

Syntactic Regularities in OWL Ontologies based on Language Abstractions

Christian Kindermann, Bijan Parsia, and Uli Sattler

University of Manchester, Oxford Rd, Manchester M13 9PL, UK

`christian.kindermann,bijan.parsia,uli.sattler@manchester.ac.uk`

Abstract. OWL ontologies are built and maintained on the basis of all sorts of methods and methodologies using a wide range of tools. As ontologies are primarily published as sets of axioms, its underlying design principles often remain opaque. However, a principled and systematic ontology design is likely to be reflected in regularities for axioms. Identifying such regularities may help to recover and unveil conscious design choices and otherwise recurring modelling practices. In this work, we propose a framework for identifying syntactic regularities for axioms in an automated manner. Using this framework, we survey ontologies indexed in BioPortal. We find that most axioms conform to a few number of syntactic regularities. Since conceptual entities are often represented by more than one axiom in OWL, we also motivate a notion for syntactic regularities over sets of axioms. For sets of axioms, an ontology’s design appears more variegated. Still, we can often identify a small number of prevalent regularities.

1 Introduction

OWL ontologies are built and maintained on the basis of all sorts of methods and methodologies using a wide range of tools. A basic level of interoperability and reuse is guaranteed by the W3C standards for the Semantic Web.¹ Here, an OWL ontology is defined to consist of two main components: a set of axioms, on the one hand, and a set of imported ontologies, on the other hand. In theory, ontology imports can be used as a basic means to structure an ontology in a meaningful way. Yet, in practice, ontology imports are not used to divide an ontology into parts of fine granularity. As a result, underlying design principles and consistently used modelling techniques remain opaque in published ontologies.

However, a principled and systematic ontology design is likely to be reflected in regularities for axioms. Identifying such regularities may help recover and unveil conscious design choices and otherwise recurring modelling practices. This arguably helps with ontology comprehension and quality assurance in practice.

In this work, we propose a framework for identifying regularities for axioms in an automated manner. The contributions are as follows: (i) we develop a formal framework for discovering syntactic regularities in OWL ontologies, (ii) we propose metrics to qualify discovered regularities and to characterise an ontology’s overall design, (iii) we survey ontologies indexed in BioPortal using our framework and metrics.

¹ <https://www.w3.org/standards/semanticweb/>

2 Preliminaries

We follow the terminology and notation introduced in [1] but adapt definitions for our purposes as needed.²

Terms and Trees We define *terms* in the usual manner by an inductively defined set $T(\Sigma, \mathcal{X})$ over a set of both ranked and unranked symbols Σ as well as a set of constant symbols \mathcal{X} of variables. A term may be viewed as a *tree* (a connected acyclic undirected graph) the leaves of which are labeled with constant symbols (or variables) and internal nodes are labeled with symbols of positive arity. In the following we do not distinguish between terms and their corresponding trees.

Substitutions A *substitution* σ is a function $\mathcal{X} \rightarrow T(\Sigma, \mathcal{X})$ where there are only finitely many variables not mapped to themselves. The *domain* of a substitution σ is the subset of variables $x \in \mathcal{X}$ such that $\sigma(x) \neq x$. The substitution $\{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$ is the identity on $\mathcal{X} \setminus \{x_1, \dots, x_n\}$ and maps $x_i \in \mathcal{X}$ on $t_i \in T(\Sigma, \mathcal{X})$ for every index $1 \leq i \leq n$. Substitutions are extended to $f(t_1, \dots, t_n) \in T(\Sigma, \mathcal{X})$ by $\sigma(f(t_1, \dots, t_n)) = f(\sigma(t_1), \dots, \sigma(t_n))$.

In the following, we do not distinguish between a substitution and its extension to $T(\Sigma, \mathcal{X})$, and as usual, we will often use substitutions in postfix notation, i.e., $t\sigma$ is the result of applying σ to the term t .

Contexts A term $C \in T(\Sigma, \mathcal{X})$ with variables is called a *context* and we write $C[x_1, \dots, x_n]$ to denote a context C with variables x_1, \dots, x_n . We write $C[t_1, \dots, t_n]$ to denote the term $C\sigma \in T(\Sigma)$ where σ is the substitution $\{x_1 \mapsto t_1, \dots, x_n \mapsto t_n\}$.

Tree Transducers We define transformations on trees as finite bottom-up tree transducers. A *tree transducer* is a tuple $A = (Q, \Sigma, \Sigma', Q_f, \Delta)$ where Σ is the input alphabet, Σ' is the output alphabet, Q is a set of (unary) states, $Q_f \subseteq Q$ is a set of final states, and Δ is a set of transduction rules for trees.

For the purpose of treating ranked and unranked symbols in a uniform manner, we describe the left-hand side of transduction rules by $f(h)$, where f is a symbol in Σ and h is a (finite) sequence of symbols. To ensure that h is of suitable length, we condition each transduction rule to a language \mathcal{H} , called a *horizontal language*³ that formulates constraints for h as needed. With this, a transduction rule in Δ is of the following form:

$$f(h) \rightarrow q(\phi(h)) \text{ subject to } h \in \mathcal{H},$$

where $f \in \Sigma, q \in Q, \mathcal{H}$ a horizontal language, and ϕ is a function $\mathcal{H} \rightarrow T(\Sigma', \mathcal{X})$. Note, that each transduction rule has a state q as the root of the tree on its right-hand side. Therefore we can assume h to be a sequence of states or the empty tree in case of

² This paper is accompanied by a technical report providing a more detailed presentation of the used notions with elaborate examples: <https://github.com/ckindermann/syntacticRegularities>.

