags consisting of 1–8 ASCII letters and/or digits.

A custom datatype with built-in datatype can do much better if it implements checking for the syntax defined in [RFC4646] and also contains data from the IANA registry. This approach was chosen and partially

implemented. However, the implementation was not finished within the scope of this thesis project.

Making the lexical space of a datatype dependent on a registry that changes over time poses a versioning problem. However, as long as implementations document when their IANA registry snapshot was taken, introducing new language tags will not be a significant problem. Introducing new values to the registry will expand the valid lexical space without making any previously conforming documents non-conforming. This is in contrast with the IRI scheme issue. Introducing support for a previously unsupported IRI scheme would narrow the valid lexical space.

#### 5.4.4 ECMAScript Regular Expressions

Web Forms 2.0 introduces an attribute called `pattern` for form controls that accept textual input. The attribute specifies a regular expression that the value of the form control must match. Web Forms 2.0 does not define its own regular expressions. Instead, ECMAScript regular expressions are used. When `^(?:` is prepended to the value of the `pattern` attribute and `)$` is appended to it, the resulting string is required to be a conforming ECMAScript regular expression. The datatype library handles the checking ECMAScript regular expression syntax by wrapping the regular expression parser of Rhino.

### 5.5 The Schematron Schema

A Schematron schema was used where RELAX NG did not work conveniently and where using custom Java code was not strictly necessary. Surprisingly, the use cases for Schematron turned out to be narrower than originally expected.

#### 5.5.1 Exclusions

While grammars are very good at enforcing parent-child relationships, they are not good at enforcing ancestor-descendant relationships, because in the absence of intersections and negations, some the productions would have to “remember” what kind of ancestors of interest there have been. Suppose element A is not allowed to have element B as a descendant. In this case, all the elements that can occur on the ancestry path from B to A would need to have two grammar productions: one that can be derived from the root without an intervening A element and one that cannot. Moreover, if there

were more such exclusions, the number of parallel grammar productions per each neutral element would double per each exclusion rule. Obviously, this could be a serious maintainability problem.

Fortunately, exclusions are extremely easy to express in Schematron. The forbidden descendant element is used as the context node for an assertion, which then states that the context node does not have an ancestor element that could not have the context element as its descendant.

For example, this Schematron pattern causes `blockquote` elements that are descendants of header elements to be reported as errors.

```
<rule context="h:blockquote">
  <report test="ancestor::h:header">
    The blockquote element cannot appear as a
    descendant of the header element.
  </report>
</rule>
```

The Schematron wrapping of the XPath expressions is rather verbose, but the expressions themselves are simple. The expression `h:blockquote` matches the `blockquote` element in the XHTML namespace. (The prefix `h` is bound to the XHTML namespace.) The rest of the rule is only applied if the context expressions matches, and in that case, the `blockquote` element is used as the XPath context for the second expression `ancestor::h:header`. This expression matches header elements (in the XHTML namespace) that are also ancestors of the context node. If the set of matching nodes is non-empty, the natural-language error message is reported.

### 5.5.2 Required Ancestors

Opposite to exclusions, there are also checks for required ancestors. Specifically, the Web Forms 2.0 repetition model requires form inputs for moving and removing repetition blocks to have a repetition block or a repetition template as an ancestor.

The checks are of the following form:

```
<rule context='h:input[@type=move-down] '>
  <assert test='ancestor::h:*[@repeat] '>
    An input element of type="move-down"
    must have a repetition block or a
    repetition template as an ancestor.
```

```
</assert>  
</rule>
```

The context matches `input` elements in the XHTML namespace that have an attribute named `type` that has the value `move-down`. The assertion test matches ancestors of the context node that are in the XHTML namespace and have an attribute named `repeat`. This rule is correct for pure (X)HTML5 documents but is not sufficient for compound documents. Support for compound documents (mixing XHTML with e.g. SVG and MathML with the Web Forms 2.0 repetition model) was left outside the scope of this Master's Thesis project.

### 5.5.3 Referential Integrity

It turns out that RELAX NG is not good for enforcing referential integrity. Enforcing referential integrity means checking cases where an attribute value is required to be a reference to the ID of another element. RELAX NG has an optional extension called RELAX NG DTD Compatibility. This extension makes it possible to check that values designated as ID references actually referred to IDs in the same document. However, that's all it can check. Moreover, enabling this extension places restrictions on the ambiguity of the schema, which makes schemata harder to write in some cases. Due to the limitations of RELAX NG DTD Compatibility, all ID-related checking was moved away from RELAX NG. Most cases are handled in Schematron.

In Schematron, referential integrity checking builds on the XPath `id()` function. The argument of the function is coerced into a string which is split on whitespace. The return value is a node set containing the elements that have an ID equal to any of the tokens split for the argument.

The use of the `id()` function has two crucial differences compared to testing equality against the values of two attributes. First, the function operates on IDness and not on the names of attributes. Second, the argument is split on whitespace to produce a list of tokens, so the function works for the case where multiple IDs are referenced (IDREFs in DTD terms).

The concept of IDness is part of XML 1.0 itself. The IDness is established by processing a DTD that declares certain attributes to have the type ID. The DTD-based type of attributes is data that is exposed through the core SAX2 interface. More precisely, it is exposed via the `getType()` methods of the `Attributes` interface. The XPath implementation determines which attributes should be considered to be of type ID by calling one of the `getType()` methods.

In the case of the HTML5 serialization, there is no DTD processing whatsoever involved. The parser simply assigns IDness to the attribute named `id` and exposes this through SAX.

In the case of the XML serialization (XHTML5), IDness that is not based on DTD processing is assigned between the parser and the validation SAX listeners. Immediately after the parser in the pipeline, there is a SAX filter that constitutes an “xml:id processor” as per [xmlid]. Immediately after the xml:id processor, there is a SAX filter that performs similar ID assignment on attributes named `id` that are not in a namespace and belong to an element that is in the XHTML namespace. This second filter could be called an “XHTML id processor”.

With ID assignment performed before the Schematron stage, Schematron can be used to check that referents are of the right kind. For example:

```
<rule context='h:input[@list]'\>
  <assert test='id(@list)/self::h:datalist or
              id(@list)/self::h:select'\>
    The list attribute of the input element must
    refer to a datalist element or to a select element.
  </assert>
</rule>
```

## 5.6 The Non-Schema-Based Checkers

It turns out, as expected, that HTML5 has conformance requirements that cannot be expressed in RELAX NG or Schematron or that would be inconvenient to express in RELAX NG or Schematron. Checking for such requirements is handled by non-schema-based checkers.

Conceptually, a non-schema-based checker listens to parse events and does whatever is necessary and computable to identify the kind of conformance requirement violations that the checker is designed to handle. In practice, a non-schema-based checker is a Java class that implements the SAX2 `ContentHandler` interface and reports to a SAX2 `ErrorHandler`. Since non-schema-based checkers are implemented in a full programming language, they can check for any machine-checkable conformance requirement.

To make the non-schema-based checkers fit into the Jing-based architecture, a wrapper class that implement the `Jing Validator` implemented. With this wrapper class, the non-schema-based checkers conform to the same interface as the schema-based validators and can be combined into the same validation pipeline. Every non-schema-based checker is also given a



URI so that they can be instantiated by URI in the extended validation user interface just like validators instantiated from schemata identified by URIs.

The organization of checks for different requirements into different non-schema-based checkers is a matter of software design. After all, the specification does not mandate a particular code organization. If every requirement is implemented as a separate checker, there will be a lot of code duplication, since many checkers need to do similar things. On the other hand, implementing everything a single checker would make the code unmaintainable. The organization below reflects the design decisions of the author.

Developing checkers to cover all of HTML5 was deemed to be out of the scope of this thesis project. Instead, a selection of prototype checkers were implemented as a proof of concept.

### 5.6.1 Table Integrity Checker

A table integrity checker was chosen as the main proof of concept non-schema-based checker. Table integrity was deemed to be the most complex of the requirements that call for a non-schema-based checker making it ideal for proving by implementation that the approach works. Also, purely schema-based validators are incapable of checking table integrity and table integrity has notable relevance to table rendering in browsers, so table integrity checking makes for a good demo.

Since the table model for (X)HTML5 was only being specified when the checker was prototyped, the checker was speculatively based on the HTML 4.01 table model and browser behavior. The differences from HTML 4.01 are that `colspan='0'` is treated as `colspan='1'` and that `headers` must refer to `th` cells. The top left corner of cells is placed in the first available slot on the row, which is browser-compatible but different from what the CSS2 specification says.

