

A peer-reviewed version of this preprint was published in PeerJ on 30 January 2017.

[View the peer-reviewed version](https://peerj.com/articles/cs-105) (peerj.com/articles/cs-105), which is the preferred citable publication unless you specifically need to cite this preprint.

Dimou A, Vahdati S, Di Iorio A, Lange C, Verborgh R, Mannens E. 2017. Challenges as enablers for high quality Linked Data: insights from the Semantic Publishing Challenge. PeerJ Computer Science 3:e105 <https://doi.org/10.7717/peerj-cs.105>

Challenges as Enablers for High Quality Linked Data: Insights from the Semantic Publishing Challenge

Anastasia Dimou ^{Corresp., 1,2}, **Sahar Vahdati** ³, **Angelo Di Iorio** ⁴, **Christoph Lange** ^{3,5}, **Ruben Verborgh** ^{1,2}, **Erik Mannens** ^{1,2}

¹ Faculty of Engineering and Architecture, Ghent University, Ghent, Belgium

² imec, Leuven, Belgium

³ Department of Intelligent Systems, University of Bonn, Bonn, Germany

⁴ Department of Computer Science and Engineering, University of Bologna, Bologna, Italy

⁵ Enterprise Information Systems, Fraunhofer IAI, Sankt Augustin, Germany

Corresponding Author: Anastasia Dimou

Email address: anastasia.dimou@ugent.be

While most challenges organized so far in the Semantic Web domain are focused on comparing tools with respect to different criteria such as their features and competencies, or exploiting semantically enriched data, the Semantic Web Evaluation Challenges series, co-located with the ESWC Semantic Web Conference, aims to compare them based on their output, namely the produced dataset. The Semantic Publishing Challenge is one of these challenges. Its goal is to involve participants in extracting data from heterogeneous sources on scholarly publications, and producing Linked Data that can be exploited by the community itself. This paper reviews lessons learned from both (i) the overall organization of the Semantic Publishing Challenge, regarding the definition of the tasks, building the input dataset and forming the evaluation, and (ii) the results produced by the participants, regarding the proposed approaches, the used tools, the preferred vocabularies and the results produced in the three editions of 2014, 2015 and 2016. We compared these lessons to other Semantic Web Evaluation challenges. In this paper, we (i) distill best practices for organizing such challenges that could be applied to similar events, and (ii) report observations on Linked Data publishing derived from the submitted solutions. We conclude that higher quality may be achieved when Linked Data is produced as a result of a challenge, because the competition becomes an incentive, while solutions become better with respect to Linked Data publishing best practices when they are evaluated against the rules of the challenge.

Challenges as Enablers for High Quality Linked Data: Insights from the Semantic Publishing Challenge

Anastasia Dimou^{1,2}, Sahar Vahdati³, Angelo Di Iorio⁴, Christoph Lange^{3,5}, Ruben Verborgh^{1,2}, Erik Mannens^{1,2}

¹ Faculty of Engineering and Architecture, Ghent University, Ghent, Belgium

² imec, Leuven, Belgium

³ Department of Intelligent Systems, University of Bonn, Bonn, Germany

⁴ Department of Computer Science and Engineering, Università di Bologna, Bologna, Italy

⁵ Enterprise Information Systems, Fraunhofer IAIS, Sankt Augustin, Germany

Corresponding Author:

Anastasia Dimou¹

Sint-Pietersnieuwstraat 41, Ghent, B-9000, Belgium

Email address: Anastasia.Dimou@UGent.be

Challenges as Enablers for High Quality Linked Data: Insights from the Semantic Publishing Challenge

Anastasia Dimou^{1,2}, Sahar Vahdati³, Angelo Di Iorio⁴, Christoph Lange^{3,5}, Ruben Verborgh^{1,2}, and Erik Mannens^{1,2}

¹Faculty of Engineering and Architecture, Ghent University, Ghent, Belgium

²imec, Leuven, Belgium

³Department of Intelligent Systems, University of Bonn, Bonn, Germany

⁴Department of Computer Science and Engineering, University of Bologna, Bologna, Italy

⁵Enterprise Information Systems, Fraunhofer IAIS, Sankt Augustin, Germany

ABSTRACT

While most challenges organized so far in the Semantic Web domain are focused on comparing tools with respect to different criteria such as their features and competencies, or exploiting semantically enriched data, the Semantic Web Evaluation Challenges series, co-located with the ESWC Semantic Web Conference, aims to compare them based on their output, namely the produced *dataset*. The Semantic Publishing Challenge is one of these challenges. Its goal is to involve participants in extracting data from heterogeneous sources on scholarly publications, and producing Linked Data that can be exploited by the community itself. This paper reviews lessons learned from both (i) the overall organization of the Semantic Publishing Challenge, regarding the definition of the tasks, building the input dataset and forming the evaluation, and (ii) the results produced by the participants, regarding the proposed approaches, the used tools, the preferred vocabularies and the results produced in the three editions of 2014, 2015 and 2016. We compared these lessons to other Semantic Web Evaluation challenges. In this paper, we (i) distill best practices for organizing such challenges that could be applied to similar events, and (ii) report observations on Linked Data publishing derived from the submitted solutions. We conclude that higher quality may be achieved when Linked Data is produced as a result of a challenge, because the competition becomes an incentive, while solutions become better with respect to Linked Data publishing best practices when they are evaluated against the rules of the challenge.

Keywords: Challenge, Semantic Publishing, Linked Data Publishing

1 INTRODUCTION

The Semantic Web aims to extend the human-readable Web by encoding the semantics of resources in a machine-comprehensible and reusable fashion. Over the past years, a growing amount of research on publishing and consuming Linked Data, i.e. data represented and made available in a way that maximizes reusability, has facilitated Semantic Web adoption. However, one of the remaining issues is lack of high quality Linked Data. A promising means to foster and accelerate the publication of such high quality Linked Data is the organization of *challenges*: competitions during which participants complete tasks with innovative solutions that are then ranked in an objective way to determine the winner. A significant number of challenges has been organized so far, including the Semantic Web Challenge¹, its Big Data Track formerly known as the Billion Triples Challenge, and the LinkedUp Challenge², to mention a few of the longest lasting. However, these challenges targeted broad application domains and were more focused on innovative ways of exploiting Semantic Web enabled tools (*Linked Data consumption*) than on the output actually produced (*Linked Data production*). Therefore, such challenges enable advancement of Semantic Web technology but overlook the possibility of also advancing Linked Datasets per se.

This paper focuses on a series of Challenges in the Semantic Publishing domain. *Semantic publishing* is defined as “the enhancement of scholarly publications by the use of modern Web standards to improve

46 interactivity, openness and usability, including the use of ontologies to encode rich semantics in the form
47 of machine-readable RDF metadata” by Shotton (2009). The 2014 Semantic Publishing Challenge, was
48 themed “Assessing the Quality of Scientific Output” (Lange and Di Iorio, 2014)³, in 2015 we mentioned
49 the techniques more explicitly by appending “. . . by Information Extraction and Interlinking” (Di Iorio
50 et al., 2015)⁴, and in 2016 we generalized to “. . . in its Ecosystem” to emphasize the multiple dimensions
51 of scientific quality and the potential impact of producing Linked Data about it (Dimou et al., 2016)⁵.

52 According to Miller and Mork (2013), extracting, annotating and sharing scientific data (by which,
53 here, we mean standalone research datasets, data inside documents, as well as metadata about datasets
54 and documents) and then building new research efforts on them, can lead to a data value chain producing
55 value for the scholar and Semantic Web community. On the one hand, the scholar community benefits
56 from a challenge that produces data, as the challenge results in more data and in data of higher quality
57 being available to the community to exploit. On the other hand, the Semantic Web community benefits:
58 participants optimize their tools towards performance in this particular challenge, but such optimisations
59 may also improve the tools in general. Once such tools are reused, any other dataset benefits from their
60 advancements, because the processes producing them has been improved. However, bootstrapping and
61 enabling such value chains is not easy.

62 In a recent publication (Vahdati et al., 2016), we discussed lessons we learned from our experience
63 in organizing the first two editions of the Semantic Publishing Challenge – mainly from the perspective
64 of how to improve the organization of further editions and of providing a better service to the scholar
65 community. The lessons are related to the challenge organization, namely defining the tasks, building
66 the input datasets and performing the evaluation, as well as lessons we learned by studying the solutions,
67 with respect to the methodologies, tools and ontologies used, and data produced by the participants. We
68 organized the third edition based on these lessons learned.

69 In this paper, we revise our lessons learned, taking into consideration experience gained by organizing
70 the challenge’s third edition, whose results validate in principle our lessons learned. We argue that
71 challenges may act as enablers for the generation of higher quality Linked Data, because of the competitive
72 aspect. However, organizing a successful challenge is not an easy task. Therefore, the goal of this paper is
73 to distill *generic best practices*, which could be applied to similar events, rendering the challenge tasks into
74 meaningful milestones for efficient Linked Data generation and publishing. To achieve that, we validated
75 the generalizability of our lessons learned against the other Semantic Web Evaluation Challenges^{6,7,8}.

76 We concluded that our lessons learned are applicable to other challenges too; thus they can be
77 considered *best practices* for organizing a challenge. Other challenge organizers may benefit from relying
78 on these best practices when organizing their own challenge. Additionally, we thoroughly analyze and
79 report best practices followed by the Linked Data that the *solutions* to our challenge’s tasks produce. Our
80 study of the different solutions provides insights regarding different approaches that address the same task,
81 namely it acts as if the challenge benchmarks those different solutions against a common problem. Last,
82 we assess based on the produced datasets how the challenge organization reinforces increasing Linked
83 Data quality in respect to the different Linked Data dimensions identified by Zaveri et al. (2016).

84 Thus, besides the scholarly community and the CEUR-WS.org open access repository, which is the
85 owner of the underlying data, the broader Linked Data community may benefit from looking into our
86 cumulative results. Other Linked Data owners may find details on different approaches dealing with the
87 same problem and the corresponding results they produce. Taking them into consideration, they can
88 determine their own approach for an equivalent case or even consider launching a corresponding challenge
89 to determine the best performing tool with respect to the desired results and consider this one for their
90 regular long term use. Moreover, other Linked Data publishers may advise the results or consider the best
91 practices as their guidelines for improving their tools and thus their results.

92 In summary, our contributions are:

- 93 • an *outline of challenges* organized in the field of Linked Data and Semantic Web technologies,
- 94 • an *exhaustive analysis* of all solutions to every task of all editions of the Semantic Publishing
95 Challenge series,
- 96 • a *systematic discussion of lessons* that we have learned from organizing the Semantic Publishing
97 Challenge, and

- 98 • a structured set of *best practices* for organizing similar challenges, resulting from validating our
99 lessons against other Semantic Web Evaluation challenges.

100 The remainder of the paper is structured as follows: Section 2 reviews related work; in particular it
101 sets the background for our study by recapitulating the Semantic Publishing Challenges run so far and
102 comparing them to related challenges. Section 3 revisits the lessons learned, taking into consideration all
103 three editions, validates them against other challenges and concludes in best practices for organizing such
104 challenges. Section 4 exhaustively and cumulatively analyses the solutions submitted to all tasks of all
105 challenges in the series. Section 5 reviews the Semantic Publishing Challenges as a means of assessing
106 the quality of data, and Section 6 summarizes our conclusions.

107 2 BACKGROUND AND RELATED WORK

108 This section sets the background of the Semantic Publishing Challenges so far. Section 2.1 summarizes
109 other challenges, mainly those run in the Semantic Web community. Then, Section 2.2 recapitulates the
110 Semantic Publishing Challenges run so far, including the definitions of their tasks, and their outcomes.

111 2.1 State of the Art on Previously Organized Challenges

112 Several related challenges were organized in the past for different purposes and application domains.
113 In this section, we summarize the most well-known, long-lasting and closely related challenges in the
114 Semantic Web field. Where applicable, we report on systematic *reviews* of challenges for lessons learned.

115 2.1.1 Ontology Matching Challenges

116 The *Ontology Matching Challenges*⁹ have been organized since 2004 by the Ontology Alignment Eval-
117 uation Initiative (OAEI)¹⁰ and co-located with several top Information Systems and Web conferences
118 such as WWW¹¹ or VLDB¹². It aims to forge a consensus for evaluating the different emerging methods
119 for schema or ontology matching. The OAEI aims to assess the strengths and weaknesses of alignment/
120 matching systems, compare the performance of techniques, and improve evaluation techniques to help
121 improving the work on ontology alignment/matching through evaluating the techniques' performances.
122 Following a similar structure as the Semantic Publishing Challenge, the OAEI challenge provides a
123 list of test ontologies as training datasets. The SEALS infrastructure¹³ to evaluate the results has been
124 made available since 2011. The results are presented during the Ontology Matching workshop, which is
125 usually co-located with the International Semantic Web Conference (ISWC¹⁴). The tests and results of
126 the challenge are published for further analysis.

127 2.1.2 Semantic Web Challenge

128 The *Semantic Web Challenge*¹⁵ aims to apply Semantic Web techniques in building online end-user
129 applications that integrate, combine and deduce information needed to assist users in performing tasks. It
130 features a track about Big Data designed to demonstrate approaches which can work on Web scale using
131 realistic Web-quality data. The Big Data Track, formerly known as the *Billion Triples Challenge (BTC)*,
132 started from 2008 mostly co-located with ISWC. The Billion Triples Challenge aimed to demonstrate
133 the capability of Semantic Web technologies to process very large and messy data as typically found on
134 the Web. The track was renamed to "Big Data Track" because very large data sets are now ubiquitous
135 and the competition was opened to broader range of researchers dealing with their own big data. The
136 functionality of submitted solutions is open but, to address real scalability issues, it forces all participants
137 to use a specific Billion Triple Challenge Dataset provided by the challenge's organizers.

138 2.1.3 Question Answering over Linked Data (QALD)

139 The *Question Answering over Linked Data (QALD)* challenge¹⁶ (Lopez et al., 2013; Unger et al., 2015)
140 focuses on answering natural language or keyword-based questions over linked datasets. Co-located with
141 the ESWC Semantic Web Conference (ESWC¹⁷) in its first two editions in 2011 and 2013, it moved to
142 the Conference and Labs of the Evaluation Forum (CLEF¹⁸) for the three following editions, to return to
143 ESWC as a part of its Semantic Web Evaluation Challenges track explained below. In all editions, a set of
144 up to 340 questions over DBpedia¹⁹ served as input; participants were expected to answer these questions.
145 The 2013 to 2016 editions had a task on multilingual questions, while from 2014, a task on hybrid question
146 answering over RDF and free text was added. Some editions considered alternative datasets, e.g., about
147 drugs or music, and had alternative sub-tasks on answering questions over interlinked datasets or finding

148 lexicalizations of ontological terms. Only few submitted solutions address the question/answering issues
149 over a distributed and large collection of interconnected datasets.

150 The first two editions of the QALD Challenge were reviewed (Lopez et al., 2013); similarly to our
151 work, this review “discuss[es] how the second evaluation addressed some of the issues and limitations
152 which arose from the first one, as well as the open issues to be addressed in future competitions”. Like
153 us, Lopez et al. present the definition of the QALD challenge’s tasks and the datasets used, and draw
154 conclusions for the subsequent evaluation of question answering systems from reviewing concrete results
155 of the first two challenge editions. Their review of related work includes a review of methods for evaluating
156 question answering systems, whereas the Semantic Publishing Challenge was created to address the lack
157 of such methods for evaluating semantic publishing tools (cf. Section 2.2). We additionally present lessons
158 learned for challenge organization (Section 3) and about semantic publishing tools (Section 4), which,
159 together, constitute the main contribution of this paper.

160 **2.1.4 LAK Challenges**

161 The *Learning Analytics and Knowledge Challenges* (LAK²⁰) use a specific dataset of structured metadata
162 from research publications in the field of learning analytics. The challenge was organized in 2011 for
163 the first time and has so far continued yearly with the LAK conference. Beyond merely publishing the
164 data, the LAK challenges encourage its innovative use and exploitation. Participants submit a meaningful
165 use case of the dataset in the scope of six topic categories, such as comparison of the LAK and EDM
166 (Educational Data Mining) communities, innovative applications to explore, navigate and visualize,
167 enrichment of the Dataset, and usage of the dataset in recommender systems. Considering that a lot of
168 information is still available only in textual form, the submitted approaches can not only deal with the
169 specific character of structured data. The aim for further challenges is to combine solutions for processing
170 both structured and unstructured information from distributed datasets.

