

# Validating RDF Data

# Synthesis Lectures on Semantic Web: Theory and Technology

## Editors

**Ying Ding, Indiana University**

**Paul Groth, Elsevier Labs**

## Founding Editor Emeritus

**James Hendler, Rensselaer Polytechnic Institute**

*Synthesis Lectures on the Semantic Web: Theory and Technology* is edited by Ying Ding of Indiana University and Paul Groth of Elsevier Labs. Whether you call it the Semantic Web, Linked Data, or Web 3.0, a new generation of Web technologies is offering major advances in the evolution of the World Wide Web. As the first generation of this technology transitions out of the laboratory, new research is exploring how the growing Web of Data will change our world. While topics such as ontology-building and logics remain vital, new areas such as the use of semantics in Web search, the linking and use of open data on the Web, and future applications that will be supported by these technologies are becoming important research areas in their own right. Whether they be scientists, engineers or practitioners, Web users increasingly need to understand not just the new technologies of the Semantic Web, but to understand the principles by which those technologies work, and the best practices for assembling systems that integrate the different languages, resources, and functionalities that will be important in keeping the Web the rapidly expanding, and constantly changing, information space that has changed our lives.

Topics to be included:

- Semantic Web Principles from linked-data to ontology design
- Key Semantic Web technologies and algorithms
- Semantic Search and language technologies
- The Emerging “Web of Data” and its use in industry, government and university applications
- Trust, Social networking and collaboration technologies for the Semantic Web
- The economics of Semantic Web application adoption and use
- Publishing and Science on the Semantic Web
- Semantic Web in health care and life sciences

[\*\*Validating RDF Data\*\*](#)

Jose Emilio Labra Gayo, Eric Prud'hommeaux, Iovka Boneva, and Dimitris Kontokostas  
2017

[\*\*Natural Language Processing for the Semantic Web\*\*](#)

Diana Maynard, Kalina Bontcheva, and Isabelle Augenstein  
2016

[\*\*The Epistemology of Intelligent Semantic Web Systems\*\*](#)

Mathieu d'Aquin and Enrico Motta  
2016

[\*\*Entity Resolution in the Web of Data\*\*](#)

Vassilis Christophides, Vasilis Efthymiou, and Kostas Stefanidis  
2015

[\*\*Library Linked Data in the Cloud: OCLC's Experiments with New Models of Resource Description\*\*](#)

Carol Jean Godby, Shenghui Wang, and Jeffrey K. Mixter  
2015

[\*\*Semantic Mining of Social Networks\*\*](#)

Jie Tang and Juanzi Li  
2015

[\*\*Social Semantic Web Mining\*\*](#)

Tope Omitola, Sebastián A. Ríos, and John G. Breslin  
2015

[\*\*Semantic Breakthrough in Drug Discovery\*\*](#)

Bin Chen, Huijun Wang, Ying Ding, and David Wild  
2014

[\*\*Semantics in Mobile Sensing\*\*](#)

Zhixian Yan and Dipanjan Chakraborty  
2014

[\*\*Provenance: An Introduction to PROV\*\*](#)

Luc Moreau and Paul Groth  
2013

[\*\*Resource-Oriented Architecture Patterns for Webs of Data\*\*](#)

Brian Sletten  
2013

[Aaron Swartz's A Programmable Web: An Unfinished Work](#)  
Aaron Swartz  
2013

[Incentive-Centric Semantic Web Application Engineering](#)  
Elena Simperl, Roberta Cuel, and Martin Stein  
2013

[Publishing and Using Cultural Heritage Linked Data on the Semantic Web](#)  
Eero Hyvönen  
2012

[VIVO: A Semantic Approach to Scholarly Networking and Discovery](#)  
Katy Börner, Michael Conlon, Jon Corson-Rikert, and Ying Ding  
2012

[Linked Data: Evolving the Web into a Global Data Space](#)  
Tom Heath and Christian Bizer  
2011

© Springer Nature Switzerland AG 2022  
Reprint of original edition © Morgan & Claypool 2018

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopy, recording, or any other except for brief quotations in printed reviews, without the prior permission of the publisher.