An HTML table has a Cartesian grid of slots. A cell can span multiple slots. Subsequent cells are moved to the right until a free slot for the top left corner of the cell is found. When cells span multiple rows, this slot allocation policy can lead to overlapping cells.

The table integrity checker consists of 7 classes: `TableChecker`, `Table`, `RowGroup`, `Cell`, `ColumnRange`, `HorizontalCellComparator` and `VerticalCellComparator`. The `TableChecker` class is the non-schema-based checker class used by other code and the rest of the classes are internal to the checker.

`TableChecker` maintains a stack of `Table` instances. When a `startElement` event for the `table` element is seen, a new `Table` instance is pushed onto the stack. Likewise, the stack is popped upon seeing an

`endElement` event for `table`. The rest of table-significant events (`startElement` and `endElement` the `col`, `colgroup`, `thead`, `tbody`, `tfoot`, `tr`, `td` and `th`) are delegated to the `Table` object at the top of the stack.

`Table` maintains state of where in the table markup (e.g. in table row inside an explicit row group like `tbody`) the parse currently is. The `Table` object only sees table-significant parse events. If an event occurs out of the permitted sequence (e.g. a cell start occurs when the state is not in a table row) the whole subtree of misplaced elements is silently ignored at that point. Reporting it as an error is left to the RELAX NG validator.

`Table` maintains a linked list `ColumnRanges`. A `ColumnRange` represents a contiguous range of columns that has been established e.g. by the `col` element but does not yet have any cells starting in the range. A `ColumnRange` knows its start and end column slot indexes. In contrast to e.g. an array of column slots, the memory usage of this data structure is proportional to the number of ranges rather than the size of the ranges. Therefore, the memory usage can be throttled by limiting the number of elements that are processed, which can indirectly be throttled by limiting how large documents in term of bytes are processed. A malicious content author cannot make the checker allocate excessive amounts of memory by declaring a column group that spans an immense number of column slots.

`Table` also has the responsibility of instantiating `RowGroups`. The `thead`, `tfoot` and `tbody` elements establish explicit row groups represented as `RowGroup` objects. In XHTML, `tr` elements may occur as children of the `table` element. (In HTML, the parser infers a wrapping `tbody`.) This case is treated as an implicit row group also represented by `RowGroup`. Since table cells are assigned to slots on a per-row group basis in (X)HTML, the `Table` only instantiates `Cell` objects and gives them to the current `RowGroup` for slot allocation.

A `Cell` knows the column index where it starts and where it ends horizontally. It also knows the row until which it spans vertically. `RowGroup` maintains a set of cells that span more than one row and that are still in effect—that is, the first row onto which they do not span has not passed yet. As with `ColumnRanges`, the data structure does not involve allocating a two-dimensional array of slots. Therefore, the memory requirements are bounded in proportion to the number of cell elements rather than the number of slots that cells are declared to span.

When a `Cell` spans more than one row, it is added to the set of `Cells` that are in effect. When a row ends, `Cells` that are no longer in effect are culled from the set. When a new row starts, the `Cells` that are in effect are sorted according to their start column index. A new `Cell` on the row get its

top left corner assigned to the first free slot on the row. The slots occupied by the `Cells` that are still in effect (i.e. span down from earlier rows) are skipped. The `RowGroup` maintains the index of the current prospective slot index and the index of the next uninspected cell still in effect as part of its state, so the assignment of `Cells` to slots can progress from left to right without having to inspect the same cell still in effect more than once per row. If the set of `Cells` in effect is not empty when the row group ends, an error is reported for each `Cell` still in the set. (Cells cannot span across row groups.)

When `RowGroup` has allocated a `Cell` to particular slots, the `Cell` is reported back to `Table`, so that the list of `ColumnRanges` without cells can be adjusted accordingly. Since a linked list of `ColumnRanges` with no cells starting in them is maintained, the possible effects are: no change, removing a single-column range, narrowing a range by one and splitting a range. If there are `ColumnRanges` left in the linked list when the table ends, the ranges without cells starting in them are reported.

The checker emits both warnings and errors. Since the specification is still a draft, it is too early to tell which conditions should be treated as errors eventually. The rest will be just warnings.

The following conditions are detected:

- A table cell is overlapped by later table cell.
- A table cell overlaps an earlier table cell. (Single overlap gets reported in both directions to show source location for both cells.)
- A table cell spans past the end of its row group.
- A row has no cells starting on it.
- The column count on a table row is greater than the column count established by `cols/colgroups`.
- The column count on a table row is less than the column count established by `cols/colgroups`.
- The headers attribute does not point to th elements in the same table. (This feature was based on speculative information and will likely have to be revisited as the specification matures.)
- A column range has no cells starting on it.
- The value of a `colspan` attribute exceeds 1000, which is a magic number in Gecko (and according to comments in Gecko source, also in IE and Opera).
- The value of a `rowspan` attribute exceeds 8190, which is a magic number in Gecko.
- The column count on a table row is greater than the column count established by the first row in the absence of `cols/colgroups`.

- The column count on a table row is less than the column count established by the first row in the absence of `cols/colgroups`.
- A `col` element causes a `span` attribute to be ignored on the parent `colgroup`. (Conforming in HTML 4 / XHTML 1.0; non-conforming in (X)HTML5. With (X)HTML5 there's also a schema-level error.)

## 5.6.2 Checking the Text Content of Specific Elements

Some new elements in HTML5 have a special format for their text content. The `time`, `meter` and `progress` elements have a specific format for their text content. HTML5-aware user agents are required to parse the text content but the content serves as a fallback in legacy user agents. RELAX NG makes it possible to constrain the text content of an element if the element only has text content and no child elements. In HTML5, the elements with a specific text content format are allowed to have child elements. Hence, merely developing a custom datatype and using it from the RELAX NG schema would not work exactly the right way.

To make it possible to move the text content checking to RELAX NG in situations where custom datatype libraries are supported by non-schema-based checkers are not and text content checking is valued higher than child elements (for example, a Java-based RELAX NG-aware XML editor), checkers for the particular text content formats were implemented as classes that implement the `Datatype` interface from the Java API for custom RELAX NG datatypes. The non-schema-based checker then uses these `Datatype` implementations to check the text content of particular elements.

The Java API for custom RELAX NG datatypes supports streaming reporting of text content to a `Datatype` implementation. The `Datatype` interface has a method called `createStreamingValidator` which returns an object that implements an interface called `DatatypeStreamingValidator`. Character data can be reported to a `DatatypeStreamingValidator` in chunks the same way as the SAX interface handles reporting of character data. The non-schema-based checker simply maintains a stack of active `DatatypeStreamingValidator`s and reports the SAX character data to every `DatatypeStreamingValidator` on the stack.

The instantiation rules for the `DatatypeStreamingValidator`s are hard-coded although the non-schema-based checker could in theory be generalized to accept parameters from the outside stating which elements require which `Datatype` to be instantiated and the corresponding `DatatypeStreamingValidator` to be pushed onto the stack. However,

hard-coding these rules was simple and quite sufficient to meet the objectives of this project.

When a `DatatypeStreamingValidator` on the stack is about to be popped—that is, the end of the element is seen—the `DatatypeStreamingValidator` is queried whether it accepts the reported content or not. If it does not accept the reported content, the non-schema-based checker reports an error.

### 5.6.3 Checking for Significant Inline Content

HTML5 introduces a concept of significant inline content. Significant inline content consists of embedded content (such as an `img` element) or significant text. Significant text is defined to be text that contains any character that is not in the Unicode categories `Zs`, `Zl`, `Zp`, `Cc` and `Cf`.

HTML 4.01 stated: “We discourage authors from using empty `P` elements. User agents should ignore empty `P` elements.” As a result, markup generators started to generate paragraphs containing a single `NO-BREAK SPACE`. In HTML5, a `NO-BREAK SPACE` does not count as significant text, since it belong in the Unicode category `Zs`. I believe this kind of conformance definition and workaround arms race is not productive. However, since the requirement was in the specification draft, a checker was developed. This non-schema-based checker is an example of a simple checker. Having also makes it possible to test the requirement of significant inline content with existing Web pages to see how realistic the requirement is.

The implementation of the checker is very simple. There is a stack corresponding to open elements. The information of which element have already had significant inline content is maintained on the stack. Starts of elements that constitute embedded content cause the current stack nodes to be marked to have significant inline content. Character data is looped over and each character is tested against an ICU4J `UnicodeSet` that represents the characters in the above-mentioned Unicode categories. For forward-compatibility, the checking supports supplementary planes, even though as of Unicode 5.0 all the characters in the Unicode categories that do not count as significant text are all on the Basic Multilingual Plane of Unicode.

