

Understanding and Improving Ontology Reasoning Efficiency through Learning and Ranking

Yong-Bin Kang^a, Shonali Krishnaswamy^a, Wudhichart Sawangphol^b, Lianli Gao^c,
Yuan-Fang Li^{d,*}

^a Department of Computer Science and Software Engineering, Swinburne University of Technology, Australia

^b Faculty of Information and Communication Technology, Mahidol University, Thailand

^c School of Computer Science and Engineering,

The University of Electronic Science and Technology of China

^d Faculty of Information Technology, Monash University, Australia

Abstract

Ontologies are the fundamental building blocks of the Semantic Web and Linked Data. Reasoning is critical to ensure the logical consistency of ontologies, and to compute inferred knowledge from an ontology. It has been shown both theoretically and empirically that, despite decades of intensive work on optimising ontology reasoning algorithms, performing core reasoning tasks on large and expressive ontologies is time-consuming and resource-intensive. In this paper, we present the meta-reasoning framework $R_2O_2^*$ to tackle the important problems of understanding the source of TBox reasoning hardness and predicting and optimising TBox reasoning efficiency by exploiting machine learning techniques. $R_2O_2^*$ combines state-of-the-art OWL 2 DL reasoners as well as an efficient OWL 2 EL reasoner as components, and predicts the most efficient one by using an ensemble of robust learning algorithms including XGBoost and Random Forests. A comprehensive evaluation on a large and carefully curated ontology corpus shows that $R_2O_2^*$ outperforms all six component reasoners as well as AutoFolio, a robust and strong algorithm selection system.

Keywords: OWL, Reasoning, Performance prediction, Ontology, Metrics, Learning, Meta-reasoning, Semantic Web

*Author for correspondence

Email addresses: ykang@swin.edu.au (Yong-Bin Kang), skrishnaswamy@swin.edu.au (Shonali Krishnaswamy), wudhichart.saw@mahidol.edu (Wudhichart Sawangphol), lianli.gao@uestc.edu.cn (Lianli Gao), yuanfang.li@monash.edu (Yuan-Fang Li)

1. Introduction

2 Ontologies are essential building blocks of the Semantic Web. Expressive ontol-
3 ogy languages OWL DL and OWL 2 DL are widely used to represent many complex
4 phenomena in a number of application domains, including bioinformatics [1], software
5 engineering [2] and data management [3–6]. In these domains, maintaining the logical
6 correctness of ontologies (i.e. consistency checking) and deducing implicit facts from
7 ontologies (i.e. classification) are both important tasks that may need to be performed
8 repeatedly. However, ontologies as expressed in common ontology languages such as
9 OWL [7] and OWL 2 [8] can be large, complex, or both. The high worst-case com-
10 plexity of these ontology languages incurs high computational costs on the above core
11 reasoning problems. Checking the logical consistency of an ontology in $\mathcal{SHOIN}(\mathbf{D})$,
12 the description logic (DL) underlying OWL DL, has NEXPTIME -complete worst-case
13 complexity [7]. The complexity of the same problem for $\mathcal{SROIQ}(\mathbf{D})$, the DL underlying
14 OWL 2 DL, is even higher (2NEXPTIME -complete) [8].

15 The past decade has seen the development of highly optimised inference algo-
16 rithms for description logics, with (hyper) tableau algorithms [9] being a leading ex-
17 emplar. A number of high-performance DL reasoners have been developed, including
18 FaCT++ [10], HermiT [11], Konclude [12], Pellet [13] and TrOWL [14]. Despite the
19 tremendous progress in both theoretical research and practical implementation, the high
20 theoretical worst-case complexity results for OWL DL and OWL 2 DL still imply that
21 core reasoning services may be computationally very expensive. It has been shown
22 empirically that reasoning on large and complex ontologies in OWL 2 DL and OWL
23 2 EL (a less expressive profile that enjoys a PTIME -complete complexity) can be very
24 time-consuming for state-of-the-art reasoners [15, 16]. Such high difficulty of reasoning
25 and the fundamental role inference plays in ontology-based applications make it highly
26 desirable to be able to accurately predict inference performance for ontologies and
27 reasoners.

28 It is well-known that worst-case complexity does not necessarily provide useful
29 insights into hardness of individual instances [17, 18]. In this context, it is noteworthy

30 that reasoner benchmarking has been conducted previously [15, 19–22]. Except the two
31 ORE competitions, these works only compared inference performance on a small set
32 of ontologies. Moreover, they did not attempt to correlate characteristics of ontologies
33 with their inference performance. Hence, they do not provide insight into what makes
34 inference difficult on a given ontology.

35 The *robustness* of ontology reasoners was recently investigated [23], with a particular
36 focus on reasoning efficiency. It was observed that given a corpus of ontologies and a
37 number of state-of-the-art reasoners, it is highly likely that one of the reasoners performs
38 sufficiently well on any given ontology in the corpus. However, this *virtual best reasoner*
39 is only found *a posteriori*, and the paper did not discuss how the best reasoner may be
40 selected automatically. It only stated that this task is not straightforward.

41 In our previous work we studied the characterisation of ontology’s design complexity
42 using metrics [24], the prediction of ontology classification efficiency [25, 26], and
43 proposed a *meta-reasoner* R_2O_2 [27]. A meta-reasoner is one that combines other
44 (component) reasoners. Given an ontology, a meta-reasoner predicts the most efficient
45 component reasoner and selects it to carry out reasoning on that ontology. In this
46 paper, we improve upon our existing work and present a learning- and ranking-based
47 framework for the understanding of sources of ontology reasoning hardness, prediction
48 of ontology reasoning time, and ultimately improving reasoning performance, under
49 a unifying meta-reasoning framework. The main contributions of this paper can be
50 summarised as follows:

- 51 • **Accurate prediction models:** A regression based prediction model is learned for
52 each of a number of state-of-the-art OWL 2 DL reasoners. Evaluated with 10-
53 fold cross validation, all the models are highly accurate, with R^2 (*coefficient of*
54 *determination*) values in [0.71, 0.95].
- 55 • **A meta-reasoning framework:** A novel meta-reasoning framework, $R_2O_2^*$, is de-
56 veloped. Building on our previous work, $R_2O_2^*$ ranks and selects OWL reasoners
57 with the aim of determining the most efficient reasoner for an unknown ontology.
58 Compared with R_2O_2 [27], $R_2O_2^*$ utilises a robust, state-of-the-art prediction model
59 XGBoost [28] based on gradient boosting [29] and an ensemble learning method
60 *stacking* that combines multiple learning algorithms to obtain better predictive per-

61 formance. $R_2O_2^*$ integrates state-of-the-art, *sound* and *complete* reasoners that are
62 both efficient and robust. Moreover, $R_2O_2^*$ also incorporates ELK [30], an efficient
63 reasoner for OWL 2 EL ontologies, to further improve reasoning efficiency for a
64 wider variety of ontologies.

65 • **Comprehensive evaluation:** A comprehensive evaluation on reasoning time has
66 been conducted on a modern, large, and carefully-curated ontology corpus from the
67 ORE 2015 reasoner competition [22]. Our evaluation shows that $R_2O_2^*$ outperforms
68 all of the six OWL 2 DL component reasoners in the evaluation. The complete
69 meta-reasoner variant, $R_2O_2^*_{(all)}$, also outperforms our previous meta-reasoners PR
70 and R_2O_2 [27]. More importantly, $R_2O_2^*$ outperforms AutoFolio [31], a state-of-
71 the-art portfolio-based algorithm selection model that has demonstrated excellent
72 performance in a number of domains.

73 • **Ontology metrics:** Furthermore, we give full definitions of a large suite of 91 on-
74 tology metrics that have been mentioned but not formally defined in our previous
75 work [24, 26]. The formal definitions provide a valuable insight into ontology en-
76 gineering and maintenance. We also identify the most important metrics that affect
77 ontology reasoning efficiency, which further informs ontology engineering practices.

78 The $R_2O_2^*$ meta-reasoner, including components that calculate metrics and train
79 prediction models and rankers, has been made freely available for wider dissemination.¹

80 The rest of the paper is organised as follows. A brief overview of background
81 knowledge of ontologies and reasoning is given in Section 2, followed by a discussion
82 of closely-related work in Section 3. The suite of all metrics that characterise the design
83 complexity of OWL ontologies are formally defined in Section 4. The four variants
84 of the meta-reasoner $R_2O_2^*$ are described in detail in Section 5. Section 6 presents
85 the evaluation framework for building prediction models and the meta-reasoner. Our
86 detailed evaluation results and analysis are presented in Section 7. Lastly, we conclude
87 the paper and discuss future work in Section 8.

¹<https://github.com/liyuanfang/r2o2-star>

88 2. Ontologies and Reasoning

89 Ontologies organise domain knowledge in a structured and logical way. Semantic
90 Web ontologies have been widely used in many different areas as a medium for knowl-
91 edge representation and data integration. Common ontology languages such as OWL
92 1 [7] and OWL 2 [8] have formal semantics defined by Description Logics (DL) [32], a
93 family of logics created specifically for the purpose of knowledge representation. In
94 simple terms, knowledge in DL is characterised by abstract *concepts*, which represent
95 sets of entities; *properties* or roles, which are binary relations between entities; and *indi-*
96 *viduals*, which represent entities themselves. Hence, a concept semantically represents
97 a set of individuals.

98 The logical nature gives rise to reasoning support for ontologies. These, among
99 others, include concept satisfiability checking, concept subsumption checking, and
100 classification. Concept satisfiability checking ensures that a concept can contain at least
101 one individual. Concept subsumption checks whether two concepts have a sub-class
102 relationship. Classification computes the subsumption relationship between all pairs of
103 (named) concepts in an ontology.

104 For example, the following axioms in the *DL syntax* describe some knowledge about
105 pizzas. Axioms (1) and (2) state that *AmericanHot* is a *Pizza*, and that it has some
106 *MozzarellaTopping*. Axiom (3) states that *MozzarellaTopping* is a *CheeseTopping*.
107 Axiom (4) states that *CheesyPizza* is exactly those *Pizza* that has some *CheeseTopping*.
108 Through classification, the concept *AmericanHot* is found to be a subclass of *CheesyPizza*,
109 as it is a *Pizza*, and that it also has some *MozzarellaTopping*, which is a type of
110 *CheesyTopping*.

$$\textit{AmericanHot} \sqsubseteq \textit{Pizza} \tag{1}$$

$$\textit{AmericanHot} \sqsubseteq \exists \textit{hasTopping.MozzarellaTopping} \tag{2}$$

$$\textit{MozzarellaTopping} \sqsubseteq \textit{CheeseTopping} \tag{3}$$

$$\textit{CheesyPizza} \equiv \textit{Pizza} \sqcap \exists \textit{hasTopping.CheeseTopping} \tag{4}$$

111 A number of different DLs have been proposed over the years. These DLs include

112 different combinations of language constructs, hence they have different expressive
113 power. As a result, they also have different worst-case complexity results for core
114 reasoning tasks. A detailed introduction to the syntax, semantics and complexity of
115 these ontology languages can be found in the literature [7, 8].

116 Optimisation of ontology reasoning algorithms has been aggressively pursued over
117 the past decades, and a number of highly optimised reasoners have been produced.
118 These include sound and complete reasoners such as FaCT++ [33], HermiT [34], Kon-
119 clude [12], and Pellet [13]; as well as the sound but incomplete reasoner TrOWL [14].
120 Despite tremendous progress in optimisation, there has been ample empirical evidence
121 of the actual *hardness* of real-world ontologies [15, 25, 35]. Therefore, efficient reason-
122 ing over large and expressive ontologies remains a computationally challenging task.
123 On the other hand, efficient reasoners dedicated to less expressive profiles have also
124 been developed. ELK [30] is such a concurrent reasoner for ontologies in the OWL 2
125 EL profile.

126 **3. Related Work**

127 Our related work is divided into the three categories according to the focuses in this
128 paper: ontology metrics, prediction models for OWL reasoners, and algorithm selection
129 and meta-reasoners.

130 ***Ontology Metrics.*** There has been research on the development of a series of metrics for
131 analysing ontology complexity [36]. The pioneering work for identifying the proposed
132 metrics in this paper is found in [24] that introduced a suite of 8 metrics with the
133 aim of characterising different aspects of ontology *design complexity*. Further, we
134 identified an additional 19 metrics that can measure different aspects of the size and
135 structural characteristics of an ontology [25]. These 27 metrics were used for predicting
136 discretised reasoning performance of reasoners. These 27 metrics were combined with
137 another set of metrics that capture ontology complexity in [26]. In total, 91 metrics
138 were collected and used to build models for predicting absolute reasoning performance
139 of reasoners. In this work, we use these 91 metrics to build prediction models and the

140 meta-reasoner $R_2O_2^*$. We also give the full definitions of all the 91 metrics, which have
141 not been previously defined formally.

142 ***Prediction of Reasoning Performance.*** Ontology reasoning tasks are hard decision
143 problems that may go beyond NP-hard. For very expressive DLs, ontology reasoning has
144 a very high worst-case complexity of 2NEXPTIME-complete [8]. For ontology reasoning
145 optimisation, the research community have been interested in benchmarking of reasoner
146 performance. In [15, 16], a number of modern reasoners were compared, and it was
147 observed that the reasoners exhibit significantly different performance characteristics,
148 thereby choosing an efficient reasoner for an ontology is a non-trivial task.

149 In a previous work [25], we developed classifiers to predict ontology classification
150 performance categories for FaCT++, HermiT, Pellet and TrOWL, using ontology metrics
151 as predictors [25]. The raw reasoning time is discretised into 5 increasingly large
152 categories. High prediction accuracy of over 80% is achieved for all the 4 reasoners.
153 Although highly accurate, the limitation of this work is that only the hardness category
154 is predicted, not the actual reasoning time. To overcome this problem, we further
155 investigated regression-based prediction models [26] to predict actual (or absolute)
156 reasoning time of reasoners. In this approach, regression analysis was applied to
157 estimate a numeric *response variable* (i.e. predicted reasoning time) from some *predictor*
158 *variables* (i.e. 91 ontology metrics). These regression models were built on a small
159 number of ontologies (i.e. 451) for 6 reasoners (FaCT++, HermiT, JFact, MORE, Pellet,
160 TrOWL). These were implemented in R^2 . In this work, we improve upon these models
161 by using a modern, carefully curated dataset of 1,920 ontologies from the ORE 2015
162 reasoner competition [22], an additional robust learning algorithm XGBoost [28], and
163 an updated list of reasoners that includes Konclude [12] and ELK [30] (for OWL 2
164 EL ontologies only) and excludes TrOWL [14] as it is an approximate thus incomplete
165 reasoner.

166 Sazonau et al. [37] proposed a *local* approach to predicting OWL reasoner effi-
167 ciency. Small subsets of a given ontology are repeatedly created, on which reasoning

²<https://www.r-project.org/>

168 is performed. Reasoning time data is then used to extrapolate a reasoner’s discretised
169 reasoning time on the whole ontology. Principal component analysis (PCA) was also
170 employed to reduce the number of features (metrics). Evaluation conducted on 357
171 ontologies and 3 reasoners shows that the local prediction method performs as well as
172 the *global* approach [25]. Moreover, they observed that the prediction model based on
173 one feature (number of axioms) has comparable performance as that using a set of 57
174 features.

175 In a similar spirit, we investigated the prediction of reasoning time of ABox-intensive
176 OWL 2 EL ontologies [38] and energy consumption of reasoning tasks on the Android
177 platform [39].

178 **Algorithm selection and meta-reasoner.** Algorithm selection [40] is the problem of
179 selecting a well-performing algorithm for a given problem instance. It has been success-
180 fully applied to machine learning, combinatorial optimisation and constraint satisfaction
181 problems [41, 42]. SATzilla [43], for instance, a portfolio-based SAT solver, has demon-
182 strated higher efficiency over single solvers. Compared to SAT, ontology languages are
183 more expressive with the inclusion of many more language constructs. As a result, it is
184 more challenging to accurately characterising ontology complexity.

