

Designing Efficient Controlled Languages for Ontologies

Camilo Thorne, Raffaella Bernardi, and Diego Calvanese

1 Introduction

The attempts made in the 70s and 80s to build natural language interfaces (NLIs) to information systems and databases turned into disappointments towards the 90s (Androutsopoulos et al., 1995). One of the reasons were the challenges posed by structural and semantic ambiguity in arbitrary natural language input. As a way to overcome the ambiguity problem, the *controlled natural language* (CL) paradigm was proposed (Huijsen, 1998; Kittredge, 2003), to build NLIs where only a restricted fragment of a natural language can be used. An important area of application of CLs is to provide front-ends to ontologies and ontology-based systems. In this setting, CLs allow the systems to parse efficiently user statements and questions. It is less clear however whether they can be understood as efficiently, in particular by ontology-based systems that need to *reason* over the semantic representations of user inputs. The present chapter intends to study the semantic complexity of CLs, together with the conditions under which reasoning with a CL can scale to very large ontologies and ontology-based systems.

By an *ontology* we mean here a conceptualization of a domain of interest, expressed as a set of logical assertions. Specifically, ontologies formulated in variants of description logics (DLs), which are fragments of first-order logics with well understood computational properties and for which logical reasoning (e.g., to detect inconsistencies in a specification) is decidable and in significant cases also computationally tractable (Baader et al., 2003). DLs provide the formal underpinning for

Camilo Thorne

Free University of Bozen-Bolzano, Bolzano, Italy, e-mail: cthorne@inf.unibz.it

Raffaella Bernardi

University of Trento, Trento, Italy, e-mail: bernardi@disi.unitn.it

Diego Calvanese

Free University of Bozen-Bolzano, Bolzano, Italy, e-mail: calvanese@inf.unibz.it

the Web Ontology Language OWL (Horrocks et al., 2003), which is the ontology specification language standardized by the W3C.¹

The present chapter specifically addresses the questions of (1) *which* should be the CL to be used to manage ontologies efficiently, and (2) *how* can it be defined. Concretely, our proposal is to determine a methodology for defining *exactly* those fragments with a *desirable computational complexity*. We use DLs as the starting point to answer Question (1), viz. which is the most suitable NL fragment, and we use categorial grammars (CGs) to provide an answer to Question (2), viz. how to capture the syntactic structures corresponding exactly to the semantic representations allowed by the chosen, efficient DL.

With respect to the kind of DL, we focus our attention on *DL-Lite*, which is a family of DLs studied in the context of ontology-based access to (relational) databases (Calvanese et al., 2007, 2011). When considering the well-known trade-off between expressive power and computational complexity of inference, *DL-Lite* is specifically optimized for efficient reasoning also in the presence of large datasets, taking into account that in ontology-based systems the size of the data (stored in relational databases or in possibly very large triple stores) largely dominates the size of the ontology's intensional descriptions. Indeed, in *DL-Lite*, reasoning is computationally tractable in general, and can actually be carried out by exploiting the query answering functionalities of the data storage layer. This can be contrasted with the computational properties of more expressive DLs, such as *SHOIN*, the DL underlying OWL, in which reasoning is computationally intractable also when the complexity is measured with respect to the size of the data only. With respect to the CG, we define a grammar that relies on the sub-categorization of syntactic constituents to capture exactly the intended logic.

We exploit the syntax-semantics interface as realized by CGs to obtain *DL-Lite* meaning representations compositionally while parsing (van Benthem, 1987; Moortgat, 1997). To this end we consider of particular value the studies carried out by Pratt and Third (2006), who have investigated the satisfiability of sets of sentences in fragments of natural language and their computational complexity, but start instead from the logic (viz., an OWL fragment) as do Kaljurand and Fuchs (2006).

The rest of the chapter is structured as follows. In Section 2, we provide an overview of controlled languages and semantic complexity, highlighting the open questions that motivate our contributions. In Sections 3 and 4, we introduce respectively the DL and the grammar we work with. In Section 5, we describe in detail how CGs can capture exactly the desired fragments of natural language. In Section 6, we show how corpora analysis can be used to justify further CL design choices. In Section 7, we provide an overview of related work, in the form of related results on CLs obtained and published elsewhere by the authors, and in the form of other CLs for ontologies that have been proposed in the literature. Finally, in Section 8, we summarize our results and outline our ongoing work regarding the computational properties of controlled languages, both declarative and interrogative.

¹ <http://www.w3.org/TR/owl2-primer/>

2 Controlled Languages and Semantic Complexity

A *controlled language* (CL) is a fragment of natural language such as English with a limited lexicon and a small set of grammar rules (Huijsen, 1998; Kittredge, 2003). Importantly, CLs are engineered to handle natural language ambiguity, so that their utterances “compile”, via, e.g., a rule-based, symbolic and compositional syntax-directed translation algorithm (in a way similar to programming languages’ compilation), in *unambiguous* logical axioms and/or queries, due to their restricted syntax and lexicon.

This tight integration with formal ontology and query languages gives rise to a more general phenomenon: the property of *semantic complexity* as defined and investigated by Pratt and Third (2006). They show that each (controlled) fragment of English generates a logic fragment: the set of its *meaning representations* (MRs); semantic complexity is then naturally defined as the computational complexity of reasoning with its MRs (i.e., the computational complexity of the associated satisfiability problem). Furthermore, they show that semantic complexity correlates with coverage by considering the impact that particular combinations of English constructs (negation, relatives, transitive verbs, etc.) have on semantic complexity, pinpointing combinations that are: (i) tractable (**PTime** semantic complexity), (ii) intractable (**NP-hard** semantic complexity), or (iii) undecidable.