The checker for significant inline content demonstrates that non-schema-based checkers for particular requirements can be very simple and straight-forward.

## 5.6.4 Unicode Normalization Checking

In order to develop a prototype checker for a potential requirement that needs checking not only in the validation pipeline but also in the parser, checking for Unicode normalization was prototyped.

[CharmodNorm] is still in the Working Draft state. Nonetheless, a checking for compliance was prototyped as if [CharmodNorm] was a normative part of (X)HTML5. This way, the feasibility of the requirements of [CharmodNorm] could be evaluated in practice and checks with parser-level impact could be prototyped.

### 5.6.4.1 Requirements

The definition for Fully-normalized Text involves checking normalization before and after parsing. That is, the source text is required to be in Unicode Normalization Form C and after parsing the “constructs” parsed out of the source are required to be in Unicode Normalization Form C and are required not to start with a “composing character” (which is not exactly the same as a “combining character” in Unicode).

In order to integrate normalization checking of the unparsed character stream into *Ælfred2*, special-case decoding for US-ASCII, ISO-8859-1, UTF-8, UTF-16 and UTF-32 was removed and all character decoding was unified to use the `java.nio.charset` framework.

The definition involves peeking underneath the parser, which might be considered a violation of abstraction layers. However, the requirement of checking the unparsed source text does have the benefit that if the source is in Unicode Normalization for C, the document cannot be accidentally broken by editing it in a normalizing text editor.

### 5.6.4.2 Interpretation

[CharmodNorm] does not define what “constructs” are in the context of XML 1.0 or HTML5. However, XML 1.1 does define what “relevant constructs” are, so that definition might be generalizable to XML 1.0 and HTML5. Unfortunately, XML 1.1 defines relevant constructs in terms of the grammar productions of XML itself instead of the significant information items that an XML processor reports to the application. As a result, the XML 1.1 definition is not particularly useful in practice.

Since XML 1.0 and HTML5 do not have a definition for “constructs”, a definition that makes sense was devised for the purpose of prototyping. Canonical XML and the SAX2 `ContentHandler` interface were used as

indicators of what the meaningful constructs in XML are once the differences in syntactic sugar are out of the way.

This gave the following definition of constructs:

- Local names of elements
- Local names of attributes
- Attribute values
- Declared namespace prefixes
- Declared namespace URIs
- Processing instruction targets
- Processing instruction data
- Concatenations of consecutive character data between element boundaries and processing instructions ignoring comments and CDATA section boundaries.

#### 5.6.4.3 Implementation

As with the checker for significant inline content, the implementation turns out to be rather simple. The checker outsources most work to ICU4J. An ICU4J `UnicodeSet` is used for testing whether a character is a composing character. The ICU4J `Normalizer` class is used for testing Unicode normalization.

Constructs that are exposed as Java `Strings` in the SAX API are very simple to check. The first character is checked against the above-mentioned `UnicodeSet` and the whole string is passed to the `Normalizer` for normalization checking.

Character data, however, is checked in a piecewise manner. Most complexity in the checker is due to trying to avoid buffering as much as possible while still using the ICU4J API unmodified. Also, dealing with the halves of a surrogate pair falling into different UTF-16 code unit buffers causes complexity by roughly doubling the lines of code compared to an implementation that was not supplementary plane-aware.

The checker tries to check as much character data as possible by passing runs of the SAX-provided buffer to ICU4J. However, normalization-sensitive data may continue over the buffer boundary, so the checker copies potentially normalization-sensitive data near the buffer boundaries to its internal buffer which it later passes to ICU4J. The ability to test for composing characters is used for finding safe points for slicing buffers. By definition, it is always safe to slice buffers so that in piecewise normalization checking each buffer slice being tested starts with a character that is not a composing character.

Unfortunately, the column and line numbers reported on errors are very inaccurate due to buffering. Due to the design of the SAX API, accurate column and line positions are unavailable within a particular buffer of character data.

The normalization checking of the source text is performed by making the parsers (both the HTML parser and the XML parser) instantiate the checker on the parser level. As a side effect of reading from the decoded character stream, each buffer of UTF-16 code units is passed to the checker the way SAX character data would. The checker has a mode flag for this usage so that the error messages make sense.

## 5.7 The HTML Parser

The fundamental idea underlying the `text/html` support of the conformance checker is that HTML5 can be treated as an alternative infoset serialization for a subset of possible XML infosets, so an HTML5 parser can appear to the XML tooling as an XML parser parsing XHTML5. That is, the HTML5 parser needs to emit SAX parse events as if it was parsing an equivalent XHTML5 document.

This approach was inspired by John Cowan's TagSoup. TagSoup itself was deemed unsuitable for this project, however. TagSoup is designed to perform fix-up on horribly malformed markup and to never report a parse error of any kind. A parser must report low-level errors and preserve high-level errors to be suitable for conformance checking. Preserving high-level errors means that the kind of errors that the RELAX NG layer, the Schematron layer and the non-schema checkers are designed to find are not hidden. Tokenization errors are low-level errors that need to be reported on the parser level.

After inspecting the TagSoup code base, it was deemed easier and less conflicting with the goal of the TagSoup project to develop a parser specifically for HTML5 conformance checking. An experimental special-purpose parser was developed speculatively before the HTML5 parsing algorithm draft was published.

The parser consists of a tokenizer that emits SAX events and SAX filters that add `endElement` events for empty elements and perform tag inference for optional tags. The tokenizer emits `startElement` and `endElement` events for start tags and end tags, respectively. As a result, the SAX events emitted by the tokenizer may violate the SAX API contract. The empty element filter adds an `endElement` event for every `startElement` event belonging to an element that does not have an end tag in HTML. The tag



inference filter adds `startElement` and `endElement` events where HTML 4.01 Strict would allow tags to be omitted. In most cases the optional tags are end tags and start tags that imply the end tag can be enumerated in a simple lookup table. For implicit start tags, case-by-case checks are used.

Since the parser was developed before the HTML5 parsing algorithm was specified, it did not make sense to implement error recovery. Instead, the parser treats most tokenization errors as fatal. Also, by design the parser is not concerned with nesting errors that are known to be caught on the RELAX NG level even though normal HTML parser have to deal with issues like block elements as children of inline elements on the parser level.

There are some minor problems related to mapping HTML5 to XHTML5. XHTML5 does not allow the character encoding to be declared using the `meta` element. The parser simply does not report the `meta` element to the validation layer in that case in order to avoid schema differences for HTML5 and XHTML5 on that point. Likewise, the `base` element is not reported to the validation layer. XML places rather arbitrary restrictions the characters that are permitted in element and attribute names. Fortunately, these restrictions are not a problem for conforming HTML5 documents. Unfortunately, XML also restricts the characters permitted in content more strictly than HTML5 does. For example, HTML5 allows the form feed character. In the current implementation, the XML restrictions apply. This issue needs to be addressed if HTML5 continues to allow control characters that XML does not. Finally, the permitted contents of comments differ in a subtly way between XML and HTML5, but this is not a problem, because the parser does not report comments to the validation layer.

The version of the parser that was implemented sufficiently demonstrates the feasibility of the chosen approach for the purpose of this master's thesis. Modifying the parser to implement the HTML5 parsing algorithm was deemed to be out of the scope of the thesis project and is discussed as possible future work.

## 5.8 Character Model Checking

[WebForms2] requires documents to conform to [Charmod]. Therefore, checking for the conformance requirements that [Charmod] places on content was investigated.

Since the requirements of [Charmod] deal with issues related to character encoding, it turns out that the best opportunity for checking whether a document conforms to [Charmod] is in the parser. Hence, the checks were implemented in the parsers and not in the validation pipeline.

It turned out that some of the requirements were not machine-checkable and, thus, had to be omitted. For example, the requirement “Specifications, software and content **MUST NOT** require or depend on a one-to-one correspondence between characters and the sounds of a language.” is not machine-checkable, because the checking software cannot tell what the author expects of correspondence to sounds and if the expectation is relied upon.

Many requirements deal with the identification of the character encoding and the use of preferred IANA-registered names. These checks were implemented as part of the code that sets up the input stream decoding in the parsers.

It turned out that in addition to the requirements related to identifying the character encoding, there were only three other machine-checkable requirements. However, all three are problematic in terms of reporting errors.

The requirement “Publicly interchanged content **SHOULD NOT** use codepoints in the private use area.” was deemed worth a warning (only once per document), because human judgment is needed in order to decide when the private use area is used in a legitimate way—e.g. to encode scripts that have not yet been assigned official Unicode code points or that will not be included in Unicode as a matter of policy (such as Tengwar or Klingon which have been constructed to support particular works of fiction). Moreover, another requirement in [Charmod] denies denying the use of the private use area.