185 AutoFolio [31] is a state-of-the-art, general-purpose algorithm selection system that
186 performs automatic algorithm selection as well as hyper-parameter tuning. In this paper
187 we use AutoFolio as a strong baseline to evaluate $R_2O_2^*$ in Section 7.2.3.

188 CHAINSAW [44] first proposed the notion of a *metareasoner* for OWL ontologies.
189 Given a *query* (i.e. reasoning task) on an ontology, CHAINSAW constructs the smallest
190 possible subset of the ontology while guaranteeing completeness of answering the query.
191 This is achieved through the extraction of *locality-based modules* [45] using atomic
192 decomposition [46]. The size of the extracted module is dependent on the reasoning
193 task. For certain tasks such as consistency checking, the entire ontology needs to
194 be extracted, hence not resulting in gains in efficiency. Also, given the potentially
195 substantial overhead of computing modules, CHAINSAW may not be competitive for
196 simpler ontologies. As a prototype reasoner, CHAINSAW uses FaCT++ version 1.5.3 as the
197 delegate (i.e. component) reasoner. In the ORE 2015 ontology reasoner competition [47],

198 CHAINSAW , CHAINSAW did not perform competitively against state-of-the-art reasoners:
199 it was ranked 10/10 for the task of OWL DL classification and 11/13 for OWL EL
200 classification.

201 WSReasoner [48] is a hybrid reasoner designed for large and complex ontologies
202 in the description logic \mathcal{ALCHOI} . Given an ontology O , WSReasoner builds two
203 approximate ontologies: a weakened version O_{wk} and a strengthened version O_{str} , both
204 of which are in the less expressive (thus less complex) logic \mathcal{ALCH} . WSReasoner
205 employs two component reasoners: a consequence-based reasoner that classifies both
206 O_{wk} and O_{str} . As reasoning over O_{str} may not be sound, WSReasoner also employs a
207 tableau-based reasoner to verify these results obtained on O_{str} . In its evaluation on a
208 number of well-known hard ontologies including DOLCE, FMA and variants of Galen,
209 WSReasoner outperforms tableau-based reasoners FaCT++, HermiT, Pellet and the
210 approximate, consequence-based reasoner TrOWL.

211 In our preliminary work [27], we proposed a meta-reasoner, R_2O_2 , that makes use of
212 regression-based prediction models of six OWL 2 DL reasoners (i.e. FaCT++, HermiT,
213 JFact, Konclude, MORE, TrOWL). R_2O_2 takes two steps in the training phase. First,
214 given training ontologies characterised by a set of metrics [26] and their reasoning
215 time by the reasoners, R_2O_2 constructs a regression-based prediction model for each
216 of the six reasoners. Second, given another set of training ontologies, a *ranking matrix*
217 is generated using the prediction models. In the ranking matrix, each row represents
218 the values of the ontology metrics and a ranking of the reasoners according to their
219 *predicted* reasoning time. Several rankers were trained on this ranking matrix to learn
220 how ontology metrics can be mapped to a relative ordering by the *predicted* performance
221 of the reasoners. In the actual reasoning (testing) phase, given an unknown ontology,
222 R_2O_2 makes performance predictions for the reasoners. It then ranks the reasoners
223 according to their predicted reasoning time. The rankings recommended by the trained
224 rankers are averaged to determine a unique rank of each reasoner. The highest ranked
225 reasoner is chosen to perform the reasoning task for the unknown ontology. The
226 evaluation on R_2O_2 [27] shows that R_2O_2 outperforms all of the six state-of-the-art
227 OWL 2 DL reasoners, including Konclude [12], the most efficient OWL 2 DL reasoner.

228 As a baseline model to evaluate R_2O_2 [27], we also constructed a portfolio-based

229 OWL reasoner PR, which always selects the most efficient reasoner for any given
230 ontology according to predicted reasoning time of all component reasoners.

231 R_2O_2 is different from PR in the following way. Instead of choosing the best
232 reasoner according to predicted reasoning time of reasoners, as in PR, R_2O_2 selects a
233 best possible reasoner from an aggregation of the *rankings* of component reasoners.

234 Recently a multi-criteria meta-reasoner Multi-RakSOR [49, 50] has been proposed
235 for reasoning about OWL 2 DL and EL ontologies. Multi-RakSOR incorporates two
236 objectives in selecting the best reasoner: reasoning efficiency and robustness. Efficiency
237 is measured by execution time. The robustness of a reasoner is measured by four ordered
238 *termination states*: (1) success (\mathcal{B}_S); (2) unexpected (\mathcal{B}_U), where the reasoning result is
239 not expected; (3) timeout (\mathcal{B}_T), where the reasoner times out on an ontology; and (4)
240 halt (\mathcal{B}_H), where the reasoner crashes. The ordering on these states is then defined to be
241 $\mathcal{B}_S < \mathcal{B}_U < \mathcal{B}_T < \mathcal{B}_H$.

242 Multi-RakSOR encompasses two main components: (1) a multi-label classifier that
243 predicts the termination state of a reasoner, and (2) a multi-target regression model that
244 predicts the ranking of the reasoners (with tie breaking) that respects the above ordering.

245 A comprehensive evaluation of Multi-RakSOR is performed on the ORE 2015
246 reasoner competition ontology dataset³, which we also use for evaluating $R_2O_2^*$, and a
247 set of 10 OWL 2 DL/EL reasoners. The paper also describes an “upgraded” version of
248 Multi-RakSOR, dubbed Meta-RakSOR, that is able to handle both OWL 2 DL (more
249 expressive) and OWL 2 EL (less expressive) ontologies. Meta-RakSOR is evaluated on
250 the task of ontology classification on two datasets: one for OWL 2 DL and for OWL
251 2 EL. The main evaluation results show that for both datasets, Meta-RakSOR has the
252 highest number of ontologies successfully reasoned over. In terms of efficiency, for
253 OWL 2 DL, Meta-RakSOR demonstrates competitive performance (but not better) than
254 the best single reasoner Konclude. For OWL 2 EL, it is shown that Meta-RakSOR ranks
255 6th of the eleven reasoners evaluated, in terms of average reasoning time.

256 The meta-reasoners/hybrid reasoners described so far are all focussed on TBox
257 reasoning (consistency checking or classification). PAGOdA [51] is a hybrid system

³<https://zenodo.org/record/50737>

258 designed for the task of *query answering* over ABox data. Employing an approach
259 similar to WSReasoner [48], PAGOdA uses an efficient reasoner, in this case a Datalog
260 reasoner, to compute a lower bound answer (sound but possibly incomplete) and an
261 upper bound answer (complete but possibly unsound). When the two answers do not
262 completely coincide, PAGOdA extracts relevant subsets from the TBox and the ABox
263 are extracted, which are used to verify the answers by a fully-fledged OWL 2 DL
264 reasoner.

265 **4. Ontology Metrics**

266 Metrics have been proposed to quantitatively measure the quality, complexity, testa-
267 bility, and maintainability of ontologies. Inspired by software metrics [52], we proposed
268 a set of *91 metrics* [24–26] for characterising the *design complexity* of ontologies. How-
269 ever, the definition of many of these metrics were not formally given. In this section,
270 we give a detailed account of this suite of 91 metrics that comprehensively characterise
271 ontologies in terms of their size and syntactic and structural complexity. These metrics
272 serve as distinctive features for learning ontology reasoning prediction models and
273 building the proposed meta-reasoning framework $R_2O_2^*$.

274 These metrics are organised by what they characterise: (1) *the ontology itself*, (2)
275 *classes*, (3) *anonymous class expressions*, and (4) *properties*. These metrics are proposed
276 with efficient computation as a key consideration. In the calculation of metrics, we
277 adopt a graph-based view of ontologies [24] to capture the complexity of ontologies and
278 generate a set of metrics.

279 The subsequent subsections present details of the metrics. Note that a metric name
280 without the percent sign (%) is a *count*, and one with it is a *ratio*. A count metric
281 shows how a component of an ontology has impact on the reasoning performance
282 by its occurrences. A *ratio metric* is used to explain the relationship between such a
283 component with respect to the overall structure of the ontology. Intuitively, a count/ratio
284 metric represents the absolute/relative value of a metric of an ontology, respectively. We
285 note that all metrics are computed on the *asserted* ontology hierarchy. In other words,
286 no reasoning is performed prior to computing these ontologies.

287 *4.1. Ontology-level Metrics (ONT)*

288 The 6 ONT metrics were previously defined [24]. Here, we define an additional 18
289 ONT metrics that have not been described previously. They measure the overall size
290 and complexity of an ontology.

291 **IND** counts the number of (named or anonymous) *individuals* in an ontology. The
292 remaining 17 metrics are collected by observing the structure (i.e. language constructs)
293 of a given ontology.

294 **GCI/HGCI**: These metrics measure the number of general concept inclusion (GCI)
295 axioms and hidden GCI (HGCI) axioms, respectively. GCI counts the number of sub-
296 sumption axioms whose subclass is a complex concept (anonymous class expression).
297 HGCI counts the number of (named) concepts that appear as a subclass in some sub-
298 sumption axioms as well as in some equivalent classes axioms. In general, the presence
299 of GCI axioms may increase reasoning complexity as they may introduce nondetermin-
300 ism [53]. A GCI axiom is *hidden* when a named class is the LHS of a subclass axiom as
301 well as an equivalent class axiom.

302 Either an equivalent class axiom or a subclass axiom where the left-hand side of the
303 subclass axiom is a named class.

304 **ESUB%/DSUB%/CSUB%**: These metrics measure ratios of subclass axioms that
305 contain (possibly nested) specific types of class expressions, including those that are
306 nested, and all subclass axioms. For these metrics, the subclasses and super classes are
307 *flattened*, and an axiom is considered to contain a specific type of expressions iff one of
308 the flattened expressions is of that type.

309 Respectively, **ESUB%**, **DSUB%** and **CSUB%** calculate the ratio of subclass axioms
310 that contain existential restrictions ($\exists R. _$), disjunctions ($_ \sqcup _$), and conjunctions
311 ($_ \sqcap _$). Additionally **CSUB%** requires that at least one of the conjuncts in the
312 subclass is an anonymous class expression. Existential restrictions (**ESUB%**) and
313 disjunctions (**DSUB%**) could generate AND-branching and OR-branching, respectively,
314 during the reasoning process. AND- and OR-branching are major sources of complexity
315 for tableau-based algorithms [54], hence their presence may negatively correlate with
316 performance. For the **CSUB%** metric, anonymous conjuncts in the subclass can generate

317 more axioms during the normalization process, hence it may increase workload for a
318 reasoner.

319 **ELCLS%/ELAX%**: These two metrics measure the ratios of (nested) class ex-
320 pressions (ELCLS%) and axioms (ELAX%), respectively, in the OWL 2 EL profile,
321 a sublanguage of OWL that is based on \mathcal{EL} [55], a description logic with efficient
322 PTIME-complete algorithms. Our intuition is that as reasoning in the EL profile is in
323 general easier, a reasoner such as MORE [56] that is able to delegate EL reasoning to an
324 efficient EL reasoner could be more efficient with ontologies with a large percentage of
325 EL expressions and axioms.

326 **HLC/HLC%**: These metrics are the count (HLC) and ratio (HLC%) of *hard*
327 language constructs, which include disjunctions, transitive properties, inverse properties
328 and property chains, in superclass expressions. These language constructs can potentially
329 negatively influence the performance of ontology reasoners.

330 **SUBCECHN/SUPCECHN**: These metrics are the count of top-level class expres-
331 sions containing *chained* (a sequence of) existential restrictions (i.e., $\exists R.C$ in the DL
332 syntax) as a subclass (SUBCECHN) or a superclass (SUPCECHN). These metrics
333 measure the impact of $\exists R.C$ expressions on the performance of reasoning as they can
334 potentially slow down the reasoning process by increasing the search space.

335 For example, suppose an ontology contains two subsumption axioms: (1) $\exists R.(A \sqcap$
336 $\exists R.(B \sqcap \exists P.C)) \sqsubseteq E$, and (2) $\exists R.(D \sqcap \exists R.(F \sqcap \exists P.G)) \sqsubseteq H$, where R, P represent
337 properties and A, \dots, H represent classes. For this ontology, SUBCECHN = 2 because
338 axiom (1)'s subclass is a chained class expression containing existential restriction
339 $\exists R.(A \sqcap \exists R.(B \sqcap \exists P.C))$ and axiom (2)'s subclass is also a chained existential restriction
340 $\exists R.(D \sqcap \exists R.(F \sqcap \exists P.G))$. On the other hand, there is no chained class expressions
341 containing existential restrictions as the superclass hence SUPCECHN = 0.

342 **DSUBCECHN/DSUPCECHN**: These count metrics calculate, in a depth-first man-
343 ner, the maximum depth of nested class expressions containing existential restrictions
344 as a subclass (DSUBCECHN) or a superclass (DSUPCECHN), with the intuition that
345 deeper subclass chains may increase reasoning time. For example, suppose an ontology
346 contains $\exists R.(A \sqcap \exists R.(B \sqcap \exists P.C)) \sqsubseteq E, D \sqcap \exists R.(F \sqcap \exists P.G) \sqsubseteq H$, where R, P represent
347 properties and A, B, C, D, E, F, G, H represent classes. DSUBCECHN is 3 because the

348 depth of nested class expressions containing existential restrictions in the first axiom
 349 $\exists R.(A \sqcap \exists R.(B \sqcap \exists P.C)) \sqsubseteq E$ is 3 and that of the second axiom $D \sqcap \exists R.(F \sqcap \exists P.G) \sqsubseteq H$
 350 is 2. In addition, there is no nested class expression containing existential restrictions in
 351 the superclass, hence DSUPCECHN is 0.

352 **SUBCCHN/SUPCCHN:** These metrics represent the number of class expressions
 353 containing *chained* conjunction expression as a subclass/superclass. For tableau-based
 354 algorithms, conjunctions of complex concepts in a subclass may not be easily normalised
 355 for some reasoners. Hence a subclass expression containing many complex class
 356 expressions may slow down the reasoning process.

357 **DSUBCCHN/DSUPDCHN:** These metrics represent maximum depth of nesting
 358 of class expressions containing disjunction expressions as a subclass/superclass. The
 359 idea of these metrics is similar to the metrics DSUBCECHN/DSUPCECHN.

360 In total there are 24 ONT metrics, which are summarised in Table A.12 in the
 361 appendix.

362 4.2. Class-level Metrics (CLS)

363 The CLS metrics capture characteristics of classes, which are first-class citizens in
 364 OWL ontologies. Five functions, NOC (number of children), NOP (number of parents),
 365 DIT (depth of inheritance tree), CID (class in-degree), and COD (class out-degree),
 366 were defined previously [24]. Each of these functions returns, for a (named or possibly
 367 nested anonymous) class expression in an ontology, a count value respectively. For a
 368 given class C , $\text{NOC}(C)$ and $\text{NOP}(C)$ return the number of direct subclasses and super
 369 classes of C in the ontology, respectively. $\text{DIT}(C)$ returns the longest path from C to \top ,
 370 the root class, in a depth-first manner. $\text{CID}(C)$ and $\text{COD}(C)$ calculate, respectively, the
 371 number of incoming and outgoing edges of C .

372 For each of these five functions, we identify three metrics: the *total*, the *average*,
 373 and the *maximum* values across all classes for a given ontology. For example for NOC,
 374 the total NOC (tNOC) is calculated by summing the NOC value for all classes in
 375 an ontology, and the average NOC (aNOC) is tNOC divided by the total number of
 376 class expressions. Similarly, the maximum NOC (mNOC) is the maximum number of
 377 children among all classes.