171 **2.1.5 LinkedUp**

172 The *LinkedUp* challenge was run by the LinkedUp project²¹ since 2014. The main purpose of the
173 project was to push educational organizations to make their data publicly available on the Web. One of
174 the activities towards this purpose was to organize the *LinkedUp Challenge*. The three editions of the
175 challenge focused on three different levels of maturity: demo prototypes and applications, innovative
176 tools and applications, and mature data-driven applications. Participants were asked to submit demos of
177 tools that analyze and/or integrate open Web data for educational purposes. For all the above challenges,
178 the participants were asked to submit a scientific paper along with their tool and dataset.

179 d’Aquin et al. (2014) present lessons learned from the LinkedUp project (Linking Web Data for
180 Education). However, their paper provides a summary of the outcomes of the project, including a
181 summary of the LinkedUp Challenge, rather than a systematically structured account of lessons learned.

182 **2.1.6 Dialog State Tracking Challenge (DSTC)**

183 The challenge series review that is most closely related to ours in its methodology has been carried out by
184 Williams et al. (2016) over a challenge series from a field of computer science that is related to semantics
185 but not to the Web: the Dialog State Tracking Challenge (DSTC²²) on “correctly inferring the state of [a]
186 conversation [...] given all of the dialog history”. Like our review, the one of DSTC is based on three
187 editions of a challenge, each of which built on its predecessor’s results, and it presents the definition of
188 the challenge’s tasks and the datasets used. Like we do in Section 4, they provide a structured overview of
189 the submissions to the DSTC challenges. However, the focus of their review is on the evolution of tools in
190 their domain of dialog state tracking, whereas our review additionally covers lessons learned for challenge
191 design (cf. Section 3), besides tools in the domain of Semantic publishing.

192 **2.1.7 Other related works**

193 There are further related works and challenges that we consider out of the scope, as they are not focused
194 on Linked Data sets. For example, the *AI Mashup Challenge*²³ as a part of the ESWC conference
195 focused on innovative mashups, i.e. web applications combining multiple services and datasets, that were
196 evaluated by a jury. Information Retrieval campaigns are a series of comparative evaluation methods that
197 originate from the 1960s and are used to compare various retrieval strategies or systems. As an example
198 of such campaigns SemEval (Semantic Evaluation)²⁴ is one of the ongoing series of evaluations of
199 computational semantic analysis systems with a focus on Textual Similarity and Question Answering and
200 Sentiment Analysis (Clough and Sanderson (2013)). The *Computational Linguistics Scientific Document*

201 *Summarization Shared Task (CL-SciSumm)*²⁵ is based on a corpus of annotated documents; tasks focus on
 202 correctly identifying the underlying text that a summary refers to, but also on generating summaries.

Table 1. Semantic Web Evaluation Challenges

Abbreviation	Challenge	Years
SemPub	Semantic Publishing Challenge	2014, 2015, 2016
CLSA	(Concept-Level) Sentiment Analysis Challenge	2014, 2015, 2016
RecSys	Linked Open Data-Enabled Recommender System Challenge	2014, 2015
OKE	Open Knowledge Extraction Challenge	2015, 2016
SAQ	Schema-agnostic Queries over Linked Data	2015
QALD	Open Challenge on Question Answering over Linked Data	2016
Top-K	Top-K Shortest Path in Large Typed RDF Graphs Challenge	2016

203 **2.1.8 Semantic Web Evaluation Challenges**

204 The *Semantic Web Evaluation Challenges*, including our Semantic Publishing Challenge, aim at de-
 205 veloping a set of common benchmarks and establish evaluation procedures, tasks and datasets in the
 206 Semantic Web field. They are organized as an official track of the ESWC Semantic Web Conference,
 207 which introduces common standards for its challenges, e.g., common deadlines for publishing the training
 208 and evaluation datasets. The purpose of the challenges is to showcase methods and tools on tasks common
 209 to the Semantic Web and adjacent disciplines, in a controlled setting involving rigorous evaluation. Each
 210 Semantic web Evaluation Challenge is briefly described here and all of them are summarized at Table 1.

211 **Concept-Level Sentiment Analysis Challenge** The *Concept-Level Sentiment Analysis Challenge*
 212 (*CLSA*) focuses on semantics as a key factor for detecting the sentiment of a text, rather than just performing
 213 a lexical analysis of text; cf. Reforgiato Recupero and Cambria (2014) and Reforgiato Recupero et al.
 214 (2015). Participants are asked to use Semantic Web technology to improve their sentiment analysis system
 215 and to measure the performance of the system²⁶ within the Sentiment Analysis track of the SEMEVAL
 216 2015 workshop²⁷. An automatic evaluation tool²⁸ was applied to the submissions; it was made available
 217 to the participants before their submission. In the second edition, participants were asked to submit a
 218 concept-level sentiment analysis engine that exploited linked datasets such as DBpedia.

219 **Linked Open Data-Enabled Recommender Systems Challenge** The *Linked Open Data-Enabled*
 220 *Recommender Systems Challenge* (Di Noia et al., 2014) was designed with two main goals: i) establish
 221 links between the two communities of recommender systems and Semantic Web, ii) develop content-based
 222 recommendation systems using interlinking and other semantic web and technologies. The first edition
 223 featured three independent tasks related to a book recommendation use case. While the first edition was
 224 successful, the second edition was canceled because it had no participants.

225 **Open Knowledge Extraction Challenge** The *Open Knowledge Extraction Challenge (OKE)* focuses on
 226 content extraction from textual data using Linked Data technology (Nuzzolese et al., 2015a). The challenge
 227 was divided into two sub-tasks²⁹ focusing on entity recognition and entity typing. The participants of
 228 the challenge were the developers of four different well-known system in this community. The three
 229 defined tasks were focused on a) entity recognition, linking and typing for knowledge base population,
 230 b) entity typing for vocabulary and knowledge Base enrichment and c) Web-scale knowledge extraction
 231 by exploiting structured annotation. The submissions were evaluated using two different methods: i)
 232 using datasets for training purposes and for evaluating the performance of the submitted approaches, ii)
 233 establishing an evaluation framework to measure the accuracy of the systems. The applications of task 1
 234 and 2 were published as web services with input/output provided in the NLP Interchange Format NIF³⁰.

235 **Schema-Agnostic Queries over Linked Data Challenge** The *Schema-Agnostic Queries over Linked*
 236 *Data Challenge (SAQ)* was designed to invite schema-agnostic query approaches and systems (Freitas
 237 and Unger, 2015). The goal of this challenge is to improve querying approaches over complex databases
 238 with large schemata and to relieve users from the need to understand the database schema. Tasks were
 239 defined for two types of queries: schema-agnostic SPARQL queries and schema-agnostic keyword-based
 240 queries. Participants were asked to submit the results together with their approach without changing the

241 query syntax but with different vocabularies and structural changes. A gold standard dataset was used to
242 measure precision, recall and F1-score.

243 **2.2 Semantic Publishing Challenge: 2014–2016**

244 In this section, we briefly summarize the history of the Semantic Publishing Challenge to provide the
245 necessary background for the following discussion. More detailed reports for each edition have been
246 published separately by Lange and Di Iorio (2014), Di Iorio et al. (2015), and Dimou et al. (2016).

247 We sought a way to challenge the semantic publishing community to accomplish tasks whose results
248 could be compared *in an objective way*. After some preliminary discussion, we focused on *information*
249 *extraction tasks*. The basic idea was to provide as input some scholarly papers – in multiple formats – and
250 some queries in natural language. Participants were asked to extract data from these papers and to publish
251 them as an RDF dataset that could be used to answer the input queries. The best performing approach
252 was identified automatically by comparing the output of the queries in the produced datasets against a
253 gold standard, and by measuring precision and recall. Our selection of queries was motivated by quality
254 assessment scenarios complementary to the traditional metrics based on counting citations: how can the
255 extracted information serve as indicators for the quality of scientific output such as publications or events.
256 The same motivation, structure and evaluation procedure have been maintained in the following years,
257 with some improvements and extensions.

258 All challenge's series' tasks (Section 2.2.1), the input to the tasks, namely the training and evaluation
259 datasets (Section 2.2.2), the output, namely the submitted solutions and the produced dataset (Section 2.2.3)
260 and how their evaluation was conducted (Section 2.2.4) are briefly explained below.

261 **2.2.1 Tasks Evolution**

262 Table 2 summarizes the tasks' full history. For each year and each task, we highlight the data source and
263 the format of the input files, along with a short description of the task and a summary on the participation.

264 **2014 edition tasks.** The first edition had two main tasks (Task 1 and Task 2) and an open task (Task 3;
265 see Lange and Di Iorio (2014) for full details and statistics of this challenge's edition).

266 For Task 1, the participants were asked to extract information from selected CEUR-WS.org workshop
267 proceedings volumes to enable the computation of indicators for the workshops' quality assessment.
268 The input files were HTML tables of content using different levels of semantic markup, as well as PDF
269 full text. The participants were asked to answer twenty queries. For Task 2, the input dataset included
270 XML-encoded research papers, derived from the PubMedCentral and Penseful Open Access archives.
271 The participants were asked to extract data about citations to assess the value of articles, for instance by
272 considering citations' position in the paper, their co-location with other citations, or their purpose. In total,
273 they were asked to answer ten queries. Dataset and queries were completely disjoint from Task 1.

274 After circulating the call for submissions, we received feedback from the community that mere
275 information extraction, even if motivated by a quality assessment use case, was not the most exciting task
276 related to the future of scholarly publishing, as it assumed a traditional publishing model. Therefore, to
277 address the challenge's primary target, i.e. 'publishing' rather than just 'metadata extraction', we widened
278 the scope by adding an open task (Task 3). Participants were asked to showcase data-driven applications
279 that would eventually support publishing. We received a good number of submissions; winners were
280 selected by a jury.

281 **2015 edition tasks.** In 2015 we were asked to include only tasks that could be evaluated in a fully
282 objective manner, and thus we discarded the 2014's edition open task (Task 3).

283 While Task 1 queries remained largely stable from 2014 to 2015, the queries for Task 2 changed. We
284 transformed Task 2 into a PDF mining task, instead of XML, and thus moved all PDF-related queries
285 there. The rationale was to differentiate tasks on the basis of the competencies and tools required to solve
286 them. Since the input format was completely new and we expected different teams to participate (as
287 actually happened), we wanted to explore new areas and potentially interesting information. In fact, we
288 asked participants to extract data not only on citations but also on affiliations and fundings. The number of
289 queries remained unchanged (ten in total). We also decided to use the same data source for both tasks, and
290 to make them interplay. CEUR-WS.org data has become the central focus of the whole Challenge, for two
291 reasons: on the one hand, the data provider (CEUR-WS.org) takes advantage of a broader community that
292 builds on its data, which, before the Semantic Publishing Challenges, had not been available as Linked

Table 2. Semantic Publishing Challenge Evolution from 2014 to 2016

		2014 edition	2015 edition	2016 edition
Task 1	Task	Extracting data on workshops history and participants	Extracting data on workshops history and participants	Extracting data on workshops history and participants
	Source	CEUR-WS.org proceedings volumes	CEUR-WS.org proceedings volumes	CEUR-WS.org proceedings volumes
	Format	HTML and PDF	HTML	HTML
	Solutions	3	4	0
	Awards	best performance innovation	best performance innovation	–
	Decision	–	chairs' assessment	chairs' assessment
Task 2	Task	Extracting data on citations	Extracting data on citations, affiliations, fundings	Extracting data on internal structure, affiliations, fundings
	Source	PubMed	CEUR-WS.org	CEUR-WS.org
	Format	XML	PDF	PDF
	Solutions	1	6	5
	Awards	–	best performance most innovative	best performance most innovative
	Decision	–	chairs' assessment	chairs' assessment
Task 3	Task	Open task: showcasing semantic publishing applications	Interlinking cross-dataset entities	Interlinking cross-dataset entities cross-task entities
	Source	–	CEUR-WS.org, Colinda DBLP, Springer LD Lancet, SWDF	CEUR-WS.org, Colinda DBLP, Springer LD
	Format	–	RDF	RDF
	Solutions	4	0	0
	Awards	most innovative (jury assessment)	–	–

293 Data³¹. On the other hand, data consumers gain the opportunity to assess the quality of scientific venues
294 by taking a deeper look into their history, as well as the quality of the publications.

295 In 2015, we also introduced a new Task 3. Instead of being an open task, Task 3 was focused
296 on interlinking the dataset produced by the winners of Task 1 from the 2014 edition of the Semantic
297 Publishing Challenge with related datasets in the Linked Data Cloud.

298 **2016 edition tasks.** The tasks of the 2016 edition were designed to ensure continuity and to allow
299 previous participants to use and refine their tools.

300 In particular, Task 1 was unchanged except for some minor details on queries. Task 2 was still on
301 PDF information extraction but queries were slightly changed: considering the interest and results of the
302 participants in the past, we did not include citations any more. Rather, we added some queries on the
303 identification of the structural components of the papers (table of contents, captions, figures and tables)
304 and maintained queries on funding agencies and projects. In total, we had ten queries in 2016 as well.

305 Task 3 remained the same but it was repurposed. Instead of only aiming for cross-dataset links
306 between the dataset produced by the Task 1 winners of the previous edition of the challenge and other,
307 external datasets, Task 3 now focused on interlinking the datasets produced by the winners of Task 1 and
308 Task 2 of the 2015 edition. Thus, the task aimed not only at *cross-dataset* but also at *cross-task* links: the
309 goal was to link entities identified in the CEUR-WS.org website with the same entities that were extracted
310 from the proceedings papers. Moreover, the number of external datasets was reduced.

311 **2.2.2 Input: Training and Evaluation Datasets**

312 In this section we give an overview of the datasets used for the above mentioned tasks. These datasets were
313 incrementally refined and, as discussed below in Section 3.2.1, some valuable indications can be taken
314 from their analysis. For each task, and for each year, we published two datasets: (i) a training dataset (TD)
315 on which the participants could test and train their extraction tools and (ii) an evaluation dataset (ED)
316 made available a few days before the final submission and used as input for the final evaluation.

317 **Training and Evaluation dataset for Task 1.** The CEUR-WS.org workshop proceedings volumes
318 served as the source for selecting the training and evaluation datasets of Task 1 in all challenge editions.
319 In this data source, which included data spanning over 20 years, workshop proceedings volumes were
320 represented in different formats and at different levels of encoding quality and semantics. An HTML 4
321 main index page³² links to all workshop proceedings volumes, which have HTML tables of contents and
322 contain PDF or PostScript full texts. A mixture of different HTML formats (no semantic markup at all,
323 different versions of microformats, RDFa) were chosen for both the training and evaluation datasets. The
324 training dataset comprised all volumes of several workshop series, including, e.g., the Linked Data on the
325 Web workshop at the WWW conference, and all workshops of some conferences, e.g., of several editions
326 of ESWC. In 2014 and 2015, the evaluation dataset was created by adding further workshops on top of
327 the training dataset. To support the evolution of extraction tools, the training datasets of 2015 and 2016
328 were based on the unions of the training and evaluation datasets of the previous years. In 2015 and 2016,
329 the Task 1 dataset of the previous year served as an input to Task 3.

330 **Training and Evaluation dataset for Task 2.** In 2014, the datasets for Task 2 included XML files
331 encoded in JATS³³ and TaxPub³⁴, an official extension of JATS customized for taxonomic treatments (Cat-
332 apano, 2010). The training dataset consisted of 150 files from 15 journals, while the evaluation dataset
333 included 400 papers and was a superset of the training dataset. In 2015, we switched to PDF information
334 extraction: the training dataset included 100 papers taken from some of the workshops analyzed in Task
335 1, while the evaluation dataset included 200 papers from randomly selected workshops (uniform to the
336 training dataset). In 2016, we reduced the number of papers increasing the cases for each query. Thus, we
337 included 50 PDF papers in the training and 40 in the evaluation dataset. Again, the papers were distributed
338 in the same way and used different styles for headers, acknowledgments and structural components.