Two requirements were ignored as too impractical when considered in the light on real-world HTML authoring practice.

The first ignored requirement reads: “Escapes **SHOULD** only be used when the characters to be expressed are not directly representable in the format or the character encoding of the document, or when the visual representation of the character is unclear.” The second one is: “Content **SHOULD** use the hexadecimal form of character escapes rather than the decimal form when there are both.”

Using the five predefined entities in XML, using the HTML5 entities from the specification or using numeric character references is harmless when it comes to the parsed document tree. In XML decimal in hexadecimal character references work equally well. In HTML, the decimal form actually works better in very old browsers. Enforcing these requirements would mean proclaiming a prevalent authoring practice non-conforming on the grounds of the aesthetic preferences of the authors of [Charmod]. Moreover, [Charmod] does not give a solid machine-checkable definition for characters whose visual representation is unclear.

## 5.9 The Front End

The conformance checker needs a user interface. The user interface should itself be conforming. To meet this goal, the user interface is generated using a SAX pipeline. The output markup is produced by a SAX serializer that produces HTML5 markup from SAX events representing an XHTML5 document. Hence, the approach taken with output is the reverse of the approach taken with HTML5 input.

This method of producing markup was chosen, because isolating the markup text generation in a single serializer class makes it easier to get markup generation right compared to more traditional methods where snippets of markup text are generated ad hoc in multiple places in the program. In particular, using an isolated serializer eliminates a whole class of markup generation problems: forgetting to escape markup-significant characters or accidentally escaping them twice. [ProducingXML]

All parts of the conformance checker report to a SAX `ErrorHandler`. The `ErrorHandler` implementation that is used emits SAX events for the error messages formatted as XHTML list items. There is no intermediate data model. The XHTML list items are streamed out as the back end progresses on processing the input document.

Expressing an XML document by writing Java method calls is not as convenient as writing tags. To avoid hand-coding large unchanging parts of the user interface markup as Java method calls, a code generator called `SaxCompiler` was developed [SaxCompiler]. `SaxCompiler` is a SAX `ContentHandler` that produces Java source for a class that will call the `ContentHandler` methods in the same sequence as the methods of `SaxCompiler` itself were called. That is, the resulting class plays back SAX method calls recorded from a real XML parse. `SaxCompiler` supports the insertion of method call back to the application that uses the generated classes and supports recoding XML fragments.

The front end is a rather straight-forward Java servlet that handles form input, runs the back end code and produces output using a SAX pipeline as described above. Since the `ErrorHandler` implementation writes to the output SAX pipeline, the view and the controller are effectively conflated. There is no need for a separate data model, as the front end is only a thing wrapper for the back end code.



# Chapter 6

## Shortcomings

The developed service has some shortcomings. The real shortcomings pertain to error messages. The less important shortcomings pertain to the software not being quite as efficient in terms of performance as it could theoretically be.

### 6.1 Non-Ideal Error Messages

The foremost shortcoming of the conformance checker is the lack of useful detail in the error messages emitted upon violations of the RELAX NG schema. This problem was anticipated before starting the project. Also, the users of development versions have identified the messages generated by the Jing validation engine as a notable problem.

Jing has a handful of different messages for RELAX NG validation failures. They are all very terse. For example: “Element *name* not allowed in this context.”, “Attribute *name* not allowed at this point; ignored.” and “Bad value for attribute *name*.” The error messages are always correct, but they do not help the user understand why something went wrong.

Moreover, since the error messages do not show the erroneous markup, it is unnecessarily hard for the user to see where the problem is.

#### 6.1.1 Bimorphic Content Models

When the permissibility of an element depends on something more complex than the parent element, the error messages may confuse the user. For example, some elements take either inline or block children but not both. Moreover, in HTML 4.01 Transitional these elements were generally allowed to take a mix of inline and block children. Consider this fragment:

```
<div><em>foo</em><p>bar</p></div>
```

When the validation engine sees the `p` element, it has already committed to a derivation in the grammar that allows `em` as a child of `div`. That derivation is the inline branch of the bimorphic content model. Hence, the `p` element is not allowed in the derivation and the validation simply states that the element `p` is not allowed there. This may be confusing, because `div` does allow `p` children in other situations. A better error message would state that `div` only takes either inline or block children.

One way of addressing this problem would be allowing a mix of inline and block children in the schema and using a non-schema-based checker for detecting and reporting the mixed use of inline and block content as children of elements that have a bimorphic content model.

### 6.1.2 Lack of Datatype Diagnostics

When the value of an attribute is not permissible according to the datatype of the attribute, the validation engine emits a message simply stating that the attribute had a “bad” value. No hint is given on why the value is bad.

The RELAX NG datatype API allows a datatype implementation to communicate diagnostic messages to the validation engine in exception messages. However, Jing does not expose these messages to the user. In theory, an ambiguous grammar could cause the validation engine to test a single attribute value against multiple datatypes, so there would not be a single diagnostic message. However, in practice with many schemata it would be quite helpful to provide the diagnostic message from at least one (usually the only) datatype that did not accept the value.

Even if diagnostic information was reported, there would be a further complication. In many cases the datatype is defined using a regular expression. If the regular expression does not match, there is no useful natural-language explanation available.

## 6.2 Poor Localizability

Although the conformance checker is carefully internationalized in the sense that it correctly handles input documents in any language and supports supplementary characters in addition to the Basic Multilingual Plane of Unicode, the messages the conformance checker reports are available in English only and there is no mechanism in place for supporting other languages.

There are a number of problems related to the translatability of the messages emitted by the software. Various libraries are used, so the ways in

which messages originate is not unified (except that all libraries emit English-language messages by default). Some libraries have no localization facilities whatsoever. Other libraries do have localization mechanisms but the mechanisms assume that the user interface language is a property of the entire Java process and discovered from the environment of the process. However, this assumption does not hold for Web applications. The locale of the server is uninteresting and instead Web applications should be able to vary the user interface language on a per-HTTP response basis. In addition to these problems, the Schematron schema contains messages that are exposed to the end user.

Instead of modifying the libraries themselves, an alternative approach to localization would be reverse templating. The English messages would be matched against known patterns that would allow the variable parts to be extracted. The variable parts could then be plugged into a translated message corresponding to the matched pattern.

In order to focus on HTML5 conformance checking instead of solving the translatability problems discussed above, translatability of the user interface was left out of the scope of this project. Moreover, other validation services whose developers have borne the burden of making the software localizable have not been met with notable enthusiasm to actually produce translations.

## 6.3 Opportunities for Optimization

Some shortcomings relate to the implementation not being as efficient in terms of performance as theoretically possible. These shortcomings are not necessarily practical problems and the achievable improvements may not be worth the effort that would be required.

### 6.3.1 RELAX NG

The Jing RELAX NG engine took its current form 2003. Back then, it was designed to be compatible with Java 1.1. Dropping support for Java 1.1 opens up opportunities for performance optimization by replacing thread-safe classes with classes that perform the same tasks but do not use thread synchronization features.

When Java 1.0 and 1.1 were designed, all the classes in the standard library were made thread-safe as a matter of policy. In retrospect, this has turned out to be a bad policy. Often, a given object is only accessed from one thread, which makes synchronized monitor entry and exit useless. When an object is shared between the threads, it is likely that more than one standard class

library object need to be mutated in an atomic operation, which means the application means to manage the synchronization anyway. Even with modern virtual machine designs that can make monitors biased towards one thread so that access from that thread does not actually cause real inter-thread synchronization, class implementations that do not use useless synchronization perform better in scenarios where thread-safety is not necessary.

For compatibility with Java 1.1, Jing uses the `Hashtable` class instead of the `HashMap` class introduced in Java 1.2. Profiling the conformance checker with the NetBeans profiler shows that `Hashtable$Entry` is the second most common object (after `char[]`) in the memory allocation statistics in terms of number of live objects of a given type. While this statistics does not indicate how often the methods of `Hashtable` are called, it is still reasonable to expect the `Hashtable` class to be used a lot. Therefore, replacing occurrences of the `Hashtable` class with the API-compatible non-synchronized `HashMap` class would likely make RELAX NG validation is slightly faster.

The two other typical candidates for replacement with non-synchronized counterparts are the `Vector` class and the `StringBuffer` class. They could be replaced with `ArrayList` (introduced in Java 1.2) and `StringBuilder` (introduced in Java 5), respectively. Even though the instances of these classes are not as common as instances of `Hashtable` and `Hashtable$Entry`, it would make sense to use the non-synchronized counterparts in these cases as well.