³ The notion of a horizontal language is closely related to the notion of *hedge languages* for unranked symbols in the literature. However, as we work with both ranked and unranked symbols, we decided against overloading this commonly used terminology to avoid confusion.

rules for leafs. To make this precise, we restrict horizontal languages \mathcal{H} as follows:

$$\begin{aligned} \mathcal{H} \subseteq \mathcal{H}_Q = \{ \varepsilon \} \cup \{ q_1(x_1) \cdots q_n(x_n) \mid & n \in \mathbb{N}, \\ & q_1, \dots, q_n \in Q, \\ & x_1, \dots, x_n \in \mathcal{X}, \\ & x_i \neq x_j \text{ for } i, j \in \{1, \dots, n\} \text{ and } i \neq j \}. \end{aligned}$$

So $h \in \mathcal{H}_Q$ is a (finite) sequence of (single variable) contexts $q_i(x_i)$. Each context in such a sequence can be uniquely identified by its variable x_i . Note, however, that the states of contexts $q_i(x_i)$ and $q_j(x_j)$ with $x_i \neq x_j$ in a given sequence can be identical. Transduction rules with $\mathcal{H} = \{ \varepsilon \}$ can be applied to the leaves of a tree will be written $a \rightarrow q(\phi(a))$ for $a \in \Sigma$.

In the following, we assume horizontal languages to be regular languages. Furthermore, we will use regular expressions to define them. We assume regular expressions using the Kleene star to introduce fresh individuals as required. For example $q(x_1)q(x_2) \in Q(x)^*$ but $q(x_1)q(x_1) \notin Q(x)^*$ (assuming $q \in Q$).

Let $t, t' \in T(\Sigma \cup \Sigma' \cup Q)$. The move relation \rightarrow_A for a tree transducer A is defined by

$$t \rightarrow_A t' \Leftrightarrow \begin{cases} \exists f(h) \rightarrow q(\phi(h)) \text{ subject to } h \in \mathcal{H} \in \Delta: \\ \exists q_1(x_1) \cdots q_n(x_n) = h \in \mathcal{H}: \\ \exists u_1, \dots, u_n \in T(\Sigma'): \\ \exists C \in \mathcal{C}(\Sigma \cup \Sigma' \cup Q): \\ t = C[f(q_1(u_1), \dots, q_n(u_n))] \wedge \\ t' = C[q(\phi(h)\{x_1 \rightarrow u_1, \dots, x_k \rightarrow u_k\})] \end{cases}$$

The reflexive and transitive closure of \rightarrow_A is \rightarrow_A^* . A transduction from a ground term $t \in T(\Sigma)$ to a ground term $t' \in T(\Sigma')$ is a sequence of move steps of the form $t \rightarrow_A^* q(t')$, such that q is a final state. The relation between terms $T(\Sigma)$ and $T(\Sigma')$ induced by A is the relation defined by

$$A^* = \{ (t, t') \mid t \rightarrow_A^* q(t'), t \in T(\Sigma), t' \in T(\Sigma'), q \in Q_f \}.$$

A tree transducer is *deterministic* if there are no rules with the same left-hand side. In the following, we assume transducers to be deterministic and we write $A(t) = t'$ if $(t, t') \in A^*$.

Graphs, Graph Minors, and Graph Isomorphisms A *graph* G is an ordered pair (V, E) where V is a set of *vertices*, and $E \subseteq \{ \{v_1, v_2\} \mid (v_1, v_2) \in V \times V \wedge v_1 \neq v_2 \}$ is a set of *edges*. A *contraction* of an edge $e = \{v_1, v_2\} \in E$ denotes the removal of e from E and the replacement of both v_1 and v_2 by a single node v' s.t. any node adjacent to either v_1 or v_2 is adjacent to v' . A *minor* of a graph is a graph obtained by repeatedly contracting edges, removing edges, or removing vertices without adjacent nodes. A *proper minor* is a minor other than the graph itself. A *graph isomorphism* between two graphs $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ is a bijection $f: V_1 \rightarrow V_2$ s.t. $\{v_1, v'_1\} \in E_1$ iff $\{f(v_1), f(v'_1)\} \in E_2$.

OWL Ontologies The Web Ontology Language (OWL) is defined defined in terms of a structural specification.⁴ This specification, however, does not define a formal lan-

⁴ <https://www.w3.org/TR/owl2-syntax/>

guage in the conventional sense. Rather, it specifies structural features that a formal language needs to represent (in whatever way) to be considered an OWL language. In this work, we pick a concrete OWL syntax, namely the functional-style syntax, as a representative for OWL.⁵ So, an *OWL axiom* is an axiom as specified by OWL’s structural specification⁶ written in functional-style syntax. An *OWL ontology* is a finite set of OWL axioms. Please note that this notion of an OWL ontology does, strictly speaking, not correspond the definition of an OWL ontology in OWL’s structural specification.

3 Language Abstraction

We understand the notion of *abstraction* intuitively as a vehicle for analysing a subject of interest w.r.t. some aspects while disregarding others. As such, an abstraction involves a conceptualisation of the subject that is based on the removal of some level of detail. This removal of information is warranted if the removed details are immaterial for the subsequent analysis. Thus, we can characterise and describe an abstraction both in terms of what kind of information is preserved and what kind of information is discarded. We use this intuitive understanding of the notion of abstraction to motivate a (formal) notion for abstracting from a formal language into another.

3.1 Abstraction for Formal Languages

Following the intuition that the notion of abstraction is based on the removal of some information, we propose a notion of abstraction for formal languages that captures this intuition by syntax-directed transformations.

Definition 1 (Language Abstraction). *An abstraction for a language \mathcal{L} over finite alphabet Σ into a language \mathcal{L}_a over finite alphabet Σ_a is defined by a tree transducer $A = (Q, \Sigma, \Sigma_a, Q_f, \Delta)$ such that:*

- (i) *for all $t \in \mathcal{L}$ there exists a $A(t) \in \mathcal{L}_a$,*
- (ii) *there exist $t, t' \in \mathcal{L}$ s.t. $t \neq t'$ with $A(t) = A(t')$,*
- (iii) *for $t \in \mathcal{L}$ there exists a minor t_m that is isomorphic to $A(t)$.*

A language abstraction can be seen as a function from one formal language into another. Note that this function involves a loss of information since there are at least two distinct elements in the original language that are indistinguishable under the abstraction (see condition (ii) in Definition 1).