Validating RDF Data

Jose Emilio Labra Gayo, Eric Prud'hommeaux, Iovka Boneva, and Dimitris Kontokostas

ISBN: 978-3-031-79477-3 paperback  
ISBN: 978-3-031-79478-0 ebook  
ISBN: 978-3-031-79478-0 e-pub

DOI 10.1007/978-3-031-79478-0

A Publication in the Springer series

*SYNTHESIS LECTURES ON SEMANTIC WEB: THEORY AND TECHNOLOGY*

Lecture #16

Series Editors: Ying Ding, *Indiana University*

Paul Groth, *Elsevier Labs*

Founding Editor Emeritus: James Hendler, *Rensselaer Polytechnic Institute*

Series ISSN

Print 2160-4711 Electronic 2160-472X

Cover art: Amy van der Hiel

# Validating RDF Data

Jose Emilio Labra Gayo

University of Oviedo

Eric Prud'hommeaux

W3C/MIT and Micelio

Iovka Boneva

University of Lille

Dimitris Kontokostas

University of Leipzig

*SYNTHESIS LECTURES ON SEMANTIC WEB: THEORY AND  
TECHNOLOGY #16*

## ABSTRACT

RDF and Linked Data have broad applicability across many fields, from aircraft manufacturing to zoology. Requirements for detecting bad data differ across communities, fields, and tasks, but nearly all involve some form of data validation. This book introduces data validation and describes its practical use in day-to-day data exchange.

The Semantic Web offers a bold, new take on how to organize, distribute, index, and share data. Using Web addresses (URIs) as identifiers for data elements enables the construction of distributed databases on a global scale. Like the Web, the Semantic Web is heralded as an information revolution, and also like the Web, it is encumbered by data quality issues. The quality of Semantic Web data is compromised by the lack of resources for data curation, for maintenance, and for developing globally applicable data models.

At the enterprise scale, these problems have conventional solutions. Master data management provides an enterprise-wide vocabulary, while constraint languages capture and enforce data structures. Filling a need long recognized by Semantic Web users, shapes languages provide models and vocabularies for expressing such structural constraints.

This book describes two technologies for RDF validation: Shape Expressions (ShEx) and Shapes Constraint Language (SHACL), the rationales for their designs, a comparison of the two, and some example applications.

## KEYWORDS

RDF, ShEx, SHACL, shape expressions, shapes constraint language, data quality, web of data, Semantic Web, linked data

# Contents

Preface .....	xv
Foreword by Phil Archer .....	xvii
Foreword by Tom Baker .....	xix
Foreword by Dan Brickley and Libby Miller .....	xi
Acknowledgments .....	xxiii
<b>1 Introduction .....</b>	<b>1</b>
1.1 RDF and the Web of Data .....	1
1.2 RDF: The Good Parts .....	1
1.3 Challenges for RDF Adoption .....	3
1.4 Structure of the Book .....	5
1.5 Conventions and Notation .....	6
<b>2 The RDF Ecosystem .....</b>	<b>9</b>
2.1 RDF History .....	9
2.2 RDF Data Model .....	10
2.3 Shared Entities and Vocabularies .....	17
2.4 Technologies Related with RDF .....	18
2.4.1 SPARQL .....	18
2.4.2 Inference Systems: RDF Schema and OWL .....	20
2.4.3 Linked Data, JSON-LD, Microdata, and RDFa .....	23
2.5 Summary .....	25
2.6 Suggested Reading .....	25
<b>3 Data Quality .....</b>	<b>27</b>
3.1 Non-RDF Schema Languages .....	28
3.1.1 UML .....	28
3.1.2 SQL and Relational Databases .....	29