378 Thus, in total, 15 CLS metrics are identified, which are shown in Table A.13 in the
379 appendix.

380 4.3. Anonymous Class Expression Metrics (ACE)

381 The ACE metrics are an important ingredient in building expressive classes. Dif-
382 ferent types of anonymous class expressions can have different impact on reasoning
383 performance. The 9 ACE count metrics have been previously defined [25], one for
384 each different type of (possibly nested) anonymous class expressions (enumerations,
385 negations, conjunctions, disjunctions, universal restrictions, existential restrictions, and
386 min/max/exact cardinality restrictions). We further define two additional ACE count
387 metrics that represent the number of *value restrictions* (VALUE, for $\exists R.\{a\}$, where a
388 is an individual) and *self references* (SELF, for $\exists R.\text{self}$). We also propose their corre-
389 sponding *ratio* metrics that measure the percentage of each count ACE metric over all
390 (possibly nested) anonymous class expressions.

391 Hence in total there are 22 ACE metrics, shown in Table A.14 in the appendix.

392 4.4. Property Metrics (PRO)

393 Additional pairs of count and ratio metrics are defined: ASYM/ASYM% (asymmet-
394 ric properties), REFLE/REFLE% (reflective properties), IRREF/IRREF% (irreflective
395 properties), and CHN/CHN% (property chains).

396 Similarly, the 6 of the 8 existing count PRO metrics [25] are augmented with
397 their corresponding ratio metrics. For example, the DTP metric counts the number of
398 datatype properties. The metric DTP% records the ratio between the number of datatype
399 properties and the total number of properties. These 6 are: OBP (object properties),
400 DTP (datatype properties), FUN (functional properties), SYM (symmetric properties),
401 TRN (transitive properties), and IFUN (inverse functional properties). Furthermore,
402 four count metrics are defined to record the number of some property axioms, including
403 SUBP (subproperties), DISP (disjoint properties), DOMN (domain), and RANG (range).

404 Finally, four additional metrics are defined to measure the *usage* of properties in an
405 ontology.

- 406 • **ELPROP%**: This metric measures the ratio, of all property axioms, the number of
 407 property axioms allowed in the OWL 2 EL profile, which include subproperty axioms,
 408 equivalent property axioms, transitive axioms, reflexive axioms, domain/range axioms,
 409 and functional data property axioms. Intuitively, the higher the ELPROP% of an
 410 ontology is, the more efficient its reasoning may be.
- 411 • **IHR, IIR, ITR**: These metrics measure the count of class axioms (e.g., subclass
 412 axioms and class/property assertions) that make use of some property in some property
 413 hierarchy (IHR), inverse properties (IIR), and transitive properties (ITR). The intuition
 414 is that the more these types of properties are used in class axioms, the more difficult
 415 reasoning may be for this ontology.

416 There are in total 30 PRO properties as summarised in Table A.15.

417 5. Meta-Reasoning Models

418 In this section, as our major contributions of this paper, we propose our meta-
 419 reasoning framework $R_2O_2^*$ and its four different *meta-reasoning models* (simply
 420 *meta-reasoners*). Each meta-reasoner recommends the most efficient reasoner for
 421 unknown ontologies using different machine learning techniques. We first introduce
 422 basic notations we use in the paper. Then, we present the details of the four meta-
 423 reasoners with their learning objectives in the training phase and their utilisation in the
 424 recommendation (or testing) phase.

425 5.1. Notation Definition

426 The following basic notations are used in the paper.

- 427 • Let $R = \{r_1, \dots, r_n\}$ be a set of n reasoners, also called *component reasoners* in the
 428 paper.
- 429 • Let $\hat{R} = \{\hat{r}_1, \dots, \hat{r}_n\}$ be a set of n *prediction models* such that \hat{r}_i predicts the reasoning
 430 time of r_i .
- 431 • Let $OM = \{om_1, \dots, om_q\}$ be a set of q ontology metrics.
- 432 • Let $O = \{o_1, \dots, o_h\}$ be a set of h ontologies that can be reasoned about by at least one
 433 reasoner in R without timing out or errors. Each ontology in O is represented using
 434 its values of ontology metrics OM in this paper.

- 435 • Given a reasoner r and an ontology o , let $\theta(r, o)$ represent the actual reasoning time
- 436 of r for o for the task of ontology classification. Similarly, let $\theta(\hat{r}, o)$ represent the
- 437 reasoning time predicted by \hat{r} for o for the same task.
- 438 • Two partitioned subsets O_{tr} and O_{te} are drawn from O for training and testing the
- 439 proposed meta-reasoners, respectively.

440 5.2. Details of the Meta-Reasoners

441 In the following, we present the details of each of our four meta-reasoners. For

442 each reasoner, we describe its learning objective and how to build it in detail. All these

443 meta-reasoners are built on the training dataset $O_{tr} \subset O$.

444 The first meta-reasoner, $R_2O_2^*_{(pt)}$, directly trains prediction models \hat{R} of R on the

445 training data O_{tr} , and uses the predicted reasoning time that has been estimated by \hat{R}

446 to find the most efficient reasoners for unknown ontologies. The underlying idea is to

447 choose a reasoner r , whose predicted reasoning time estimated by \hat{r} is the most efficient

448 among \hat{R} , as the most efficient reasoner for an unknown ontology. The second meta-

449 reasoner, $R_2O_2^*_{(rk)}$, trains a ranking algorithm to learn the rankings of the reasoners

450 in R according to the actual reasoning time of the training data O_{tr} , and uses it to

451 predict the best ranked reasoners for unknown ontologies. The third meta-reasoner,

452 $R_2O_2^*_{(mc)}$, trains a classifier that learns the most efficient reasoner on O_{tr} , and uses

453 it to directly predict the most efficient reasoners for unknown ontologies. The fourth

454 meta-reasoner, $R_2O_2^*_{(all)}$, is an ensemble classifier that uses the predictions of the above

455 three meta-reasoners.

456 5.2.1. Meta-reasoner based on the direct use of predicted reasoning time: $R_2O_2^*_{(pt)}$

457 The meta-reasoner $R_2O_2^*_{(pt)}$ aims to recommend the most efficient reasoners for

458 unknown ontologies based on the *direct use of the predicted reasoning time* of all

459 reasoners in R . Therefore, for all reasoners in R , building their corresponding prediction

460 models \hat{R} is essential prior to making use of $R_2O_2^*_{(pt)}$ for determining such most efficient

461 reasoners. $R_2O_2^*_{(pt)}$ is similar to the non-ranking portfolio reasoner PR described in our

462 preliminary work [27] in that it leverages predicted reasoning time of \hat{R} . The difference

463 is that $R_2O_2^*_{(pt)}$ uses an ensemble regression model instead of using a random forest

464 regression algorithm as PR [27]. Also, $R_2O_2^*_{(pt)}$ incorporates the average rankings of
 465 the reasoners on the training data when there are more than two reasoners that were
 466 chosen as the most efficient (i.e. their predicted reasoning time is the same), while PR
 467 chooses one in random. The details of resolving a tie-breaking method $R_2O_2^*_{(pt)}$ is
 468 explained below.

469 Consequently, the effectiveness of $R_2O_2^*_{(pt)}$ relies mainly on the accuracy of the
 470 prediction models in \hat{R} . In order to build such prediction models, we use *stacking* [57],
 471 an ensemble learning technique to combine multiple classification (or regression) models
 472 in which (1) *base learners (or regression models)* (or level-0 models) are trained on the
 473 training data O_{tr} , and (2) the outputs of the base learners are combined using a *meta-*
 474 *learner (or meta-regression models)* (level-1 model), in our context. More specifically,
 475 for each reasoner, each learner is trained to learn a mapping function from the values
 476 of ontology metrics on the training data O_{tr} to their actual reasoning time. Then, a
 477 meta-learner is trained to learn a mapping function from the predicted outputs of O_{tr} ,
 478 which have been estimated by the base learners, to actual reasoning time.

479 Here, our aim is to use the decisions of the individual base learners that employ
 480 different learning criteria, and to combine their decisions to outperform each individual
 481 base learner using a meta-learner.

More formally, to build each prediction model $\hat{r}_k \in \hat{R}$ for reasoner $r_k \in R$, we
 represent each ontology $o_i \in O_{tr}$ as follows:

$$o_i = \underbrace{om_{i,1}, \dots, om_{i,q}}_{\text{ontology metrics}}, \underbrace{\theta(r_k, o_i)}_{\text{actual reasoning time}} \quad (5)$$

482 where $om_{i,j}$ is the value of the j -th ontology metric om_j of ontology o_i , and $\theta(r_k, o_i)$
 483 denotes the actual reasoning time of r_k on ontology o_i .

484 Using the above representation scheme, for each reasoner, we train k base learners
 485 (level-0 models) on O_{tr} . Then, we generate level-1 data obtained from the predictions
 486 of the k base learner over the instances in O_{tr} . The level-1 data have k attributes whose
 487 values are the predictions (i.e. predicted time) of the k base learners for every instance
 488 in O_{tr} . Thus, each training example for a meta-learner (level-1 model) will be composed
 489 of k attributes (e.g. k predictions from the k base learners) and the target which is the

490 actual reasoning time for every instance in O_{tr} . Once the level-1 data have been built
491 from all instances in O_{tr} , any learning regression models can be used to generate the
492 meta-learner. In this paper we choose $k = 2$.

493 In this paper, we use two robust base learners: (1) the *random forest regression*
494 *algorithm* and (2) the *XGBoost (eXtreme Gradient Boosting)* algorithm [28]:

- 495 • **Random forest regression algorithm:** As a base learner, we build a regression
496 model using *random forest regression algorithm*, an efficient and robust learning
497 model, which has produced good predictive performance in our previous work [26].
498 In our context, the random forest model combines a number of *decision trees*, each of
499 which is trained using a subset of training ontologies, to build a prediction model for
500 a given reasoner.
- 501 • **XGBoost:** As another base learner, we use the state-of-the-art learning algorithm,
502 XGBoost [28], which has recently shown dominant performance on a number of
503 Kaggle competitions for structured or tabular data. It is an implementation of gradient
504 boosted decision trees designed for achieving better computational efficiency and
505 prediction performance.

506 We again consider random forest and XGBoost as candidates of a meta-learner
507 (level-1 model) due to their strong predictive performance. These level-1 candidates
508 are denoted by meta-RF and meta-XGBoost in this paper. As can be seen in Sec-
509 tion 7.1.1, we eventually choose meta-XGBoost as our meta-learner given its best
510 overall performance.

511 Once we train the meta-regression model on O_{tr} , given an unknown ontology in
512 the testing data O_{te} , the meta-reasoner $R_2O_2^*_{(pt)}$ will recommend a reasoner whose
513 predicted reasoning time is the fastest among all of the predictions of the prediction
514 models in \hat{R} as the most efficient reasoner.

515 If more than two reasoners are chosen as the most efficient, a tie-breaking method is
516 also applied to select one of them. This method takes into consideration the *precision*
517 *at 1 (P@1)* of the reasoners in R that are measured on the training data O_{tr} . In this
518 context, for each reasoner in R , its P@1 is measured by the proportion that the reasoner
519 is the most efficient across all instances in O_{tr} . Our tie-breaking method chooses the
520 reasoner with the highest P@1. This tie-breaking method is applied to all the other

521 meta-reasoners in our $R_2O_2^*$.

522 5.2.2. Meta-reasoner based on ranking algorithm: $R_2O_2^{*(rk)}$

523 $R_2O_2^{*(rk)}$ is a meta-reasoner that learns the rankings of the reasoners in R . During
 524 the training phase, it trains a *ranking algorithm* (simply *ranker*) that learns the rankings
 525 of the reasoners on O_{tr} in terms of their reasoning time. Once the ranker is trained,
 526 given an unknown ontology in O_{te} , $R_2O_2^{*(rk)}$ ranks and recommends the most efficient
 527 reasoner for that ontology.

528 $R_2O_2^{*(rk)}$ follows a similar spirit of R_2O_2 [27] in that $R_2O_2^{*(rk)}$ uses a *ranking*
 529 *matrix* and uses a ranker. The difference is that $R_2O_2^{*(rk)}$ uses a single ranker, rather
 530 than aggregating multiple rankers as R_2O_2 , to recommend the most efficient reasoner
 531 for an unknown ontology. Given a pool of rankers using different criteria for learning,
 532 we have chosen the one with the best ranking performance through our experiments
 533 which will be further discussed in Section 7.1.1.

Given the training data O_{tr} , we generate a *ranking matrix*, where each row represents
 the values of ontology metrics and the rankings of reasoners R according to their actual
 reasoning time. Then, a ranker is trained on this matrix to learn how the characteristics
 of the ontologies in O_{tr} can be optimally mapped to the relative ordering of the reasoning
 performance of the reasoners in R . Initially, we build an $|O_{tr}| \times (q + n)$ data matrix \mathbf{M}_d
 (recall that $q = |OM|$, $n = |R| = |\hat{R}|$), where row i represents an ontology $o_i \in O_{tr}$ and
 the actual reasoning time of the reasoners in R for o_i :

$$o_i = \underbrace{om_{i,1}, \dots, om_{i,q}}_{\text{ontology metrics}}, \underbrace{\theta(r_1, o_i), \dots, \theta(r_n, o_i)}_{\text{actual reasoning time}}, \quad (6)$$

where $om_{i,j}$ is the value of the j -th ontology metric om_j of o_i , and $\theta(r_s, o_i)$ denotes r_s 's
 actual reasoning time for o_i . From \mathbf{M}_d , we build the corresponding $|O_{tr}| \times (q + n)$ ranking
 matrix \mathbf{M}_r , where row i is represented as:

$$o_i = \underbrace{om_{i,1}, \dots, om_{i,q}}_{\text{ontology metrics}}, \underbrace{\pi(r_1, o_i), \dots, \pi(r_n, o_i)}_{\text{ranking of reasoners}}, \quad (7)$$

534 where $\pi(r_s, o_i)$ denotes the *rank* of $r_{s[1..n]} \in R$ for o_i determined by $\theta(r_s, o_i)$. On \mathbf{M}_r , the
 535 more efficient a reasoner is, the higher ranked it is (the smaller the rank number). To

536 illustrate this, suppose there are 3 reasoners $\{r_1, r_2, r_3\}$, and their actual reasoning time
 537 for an ontology o_i is 100s, 90s, and 10s, respectively, i.e., $(\theta(r_1, o_i), \theta(r_2, o_i), \theta(r_3, o_i))$
 538 $= (100s, 90s, 10s)$. Thus, the ranking is $(\pi(r_1, o_i), \pi(r_2, o_i), \pi(r_3, o_i)) = (3, 2, 1)$. If the
 539 reasoning time is (10s, 10s, 100s) instead, the ranking will be (1, 1, 3).

540 In summary, the goal is to learn rankings of all reasoners in R on the ranking matrix
 541 \mathbf{M}_r , and to predict a ranking of the reasoners for an unseen ontology. The top-ranked
 542 reasoner will be chosen by the meta-reasoner $R_2O_2^*_{(rk)}$ to be the most efficient reasoner
 543 to reason over the ontology. Comparing to $R_2O_2^*_{(pt)}$, the main feature of $R_2O_2^*_{(rk)}$
 544 stems from that it is built on the ranking matrix that uses rankings of the reasoners in R ,
 545 rather than the direct use of the predicted reasoning time estimated from \hat{R} .

546 5.2.3. Meta-reasoner based on multi-class classification: $R_2O_2^*_{(mc)}$

547 $R_2O_2^*_{(mc)}$ formulates the learning problem into a multi-class classification problem,
 548 where its goal is to classify an ontology with one of the reasoners that is able to reason
 549 about the ontology the most efficiently. During the training phase, $R_2O_2^*_{(mc)}$ learns the
 550 most efficient reasoner for each ontology on the training data O_{tr} . The most efficient
 551 reasoner is determined by means of the actual reasoning time of the reasoners in R ,
 552 meaning that the fastest reasoner is chosen as the most efficient reasoner.