In our work, we extend both their methodology and their results: their methodology, by using categorial grammars to “reverse engineer” English controlled fragments from logics that exhibit desirable computational properties; their results, by considering semantic *data complexity*, viz., the semantic complexity of CLs for ontologies measured only in the size of the (typically very large) data repositories they are meant to manage as opposed to the size of the complete logical specification derived from the natural language utterances. More precisely, we (i) consider logic constructs that give rise to tractable data complexity, (ii) pinpoint those structures of English that map (following Montagovian semantics) into those logic constructs, and (iii) propose grammars that generate such structures. In this way one can determine the best trade-off between coverage and tractability holding for NLI systems on ontologies and ontology-based systems.

Pratt and Third’s Fragments of English. The work of Pratt and Third provides hints on how to determine which fragments hold the right expressiveness for ontology-based systems, via their notion of semantic complexity. We give now a brief overview of their *controlled fragments of English* (cf., Pratt and Third, 2006), which are subsets of standard English meant to capture some simple, albeit for our purpose important, structures of English.

The fragments of Pratt and Third are built incrementally, starting with copula, nouns, negation, and the universal and existential quantifiers, and extending later coverage to larger portions of English – relative constructions, ditransitive verbs, and anaphora, as summarized in Table 1. The fragments are named after such combinations, COP if their sentences contain only the copula, COP+TV if they contain in addition transitive verbs, and COP+TV+DTV if they contain both transitive and

Table 1 Fragments of English studied by Pratt and Third (2006).

| Fragment | Coverage | Semantic Complexity |
|----------------|---|---------------------------|
| COP | Copula, common, proper nouns, negation, universal and existential quantifiers | PTime |
| COP+TV+DTV | COP+transitive verbs (“reads”) + ditransitive verbs (“gives”) | PTime |
| COP+Rel | COP+relative pronouns (“who”, “that”, “which”) | NP -complete |
| COP+Rel+TV | COP+Rel+transitive verbs | ExpTime -complete |
| COP+Rel+TV+DTV | COP+Rel+TV+ditransitive verbs | NExpTime -complete |
| COP+Rel+TV+RA | COP+Rel+TV+restr. anaphora (“him”, “she”, “itself” with <i>bounded</i> anaphoric co-references) | NExpTime -complete |
| COP+Rel+TV+GA | COP+Rel+TV+gen. anaphora (<i>unbounded</i> anaphoric pronouns) | undecidable |

ditransitive verbs. Further differences are due to the presence in the lexicon of the relative pronoun (Rel) and of anaphora in a general (GA) or restricted form (RA).

Each NL construct has a MR introducing an n -ary predicate or a logical operation in First Order Logic (FO): The MRs of relatives (e.g., “who”) introduce conjunction (\wedge); negations (e.g., “no”, “not”) introduce logical negation (\neg); intransitive verbs (e.g., “runs”) and nouns (e.g., “man”) correspond to unary predicates; transitive verbs (e.g., “loves”) correspond to binary predicates, and ditransitive verbs (e.g., “sells to”) to ternary predicates; universal quantifiers (“every”, “all”, “everyone”) to universal quantification (\forall), and existentials (“some”, “someone”) to existential quantification (\exists).

Example 1. COP and COP+TV+DTV generate English utterances such as:

- (1) Some people are weak.

$$[\exists x(\text{People}(x) \wedge \text{Weak}(x))]$$

- (2) Every husband has a wife.

$$[\forall x(\text{Husband}(x) \rightarrow \exists y(\text{Wife}(y) \wedge \text{Has}(x, y)))]$$

- (3) Every salesman sells some item to some customer.

$$[\forall x(\text{Salesman}(x) \rightarrow \exists y(\text{Customer}(y) \wedge \exists z(\text{Item}(z) \wedge \text{Sells}(x, z, y)))]$$

Note that in (2) and (3) above, other translations might be possible due to NL ambiguity. However, these are discarded by the grammar, which follows only the surface order of constituents. ♣

Boolean- and non-Boolean-closed fragments. As shown in Table 1, where we report the results of Pratt and Third (2006), the most expressive fragment of English they consider is undecidable. As a matter of fact, only the first two fragments,

COP and COP+TV+DTV, are tractable, i.e. have **PTime** semantic complexity. Notice that as soon as we add rules dealing with the relative clause we lose tractability. COP+Rel (i.e., COP with relative clauses) is already **NP**-hard. This is because relatives express conjunctions which, together with negation, generate logics (i.e., fragments of FO) that contain the propositional calculus (for which reasoning is **NP**-complete). In other words, COP+Rel and all the fragments containing it are “Boolean-closed”, and allow negation to be freely combined with conjunction and relatives. Instead, COP and COP+TV+DTV are “non-Boolean-closed”. The challenge that we face here is to develop a methodology for defining tractable, “non-Boolean-closed” CLs that capture tractable ontology languages.

3 *DL-Lite* and its Computational Properties

Description logics (DLs) (Baader et al., 2003) are the logics, typically fragments of FO, that provide the formal underpinning to ontologies and the Semantic Web (Horrocks et al., 2003). They allow one to structure the domain of interest by means of *concepts*, denoting sets of objects, and *roles*, denoting binary relations between (instances of) concepts. Complex concept and role expressions are constructed starting from a set of atomic concepts and roles by applying suitable constructs. The domain of interest is then represented by means of a DL knowledge base, consisting of a TBox (for “terminological box”), storing intensional information, and an ABox (for “assertional box”), storing extensional information about individual objects of the domain of interest.

We focus our attention on *DL-Lite* (Calvanese et al., 2007, 2011), a family of DLs specifically tailored to manage large amounts of data efficiently. Specifically, we consider variants of *DL-Lite* in which the TBox is constituted by a set of *inclusion assertions* of the form

$$Cl \sqsubseteq Cr$$

where *Cl* and *Cr* denote concepts that may occur respectively on the left and right-hand side of inclusion assertions. The form of such concepts depends on the specific variant of *DL-Lite*. Here, we consider two variants, called *DL-Lite_{core}* and *DL-Lite_{R,□}*, which we define below. In fact, *DL-Lite_{core}* represents a core part shared by all logics of the *DL-Lite* family.