3.1.3	XML .....	31
3.1.4	JSON .....	37
3.1.5	CSV .....	39
3.2	Understanding the RDF Validation Problem .....	40
3.3	Previous RDF Validation Approaches .....	45
3.3.1	Query-based Validation .....	45
3.3.2	Inference-based Approaches .....	47
3.3.3	Structural Languages .....	48
3.4	Validation Requirements .....	48
3.4.1	General Requirements .....	49
3.4.2	Graph-based Requirements .....	49
3.4.3	RDF Data Model Requirements .....	50
3.4.4	Data-modeling-based Requirements .....	50
3.4.5	Expressiveness of Schema Language .....	51
3.4.6	Validation Invocation Requirements .....	52
3.4.7	Usability Requirements .....	52
3.5	Summary .....	52
3.6	Suggested Reading .....	53
<b>4</b>	<b>Shape Expressions .....</b>	<b>55</b>
4.1	Use of ShEx .....	55
4.2	First Example .....	56
4.3	ShEx implementations .....	58
4.4	The Shape Expressions Language .....	59
4.4.1	Shape Expressions Compact Syntax .....	59
4.4.2	Invoking Validation .....	60
4.4.3	Structure of Shape Expressions .....	63
4.4.4	Start Shape Expression .....	64
4.5	Node Constraints .....	65
4.5.1	Node kinds .....	67
4.5.2	Datatypes .....	68
4.5.3	Facets on Literals .....	70
4.5.4	Value Sets .....	73
4.6	Shapes .....	78
4.6.1	Triple Constraints .....	79
4.6.2	Groupings .....	80
4.6.3	Cardinalities .....	80

4.6.4	Choices . . . . .	82
4.6.5	Nested Shapes . . . . .	84
4.6.6	Inverse Triple Constraints . . . . .	85
4.6.7	Repeated Properties . . . . .	86
4.6.8	Permitting other Triples . . . . .	87
4.7	References . . . . .	90
4.7.1	Shape References . . . . .	90
4.7.2	Recursion and Cyclic References . . . . .	91
4.7.3	External Shapes . . . . .	92
4.7.4	Labeled Triple Expression . . . . .	93
4.7.5	Annotations . . . . .	94
4.8	Logical Operators . . . . .	95
4.8.1	Conjunction . . . . .	95
4.8.2	Disjunction . . . . .	98
4.8.3	Negation . . . . .	101
4.9	Shape Maps . . . . .	105
4.9.1	Fixed Shape Maps . . . . .	105
4.9.2	Query Shape Maps . . . . .	106
4.9.3	Result Shape Maps . . . . .	108
4.9.4	JSON Representation . . . . .	109
4.9.5	Chaining Validation Workflows . . . . .	110
4.10	Semantic Actions . . . . .	110
4.11	ShEx and Inference . . . . .	111
4.12	Importing schemas . . . . .	113
4.13	RDF and JSON-LD Syntax . . . . .	114
4.14	Summary . . . . .	116
4.15	Suggested Reading . . . . .	116
<b>5</b>	<b>SHACL . . . . .</b>	<b>119</b>
5.1	Simple Example . . . . .	119
5.2	SHACL Implementations . . . . .	122
5.3	Basic Definitions: Shapes Graphs, Node, and Property Shapes . . . . .	124
5.4	Importing other Shapes Graphs . . . . .	125
5.5	Validation Report . . . . .	126
5.6	Shapes . . . . .	129
5.6.1	Node shapes . . . . .	129