### 6.3.2 Schematron

The Schematron implementation in Jing is based on XSLT. This is natural, considering that Schematron has been designed to be easily implementable as an XSLT transformation on the document being validated. Inside the transformer, the tree is built from the SAX parse events. The Schematron assertions fire when the entire document has been reported to the XSLT transformer.

Also, Jing creates a short-lived helper thread that sleeps when the main thread runs for fitting together API calls whose blocking behavior makes it impossible to use them from one thread. The helper thread pretends to call in to an XML parser and then blocks unblocking the main thread and allowing the parse events from the main thread to be reported as if coming from the XML parser the that the helper thread pretended to call.

The implementation approach has many points where it could be improved.



First, the helper thread can be eliminated by using an API that exposes the XSLT engine as a SAX ContentHandler instead of insisting on the the XSLT engine initiating the parse. This approach was prototyped and, indeed, it was possible to eliminate the helper thread and the overhead associated with creating and destroying it. cursory testing locally without network suggested that this improved the throughput (number of requests per unit of time) of the system by about 1%, but due to the variation between test runs the figure should be considered inaccurate.

Second, even though Schematron is designed to be implementable in XSLT, running an XSLT transformation on a full XSLT implementation is more complex than what would be minimally required to implement Schematron. The XSLT transformation spends time creating a report document which is then converted to calls to the SAX ErrorHandler. A Schematron implementation without XSLT could run XPath on a tree model and produce error messages as necessary without creating report document.

Third, Schematron implemented by running XPath expressions on a full document tree causes all the messages to appear after the entire document has been parsed and the tree built. In some cases, messages could logically be triggered much sooner. For example, in the case of exclusions as soon as an element is seen with a forbidden parent on the stack of open elements, an error message could be produced. In theory, a streaming XPath matcher could both produce errors so little in the parse and consume less memory. However, implementing such a streaming matcher for this project in particular would have been (and still would be) too great a distraction from the main goals of the project. If a streaming Schematron implementation was available off-the-shelf, using it would be worthwhile.

Fourth, the project ended up using Schematron only for two simple things: exclusions and referential integrity checking. For just these two purposes, Schematron in general and XSLT-based Schematron in particular is unnecessarily heavy. An extremely simple hand-crafted non-schema checker could replace the Schematron part of the system. A rough estimate based on the bench marking the throughput of the system with and without Schematron part suggests that the throughput of the system could increase by about 5% if the Schematron part was replaced with a hand-crafted non-schema checker. Moreover, such a checker would make it extremely easy to emit errors related to exclusions as soon as logically possible.

In summary, the best way to optimize the performance of the Schematron part would be to treat it as a rapid prototype and replace it with hand-crafted code once the HTML5 language requirements have stabilized and there is less need for the conformance checker to be easily modifiable.



# Chapter 7

## Applicability in Other Contexts

### 7.1 RELAX NG-Guided Autocompletion in Editors

Early on in the project it was assumed that the RELAX NG schema would be directly usable in RELAX NG-aware XML editors for guiding the auto-completion. RELAX NG-aware editors include nxml-mode for Emacs, oXygen XML and Etna. Unfortunately, during the course of the project it became apparent that in some cases RELAX NG could in theory express a constraint, but expressing it in Schematron or in custom code would be easier and would provide better error messages. Exclusions of the foremost class of constraints for which this is the case. Moreover, the use of a custom datatype library makes the schema less portable.

If RELAX NG could combine pattern with intersection and negation operators, writing a schema for auto-completion would be easier. However, it would not solve the problem that it would be hard for a generic RELAX NG validator to generate better error messages for exclusions than what hand-crafted messages in Schematron can provide trivially.

### 7.2 Content Management Systems

Many content management systems in use today do not properly check the input they accept. This leads to a situation commonly referred to as garbage in garbage out. The content management systems serve erroneous markup if erroneous markup has been entered into the system. The back end of the conformance checker implemented in this project could be used in content management systems to check input.

Content management systems written in Java could easily integrate the back end of the conformance checker. However, to support other programming languages the conformance checker would need to expose a remote interface that could be used from other programming languages with minimal client code. In practice, it would make sense to implement such an interface as a Web service following the REST architectural style.

# Chapter 8

## Future Work

Developing a full HTML5 conformance checker is too broad a task for a master's thesis project. For this reason, the software was developed to a point where the feasibility of implementation is demonstrated in all areas but every area was not pushed to completion.

### 8.1 Open Up

Even though the software developed in this project is Free Software / Open Source, it has not been developed in a way that would make it easily approachable to potential contributors. Perhaps the most pressing need for a change in order to move the software forward after the completion of this thesis is moving the software to a public version control system and making building and deploying the software easy.

### 8.2 The HTML5 Parsing Algorithm

The HTML parser was implemented speculatively before the HTML5 parsing algorithm had been defined. The parser needs to be revised to implement the HTML5 parsing algorithm.

The parser is designed for the SAX API, which is a streaming API. The HTML5 parsing algorithm has error recovery features that requires the parser to mutate parts of the parse tree that have already been built in earlier stages of the parse. This is incompatible with SAX. In order not to compromise streamability, the revised parser would have to treat errors that require SAX-incompatible recovery as fatal errors. This is allowed by the specification.

To make the parser reusable as a general-purpose HTML5 for other Java programs, it would be desirable to also implement the full HTML5 parsing algorithm including the SAX-incompatible parts. This would require a tree building layer as an alternative to the tag inference filter that does not build a tree. There could be an interface for pluggable tree builders. Tree builders could be provided for DOM, XOM and a special-purpose tree designed for efficient SAX event replay. A tree designed for SAX replay could store attributes as objects implementing the SAX `Attributes` interface and could store character data in `char` arrays as opposed to `Strings`.

A preliminary review of the HTML5 parsing algorithm indicates that tokenizer would not need to be completely rewritten even though the tokenizer maintains its state implicitly in the runtime stack and the HTML5 parsing algorithm maintains state explicitly and on the first sight appears to allow state transitions that do not appear to correspond to normal returns on the runtime stack. However, on a closer inspection the abnormal state transitions are always abrupt returns to the main loop and could be implemented as an exception caught in the main loop.

### 8.3 Tracking the Specification

Since the conformance checker and HTML5 itself have been developed in parallel, the conformance checker has been almost constantly more or less out of sync with the specification. Ideally, future development should track the specification on a near-daily basis instead of major synchronization work every few months.

Moreover, during the work on this project, various small issues with no clear answer in the specification were discovered. Once the issues are clarified in the specification, the software needs to be updated accordingly.

### 8.4 RELAX NG Message Improvements

The foremost problem with the RELAX NG-based approach to HTML5 conformance checking is that the error messages that a generic validator can realistically generate cannot be as good as messages written by a human for specific situations. Moreover, the actual messages emitted by Jing are not as good as they could be.

The first potential improvement is exposing the diagnostic messages from datatype libraries to the end user. Another relatively simple improvement would be mentioning the parent element of whatever element is

deemed to be forbidden in a given context. With RELAX NG, the parent is not sufficient for explaining the situation, but with HTML5 it usually is. No assessment was made to determine the feasibility of these changes, but the changes do seem simple.

If the validation engine reported the misplaced element and its parent, an HTML5-aware error message decorator could add natural-language descriptions about the content model of the parent and the allowed contexts of the child.

A more involved change would involve giving a hint of what kind of elements would have been allowed in place of the element that was not allowed. Even if querying the schema in this way was possible to add to Jing, there would be additional complication that RELAX NG-based expectations could still be wrong according to e.g. exclusions expressed in Schematron. For this reason, it may be better to focus on the first two improvement ideas.

## 8.5 More Non-Schema-Based Checkers

The table integrity checker was the most complex non-schema-based checker that is needed. Therefore, there is not work of comparable complexity left to be implemented. However, in terms of quantity, there are many simple and small requirements scattered around [WebApps] and [WebForms2] that need to be addressed using non-schema-based checkers or alternatively, in some cases, Schematron. Making sure that all these requirements have been identified, writing test cases for the requirements and implementing checking for each requirement may well end up being more time consuming than the development of the table integrity checker.

Identifying all these requirements requires particular attention, as the requirements are scattered around the specifications and many times only get a brief mention in passing. The following were identified as unimplemented features requiring non-RELAX NG checking:

- Emit a warning if there is no selected radio button in a radio button group.
- Emit an error if there are more than one radio button selected in a radio button group.
- Emit an error if a form field name starts with `Ecom_` and the name is not listed in [RFC3106].
- Emit an error if a form attribute is non-empty but does not point to a form element by ID.