Definition 1 (*DL-Lite_{core}* and *DL-Lite_{R,□}*). In *DL-Lite_{core}*, *Cl* and *Cr* are defined as follows:²

$$Cl \longrightarrow B \mid \exists R \qquad Cr \longrightarrow B \mid \neg B \mid \exists R \mid \neg \exists R$$

where *B* denotes an atomic concept, and *R* denotes an atomic role. In *DL-Lite_{R,□}*, in addition to the clauses of *DL-Lite_{core}*, we have also:

² We have omitted *inverse* roles R^- from the DLs to simplify the presentation of the main idea we are investigating.

$$\begin{array}{ll}
MScStudent \sqsubseteq Student & Student \sqcap Busy \sqsubseteq Works \\
MScStudent \sqsubseteq Works & Student \sqcap \exists Reads \sqsubseteq Works \\
MScStudent \sqsubseteq \neg BScStudent & \exists Reads \sqcap \exists Writes \sqsubseteq Works \\
\exists Reads \sqsubseteq Works & Student \sqsubseteq \exists Reads.Book \\
Student \sqsubseteq \exists Reads &
\end{array}$$

Fig. 1 An example $DL-Lite_{core}$ TBox (left part), and some additional $DL-Lite_{R,\sqcap}$ assertions (right part).

$$Cl \longrightarrow Cl_1 \sqcap Cl_2 \qquad Cr \longrightarrow \exists R.B$$

where R denotes again an atomic role. ♠

The \sqcap construct denotes conjunction, and \neg negation (or complement). The $\exists R$ construct is called *unqualified existential quantification*, and intuitively denotes the *domain* of role R , i.e., the set of objects that are connected through role R to some (not further specified) object.³ Finally, the $\exists R.Cr$ construct, called *qualified existential quantification*, allows one to further qualify the object connected through role R as an instance of concept Cr .

As an example, consider the $DL-Lite_{core}$ TBox depicted in the left part of Figure 1, which makes use of various concepts ($Student$, $MScStudent$, $BScStudent$, $Works$) and roles ($Reads$, $Writes$) to express some simple knowledge about the student domain. Specifically, the TBox assertions state that every MSc-student is a student, that MSc-students work, and that no MSc-student is a BSc-student, i.e., the two concepts are *disjoint*. Note that in $DL-Lite$, negation is used only to express disjointness, as in the statement in Figure 1. Additionally, making use of unqualified existential quantification, we can express that everyone who reads something (i.e., is in the domain of the $Reads$ role) works, and that every student reads something. The latter is also called a *participation constraint*, since it forces instances of $Student$ to participate in the $Reads$ role. In the right part of Figure 1, we have shown also some $DL-Lite_{R,\sqcap}$ inclusion assertions, which make use of conjunction in the left-hand side to express that busy students work, that students who read something work, and that everyone who reads something and writes something works. Finally, to express that every student reads some book, we can make use of qualified existential quantification (allowed to appear only in the right-hand side of inclusion assertions).

To formally specify the semantics of $DL-Lite$, we provide its standard translation to FO. Specifically, we map each concept C (we use C to denote an arbitrary concept, constructed applying the rules above) to a FO formula $\varphi(C, x)$ with one free variable x (i.e., a unary formula), and each role R to a binary formula $\varphi(R, x, y)$ as follows:

³ Instead, $\exists R^-$, for an inverse role R^- , denotes the *range* of role R .

Table 2 Combined complexity and data complexity of consistency in different DLs.

| DL | Combined Complexity | Data Complexity |
|-------------------------------|--------------------------|-----------------------|
| <i>DL-Lite_{core}</i> | in NLogSpace | AC⁰ |
| <i>DL-Lite_{R,∩}</i> | PTime -complete | AC⁰ |
| <i>ALC</i> | ExpTime -complete | coNP -complete |
| <i>SHOIN</i> | NExpTime -hard | coNP -hard |

$$\begin{aligned}
\varphi(B, x) &= \mathbb{B}(x) & \varphi(R, x, y) &= \mathbb{R}(x, y) \\
\varphi(\neg C, x) &= \neg \varphi(C, x) & \varphi(\exists R, x) &= \exists y \varphi(R, x, y) \\
\varphi(C_1 \sqcap C_2, x) &= \varphi(C_1, x) \wedge \varphi(C_2, x) & \varphi(\exists R.C, x) &= \exists y \varphi(R, x, y) \wedge \varphi(C, y)
\end{aligned}$$

In the translation of $\exists R.C$, the variable y is considered to be a fresh variable. An inclusion assertion $Cl \sqsubseteq Cr$ of the TBox corresponds then to the universally quantified FO sentence $\forall x. \varphi(Cl, x) \rightarrow \varphi(Cr, x)$.

We observe that the above translation actually generates a formula in the guarded fragment of FO. This holds not only for *DL-Lite* but for many other expressive DLs as well, and accounts for the good computational properties of such logics (Baader et al., 2003).

Finally, in *DL-Lite*, an ABox is constituted by a set of *assertions on individuals*, of the form $B(a)$ or $R(a, b)$, where B and R denote respectively an atomic concept and role, and a, b denote constants. As in FO, each constant is interpreted as an element of the interpretation domain. The above ABox assertions correspond to the analogous FO facts, or, by resorting to the above mapping, to $\varphi(B, x)(a)$ and $\varphi(R, x, y)(a, b)$, respectively. A *DL-Lite knowledge base* is simply a pair $(Tbox, Abox)$, where $Tbox$ is a TBox and $Abox$ an ABox. A *model* of such knowledge base is a FO interpretation in which the (closed) FO formulae resulting from the translation of all assertions in $Tbox \cup Abox$ evaluates to true.