5.6.2	Property Shapes . . . . .	130
5.6.3	Constraint Components . . . . .	131
5.6.4	Human Friendly Messages . . . . .	133
5.6.5	Declaring Shape Severities . . . . .	134
5.6.6	Deactivating Shapes . . . . .	135
5.7	Target Declarations . . . . .	137
5.7.1	Target Node . . . . .	137
5.7.2	Target Class . . . . .	138
5.7.3	Implicit Class Target . . . . .	139
5.7.4	Target Subjects Of . . . . .	140
5.7.5	Target Objects Of . . . . .	141
5.8	Cardinality . . . . .	141
5.9	Constraints on Values . . . . .	142
5.9.1	Datatypes . . . . .	142
5.9.2	Class of Values . . . . .	145
5.9.3	Node Kinds . . . . .	146
5.9.4	Sets of Values . . . . .	147
5.9.5	Specific Value . . . . .	148
5.10	Datatype Facets . . . . .	148
5.10.1	Value Ranges . . . . .	149
5.10.2	String-based Constraints . . . . .	149
5.10.3	Language-based Constraints . . . . .	151
5.11	Logical Constraints: and, or, not, xone . . . . .	154
5.11.1	AND . . . . .	154
5.11.2	OR . . . . .	157
5.11.3	Exactly One . . . . .	159
5.11.4	Not . . . . .	162
5.11.5	Combining Logical Operators . . . . .	162
5.12	Shape-based Constraints . . . . .	164
5.12.1	Shape References and Recursion . . . . .	166
5.12.2	Qualified Value Shapes . . . . .	174
5.13	Closed Shapes . . . . .	177
5.14	Property Pair Constraints . . . . .	180
5.15	Non-validating SHACL Properties . . . . .	182
5.16	SHACL-SPARQL . . . . .	184
5.16.1	SPARQL Constraints . . . . .	184

5.16.2 SPARQL-based Constraint Components .....	185
5.17 SHACL and Inference Systems .....	188
5.18 SHACL Compact Syntax .....	190
5.19 SHACL Rules and Advanced Features .....	190
5.20 SHACL Javascript .....	193
5.21 Summary .....	194
5.22 Suggested Reading .....	194
<b>6 Applications .....</b>	<b>195</b>
6.1 Describing a Linked Data Portal .....	195
6.1.1 WebIndex in ShEx .....	197
6.1.2 WebIndex in SHACL .....	200
6.2 Describing Clinical Records—FHIR .....	204
6.2.1 FHIR as Linked Data .....	206
6.2.2 Consistency constraints .....	206
6.2.3 FHIR/RDF Development .....	209
6.2.4 Generic Properties .....	210
6.3 Springer Nature SciGraph .....	212
6.4 DBpedia Validation Use Cases .....	213
6.4.1 Ontology-based Validation .....	213
6.4.2 RDF Mappings Validation .....	214
6.4.3 Validating Link Contributions with SHACL .....	215
6.4.4 Ontology Validation with SHACL .....	216
6.5 ShEx for ShEx .....	219
6.6 SHACL in SHACL .....	225
6.7 Summary .....	230
6.8 Suggested Reading .....	231
<b>7 Comparing ShEx and SHACL .....</b>	<b>233</b>
7.1 Common Features .....	233
7.2 Syntactic Differences .....	237
7.3 Foundation: Schema vs. Constraints .....	239
7.4 Invoking Validation .....	240
7.5 Modularization and Reusability .....	242
7.6 Shapes, Classes, and Inference .....	244
7.7 Violation Reporting and Severities .....	246

7.8	Default Cardinalities .....	246
7.9	Property Paths .....	247
7.10	Recursion .....	248
7.11	Property Pair Constraints and Uniqueness .....	250
7.12	Repeated Properties .....	251
7.13	Exactly One and Alternatives .....	254
7.14	Treatment of Closed Shapes .....	257
7.15	Stems and Stem Ranges .....	259
7.16	Annotations .....	260
7.17	Semantics and Complexity .....	261
7.18	Extension Mechanisms .....	262
7.19	Conclusions and Outlook .....	263
7.20	Summary .....	266
7.21	Suggested Reading .....	266
<b>A</b>	<b>WebIndex in ShEx .....</b>	<b>267</b>
<b>B</b>	<b>WebIndex in SHACL .....</b>	<b>269</b>
<b>C</b>	<b>ShEx in ShEx .....</b>	<b>275</b>
<b>D</b>	<b>SHACL in SHACL .....</b>	<b>279</b>
	<b>Bibliography .....</b>	<b>285</b>
	<b>Authors' Biographies .....</b>	<b>295</b>