339 **Training and Evaluation dataset for Task 3.** The training dataset for Task 3 consists of the CEUR-
340 WS.org dataset produced by the 2014 winning tool of Task 1³⁵, COLINDA³⁶, DBLP³⁷, Lancet³⁸, SWDF³⁹,
341 and Springer LD⁴⁰ in 2015 and the CEUR-WS.org datasets produced by the 2015 winning tools of Task
342 1⁴¹ and Task 2⁴², of COLINDA, DBLP, and Springer LD in 2016.

Table 3. Task 1 solutions: their primary analysis methods, methodologies, implementations basis and evaluation results.

	Solution 1.1	Solution 1.2	Solution 1.3	Solution 1.4
Publications	Kolchin et al. (2015) Kolchin and Kozlov (2014)	Heyvaert et al. (2015) Dimou et al. (2014)	Ronzano et al. (2015) Ronzano et al. (2014)	Milicka and Burget (2015) –
Primary analysis				
structure-based		✓	✓	
syntactic-based	✓			✓
linguistic-based			✓	
layout-based				✓
Methodology				
method	Crawling	Generic solution for abstracted mappings	Linguistic and structural analysis	Visual layout multi-aspect content analysis
case-specific	✓		✓ (partly)	✓ (partly)
template-based	✓	✓		
NLP/NER			✓	✓
Implementation				
basis	n/a	RML	GATE	FITLayout
language	Python	Java	Java	Java, HTML
rules language	XPath	RML, CSS	JAPE	HTML, CSS
code/rule separation		✓	✓	
regular expressions	✓	✓	✓	✓
external services			✓	✓
open source	✓	✓		✓
license	MIT	MIT	–	GPL-3.0
Evaluation				
precision improvement	11.1%	11.4%	10.7%	–
recall improvement	11.3%	11.3%	10.9%	–
best performing	✓ (2014)			✓ (2015)
most innovative			✓ (2014)	✓ (2015)

Table 4. Task 1 and 2 solutions: the vocabularies used to annotate the data.

	Sol 1.1	Sol 1.2	Sol 1.3	Sol 1.4	Sol 2.1	Sol 2.2	Sol 2.3	Sol 2.4	Sol 2.5	Sol 2.6	Sol 2.7	Sol 2.8
bibo ⁹³	✓	✓		✓				✓	✓			
co ⁴³			✓							✓		
DBO ^{4.2.2}	✓		✓	✓		✓			✓			
DC ¹⁰²	✓	✓	✓	✓	✓				✓			✓
DCterms ¹⁰³	✓			✓	✓			✓		✓		
event ¹⁰⁷		✓							✓			
FOAF ¹⁰⁴	✓	✓	✓	✓				✓	✓	✓		✓
schema ¹⁰⁹						✓		✓				
SKOS ⁴⁴	✓											
SPAR ⁹⁵		✓	✓				✓	✓		✓		✓
BiRO			✓							✓		
CiTO												✓
DoCO							✓	✓				✓
FaBiO		✓	✓				✓	✓		✓		
FRAPO							✓	✓				
FRBR			✓									
PRO			✓				✓	✓		✓		
SWC ⁹⁴	✓			✓				✓				
SWRC ⁹²	✓	✓	✓	✓					✓	✓		
timeline ¹⁰⁸	✓			✓								
vcard ¹⁰⁶			✓	✓	✓							
custom								✓	✓		✓	✓

Table 5. Statistics about the model (Task 1 – 2014 and 2015 editions)

year	Solution 1.1		Solution 1.2		Solution 1.3		Solution 1.4
	2014	2015	2014	2015	2014	2015	2015
Conferences	swc:OrganizedEvent	swc:OrganizedEvent	swc:Event	bibo:Conference	swrc:Event	swrc:Conference	swrc:ConferenceEvent
Workshops	bibo:Workshop	bibo:Workshop	swc:Event	bibo:Workshop	swrc:Event	swrc:Workshop	swrc:Section
Proceedings	swrc:Proceedings	bibo:Proceeding	bibo:Volume	bibo:Proceeding	swrc:Proceedings	swrc:Proceedings	swrc:Proceedings
Papers	swrc:InProceedings	swrc:InProceedings, foaf:Document	bibo:Article	swrc:InProceedings	swrc:Publication	swrc:Publication	swc:Paper
Persons	foaf:Agent	foaf:Person	foaf:Person	foaf:Person	foaf:Person	foaf:Person	foaf:Person

343 **2.2.3 Output: Solutions and Datasets produced**

344 There were four distinct solutions in total for Task 1 during the three editions of the challenge, eight
 345 distinct solutions in total for Task 2 and none for Task 3 during the last two editions. All solutions for
 346 each task are briefly summarized here.

347 **Task 1.** There were four distinct solutions proposed to address Task 1 in 2014 and 2015 editions of the
 348 challenge. Three participated in both editions, whereas the fourth solution participated only in 2015. All
 349 solutions are briefly introduced here and summarized in Table 3, Table 4, Table 5, Table 6, and Table 7.
 350 Table 3 provides details about the methodologies, approach and implementation each solution followed.
 351 Table 4 summarizes the model and vocabularies/ontologies each solution used (both for Task 1 and Task
 352 2), whereas Table 7 provides statistics regarding the dataset schema/entities and triples/size each solution
 353 produced (again both for Task 1 and Task 2). Last, Table 5 summarizes the data model each solution
 354 considered and Table 6 the number of instances extracted and annotated per concept for each solution.

Table 6. Number of entities per concept for each solution (Task 1 – 2014 and 2015 editions)

year	Solution 1.1		Solution 1.2		Solution 1.3		Solution 1.4
	2014	2015	2014	2015	2014	2015	2015
Conferences	21	46		46		5	47
Workshops	132	252	14	1,393	1,516	127	198
Proceedings	126	243	65	1,392	124	202	1,353
Papers	1,634	3,801	971	2,452	1,110	720	2,470
Persons	2,854	6,700	202	6,414	2,794	3,402	11,034

Table 7. Statistics about the produced dataset (Task 1 – 2014 and 2015 editions)

year	Solution 1.1		Solution 1.2		Solution 1.3		Solution 1.4
	2014	2015	2014	2015	2014	2015	2015
dataset size	1.5M	25M	1.7M	7.2M	2.7M	9.1M	9.7M
# triples	32,088	177,752	14,178	58,858	60,130	62,231	79,444
# entities	4,770	11,428	1,258	11,803	9,691	11,656	19,090
# properties	60	46	43	23	45	48	23
# classes	8	30	5	10	10	19	6

355 **Solution 1.1** Kolchin et al. (2015) and Kolchin and Kozlov (2014) presented a case-specific crawling
 356 based approach for addressing Task 1. It relies on an extensible template-dependent crawler that uses
 357 sets of special predefined templates based on XPath and regular expressions to extract the content from
 358 HTML and convert it in RDF. The RDF is then processed to merge resources using fuzzy-matching. The
 359 use of the crawler turns the system tolerant to invalid HTML pages. This solution improved its precision
 360 in 2015 as well the richness of the data model.

361 **Solution 1.2** Heyvaert et al. (2015) and Dimou et al. (2014) exploited a generic tool for generating RDF
 362 data from heterogeneous data. It uses the RDF Mapping Language (RML)⁴⁵ to define how data extracted
 363 from CEUR-WS.org Web pages should be semantically annotated. RML extends R2RML⁴⁶ to express
 364 mapping rules from heterogeneous data to RDF. CSS3 selectors⁴⁷ are considered to extract the data from
 365 the HTML pages. The RML mapping rules are parsed and executed by the RML Processor⁴⁸. In 2015 the
 366 solution reconsidered its data model and was extended to validate both the mapping documents and the
 367 final RDF, resulting in an overall improved quality dataset.

368 **Solution 1.3** Ronzano et al. (2015, 2014) designed a case-specific solution that relies on chunk-based
 369 and sentence-based Support Vector Machine (SVM) classifiers which are exploited to semantically
 370 characterize parts of CEUR-WS.org proceedings textual contents. Thanks to a pipeline of text analysis
 371 components based on the GATE Text Engineering Framework⁴⁹, each HTML page is characterized
 372 by structural and linguistic features: these features are then exploited to train the classifiers on the
 373 ground-truth provided by the subset of CEUR-WS.org proceedings with microformat annotations. A
 374 heuristic-based annotation sanitizer is applied to fix classifiers imperfections and interlink annotations.
 375 The produced dataset is also extended with information retrieved from external resources.

376 **Solution 1.4** Milicka and Burget (2015) presented an application of the FITLayout framework⁵⁰. This
 377 solution participated in the Semantic Publishing Challenge only in 2015. It combines different page
 378 analysis methods, i.e. layout analysis and visual and textual feature classification to analyze the rendered
 379 pages, rather than their code. The solution is quite generic but requires domain/case-specific actions in
 380 certain phases (model building step).

381 **Task 2** There were eight distinct solutions proposed to address Task 2 in the 2015 and 2016 editions
 382 of the challenge. Three participated in both editions, three only in 2015 and two only in 2016. As the
 383 definition of Task 2 changed fundamentally from 2014 to 2015, the only solution submitted for Task 2 in
 384 2014 (Bertin and Atanassova, 2014) is not comparable to the 2015 and 2016 solutions and therefore not
 385 discussed here. All solutions for Task 2 – except for the one of 2014 – are briefly introduced here and
 386 summarized in Table 4, Table 8, Table 9, Table 10 and Table 11. Table 9 and Table 10 provide details
 387 about the methodologies and approach each solution followed. Table 11 summarizes details regarding
 388 the implementation and its components each solution employed to address Task 2. Table 4 summarizes
 389 the model and vocabularies/ontologies each solution used (both for Task 1 and Task 2), whereas Table 8
 390 provides statistics regarding the dataset schema/entities and triples/size each solution produced (again
 391 both for Task 1 and Task 2).

Table 8. Statistics about the produced dataset (Task 2 – 2015 and 2016 editions)

	Sol 2.1	Sol 2.2		Sol 2.3	Sol 2.4	Sol 2.5	Sol 2.6	Sol 2.7	Sol 2.8
year	2015	2015	2016	2016	2015	2015	2015	2016	2016
dataset size	2.6M	1.5M	285	184K	3.6M	2.4M	17M	152	235
# triples	21,681	10,730	2,143	1,628	15,242	12,375	98,961	1,126	1,816
# entities	4,581	1,300	334	257	3,249	2,978	19,487	659	829
# properties	12	23	23	15	19	21	36	571	23

Table 9. Task 2 solutions: their primary analysis methods, their methodologies (i) in general as well as with respect to (ii) extraction, (iii) text recognition and (iv) use of machine learning techniques, and evaluation results.

	Solution 2.1	Solution 2.2	Solution 2.3	Solution 2.4	Solution 2.5	Solution 2.6	Solution 2.7	Solution 2.8
Publications	Tkaczyk (2015)	Klampfl (2016)	Nuzzolese (2016)	Sateli (2016)	Kovriguina (2015)	Ronzano (2015)	Ahmad (2016)	Ramesh (2016)
	–	Klampfl (2015)	Nuzzolese (2015)	Sateli (2015)	–	–	–	–
Primary Analysis								
structure-based	✓	✓				✓	✓	✓
linguistic-based		✓	✓	✓	✓	✓	✓	
presentation-based	✓				✓	✓		✓
Methodology								
workflow	parallel pipelines	parallel pipelines	single pipeline	iterative approach	single pipeline	single pipeline	single pipeline	layered approach
external services		✓	✓	✓		✓		
Extraction								
PDF-to-XML	✓	✓		✓ (2016)			✓	✓
PDF-to-HTML					✓			
PDF-to-text			✓	✓ (2015)		✓	✓	
Machine Learning								
supervised	✓	✓	✓	✓		✓		✓
unsupervised	✓	✓						
CRF	✓	✓						✓
Text recognition								
NLP/NER		✓	✓	✓	✓	✓		
heuristics	✓	✓	✓	✓	✓	✓	✓	✓
regEx	✓	✓	✓	✓	✓	✓		✓
Evaluation								
best performing	✓ (2015)						✓ (2016)	
most innovative		✓ (2016)		✓ (2015)				

Table 10. Task 2 solutions: how they address different subtasks to accomplish Task 2. **n/a** stands for subtasks that were not required the year the solution participated in the challenge. **X** stands for subtasks that were not addressed by a certain solution.

Information to extract	Solution 2.1	Solution 2.2	Solution 2.3	Solution 2.4	Solution 2.5	Solution 2.6	Solution 2.7	Solution 2.8
document structure	enhanced docstrum	max entropy, merge & split, clustering	NLP to break the text down in sections & sentences	span between Gazetteer's segment headers	font characteristics, text position	rule-based iterative PDF analysis	heuristics on titles, capital-case and style	level I & II CRF
fragments' classification	SVM	supervised ML	Stanford CoreNLP & NLTK	Gazetteer	font-based blocks & sorting	structural features, chunk- & sentence-based SVM	pattern-matching	level II CRF
authors	SVM (LibSVM)	unsupervised ML & classification	heuristics, NER, CoreNLP	Gazetteer's person first names	e-mail 1st part frequent patterns & string comparison	layout info, ANNIE, external repos	from plain text: start/end identifiers return character	level III CRF
affiliations	CRF	unsupervised ML & classification	NER, statistical rules, patterns	organizations names rules patterns	e-mail 2nd part frequent patterns & string comparison	ANNIE, external repos	from plain text: start/end identifiers return character	level III CRF, affiliation markers, POS, NER
funding	X	NER, sequence classification	'Acknowledgments' section, regEx, number or identifier	'Acknowledgments' section, upper-initial word token or name of organization	'Acknowledgments' section, string-matching: 'support/fundl sponsor', etc.	manual JAPE grammars	'Acknowledgments' section, string matching: 'the'... 'project', etc.	level II CRF
references	CRF	geometrical block segmentation	ParseCit CrossRef	hand-crafting rules for multiple cases	Heuristics on 'References' section	external services	n/a	level III CRF (even though n/a in 2016)
ontologies	X	n/a	match named entities to indexed ontologies	root tokens of ontology names	'Abstract' stop-list of acronyms	JAPE grammars	n/a	n/a
tables & figures	n/a	max entropy, merge & split	X	'Table' 'Figure Fig' trigger words	n/a	n/a	heuristics on captions, string matching	level II CRF
supplementary material	n/a	max entropy, merge & split	X	heuristics on links	n/a	n/a	heuristics on links and string matching	X

Table 11. Implementation details for Task 2 solutions

	Solution 2.1	Solution 2.2	Solution 2.3	Solution 2.4	Solution 2.5	Solution 2.6	Solution 2.7	Solution 2.8
Implementation language								
C++								✓
Java	✓	✓	✓	✓		✓	✓	✓
Python		✓	✓		✓			
PDF character extraction								
Apache PDFBox ⁵¹		✓					✓	✓
iText ⁵²	✓							
Poppler ⁵³						✓		
PDFMiner ⁵⁴			✓		✓			
PDFX ⁵⁵				✓ (2016)		✓	✓	
Xpdf ⁵⁶				✓ (2015)				
Intermediate representation								
HTML					✓			
JSON			✓					
text		✓		✓	✓		✓	
XML	✓ (NLM JATS)					✓	✓	✓ (NLM JATS)
External components								
CrossRef API			✓			✓		
DBpedia Spotlight ⁵⁷				✓		✓		
GATE		✓		✓		✓		
ANNIE ⁵⁸				✓		✓		
FreeCite			✓			✓		
others	GRMM ⁵⁹ , LibSVM ⁶⁰ , Mallet ⁶¹	crfsuite ⁶² , OpenNLP ⁶³ , ParsCit ⁶⁴ ,	FRED, Stanford ⁶⁵ , CoreNLP, NLTK ⁶⁶ , (WordNet ⁶⁷ , BabelNet ⁶⁸)	DBpedia SPARQL end-point	Grab spider ⁶⁹ , BeautifulSoup ⁷⁰	Bibsonomy ⁷¹ , FundRef ⁷²	EDITpad Pro ⁷³	Stanford ⁷⁴ , NERTagger, CRF++ ⁷⁵ , CoNLL ⁷⁶ , JATS2RDF ⁷⁷
(Open Source) License	AGPL-3.0	AGPL-3.0	not specified	LGPL-3.0 ⁷⁸	MIT	not specified	not specified	not specified

392 **Solution 2.1** Tkaczyk and Bolikowski (2015) relied on CERMINE⁷⁹, an open source system for
393 extracting structured metadata and references from scientific publications published as PDF files. It has a
394 loosely captured architecture and a modular workflow based on supervised and unsupervised machine-
395 learning techniques, which simplifies the system's adaptation to new document layouts and styles. It
396 employs an enhanced Docstrum algorithm for page segmentation to obtain the document's hierarchical
397 structure, Support Vector Machines (SVM) to classify its zones, heuristics and regular expressions for
398 individual and Conditional Random Fields (CRF) for affiliation parsing and thus to identify organization,
399 address and country in affiliation. Last, K-Means clustering was used for reference extraction to divide
400 references zones into individual reference strings.