More formally, to build $R_2O_2^*_{(mc)}$, we represent each ontology $o_i \in O_{tr}$ as follows:

$$o_i = \underbrace{om_{i,1}, \dots, om_{i,q}}_{\text{ontology metrics}}, \quad \underbrace{\gamma_i}_{\text{the most efficient reasoner}} \quad (8)$$

553 where $om_{i,j}$ is the value of the j -th ontology metric om_j of o_i , and γ_i denotes the actually
 554 most efficient reasoner for o_i . If there is an ontology $o_i \in O_{tr}$ that has k -reasoners
 555 (where $k > 1$) that show the equivalently most efficient reasoning time, we generate
 556 k instances with k most efficient reasoners for o_i . For example, given on ontology o_i ,
 557 suppose that there are two most efficient reasoners: Konclude and Pellet. We then
 558 generate two instances for o_i as follows: (1): " $om_{i,1}, \dots, om_{i,q}, \text{Konclude}$ ", and (2)
 559 " $om_{i,1}, \dots, om_{i,q}, \text{Pellet}$ ".

560 Using the above representation scheme, we train a classifier on the training data O_{tr} .
 561 In our experiments, we have considered two classifiers: (1) *random forest algorithm*
 562 and (2) *XGBoost*, because of their robust classification performance as in the case of the

563 meta-reasoner $R_2O_2^*_{(pt)}$. Based on our cross-validation on O_{tr} , we eventually choose
 564 XGBoost which will be further discussed in Section 7.1.1.

565 In comparison with the meta-reasoner $R_2O_2^*_{(pt)}$, $R_2O_2^*_{(mc)}$ does not directly use
 566 the predicted reasoning time of the reasoners in R , that is, it does not rely on \hat{R} . Rather
 567 it learns which reasoners have been the most efficient reasoners for ontologies in
 568 the training data O_{tr} . The learning goal of $R_2O_2^*_{(mc)}$ is similar to the meta-reasoner
 569 $R_2O_2^*_{(rk)}$ in that its learning is based on the ranks of the reasoners in R on the training
 570 data O_{tr} . However, $R_2O_2^*_{(mc)}$ differs in that it learns the most efficient reasoner only,
 571 not the rankings of all reasoners in R as $R_2O_2^*_{(rk)}$.

572 5.2.4. Ensemble meta-reasoner: $R_2O_2^*_{(all)}$

573 $R_2O_2^*_{(all)}$ is a stacking classifier that learns from the predictions of the above three
 574 meta-reasoners: (1) $R_2O_2^*_{(pt)}$ that directly uses the predicted reasoning time of reasoners
 575 for the training set O_{tr} , (2) $R_2O_2^*_{(rk)}$ that learns the rankings of reasoners by means
 576 of their actual reasoning time for O_{tr} , and (3) $R_2O_2^*_{(mc)}$ that learns the most efficient
 577 reasoners for ontologies in O_{tr} . In other words, in $R_2O_2^*_{(all)}$, those three meta-reasoners
 578 can be seen as base classifiers (i.e. level-0 models), and $R_2O_2^*_{(all)}$ learns a meta-classifier
 579 (i.e. level-1 model) from the predictions of the base classifiers.

Thus, $R_2O_2^*_{(all)}$ trains a meta-classifier on the training data O_{tr} , where each ontology
 $o_i \in O_{tr}$ is represented as follows:

$$o_i = \underbrace{R_2O_2^*_{(pt)}(o_i), R_2O_2^*_{(rk)}(o_i), R_2O_2^*_{(mc)}(o_i)}_{\text{predicted reasoners of the meta-reasoners}}, \underbrace{\gamma_i}_{\text{the most efficient reasoner}} \quad (9)$$

580 where $R_2O_2^*_{(pt)}(o_i)$, $R_2O_2^*_{(rk)}(o_i)$ and $R_2O_2^*_{(mc)}(o_i)$ are the predicted most efficient
 581 reasoners of $R_2O_2^*_{(pt)}$, $R_2O_2^*_{(rk)}$, $R_2O_2^*_{(mc)}$, respectively, for o_i . As in Equation 8, γ_i
 582 denotes the actually most efficient reasoner for o_i .

583 As candidates for a meta-classifier, we also consider the same learning algorithms
 584 used to built the meta-reasoner $R_2O_2^*_{(pt)}$ because of the same reason: their proven,
 585 robust predictive performance: (1) random forest algorithm and (2) XGBoost. The one
 586 with the better performance from these two candidates is chosen, where the performance
 587 is measured via cross-validation on the training data O_{tr} . Eventually, XGBoost is chosen
 588 as our meta-classifier which will be presented in Section 7.1.1.

589 Thus, $R_2O_2^*_{(all)}$ trains the meta-classifier that learns relationships between the
 590 predicted reasoners of the previous three meta-reasoners and the actually most efficient
 591 reasoners on the training data O_{tr} , and then make predictions about the most efficient
 592 reasoners for unknown reasoners for the testing ontologies O_{te} .

593 5.2.5. Meta-reasoners with ELK

594 In the previous sections, we have presented the four different meta-reasoners in
 595 $R_2O_2^*$. Note that all of these meta-reasoners utilise a number of high-performance DL
 596 reasoners. Here, an interesting question is whether incorporating an EL reasoner into
 597 the meta-reasoners, such as ELK [30], can further improve reasoning efficiency for the
 598 less complex OWL 2 EL ontologies.

599 To address this question, we augment our meta-reasoners to incorporate an EL
 600 reasoner, ELK, that is designed to support the less expressive OWL 2 EL profile. ELK
 601 does not support reasoning over non-EL ontologies. Thus, the meta-reasoners only
 602 incorporate ELK when predicting the most efficient reasoners for unknown OWL 2 EL
 603 ontologies.

604 First, we train a prediction model for ELK on the training data O_{tr} following the
 605 steps presented in Section 5.2.1. Then, given an ontology $o \in O_{te}$ whose reasoning time
 606 is unknown, each meta-reasoner (representatively denoted by $R_2O_2^*$) is extended by
 607 taking the following further steps to find the most efficient reasoner to reason about o :

- 608 1. Check whether o is an EL ontology or not.
2. If o is an EL ontology, compare the predicted reasoning time and the most
 efficient reasoner determined by $R_2O_2^*$. Then, we choose the fastest one as the
 most efficient reasoner. Formally,

$$\hat{\gamma}(o) = \underset{\hat{r}_i \in R_2O_2^*, \hat{r}_{elk}}{\operatorname{argmin}} \theta(\hat{r}_i, o) \quad (10)$$

609 where $\hat{\gamma}(o)$ is the most efficient reasoner for o determined eventually; $\theta(\hat{r}_i, o)$ is
 610 the predicted reasoning time of the reasoner r_i for o ; and \hat{r}_{elk} is the prediction
 611 model of ELK.

612 3. If o is not an EL ontology, follow the recommendation of $R_2O_2^*$ not considering
613 r_{elk} .

614 Thus, in our extension for EL ontologies, for each meta-reasoner, the main idea
615 is to simply incorporate predicted reasoning time of ELK and the most efficient rea-
616 soner determined by the meta-reasoner, to compare their reasoning time, and finally to
617 recommend the one with the most efficient predicted reasoning time.

618 Here we summarise the main differences between our meta-reasoning framework
619 $R_2O_2^*$ and our preliminary meta-reasoner R_2O_2 and the portfolio-based meta-reasoner
620 PR, which were described in detail in Section 3.

- 621 1. As a framework (but not individual meta-reasoners), $R_2O_2^*$ adapts and incorpo-
622 rates both R_2O_2 (as $R_2O_2^*_{(rk)}$) and PR (as $R_2O_2^*_{(pt)}$).
- 623 2. $R_2O_2^*$ incorporates more advanced prediction models (Random Forests and
624 XGBoost) through *ensembling*.
- 625 3. $R_2O_2^*$ generates a ranking matrix using the actual reasoning time of the compo-
626 nent reasoners. On the other hand, R_2O_2 built a ranking matrix using the predicted
627 reasoning time of the reasoners, where such time was estimated by prediction
628 models.
- 629 4. $R_2O_2^*$ incorporates a single best ranker instead of using an aggregation of multiple
630 rankers as R_2O_2 . The single best ranker is determined empirically through cross-
631 validation on the training data.
- 632 5. $R_2O_2^*$ invokes the efficient reasoner ELK [30] directly for OWL 2 EL ontologies.
- 633 6. $R_2O_2^*$ is built using a different mix of component reasoners that includes Pellet
634 (which R_2O_2 does not include) but excludes TrOWL, as TrOWL is an approximate,
635 therefore incomplete reasoner.
- 636 7. Finally, $R_2O_2^*$ is built and evaluated using a more modern set of ontologies and
637 more recent versions of reasoners.

638 In the following sections, we discuss our evaluation framework and results that
639 compare the performance of all the proposed meta-reasoners (with and without ELK)
640 using the ORE 2015 competition corpus [22].

641 **6. Evaluation Framework**

642 In this section, we describe the evaluation framework used for evaluating the pro-
643 posed prediction models and meta-reasoners, including details of the reasoners, ontolo-
644 gies and the evaluation environment. The notations used in this section follow those
645 defined in Section 5.1.

646 **Reasoners:** Six state-of-the art OWL 2 DL reasoners that participated in ORE the
647 2015 reasoner competition [22] are used as component reasoners (simply reasoners)⁴:
648 FaCT++ [10], HermiT [34], JFact,⁵ Konclude [12], MORE [56] (with HermiT as the
649 underlying OWL 2 DL reasoner), and Pellet [13]. Besides these six OWL DL reasoners,
650 we also incorporate ELK [30], the efficient reasoner for the less expressive OWL EL
651 profile. The versions of the reasoners are the same as those in ORE 2015. As described
652 in Section 5, we build a prediction model for each reasoner, which is one of the key
653 components in the meta-reasoner framework $R_2O_2^*$.

654 **Target Reasoning Task:** For the ontology reasoning task, we choose ontology *clas-*
655 *sification*. The actual reasoning time (wall-time) of each ontology in the dataset was
656 measured on a high-performance server running CentOS Linux 7.4 (Core) and Java 1.8
657 on single-core Intel Gold 6140 each at 2.3 GHz, with a maximum of 10 GB memory
658 allocated to the reasoner. A timeout of *30 minutes* of wall-time is imposed on each
659 (reasoner, ontology) pair.

660 To ensure consistency of the evaluation across reasoners, the ORE 2015 competition
661 framework⁶ is used to invoke all reasoners and record their reasoning time. For each
662 ontology, the competition framework converts it into the OWL functional syntax (FSS),

⁴Chainsaw [44] and Racer [58] (two OWL 2 DL reasoners that participated in ORE 2015) are excluded due to reasoning errors in an excessive number of ontologies.

⁵<http://jfact.sourceforge.net>

⁶<https://github.com/andreas-steigmiller/ore-competition-framework>

663 invokes reasoners for classification using a Bash shell script, and records reasoning time
664 using the GNU `time` command. We follow the framework and include ontology loading
665 time as part of classification time.

666 **Ontologies:** We collected all 1,920 ontologies from the ORE 2015 reasoner compe-
667 tition [22]⁷. Prior to building our proposed meta-reasoners, we performed the three
668 preprocessing steps on the 1,920 ontologies, following the steps in [26, 27]. The aim
669 of these steps is to remove duplicate ontologies that may exist in the given ontology
670 collection, normalise their metric values to avoid the high skewness of values of ontol-
671 ogy metrics *OM*, and remove ontology metrics that may influence learning a prediction
672 model with lower prediction accuracy:

- 673 1. **Cleansing:** Since the given ontology collection has been obtained from multiple
674 repositories, it may contain duplicates. All but one ontology is removed from
675 each set of ontologies with duplicate metric values. After removing duplicates,
676 1,760 ontologies remained.
- 677 2. **Normalisation:** In the given 1,760 ontology collection, values of some of the
678 metrics span a large range and are very skewed as discussed in [26]. We apply a
679 commonly-used log-transformation on the metric values of the ontologies that are
680 greater than 10. The log-transformation is also performed on reasoning time.
- 681 3. **Metric removal:** It is a widely used practice to remove features that have near-zero
682 variance values and features that are highly correlated (with respect to the dataset).
683 In this paper we follow this practice. Following our previous work [26, 27], we
684 consider two metrics with correlation coefficients above 0.9 to be highly correlated.
685 To observe a better generalised distribution of these metrics, we measure their
686 correlation on a larger ontology collection, the one used in the ORE 2014 reasoner
687 competition [21]. It contains 16,555 ontologies, which are split into four groups
688 by percentiles of file size. Ontologies are randomly sampled from within these
689 groups. This is to ensure that files of different sizes are sufficiently represented.

⁷<http://owl.cs.manchester.ac.uk/publications/supporting-material/ore-2015-report/>

690 Given correlation calculated from this ontology collection, we remove all but
691 29 metrics: 12 ONT metrics (SOV, ENR, EOG, CYC, RCH, IND, ESUB%,
692 ELCLS%, ELAX%, HLC, HLC%, SUPCECHN); 7 CLS metrics (tNOC, aNOC,
693 aCID, mCID, tCOD, aCOD, aNOP); 4 ACE metrics (CONJ%, UF%, EF, EF%);
694 and 6 PRO metrics (OBP%, DTP%, FUN%, CHN%, ELPROP%, IHR).

695 Finally, as our ontology collection (O), we used 1,760 ontologies, where each
696 ontology is represented by the 29 metrics and the metric values are log-scaled.

697 Table 1 shows, for each reasoner, the total number of ontologies it successfully
698 handled, that result in an error, and that time out, and brief statistics of reasoning time
699 (in seconds, and excluding those ontologies that time out). It can be observed that
700 the reasoning time spans a large range for all the reasoners, and that Konclude is the
701 most efficient as well as most robust reasoner (the least number of error and timeout
702 ontologies). In total, 1,269 ontologies (excluding timeout ontologies) were reasoned
703 about successfully (no error, no timeout) by all the six reasoners, and 1,390 ontologies
704 did not result in a runtime error (timeout ontologies are included in this case). Note that
705 timeout ontologies are used to evaluate our meta-reasoners in one set of experiments.
706 Figure 1 in Section 7 depicts the performance characteristics of the component reasoners
707 in more details in violin plots. Of the 1,760 ontologies, 761 are in the OWL 2 EL profile.
708 Hence, we performed classification on these ontologies using ELK.

709 **Evaluation method:** We evaluate our meta-reasoners with the state-of-the-art algorithm
710 selection framework AutoFolio [31] that is configured to minimise runtime. As discussed
711 earlier in Section 5, meta-reasoners PR and R_2O_2 described in our previous work [27]
712 are similar to $R_2O_2^*_{(pt)}$ and $R_2O_2^*_{(rk)}$, respectively. Hence this comparison also assesses
713 the performance of $R_2O_2^*_{(all)}$ against our previous models. We also attempted to
714 evaluate our meta-reasoners against a recently proposed multi-criteria meta-reasoner
715 Meta-RaksOR [49, 50], which has dual optimisation objectives of reasoning correctness
716 as well as efficiency. However, we are unable to evaluate Meta-RaksOR due to two
717 reasons. Firstly, for the OWL 2 DL classification task, Meta-RaksOR incorporates
718 eight component reasoners but our meta-reasoners only incorporates six. Among the
719 two reasoners that are not considered by our meta-reasoners, TrOWL [14], which is

Table 1: A summary of statistics of the deduplicated ORE 2015 competition dataset containing a total of 1,760 unique ontologies. Note the reasoning time is measured in seconds and without considering timeout ontologies.