Also note that a language abstraction allows for the removal of two distinct kinds of information. On the one hand, an abstraction can remove structural information by transforming a tree following operations of graph minors. And on the other hand, lexical information can be removed by mapping to different symbols in the original language to the same symbol.

⁵ Note that the particular choice of a representative language for OWL is immaterial for the remainder of this work. In the following, we discuss OWL axioms exclusively on the basis of OWL’s structural specification.

⁶ <https://www.w3.org/TR/owl2-syntax/#Axioms>

Example 1 (Language Abstraction for OWL). Consider the axioms

$$\text{SubClassOf}(A_1, \text{ObjectMinCardinality}(2, S, A_2)) \text{ and} \\ \text{SubClassOf}(B_1, \text{ObjectSomeValuesFrom}(R, B_2)).$$

Clearly, these axioms differ w.r.t. their syntactical structure. However, they are similar in the sense that they both involve an object property restriction on the right-hand side of an `SubClassOf` axiom. This similarity can be made explicit by abstracting over the particular kind of object property restrictions. Using a language abstraction over OWL with the following transduction rules:⁷

- (i) $x \rightarrow q_{iri}(x)$ for $x \in \{A_1, A_2, B_1, B_2, S, R\}$
- (ii) $x \rightarrow q_N(x)$ for $x \in \mathbb{N}$
- (iii) $\text{ObjectSomeValuesFrom}(q_{iri}(x), q_{iri}(y)) \rightarrow q(\text{Restriction}(x, y))$
- (iv) $\text{ObjectMinCardinality}(q_N(x), q_{iri}(y), q_{iri}(z)) \rightarrow q(\text{Restriction}(y, z))$
- (v) $\text{SubClassOf}(q(x), q(y)) \rightarrow q_{final}(\text{SubClassOf}(x, y))$

the two axioms may be transformed into the (abstract) axioms

$$\text{SubClassOf}(A_1, \text{Restriction}(S, A_2)) \text{ and} \\ \text{SubClassOf}(B_1, \text{Restriction}(R, B_2)).$$

that only differ w.r.t. their signature (which could also be abstracted over, e.g., by the introduction of variables).

3.2 Abstraction Algebra

In this section, we briefly touch on the subject of relations between language abstractions. For example, two language abstractions can be related by some form of generalisation or specialisation or they can be completely orthogonal. To discuss such relations we propose the idea of an algebra for language abstractions. Note, however, that the following presentation is not intended to give a exhaustive or otherwise complete account of relationships for language abstractions.

Definition 2 (Abstraction by Ground Generalisation). An abstraction by ground generalisation w.r.t. a set of constants $C \subseteq \Sigma$ is a language abstraction $A = (Q, \Sigma, \Sigma \cup \mathcal{X}, Q_f, \Delta)$ s.t.

- (i) for $c \in C$ there exists a rule $c \rightarrow q(\phi(c)) \in \Delta$, where $\phi(c) = x$ maps a constant to a fresh variable, and
- (ii) for $s \notin C$ there exists a rule $f(h) \rightarrow q(\phi(h)) \in \Delta$, where $\phi(h) = f(h)$.

An Abstraction by Ground Generalisation removes lexical information by replacing constant symbols with variables.⁸ However, all of the syntactical structure pertaining to non-constant symbols is preseved.

⁷ We specify transduction rules without strictly following the correct formal notation. We refer the reader to the technical report for a fully worked-out and formally correct presentation.

⁸ One can argue that the introduction of different variables for the same constant symbol at different positions in a tree removes structural information. Indeed, one can easily define an abstraction by *renaming* that uniformly replaces a constant at different positions in a tree with the same variable.

A straightforward abstraction that does not preserve the structure of non-constant symbols, simply reduces the number of its children, i.e., immediate subterms.

Definition 3 (Abstraction by Projection). A language abstraction $A = (Q, \Sigma, \Sigma_a, Q_f, \Delta)$ is a projection for symbol $f \in \Sigma$ if there exists a rule $f(h) \rightarrow q(\phi(h)) \in \Delta$ s.t. $\phi(h) = f'(\phi(h))$ and $|\phi(h)| \leq |h|$.

A special case of an abstraction by projection is an abstraction that removes selected symbols together with all of their subterms.

Definition 4 (Abstraction by Selection). A language abstraction $A = (Q, \Sigma, \Sigma_a, Q_f, \Delta)$ is a selection for symbols $S \subseteq \Sigma$ if for all $f \notin S$ there exists a rule $f(h) \rightarrow q(\phi(f(h))) \in \Delta$ and $\phi(f(h)) = \varepsilon$.

4 Measuring Syntactic Regularity

We are interested in syntactic *regularities* for axioms in ontologies. We understand syntactic regularities in terms of shared syntactic properties between axioms. Such shared properties can be made precise and explicit by language abstractions as defined in Section 3. The basic idea being that shared syntactic properties are *preserved* under language abstractions while differences are *abstracted away*.

Definition 5 (Syntactic Regularity). Let A be a language abstraction over language \mathcal{L} . A syntactic regularity is an equivalence class $[t]_A = \{t' \in \mathcal{L} \mid A(t') = A(t)\}$ for a term $t \in \mathcal{L}$ induced by A .

Elements of a syntactic regularity are (syntactically) indistinguishable under the used abstraction. These elements necessarily share some syntactic aspects due to constraint (iii) in Definition 1 for language abstractions.

We extend the notion of syntactic regularities for elements of a language to *sets* of elements in a straightforward manner.