401 **Solution 2.2** Klampfl and Kern (2015, 2016) implemented a processing pipeline that analyzes a PDF
402 document structure incorporating a diverse set of machine learning techniques. To be more precise, they
403 employ unsupervised machine learning techniques (Merge-&-Split algorithm) to extract text blocks and
404 supervised (Max Entropy and Beam search) to extend the document's structure analysis and identify sec-
405 tions and captions. They combine the above with clustering techniques to obtain the article's hierarchical
406 table of content and classify blocks into different meta-data categories. Heuristics are applied to detect the
407 reference section and sequence classification to categorize the tokens of individual references to strings.
408 Last, Named Entity Recognition (NER) is used to extract references to grants, funding agencies, projects,
409 figure and table captions.

410 **Solution 2.3** Nuzzolese et al. (2015b, 2016) relied on the Metadata And Citations Jailbreaker (MACJa –
411 IPA) in 2015, which was extended to the Article Content Miner (ACM) in 2016. The tool integrates hybrid
412 techniques based on Natural Language Processing (NLP, Combinatory Categorical Grammar, Discourse
413 Representation Theory, Linguistic Frames), Discourse Reference Extraction and Linking, and Topic
414 Extraction. It also employs heuristics to exploit existing lexical resources and gazetteers to generate
415 representation structures. Moreover, it incorporates FRED⁸⁰, a novel machine reader, and includes
416 modules to query external services to enhance and validate data.

417 **Solution 2.4** Sateli and Witte (2015, 2016), relying on LODEXporter⁸¹, proposed an iterative rule-based
418 pattern matching approach. The system is composed of two modules: (i) a text mining pipeline based
419 on the GATE framework to extract structural and semantic entities. It leverages existing NER-based
420 text mining tools to extract both structural and semantic elements, employing post-processing heuristics
421 to detect or correct the authors affiliations in a fuzzy manner, and (ii) a LOD exporter, to translate the
422 document annotations into RDF according to custom rules.

423 **Solution 2.5** Kovriguina et al. (2015) relies on a rule-based and pattern matching approach, implemented
424 in Python. Some external services are employed for improving the quality of the results (for instance,
425 DBLP for validating author's data), as well as regular expressions, NLP methods and heuristics for HTML
426 document style and standard bibliographic description. It also relies on an external tool to extract the
427 plain text from PDFs.

428 **Solution 2.6** Ronzano et al. (2015) extended their framework used for Task 1 (and indicated as Solution
429 1.3 above) to extract data from PDF as well. Their linear pipeline includes text processing and entity
430 recognition modules. It employs external services for mining PDF articles and heuristics to validate,
431 refine, sanitize and normalize the data. Moreover, linguistic and structural analysis based on chunk-based
432 & sentence-based SVM classifiers are employed, as well as enrichment by linking with external resources
433 such as Bibsonomy, DBpedia Spotlight, DBLP, CrossRef, FundRef & FreeCite.

434 **Solution 2.7** Ahmad et al. (2016) proposed a heuristic-based approach that uses a combination of
435 tag-/rule-based and plain text information extraction techniques combined with generic heuristics and
436 patterns (regular expressions). Their approach identifies patterns and rules from integrated formats.

437 **Solution 2.8** Ramesh et al. (2016) proposed a solution based on a sequential three-level Conditional
438 Random Fields (CRF) supervised learning approach. Their approach follows the same feature list as
439 Klampfl and Kern (2015). However, they extract PDF to an XML that conforms to the NLM JATS DTD,
440 and generate RDF using an XSLT transformation tool dedicated for JATS.

441 **2.2.4 Tasks Evaluation**

442 The evaluation of the submitted solutions was conducted in a transparent and objective way by measuring
443 precision and recall. To perform the evaluation, we relied on (i) a gold standard and (ii) an evaluation tool
444 which was developed to automate the procedure.

445 **Gold standard** The gold standard used for each task's evaluation was generated *manually*. It consisted
446 of a set of CSV files, each corresponding to the output of one of the queries used for the evaluation.
447 Each file was built after checking the original sources – for instance HTML proceedings in case of
448 Task 1 and PDF papers for Task 2 – and looking for the output of the corresponding query; then, it
449 was double-checked by the organizers. Furthermore, we also made available the gold standard to the
450 participants (after their submission) so as they have the chance to report inaccuracies or inconsistencies.
451 The final manually-checked version of the CSV files was used as input for the evaluation tool.

452 **Evaluation tool** The evaluation tool⁸² compares the queries output provided by the participants (in
453 CSV) against the gold standard and measures precision and recall. It was not made available to the
454 participants after the 2014 edition, it was only made available after the 2015 edition, while it was made
455 available already by the end of the training for the 2016 edition. This not only increased transparency but
456 also allowed participants to refine their tools and address output imperfections, increasing this way the
457 quality of their results.

458 **3 BEST PRACTICES FOR CHALLENGE ORGANIZATION**

459 In this section we discuss lessons learned from our experience in organizing the challenge and from
460 (even unexpected) aspects that emerged while running the challenge. This section presents the lessons
461 learned by looking at the solutions and data produced by the participants. We have grouped the lessons in
462 categories for clarity, even though there is some overlap between them.

463 Moreover, we validated our lessons learned with respect to other Semantic Web Evaluation Challenges,
464 aiming to assess whether the lessons learned from the Semantic Publishing Challenge are transferable to
465 their settings too. Besides the Semantic Publishing Challenge, another five challenges are organized in
466 the frame of the Semantic Web Evaluation Challenges track at the ESWC Semantic Web Conference (cf.
467 Section 2.1). To validate our challenge's lessons learned, we conducted a survey, which we circulated
468 among the organizers of the different Semantic Web Evaluation Challenges. One organizer per challenge
469 filled in the questionnaire, providing representative answers for the respective challenge. Based on our
470 survey's results, we distill generic best practices that could be applied to similar events. Our lessons
471 learned are outlined in this section, together with their validation based on the other challenges, as well as
472 the corresponding distilled best practices.

473 **3.1 Lessons learned from defining tasks**

474 For the Semantic Publishing Challenge, it was difficult to define appealing tasks that bridge the gap be-
475 tween building up initial datasets and exploring possibilities for innovative semantic publishing. Therefore,
476 as discussed in Section 2.2, we refined the challenge's tasks over the years according to the participants'
477 and organizers' feedback.

478 **3.1.1 Task continuity**

479 **Lesson:** In the case of the Semantic Publishing Challenge, the first edition's tasks were well perceived
480 by potential participants and all of them had submissions. In the second edition (2015), in fact, the
481 challenge was re-organized aiming at committing participants to re-submitting overall improved versions
482 of their first edition's submissions. Results were positive, as the majority of the participants of the first
483 edition competed in the second one too. Therefore, task continuity is a key aspect of the Semantic
484 Publishing Challenge, whose tasks in every year are broadly the same as the previous year's edition,
485 allowing participants to reuse their tools to adapt to the new call after some tuning.

486 **Validation:** Three of the other four Semantic Web Evaluation challenges have also been organized for
487 several times. Table 1 shows the sustainability of the challenges considering recency and regularity of
488 revisions over their lifetimes. Task continuity was embraced in all challenges by their participants, who
489 not only resubmitted their solutions but also showed continuously improved performance for all three
490 challenges that had multiple editions, according to the organizers' answers to our survey.

491 **Best Practice:** Tasks should be continued over the course of different editions. Nevertheless, they
492 should be adjusted to pose new challenges that allow the authors of previous editions' submissions
493 to participate again in the challenge, thus offering them incentives to improve their solution, without
494 excluding though new submissions at the same time.

495 **3.1.2 Distinct Tasks**

496 **Lesson:** The initial goal of the Semantic Publishing Challenge was to explore a larger amount of
497 information derived from CEUR-WS.org data and to offer a broad spectrum of alternative options for
498 potential participants but, in retrospect, such heterogeneity proved to become a limitation. One of the
499 main problems we faced was that some of the queries classified under the same task were cumbersome
500 for the participants. For instance, in particular the submissions to Task 2 – extraction from XML and
501 PDF – showed an unexpectedly low performance. The main reason, in our opinion, is that the task was
502 actually composed of two sub-tasks that required different tools and technologies: some queries required
503 participants to basically map data from XML/PDF to RDF, while the others required additional processing
504 of the content. Potential participants were discouraged to participate as they only felt competitive for
505 the one and not for the other. A sharper distinction between tasks would have been more appropriate.
506 In particular, it is important to separate tasks on plain data extraction from those on natural language
507 processing and semantic analysis.

508 **Validation:** According to the results of our survey, the Semantic Web Evaluation challenges were
509 designed with more than one task, more precisely, on average three tasks per challenge. In addition, all
510 the individual tasks of the challenges were defined related to each other but independently at the same
511 time, so that participants could take part in all or some of the tasks. Nevertheless, only two challenges
512 had submissions for all tasks, while three out of five challenges lacked submissions only for one task.
513 All challenges though, according to our survey, split the tasks considering the required competencies to
514 accomplish them. Three out of five challenges even distinguish the training dataset used by each task
515 to render the different tasks even more distinct. This contributes to enabling participation in certain
516 tasks, while more challenging tasks or tasks of different nature are isolated. Thus, participants are not
517 discouraged from participating if they are not competent for these parts; they can still participate in the
518 tasks where they feel competent.

519 **Best Practice:** Splitting tasks with a clear and sharp distinction of the competencies required to
520 accomplish them is a key success factor. Task should be defined taking into consideration the technology,
521 tools and skills required to accomplish them.

522 **3.1.3 Participants involvement**

523 **Lesson:** One of the incentives of the challenge's successive editions was to involve participants in the
524 tasks' definition, because potential tasks or obstacles might be identified more easily, if not intuitively, by
525 them. However, even though we collected feedback from previous years' participants when designing the
526 tasks, we noticed that such a preliminary phase was not given enough attention. Even though participants
527 provided feedback immediately after the challenge was completed they were not equally eager to give
528 feedback when they were asked just before the new edition was launched. Talking to participants, in fact,
529 helped us to identify alternative tasks.

530 **Validation:** It is common practice that challenge organizers ask for the participants' feedback. According
531 to our survey three out of four challenges (including Semantic Publishing Challenge) which had more
532 than one submission took into consideration the participants' feedback to adjust the tasks or to define new.

533 **Best Practice:** Exploiting participants feedback and involving them in the task definition creating a
534 direct link between different editions is a key success factor. The participants' early feedback can help to
535 identify practical needs and correspondingly shape and adjust tasks. Tasks proposed or emerged from the
536 community can be turned into an incentive to participate.

537 **3.1.4 Community traction**

538 **Lesson:** Although the challenge was open to everyone from industry and academia, we originally
539 expected participants from the Semantic Web community. However, the submitted solutions include
540 participants with completely different research focus areas, even without any Semantic Web background.
541 This changed our perception of the core communities in the challenge. In future, one might therefore
542 consider defining a cross-domain task, e.g., using a dataset of publications from the biomedical domain.

543 **Validation:** Evaluating the scientific profiles of participants and the submitted solutions highlights the
544 diversity of professions. The participants of Task 2 are mainly active researchers in the fields of NLP
545 (Natural Language Processing), Text Mining, and Information Retrieval. Submissions to Task 1 are
546 mostly from the Linked Data and semantic publishing communities, addressing various subjects of interest
547 such as User Modeling, Library Science, and Artificial Intelligence. This diversity of professions was
548 acknowledged while inviting the members of the challenge’s program committee, and during the process
549 of assigning them as reviewers to submissions.

550 **Best Practice:** Defining independent tasks and using datasets related to other fields of study can build a
551 bridge across disciplines. The use case dataset contains data about computer science publications, and the
552 super-event of the Semantic Publishing Challenge series, the ESWC conference, is highly ranked, and
553 thus of potential interest to a wide audience, but focused on a dedicated sub-field of computer science.
554 This choice of subject potentially restricts the target audience and the publicity of the challenge; however,
555 with a slight shift of any of these, it becomes possible to involve other research communities.

556 **3.2 Lessons learned from building training and evaluation datasets**

557 The training and output dataset definition are also crucial parts when organizing a challenge. In the
558 Semantic Publishing Challenge case, we experimented with (i) maintaining the same training and output
559 dataset, as well as the same tasks, as in the case of Task 1, and (ii) modifying the dataset but keeping
560 almost the same tasks, as in the case of Task 2 and 3. This way, we bridged the gap between building up
561 initial datasets and exploring possibilities for innovative semantic publishing. As mentioned in Section 2.2,
562 we refined both the datasets and their corresponding tasks over the years according to the participants’
563 and organizers’ feedback.

564 **3.2.1 Dataset continuity**

565 **Lesson:** We noticed benefits of not only continuing the same tasks but also using the same datasets
566 across multiple editions of the challenge. In Task 1 of each edition, we evolved training and evaluation
567 datasets based on the same data source over the three years. Participants were able to reuse their existing
568 tools and extend the previously-created knowledge-bases with limited effort. However, for the other tasks,
569 whose datasets were not equally stable, we had to rebuild the competition every year without being able to
570 exploit the past experience. Once solutions were submitted for Task 2 though and it was repeated with the
571 same dataset in 2016 as in 2015, the Semantic Publishing Challenge immediately gained corresponding
572 profit as for Task 1, as the majority of the submitted solutions were resubmitted. This did not happen with
573 Task 3, which did not gain traction in the first place and changing the training dataset and tasks did not
574 attract submissions. Therefore, the “continuity” lesson is equally applicable to tasks as well as to datasets.

575 **Validation:** Dataset continuity is not as persistent as task continuity for most challenges, but it still
576 occurs. To be more precise, most challenges in principle reuse the same datasets across different editions:
577 two of the four Semantic Web Evaluation challenges with multiple editions reused the same dataset, while
578 the other two did the same except for one of their editions, where a different dataset was considered, albeit
579 one of the same nature.

580 **Best Practice:** Same datasets should be continuously reused over the course of different editions.
581 Nevertheless, eventually substituting them by another dataset of the same nature, where the same tasks
582 and tools are equally applicable, does not harm the challenge.

583 **3.2.2 Single dataset for all tasks**

584 **Lesson:** Similarly, we observed that it is valuable to use the same dataset for multiple tasks. For
585 instance, in the Semantic Web Challenge case, completely different datasets were used for Task 1 and 2
586 for the first edition, but complementary datasets were used for the same tasks during the second and third
587 edition, while Task 3 considered the previous year’s output of Task 1.

588 The participants can extend their existing tools to compete for different tasks, with limited effort. This
589 also opens new perspectives for future collaboration: participants’ work could be extended and integrated
590 in a shared effort for producing useful data. It is also worth highlighting the importance of such uniformity
591 for the organizers. It reduces the time needed to prepare and validate data, as well as the risk of errors and
592 imperfections. Last but not least, it enables designing interconnected tasks and producing richer output.

593 **Validation:** All four Semantic Web Evaluation challenges with multiple editions used the same dataset
594 or subsets of it for all different tasks of the challenge.

595 **Best Practice:** It is clearly beneficial for the challenge to consider the same dataset for all tasks.

596 **3.2.3 Exhaustive output dataset description**

597 **Lesson:** An aspect that was underestimated in the first editions of the Semantic Publishing Challenge
598 was the training and output dataset description. While we completely listed all data sources, we did not
599 provide enough information on the expected output: we went into details for the most relevant and critical
600 examples, but we did not provide the exact expected output for all cases in the training dataset. Such
601 information should have been provided, as it directly impacts the quality of the submissions and helps
602 participants to refine their tools.

603 **Validation:** According to the survey results, the other Semantic Web Evaluation challenges seem to
604 share the same principle about the exhaustive description of the expected output dataset. To be more
605 precise, only one of the Semantic Web Evaluation challenges does not provide a detailed and exhaustive
606 description of the expected output.

607 **Best Practice:** Exhaustive and detailed description of both the training and evaluation dataset is required,
608 as it affects the submissions' quality and helps participants to refine their tools.