Reasoner	No. of Ontologies			Reasoning time (s)			
	Successful	Error	Timeout	Min	Max	Median	Mean
FaCT++	1,461	109	190	0.53	1,638.6	1.4	72.4
HermiT	1,658	51	51	0.70	1,622.4	3.6	35.8
JFact	1,292	161	307	1.03	1,788.6	3.4	72.8
Konclude	1,737	8	15	0.03	1,087.2	0.3	3.7
MORe	1,706	18	36	2.00	1,684.5	4.5	87.1
Pellet	1,477	114	169	1.01	1,773.9	3.5	39.5

720 an approximate hence incomplete reasoner, and RACER [59], which did not execute
721 properly on our evaluation hardware. Secondly, even though Meta-RakSOR has released
722 source for *running* the meta-reasoner⁸, it however does not include the source code
723 of Meta-RakSOR itself. Therefore we are unable to modify it and compare with it.
724 However, as can be seen from Table 1 of Meta-RakSOR [50], Meta-RakSOR does
725 not outperform Konclude on average reasoning time, it is thus not unreasonable to
726 hypothesise that our meta-reasoners would outperform Meta-RakSOR, as our meta-
727 reasoners outperform Konclude.

728 In our evaluation we retain timeout ontologies to realistically assess performance
729 of all reasoners. We assess the impact of those ontologies that result in a runtime
730 error in two different experiments. In the first experiment (hereinafter referred as
731 **ErrorsRemoved**), the error ontologies for each reasoner are removed. In the second
732 experiment (hereinafter referred as **ErrorsReplaced**), the error ontologies for each
733 reasoner are treated as they timeout. Note that in ErrorsRemoved, it may be the case
734 that a reasoner that strictly conforms to OWL semantics may throw many runtime errors

⁸<https://github.com/Alaya2016/Multi-RakSORDemo>

735 on a corpus of ontologies as it may reject non-conformant constructs (e.g. imaginary
736 numbers). As such, such a reasoner may turn out to be efficient in the experiment
737 **ErrorsRemoved** than in **ErrorsReplaced**.

738 In each experiment, standard 10-fold cross validation is performed to adequately
739 assess the performance of the meta-reasoners. That is, we take the following steps: (1)
740 shuffle the ontologies O randomly; (2) split O into 10 subsets; (3) for each subset, take
741 the subset (i.e. 10%) as test set (O_{te}), and take the remaining subsets (i.e. 90%) as a
742 training set (i.e. O_{tr}); (4) fit a model on the training set and evaluate it on the test set;
743 and (5) average the evaluation scores of the model across 10-times.

744 As evaluation metrics, average runtime (i.e. reasoning time) is used as the main
745 evaluation metric for the meta-reasoner. To evaluate our runtime prediction (regression)
746 models, we used the standard ‘coefficient of determination’ (R^2) as used in [26]. R^2
747 denotes the proportion of the variation in the target variable (i.e. reasoning time) that
748 can be explained by each prediction model $\hat{r} \in \hat{R}$. The higher the R^2 value is, the more
749 accurate the model is. Moreover, we evaluate our meta-reasoners using the standard
750 metric precision at 1 ($P@1$), as the meta-reasoners require the predicted best reasoner.
751 The performance number reported in the rest of the section is the average on the test set
752 over the 10 folds for each experiment.

753 In the experiment **ErrorsReplaced**, in each fold of the cross-validation, 90% of
754 1,760 ontologies ($\approx 1,584$) are used for training, and the rest 10% (≈ 176) are used for
755 testing. In the experiment **ErrorsRemoved**, the 1,389 ontologies are randomly divided
756 into the training set ($\approx 1,250$) and the test set (≈ 139) in each fold.

757 **7. Evaluation Results and Analysis**

758 In this section we present the evaluation results and their detailed analysis. The
759 evaluation is conducted in two parts. In Section 7.1, we present the performance
760 evaluation of the prediction models used in $R_2O_2^*$ on the training set O_{tr} . In Section 7.2,
761 we present the overall evaluation results, obtained on the test set O_{te} , comparing $R_2O_2^*$
762 with component reasoners and AutoFolio.

763 *7.1. Performance Evaluation of the Key Learning Components in $R_2O_2^*$*

764 Here, we present the performance evaluation of the prediction models (regression
765 models, classifiers and rankers) used in our meta-reasoning framework $R_2O_2^*$, obtained
766 through 10-fold cross-validation on the training set O_{tr} .

767 *7.1.1. Performance of regression models in $R_2O_2^*_{(pt)}$*

768 Here, our goal is to assess the generalizability of the prediction models \hat{r}_i of rea-
769 soner $r_i \in R$. As described in Section 5.2.1, these prediction models are central to
770 $R_2O_2^*_{(pt)}$. Performance (i.e. generalizability) was measured in terms of the coefficient
771 of determination (R^2), as introduced in Section 6.

772 As presented in Section 5.2.1, we use a stacking approach to build a prediction
773 model with two base regression models (level-0 models): random forest regression
774 algorithm and XGBoost [28]. As a meta-regression model (level-1 model), we also used
775 these two algorithms.

776 We employed the widely-used Weka framework⁹ as our evaluation environment.
777 For ease of experimentation we chose Weka’s version of the random forest algorithm.
778 For XGboost, we used Weka-XGBoost¹⁰ that can easily interface with Weka.

779 For the random forest algorithm, we use the default configuration in Weka. For
780 XGBoost, we set the following parameters keeping all the others fixed as default
781 throughout this paper: num_round = 50 (the number of rounds for boosting), eta =
782 0.1 (learning (or shrinkage) parameter that controls how much information from a new
783 tree will be used in the Boosting), max_depth = 10 (controls the maximum depth of
784 the trees: deeper trees have more terminal nodes and fit more data), sub_sample =
785 0.5 (determines if we are estimating a Boosting or a Stochastic Boosting. A value 1
786 represents the regular boosting, and a value between 0 and 1 is for the stochastic case.
787 The stochastic Boosting uses only a fraction of the data to grow each tree. For example,
788 if we use 0.5 each tree will sample 50% of the data to grow). Note that these parameter
789 values were chosen empirically. XGBoost uses multiple parameters and determining

⁹<https://www.cs.waikato.ac.nz/ml/weka/>.

¹⁰<https://github.com/SigDelta/weka-xgboost>.

790 optimal parameter values is beyond the scope of this paper.

791 Tables 2 and 3 show the R^2 values of the 7 prediction models obtained from cross-
 792 validation on the training data O_{tr} in the two experiments: **ErrorsRemoved** and **Er-**
 793 **rorsReplaced**. Note that we have compared the prediction performance of 4 regression
 794 models as candidates as a prediction model for each reasoner: (1) random forest regres-
 795 sion algorithm (denoted by RF), (2) XGBoost, (3) a stacking meta-regression model
 796 using random forest regression algorithm (denoted by meta-RF), and (4) a stacking
 797 meta-regression model using XGBoost (denoted by meta-XGBoost). R^2 denotes the pro-
 798 portion of the variation in the target variable (i.e. reasoning time) that can be explained
 799 by the model. For example, 0.853 in the model $\hat{r}_{\text{FaCT++}}$, implemented by meta-XGBoost,
 800 indicates that 85.3% of the variation in the reasoning time can be accounted for by
 801 meta-XGBoost. In Table 2, both RF and meta-XGBoost show the best prediction perfor-
 802 mance whereas in Table 3, meta-XGBoost has the highest R^2 value. Thus, we choose
 803 meta-XGBoost to implement the meta-reasoner $R_2O_2^*_{(pt)}$.

Table 2: A summary of prediction model performance as measured by R^2 on the dataset **ErrorsRemoved**.

Model	RF	XGBoost	meta-RF	meta-XGBoost
$\hat{r}_{\text{FaCT++}}$	0.852	0.852	0.852	0.853
\hat{r}_{HermiT}	0.825	0.824	0.831	0.833
\hat{r}_{JFact}	0.904	0.903	0.906	0.907
$\hat{r}_{\text{Konclude}}$	0.910	0.905	0.909	0.909
\hat{r}_{MORe}	0.723	0.705	0.707	0.708
\hat{r}_{Pellet}	0.770	0.763	0.765	0.766
\hat{r}_{Elk}	0.948	0.947	0.949	0.950
Mean	0.847	0.843	0.846	0.847

804 The above observations provide insight into how well the 7 prediction models can fit
 805 the given data. The similar values of R^2 with the ones in [26] (i.e. averaged $R^2 = 0.869$)
 806 suggest a good generalizability of the models.

Table 3: A summary of prediction model performance as measured by R^2 on the dataset **ErrorsReplaced**.

Model	RF	XGBoost	meta-RF	meta-XGBoost
$\hat{r}_{\text{FaCT}++}$	0.835	0.832	0.834	0.835
\hat{r}_{HermiT}	0.807	0.807	0.812	0.814
\hat{r}_{JFact}	0.843	0.833	0.839	0.840
$\hat{r}_{\text{Konclude}}$	0.909	0.909	0.912	0.913
\hat{r}_{MORe}	0.757	0.744	0.750	0.756
\hat{r}_{Pellet}	0.727	0.724	0.719	0.723
\hat{r}_{Elk}	0.939	0.936	0.939	0.940
Mean	0.831	0.826	0.829	0.832

807 *7.1.2. Performance of rankers in $R_2O_2^*_{(\text{rk})}$*

808 As explained in Section 5.2.2, our meta-reasoner, $R_2O_2^*_{(\text{rk})}$, incorporates a single
 809 ranker which differs from our previous approach in R_2O_2 [27] that uses an aggregation of
 810 multiple rankers. To implement $R_2O_2^*_{(\text{rk})}$, we initially considered the 5 rankers¹¹ that are
 811 used in R_2O_2 [27] on the ranking matrix \mathbf{M}_r : (1) *KNNRanker* that aggregates rankings
 812 using the nearest neighbors, (2) *PCTRanker* that is based on predictive clustering tree for
 813 ranking, (3) *RPCRanker* that is based on a multiple binary pairwise classifier to construct
 814 a ranking model, (4) *RegRanker* where the ranking problem is cast to a multi-target
 815 regression problem, where the rank position values of each reasoner are the targets in
 816 the multi-target regression setting, and (5) *ARFRanker* that is based on using random
 817 forests for ensembling multiple binary approximated ranking trees. Then, we choose
 818 the one showing the best performance in 10-fold cross-validation on the training data
 819 O_{tr} . For each ranker, its performance was measured by precision at 1 (denoted by P@1)
 820 that measures the proportion of the reasoners that are correctly recommended by the
 821 ranker as most efficient reasoner.

822 Table 4 shows the performance of the 5 rankers in terms of P@1 on both datasets,

¹¹These rankers are available in <http://www.quansun.com>, Readers interested in the details of the rankers are referred to [60].

823 **ErrorsRemoved** and **ErrorsReplaced**. For each ranker, P@1 was measured using
824 10-fold cross-validation from the ranking matrix generated from the training data O_{tr} .
825 The ranking matrix formation was presented in Equation 7 in Section 5.2.2. As seen in
826 the table, all 5 rankers showed high performance, achieving more than 90% of the P@1
827 values. ARFRanker shows the best performance (denoted in bold): 0.978 and 0.955 on
828 **ErrorsRemoved** and **ErrorsReplaced**, respectively. Consequently, to implement our
829 meta-reasoner $R_2O_2^*_{(rk)}$, we use ARFRanker.

Table 4: A summary of performance assessment of the 5 rankers in terms of P@1.

Ranker	ErrorsRemoved	ErrorsReplaced
KNNRanker	0.973	0.952
PCTRanker	0.974	0.944
RPCRanker	0.970	0.947
RegRanker	0.967	0.947
ARFRanker	0.978	0.955

830 7.1.3. Performance of the classifiers in $R_2O_2^*_{(mc)}$

831 As presented in Section 5.2.3, the meta-reasoner: $R_2O_2^*_{(mc)}$ learns the most efficient
832 reasoner on the training data, and predicts the most likely efficient reasoner for an
833 unknown ontology. Using the ontology representation scheme in Equation 5, we
834 considered two classifiers: random forest algorithm and XGBoost.

835 Table 5 shows the prediction performance in terms of classification accuracy on
836 both datasets, **ErrorsRemoved** and **ErrorsReplaced**. As can be seen, the prediction
837 performance is very similar between RF and XGBoost where the best one on each
838 dataset is denoted in bold. In our evaluation, we use XGBoost as it shows the better
839 performance than RF overall.

840 In the following subsection, we present and analyse the evaluation results of our
841 proposed four meta-reasoners on the testing data O_{te} .

Table 5: A summary of performance of the two prediction models for meta-reasoner $R_2O_2^*_{(mc)}$ in terms of accuracy.

Model	ErrorsRemoved	ErrorsReplaced
RF	0.984	0.944
XGBoost	0.984	0.945

842 *7.2. Performance Evaluation of $R_2O_2^*$ on Reasoning Efficiency*

843 In this subsection we discuss the evaluation results comparing our meta-reasoners
 844 with the various component reasoners as well as AutoFolio [31], a state-of-the-art
 845 algorithm selection framework as a strong and robust baseline, following the evaluation
 846 framework presented in Section 6.

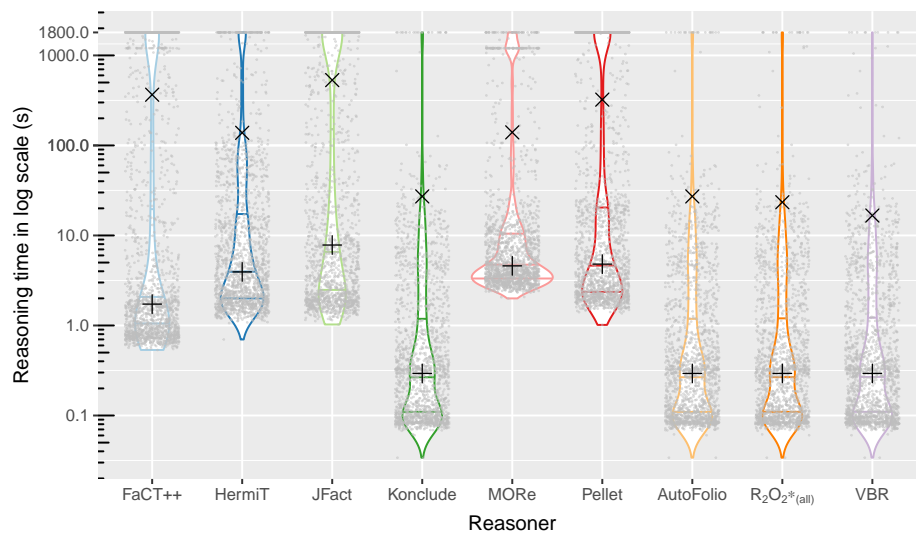


Figure 1: A summary of reasoning time characteristics, in seconds and log-scale, of various reasoners in violin plots.

847 To summarise the performance characteristics of the various component and meta-
 848 reasoners, Figure 1 shows a violin plot of the reasoning time of the reasoners on log
 849 scale, for the experiment **ErrorsReplaced**. A violin plot is a combination of a boxplot
 850 and a mirrored kernel density plot. As a result, a violin plot visualises the underlying

851 distribution that boxplot does not show. Each shape in Figure 1 contains the following
852 components.

- 853 • The (mirrored) violin itself shows the distribution of reasoning time.
- 854 • The cross (×) in the middle shows the mean reasoning time of the reasoner.
- 855 • The plus symbol (+) in the middle shows the median reasoning time of the reasoner.
- 856 • The three horizontal lines within each shape shows the 25%, 50%, and 75% of data,
857 respectively.
- 858 • The grey dots represent the actual reasoning time of all ontologies.

859 As can be seen from the figure, the dominance of Konclude over the other five
860 component reasoners is evident. $R_2O_2^*_{(all)}$ shares similar performance characteristics
861 with AutoFolio and VBR, but with a lower mean reasoning time than AutoFolio. The
862 term VBR stands for the virtual best reasoner, which exhibits the optimal efficiency.
863 Even VBR times out on a number of ontologies (grey dots on 1800.0 at the top of the
864 plot), showing the challenging nature of ontology reasoning. The remainder of the
865 section will present more details and discussions on these performance comparisons.