To study efficiency we consider the computational reasoning problems relevant to DL ontologies and knowledge bases. The key problem, to which most other ones can be reduced, is the problem of *knowledge base consistency*, in which, given a knowledge base $(Tbox, Abox)$, we ask whether it has a model. Following Vardi (1982), when we consider the computational complexity measured only in terms of the *size* of the ABox (defined as the number of constants the ABox contains), we speak about *data complexity*. When instead the complexity is measured in terms of the size of the whole input, we speak of *combined complexity*. A DL can be considered as “efficient” for ontology-based data management, whenever such complexity is *tractable* (in **PTime**).

It turns out that *DL-Lite*, and in particular *DL-Lite_{core}* and *DL-Lite_{R,∩}*, are “optimally efficient”, in the sense that their data complexity is even lower. Indeed, relatively to consistency, the problem we are interested in this chapter, they are in **AC⁰**⁴

⁴ The class **AC⁰**, is a complexity class strictly contained in (and hence easier than) **PTime**. SQL query evaluation in relational databases is in **AC⁰** in data complexity, which accounts for the efficiency of database management systems in dealing with large amounts of data.

in data complexity and in **PTime** in combined complexity.⁵ The DLs in the *DL-Lite* family are essentially the maximal DLs that exhibit such nice computational properties (Calvanese et al., 2013; Artale et al., 2009). This is a consequence of suitable syntactic restrictions that have been imposed in such logics:

- Concepts are not closed under Boolean operations: negation is restricted to basic concepts within the scope of a *right Cr*, and the use of disjunction is ruled out.
- Value restriction, a typical DL construct corresponding to a form of universal quantification, is not allowed, and the use of qualified existential quantification is restricted to the right-hand side of inclusion assertions.

These restrictions ensure that the *DL-Lite* logics are contained in the Horn fragment of FO. The *DL-Lite* constructs are nevertheless sufficiently expressive to cover the main features of conceptual modeling languages such as UML class diagrams and of concept hierarchies in ontologies and ontology-based systems. This is important, since it implies that in practice reasoning does indeed scale to very large ontologies that can capture several naturally arising domains of interest. This has to be compared with the much higher computational complexity of more expressive DLs. For illustration, consider in Table 2 the complexity of the DL *ALC*, which is the smallest logic containing the *DL-Lite* logics that we have considered here⁶ and closed under Boolean operations. Both for *ALC* and for *SHOIN*, the DL that underpins OWL DL, reasoning is **coNP**-hard in data complexity, and provably exponential in combined complexity. Notice that as soon as a DL becomes closed under Boolean operations, it is intractable, and hence reasoning does not really scale well with data growth.

We are interested in studying the linguistic structures that correspond to the *DL-Lite* constructs. In what follows (Section 5 below), we will look at straightforward ways to express them in natural language.

4 Categorical Grammars

As most of the linguistically motivated formal grammars currently in use, categorical grammars (CGs) are a class (or family of classes) of lexicalized grammars, i.e., grammars where the lexicon carries most of the information about how words can be assembled to form grammatical structures. In this framework, syntactic categories are seen as *formulas* and their category forming operators as connectives, i.e., *logical constants*. In addition, the Curry-Howard correspondence ensures the Montagovian homomorphism, a.k.a. syntax-semantics interface, between the (logical) calculus of syntactic categories and FO MRs (van Benthem, 1987).

⁵ Notice that Pratt and Third's complexity results do not distinguish between data and combined complexity.

⁶ All *DL-Lite* logics include also the *inverse role* constructor, which cannot be captured in *ALC*. Moreover, some *DL-Lite* variants use (complex) role inclusions, which also would lead the logic outside the scope of *ALC*.

The peculiarity of CGs is that word assembly is carried out by natural deduction logical rules (that take care of natural language syntax); such natural deduction rules are coupled with (via the Curry-Howard correspondence) λ -calculus operations dealing with the FO meaning assembly, via the intermediate λ -FO formalism, viz., FO extended with (typed) λ -calculus λ -abstractions and λ -applications. In so doing, CGs capture better and more elegantly the tight correspondence between syntax and semantics of NL and its fragments than other equivalent grammatical formalisms such as semantically-enriched context-free grammars or some simple kinds of definite clause grammars.

This aspect of the formalism significantly simplifies the implementation task, since one has to focus only on the construction of the lexicon and can rely on any existing parser for the calculus. Information both about the syntactic structure where the word could occur and its meaning are stored in the lexicon. As derivation or logical deduction rules, we use the product free version of the (non associative) Lambek calculus. (Lambek, 1958; Moortgat, 1997)⁷

Definition 2 (Term Labeled Lexicon, Categorical Grammar). A (syntactic) *category* A is defined as follows

$$A \longrightarrow np \mid n \mid s \mid A_1 \backslash A_2 \mid A_2 / A_1$$

where np (*noun phrases*), n (*nouns*) and s (*complete sentences*) are *atomic* categories. Complex categories are built out of atomic categories by means of the directional *left* and *right* functional connectives \backslash and $/$ ($A_1 \backslash A_2$, resp. A_2 / A_1 , applied to a category A_1 situated to its left, resp. its right, yield category A_2). We denote by CAT the set of all such categories and by ATOM the set $\{np, n, s\}$.

We map each syntactic category A to a (semantic) *type* $typ(A)$ as follows:

$$\begin{aligned} typ(np) &= e; & typ(s) &= t; & typ(n) &= (e, t), \\ typ(A_1/A_2) &= (typ(A_2), typ(A_1)); & typ(A_2 \backslash A_1) &= (typ(A_2), typ(A_1)). \end{aligned}$$

where the atomic types are e (*entities*) and t (*Booleans*), and (τ, τ') denotes the *functional type* (the type of functions from τ into τ').