- Emit an error if a `label` element has more than one form control descendant or if a `label` element has the `for` attribute and a form control as a descendant.
- Emit an error if the form submission method is `get`, the form is designated to submit to an `http` URI and the `accept-charset` attribute designates an encoding other than `US-ASCII` or `UTF-8`.
- Emit a warning if the `accept-charset` attribute is used on a form element that uses the XML submission format.
- Emit an error if there are non-unique term definitions using the `dfn` element.
- Emit an error if there is an `abbr` element that has neither a `title` attribute nor a corresponding `dfn` element.
- Emit an error if the attribute values on the `meter` element do not satisfy the expected inequalities.
- Emit an error if the value of the `value` attribute on the `progress` element is greater than the value of the `max` attribute.
- Emit an error if IDs are not unique.
- Support checking for `rel` and `class` attribute values in a way that updates the registered permissible values from an external service dynamically.

Of the features listed above, the last one will likely take the most effort to implement.

A relatively complex component similar to a non-schema-based checker is needed for showing the document outline. However, strictly speaking, such a component would not be checking machine-checkable conformance requirements and, therefore, is not included in the list above.

## 8.6 Assistance for Checking Human-Checkable Requirements

The conformance checker is unable to check for conformance requirements that require human judgment. However, the conformance checker could make it easier for human users to check such requirements.

For example, a machine cannot check if the document outline produced by the HTML5 outline algorithm makes sense. However, a machine could generate the outline and show it to the human user for review.



## 8.7 Web Service

## 8.8 Embedded MathML and SVG

Browsers from the vendors involved with the WHATWG

## 8.9 Showing the Erroneous Source Markup

The error messages give the line and column for errors, but the output does not show the actual erroneous part of the source markup. This makes it harder to see where the problems are.

Future development could include a class for collecting the character-decoded and line-identified source code. This would require changes to both the HTML and XML parsers. The parsers would need to report the unparsed source to a data holder class at the point where the byte stream has already been decoded to UTF-16 and the line boundaries have been identified.

The data holder class could be used for extracting markup snippets corresponding to each message. Also, the data could be used for formatting the entire source markup in the conformance checker output in a way that allowed linking to the points that involve errors.



# Chapter 9

## Conclusions

In this thesis the implementation of an HTML5 conformance checker was examined. The conformance checker was built around a RELAX NG validator.

The prior expectation related to the expressiveness of RELAX NG was that RELAX NG alone would not be sufficient but would have to be augmented with Schematron and hand-crafted custom code. Also, it was expected that the convenience of using RELAX NG would have the cost of making error messages worse than what they could be if they were hand-crafted a case-by-case basis.

The prior expectation related to the `text/html` serialization of HTML5 was that it could be treated as an alternative infoset serialization for a subset of possible XML infosets and, therefore, XML tools would be applicable if the parser for `text/html` exposed the result of the parse in a way an XML parser would expose the result of parsing an equivalent XML document.

It was found, through implementation experience, that these prior expectations were correct.

It was expected that it would make sense to express as RELAX NG virtually all the HTML5 conformance requirements that would be possible to express as RELAX NG. This expectation turned out to be incorrect. It turned out that especially in cases of exclusions implementation using Schematron (or custom Java code) was by far easier than expressing the exclusions in the RELAX NG grammar. Moreover, expressing exclusions in Schematron (or custom Java code) also made it possible to give more sensible error messages than what a RELAX NG validator would give if the exclusions were woven into the RELAX NG grammar.

It would be fair to say that for purposes of validation, when Schematron (or hand-crafted code) can be used alongside RELAX NG, RELAX NG should be used for expressing the general rules and Schematron should be used for expression the exceptions to the general rules. Trying to factor the

exceptions to the general rules into the RELAX NG schema is bad both for schema maintainability and for the error messages that are generated. That is, the full formal power of RELAX NG cannot be put to use, because the generation of useful error messages cannot keep up.

This has reusability implications for the RELAX NG schema. Since the Schematron schema and the custom Java code do not just improve on things that RELAX NG cannot express, using the RELAX NG schema alone in, for example, an editing system that only supports RELAX NG would mean losing some exclusion features that theoretically could be expressed in RELAX NG. Therefore, in such cases it may be worth while to use code generation together with the schema from this project to produce a schema that incorporates the exclusions in RELAX NG.

Moreover, it was unexpected that the Schematron part turned out to be constrained to exclusions and referential integrity checking. Indeed, thinking of the Schematron part as a rapid prototype of a certain kind of a non-schema checker came relatively late in the project.

Even though the applicability of schema languages was overestimated before the project, the overall hybrid implementation approach worked out fine.

# References

## [Cascading]

*Cascading Style Sheets*. Håkon Wium Lie. PhD thesis, University of Oslo, 2005.

<http://people.opera.com/howcome/2006/phd/> (referenced: 2007-02-26)

## [cdf-ws-minutes2]

*W3C Workshop on Web Applications and Compound Documents (Day 2) Jun 2, 2004*. Leigh Klotz, editor. W3C, 2004.

<http://www.w3.org/2004/04/webapps-cdf-ws/minutes-20040602.html>  
(referenced: 2006-10-18)

## [Charmod]

*Character Model for the World Wide Web 1.0: Fundamentals*. Martin J. Dürst, François Yergeau, Richard Ishida, Misha Wolf and Tex Texin, editors. W3C, 2005.

<http://www.w3.org/TR/2005/REC-charmod-20050215/>

## [CharmodNorm]

*Character Model for the World Wide Web 1.0: Normalization*, working draft. François Yergeau, Martin J. Dürst, Richard Ishida, Addison Phillips, Misha Wolf and Tex Texin, editors. W3C, 2005.

<http://www.w3.org/TR/2005/WD-charmod-norm-20051027/>

## [Compact]

*RELAX NG Compact Syntax*. James Clark, editor. OASIS, 2002.

<http://relaxng.org/compact-20021121.html>

## [CompactXSD]

*A Compact Syntax for XML Schema*. Kilian Stillhard. Master's thesis, Swiss Federal Institute of Technology Zurich, 2003.

<http://dret.net/netdret/docs/da-ws2002-stillhard.pdf>

## [DatatypeAPI]

*RELAX NG Datatype Interface*. James Clark and Kohsuke Kawaguchi.

<http://relaxng.sourceforge.net/datatype/java/apiDocs/>

## [Derivative]

*An algorithm for RELAX NG validation.* James Clark. 2002.  
<http://www.thaiopensource.com/relaxng/derivative.html>  
(referenced: 2007-02-27)

## [DocBook]

*The DocBook Schema*, working draft. Norman Walsh, editor. OASIS, 2006.  
<http://docbook.org/specs/docbook-5.0b8-spec-wd-01.html>

## [DTDDCompat]

*RELAX NG DTD Compatibility.* James Clark and Makoto Murata, editors. OASIS, 2001.  
<http://relaxng.org/compatibility-20011203.html>

## [EarlyHistory]

*The Early History of HTML.* Sean B. Palmer.  
<http://infomesh.net/html/history/early/> (referenced: 2006-10-05)

## [Frames]

*Why Frames Suck (Most of the Time).* Jakob Nielsen. 1996.  
<http://www.useit.com/alertbox/9612.html> (referenced: 2006-10-09)

## [Freddy]

*Can Blind Freddy see a pattern here?.* Rick Jelliffe. 2005.  
<http://lists.xml.org/archives/xml-dev/200507/msg00057.html>

## [Handbook]

*The SGML Handbook.* Charles F. Goldfarb. Oxford University Press, 1991.  
ISBN: 0-19-853737-9.

## [Harmful]

*Sending XHTML as text/html Considered Harmful.* Ian Hickson.  
<http://www.hixie.ch/advocacy/xhtml> (referenced: 2006-10-14)

## [HTML30]

*HyperText Markup Language Specification Version 3.0*, working draft. Dave Raggett, editor. IETF, 1995.  
<http://www.w3.org/MarkUp/html3/html3.txt>

## [HTML32]

*HTML 3.2 Reference Specification.* Dave Raggett. W3C, 1997.  
<http://www.w3.org/TR/REC-html32>

## [HTML40]

*HTML 4.0 Specification.* Dave Raggett, Arnaud Le Hors and Ian Jacobs, editors. W3C, 1997.  
<http://www.w3.org/TR/REC-html40-971218/>

## [HTML401]

*HTML 4.01 Specification*. Dave Raggett, Arnaud Le Hors and Ian Jacobs, editors. W3C, 1999.

<http://www.w3.org/TR/1999/REC-html401-19991224/>

## [HTML40rev]

*HTML 4.0 Specification*. Dave Raggett, Arnaud Le Hors and Ian Jacobs, editors. W3C, 1998.