866 7.2.1. Meta-reasoner time overhead

867 The four different variants of $R_2O_2^*$ all require some additional tasks at both
868 training time and test time. At training time, overall across the 10-fold cross validation,
869 $R_2O_2^*_{(pt)}$ needs to learn regression models (Section 5.2.1); $R_2O_2^*_{(rk)}$ needs to learn
870 rankers (Section 5.2.2), $R_2O_2^*_{(mc)}$ needs to learn a multi-class classifier (Section 5.2.3);
871 and $R_2O_2^*_{(all)}$ ensembles all the above (Section 5.2.4), hence needing to learn all those
872 models. At testing time, each of the meta-reasoners will need to apply these models.

873 We have calculated the time overhead of learning and applying these models. At
874 training time, building a regression model for a reasoner takes 1–2.5 sec (stacking
875 model that combines RF and XGBoost), building a ranker takes 0.3–0.6 sec, building a
876 multi-class classifier takes 0.3–0.5 sec, and building a final stacking model ($R_2O_2^*_{(all)}$)
877 takes an additional 0.1 sec. At testing time, making prediction for a given ontology by
878 $R_2O_2^*_{(all)}$ takes a negligible < 0.5 millisecond.

879 We note that the time overhead at training time does not affect $R_2O_2^*$'s performance
880 as an OWL reasoner. It is only the overhead at testing time that does, as for a new

881 ontology, predictions need to be made for the various models. In all the experiments
882 below in the rest of this section, $R_2O_2^*$'s reported reasoning time already includes the
883 testing-time time overhead.

884 7.2.2. Comparison with our meta-reasoners and component reasoners

885 We now evaluate our meta-reasoners against the six component reasoners, Autofolio
886 and VBR in detail. The evaluation is conducted for the two experiments described
887 in the previous section, **ErrorsRemoved**, in which error ontologies are removed, and
888 **ErrorsReplaced**, in which error ontologies are treated as they time out. Furthermore,
889 for each experiment, we experiment with the inclusion/exclusion of ELK in our meta-
890 reasoners, as described in Section 5. By incorporating ELK, our meta-reasoners can
891 handle the simpler OWL 2 EL ontologies more efficiently, which will improve the
892 overall performance.

893 The reasoners are evaluated on two metrics: (1) average precision at 1 (P@1), which
894 measures whether a reasoner is the most efficient on a given test collection O_{te} across 10-
895 fold cross validation, and (2) average runtime (Avg. runtime), which measures the mean
896 reasoning time on a given test collection O_{tr} in seconds across 10-fold cross-validation.
897 The results presented in the rest of this subsection are the average of those obtained
898 on the test set in each of the 10 folds. Note that in each fold, the test collection O_{te} is
899 differently chosen from the total ontology collection O . A more detailed description
900 about our evaluation framework is found in Section 6.

901 Table 6 and Table 7 show evaluation results for the two experiments, **ErrorsRe-**
902 **moved** and **ErrorsReplaced**, respectively. We note that in the case where ELK is
903 included in our meta-reasoners, the P@1 and Avg. runtime values for ELK are not
904 recorded over the entire test set, but only the subset of OWL 2 EL ontologies. Hence a
905 comparison between ELK and the other reasoners is not meaningful. Note that as we
906 discussed in Section 5.2, $R_2O_2^*_{(pt)}$ follows a similar spirit of the non-ranking portfolio
907 reasoner PR [27], and $R_2O_2^*_{(rk)}$ follows a similar spirit of R_2O_2 [27]. A number of
908 important observations can be made from these tables.

- 909 • In the experiment **ErrorsRemoved**, in all but one case, the best variant of our meta-
910 reasoners outperforms all component reasoners in terms of P@1. Konclude has

Table 6: Performance evaluation in the experiment **ErrorsRemoved**, in which error ontologies are removed. Our meta-reasoners are compared with component reasoners and the virtual best reasoner (VBR). In each column, the best results are **bolded** and the second best results are underlined.

Reasoner	Without ELK		With ELK	
	P@1	Avg. runtime	P@1	Avg. runtime
FaCT++	0.36%	201.62	0.36%	201.62
Hermit	0.14%	54.17	0.14%	54.17
JFact	0.07%	308.47	0.07%	308.47
Konclude	98.78%	<u>8.58</u>	97.26%	8.58
MORe	0.79%	89.93	0.58%	89.93
Pellet	0.22%	164.30	0.22%	164.30
ELK	-	-	1.73%	0.69
$R_2O_2^*_{(pt)}$	<u>98.56%</u>	8.60	98.20%	8.53
$R_2O_2^*_{(rk)}$	98.78%	<u>8.58</u>	<u>98.42%</u>	<u>8.51</u>
$R_2O_2^*_{(mc)}$	98.78%	7.48	98.49%	7.40
$R_2O_2^*_{(all)}$	98.78%	7.48	98.49%	7.40
VBR	100%	4.50	100%	4.43

911 the same performance (98.78%) as our meta-reasoners ($R_2O_2^*_{(rk)}$, $R_2O_2^*_{(mc)}$, and
 912 $R_2O_2^*_{(all)}$) when ELK is not considered.

- 913 • In both experiments **ErrorsRemoved** and **ErrorsReplaced**, in all cases, the best
 914 variant of our meta-reasoners outperforms all component reasoners in terms of average
 915 runtime, including the highly efficient, parallelising reasoner Konclude.
- 916 • Out of all the variants of our meta-reasoners, $R_2O_2^*_{(all)}$ exhibits overall best efficiency.
 917 Of all the four experimental setups, $R_2O_2^*_{(all)}$ exhibits the best performance of three,
 918 and second best in the other one. $R_2O_2^*_{(all)}$ achieves a speedup of at least 1.10x
 919 (over Konclude in experiment **ErrorsReplaced**) and at most 41.69x (over JFact
 920 in experiment **ErrorsRemoved**). Its best efficiency is the result of the stacking
 921 technique employed in $R_2O_2^*_{(all)}$.

Table 7: Performance evaluation of the experiment **ErrorsReplaced**, in which error ontologies are replaced by timeouts (30 minutes). Our meta-reasoners is compared with component reasoners and the virtual best reasoner (VBR). In each column, the best results are **bolded** and the second best results are underlined.

Reasoner	Without ELK		With ELK	
	P@1	Avg. runtime	P@1	Avg. runtime
FaCT++	1.31%	365.92	1.25%	365.92
Hermit	0.68%	138.06	0.68%	138.06
JFact	0.57%	532.11	0.57%	532.11
Konclude	97.22%	27.15	94.55%	27.15
MORe	2.16%	139.68	1.48%	139.68
Pellet	0.91%	322.59	0.85%	322.58
ELK	-	-	3.47%	0.97
$R_2O_2^*_{(pt)}$	97.10%	24.50	96.82%	<u>23.89</u>
$R_2O_2^*_{(rk)}$	97.73%	25.72	<u>97.05%</u>	24.45
$R_2O_2^*_{(mc)}$	97.39%	26.93	96.82%	25.66
$R_2O_2^*_{(all)}$	<u>97.56%</u>	<u>24.67</u>	97.16%	23.39
VBR	100%	16.65	100%	16.22

- The inclusion of ELK in our meta-reasoners indeed improves reasoning efficiency. Even though ELK is only able to handle the less expressive OWL 2 EL profile, our meta-reasoners are able to take advantage of its efficiency in handling such ontologies.

Table 1 in Section 7 shows the vast efficiency dominance of Konclude over the other component reasoners. To better understand the reasoners' performance, we divide the ORE 2015 dataset into four bins of *discretised reasoning time*. The discretisation is performed on the *best reasoning time* (in seconds) of the virtual best reasoner (VBR), into four bins: 'A' (0, 1), 'B' [1, 10), 'C' [10, 100), and 'D' [100, 1,800], which represent ontologies with increasing difficulty.

Tables 8 and 9 below summarise, in each bin and the entire ontology collection O , the percentage of each component reasoner being the most efficient. The component

Table 8: Comparison of percentage of each component reasoner being the most efficient in each bin and overall in experiment **ErrorsRemoved** (Rea: Reasoner, E: ELK, F: FaCT++, H: HermiT, J: JFact, K: Konclude, M: MORE, P: Pellet).

Rank	Overall		A		B		C		D	
	%	Rea	%	Rea	%	Rea	%	Rea	%	Rea
1	96.92	K	100	K	88.62	K	80.56	K	27.27	F,K
2	1.72	E	0	E,F,H,J,M,P	8.94	E	11.11	M	18.18	M
3	0.57	M	-	-	0.81	M,P	5.56	E	9.09	H,J,P
4	0.36	F	-	-	0.41	F,H	2.78	F	0	E
5	0.22	P	-	-	0	J	0	H,J,P	-	-
6	0.14	H	-	-	-	-	-	-	-	-
7	0.07	J	-	-	-	-	-	-	-	-

Table 9: Comparison of percentage of each component reasoner being the most efficient in each bin and overall in experiment **ErrorsReplaced** (Rea: Reasoner, E: ELK, F: FaCT++, H: HermiT, J: JFact, K: Konclude, M: MORE, P: Pellet).

Rank	Overall		A		B		C		D	
	%	Rea	%	Rea	%	Rea	%	Rea	%	Rea
1	91.93	K	100	K	84.41	K	67.31	K	20.83	M
2	3.37	E	0	E,F,H,J,M,P	11.18	E	22.12	E	18.06	F,K
3	1.44	M	-	-	2.35	F	7.69	M	15.28	H
4	1.22	F	-	-	0.88	M,P	1.92	P	13.89	J,P
5	0.83	P	-	-	0.29	H	0.96	F	0	E
6	0.66	H	-	-	0	J	0	H,J	-	-
7	0.55	J	-	-	-	-	-	-	-	-

933 reasoners are ordered by their percentages of being the most efficient, from highest to
934 lowest. Note ELK only performs reasoning on the subset of OWL 2 EL ontologies.
935 Note that the % values in the tables are averaged from the test sets of 10-fold cross
936 validation in the above two experimental setups. Also note that the % values in Tables 8

937 and 9 are different from the P@1 values in Table 6 and Table 7. The % values show
938 the percentage being the most efficient across all the component reasoners. Thus, the
939 sum of the % values in each of the columns - 'Overall', 'A', 'B', 'C' and 'D' is 1. From
940 these tables, we note a number of interesting observations:

- 941 • Konclude’s dominance is again evident, as it is the most efficient reasoner for all
942 bins except only in bin D, experiment **ErrorsReplaced**, where MORE is the most
943 efficient.
- 944 • Despite its dominance, Konclude does not dominate the most challenging category,
945 bin D. In experiment **ErrorsRemoved** both FaCT++ and Konclude are the most
946 efficient, and in experiment **ErrorsReplaced** MORE is the most efficient.
- 947 • For the most difficult category, bin D, many component reasoners are the most efficient
948 for a considerable percentage. In other words, the dominance of any component
949 reasoner is much less pronounced. This shows the challenging nature of algorithm
950 selection in the ontology reasoning context, where for the most challenging instances,
951 a large number of choices are possible. This observation also indicates considerable
952 room for further investigation.

953 7.2.3. Comparison with AutoFolio

954 In this section, we evaluate $R_2O_2^*_{(all)}$ against AutoFolio for the ontology classifica-
955 tion problem¹². We have chosen $R_2O_2^*_{(all)}$ for comparison with AutoFolio as it shows
956 the best performance overall in our evaluation as discussed in the previous section.
957 AutoFolio and $R_2O_2^*_{(all)}$ use the same set of metrics, and are trained and tested on the
958 same data splits and the same set of six reasoners. Hence, the comparison is fair. Note
959 that AutoFolio makes use of functionally equivalent components hence it cannot take
960 advantage of ELK. As described in Section 6, 10-fold cross validation is carried out for
961 both experiments. In each fold, an AutoFolio model is trained on the training set, and
962 then evaluated on the test set.

963 Table 10 summarises the mean reasoning time for the two experiments. As can
964 be observed in the table, $R_2O_2^*_{(all)}$ outperforms AutoFolio in both experiments. The

¹²AutoFolio is available at <https://github.com/mlindauer/AutoFolio>.

965 best performance is achieved when ELK is incorporated in $R_2O_2^*_{(all)}$, where $R_2O_2^*_{(all)}$
 966 outperforms AutoFolio by 15.95% and 16.08% respectively. However, even without
 967 ELK, $R_2O_2^*_{(all)}$ also outperforms AutoFolio, by 14.71% and 10.05% respectively. These
 968 results demonstrate the effectiveness of $R_2O_2^*_{(all)}$ as AutoFolio represents the state-of-
 969 the-art method in automated algorithm selection.

Table 10: Summary of mean reasoning performance (in seconds) comparison between AutoFolio and $R_2O_2^*_{(all)}$ ($R_2O_2^*_{(all)}$, both with and without ELK). In each row, best performance in the test set is highlighted in **bold**, and second best performance in the test set is underlined.

Experiment	AutoFolio	$R_2O_2^*_{(all)}$ (without ELK)	$R_2O_2^*_{(all)}$ (with ELK)
ErrorsRemoved	8.58	<u>7.48</u>	7.40
ErrorsReplaced	27.15	<u>24.67</u>	23.39

970 We further analysed the component reasoners selected by AutoFolio and $R_2O_2^*_{(all)}$
 971 in both experiments. In **ErrorsRemoved**, AutoFolio exclusively selects Konclude for
 972 all 1,389 instances. In **ErrorsReplaced**, AutoFolio selects MORE in 11 (0.625%) of
 973 the 1,760 instances, and Konclude 1,749 (99.375%) instances.

974 On the other hand, $R_2O_2^*_{(all)}$ is able to select a more diverse set of efficient reasoners.
 975 Table 11 shows, with ELK included, the number and percentage of times each compo-
 976 nent reasoner is selected by $R_2O_2^*_{(all)}$ for each bin as well as the entire dataset. These
 977 results provide additional evidence that always selecting the most-efficient-on-average
 978 reasoner(s) is not the most optimal approach. For example, Table 1 shows that the
 979 mean reasoning time of MORE is more than 20x slower than Konclude. However, in
 980 experiments **ErrorsReplaced**, MORE is selected almost 30% among the most diffi-
 981 cult ontologies (bin D). Thus, this analysis further validates the effectiveness of the
 982 sophisticated meta-reasoning framework $R_2O_2^*_{(all)}$.

983 7.3. Key Metrics Identification

984 Lastly, we investigate the identification of each of the 29 metrics' influence on the
 985 performance of our meta-reasoners. This will provide insight into the contribution of
 986 individual metrics we have chosen (i.e. 29 metrics) on their performance. We measure

Table 11: The number and percentage of each component reasoner being selected by $R_2O_2^*$ _(all) in both experiments, with ELK included. Percentages are calculated column-wise.

	Overall (%)	A (%)	B (%)	C (%)	D (%)
ErrorsRemoved					
Konclude	1,354 (97.48)	1,101 (100)	219 (88.66)	28 (82.35)	6 (85.71)
MORe	1 (0.07)	0	0	0	1 (14.29)
ELK	34 (2.45)	0	28 (11.34)	6 (17.65)	0
ErrorsReplaced					
FaCT++	6 (0.34)	0	6 (1.80)	0	0
Konclude	1,682 (95.57)	1,294 (100)	291 (87.12)	79 (75.24)	18 (66.67)
MORe	13 (0.74)	0	0	5 (4.76)	8 (29.63)
Pellet	2 (0.11)	0	0	1 (0.95)	1 (3.70)
ELK	57 (3.24)	0	37 (11.08)	20 (19.05)	0

987 the relative metric importance by using the feature importance values provided by the
988 XGBoost classifier used in $R_2O_2^*$ _(mc). Figure 2 shows the most important metrics for
989 the prediction task of $R_2O_2^*$ _(mc), which is to predict the most efficient reasoner, in both
990 experiments. The weights are estimated by calculating the average of importance of the
991 metrics through 10-fold cross validation and normalised into [0, 1]. The weight of each
992 metric is measured based on the number of times it is used to split the data across all
993 trees in XGBoost.