Given a set Σ of natural language basic expressions (i.e., a natural language vocabulary), a *term labeled categorical lexicon* is a relation,

$$\text{LEX} \subseteq \Sigma \times (\text{CAT} \times \text{TERM}) \quad \text{s.t.,} \quad \text{if } (w, (A, \alpha)) \in \text{LEX}, \text{ then } \alpha \in \text{TERM}_{typ(A)}$$

where TERM is the set of all lambda terms and $\text{TERM}_{typ(A)}$ denotes the set of lambda terms whose type is mapped to the category A .

Given a term labeled lexicon LEX, a *categorical grammar* is any finite subset $G \subseteq \text{LEX}$. ♠

This constraint on lexical entries categories and terms enforces the following requirement: if the expression (or word) w is assigned the syntactic category A and

⁷ The lexicon we present in this article has been tested using the GRAIL parser (Moot, 1998), based on the Lambek calculus.

the term α , then the term α must be of a type appropriate for the category A . We assign lambda terms whose body is a FO formula, viz., λ -FO terms. We look at the determiner *every*, by means of example, since it has a crucial role in our grammar. The reader is referred to work by Keenan and Faltz (1985) and van Eijck (1985) for an in-depth explanation of this example in particular and the relationships between CGs and λ -FO in general.

Example 2 (Determiner). The meaning of “every NOUN” (e.g., “every man”) is the set of those properties that “every NOUN” (e.g., “man”) has

$$\llbracket \text{every NOUN} \rrbracket = \{X \mid \llbracket \text{NOUN} \rrbracket \subseteq X\}.$$

In a functional perspective, the determiner “every” is seen as a two-argument function taking a noun and a verb phrase (a property) as arguments. The syntactic category expressing this functional view as well as word order is the following

$$(s/(np\s))/n$$

where the n is the first argument that must occur on the right of “every” and $np\s$, i.e., a verb phrase, is its second argument to occur still on the right of “every NOUN” (viz. “every NOUN VERB_PHRASE”). The typed lambda term (according to generalized quantifier theory, see Barwise and Cooper (1980)) corresponding to this syntactic category is: $\lambda Y_{(e,t)}. \lambda X_{(e,t)}. \forall x_e (Y(x) \rightarrow X(x))$. In the following, we will not use types on lambda terms unless necessary. ♣

An important feature of CGs is their “parsing as deduction” approach, which reduces the problem of checking whether a linguistic string is grammatical to the problem of proving that the string is of a certain syntactic category. More precisely, instead of directly recognizing linguistic word strings $w_1 \cdots w_n$, we work on the corresponding set of Lambek calculus formulas: to each lexicon entry $(w_i, (A_i, \alpha_i))$, for $i \in \{1, \dots, n\}$ we associate a (Lambek calculus) *sequent* $A_i \vdash A_i : \alpha_i$; thereafter, following the inference rules of the calculus, a proof (a tree-shaped derivation) of a sequent $\Gamma \vdash s : \phi$, with ϕ of type t (a λ -FO formula) is constructed. More formally:

Definition 3 (Recognized Language). Given a categorial grammar G the *language recognized by G* , denoted $L(G)$ is the set of all word strings $w_1 \cdots w_n$ such that the sequent $\Gamma \vdash s : \phi$, has a proof in the Lambek calculus; where Γ consists of a set $\{A_1 : \alpha_1, \dots, A_n : \alpha_n\}$ of pairs of categories and terms as defined in the term labeled lexicon $\{(w_i, (A_i, \alpha_i)) \mid i = 1, \dots, n\}$, and ϕ is a λ -FO formula (a term of type t). ♠

As by-product of the derivation one derives also the MR of the structure assigned to the string, i.e., the λ -FO term ϕ which after reduction gives rise to a FO closed formula or sentence. As such, NLs (and fragments thereof) recognized by a CG that does not cover purely higher-order NL constructs such as, e.g., the second-order determiner “most”, can induce (in a way similar, though more general, to Pratt and Third’s fragments) a fragment of FO: the set of all the first-order MRs associated with its (grammatical) complete sentences. We will exploit this particular feature of the formalism to define a CL in the next section that generates the *DL-Lite* logics.

5 Lite English and its Grammar CG-lite

As mentioned above, the goal of our methodology is to define CLs for ontologies that are *efficient*, i.e., tractable w.r.t. semantic data complexity. We propose to this end to define them vis-à-vis those ontology constructs that give rise to tractable data complexity. More precisely, we propose to identify English syntactic categories that lexically control the restrictions imposed by the *DL-Lite* constructs. Such categories will naturally induce a CG (i.e., a term-labeled categorial lexicon) expressing *exactly* the *DL-Lite* family of logics as described earlier. In this section we outline such syntactic categories and how they were obtained. We proceed in three steps. Firstly, we outline the key constraints to be satisfied for a CL to induce *DL-Lite*. Secondly, we provide a sample CG (a finite term-labeled lexicon). Thirdly, we describe the main features of the fragment thus generated. Notice also that the methodology proposed is not, per se, grammar dependent, since our CLs can be equally, although less succinctly and not as elegantly, defined using semantically-enriched context-free grammars as we did in some previous work (Thorne, 2010, Ch. 4). We call Lite-English the resulting CL, and CG-lite its CG.⁸

5.1 Fragment of Natural Language for DL-Lite

The constraints expressed in the TBox are universally quantified FO sentences. They are of the form $Cl \sqsubseteq Cr$, which translates into FO as $\forall x. \varphi(Cl, x) \rightarrow \varphi(Cr, x)$ and can be expressed by the following NL sentence patterns:

- (a) [Every $\underbrace{\text{NOUN}}_{Cl}$ $\underbrace{\text{VERB_PHRASE}}_{Cr}$]
 (b) [[Everyone $\underbrace{[\text{who VERB_PHRASE}]}_{Cl}$] $\underbrace{\text{VERB_PHRASE}}_{Cr}$]

The determiner “every” and the quantifier phrase “everyone” play a crucial role in determining the linguistic structures that belong to the natural language fragment corresponding to a *DL-Lite* TBox. In the following, we zoom into the `NOUN` and `VERB_PHRASE` constituents. In other words, we spell out how *DL-Lite* *Cl* and *Cr* concepts can be expressed in English. In doing so, we follow Definition 1.