<http://www.w3.org/TR/1998/REC-html40-19980424/>

## [HTMLplus]

*HTML+ (Hypertext markup format)*, working draft. Dave Raggett. 1993.

[http://www.w3.org/MarkUp/HTMLPlus/htmlplus\\_1.html](http://www.w3.org/MarkUp/HTMLPlus/htmlplus_1.html)

## [IIIR-HTML]

*Hypertext Markup Language (HTML)*, working draft. Tim Berners-Lee and Daniel Connolly, editors. IETF, 1993.

<http://www.w3.org/MarkUp/draft-ietf-iiir-html-01.txt>

## [IntroXML]

*An Introduction to XML and Web Technologies*. Anders Møller and Michael I. Schwartzbach. Addison-Wesley, 2006. ISBN: 0321269667.

<http://www.brics.dk/ixwt/>

## [ISO15445]

ISO/IEC 15445:2000(E), *Information technology – Document description and processing languages – HyperText Markup Language (HTML)*. ISO, 2000.

<http://purl.org/NET/ISO+IEC.15445/15445.html>

## [ISO15445TC1]

ISO/IEC 15445:2000(E) TC1, *Information technology – Document description and processing languages – HyperText Markup Language (HTML), Technical Corrigendum 1*. ISO, 2002.

<http://purl.org/NET/ISO+IEC.15445/TC1.html>

## [ISO19757-2]

ISO/IEC 19757-2:2003(E), *Information technology – Document Schema Definition Language (DSDL) – Part 2: Regular-grammar-based validation – RELAX NG*. ISO, 2003.

[http://standards.iso.org/ittf/PubliclyAvailableStandards/c037605\\_ISO\\_IEC\\_19757-2\\_2003\(E\).zip](http://standards.iso.org/ittf/PubliclyAvailableStandards/c037605_ISO_IEC_19757-2_2003(E).zip)

## [ISO19757-2Amd1]

ISO/IEC 19757-2:2003 / Amd 1:2006(E), *Information technology – Document Schema Definition Language (DSDL) – Part 2: Regular-grammar-based validation – RELAX NG – Amendment 1: Compact Syntax*. ISO, 2006.

[http://standards.iso.org/ittf/PubliclyAvailableStandards/c040774\\_ISO\\_IEC\\_19757-2\\_2003\\_Amd\\_1\\_2006\(E\).zip](http://standards.iso.org/ittf/PubliclyAvailableStandards/c040774_ISO_IEC_19757-2_2003_Amd_1_2006(E).zip)

## [ISO19757-3]

ISO/IEC 19757-3:2006(E), *Information technology – Document Schema Definition Language (DSDL) – Part 3: Rule-based validation – Schematron*. ISO, 2006.

[http://standards.iso.org/ittf/PubliclyAvailableStandards/c040833\\_ISO\\_IEC\\_19757-3\\_2006\(E\).zip](http://standards.iso.org/ittf/PubliclyAvailableStandards/c040833_ISO_IEC_19757-3_2006(E).zip)

## [ISO8879]

ISO 8879:1986, *Information processing – Text and office systems – Standard Generalized Markup Language (SGML)*. ISO, 1986.

## [JointPosition]

*Position Paper for the W3C Workshop on Web Applications and Compound Documents*. The Mozilla Foundation and Opera Software, 2004.

<http://www.w3.org/2004/04/webapps-cdf-ws/papers/opera.html>  
(referenced: 2006-10-16)

## [M12N]

*Modularization of XHTML™*. Murray Altheim, Frank Boumphrey, Sam Dooley, Shane McCarron, Sebastian Schnitzenbaumer and Ted Wugofski, editors. W3C, 2001.

<http://www.w3.org/TR/2001/REC-xhtml-modularization-20010410/>

## [M12N-RNG]

*Modularization of XHTML in RELAX NG*. James Clark. Thai Open Source Software Center Ltd, 2003.

<http://www.thaiopensource.com/relaxng/xhtml/>

## [M12N11]

*XHTML™ Modularization 1.1*. Murray Altheim, Frank Boumphrey, Sam Dooley, Shane McCarron, Sebastian Schnitzenbaumer, Ted Wugofski, Daniel Austin, Subramanian Peruvemba and Masayasu Ishikawa, editors. W3C, 2006.

<http://www.w3.org/TR/2006/PR-xhtml-modularization-20060213/>

## [mod\_validator]

*mod\_validator*. Web Thing.

[http://apache.webthing.com/mod\\_validator/](http://apache.webthing.com/mod_validator/) (referenced: 2007-01-30)



## [MozFAQ]

*Mozilla Web Author FAQ*. Henri Sivonen. 2005.  
<http://www.mozilla.org/docs/web-developer/faq.html>  
(referenced: 2006-10-18)

## [MS-WebApps]

*Paper for participation in the W3C Workshop on Web Applications and Compound Documents*. Alex Hopmann and Michael Wallent. Microsoft, 2004.  
<http://www.w3.org/2004/04/webapps-cdf-ws/papers/microsoft.html>  
(referenced: 2006-10-16)

## [ProducingXML]

*HOWTO Avoid Being Called a Bozo When Producing XML*. Henri Sivonen. 2005.  
<http://hsivonen.iki.fi/producing-xml/> (referenced: 2007-02-22)

## [Raggett]

*Raggett on HTML 4*. Dave Raggett, Jenny Lam, Ian Alexander and Michael Kmieć. Addison Wesley Longman, 1998. ISBN: 0-201-17805-2.  
<http://www.w3.org/People/Raggett/book4/ch02.html>

## [Relaxed]

*Relaxed: on the way towards true validation of compound documents*. Jirka Kosek and Petr Nálezka. In *WWW '06: Proceedings of the 15th international conference on World Wide Web*, pages 427–436. ACM Press, 2006. ISBN: 1-59593-323-9.  
<http://doi.acm.org/10.1145/1135777.1135841>

## [Relaxtron]

*Combining RELAX NG and Schematron*. Eddie Robertsson. O'Reilly Media, Inc., 2004.  
<http://www.xml.com/pub/a/2004/02/11/relaxtron.html>

## [RFC1738]

RFC 1738, *Uniform Resource Locators (URL)*. Tim Berners-Lee, Larry Masinter and Mark McCahill, editors. IETF, 1994.  
<http://ietf.org/rfc/rfc1738>

## [RFC1866]

RFC 1866, *Hypertext Markup Language - 2.0*. Tim Berners-Lee and Dan Connolly, editors. IETF, 1995.  
<http://ietf.org/rfc/rfc1866>

## [RFC1942]

RFC 1942, *HTML Tables*. Dave Raggett. IETF, 1996.  
<http://ietf.org/rfc/rfc1942>

## [RFC2070]

RFC 2070, *Internationalization of the Hypertext Markup Language*. François Yergeau, Gavin Thomas Nicol, Glenn Adams and Martin J. Duerst. IETF, 1997.

<http://ietf.org/rfc/rfc2070>

## [RFC2368]

RFC 2368, *The mailto URL scheme*. Paul E. Hoffman, Larry Masinter and Jamie Zawinski. IETF, 1998.

<http://ietf.org/rfc/rfc2368>

## [RFC2397]

RFC 2397, *The "data" URL scheme*. Larry Masinter. IETF, 1998.

<http://ietf.org/rfc/rfc2397>

## [RFC2616]

RFC 2616, *Hypertext Transfer Protocol – HTTP/1.1*. Roy T. Fielding, James Gettys, Jeffrey C. Mogul, Henrik Frystyk Nielsen, Larry Masinter, Paul J. Leach and Tim Berners-Lee. IETF, 1999.

<http://ietf.org/rfc/rfc2616>

## [RFC2818]

RFC 2818, *HTTP Over TLS*. Eric Rescorla. IETF, 2000.

<http://ietf.org/rfc/rfc2818>

## [RFC3066]

RFC 3066, *Tags for the Identification of Languages*. Harald Tveit Alvestrand. IETF, 2001.

<http://ietf.org/rfc/rfc3066>

## [RFC3106]

RFC 3106, *ECML v1.1: Field Specifications for E-Commerce*. Donald E. Eastlake and Ted Goldstein. IETF, 2001.

<http://www.ietf.org/rfc/rfc3106>

## [RFC4287]

RFC 4287, *The Atom Syndication Format*. Mark Nottingham and Robert Sayre, editors. IETF, 2005.

<http://ietf.org/rfc/rfc4287>

## [RFC4646]

RFC 4646, *Tags for Identifying Languages*. Addison Phillips and Mark Davis, editors. IETF, 2006.

<http://ietf.org/rfc/rfc4646>

## [RFC4647]

RFC 4647, *Matching of Language Tags*. Addison Phillips and Mark Davis, editors. IETF, 2006.