609 **3.3 Lessons learned from evaluating results**

610 All three editions of the Semantic Publishing Challenge shared the same evaluation procedure (see
611 Section 2.2.4 for details). However, it presented some weaknesses, especially in the first two editions,
612 which we subsequently addressed. Three lessons are derived from the issues that are explained below.

613 **3.3.1 Entire dataset evaluation**

614 **Lesson:** Even though we asked participants to run their tools on the entire evaluation dataset, we
615 considered only a subset for the final evaluation. The subset has been randomly selected from clusters
616 representing different cases, which participants were required to address. On the one hand, since the
617 subset was representative of these cases, we received a fair indication of each tool capabilities. On the
618 other hand, some submissions were penalized as their tool could have worked well on other values, which
619 were not taken into account for the evaluation. In the second edition, we tried to resolve this issue by
620 increasing the number of evaluation queries, without reaching the desired results though, but causing
621 instead some additional overhead to the participants. In the third edition, we reduced the number of
622 evaluation queries, but we radically increased their coverage to assure that the greatest part of the dataset
623 (or even the whole dataset) is covered.

624 **Validation:** Our lesson learned was validated by our survey in this case too. Only one of the Semantic
625 Web Evaluation challenges does not take into consideration the entire dataset for the evaluation.

626 **Best Practice:** The evaluation method should cover the entire evaluation dataset to be fair, to avoid bias
627 and to reinforce submissions to maintain a high quality across the entire dataset.

628 **3.3.2 Disjoint training and evaluation dataset**

629 **Lesson:** During the first two editions of the Semantic Publishing Challenge, the evaluation dataset was
630 a superset of the training one. This may have resulted in some over-training of the tools, and caused
631 imbalance in the evaluation, as certain tools performed very well for the training dataset but not for the
632 entire dataset. In an effort to avoid this, we made the training and evaluation datasets disjoint for the third
633 edition of the Semantic Publishing Challenge. It is more appropriate to use completely disjoint datasets,
634 as a solution to avoid over-trained tools.

635 **Validation:** Our lesson learned regarding disjoint training and evaluation datasets was validated by the
636 other challenge organizers. Only one of the Semantic Web Evaluation challenges considers an evaluation
637 dataset which is a subset of the training dataset. All the others consider disjoint training and evaluation
638 datasets.

639 **Best Practice:** The training and evaluation dataset should be disjoint to avoid over-trained tools.

640 **3.3.3 Available evaluation tool**

641 **Lesson:** The evaluation was totally transparent and all participants received detailed feedback about
642 their scores, together with links to the open source tool used for the final evaluation. However we were
643 able to release the evaluation tool only after the challenge for the last two editions. The evaluation tool
644 was not made available after the 2014 edition, it was only made available after the 2015 edition, while
645 it was made available by the end of the training for the 2016 edition. It is instead more meaningful to
646 make it available during the training phase, as we did for the challenge's third edition. Participants can
647 then refine their tool and improve the overall quality of their output. Moreover, such an approach reduces
648 the (negative) impact of output imperfections. Though the content under evaluation was normalized and
649 minor differences were not considered as errors, some imperfections were not expected and were not
650 handled in advance. Some participants, for instance, produced CSV files with columns in a different
651 order or with minor differences in the IRI structure. These all could have been avoided if participants
652 had received feedback during the training phase, with the evaluation tool available as a downloadable
653 stand-alone application or as a service.

654 **Validation:** Our lesson learned regarding the availability of the evaluation tool was also validated by our
655 survey. To be more precise, all the Semantic Web Evaluation challenges make the evaluation tool available
656 to the challenge participants. There is only one that does not, but only because there is no evaluation tool.

657 **Best Practice:** The evaluation tool should be made available to the participants as early as possible
658 while the participants are still working with the training dataset and fine tuning their approaches.

659 **3.4 Lessons learned from expected output use and synergies**

660 In all three editions of the Semantic Publishing Challenge, the potential use of the expected output was
661 clearly stated in the call, but not the output dataset license; it was up to the participants to choose one.
662 Moreover, the challenge was disseminated and supported thanks to synergies with other events. In this
663 section, we outline lessons learned regarding how the expected use of the challenge output and synergies
664 reflect on the challenge perspective, also on the participants and their submissions.

665 **3.4.1 Expected output use**

666 **Lesson:** The uppermost goal of the Semantic Publishing challenge was to obtain the best output dataset.
667 To achieve that, it is required to identify the best performing tool, namely the tool that actually produces
668 the best output dataset. This tool – or a refined version – is subsequently used to generate the RDF
669 representation of the whole CEUR-WS.org corpus⁸³. The fact that the submitted tools are expected to be
670 reused becomes a critical issue: participants' submission should not only target the challenge, but they
671 should produce an output that is directly reusable. Therefore, it is in fact critical to state how the results of
672 the challenge will be eventually used, in order to encourage and motivate participants.

673 **Validation:** Three out of the other four Semantic Web Evaluation challenges do clearly mention the
674 expected output use, as the Semantic Publishing Challenge does too.

675 **Best Practice:** The expected output use and conditions should be explicitly specified in advance.

676 **3.4.2 License**

677 **Lesson:** The incentive to organize the Semantic Publishing Challenge was to reuse the output dataset.
678 Thus, having the permission to do so, which is specified through the dataset license, but also to reuse
679 the tool that produces this output to systematically generate the CEUR-WS.org dataset, is of crucial
680 importance. Particular attention should be given to the licensing of the output produced by the participants.
681 We did not explicitly say which license the submitted solutions should have: we just requested from
682 participants to use an open license on data (at least as permissive as the source of data) and we encouraged
683 open-source licenses on the tools (but not mandatory). Most of the participants did not declare which
684 exact license applies to their data. This is an obstacle for its reusability: especially when data come from
685 heterogeneous sources (e.g., paper full texts copyrighted by the individual authors, as well as metadata
686 copyrighted by the workshops' chairs) and are heterogeneous in content and format, as in the case of
687 CEUR-WS.org, it is very important to provide an explicit representation of the licensing information.

688 **Validation:** Like the Semantic Publishing Challenge, none of the other Semantic Web Evaluation
689 Challenges specified the tool or output dataset license. As a result, none of the submitted solutions
690 provided any licensing information, apart from one challenge where some of the submitted solutions
691 provided licensing information. Even though all Semantic Web Evaluation Challenges follow the same
692 practice of not specifying the output dataset potential license, it becomes obvious based on the results that
693 explicitly specifying it is important if the challenge output is desired to be reused.

694 **Best Practice:** The output dataset license should be explicitly requested to be provided for each one
695 of the submitted solutions. Moreover, participants should be advised to respectively specify their tools'
696 licensing information, to enable inference of their potential re-usability.

697 **3.4.3 Conflicts and synergies**

698 **Lesson:** Based on our experience from organizing three editions of the Semantic Publishing Challenge,
699 we realized that the dissemination should happen in a targeted way. To this extent, other events thematically
700 relevant to the challenge are considered important synergies that contribute to generating interest and
701 identifying potential participants: For instance, in the Semantic Publishing Challenge case the fact that
702 the SePublica 2014 workshop on semantic publishing was organized at ESWC 2014 reflected positively
703 on our challenge, since we had fruitful discussions with its participants. Moreover, the fact that results
704 from the first two editions of the Semantic Publishing Challenge (Vahdati et al., 2016) were presented
705 at the SAVE-SD workshop on semantics, analytics, visualization and enhancement of scholarly data
706 (SAVE-SD2016⁸⁴), which was co-located with WWW 2016, contributed to the challenge dissemination's
707 and in particular to an audience both thematically and technologically relevant to the challenge. To
708 the contrary, in 2015, we introduced a task on interlinking and realized possible conflicts with other
709 challenges, like OAEI (Ontology Alignment Evaluation Initiative), which may have resulted in the lack of
710 participation to Task 3 – even though Task 3 did not intend to cover the specialized scope of OAEI, but
711 rather put the interlinking task into the scope of a certain use case that merely served in aligning the tasks'
712 outputs among each other and with other datasets in the LOD Cloud. Therefore, we concluded that it is
713 important not only to generate interest but also to identify and avoid potential conflicts.

714 **Validation:** All Semantic Web Evaluation challenges collaborate with the ESWC conference, as they
715 are co-located with this event. Besides the main conference, which drives the challenges, it appears that
716 most of them, and in particular the most long-standing ones, also collaborate with other events and, in
717 particular, with other workshops. For instance, the QALD challenge collaborates with the CLEF QA
718 track⁸⁵, and the challenge on Semantic Sentiment Analysis collaborates with the workshop on Semantic
719 Sentiment Analysis⁸⁶, which is also co-organized with ESWC. Last, the OKE challenge collaborates with
720 the Linked Data for Information Extraction workshop (LD4IE)⁸⁷ which, in turn, is co-located with ISWC.
721 According to our survey, none of the other challenges experienced conflicts with further challenges.

722 **Best Practice:** Establish synergies with other events that are thematically and/or technologically relevant
723 to reinforce dissemination and to identify potential participants.

724 **4 CHALLENGE SOLUTIONS ANALYSIS**

725 In this section, we discuss observations from the participants' solutions and derive corresponding conclu-
726 sions that can be used in the Linked Data publishing domain. We group the lessons into four categories:
727 tools, ontologies, data and evaluation process, even though there is some overlap between these aspects.

728 **4.1 Lessons learned from the tools**

729 Valuable indications can be derived by looking at the tools implemented by the participants. In particular,
730 we focus on the software used to address Tasks 1 and 2.

731 **4.1.1 Primary Analysis.**

732 **Observation:** The Semantic Publishing Challenge tasks could be addressed by both generic and ad-hoc
733 solutions, as well as different methodologies and approaches; nevertheless, solutions tend to converge.

734 For Task 1, two out of four solutions primarily consisted of a tool developed specifically for this task,
735 whereas the other two solutions only required task-specific templates or rules to be used within their
736 otherwise generic implementations. In the latter case, Solution 1.2 abstracts the extraction rules from the
737 implementation, whereas Solution 1.4 keeps them inline with the implementation. Those two solutions

738 are generic enough to be adapted even to other domains. Even though solutions were methodologically
739 different, four approaches for dealing with the HTML pages prevailed: (i) *structure-based* (relying on
740 the HTML code/structure), (ii) *layout-based* (relying on the Web page layout), (iii) *linguistic-based*, and
741 (iv) *presentation-based*. **Most tools relied on structured-/layout-based approach** (three out of four)
742 and only one on a partially linguistic-based approach (Solution 1.3).

743 As far as Task 2 is concerned, there were different methodologies and approaches combined in different
744 ways. The overall picture is summarized in Table 9 and Table 10. The nature of the task influenced
745 the proposed solutions. In fact the task was composed of two subtasks: (i) identifying the structural
746 components of the PDF papers and (ii) processing the extracted text. Thus, **some solutions mainly**
747 **focused on structure-based analysis** (five out of eight); others gave more relevance to the *linguistic-*
748 *based* analysis (three out of eight) for their primary analysis. Last, up to four used the *linguistic-*
749 *based analysis to complement their primary approach*, while two solutions also used formatting styles/rules
750 to increase the quality of their output (*style-based* analysis).

751 We also observed that most solutions implemented a **modular pipeline**. In particular, the solutions
752 that followed a structure-based analysis had a workflow with a single pipeline, whereas **linguistic-based**
753 **approaches required parallel or iterative pipelines to address different aspects of the solution and**
754 **to increase performance**. It is also worth mentioning that two solutions over eight, one being the 2015
755 most innovative solution, adopted an iterative approach. One of them iterates over the same analysis
756 multiple times to refine the results (Solution 2.4); the other one (Solution 2.8) adopted a layered approach,
757 in which each iteration adds new information to the previously-produced output.

758 **Conclusion:** The solutions were methodologically different among each other, and modular and hybrid
759 solutions prevailed compared to case-specific ones. This is important as case-specific solutions do not
760 extend beyond the scope of challenges, but generic ones do. It is interesting to note that both 2015 and
761 2016 the best solutions for Task 2 relied primarily on structure analysis, whereas the most innovative
762 solutions focused on linguistic analysis. This might indicate that further research on linguistic approaches
763 might bring interesting results for optimizing the output of such tasks. A deep analysis of the structure,
764 in fact, made participants capture more information; on the other hand, these approaches were quite
765 straightforward and less innovative. It is interesting, though, to note here that the best performing tool
766 of 2016 grounded its structured-based approach on a prior linguistic analysis, whereas most solutions
767 grounded their linguistic analysis on a prior structure analysis. Thus, hybrid solutions are obviously
768 required but their execution order should not be taken for granted. It is also worth discussing the recall
769 score of the linguistic-based tools: these tools most probably suffer from noisy text extraction. In fact
770 the three solutions (Solution 2.2, Solution 2.3 and Solution 2.4) that mainly rely on linguistic analysis
771 achieved the lowest recall scores both in 2015 and 2016 editions, even though they showed significant
772 improvement in the latter edition.

773 Similarly, the tool that relied on a linguistic analysis for Task 1 showed significantly lower precision
774 and recall, compared to the other tools, indicating that linguistic-based solutions are not enough, if not
775 supported by a precise structure analysis. Even though the linguistic-based approach was considered a
776 rather innovative way of dealing with Task 1, the evaluation showed that a linguistic-based analysis might
777 not be able to perform as well as a structure-based one.

778 **4.1.2 Methodologies: extraction, intermediate format and machine learning**

779 **Observation:** Diverse methodologies were employed by the participants to extract and analyze content.
780 There were no prevalent approaches, but some tendencies were observed.

781 For Task 1, three out of four solutions considered **rules to extract data from the HTML pages**; two
782 of them considered CSS to define the rules, while the other one, which relied on linguistic-based analysis,
783 considered JAPE; the latter solution was based on crawling. Last, **all solutions used regular expressions**
784 at some point of their workflow.

785 For Task 2, half of the solutions in 2015 but only two out of five in 2016 extracted the text from PDF
786 documents and turned it into plain text. On the contrary, **the majority extracted the text from the PDF**
787 **files but turned it into XML** (two out of six solutions in 2015 and four out of five in 2016). There was
788 only one solution that used HTML as intermediate format. We noted that, **both in 2015 and 2016, the**
789 **best performing solutions relied on a PDF-to-XML extraction**. Moreover, one solution changed from
790 PDF-to-text to PDF-to-XML and indeed performed better in 2016, but we cannot state with high certainty
791 if this was the determining factor. Besides extraction, as far as text analysis is concerned, five solutions

792 in 2015 and four in 2016 relied on supervised Machine Learning. Only two solutions in 2015 and one
793 in 2016 (the same as in 2015) additionally relied on unsupervised Machine Learning to address Task
794 2. Last, **all solutions employed heuristics and regular expressions**. Five out of six solutions in 2015
795 employed Natural Language Processing (NLP) and Named Entity Recognition (NER), and those that also
796 participated in 2015 kept NLP/NER in their workflows in 2016.

797 **Conclusion:** Solutions based on supervised Machine Learning were awarded as the most innovative
798 both in 2015 and in 2016. Therefore, it seems that there is potential on experimenting with supervised
799 Machine Learning approaches to address such a task. Nevertheless, even though the best performing
800 solution in 2015 did use supervised Machine Learning, it is not the case for 2016, which makes us
801 conclude that fundamentally alternative solutions might show good results too. Overall, there is potential
802 for improvement and plenty alternative methodologies can be investigated. The intermediate format used
803 by each solution, on the other hand, had no relevant impact on the final results.

804 4.1.3 Source tools

805 **Observation:** The Semantic Publishing Challenge call did not prescribe (i) the implementation language,
806 (ii) the license, as well as whether the tools should (iii) reuse existing components or external services,
807 and (iv) be open-sourced or not. The participants were allowed to follow their preferred approaches.

808 Three out of four Task 1 solutions, as shown in Table 3, and seven out of eight Task 2 solutions,
809 as shown in Table 11, **primarily relied on Java-based implementations**. In both cases, the remaining
810 solution relied on Python. Two out of eight solutions for Task 2 complemented their Java-based imple-
811 mentations with Python-based parts. Moreover, as it is observed in Table 3, for Task 1, three out of four
812 solutions **relied on tools totally open-sourced**, while the fourth one, the one that addressed both Task 1
813 and Task 2, **relied on a stack of tools which are open-sourced**, but the workflow used was not. This is
814 also observed in most tools for Task 2, as shown in Table 11 (six out of the eight solutions).