First of all, a *Cl* or a *Cr* could be an atomic concept *A*. An atomic concept *A* corresponds to a unary predicate, which following standard formal semantic theory can be expressed either by a noun such as “student” (see (4) below), or an intransitive verb such as “work” (see (5) below).

The introduction of negation $\neg A$ on atomic concepts *A*, however, can occur only in a *Cr* and can thus be expressed only by a *predicate* `VERB_PHRASE` such as “is not a BSc-student” (6), or “does not work” (7).

⁸ We refer the reader to the appendix for the formal proofs of the claims made in this section.

The introduction of the $\exists R$ in a *Cl* can be performed by means of the quantifier phrase “everyone” followed by the relative pronoun “who” (9) (or by the conjunction that would correspond to the use of \sqcap on the *Cl* part allowed in the *DL-Lite_{R, \sqcap}* fragment, see (16) below).

- | | |
|---|--|
| (4) Every MSc-student is a student. | $[MScStudent \sqsubseteq Student]$ |
| (5) Every MSc-student works. | $[MScStudent \sqsubseteq Works]$ |
| (6) Every MSc-student is not a BSc-student. | $[MScStudent \sqsubseteq \neg BScStudent]$ |
| (7) Every BSc-student does not work. | $[BScStudent \sqsubseteq \neg Works]$ |
| (8) Everyone who learns works. | $[Learns \sqsubseteq Works]$ |
| (9) Everyone who reads something works. | $[\exists Reads \sqsubseteq Works]$ |

On the other hand, the introduction of $\exists R$ on the *Cr* part corresponds to the use of a transitive verb followed by an existential quantifier phrase, “something” (10), and its negation to the use of “does not” to negate such construction (11).

- | | |
|---|--|
| (10) Every student reads something. | $[Student \sqsubseteq \exists Reads]$ |
| (11) Every student does not read something. | $[Student \sqsubseteq \neg \exists Reads]$ |

Note that, as the *DL-Lite* clause shows, the only reading of the ambiguous sentence in (11) is the one with *every* having wide scope and *something* being in the scope of *not*⁹.

Also, the *VERB_PHRASE* in (a) and the second *VERB_PHRASE* in (b) (i.e., the *VERB_PHRASE* of the main clause expressing a *DL-Lite Cr* concept) can be of any of the structures in (4)–(11). On the other hand, the first *VERB_PHRASE* in (b) (i.e., the *VERB_PHRASE* of the relative clause expressing a *DL-Lite Cl* concept) cannot contain negation: for it only the cases 5-4 above hold.

When we move to *DL-Lite_{R, \sqcap}*, the addition of the conjunction in the *Cl* corresponds to the use of adjective (12), or relative clauses modifying the noun quantified by “every” (13-15), or the “and” coordinating two VPs (16).

- | | |
|---|---|
| (12) Every nice student works. | $[Student \sqcap Nice \sqsubseteq Works]$ |
| (13) Every student who learns works. | $[Student \sqcap Learns \sqsubseteq Works]$ |
| (14) Every student who is a BSc-student works. | $[Student \sqcap BScStudent \sqsubseteq Works]$ |
| (15) Every student who reads something works. | $[Student \sqcap \exists Reads \sqsubseteq Works]$ |
| (16) Everyone who reads something and writes something works. | $[\exists Reads \sqcap \exists Writes \sqsubseteq Works]$ |

Furthermore, the introduction of the qualified existential on the *Cr* is performed by the determiner “a” (17).

- | | |
|----------------------------------|--|
| (17) Every student reads a book. | $[Student \sqsubseteq \exists Reads.Book]$ |
|----------------------------------|--|

⁹ For ease of explanation we do not consider the distinction between *something* and the negative polarity item *anything*. This distinction could be incorporated into the fragment, as studied by Bernardi (2002).

Non-Boolean-closedness (tractability) of the fragment. An important point to emphasize is the presence of the relative pronoun in the above fragment of sentences. Pratt and Third have shown how the uncontrolled use of such expression leads to **NP**-complete fragments when allowing the use only of the copula, or even to **ExpTime**-completeness when adding transitive verbs. Below, we will show how relative pronouns can be used in a controlled grammar while preserving tractability of inferences.

5.2 Expressing *DL-Lite_{core}*

We start again by looking at the main syntactic constraints over *DL-Lite_{core}* concepts and consider, in particular, the two constraints regarding the use of negation:

1. negation of atomic concepts can occur in a *Cr* but not in a *Cl*: $Cl \longrightarrow B, Cr \longrightarrow B \mid \neg B$;
2. an unqualified existential can occur both in a *Cl* and a *Cr*, but its negation can occur only in *Crs*: $Cl \longrightarrow \exists R, Cr \longrightarrow \exists R \mid \neg \exists R$.

As we anticipated before, *Cl* and *Cr* concepts correspond respectively to the so-called “restrictive scope” (the subject NOUN constituent), and “nuclear scope” (the predicate VERB_PHRASE constituent) of the sentence-building DET *every*. We need to constrain the linguistic structures that occur within them. In particular, we need to block the occurrences of negation within *Cl*s and express the fact that NOT cannot outscope any VERB_PHRASE that occurs within the restrictive scope of the determiner *every*. As emphasized by Bernardi (2002), in CGs scope is determined by the sentential categories *s* that arise from complex CG syntactic categories. Different (possibly mutually exclusive) scope distributions can be enforced by multiplying sentential categories via *sentential levels*, and exploiting the derivability relations (and restrictions) among CG categories. It suffices to provide the intuition behind the proposed solution without going into its details: a complex category $A_1 \setminus A_2$, can be applied to either category A_1 or to a category A_3 that derives A_2 ($A_3 \Rightarrow A_1$). In our case, \Rightarrow is the *derivability* relation of the logical grammar we use.