<http://ietf.org/rfc/rfc4647>

## [RNG-XSD]

*Guidelines for using W3C XML Schema Datatypes with RELAX NG*. James Clark and Kohsuke Kawaguchi, editors. OASIS, 2001.

<http://relaxng.org/xsd-20010907.html>

## [RNGdesign]

*The Design of RELAX NG*. James Clark.

<http://www.thaiopensource.com/relaxng/design.html>

## [Ruby]

*Ruby Annotation*. Marcin Sawicki, Michel Suignard, Masayasu Ishikawa, Martin Dürst and Tex Texin, editors. W3C, 2001.

<http://www.w3.org/TR/2001/REC-ruby-20010531/>

## [SaxCompiler]

*SaxCompiler*. Henri Sivonen. 2005.

<http://hsivonen.iki.fi/saxcompiler/> (referenced: 2007-02-22)

## [Schematron15]

*The Schematron Assertion Language 1.5*. Rick Jelliffe. Academia Sinica Computing Centre, 2002.

<http://xml.ascc.net/resource/schematron/Schematron2000.html>

## [SchematronOld]

*The Schematron – An XML Structure Validation Language using Patterns in Trees*. Rick Jelliffe. Academia Sinica Computing Centre, 2001.

<http://xml.ascc.net/resource/schematron/old-index.html>

(referenced: 2006-09-25)

## [SchemaUE]

*W3C Workshop on XML Schema 1.0 User Experiences and Interoperability*. W3C, 2005.

<http://www.w3.org/2005/03/xml-schema-user-program.html>

## [Several]

*Re: [whatwg] several messages about HTML5*. Ian Hickson. 2007.

<http://lists.whatwg.org/pipermail/whatwg-whatwg.org/>

2007-February/009517.html (referenced: 2007-02-27)

## [Stats]

*Web Authoring Statistics*. Google, 2005.

<http://code.google.com/webstats/index.html> (referenced: 2007-02-26)

## [Taxonomy]

*Taxonomy of XML schema languages using formal language theory*. Makoto Murata, Dongwon Lee, Murali Mani and Kohsuke Kawaguchi. In *ACM Trans. Inter. Tech.*, volume 5, number 4, pages 660–704. ACM Press, 2005. ISSN: 1533-5399.

<http://doi.acm.org/10.1145/1111627.1111631>

## [ToBeDeleted]

*<draft-ietf-iiir-html-01.txt, .ps> to be deleted..* W. Eliot Kimber. 1994.

<http://1997.webhistory.org/www.lists/www-talk.1994q1/0573.html>

## [Understanding]

*Understanding HTML, XML and XHTML*. Maciej Stachowiak. 2006.

<http://webkit.org/blog/?p=68> (referenced: 2006-10-14)

## [ValetMode]

*Parse Modes - Page Valet Help*. Web Thing.

<http://valet.webthing.com/page/parsemode.html>

(referenced: 2007-01-30)

## [Validace]

*Doplňková validace HTML a XHTML dokumentů*. Petr Nálevka. Bachelor's thesis, University of Economics, Prague, 2005.

[http://relaxed.sourceforge.net/thesis\\_cz.html](http://relaxed.sourceforge.net/thesis_cz.html)

## [ValidationProbs]

*An Experimental Study on Validation Problems with Existing HTML Webpages*. Shan Chen, Dan Hong and Vincent Y. Shen. In *Proc. 2005 International Conference on Internet Computing (ICOMP'05)*, pages 373–379. , 2005.

<http://www.cs.ust.hk/faculty/shen/100.pdf>

## [ValidatorAbout]

*About the Validation Service*. Henri Sivonen. 2007.

<http://hsivonen.iki.fi/validator-about/> (referenced: 2007-02-27)

## [WaterlooGML]

*Waterloo SCRIPT GML User's Guide*. University of Waterloo, 1988.

<http://www.uga.edu/~ucns/std docs/script-gmlref-tso.txt>

## [WCAG]

*Web Content Accessibility Guidelines 1.0*. Wendy Chisholm, Gregg Vanderheiden and Ian Jacobs, editors. W3C, 1999.

<http://www.w3.org/TR/1999/WAI-WEBCONTENT-19990505/>

## [WDG]

*WDG HTML Validator*. Liam Quinn.

<http://www.htmlhelp.com/tools/validator/> (referenced: 2007-01-25)

## [WDG1998]

*What Makes the WDG HTML Validator Special*. Liam Quinn. 1998.

<http://web.archive.org/web/19990128203022/htmlhelp.com/tools/validator/differences.html>

## [WDG2007]

*How the WDG HTML Validator differs from others*. Liam Quinn.

<http://www.htmlhelp.com/tools/validator/differences.html.en>  
(referenced: 2007-01-25)

## [WebApps]

*Web Applications 1.0*, working draft. Ian Hickson, editor. WHATWG, 2006.

<http://whatwg.org/specs/web-apps/current-work/>  
(referenced: 2006-10-19)

## [WebForms2]

*Web Forms 2.0*, working draft. Ian Hickson, editor. WHATWG, 2006.

<http://whatwg.org/specs/web-forms/current-work/>  
(referenced: 2006-10-19)

## [WHAT-Ann]

*WHAT open mailing list announcement*. Ian Hickson. WHATWG, 2004.

<http://whatwg.org/news/start> (referenced: 2006-10-18)

## [WHAT-Charter]

*Web Hypertext Application Technology Working Group Charter*. WHATWG.

<http://whatwg.org/charter> (referenced: 2006-10-19)

## [Whattf]

*RELAX NG Schema for (X)HTML 5*. Erika Etemad and Henri Sivonen. 2007.

<http://syntax.whattf.org/> (referenced: 2007-02-27)

## [Wilson]

*Jon Udell: Chris Wilson on IE7, Ajax, and web standards.* Jon Udell and Chris Wilson. Microsoft, 2007.

<http://channel9.msdn.com/podcasts/>

MSConversations\_wilson\_ch9.mp3 (referenced: 2007-02-26)

## [XHTML-MP]

*XHTML Mobile Profile.* WAP Forum, 2001.

<http://www.openmobilealliance.org/tech/affiliates/wap/wap-277-xhtmlmp-20011029-a.pdf>

## [XHTML10]

*XHTML™ 1.0: The Extensible HyperText Markup Language.* Steven Pemberton et al. W3C, 2000.

<http://www.w3.org/TR/2000/REC-xhtml1-20000126/>

## [XHTML10XSD]

*XHTML™ 1.0 in XML Schema.* Masayasu Ishikawa, editor. W3C, 2002.

<http://www.w3.org/TR/2002/NOTE-xhtml1-schema-20020902/>

## [XHTML11]

*XHTML™ 1.1 – Module-based XHTML.* Murray Altheim and Shane McCarron, editors. W3C, 2001.

<http://www.w3.org/TR/2001/REC-xhtml11-20010531/>

## [XHTMLBasic]

*XHTML™ Basic.* Mark Baker, Masayasu Ishikawa, Shinichi Matsui, Peter Stark, Ted Wugofski and Toshihiko Yamakami, editors. W3C, 2000.

<http://www.w3.org/TR/2000/REC-xhtml-basic-20001219/>

## [XML]

*Extensible Markup Language (XML) 1.0.* Tim Bray, Jean Paoli and C. M. Sperberg-McQueen, editors. W3C, 1998.

<http://www.w3.org/TR/1998/REC-xml-19980210>

## [xmlid]

*xml:id Version 1.0.* Jonathan Marsh, Daniel Veillard and Norman Walsh, editors. W3C, 2005.

<http://www.w3.org/TR/2005/REC-xml-id-20050909/>

## [XPath]

*XML Path Language (XPath).* James Clark and Steve DeRose, editors. W3C, 1999.

<http://www.w3.org/TR/xpath>

## [XSDDatatypes]

*XML Schema Part 2: Datatypes Second Edition*. Paul V. Biron and Ashok Malhotra, editors. W3C, 2004.

<http://www.w3.org/TR/2004/REC-xmlschema-2-20041028/>

## [XSDDatatypes11WD]

*XML Schema 1.1 Part 2: Datatypes*, working draft. David Peterson, Paul V. Biron, Ashok Malhotra and C. M. Sperberg-McQueen, editors. W3C, 2006.

<http://www.w3.org/TR/2006/WD-xmlschema11-2-20060217/>

## [XSDDatatypesFE]

*XML Schema Part 2: Datatypes*. Paul V. Biron and Ashok Malhotra, editors. W3C, 2001.

<http://www.w3.org/TR/2001/REC-xmlschema-2-20010502/>