815 **MIT⁸⁸ was the most popular license**, with half solutions for Task 1 using it and one out of eight
816 solutions for Task 2, followed by AGPL-3.0⁸⁹, with two out of eight solutions for Task 2 using it. Last,
817 **half of the solutions incorporated external services** to accomplish the tasks (two out of four for Task 1
818 and four out of eight for Task 2). The one of the two solutions for Task 1 that used external services was
819 the one that participated both in Task 1 and Task 2. GATE, DBpedia, CrossRef API⁹⁰, and FreeCite⁹¹ are
820 the most used external services.

821 **Conclusion:** Open-sourced tools prevailed over closed-sourced ones. None of the participants used a
822 totally closed or proprietary software. Most of the them used an open license, and Java and Python based
823 implementations prevailed both for Task 1 and Task 2. The integration of external services was also a
824 valuable solution for the participants.

825 4.2 Lessons learned from models and ontologies

826 In this section, we discuss the different solutions with respect to the data model, the vocabularies and the
827 way they used them to annotate the data.

828 4.2.1 Data model

829 **Observation:** All Task 1 solutions tend to converge regarding the data model, identifying the
830 same core concepts: *Conference*, *Workshop*, *Proceedings*, *Papers*, and *Person*. A few solutions covered
831 more details, for instance, Solution 1.1 identified also the concepts of *Invited Papers* and *Proceedings*
832 *Chair*, while Solution 1.3 captured different types of sessions by identifying additionally the concepts of
833 *Session*, *Keynote Session*, *Invited Session* and *Poster Session*, as well as the concepts of *Organization* and
834 *Topic*. In particular for Task 1, Solution 1.4 domain modeling was inspired by the model used in Solution
835 1.1, with some simplifications, a practice commonly observed in real Linked Data set modeling.

836 In contrast, **Task 2 solutions used more heterogeneous data models**. There are six high-level
837 properties identified by all solutions: *identifier*, *type*, *title*, *authors*, *affiliation* and *country*. Other entities
838 were instead described in different ways and with different granularity. That happened, for instance, to
839 the entities *organization*, *funding agency* and *grant*. In certain cases they are identified as separate entities
840 and in other cases their details constitute part of other entities descriptions (and are expressed as data or
841 object properties). The coverage of the data models was also heterogeneous: for the 2016 edition, for
842 instance, not all solutions identify the *sections* and capture the notion of caption of *figures* and *tables*.

843 **Conclusion:** Based on the aforementioned, we observe a trend of converging in respect to the model
844 the CEUR-WS.org dataset should have according to the submitted solutions. Most solutions converge on
845 the main identified concepts in the data (*Conference, Workshop, Proceedings, Paper* and *Person*) and on
846 the CEUR-WS.org dataset's graph at least for Task 1, namely the publications' metadata. The way the
847 tasks and their corresponding queries are described contributes towards this direction.

848 4.2.2 Vocabularies

849 **Observation:** There is a wide range of vocabularies and ontologies that can be used to annotate scholarly
850 data. Most of the solutions preferred to **(re)use almost the same existing ontologies and vocabularies**,
851 as summarized in Table 4. Six out of twelve solutions for both Task 1 and 2 used the Semantic Web for
852 Research Communities (*swrc*) vocabulary⁹², five used the Bibliographic Ontology (*bibo*) vocabulary⁹³
853 and three used the Semantic Web Conference (*swc*) vocabulary⁹⁴. Moreover, six solutions used one or
854 more vocabularies of the Semantic Publishing and Referencing Ontologies⁹⁵ (*SPAR*). In particular, five
855 solutions used the FRBR-aligned Bibliographic Ontology⁹⁶ (*FaBiO*) ontology, three the Publishing Roles
856 Ontology⁹⁷ (*PRO*), three the Document Components Ontology⁹⁸ (*DoCO*), two the Bibliographic Reference
857 Ontology⁹⁹ (*BiRO*), two the Funding, Research Administration and Projects Ontology¹⁰⁰ (*FRAPO*) and
858 one the Functional Requirements for Bibliographic Records¹⁰¹ (*FRBR*). Besides the domain-specific
859 vocabularies and ontologies, eight solutions used the Dublin Core vocabulary (*dc*¹⁰² and *dcterms*¹⁰³),
860 eight the Friend of a Friend vocabulary¹⁰⁴ (*foaf*), five solutions used the DBpedia ontology¹⁰⁵ (*dbo*), three
861 the VCard¹⁰⁶ (*vcard*) and two the *event*¹⁰⁷ and *timeline*¹⁰⁸ ontologies and *schema.org*¹⁰⁹. Last, there were
862 four solutions that **used their own custom vocabularies, in combination with existing ones** in most
863 cases, but only one used barely its custom vocabulary.

864 In contrast to **Task 1 solutions, which all converged on using same vocabularies and ontologies**
865 intuitively, **Task 2 solutions reused a wider range and relatively different vocabularies and ontolo-**
866 **gies** to annotate same entities appearing in the same data, which is extracted from PDF documents. This
867 is a consequence of the rather diverse data models considered by different solutions. Interestingly, most
868 Task 2 solutions use sub-ontologies of the *SPAR* ontologies family. Last, most solutions reuse the three
869 most popular vocabularies in the education field according to Schmachtenberg et al. (2014). The general
870 purpose vocabularies – such as FOAF – used by the participants are also listed high in the same ranking.

871 **Conclusion:** It is evident that the spirit of vocabulary reuse gains traction. However, it is interesting that
872 different solutions used the same ontologies to annotate the same data differently (see also Section 4.2.3).

873 4.2.3 Annotations

874 **Observation:** Even though **all solutions used almost the same vocabularies, not all of them used**
875 **the same vocabulary terms to annotate the same entities**. As far as Task 1 is concerned, all solutions
876 only converged on annotating *Persons* using the *foaf:Person* class. For the other main concepts the
877 situation was heterogeneous, as reported in Table 6. A few of them also explicitly annotated *Persons* using
878 the *foaf:Agent* class, even though *foaf:Person* is a subclass of *foaf:Agent*. *foaf:Agent*
879 was also used by one of the solutions during the first edition, but it was then replaced by the more explicit
880 *foaf:Person*. The *Conference* concept was well-captured by all solutions.

881 It is interesting to note that, **for the first edition, most solutions used relatively generic vocabulary**
882 **terms**, e.g., *swrc:Event*, *swc:Event* or *swc:OrganizedEvent* to annotate the data. How-
883 ever, **in the second edition, most solutions preferred to use more explicit vocabulary terms for the**
884 **same concept**, e.g., *swrc:Conference* and *bibo:Conference*, while they also maintained the
885 more generic vocabulary terms for events. The same occurred with the *Paper* concept. The 2014 edi-
886 tion datasets were annotated using more generic vocabulary terms, e.g., *swrc:Publication* or even
887 *foaf:Document*, whereas in 2015 more explicit terms were preferred, such as *swrc:InProceedings*
888 or *bibo:Article*. In particular *swrc:InProceedings* was adopted by three out of four solutions.

889 In contrast to **Task 1 solutions, which focus on identifying and describing concrete entities, Task**
890 **2 solutions mainly focus on capturing their properties**. This is also evident from the fact that Task
891 2 solutions rarely provide the entities' types, whereas Task 1 solutions always do, even though this
892 information could be inferred from the properties used. Moreover, Task 2 solutions generate much fewer
893 entities than Task 1 solutions. **All Task 2 solutions use approximately the same number of properties**.
894 It is interesting though to note that solutions that follow in principle the linguistic approach tend to use
895 more predicates, which are more explicit and more descriptive too.

896 **All solutions have approximately the same number of predicates, but their precision is still not**
897 **accurate.** Only one of Task 2 solutions (Solution 2.7) has a significantly higher number of predicates com-
898 pared to the other solutions. This occurs because different URIs are used for the same relationships appear-
899 ing in different files to annotate the data. For instance, the *section-title* property appears with 37 different
900 URIs, such as the following: `<http://ceur-ws.org/Vol-1558/paper5#section-title>`,
901 or `<http://ceur-ws.org/Vol-1303/paper_4#section-title>`. However, such a choice
902 prevents easily identifying same relationships.

903 **DCMI is the vocabulary most frequently used by all solutions for annotating the identifier and**
904 **the title.** *RDF(S)* is also used for the *title* (represented as `rdfs:label`), as well as for the entities'
905 *types*. For the remaining properties, a wide range of different vocabularies are considered, but they do not
906 converge on their choices. Indicatively: one of the solutions considers `schema:mentions` to describe
907 a citation, whereas other solutions consider `bibo:cites` or `biro:references`. In the same context,
908 some solutions associate authors to papers with the `dcterms:creator` property, whereas others con-
909 sider `foaf:maker`. Moreover, some solutions indicate the affiliation using the `swrc:affiliation`
910 property, whereas others use `pro:relatesToOrganization`, or some solutions represent the publi-
911 cation year using `swrc:year`, whereas others use `fabio:hasPublicationYear`. Last, it is inter-
912 esting to note that solutions may even use vocabulary terms that do not exist, such as `swrc:Section`.

913 **Conclusion:** On the one hand, the more familiar the data publishers get with the data, the more explicit
914 they become with the annotations they use and the more they converge on the choices they make. On the
915 other hand, the way different solutions extract particular properties reflects on the final data model.

916 4.3 Lessons learned from submitted RDF datasets

917 In this section, we discuss the different solutions with respect to the RDF dataset they produce.

918 4.3.1 Successive submissions improvements

919 **Observation:** From the first edition to the second edition of the Semantic Publishing Challenge, we
920 noticed that the **participants who re-submitted their solutions had improved the overall dataset**, not
921 only the parts useful to answer the queries. For instance, all three solutions of Task 1 that had participated
922 in both the 2014 and the 2015 editions modified the way they represented their data, and this resulted in
923 corresponding improvements to the overall dataset.

924 Indicatively, as far as Task 1 is concerned, Solution 1.2 addressed a number of shortcomings the
925 previous tool's version had, in particular regarding data transformations, which might have influenced
926 their precision improvement. Heyvaert et al. (2015) also assessed their mappings' quality to verify the
927 schema is valid with respect to the used vocabularies and ontologies. To address the same issue and avoid
928 inconsistencies in their dataset, Solution 1.1 preferred to align different ontologies' classes and properties,
929 e.g., aligning BIBO to the SWRC ontologies, as SWC already has some dependencies on SWRC.

930 As far as Task 2 is concerned, some parts of Solution 2.2, for instance, were changed for participating
931 in the 2016 edition. The authors employed different processing steps of their tool, which were not used in
932 the previous edition, e.g., processing section headings, hierarchy and captions, but they also introduced
933 novel aspects driven by the challenge tasks and queries, e.g., extracting links from supplementary material.
934 Among the changes of Solution 2.4, it was the PDF extraction tool used, whose change might have partially
935 contributed to their recall improvement, while a number of additional or new conditional heuristics most
936 probably led to their precision improvement. Overall, it was observed that improvements to extraction
937 might reflect on the solutions' recall, whereas improvements to text analysis on their precision.

938 **Conclusion:** The improvement of the dataset was evident on some aspects and indeed the results were
939 satisfying, but we still see room for improvement. It is interesting though to note that solutions did not
940 remain focused on improving just the *data extraction* parts of the challenge, but also the *data modeling*,
941 even though the latter is not directly assessed by the challenge.

942 4.3.2 Dataset Structure

943 **Observation:** **The different solutions differ significantly with respect to the size of the produced**
944 **dataset.** This happens for different reasons. Solution 1.1 shows an extraordinary number of triples com-
945 pared to other solutions. This occurs to a certain extent because each concept is annotated with at least two
946 classes, making one fourth of the dataset to be type declarations. Moreover, they include even annotations
947 that indicate the type of the resource or property on a very low level, namely they use `rdfs:Class`,

948 `rdfs:Property`, as well as `owl:ObjectProperty` or `owl:AnnotationProperty` etc., which
949 counts for almost 2,000 triples of the total dataset. Solution 1.4 also shows a high number of triples.
950 This occurs because the same dataset contains triples describing the structure of the HTML page,
951 as well as triples describing the actual content of the pages. Nevertheless, the main reason that
952 causes the flow of triples is the fact that a new URI is generated each time a concept appears in
953 one of the CEUR-WS.org volumes. For instance, the person *Ruben Verborgh* appears to have 9
954 URIs, e.g., `<http://ceur-ws.org/Vol-1034/#RubeniVerborgh>` for the Vol-1034 proceed-
955 ings or `<http://ceur-ws.org/Vol-1184/#RubeniVerborgh>` for the Vol-1184 proceedings.
956 The person *Christoph Lange* appears to have 15 distinct URIs, e.g., for Vol-360 proceedings, the
957 `<http://ceur-ws.org/Vol-360/#ChristophiLange>`, or for Vol-1184 proceedings, the
958 `<http://ceur-ws.org/Vol-1184/#ChristophiLange>`¹¹⁰. Solutions 1.2 and 1.3 are ap-
959 proximately at the same number of triples both for the 2014 and the 2015 editions.

960 **Conclusion:** There is a very high heterogeneity in the produced datasets; although solutions tend to
961 agree on used vocabularies, their design choices are very different and, as a consequence, the number and
962 organization of the triples is very heterogeneous.

963 4.3.3 Coverage

964 **Observation:** We further noticed that **solutions rarely agree upon the extracted information**. For
965 instance, some skip the extraction of wrong data or certain other information. Overall, we observed
966 significant differences with respect to the number of identified entities per category. The results for Task 1
967 are summarized in Table 7 and Table 6, while the results for Task 2 are summarized in Table 8.

968 **Produced datasets were very heterogeneous in term of size, number of triples and entities.** As
969 far as Task 1 is concerned, apparently, Solution 1.1 and Solution 1.3 used the individual pages to identify
970 the proceedings, whereas Solution 1.2 and Solution 1.4 used the index page to identify the proceedings,
971 this is the reason that there is so big difference in the number of *Proceedings* entities. The number of
972 identified papers is also significantly different among the different solutions, but in the *Persons* case we
973 observe the greatest variation in terms of numbers because of **different practices of assigning URIs**; a
974 few solutions reuse URIs across different proceedings volumes, others do not.¹¹¹

975 As far as Task 2 is concerned, **solutions tend to omit certain subtasks and to optimize their**
976 **performance on others due to the nature of the task** – queries were quite heterogeneous, with a clear
977 distinction, for instance, between the analysis of the structural components and of the textual content
978 of the papers. For instance, in 2015, the best performing solution focused on precisely addressing the
979 subtasks which were related to the document structure and totally omitted queries related to funding and
980 ontologies, as shown in Table 10. Similarly, in 2016, certain solutions completely omitted the queries that
981 were related to supplementary material or tables and pictures captions. Consequently, the dataset size, as
982 well as the number of triples and entities significantly diverge among the solutions.

983 **Conclusion:** The datasets' heterogeneity is also evident in the amount and type of information each
984 dataset provides. However, the more the solutions improve, the more the solutions converge at least
985 regarding the number of retrieved and/or distinctly identified entities.

986 4.4 Lessons learned from the solutions with respect to the evaluation

987 In this section, we discuss the different solutions with respect to the dataset evaluation.

988 4.4.1 Ranking

989 **Observation:** For Task 1, in 2015 the performance ranking of the three tools evolved from 2014 has not
990 changed but their performance has improved except for Solution 1.1, which improved precision but recall
991 remain the same. Disregarding the two queries that were new in 2015, Solution 1.1, which had won the
992 best performance award in 2014, performs almost as well as Solution 1.4.

993 The trend was slightly different for Task 2: all tools participating in the Challenge for the second time
994 increased their performance, but the overall ranking changed. Solution 2.4 obtained a higher score than
995 Solution 2.2 in 2016, contrarily to what happened in 2015. The position of Solution 2.3 was stable.

996 **Conclusion:** Continuity helps participants to improve their tools; the overall ranking keeps stable if
997 the tasks (and queries) are kept stable; adjustments to the tasks (and queries) may impact the ranking,
998 favoring one team more than another.