994 As can be seen from the figure, different metrics are important for each experiment.
995 However, there are some common important metrics. Two metrics, SOV and mCID,
996 are common to both experiments' top 5 metrics. Six metrics, IND, SOV, mCID, aCID,
997 RHLC, and EOG, are common to both experiments' top 10 metrics. For both exper-
998 iments, OBP% is the least important among the 29 metrics. Curiously, DTP% is the
999 most important for experiment **ErrorsRemoved**, while it is the second least important
1000 for **ErrorsReplaced**.

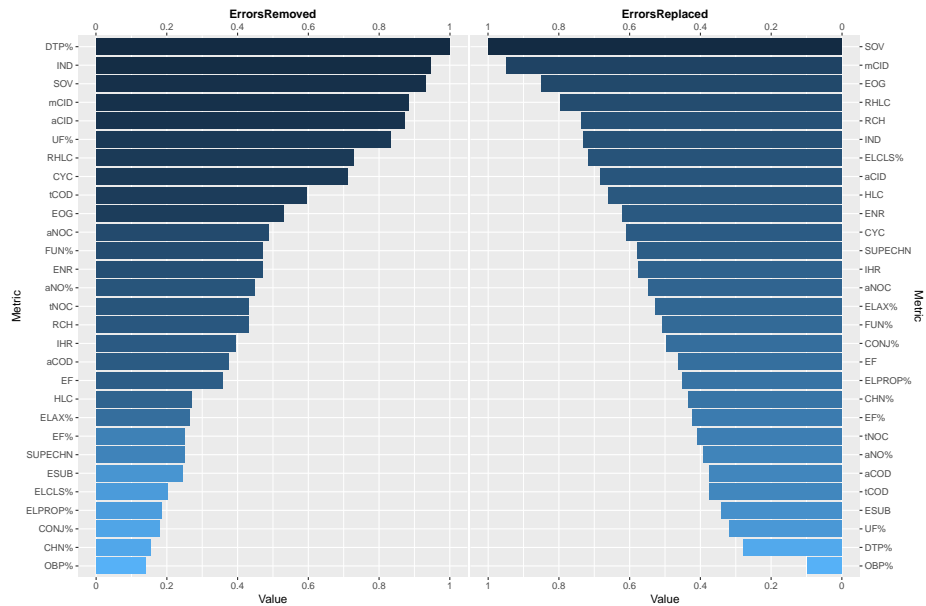


Figure 2: Importance of ontology metrics for $R_2O_2^*(mc)$ in both experiments. For each experiment the metrics are ordered in descending order by their importance values.

1001 8. Conclusion

1002 Reasoning support for OWL ontologies is essential for ensuring the correctness of
 1003 ontologies, and for inferring implicit knowledge from them. For an expressive ontology
 1004 language such as *SHOIN(D)* and *SROIQ(D)*, worst-case complexity is very high.
 1005 Moreover, ample empirical evaluation has also confirmed the hardness of actual, real-
 1006 world ontologies, even on state-of-the-art ontology reasoners such as Konclude and
 1007 HermiT.

1008 This paper presented $R_2O_2^*$, a novel, robust meta-reasoning framework that au-
 1009 tomatically ranks component reasoners by efficiency and selects the one that is most
 1010 likely the most efficient for any given ontology. The $R_2O_2^*$ framework comprises a
 1011 number of novel contributions: (1) we learn regression models that accurately predict
 1012 reasoning time for a number of state-of-the-art ontology reasoners; and (2) we propose
 1013 a learning- and ranking-based meta-reasoner that ensembles base prediction models
 1014 and thus combines component reasoners based on their predicted reasoning efficiency,

1015 and (3) we formally define a large suite of syntactic and structural metrics that describe
1016 ontologies.

1017 We performed a comprehensive evaluation on six state-of-the-art OWL 2 DL rea-
1018 soners and a large corpus of carefully curated ontologies. Our evaluation shows that
1019 $R_2O_2^*$ significantly outperforms all six component reasoners as well as AutoFolio, a
1020 strong, general-purpose, and state-of-the-art algorithm selection system. Compared to
1021 component reasoners, $R_2O_2^*$ achieves a speedup of at least 1.10x (over Konclude) and
1022 up to 41.69x (over JFact).

1023 Extending our $R_2O_2^*$ meta-reasoning framework to a multi-criteria setting as done
1024 in Multi-RakSOR [49, 50] is a natural next step. We plan to extend our methodology
1025 to support ABox reasoning, and investigate support of other non-standard reasoning
1026 problems. Continuing on our work on ABox-intensive EL ontologies [38] and reasoning
1027 on the Android platform [39], we will further investigate sources of ABox reasoning
1028 hardness by studying structural and syntactic properties of ABoxes. Performance pre-
1029 diction and optimisation utilising machine learning techniques is particularly interesting
1030 and relevant in the context of ontology-based data access (OBDA) [5], where a large
1031 database is enhanced by an ontology, and (conjunctive) query answering on the database
1032 requires ontology reasoning [61]. We will also investigate the generation of synthetic,
1033 yet realistic benchmark ontologies (TBoxes) and instances (ABoxes) to assist in the eval-
1034 uation and optimisation of reasoners. Finally, investigating the correlation of efficiency
1035 between different reasoning tasks by a same reasoner is also an interesting problem
1036 worthy of investigation.

1037 **References**

- 1038 [1] M. Ashburner, C. A. Ball, J. A. Blake, et al., Gene Ontology: Tool for the Unifica-
1039 tion of Biology, *Nat Genet* 25 (1) (2000) 25–29. doi:10.1038/75556.
1040 URL <http://dx.doi.org/10.1038/75556>
- 1041 [2] J. Z. Pan, S. Staab, U. Aßmann, J. Ebert, Y. Zhao (Eds.), *Ontology-Driven Software*
1042 *Development*, Springer, 2013.

- 1043 [3] M. Lenzerini, Data integration: a theoretical perspective, in: Proceedings of
1044 the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of
1045 database systems, PODS '02, ACM, New York, NY, USA, 2002, pp. 233–246.
1046 doi:<http://doi.acm.org/10.1145/543613.543644>.
1047 URL <http://doi.acm.org/10.1145/543613.543644>
- 1048 [4] A. Cali, D. Calvanese, G. De Giacomo, M. Lenzerini, Data integration under
1049 integrity constraints, *Information Systems* 29 (2) (2004) 147–163.
- 1050 [5] R. Kontchakov, C. Lutz, D. Toman, F. Wolter, M. Zakharyashev, The combined
1051 approach to ontology-based data access, in: *IJCAI*, Vol. 11, 2011, pp. 2656–2661.
- 1052 [6] Y.-F. Li, G. Kennedy, F. Ngoran, P. Wu, J. Hunter, An ontology-centric architecture
1053 for extensible scientific data management systems, *Future Gener. Comput. Syst.*
1054 29 (2) (2013) 641–653. doi:10.1016/j.future.2011.06.007.
1055 URL <http://dx.doi.org/10.1016/j.future.2011.06.007>
- 1056 [7] I. Horrocks, P. F. Patel-Schneider, F. van Harmelen, From *SHIQ* and RDF to
1057 OWL: The Making of a Web Ontology Language, *Journal of Web Semantics* 1 (1)
1058 (2003) 7–26.
- 1059 [8] B. Cuenca-Grau, I. Horrocks, B. Motik, B. Parsia, P. Patel-Schneider, U. Sattler,
1060 OWL 2: The next step for OWL, *Journal of Web Semantics: Science, Services and*
1061 *Agents on the World Wide Web* 6 (2008) 309–322. doi:10.1016/j.websem.
1062 2008.05.001.
1063 URL <http://portal.acm.org/citation.cfm?id=1464505.1464604>
- 1064 [9] F. Baader, U. Sattler, An overview of tableau algorithms for description logics,
1065 *Studia Logica* 69 (1) (2001) 5–40.
- 1066 [10] D. Tsarkov, I. Horrocks, FaCT++ description logic reasoner: System description,
1067 in: *Proc. of the Int. Joint Conf. on Automated Reasoning (IJCAR 2006)*, Springer,
1068 2006, pp. 292–297.

- 1069 [11] R. Shearer, B. Motik, I. Horrocks, Hermit: A Highly-Efficient OWL Reasoner,
1070 in: Proceedings of the 5th International Workshop on OWL: Experiences and
1071 Directions (OWLED 2008), 2008.
- 1072 [12] A. Steigmiller, T. Liebig, B. Glimm, Konclude: System description, J. Web Sem.
1073 27 (2014) 78–85. doi:10.1016/j.websem.2014.06.003.
1074 URL <http://dx.doi.org/10.1016/j.websem.2014.06.003>
- 1075 [13] E. Sirin, B. Parsia, B. Cuenca-Grau, A. Kalyanpur, Y. Katz, Pellet: A practical
1076 OWL-DL reasoner, Web Semantics: Science, Services and Agents on the World
1077 Wide Web 5 (2) (2007) 51–53. doi:10.1016/j.websem.2007.03.004.
1078 URL <http://dx.doi.org/10.1016/j.websem.2007.03.004>
- 1079 [14] J. Z. Pan, Y. Ren, Y. Zhao, Tractable approximate deduction for OWL, Artificial
1080 Intelligence 235 (2016) 95–155.
- 1081 [15] K. Dentler, R. Cornet, A. ten Teije, N. de Keizer, Comparison of reasoners for
1082 large ontologies in the OWL 2 EL profile, Semantic Web Journal 2 (2) (2011)
1083 71–87.
- 1084 [16] Y.-B. Kang, Y.-F. Li, S. Krishnaswamy, A rigorous characterization of reasoning
1085 performance – a tale of four reasoners, in: Proceedings of the 1st International
1086 Workshop on OWL Reasoner Evaluation (ORE-2012), 2012.
- 1087 [17] M. Y. Vardi, Boolean satisfiability: Theory and engineering, Commun. ACM 57 (3)
1088 (2014) 5–5. doi:10.1145/2578043.
1089 URL <http://doi.acm.org/10.1145/2578043>
- 1090 [18] K. Leyton-Brown, H. H. Hoos, F. Hutter, L. Xu, Understanding the empirical
1091 hardness of *NP*-complete problems, Commun. ACM 57 (5) (2014) 98–107. doi:
1092 10.1145/2594413.2594424.
1093 URL <http://doi.acm.org/10.1145/2594413.2594424>
- 1094 [19] J. Z. Pan, Benchmarking DL reasoners using realistic ontologies, in: B. Cuenca-
1095 Grau, I. Horrocks, B. Parsia, P. F. Patel-Schneider (Eds.), OWLED, Vol. 188 of
1096 CEUR Workshop Proceedings, CEUR-WS.org, 2005.

- 1097 [20] T. Gardiner, I. Horrocks, D. Tsarkov, Automated benchmarking of description
1098 logic reasoners, in: Proceedings of the 2006 International Workshop on Description
1099 Logics (DL2006), 2006.
- 1100 [21] S. Bail, B. Glimm, E. Jiménez-Ruiz, N. Matentzoglou, B. Parsia, A. Steigmiller,
1101 Summary ORE 2014 competition, in: Proceedings of the 3rd International Work-
1102 shop on OWL Reasoner Evaluation (ORE 2014), Vol. 1207, CEUR Workshop
1103 Proceedings, 2014, pp. iv–vii.
- 1104 [22] B. Parsia, N. Matentzoglou, R. S. Gonçalves, B. Glimm, A. Steigmiller, The
1105 OWL reasoner evaluation (ORE) 2015 competition report, *Journal of Automated*
1106 *Reasoning* 59 (4) (2017) 455–482.
- 1107 [23] R. S. Gonçalves, N. Matentzoglou, B. Parsia, U. Sattler, The empirical robustness
1108 of description logic classification, in: T. Eiter, B. Glimm, Y. Kazakov, M. Krötzsch
1109 (Eds.), *Description Logics*, Vol. 1014 of CEUR Workshop Proceedings, CEUR-
1110 WS.org, 2013, pp. 197–208.
- 1111 [24] H. Zhang, Y.-F. Li, H. B. K. Tan, Measuring design complexity of Semantic
1112 Web ontologies, *Journal of Systems and Software* 83 (5) (2010) 803–814.
1113 doi:DOI:10.1016/j.jss.2009.11.735.
1114 URL [http://www.sciencedirect.com/science/article/
1115 B6V0N-4XV06RX-2/2/1024981a176a351b9dbb0cc4c487c17e](http://www.sciencedirect.com/science/article/B6V0N-4XV06RX-2/2/1024981a176a351b9dbb0cc4c487c17e)
- 1116 [25] Y.-B. Kang, Y.-F. Li, S. Krishnaswamy, Predicting reasoning performance using
1117 ontology metrics, in: Cudré-Mauroux et al. [62], pp. 198–214.
- 1118 [26] Y.-B. Kang, J. Z. Pan, S. Krishnaswamy, W. Sawangphol, Y.-F. Li, How long will
1119 it take? Accurate prediction of ontology reasoning performance, in: C. E. Brodley,
1120 P. Stone (Eds.), *AAAI*, AAAI Press, 2014, pp. 80–86.
- 1121 [27] Y.-B. Kang, S. Krishnaswamy, Y.-F. Li, R2O2: an efficient ranking-based reasoner
1122 for OWL ontologies, in: M. Arenas, Ó. Corcho, E. Simperl, M. Strohmaier,
1123 M. d’Aquin, K. Srinivas, P. T. Groth, M. Dumontier, J. Heflin, K. Thirunarayan,
1124 S. Staab (Eds.), *The Semantic Web - ISWC 2015 - 14th International Semantic*

- 1125 Web Conference, Part I, Vol. 9366 of Lecture Notes in Computer Science, Springer,
1126 2015, pp. 322–338. doi:10.1007/978-3-319-25007-6_19.
1127 URL http://dx.doi.org/10.1007/978-3-319-25007-6_19
- 1128 [28] T. Chen, C. Guestrin, XGBoost: A scalable tree boosting system, in: Proceedings
1129 of the 22nd ACM SIGKDD International Conference on Knowledge Discovery
1130 and Data Mining, KDD '16, ACM, New York, NY, USA, 2016, pp. 785–794.
1131 doi:10.1145/2939672.2939785.
1132 URL <http://doi.acm.org/10.1145/2939672.2939785>
- 1133 [29] J. H. Friedman, Greedy function approximation: a gradient boosting machine,
1134 *Annals of statistics* (2001) 1189–1232.
- 1135 [30] Y. Kazakov, M. Krötzsch, F. Simancik, The incredible ELK - from polynomial
1136 procedures to efficient reasoning with \mathcal{EL} ontologies, *J. Autom. Reasoning* 53 (1)
1137 (2014) 1–61. doi:10.1007/s10817-013-9296-3.
1138 URL <http://dx.doi.org/10.1007/s10817-013-9296-3>
- 1139 [31] M. Lindauer, H. H. Hoos, F. Hutter, T. Schaub, AutoFolio: An automatically
1140 configured algorithm selector, *Journal of Artificial Intelligence Research* 53 (2015)
1141 745–778.
- 1142 [32] F. Baader, W. Nutt, Basic description logics, in: F. Baader, D. Calvanese,
1143 D. McGuinness, D. Nardi, P. Patel-Schneider (Eds.), *The description logic hand-*
1144 *book: theory, implementation, and applications*, Cambridge University Press, 2003,
1145 pp. 43–95.
- 1146 [33] I. Horrocks, The FaCT system, *Tableaux'98*, LNCS 1397 (1998) 307–312.
- 1147 [34] B. Glimm, I. Horrocks, B. Motik, G. Stoilos, Z. Wang, Hermit: An OWL
1148 2 reasoner, *J. Autom. Reason.* 53 (3) (2014) 245–269. doi:10.1007/
1149 s10817-014-9305-1.
1150 URL <http://dx.doi.org/10.1007/s10817-014-9305-1>
- 1151 [35] R. S. Gonçalves, B. Parsia, U. Sattler, Performance heterogeneity and approximate
1152 reasoning in description logic ontologies, in: Cudré-Mauroux et al. [62], pp. 82–98.