Definition 6 (Syntactic Regularity Sets). Let $\mathcal{S} = \{S_1, \dots, S_n\}$ be a family of sets over a language \mathcal{L} . A syntactic regularity set $[S_i]_A$ for $S_i \in \mathcal{S}$ w.r.t. a language abstraction A is defined by:

$$[S_i]_A = \{S' \in \mathcal{S} \mid \begin{array}{l} \text{there exist a sequence } A(t_1) \cdots A(t_n) \text{ with } \{t_1, \dots, t_n\} = S_i, \\ \text{there exist a sequence } A(t'_1) \cdots A(t'_n) \text{ with } \{t'_1, \dots, t'_n\} = S', \\ \text{s.t. } A(t_1) \cdots A(t_n) = A(t'_1) \cdots A(t'_n) \end{array}\}.$$

The amount of distinct syntactic regularities and their respective size shed light on whether the ontology is designed in a uniform manner. We refer to regularities of comparatively large size as *prevalent*.

Definition 7 (Prevalent Syntactic Regularity). A syntactic regularity $[t]_A$ in a language \mathcal{L} w.r.t. a language abstraction A is called *prevalent* if $\sup([t]_A) \geq \theta$, where $\sup: \mathcal{P}(\mathcal{O}) \rightarrow [0, 1]$ is a measure for subsets of \mathcal{L} and $\theta \in [0, 1]$ is a threshold parameter.

Note that an ontology, can contain both prevalent and non-prevalent syntactic regularities. Also note, that the existence of prevalent or non-prevalent syntactic regularities is not sufficient to consider an ontology as a whole syntactically regular or irregular. Rather, an ontology is syntactically regular, if it consists *to a large extent* of syntactic regularities. Otherwise, it is irregular.

Definition 8 (*k*-Regularity). A (finite) language \mathcal{L} is *k*-regular ($k \in \mathbb{N}$) w.r.t. a language abstraction A if there exist $t_1, \dots, t_k \in \mathcal{L}$ s.t. $\sup([t_1]_A, \dots, [t_k]_A) \geq \theta$, where $\sup: \mathcal{P}(\mathcal{O}) \rightarrow [0, 1]$ is a measure for subsets of \mathcal{L} and $\theta \in [0, 1]$ is a threshold parameter. Otherwise, \mathcal{O} is *k*-irregular.

Note that the two notions of prevalent syntactical regularity (cf. Definition 7) and *k*-regularity (cf. Definition 8) are independent from each other. While both definitions involve a measure function and a threshold parameter, we do not impose a formal relationship between them.

5 Methods

5.1 Research Questions

The notion of syntactic regularities introduced in Section 4 raises a number of research questions for OWL ontologies. Here, we distinguish between two broad categories of such questions. On the one hand, we are interested in the way OWL’s language constructors are used and composed in practice. On the other hand, we are interested in identifying reoccurring structural components beyond the use of language constructors that may reveal design decisions or common modelling practices in ontology engineering.

For the purpose of developing a first understanding of syntactic regularities in OWL ontologies, we focus on class expression axioms.⁹ This choice is motivated by the observation that OWL ontologies tend to contain large numbers of class expression axioms [6]. Also, we explore notions of syntactic regularity for *sets* of axioms on the basis of the widely-used notion of *class frames* that is used in both practice and theory as evidenced by popular ontology engineering tools and research [2,10,5].

5.2 Experimental Design

We use the following three language abstractions for class expression axioms for our empirical investigation: (i) a ground generalisation over an ontology’s class names, property names, and individual names, (ii) a ground generalisation over all constant symbols in OWL, (iii) a selection for class constructors. In the following, we refer to these three language abstractions as signature generalisation (SG), ground generalisation (GG), and class constructor preservation (CCP) respectively.

Using these three language abstractions, we conduct three distinct experiments. For each of these experiments, we distinguish between two experimental conditions: (a) syntactic regularities for *axioms* and (b) syntactic regularities for *sets of axioms*.

⁹ https://www.w3.org/TR/owl2-syntax/#Class_Expression_Axioms

While there are numerous ways sets of axioms over an ontology can be defined, we use the notion of class frames to induce a family of sets over an ontology. In particular, for a class C occurring in an ontology \mathcal{O} , we define its corresponding class frame $CF(C, \mathcal{O})$ by

$$CF(C, \mathcal{O}) = \{ \text{SubClassOf}(C, C') \in \mathcal{O} \} \cup \\ \{ \text{EquivalentClasses}(C, C'_1, \dots, C'_n) \in \mathcal{O} \} \cup \\ \{ \text{DisjointClasses}(C, C'_1, \dots, C'_n) \in \mathcal{O} \} \cup \\ \{ \text{DisjointUnion}(C, C'_1, \dots, C'_n) \in \mathcal{O} \}$$

With this, we specify the three aforementioned experiments as follows.

Syntactic Diverseness We determine the syntactic diverseness w.r.t. to a given language abstraction in an ontology by counting the number of syntactic regularities embodied in an ontology. A large number of syntactic regularities suggests that an ontology is syntactically diverse while a small number indicates syntactical homogeneity.

Prevalent Regularities We investigate the extent to which syntactic regularities in ontologies are prevalent. We measure the prevalence of a syntactic regularity with the measuring functions $\text{sup}(x) \mapsto \frac{|x|}{|\Omega|}$ in condition (a) for axioms and $\text{sup}(x) \mapsto \frac{|\{\alpha \in \Omega \mid \exists f \in x : \alpha \in f\}|}{|\Omega|}$ in condition (b) for class frames. Here, Ω denotes the set of all class expression axioms.

We decide whether a given syntactic regularity is prevalent or not by using the threshold parameter $\theta = 0.1$. So intuitively, a syntactic regularity is *prevalent* if it accounts for at least 10% of all class expression axioms in a given ontology. We report on the number of prevalent syntactic regularities and give an account of their syntactic properties.