999 **4.4.2 New and legacy solutions**

1000 **Observation:** Task 1 participants both in 2014 and 2015 had an improved version of different aspects of
1001 their solution, which resulted in correspondingly improved versions of the final dataset. The new Solution
1002 1.4, which introduced a fundamentally new approach, achieved equally good results as the best solution
1003 of 2014. The same trend was evident in Task 2, with a general improvement of all solutions that were
1004 re-proposed for the second year (2015 and 2016).

1005 **Conclusion:** Legacy solutions might be able to improve and bring stable and good results, however
1006 there is still room for improvement and mainly for fundamentally new ideas that surpass problems that
1007 legacy solutions cannot deal with.

1008 **4.4.3 Equal chances**

1009 **Observation:** Solution 1.1, the winners of Task 1 in 2014, participated in 2015 with an improved
1010 version but did not win. The 2015 winner was a new tool with a brand new approach (Solution 1.4). The
1011 same happened for Task 2: in 2016, one winner (Solution 2.7) was a brand-new solution, the other one
1012 (Solution 2.2) was an extension and improvement of a legacy solution but did not win the year before.

1013 **Conclusion:** The winners were not the same in subsequent versions of the challenge: creativity won.

1014 **5 DISCUSSION: CHALLENGE IMPACT ON LINKED DATA QUALITY**

1015 In Section 1 we motivated the Semantic Publishing Challenge as a means of producing high-quality
1016 Linked Data. In this section, we assess the potential impact of the challenge on the quality of the Linked
1017 Data produced. To be more precise, the quality of the Linked Data produced by the tools submitted has
1018 been assessed by comparing the output of a number of prescribed queries against our gold standard and
1019 measuring precision and recall, as explained in Section 2.2.4. Assessing the quality of Linked Data by
1020 running queries over it is a common approach, as the comparison of tools by Zaveri et al. (2016) confirms,
1021 whose recent survey we refer to for a comprehensive review of the state of the art regarding Linked Data
1022 quality assessment. Therefore, a challenge designed as the Semantic Publishing Challenge could act as a
1023 means to assess the Linked Data quality, and, the better the results, the higher the Linked Data quality is
1024 expected to be.

1025 The specific quality metrics that our evaluation setup assesses can be connected to the general quality
1026 dimensions (accessibility, intrinsic, contextual and representational) and certain of their corresponding
1027 metrics, as they are identified by Zaveri et al. (2016). Moreover, few other quality dimensions' metrics
1028 that are not covered by the challenge's evaluation are assessed in the frame of this review. Note that some
1029 metrics are applicable for all tasks, whereas others are only for a certain task.

1030 **5.1 Accessibility dimensions**

1031 The accessibility dimensions involve aspects related to the Linked Data access, authenticity and re-
1032 trieval (Zaveri et al., 2016). Our challenge required participants to make their data available, forcing this
1033 way the solutions to cover the availability dimension. Making the data available as an RDF dump was
1034 the minimum requirement set by the challenge, covering this way the accessibility of the RDF dumps
1035 metric. Participants were also encouraged to publish their data via other Triple Pattern Fragment (TPF)
1036 interfaces, such as SPARQL endpoints, but assessing its availability was not part of the challenge's
1037 evaluation. Moreover, participants were encouraged to publish their data using a certain license, without
1038 being a requirement though, boosting this way the licensing dimension (the corresponding detailed
1039 discussion is available in Section 3.4.2). While the aforementioned referred to all challenge's tasks, the
1040 interlinking dimension was only promoted by Task 3, which, after all, is its actual goal. Overall, even
1041 though the submitted solutions only made their datasets available as RDF dumps and did not specify the
1042 license, the challenge achieved to enable solutions to achieve the minimum requirement of making the
1043 produced datasets accessible. It is evident that, if the challenge had turned high values w.r.t. each of the
1044 aforementioned metrics mandatory, the produced dataset accessibility would have been increased.

1045 **5.2 Intrinsic dimensions**

1046 According to Zaveri et al. (2016), the intrinsic dimensions focus on whether the information correctly,
1047 compactly and completely represents the real world and is logically consistent in itself. As the Semantic
1048 Publishing Challenge requires SPARQL queries to be executed against the Linked Data produced by the

1049 different solutions, the syntactic validity of the dataset is a prerequisite, boosting this way the metrics for
1050 syntax error free documents and the absence of malformed datatypes. While our challenge evaluation
1051 covers well the syntactic validity, the semantic accuracy is not evaluated. Nevertheless, the metric which
1052 is related to the misuse of properties is discussed and assessed in a qualitative way in Section 4.2.3 of this
1053 paper, but it is not assessed quantitatively. Similarly, the population completeness, i.e., the percentage
1054 of real-world objects of a particular type that are represented in a dataset, is indirectly evaluated on the
1055 side. Namely, it is not thoroughly assessed if all real-world entities appear, but to successfully answer
1056 the evaluation queries, the population completeness is prerequisite. Moreover, a comparative evaluation
1057 of the population completeness is performed in this work (see more detailed discussion at Section 4.3.3
1058 and Table 7, Table 8). Last, even though the solutions' dataset consistency dimension could have been
1059 evaluated and shed more light to their quality, it was not done by any of the challenge's series so far. All
1060 in all, as the challenge was not focused on assessing the dataset quality, certain metrics of the intrinsic
1061 dimension were not covered intentionally, others were indirectly assessed, while a few others were only
1062 discussed in this paper. Nevertheless, if it had been intended, the challenge could have covered even more
1063 metrics of the intrinsic dimension and could have reinforced the datasets quality even more.

1064 **5.3 Contextual dimensions**

1065 The contextual dimensions highly depend on the context of the task at hand. In the case of relevancy
1066 dimension, the Semantic Publishing Challenge did not perform any relevant evaluation. Nevertheless,
1067 in this paper the coverage metric is addressed. To be more precise, in Section 4.3.3, the coverage is
1068 thoroughly discussed. The Semantic Publishing Challenge does contribute to the timeliness dimension.
1069 To be more precise, thanks to its continuity, it is assured that at least every year the challenge is organized,
1070 a new dataset for the underlying CEUR-WS.org data is generated, boosting the freshness metric. In
1071 particular the final extraction has to be made from the evaluation dataset published a few days before the
1072 final submission deadline. As a conclusion, the challenge succeeded in indirectly promoting the coverage
1073 and timeliness dimensions; however, there is potential for other dimensions to be covered as well.

1074 **5.4 Representational dimension**

1075 The representational dimension captures aspects related to the data design (Zaveri et al., 2016). As far as
1076 the interoperability dimension is concerned, the Semantic Publishing Challenge promotes the reuse of
1077 existing terms and vocabularies and, as shown in Table 4 and discussed in Section 4.2.3, the Semantic
1078 Publishing Challenge achieves its goal of promoting the re-use of existing vocabularies, even though the
1079 corresponding metric is not evaluated automatically. Moreover, thanks to Task 3, the Semantic Publishing
1080 Challenge also promotes the re-use of existing terms. Even though it failed to attract participation, it is
1081 proven that such a task contributes into increasing the overall dataset quality. Thus, the challenge enables
1082 the produced datasets to cover even the representational quality dimension.

1083 **6 CONCLUSIONS**

1084 One of the objectives of the Semantic Publishing Challenge is to produce Linked Data that contributes to
1085 improving scholarly communication. Nevertheless, the lessons learned from organizing this challenge are
1086 not only applicable in the case of a challenge on Semantic Publishing but in the case of other challenges
1087 too. Therefore, this work shed light not only on the three editions of this challenge organized by ourselves
1088 and distilled lessons learned from our experience, but we have also validated them against other challenges
1089 and concluded on general best practices for organizing such challenges. In a nutshell, continuity both
1090 in terms of the dataset and in terms of the tasks is important. Nevertheless, tasks should remain distinct,
1091 but they should refer to the same training and evaluation dataset, while participants' feedback should
1092 be taken into consideration to define or refine the tasks. Regarding the output, the larger the evaluation
1093 dataset is and the less overlapping with the training dataset, the best it is for verifying high coverage. The
1094 sooner the evaluation tool is made available, the better it is for the quality of the final output. Finally, it is
1095 a critical incentive for the participants to know how their output is intended to be reused.

1096 Besides the challenge's organizational aspects, we looked for evidence from the solutions proposed
1097 by the participants. Therefore, we analyzed them, reported our observations and came up with different
1098 conclusions related to Linked Data publishing practices followed by different participants. There are
1099 several positive aspects, among them the high participation and the quality of the produced results. This
1100 work allowed us to share those observations on semantifying scholarly data, using different ontological

1101 models, refining and extending existing datasets. Even though the Semantic Publishing Challenge focuses
1102 on scholarly data, the conclusions we draw based on our analysis are of interest for the entire community
1103 that publishes Linked Data. The possibility of sharing knowledge and solutions among participants was
1104 another key factor of the Semantic Publishing Challenge. In a nutshell, most solutions relied on generic
1105 and open-sourced tools, which allows and enables their reuse for corresponding cases. Solutions, and
1106 thus the tools that produce them, have improved from one edition to the other. Even though different
1107 methodologies were followed, there are certain prevailing approaches – based on structure/layout or on
1108 linguistics – which were instantiated in different ways. Despite the fact that tools diverge, the produced
1109 data model and final annotations converge, as solutions become more mature from one edition to the other,
1110 while well-known vocabularies are reused.

1111 Last, we assessed how the challenge's organization reflects on the submitted solutions' output, namely
1112 how the challenge's organization affects the datasets' quality. We showed that indeed the challenge's
1113 organization may have a positive impact on increasing the quality of the Linked Data produced.

1114 REFERENCES

- 1115 Ahmad, R., Afzal, M. T., and Qadir, M. A. (2016). Information Extraction for PDF Sources based on
1116 Rule-based System using Integrated Formats. In Harald Sack and Stefan Dietze and Anna Tordai and
1117 Christoph Lange (2016).
- 1118 Bertin, M. and Atanassova, I. (2014). Extraction and Characterization of Citations in Scientific Papers. In
1119 Presutti, Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván and Di Iorio, Angelo
1120 and Di Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego and Tordai, Anna (2014),
1121 pages 120–126.
- 1122 Catapano, T. (2010). TaxPub: An Extension of the NLM/NCBI Journal Publishing DTD for Taxonomic
1123 Descriptions. *Journal Article Tag Suite Conference (JATS-Con) Proceedings 2010*.
- 1124 Clough, P. and Sanderson, M. (2013). Evaluating the performance of information retrieval systems using
1125 test collections. *IR (Information Research)*, 18(2):247–375.
- 1126 d'Aquin, M., Drachsler, H., Dietze, S., Herder, E., Parodi, E., and Guy, M. (2014). Lessons Learnt from
1127 LinkedUp – Linking Web Data for Education. In *Multidisciplinary Academic Conference on Education,
1128 Teaching and E-learning*, pages 80–86.
- 1129 Di Iorio, A., Lange, C., Dimou, A., and Vahdati, S. (2015). Semantic Publishing Challenge – Assessing
1130 the Quality of Scientific Output by Information Extraction and Interlinking. In Fabien Gandon and
1131 Elena Cabrio and Milan Stankovic and Antoine Zimmermann (2015), pages 65–80.
- 1132 Di Noia, T., Cantador, I., and Ostuni, V. C. (2014). Linked Open Data-Enabled Recommender Systems:
1133 ESWC 2014 Challenge on Book Recommendation. In Presutti, Valentina and Stankovic, Milan and
1134 Cambria, Erik and Cantador, Iván and Di Iorio, Angelo and Di Noia, Tommaso and Lange, Christoph
1135 and Reforgiato Recupero, Diego and Tordai, Anna (2014), pages 129–143.
- 1136 Dimou, A., Di Iorio, A., Lange, C., and Vahdati, S. (2016). Semantic Publishing Challenge – Assessing
1137 Quality Scientific Output its Ecosystem. In Harald Sack and Stefan Dietze and Anna Tordai and
1138 Christoph Lange (2016).
- 1139 Dimou, A., Vander Sande, M., Colpaert, P., De Vocht, L., Verborgh, R., Mannens, E., and Van de Walle,
1140 R. (2014). Extraction and Semantic Annotation of Workshop Proceedings in HTML using RML. In
1141 Presutti, Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván and Di Iorio, Angelo
1142 and Di Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego and Tordai, Anna (2014),
1143 pages 114–119.
- 1144 Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine Zimmermann, editor (2015). *Semantic
1145 Web Evaluation Challenges: Second SemWebEval Challenge at ESWC 2015, Portorož, Slovenia,
1146 May 31–June 4, 2015, Revised Selected Papers*, number 548 in Communications in Computer and
1147 Information Science, Cham. Springer International Publishing.
- 1148 Freitas, A. and Unger, C. (2015). The Schema-Agnostic Queries (SAQ-2015) Semantic Web Challenge:
1149 Task Description. In Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine Zimmermann
1150 (2015), pages 191–198.
- 1151 Harald Sack and Stefan Dietze and Anna Tordai and Christoph Lange, editor (2016). *The Semantic Web:
1152 ESWC 2016 Challenges, Anissaras, Crete, Greece, May 29–June 2, 2016, Revised Selected Papers*,
1153 number 641 in Communications in Computer and Information Science, Cham. Springer International
1154 Publishing.

- 1155 Heyvaert, P., Dimou, A., Verborgh, R., Mannens, E., and Van de Walle, R. (2015). Semantically
1156 Annotating CEUR-WS Workshop Proceedings with RML. In Fabien Gandon and Elena Cabrio and
1157 Milan Stankovic and Antoine Zimmermann (2015), pages 165–176.
- 1158 Klampfl, S. and Kern, R. (2015). Machine Learning Techniques for Automatically Extracting Contextual
1159 Information from Scientific Publications. In Fabien Gandon and Elena Cabrio and Milan Stankovic
1160 and Antoine Zimmermann (2015), pages 105–116.
- 1161 Klampfl, S. and Kern, R. (2016). Reconstructing the Logical Structure of a Scientific Publication using
1162 Machine Learning. In Harald Sack and Stefan Dietze and Anna Tordai and Christoph Lange (2016).
- 1163 Kolchin, M., Cherny, E., Kozlov, F., Shipilo, A., and Kovriguina, L. (2015). CEUR-WS-LOD: Conversion
1164 of CEUR-WS Workshops to Linked Data. In Fabien Gandon and Elena Cabrio and Milan Stankovic
1165 and Antoine Zimmermann (2015), pages 142–152.
- 1166 Kolchin, M. and Kozlov, F. (2014). A Template-Based Information Extraction from Web Sites with
1167 Unstable Markup. In Presutti, Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván
1168 and Di Iorio, Angelo and Di Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego
1169 and Tordai, Anna (2014), pages 89–94.
- 1170 Kovriguina, L., Shipilo, A., Kozlov, F., Kolchin, M., and Cherny, E. (2015). Metadata Extraction from
1171 Conference Proceedings Using Template-Based Approach. In Fabien Gandon and Elena Cabrio and
1172 Milan Stankovic and Antoine Zimmermann (2015), pages 153–164.
- 1173 Lange, C. and Di Iorio, A. (2014). Semantic Publishing Challenge – Assessing the Quality of Scientific
1174 Output. In Presutti, Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván and Di Iorio,
1175 Angelo and Di Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego and Tordai,
1176 Anna (2014), pages 61–76.
- 1177 Lopez, V., Unger, C., Cimiano, P., and Motta, E. (2013). Evaluating question answering over linked data.
1178 *Web Semantics: Science Services And Agents On The World Wide Web*, 21:3–13.
- 1179 Milicka, M. and Burget, R. (2015). Information Extraction from Web Sources Based on Multi-aspect
1180 Content Analysis. In Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine Zimmermann
1181 (2015), pages 81–92.
- 1182 Miller, H. G. and Mork, P. (2013). From Data to Decisions: A Value Chain for Big Data. *IT Professional*,
1183 15(1):57–59.
- 1184 Nuzzolese, A. G., Gentile, A. L., Presutti, V., Gangemi, A., Garigliotti, D., and Navigli, R. (2015a).
1185 Open Knowledge Extraction Challenge. In Fabien Gandon and Elena Cabrio and Milan Stankovic and
1186 Antoine Zimmermann (2015), pages 3–15.
- 1187 Nuzzolese, A. G., Peroni, S., and Recupero, D. R. (2016). ACM: Article Content Miner for Assessing the
1188 Quality of Scientific Output. In Harald Sack and Stefan Dietze and Anna Tordai and Christoph Lange
1189 (2016).
- 1190 Nuzzolese, A. G., Peroni, S., and Reforgiato Recupero, D. (2015b). MACJa: Metadata and Citations
1191 Jailbreaker. In Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine Zimmermann
1192 (2015), pages 117–128.
- 1193 Presutti, Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván and Di Iorio, Angelo and
1194 Di Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego and Tordai, Anna, editor
1195 (2014). *Semantic Web Evaluation Challenge: SemWebEval 2014 at ESWC 2014, Anissaras, Crete,*
1196 *Greece, May 25–29, 2014, Revised Selected Papers*, number 457 in Communications in Computer and
1197 Information Science, Cham. Springer International Publishing.
- 1198 Ramesh, S. H., Dhar, A., Kumar, R. R., Anjaly, V., Sarath, K., Pearce, J., and Sundaresan, K. (2016).
1199 Automatically Identify and Label Sections in Scientific Journals using Conditional Random Fields. In
1200 Harald Sack and Stefan Dietze and Anna Tordai and Christoph Lange (2016).
- 1201 Reforgiato Recupero, D. and Cambria, E. (2014). ESWC’14 Challenge on Concept-Level Sentiment
1202 Analysis. In Presutti, Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván and Di
1203 Iorio, Angelo and Di Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego and
1204 Tordai, Anna (2014), pages 3–20.
- 1205 Reforgiato Recupero, D., Dragoni, M., and Presutti, V. (2015). ESWC 15 Challenge on Concept-
1206 Level Sentiment Analysis. In Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine
1207 Zimmermann (2015), pages 211–222.
- 1208 Ronzano, F., del Bosque, G. C., and Saggion, H. (2014). Semantify CEUR-WS Proceedings: Towards
1209 the Automatic Generation of Highly Descriptive Scholarly Publishing Linked Datasets. In Presutti,