We mark the structures that express *DL-Lite Cl*s and *Crs* and those that are negative or positive, by means of the four *sentential levels* s_{cl} , s_{cr} , s_{\neg} , and s , respectively, and establish the derivability relation below (we rule out any other derivability relations between atomic categorial formulas).¹⁰ These sentential levels state that a negated sentence can be in the *Cr* construct ($s_{\neg} \Rightarrow s_{cr}$) while it cannot be in the *Cl* part ($s_{\neg} \not\Rightarrow s_{cl}$) and a positive sentence can be in both ($s \Rightarrow s_{cl}, s \Rightarrow s_{cr}$):

$$s_{\neg} \not\Rightarrow s_{cl}, \quad s_{\neg} \Rightarrow s_{cr}, \quad s \Rightarrow s_{cl}, \quad s \Rightarrow s_{cr}, \quad \text{and} \quad s_{cl} \not\Rightarrow s_{cr}.$$

¹⁰ We actually use residuated unary operators to carry out these derivability relations (Kurtonina and Moortgat, 1995) exploiting their logical properties: $\diamond_j \square_j s \Rightarrow s \Rightarrow \square_i \diamond_i s$ etc. Examples of residuated unary operators are “possibility in the past” and “necessity in the future”.

Note that this induces a derivability relation between complex categories built with or containing these atomic sentential categories; for instance, from $s \Rightarrow s_{cr}$ it follows that $np \setminus s \Rightarrow np \setminus s_{cr}$. Besides these sentential levels, as we will show below, we use two other sentential levels: one to mark TBox sentences (s_{tb}) and one to mark constituents built by the relative pronoun who (s_{who}). All the constraints on these sentential levels are lexically anchored by means of the lexical assignments below.

Example 3 (Lexicon for DL-Lite_{core}). The lexicon entries to use are as below.¹¹ The content words (intransitive verbs and nouns) are only given by way of example.

- Every $\in (s_{tb}/(np \setminus s_{cr}))/n: \lambda X.\lambda Y.\forall x.(X(x) \rightarrow Y(x))$
- is a $\in (np \setminus s)/n: \lambda X.\lambda z.X(z)$
- is not a $\in (np \setminus s_{-})/n: \lambda X.\lambda z.\neg X(z)$
- does not $\in (np \setminus s_{-})/(np \setminus s): \lambda X.\lambda z.\neg X(z)$
- works $\in np \setminus s: \lambda z.WORKS(z)$
- learns $\in np \setminus s: \lambda z.LEARNS(z)$
- student $\in n: \lambda z.STUDENT(z)$
- MSc-student $\in n: \lambda z.MScSTUDENT(z)$
- BSc-student $\in n: \lambda z.BScSTUDENT(z)$
- everyone: $(s_{tb}/(np \setminus s_{cr}))/((np \setminus s_{who})): \lambda X.\lambda Y.\forall x.(X(x) \rightarrow Y(x))$
- who: $(np \setminus s_{who})/(np \setminus s_{ci}): \lambda P.\lambda z.P(z)$
- something: $((np \setminus s_{\exists})/np) \setminus (np \setminus s): \lambda Z.\lambda y.\exists x.Z(y, x)$
- reads: $(np \setminus s_{\exists})/np: \lambda x.\lambda z.READS(z, x)$ ♣

A. Using universal quantification to express concept subsumption. Notice that in Example 3 the categories assigned to *every* and *everyone* rule out the possibility for them to occur in object position –they can only be in subject position. Moreover, since they are the only entries yielding a TBox sentence (s_{tb}), only sentences starting with them will be considered as grammatical. The negation brings sentences to the negative sentential level, and once they are there, they are blocked from occurring in the restrictive scope of *every* and *everyone*.

B. Using existential quantification, relatives, and conjunction to express existentially qualified roles and their conjunctions. Since in the fragment described by Example 3, we do not have the \sqcap on the *CI*, the introduction of the unqualified existential $\exists R$ in it can be performed only by means of the quantifier *everyone* followed by the relative pronoun “who” and a transitive verb composed with *something*. The introduction of $\exists R$ on the *Cr* corresponds to the use of a transitive verb followed by an existential quantifier, *something*. The lexical entries for *everyone*, *who*, *something*, and *reads* above account for these facts. The need of the s_{who} categories is due to the fact that *everyone* must be followed by a relative clause, i.e., sentences

¹¹ Notice, in the present work we do not handle features of any sort (morphological etc). Their usage will make the lexical entries more complex but won't have any effect on the main idea we are presenting.

like *everyone left* or *everyone walks and speaks* cannot be part of the grammar. Similarly, transitive verbs can occur on the *Cr* part but only if followed by *something*, hence we use the category s_{\exists} to guarantee this requirement.¹² Finally, the category assigned to “something” is such that it can occur only in object position.

C. Controlling the behavior of negation. As the reader can see, negation in Example 3 can only occur within a `VERB_PHRASE` expressing a *Cr*. The reader can gain a better understanding of the mechanisms involved by checking how our sample lexicon, combined with the constraints CG-lite imposes over its sentential levels, ensures the ungrammaticality of the sentences below (blocked by $s_{\neg} \not\rightarrow s_{cl}$). Such sentences generate MRs that are not *DL-Lite* expressible:

- (18) Everyone who does not read something works [$\neg\exists$ Reads \sqsubseteq Works]
 (19) Everyone who is not a BSc-student works. [\neg BScStudent \sqsubseteq Works]

D. Expressing ABoxes. The fragment of sentences whose meaning representation belongs to a *DL-Lite_{core}* ABox is rather easy to build since an ABox consists of a conjunction of (ground) unary and binary logical atoms. In other words, the lexicon is built only with nouns, intransitive verbs, the copula (i.e., unary predicates), transitive verbs (i.e., binary predicates), individual names and adjectives.