Ontology Coverage with Regularities We analyse an ontology's k -regularity w.r.t. a threshold parameter $\theta = 0.9$ and use the measuring functions and $\text{sup}(x_1, \dots, x_k) \mapsto \frac{|\bigcup_k x|}{|C|}$ for condition (a), and $\text{sup}(x_1, \dots, x_k) \mapsto \frac{|\bigcup(\bigcup_k x)|}{|\Omega|}$ for condition (b). For condition (a) we determine the smallest k and in condition (b) we approximate k by iteratively considering the k -most prevalent class frame regularities.

5.3 Ontology Corpus

We work with a recent (June 2020) snapshot of BioPortal created in the same way¹⁰ as described in [7]. The data set of ontologies with their imports closure merged in encompasses a total of 622 ontologies. In experiments, we distinguish between three kinds of ontologies. First, ontologies that consist of atomic axioms only, i.e., `SubClassOf` and `EquivalentClasses` axioms that have only named classes as arguments. Second, ontologies expressible in \mathcal{EL}^{++} . And third, ontologies not expressible in \mathcal{EL}^{++} . We refer to these three kinds of ontologies as *atomic*, \mathcal{EL}^{++} , and *rich* ontologies respectively. Figure 1 shows the size of an ontology's TBox as well as the size of its subset of class expression axioms.

¹⁰ <https://github.com/matentzn/bioportal.download>

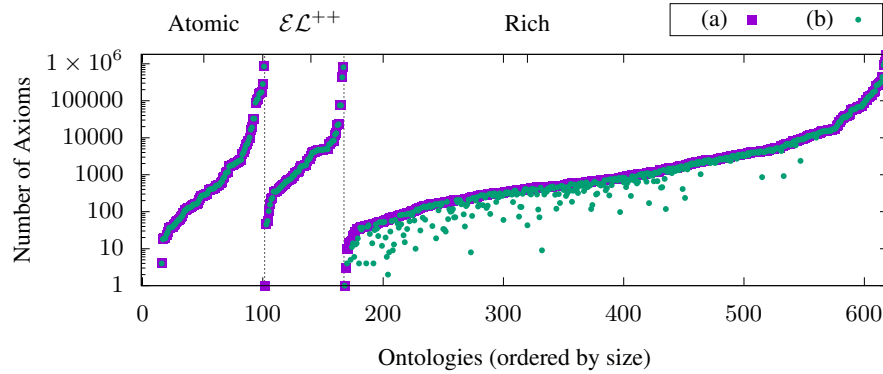


Fig. 1: Number of TBox axioms (a) and Class Expression Axioms (b).

6 Results

We present results for the three experiments as specified in Section 5.2 in distinct subsections. Each subsection is further subdivided according to the two experimental conditions for syntactic regularities for (a) axioms and (b) sets of axioms. In each case, we report on findings w.r.t. to the three language abstractions (i)–(iii).

6.1 Experiment 1: Syntactic Diverseness

(a) Axioms Figure 2 shows the (absolute) counts of syntactic regularities for axioms in BioPortal ontologies. We observe that the number of regularities (positively) correlates with both the used language profile (atomic, \mathcal{EL}^{++} , and rich) as well as the size of an ontology. However, most ontologies contain at most 100 syntactic regularities for all three language abstractions. In fact, it appears that the language abstractions, SG, GG, and CCP often coincide in terms of the number of syntactic regularities. This happens, for example, when an ontology does not contain axioms that the three language abstractions treat differently or if syntactic features that are treated differently are used in a homogeneous manner.

Yet, we also observe ontologies with large difference between the number of syntactic regularities for different language abstractions. For example, the Orphanet Rare Disease Ontology at index 166, contains a large number of syntactic regularities w.r.t. SG. This is due to the use of literal value restrictions, e.g., `DataHasValue(Orphanet_C029 "1.9"^^xsd:float)`. Since SG does not abstract over literals, axioms with different literals, e.g., different (float) numbers, belong to different syntactic regularities. The number of regularities w.r.t. GG is comparatively small because GG abstracts over differences between literal values. For example, the literal value restriction for Orphanet_C029 and `DataHasValue(Orphanet_C024 "50.0"^^xsd:float)` belong to different syntactic regularities w.r.t SG but to the same w.r.t. GG.

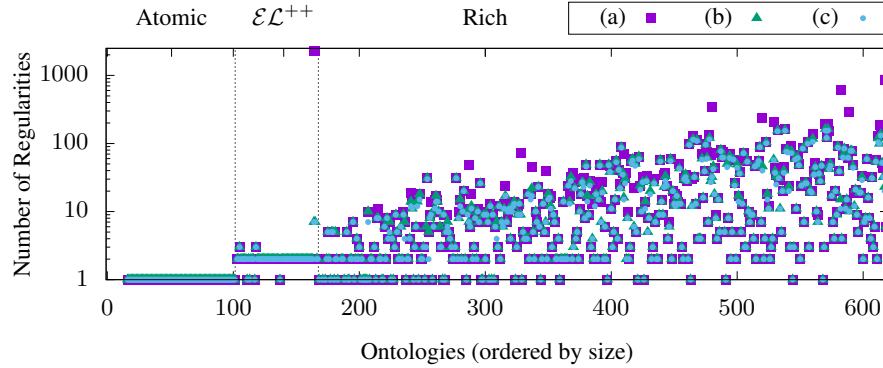


Fig. 2: Syntactic Diverseness for axioms w.r.t. (a) SG, (b) GG, and (c) CCP.

(b) Class Frames Figure 3 shows the (absolute) counts of syntactic regularities for class frames in BioPortal ontologies. We make the same observations as for syntactic regularities for axioms. Namely, the number of regularities correlates positively with both an ontology's language profile as well as its size. Also, the three language abstractions SG, GG, and CCP generally give rise to the same number of syntactic regularities. However, we need to point out the overall increased number of regularities compared to the case of (a) axioms. This is due to the fact that axioms of the same syntactic regularity can be combined in various ways to define different syntactic regularities for class frames. For example, consider the ontology

$$\mathcal{O} = \{\text{SubClassOf}(A, B), \text{SubClassOf}(A, B), \text{SubClassOf}(C, D)\}.$$

Then, there is only one syntactic regularity for axioms w.r.t. GG, but two for class frames.