- 1210 Valentina and Stankovic, Milan and Cambria, Erik and Cantador, Iván and Di Iorio, Angelo and Di
1211 Noia, Tommaso and Lange, Christoph and Reforgiato Recupero, Diego and Tordai, Anna (2014), pages
1212 83–88.
- 1213 Ronzano, F., Fisas, B., del Bosque, G. C., and Saggion, H. (2015). On the Automated Generation of
1214 Scholarly Publishing Linked Datasets: The Case of CEUR-WS Proceedings. In Fabien Gandon and
1215 Elena Cabrio and Milan Stankovic and Antoine Zimmermann (2015), pages 177–188.
- 1216 Sateli, B. and Witte, R. (2015). Automatic Construction of a Semantic Knowledge Base from CEUR
1217 Workshop Proceedings. In Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine
1218 Zimmermann (2015), pages 129–141.
- 1219 Sateli, B. and Witte, R. (2016). An Automatic Workflow for the Formalization of Scholarly Articles’
1220 Structural and Semantic Elements. In Harald Sack and Stefan Dietze and Anna Tordai and Christoph
1221 Lange (2016).
- 1222 Schmachtenberg, M., Bizer, C., and Paulheim, H. (2014). *Adoption of the Linked Data Best Practices in*
1223 *Different Topical Domains*, pages 245–260. Springer International Publishing, Cham.
- 1224 Shotton, D. (2009). Semantic publishing: the coming revolution in scientific journal publishing. *Learned*
1225 *Publishing*, 22(2):85–94.
- 1226 Tkaczyk, D. and Bolikowski, Ł. (2015). Extracting Contextual Information from Scientific Literature
1227 Using CERMINE System. In Fabien Gandon and Elena Cabrio and Milan Stankovic and Antoine
1228 Zimmermann (2015), pages 93–104.
- 1229 Unger, C., Forascu, C., Lopez, V., Ngomo, A.-C. N., Cabrio, E., Cimiano, P., and Walter, S. (2015).
1230 Question answering over linked data (QALD-5). In *CLEF 2015 Working Notes*.
- 1231 Vahdati, S., Dimou, A., Lange, C., and Di Iorio, A. (2016). Semantic Publishing Challenge: Bootstrapping
1232 a Value Chain for Scientific Data. In Alejandra Gonzalez-Beltran and Francesco Osborne and Silvio
1233 Peroni, editor, *Semantics, Analytics, Visualisation: Enhancing Scholarly Data*, Lecture Notes in
1234 Computer Science, Heidelberg. Springer.
- 1235 Williams, J. D., Raux, A., and Henderson, M. (2016). The Dialog State Tracking Challenge Series: A
1236 Review. *Dialogue & Discourse*, 7(3):4–33.
- 1237 Zaveri, A., Rula, A., Maurino, A., Pietrobon, R., Lehmann, J., and Auer, S. (2016). Quality Assessment
1238 for Linked Data: A Survey. *Semantic Web Journal*, 7(1):63–93.

1239 NOTES

- 1240 ¹Semantic Web Challenge; see <http://challenge.semanticweb.org/>
- 1241 ²LinkedUp Challenge <http://linkedup-challenge.org/>
- 1242 ³2014 SemPub Challenge, <http://2014.eswc-conferences.org/semantic-publishing-challenge.html>
- 1243 ⁴2015 SemPub Challenge, <http://2015.eswc-conferences.org/important-dates/call-SemPub>
- 1244 ⁵2016 SemPub Challenge, <http://2016.eswc-conferences.org/assessing-quality-scientific-output-its-ecosy>
- 1245 ⁶2014 Semantic Web Evaluation Challenges, <http://2014.eswc-conferences.org/important-dates/call-challenges.html>
- 1246 ⁷2015 Semantic Web Evaluation Challenges, <http://2015.eswc-conferences.org/call-challenges>
- 1247 ⁸2016 Semantic Web Evaluation Challenges, <http://2016.eswc-conferences.org/call-challenges>
- 1248 ⁹Ontology Matching Challenges <http://ontologymatching.org/>
- 1249 ¹⁰Ontology Alignment Evaluation Initiative <http://oei.ontologymatching.org/>
- 1250 ¹¹World Wide Web Conferences, https://en.wikipedia.org/wiki/International_World_Wide_Web_Conference
- 1251 ¹²Very Large Databases Conferences, <https://en.wikipedia.org/wiki/VLDB>
- 1252 ¹³SEALS infrastructure, <http://oei.ontologymatching.org/2016/seals-eval.html>
- 1253 ¹⁴ISWC Conferences, <http://swsa.semanticweb.org/content/international-semantic-web-conference-iswc>
- 1254 ¹⁵Semantic Web Challenge <http://challenge.semanticweb.org/>
- 1255 ¹⁶QALD Challenge, <http://qald.sebastianwalter.org/>
- 1256 ¹⁷ESWC Conferences, <http://eswc-conferences.org/>
- 1257 ¹⁸CLEF, https://en.wikipedia.org/wiki/Conference_and_Labs_of_the_Evaluation_Forum
- 1258 ¹⁹DBpedia, <https://dbpedia.org>
- 1259 ²⁰LAK Challenges; see http://meco.l3s.uni-hannover.de:9080/wp2/?page_id=18
- 1260 ²¹Linking Data for Education, <http://linkedup-project.eu/>
- 1261 ²²DSTC, <http://workshop.colips.org/dstc5/>
- 1262 ²³AI Mashup Challenge, <http://aimashup.org/>
- 1263 ²⁴SemEval campaigns, <http://alt.qcri.org/semeval2016/>
- 1264 ²⁵CL-SciSumm, <http://wing.comp.nus.edu.sg/cl-scisumm2016/>
- 1265 ²⁶<http://alt.qcri.org/semeval2015/task12/>
- 1266 ²⁷SEMEVAL 2015 workshop, <http://alt.qcri.org/semeval2015/>
- 1267 ²⁸ESWC-CLSA 2015, <https://github.com/diegoref/ESWC-CLSA>

- 1270 ²⁹OKE Challenge 2016, <https://github.com/anuzzolese/oke-challenge-2016#tasks-overview>
- 1271 ³⁰NIF, <http://persistence.uni-leipzig.org/nlp2rdf/>
- 1272 ³¹On a more pragmatic level, a further reason was that one of the challenge organizers, Christoph Lange, has been technical
- 1273 editor of CEUR-WS.org since 2013 and thus has (i) the mandate to advance this publication service technically, and (ii) a deep
- 1274 understanding of the data.
- 1275 ³²CEUR-WS, <http://ceur-ws.org/>
- 1276 ³³JATS, <http://jats.nlm.nih.gov/>
- 1277 ³⁴TaxPub, <https://github.com/plazi/TaxPub>
- 1278 ³⁵2014 CEUR-WS dataset, <https://github.com/ceurws/lod/blob/master/data/ceur-ws.ttl>
- 1279 ³⁶COLINDA, <http://www.colinda.org/>
- 1280 ³⁷DBLP, <http://dblp.l3s.de/dblp++.php>
- 1281 ³⁸Lancet, <http://www.semanticlancet.eu/>
- 1282 ³⁹SWDF, <http://data.semanticweb.org/>
- 1283 ⁴⁰Springer LD, <http://lod.springer.com/>
- 1284 ⁴¹2015 CEUR-WS Task 1 dataset, <http://rml.io/data/SPC2016/CEUR-WS/CEUR-WStask1.rdf.gz>
- 1285 ⁴²2015 CEUR-WS Task 2 dataset, <http://rml.io/data/SPC2016/CEUR-WS/CEUR-WStask2.rdf.gz>
- 1286 ⁴³Collections Ontology, <http://purl.org/co/>
- 1287 ⁴⁴SKOS, <http://www.w3.org/2004/02/skos/core#>
- 1288 ⁴⁵RML, <http://rml.io>
- 1289 ⁴⁶R2RML, <https://www.w3.org/TR/r2rml/>
- 1290 ⁴⁷CSS3, <https://www.w3.org/TR/selectors/>
- 1291 ⁴⁸RMLProcessor, <https://github.com/RMLio/RML-Mapper>
- 1292 ⁴⁹GATE, <https://gate.ac.uk/>
- 1293 ⁵⁰FITLayout framework, <http://www.fit.vutbr.cz/~burgetr/FITLayout/>
- 1294 ⁵¹Apache PDFBox, <https://pdfbox.apache.org/>
- 1295 ⁵²iText, <http://itextpdf.com/>
- 1296 ⁵³Poppler, <https://poppler.freedesktop.org/>
- 1297 ⁵⁴PDFMiner, <http://www.unixuser.org/~euske/python/pdfminer/>
- 1298 ⁵⁵PDFX, <http://cs.unibo.it/save-sd/2016/papers/html/pdfx.cs.man.ac.uk>
- 1299 ⁵⁶Xpdf, <http://www.foolabs.com/xpdf/>
- 1300 ⁵⁷DBpedia Spotlight, <http://spotlight.dbpedia.org/>
- 1301 ⁵⁸ANNIE, <https://gate.ac.uk/sale/tao/splitch6.html>
- 1302 ⁵⁹GRMM, <http://mallet.cs.umass.edu/grmm/>
- 1303 ⁶⁰LibSVM, <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- 1304 ⁶¹Mallet, <http://mallet.cs.umass.edu/>
- 1305 ⁶²crfsuite, <http://www.chokkan.org/software/crfsuite/>
- 1306 ⁶³OpenNLP, <https://opennlp.apache.org/>
- 1307 ⁶⁴ParsCit, <http://wing.comp.nus.edu.sg/parsCit/>
- 1308 ⁶⁵Stanford CoreNLP, <http://stanfordnlp.github.io/CoreNLP/>
- 1309 ⁶⁶NLTK, <http://www.nltk.org/>
- 1310 ⁶⁷WordNet, <https://wordnet.princeton.edu/>
- 1311 ⁶⁸BabelNet, <http://babelnet.org/>
- 1312 ⁶⁹Grab spider, <http://grablib.org/>
- 1313 ⁷⁰BeautifulSoup, <http://www.crummy.com/software/BeautifulSoup/>
- 1314 ⁷¹Bibsonomy, <http://www.bibsonomy.org/help/doc/api.html>
- 1315 ⁷²FundRef, <http://www.crossref.org/fundingdata/>
- 1316 ⁷³EDITpad Pro, <https://www.editpadpro.com/>
- 1317 ⁷⁴Stanford NERTagger, <http://nlp.stanford.edu/software/CRF-NER.shtml>
- 1318 ⁷⁵CRF++, <https://taku910.github.io/crfpp/>
- 1319 ⁷⁶CoNLL, <http://www.cnts.ua.ac.be/conll2000/chunking/>
- 1320 ⁷⁷JATS2RDF, <https://github.com/Klortho/eutils-org/wiki/JATS2RDF>
- 1321 ⁷⁸LGPL-3.0, <https://opensource.org/licenses/lgpl-3.0.html>
- 1322 ⁷⁹CERMINE, <http://cermine.ceon.pl/>
- 1323 ⁸⁰FRED, <http://wit.istc.cnr.it/stlab-tools/fred>
- 1324 ⁸¹LODeXporter, <http://www.semanticsoftware.info/lodexporter>
- 1325 ⁸²SemPubEvaluator, <https://github.com/angelobo/SemPubEvaluator>
- 1326 ⁸³The extraction tool's integration in the CEUR-WS.org production workflow is still in progress but expected to conclude in 2016.
- 1327 ⁸⁴SAVE-SD2016 Workshop, <http://cs.unibo.it/save-sd/2016/>
- 1328 ⁸⁵CLEF QA track, <http://nlp.uned.es/clef-qa/>
- 1329 ⁸⁶Semantic Sentiment Analysis Workshop, <http://www.maurodragoni.com/research/opinionmining/events/>
- 1330 ⁸⁷LD4IE2016 Workshop, <http://web.informatik.uni-mannheim.de/ld4ie2016/LD4IE2016/Overview.html>
- 1331 ⁸⁸MIT, <http://opensource.org/licenses/mit-license.html>
- 1332 ⁸⁹AGPL-3.0, <https://www.gnu.org/licenses/agpl-3.0.en.html>
- 1333 ⁹⁰CrossRef API, <http://api.crossref.org/>
- 1334 ⁹¹FreeCite, <http://freecite.library.brown.edu/>
- 1335 ⁹²SWRC, <http://swrc.ontoware.org/ontology#>
- 1336 ⁹³bibo, <http://purl.org/ontology/bibo/>
- 1337 ⁹⁴SWC, <http://data.semanticweb.org/ns/swc/ontology#>
- 1338 ⁹⁵SPAR, <http://www.sparontologies.net/>
- 1339 ⁹⁶FaBiO, <http://purl.org/spar/fabio/>

1340 ⁹⁷PRO, <http://purl.org/spar/pro/>

1341 ⁹⁸DoCO, <http://purl.org/spar/doco/>

1342 ⁹⁹BiRO, <http://purl.org/spar/biro/>

1343 ¹⁰⁰FRAPO, <http://purl.org/cerif/frapo/>

1344 ¹⁰¹FRBR, <http://purl.org/spar/frbr/>

1345 ¹⁰²DC, <http://purl.org/dc/elements/1.1/>

1346 ¹⁰³DCTerms, <http://purl.org/dc/terms/>

1347 ¹⁰⁴FOAF, <http://xmlns.com/foaf/0.1/>

1348 ¹⁰⁵DBO, <http://dbpedia.org/ontology/>

1349 ¹⁰⁶VCard, <http://www.w3.org/2006/vcard/ns#>

1350 ¹⁰⁷event ontology, <http://purl.org/NET/c4dm/event.owl#>

1351 ¹⁰⁸timeline ontology, <http://purl.org/NET/c4dm/timeline.owl#>

1352 ¹⁰⁹Schema.org, <http://schema.org>

1353 ¹¹⁰The definition of Task 1 was not explicit with regard to whether different persons with the same name (within or across different
1354 workshops proceedings volumes) should be assumed to be the same person or not. Our current work towards the release of a
1355 consolidated CEUR-WS.org dataset shows that the far majority of same names refers to the same person, which is plausible as
1356 CEUR-WS.org focuses on the relatively small computer science community. However, a general solution would be wrong to simply
1357 assume that same names mean same persons, whereas a full disambiguation of names would require a lot of information to be taken
1358 into account beyond the proceedings' tables of content: the title pages of the PDF papers plus possibly external resources.

1359 ¹¹¹Our instructions did not prescribe whether or not participants should assume persons with the same name to be the same. In
1360 the reality of the CEUR-WS.org data, there are very few cases in which the same name refers to two different persons, as the data
1361 covers the relatively small domain of computer science researchers.