# Preface

This book describes two languages for implementing constraints on RDF data, describing the main features of both Shape Expressions (ShEx) and Shapes Constraint Language (SHACL) from a user perspective, and also offering a comparison of the technologies. Throughout this book, we develop a small number of examples that typify validation requirements and demonstrate how they can be met with ShEx and SHACL. The book is not intended to be a formal specification of the languages, for which the interested reader can consult the corresponding reference documents, but rather, it is meant to serve as an introduction to the technologies with some background about the rationale of their design and some points of comparison.

Chapter 1 provides a brief introduction to the topic. Chapter 2 presents a short overview of the RDF data model and RDF-related technologies; this chapter could be skipped by any reader who already knows RDF or Turtle. Chapter 3 helps the reader to understand what to expect from data validation. It describes the problem of RDF validation and some approaches that have been proposed. This book specifically reviews two of these approaches in further detail: ShEx (Chapter 4) and SHACL (Chapter 5). These chapters describe each language and provide a practical introduction using examples. Following the discussion of both languages, Chapter 6 presents some applications using either ShEx, SHACL, or both. Finally, Chapter 7 compares ShEx and SHACL and offers some conclusions.

The goal of this book is to serve as a practical introduction to ShEx and SHACL using examples. While we omitted formal definitions or specifications, references for further reading can be found at the end of each chapter. We give a quick overview of some background and related technologies so that readers without RDF knowledge can follow the book's contents. Also, it is not necessary to have any prior knowledge of programming or ontologies to understand RDF validation technologies. The intended audience is anyone interested in data representation and quality.

Jose Emilio Labra Gayo, Eric Prud'hommeaux, Iovka Boneva, and Dimitris Kontokostas  
July 2017

# Foreword by Phil Archer

“Anyone can say anything about anything,” says the mantra for the Semantic Web. More formally, the Semantic Web adopts the Open World Assumption: just because your data encodes a set of facts, that doesn’t mean there aren’t other facts stated elsewhere about the same thing. All of which is fine and part of the design of RDF which supports the creation of a graph at Web scale, but in a lot of practical applications you just need to know whether the triples you’ve ingested match what you were expecting; you need validation. You might think of it as a defined subset of the whole graph, or maybe a profile, providing a huge boost to interoperability between disparate systems. If you can validate the data you’ve received then you can process it with confidence, using more terse code, perhaps with more performant queries. I don’t accept that RDF is hard, certainly no harder than any other Web technology; what is hard is thinking in graphs. Keeping in your head that this node supports these properties and has relationships with those other nodes becomes complex for anything other than trivial datasets. The validation techniques set out in this book provide a means to tame that complexity, to set out for humans and machines exactly what the structure of the data is or should be. That’s got to be helpful and, incidentally, ties in with new work now under way at W3C on dataset exchange. In my role at W3C I watched as the SHACL and ShEx camps tried hard to converge on a single method: they couldn’t, hence the two different approaches. Both are described in detail here with copious examples, which is just what you need to get started. How can you choose between the two methods? Chapter 7 gives a detailed comparison and allows you to make your own choice. Whichever you choose, this is the book you need to make sense of RDF validation.

Phil Archer, Former W3C Data Strategist  
July 2017

# Foreword by Tom Baker

The technologies described here meet a need first recognized, albeit dimly, two decades ago. Rewind to circa 2000, when the parallel development of two W3C specifications, XML Schema and RDF Schema, both called “schema languages” but with radically different uses, caused some confusion.