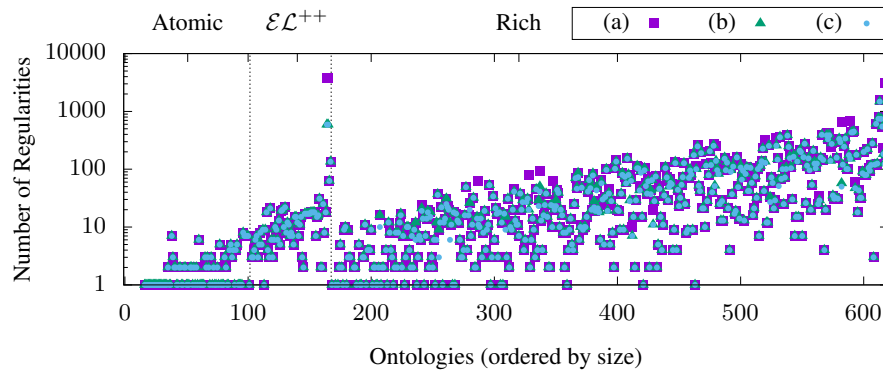


Fig. 3: Syntactic Diverseness for class frames w.r.t. (a) SG, (b) GG, and (c) CCP.

6.2 Experiment 2: Prevalent Regularities

(a) Axioms In light of only one or two syntactic regularities for almost all atomic and \mathcal{EL}^{++} ontologies (see Figure 2), it comes to no surprise that these regularities are also often prevalent, i.e., they make up at least 10% of the class expression axioms in an ontology. There is only one \mathcal{EL}^{++} with more than two prevalent regularities, namely the Orphanet Rare Disease Ontology with four regularities w.r.t. GG and CCP (recall that SG gives rise to a large number of syntactic regularities due to literal value restrictions. However, these regularities contain relatively few axioms).

Similarly, the three language abstractions give rise to at most two prevalent regularities for rich ontologies as well. This is also to be expected as (atomic) axioms from the class hierarchy often make up large a proportion of an ontology’s TBox [4]. Yet, about 15% of rich ontologies (68 out of 451) exhibit more than two prevalent regularities w.r.t. at least one of the three used language abstractions; but only about 4% (17 out of 451) exhibit more than three prevalent regularities.

Otherwise, there is no ontology with more than four prevalent regularities with the exception of only one, namely the “Neomark Oral Cancer Ontology (version 3).” Note that this ontology only contains 236 class expression axioms. So any syntactic regularity with more than 23 axioms is considered prevalent. In fact, rich ontologies with more than three prevalent regularities tend to have a “small” TBox. In particular, there are only four ontologies with at least 1000 TBox axioms that exhibit more than 3 prevalent regularities.

(b) Class Frames For class frames, we find slightly more prevalent syntactic regularities compared to regularities for (single) axioms. More than 90% ontologies exhibit at most three prevalent regularities for all language abstractions (573 out of 619). Yet, there are a few ontologies in which up to 8 regularities are prevalent. This is in line with the overall larger number of syntactic regularities for class frames (cf. Section 6.1). However, it is important to keep in mind that our notion of class frames does not induce a partition of an ontology; neither do syntactic regularities for class frames. In fact, both `EquivalentClasses` and `DisjointClasses` axioms can be part of multiple class frames associated with different classes. This can (in theory) lead to a large number of prevalent regularities. Consider the following example ontology:

$$\begin{aligned} \mathcal{O} = \{ & \text{DisjointClasses}(A, B), \text{DisjointClasses}(B, C), \\ & \text{DisjointClasses}(A, C), \text{DisjointClasses}(B, D), \\ & \text{DisjointClasses}(A, D), \text{DisjointClasses}(C, D), \\ & \text{SubClassOf}(A, \exists R.X) \quad \text{SubClassOf}(B, X \sqcap Y), \\ & \text{SubClassOf}(C, \forall R.X) \quad \text{SubClassOf}(D, X \sqcup Y) \} \end{aligned}$$

Here, the four classes A, B, C, and D all give rise to a class frame involving four axioms. Since all class frames belong to different regularities (due to the `SubClassOf` axioms) and \mathcal{O} only contains 10 axioms, all four regularities would be considered prevalent w.r.t. a threshold of 40%. However, if the class frames would induce a partition of \mathcal{O} , then there should be at most two syntactic regularities that cover 40% of class expressions in \mathcal{O} .

6.3 Experiment 3: Ontology Coverage

(a) Axioms Most ontologies contain only one or two prevalent regularities (cf. Section 6.1). Yet, the total number of syntactic regularities is larger than 10 for many rich ontologies and even larger 100 in some cases (cf. Section 6.1). In Figure 4 we show the number of syntactic regularities that need to be merged to cover 90% of class expression axioms in an ontology. We observe that two regularities are often sufficient. Furthermore, we note that cases in which more than 10 regularities are needed are rare.

This suggests that ontologies, for the most part, are built on the basis of (syntactically) homogeneous set of axioms. As argued before, this observation can (to large extents) be attributed to the predominant role of atomic subsumption axioms that are used to represent an ontology’s class hierarchy. Also, it is important to keep in mind that most entities in OWL are represented by a *combination* of a set of axioms. Therefore, it is not warranted to assume an ontology to follow an overall homogeneous design only because it is predominantly built on the basis of axioms of a certain form.