- 1153 [36] D. Zhang, C. Ye, Z. Yang, An evaluation method for ontology complexity analysis
1154 in ontology evolution, in: *Managing Knowledge in a World of Networks*, Vol.
1155 4248, 2006, pp. 214–221.
- 1156 [37] V. Sazonau, U. Sattler, G. Brown, Predicting OWL reasoners: Locally or globally?,
1157 in: M. Bienvenu, M. Ortiz, R. Rosati, M. Simkus (Eds.), *Informal Proceedings*
1158 *of the 27th International Workshop on Description Logics*, Vienna, Austria, July
1159 17-20, 2014., Vol. 1193 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2014,
1160 pp. 713–724.
1161 URL http://ceur-ws.org/Vol-1193/paper_12.pdf
- 1162 [38] J. Z. Pan, C. Bobed, I. Guclu, F. Bobillo, M. J. Kollingbaum, E. Mena, Y.-F.
1163 Li, Predicting reasoner performance on ABox intensive OWL 2 EL ontologies,
1164 *International Journal on Semantic Web and Information Systems (IJSWIS)* 14 (1).
- 1165 [39] I. Guclu, Y.-F. Li, J. Z. Pan, M. J. Kollingbaum, Predicting energy consumption of
1166 ontology reasoning over mobile devices, in: P. T. Groth, E. Simperl, A. J. G. Gray,
1167 M. Sabou, M. Krötzsch, F. Lécué, F. Flöck, Y. Gil (Eds.), *The Semantic Web -*
1168 *ISWC 2016 - 15th International Semantic Web Conference*, Kobe, Japan, October
1169 17-21, 2016, *Proceedings, Part I*, Vol. 9981 of *Lecture Notes in Computer Science*,
1170 2016, pp. 289–304. doi:10.1007/978-3-319-46523-4_18.
1171 URL http://dx.doi.org/10.1007/978-3-319-46523-4_18
- 1172 [40] J. R. Rice, The algorithm selection problem, in: *Advances in computers*, Vol. 15,
1173 Elsevier, 1976, pp. 65–118.
- 1174 [41] K. A. Smith-Miles, Cross-disciplinary perspectives on meta-learning for algorithm
1175 selection, *ACM Computing Surveys (CSUR)* 41 (1) (2009) 6.
- 1176 [42] F. Hutter, L. Xu, H. H. Hoos, K. Leyton-Brown, Algorithm runtime predic-
1177 tion: Methods & evaluation, *Artif. Intell.* 206 (2014) 79–111. doi:10.1016/j.
1178 *artint*.2013.10.003.
1179 URL <http://dx.doi.org/10.1016/j.artint.2013.10.003>

- 1180 [43] L. Xu, F. Hutter, H. H. Hoos, K. Leyton-Brown, SATzilla: Portfolio-based algo-
1181 rithm selection for SAT, *J. Artif. Int. Res.* 32 (1) (2008) 565–606.
1182 URL <http://dl.acm.org/citation.cfm?id=1622673.1622687>
- 1183 [44] D. Tsarkov, I. Palmisano, Chainsaw: a metareasoner for large ontologies, in: I. Hor-
1184 rocks, M. Yatskevich, E. Jiménez-Ruiz (Eds.), *Proceedings of the 1st International*
1185 *Workshop on OWL Reasoner Evaluation (ORE-2012)*, Manchester, UK, July 1st,
1186 2012, Vol. 858 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2012.
1187 URL http://ceur-ws.org/Vol-858/ore2012_paper2.pdf
- 1188 [45] B. Cuenca-Grau, I. Horrocks, Y. Kazakov, U. Sattler, Modular reuse of ontologies:
1189 Theory and practice, *J. Artif. Intell. Res. (JAIR)* 31 (2008) 273–318.
- 1190 [46] C. Del Vescovo, B. Parsia, U. Sattler, T. Schneider, The modular struc-
1191 ture of an ontology: Atomic decomposition, in: *Proceedings of the Twenty-*
1192 *Second International Joint Conference on Artificial Intelligence - Volume Vol-*
1193 *ume Three, IJCAI'11, AAAI Press*, 2011, pp. 2232–2237. doi:10.5591/
1194 978-1-57735-516-8/IJCAI11-372.
1195 URL <http://dx.doi.org/10.5591/978-1-57735-516-8/IJCAI11-372>
- 1196 [47] B. Parsia, N. Matentzoglou, R. S. Gonçalves, B. Glimm, A. Steigmiller, The OWL
1197 reasoner evaluation (ORE) 2015 competition report, in: T. Liebig, A. Fokoue
1198 (Eds.), *Proceedings of the 11th International Workshop on Scalable Semantic*
1199 *Web Knowledge Base Systems co-located with 14th International Semantic Web*
1200 *Conference (ISWC 2015)*, Bethlehem, PA, USA, October 11, 2015., Vol. 1457 of
1201 *CEUR Workshop Proceedings*, CEUR-WS.org, 2015, pp. 2–15.
1202 URL http://ceur-ws.org/Vol-1457/SSWS2015_paper1.pdf
- 1203 [48] W. Song, B. Spencer, W. Du, WSReasoner: A prototype hybrid reasoner for
1204 ALCHOI ontology classification using a weakening and strengthening approach,
1205 in: *OWL Reasoner Evaluation Workshop (ORE 2012)*, 2012, p. 1.
- 1206 [49] N. Alaya, S. B. Yahia, M. Lamolle, Raksor: Ranking of ontology reasoners based
1207 on predicted performances, in: *28th IEEE International Conference on Tools with*

- 1208 Artificial Intelligence, ICTAI 2016, San Jose, CA, USA, November 6-8, 2016,
1209 IEEE Computer Society, 2016, pp. 1076–1083. doi:10.1109/ICTAI.2016.
1210 0165.
1211 URL <https://doi.org/10.1109/ICTAI.2016.0165>
- 1212 [50] N. Alaya, M. Lamolle, S. B. Yahia, Multi-label based learning for better multi-
1213 criteria ranking of ontology reasoners, in: C. d’Amato, M. Fernández, V. A. M.
1214 Tamma, F. Lécué, P. Cudré-Mauroux, J. F. Sequeda, C. Lange, J. Heflin (Eds.),
1215 The Semantic Web - ISWC 2017 - 16th International Semantic Web Confer-
1216 ence, Vienna, Austria, October 21-25, 2017, Proceedings, Part I, Vol. 10587 of
1217 Lecture Notes in Computer Science, Springer, 2017, pp. 3–19. doi:10.1007/
1218 978-3-319-68288-4_1.
1219 URL https://doi.org/10.1007/978-3-319-68288-4_1
- 1220 [51] Y. Zhou, B. C. Grau, Y. Nenov, M. Kaminski, I. Horrocks, Pagoda: Pay-as-
1221 you-go ontology query answering using a datalog reasoner, *Journal of Artificial*
1222 *Intelligence Research* 54 (2015) 309–367.
- 1223 [52] N. Fenton, J. Bieman, *Software Metrics: A Rigorous and Practical Approach*,
1224 Third Edition, 3rd Edition, CRC Press, Inc., Boca Raton, FL, USA, 2014.
- 1225 [53] B. Motik, R. Shearer, I. Horrocks, Optimized reasoning in description logics using
1226 hypertableaux, in: *International Conference on Automated Deduction*, Springer,
1227 2007, pp. 67–83.
- 1228 [54] F. M. Donini, Complexity of reasoning, in: *The Description Logic Handbook:*
1229 *Theory, Implementation, and Applications*, 2003, pp. 96–136.
- 1230 [55] F. Baader, S. Brandt, C. Lutz, Pushing the \mathcal{EL} envelope, in: L. P. Kaelbling,
1231 A. Saffiotti (Eds.), *Proceedings of IJCAI 2005*, Professional Book Center, 2005,
1232 pp. 364–369.
1233 URL [http://dblp.uni-trier.de/db/conf/ijcai/ijcai2005.html#](http://dblp.uni-trier.de/db/conf/ijcai/ijcai2005.html#BaaderBL05)
1234 [BaaderBL05](http://dblp.uni-trier.de/db/conf/ijcai/ijcai2005.html#BaaderBL05)

- 1235 [56] A. A. Romero, B. Cuenca-Grau, I. Horrocks, MORE: Modular combination of
1236 OWL reasoners for ontology classification, in: Cudré-Mauroux et al. [62], pp.
1237 1–16.
- 1238 [57] M. P. Sesmero, A. I. Ledezma, A. Sanchis, Generating ensembles of heterogeneous
1239 classifiers using stacked generalization, *Wiley Int. Rev. Data Min. and Knowl.*
1240 *Disc.* 5 (1) (2015) 21–34. doi:10.1002/widm.1143.
1241 URL <http://dx.doi.org/10.1002/widm.1143>
- 1242 [58] V. Haarslev, K. Hidde, R. Möller, M. Wessel, The RacerPro knowledge representa-
1243 tion and reasoning system, *Semantic Web* 3 (3) (2012) 267–277.
- 1244 [59] V. Haarslev, R. Möller, RACER system description, in: *Proceedings of Automated*
1245 *Reasoning : First International Joint Conference*, no. 2083 in *Lecture Notes in*
1246 *Computer Science*, Siena, Italy, 2001, pp. 701–706.
- 1247 [60] Q. Sun, B. Pfahringer, Pairwise meta-rules for better meta-learning-based al-
1248 gorithm ranking, *Machine Learning* 93 (1) (2013) 141–161. doi:10.1007/
1249 s10994-013-5387-y.
- 1250 [61] A. Cali, G. Gottlob, A. Pieris, Towards more expressive ontology languages: The
1251 query answering problem, *Artificial Intelligence* 193 (2012) 87–128.
- 1252 [62] P. Cudré-Mauroux, J. Heflin, E. Sirin, T. Tudorache, J. Euzenat, M. Hauswirth, J. X.
1253 Parreira, J. Hendler, G. Schreiber, A. Bernstein, E. Blomqvist (Eds.), *The Semantic*
1254 *Web - ISWC 2012 - 11th International Semantic Web Conference*, Boston, MA,
1255 USA, November 11-15, 2012, *Proceedings, Part I*, Vol. 7649 of *Lecture Notes in*
1256 *Computer Science*, Springer, 2012.

Table A.12: Definitions of the 24 ontology-level (ONT) metrics. Note that “_” represents “don’t cares”.

Metric	Definition
SOV	No. of named entities (classes, properties & individuals)
ENR	Ratio between number of (named and anonymous) entities and number of edges
TIP	Difference between number of subclass axioms and number of (named or anonymous) classes
EOG	The <i>entropy</i> of the ontology graph, measuring the diversity of the edge distribution
CYC	The <i>cyclomatic complexity</i> of the ontology, measuring the number of linearly independent paths
RCH	The ratio of (possibly nested) anonymous class expressions and all (named or anonymous) class expressions
IND	No. of (named or anonymous) individuals
GCI	No. of GCI axioms
HGCI	No. of hidden GCI axioms
ESUB%	Ratio of subclass axioms that contain (nested) existential restrictions ($\exists R. _$)
DSUB%	Ratio of subclass axioms that contain (nested) class unions ($_ \sqcup _$)
CSUB%	Ratio of subclass axioms that contain (nested) class intersections ($_ \sqcap _$) and the subclass is anonymous
ELCLS%	Ratio of (nested) class expressions that are in OWL 2 EL profile
ELAX%	Ratio of subclass or equivalent class axioms that only contain expressions in the OWL 2 EL profile

Table A.12: Definitions of the 24 ontology-level (ONT) metrics. Note that “_” represents “don’t cares”.

Metric	Definition
HLC	Count hard language constructs containing in superclass expressions
HLC%	Ratio of HLC and subclass axioms
SUBCECHN	No. of top-level subclass expressions that contain chained existential restrictions
SUPCECHN	No. of top-level superclass expressions that contain chained existential restrictions
DSUBCECHN	Max depth of chained existential restrictions in a subclass expression
DSUPCECHN	Max depth of chained existential restrictions in a superclass expression
SUBCCHN	No. of top-level subclass expressions that contain chains of conjunctions
SUPCCHN	No. of top-level superclass expressions that contain chains of conjunctions
DSUBCCHN	Max depth of nested conjunctions in a subclass expression
DSUPCCHN	Max depth of nested conjunctions in a superclass expression

Table A.13: Definitions of the 15 class-level (CLS) metrics. Note that C represents a (named or possibly nested) class expression in an ontology, and N_C denotes the total number of (named or possibly nested) class expressions in a given ontology.

Metric	Definition
tNOC	$\sum_C NOC(C)$
aNOC	$\frac{tNOC}{N_C}$
mNOC	$\max_C NOC(C)$
tNOP	$\sum_C NOP(C)$
aNOP	$\frac{tNOP}{N_C}$
mNOP	$\max_C NOP(C)$
tCID/tCOD	\sum_C incoming/outgoing edges of C
aCID/aCOD	$aCID = \frac{tCID}{N_C}, aCOD = \frac{tCOD}{N_C}$
mCID/mCOD	$mCID = \max_C CID(C), mCOD = \max_C COD(C)$
tDIT	\sum_C distance of C from owl:Thing in a depth-first manner
aDIT	$\frac{tDIT}{N_C}$
mDIT	$\max_C DIT(C)$

Table A.14: Definition of the 22 (possibly nested) anonymous class expression (ACE) metrics. Note that each row represents a count metric and a ratio metric (represented by [%]) for the same type of class expressions.

Metric	Definition
ENUM[%]	For enumerations/nominals ($\{a, b, c\}$, where a, b, c are individuals)
NEG[%]	For class negations ($\neg C$)
CONJ[%]	For class intersections (conjunctions, $C_1 \sqcap C_2$)
DISJ[%]	For class unions (disjunctions, $C_1 \sqcup C_2$)
UF[%]	For universal restrictions ($\forall R.C$)
EF[%]	For existential restrictions ($\exists R.C$)
VALUE[%]	For value restrictions ($\exists R.\{a\}$), where a is an individual
SELF[%]	For self references
MNCAR[%]	For min cardinality restrictions ($\geq n R.C$)
MXCAR[%]	For max cardinality restrictions ($\leq n R.C$)
CAR[%]	For (exact) cardinality restrictions ($= n R.C$)

Table A.15: Definitions of the 30 property-level (PRO) metrics.

Metric	Definition
OBP[%]	Count and ratio of object-properties
DTP[%]	Count and ratio of datatype-properties
FUN[%]	Count and ratio of functional properties
SYM[%]	Count and ratio of symmetric properties
TRN[%]	Count and ratio of transitive properties
IFUN[%]	Count and ratio of inverse functional properties
ASYM[%]	Count and ratio of asymmetric properties
REFLE[%]	Count and ratio of reflexive properties
IRREF[%]	Count and ratio of irreflexive properties
CHN[%]	Count and ratio of property chain axioms
SUBP	Count of subproperty axioms
EQVP	Count of equivalent property axioms
DISP	Count of disjoint property axioms
INV	Count of inverse property axioms
DOMN	Count of domain axioms
RANG	Count of range axioms
ELPROP%	Ratio of property axioms that are in the OWL 2 EL profile
IHR	Count of class axioms that involve property hierarchies
IIR	Count of class axioms that involve inverse properties
ITR	Count of class axioms that involve transitive properties