This confusion permeated our early discussions about the Dublin Core. Was it an XML format, an RDF vocabulary, or somehow both? Could metadata just follow arbitrary formats or did it need a data model? In 2000, the Dublin Core community turned to “application profiles” as a way to mix and match multiple vocabularies to meet specific needs, and the idea was an instant hit even if people disagreed about their use. Were they more for validating data, or more about finding rough consensus on a metadata model within a community of practice? Attempts to bridge the XML and RDF mindsets in the DCMI community, notably with a Description Set Profile constraint language for validating RDF-based metadata (2008), never quite caught on. Perhaps the idea needed a bigger push?

Fast-forward to 2013, when W3C convened a workshop on RDF validation which revealed that many communities had been circling around the same issues, and which ultimately led to the results described here [82]. This book focuses on data validation, an addition to the Semantic Web stack that is long overdue. But from a DCMI perspective, the ideas for future work outlined in its Conclusion are just as exciting: the prospect of using ShEx- or SHACL-based application profiles to map and convert between data models, size up aggregated datasets, or provide nuanced feedback to data providers on quality issues. ShEx and SHACL, finally production-ready, are full of potential for further evolution.

Tom Baker, Dublin Core Metadata Initiative  
July 2017

# Foreword by Dan Brickley and Libby Miller

People think RDF is a pain because it is complicated. The truth is even worse. RDF is painfully simplistic, but it allows you to work with real-world data and problems that are horribly complicated. While you can avoid RDF, it is harder to avoid complicated data and complicated computer problems. RDF brings together data across application boundaries and imposes no discipline on mandatory or expected structures. This can make working with RDF data frustrating. Its schema and ontology languages can help define the meaning of RDF content but, again, can't express rules about how actual data records should look. The contents of this book are nearly 20 years too late, but better now than never. Recent developments around RDF validation have finally made it easier to record, exchange, and understand rules about validating and otherwise checking RDF data. Who knows what wonders await us in another 20 years.

Dan Brickley, Schema.org and Google

Libby Miller, BBC

July 2017

# Acknowledgments

This book started as an RDF validation tutorial that we prepared for the International Semantic Web Conference (ISWC) 2016. The slides are still available at [http://weso.github.io/RDF\\_Validation\\_ESWC16/](http://weso.github.io/RDF_Validation_ESWC16/). After the tutorial, the editors invited us to extend the slides to write a full book on the subject.

We want to thank Harold Solbrig for his collaboration with the tutorial and his support and participation in the development of ShEx.

Although the task of extending the tutorial seemed easy at first, writing a book on this subject was much more daunting than expected. One challenge was to update the material to the changes that ShEx and SHACL were experimenting with during the W3C process.

There were different points of view on how to tackle the RDF validation problem that led to the appearance of two camps, which at the time of this writing are represented in two W3C groups: the Shape Expressions community group and the SHACL community group. We would like to acknowledge the people participating in those groups as well as the people that were part of the W3C Data Shapes Working group.

During the development of the book we were also updating our implementations. We also want to thank all the people that have been using our implementations and online demos for their patience when something was not working as expected and for submitting issues or suggesting improvements to the tools.

We want to thank Amy van der Hiel for the design of the cover art, as well as Dan Brickley and Tom Baker for their comments and suggestions.

We are also grateful to Vladimir Alexiev for his extensive review and proofreading of the book.

Jose Emilio Labra Gayo wants to thank his colleagues from Oviedo3 and the WESO research group. Special thanks to Edita for her support and patience. He dedicates the book to Sergio, Nuria, and Alex.

Dimitris wants to thank Sebastian Hellmann and the University of Leipzig for supporting his participation in the Shape Expressions and Data Shapes groups as well as his family for their continuous support.

## **xxiv ACKNOWLEDGMENTS**

Eric wants to dedicate this book to his family, who put up with him writing when he should be playing. He also has great gratitude to Gregg Kellogg, who always takes time to hack with him, and Tom Baker, who always takes time to write with him.

Jose Emilio Labra Gayo, Eric Prud'hommeaux, Iovka Boneva, and Dimitris Kontokostas  
July 2017