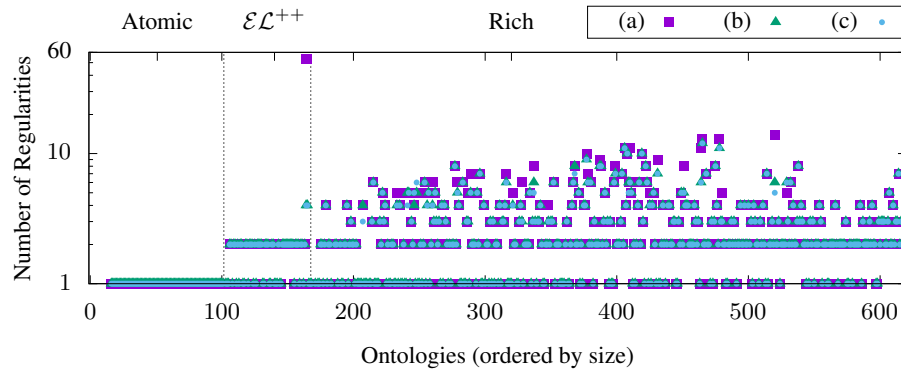


Fig. 4: Number of axiom regularities needed to cover 90% of an ontology’s class expression axioms w.r.t. (a) SG, (b) GG, and (c) CCP.

To gain a better understanding of an ontology’s composition w.r.t. syntactic regularities, we inspect what proportion of an ontology’s class expression axioms are covered by the three most prevalent regularities w.r.t. SG. We only discuss SG because GG and CCP abstract over a superset of syntactical features compared to SG and thus give rise to fewer regularities (of bigger size). While 4 already shows that many ontologies’ class expression axioms can be covered to 90% by only three regularities, Figure 5 shows that 50% of class expression axioms can be covered by the three most prevalent regularities in all ontologies. There are only 17 ontologies whose class expression axioms are covered to less than 70% by the three most prevalent regularities.

Note that some ontologies contain no class expression axioms. Such ontologies are represented in Figure 5 with a zero percentage and we excluded them for the above counts as it is a question of taste whether an empty set is covered to 0% or 100%.

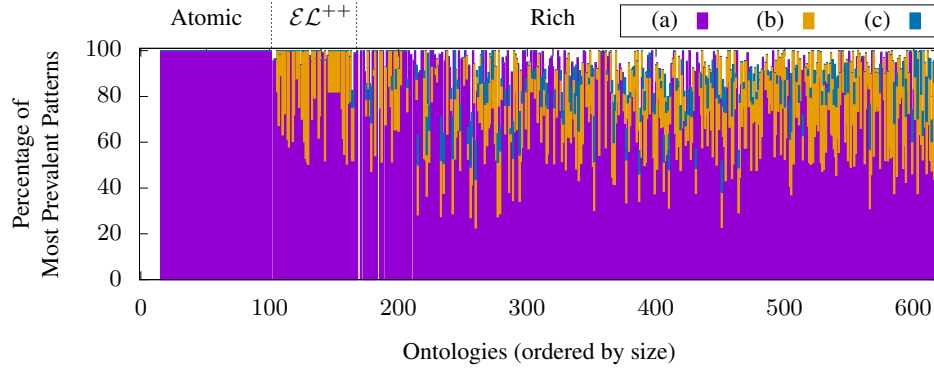


Fig. 5: Three most prevalent syntactic regularities w.r.t. SG: (a) depicts the most prevalent, (b) the second most prevalent, and (c) the third most prevalent.

(b) Class Frames For class frames, we find that class expression axioms in atomic and \mathcal{EL}^{++} ontologies can also be covered to 90% by just 10 regularities (with a few exceptions). However, for rich ontologies we observe that the number of needed regularities increases with the size of an ontology and often exceeds 10 considerably.

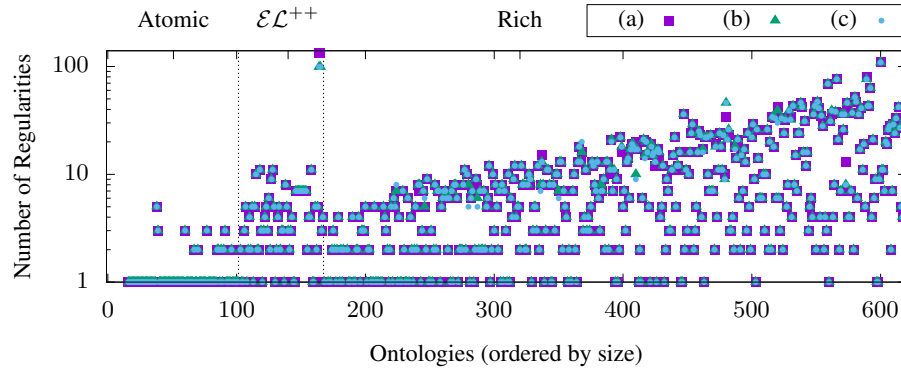


Fig. 6: Number of class frame regularities needed to cover 90% of an ontology's class expression axioms w.r.t. (a) SG, (b) GG, and (c) CCP.

As before, we inspect the three most prevalent syntactic regularities w.r.t. SG and determine what proportion of an ontology's class expression axioms they already cover. Figure 7 reveals that the three most prevalent regularities for class frames are not as prominent as for axioms. While there are ontologies for which the three most preva-

lent regularities are sufficient to cover 80% of their class expressions axioms, there are equally as many ontologies for which this does not hold.

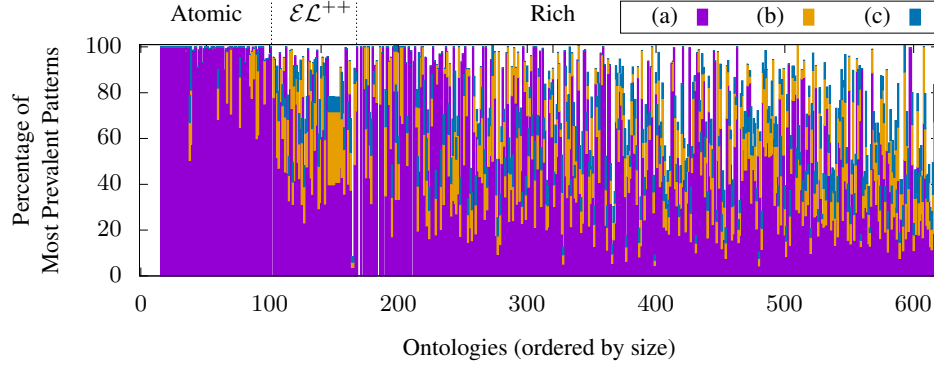


Fig. 7: Three most prevalent syntactic regularities w.r.t. SG. (a) depicts the most prevalent, (b) the second most prevalent, and (c) the third most prevalent.

7 Related Work

A range of approaches for discovering syntax-based regularities in OWL ontologies have been proposed. First, agglomerative clustering has been used to identify commonalities between similar entities [8]. The used notion of a *replacement function* can be interpreted as a language abstraction. However, its use in the similarity distance measure effectively reduces a set of axioms to a set of *types* of axioms. While the resulting loss of information can be seen as a form of abstraction, it is not directly comparable to the language abstractions motivated in this work.

Second, frequent subtree mining over OWL axioms has been motivated by interpreting them as (syntax) trees [5]. Furthermore, by using frequent itemset mining over identified regularities for axioms, regularities for sets of axioms are motivated. While frequent subtree mining aims to identify frequent structures, syntactic regularities based on language abstractions *do not*. Language abstractions are primarily concerned with syntactic properties and the corresponding notion for regularities is independent from any notion of frequency. The qualification of regularities in terms of being prevalent, i.e., capturing frequently occurring syntactic structures, is a separate level of analysis.

Third, lexical naming conventions for OWL classes in combination with structural relations between them in an ontology’s class hierarchy have been suggested to provide useful information w.r.t. an ontology’s underlying design [12].

Besides, approaches for discovering regularities from a given ontology alone, there also exists target-oriented approaches that aim to detect some predefined notion of patterns [3,9,11]. Otherwise, there are surveys of ontologies discussing syntactic properties such as the prevalence of, for example, logical constructors [6,13].

8 Conclusion

We have proposed a formal framework for identifying and characterising syntactic regularities in ontologies. We used this framework to survey a large corpus of actively maintained ontologies and found prevalent regularities that may prove to be useful in terms of ontology comprehension and maintenance. While the used language abstractions in our empirical survey are somewhat coarse-grained, they can easily be refined or otherwise modified as needed. Overall, this work is motivated by the need for data-driven methods to extract and characterise common or recurring modelling practices in published ontologies. The hope is that such an automated information extraction will help with the development of high-quality design patterns both in theory and practice.

References

1. Comon, H.: Tree automata techniques and applications. <http://www.grappa.univ-lille3.fr/tata> (1997)
2. Horridge, M., Patel-Schneider, P.F.: Manchester syntax for OWL 1.1. In: OWLED (Spring). CEUR Workshop Proceedings, vol. 496. CEUR-WS.org (2008)
3. Kindermann, C., Parsia, B., Sattler, U.: Detecting influences of ontology design patterns in biomedical ontologies. In: ISWC (1). Lecture Notes in Computer Science, vol. 11778, pp. 311–328. Springer (2019)
4. Kindermann, C., Parsia, B., Sattler, U.: Prevalence and effects of class hierarchy precompilation in biomedical ontologies. In: International Semantic Web Conference (1) (2020)
5. Lawrynowicz, A., Potoniec, J., Robaczyk, M., Tudorache, T.: Discovery of emerging design patterns in ontologies using tree mining. *Semantic Web* 9(4), 517–544 (2018)
6. Matentzoglou, N., Bail, S., Parsia, B.: A snapshot of the OWL web. In: International Semantic Web Conference (1). Lecture Notes in Computer Science, vol. 8218, pp. 331–346. Springer (2013)
7. Matentzoglou, N., Parsia, B.: Biportal snapshot 30.03.2017 (Mar 2017). <https://doi.org/10.5281/zenodo.439510>, <https://doi.org/10.5281/zenodo.439510>
8. Mikroyannidi, E., Manaf, N.A.A., Iannone, L., Stevens, R.: Analysing syntactic regularities in ontologies. In: Klinov, P., Horridge, M. (eds.) Proceedings of OWL: Experiences and Directions Workshop 2012, Heraklion, Crete, Greece, May 27–28, 2012. CEUR Workshop Proceedings, vol. 849. CEUR-WS.org (2012), http://ceur-ws.org/Vol-849/paper_11.pdf
9. Mortensen, J., Horridge, M., Musen, M.A., Noy, N.F.: Modest use of ontology design patterns in a repository of biomedical ontologies. In: WOP. CEUR Workshop Proceedings, vol. 929. CEUR-WS.org (2012)
10. Noy, N.F., Musen, M.A., Jr., J.L.V.M., Rosse, C.: Pushing the envelope: challenges in a frame-based representation of human anatomy. *Data Knowl. Eng.* 48(3), 335–359 (2004)
11. Sváb-Zamazal, O., Scharffe, F., Svátek, V.: Preliminary results of logical ontology pattern detection using SPARQL and lexical heuristics. In: WOP. CEUR Workshop Proceedings, vol. 516. CEUR-WS.org (2009)
12. Sváb-Zamazal, O., Svátek, V.: Analysing ontological structures through name pattern tracking. In: EKAW. Lecture Notes in Computer Science, vol. 5268, pp. 213–228. Springer (2008)
13. Wang, T.D., Parsia, B., Hendler, J.A.: A survey of the web ontology landscape. In: International Semantic Web Conference. Lecture Notes in Computer Science, vol. 4273, pp. 682–694. Springer (2006)