5.3 Expressing *DL-Lite_{R,□}*

We now move to *DL-Lite_{R,□}*, and account for the following additions

1. conjunctions are allowed in *Cls*: $Cl \rightarrow Cl_1 \sqcap Cl_2$;
2. the qualified existential can occur in *Crs*: $Cr \rightarrow \exists R.B$.

*Example 4 (Lexicon extension for *DL-Lite_{R,□}*).* In order to move to *DL-Lite_{R,□}*, we need to add into the lexicon the following lexical entries. The (intersective, qualitative) adjective *nice* is given only by way of example.

- nice: $n_{cl}/n_{cl}, \lambda X.\lambda z.(X(z) \wedge \text{Nice}(z))$
- who: $(n_{cl} \setminus n_{cl})/(np \setminus s_{cl}): \lambda X.\lambda Y.\lambda z.(X(x) \wedge Y(z))$
- and: $((np \setminus s_{cl}) \setminus (np \setminus s_{cl}))/(np \setminus s_{cl}): \lambda X.\lambda Y.\lambda z.(X(z) \wedge Y(z))$
- a: $((np \setminus s_{\exists})/np) \setminus (np \setminus s_{cr})/n: \lambda Y.\lambda Z.\lambda y.\exists x.(Z(y,x) \wedge Y(x))$

Again, we use sentential levels to control the occurrence of these constructs. The extended lexicon accounts also for the structures in (12)–(17). ♣

A. Controlling the interaction of conjunction and negation. Notice the need of having a conjunction operating at the sentential level s_{cl} : this blocks the composition of negation (*does not*) with a verb phrase built with an *and* that would

¹² Since we have neither *np* nor *np/n* entries we could also avoid the use of this extra sentential level s_{\exists} in the example we are considering.

wrongly give or recognize: *does not walk and speak* with *not* outscoping *and*; such constituent would yield the MR $\lambda z. \neg(\text{Walk}(z) \wedge \text{Speak}(z))$ that is not *DL-Lite* expressible, and would moreover give rise to intractable data complexity (Calvanese et al., 2013). For similar reasons we have to block the composition of *is not a* with a noun phrase built using an intersective adjective. The resulting NOUN_PHRASE constituent would yield non-*DL-Lite*-expressible λ -FO formulas where negation outscopes conjunction; e.g., a phrase like *is not a nice student* with MR $\lambda z. \neg(\text{Nice}(z) \wedge \text{Student}(z))$. The introduction of the category n_{cl} with $n \Rightarrow n_{cl}$ makes such phrases ungrammatical.

B. Qualified existential restrictions and recursive constituents. We have considered a DL, *DL-Lite*_{R, \sqcap} , with qualified existentials of the form $\exists R.A$. Hence the argument taken by the determiner a can only be a bare noun n . Finally, notice that the lexical entries for the adjective, conjunction, and qualified existential bring recursion into the language.

6 Distribution of Boolean- and non-Boolean-closed Fragments

As we have shown, reverse-engineering efficient CLs from ontologies is a promising path. Further, as shown by Thorne (2010), our methodology can be easily extended to define interrogative CLs with tractable data complexity. The question however remains as to how to identify CLs that, while enjoying the properties we desire them to have (express ontology and query languages, give rise to at most **PTime** data complexity), remain appealing to users.

In this section we propose a distributional methodology which may help in identifying desirable English constructs by focusing on their frequency in both interrogative and declarative English corpora. We believe that this method can yield techniques to pinpoint, in particular, CLs that may offer good trade-offs between coverage and semantic complexity. The intuition behind being that when we trade-off language coverage for performance (i.e., to attain tractable data complexity) in CLs, it makes sense to cover constructs that are frequently used and thus preferred by speakers. Specifically, we study the co-occurrence of crucial (for semantic complexity) logic constructs: negations, conjunctions, disjunctions, and universal and existential quantification, in English questions and sentences.

To obtain a representative sample we considered corpora of multiple domains and with sentences of arbitrary type (declarative and interrogative), since, when managing an ontology and/or an ontology-based system, users are required not only to assert but also to update and query (intensional and extensional) information belonging to different domains. We thus considered: (i) a subset (A: press articles) of the Brown corpus¹³; (ii) a subset of one (Geoquery880) of the Geoquery cor-

¹³ http://nltk.googlecode.com/svn/trunk/nltk_data/index.xml

Table 3 Corpora analyzed in this chapter.

| Corpus | Size | Domain | Sentence type |
|---------------------|------------------|--------------|---------------------------|
| Brown corpus subset | 19,741 sentences | Open (news) | Declarative ¹⁷ |
| Geoquery corpus | 364 questions | Geographical | Interrogative |
| Clinical questions | 12,189 questions | Clinical | Interrogative |
| TREC 2008 | 436 questions | Open | Interrogative |

pora¹⁴; (iii) a corpus of clinical questions¹⁵; and (iv) a sample from the TREC 2008 corpus¹⁶. Table 3 summarizes their main features.

To this end we exploited the availability of wide-coverage (statistical) deep semantic parsers such as Boxer, by Bos (2008), which output first-order MRs. We checked, for each such MR, the co-occurrence of a subset of the set $\{\forall, \exists, \neg, \wedge, \vee\}$ of FO operators (and only of that subset). Each such subset identifies MRs belonging, modulo logical equivalence, to a distinct fragment of FO. For instance, the combination $\{\forall, \exists, \wedge, \vee\}$ identifies MRs from the so-called positive fragment of FO. But it also identifies the class of corpora sentences that give rise to such MRs, and approximates the (controlled) fragment whose formal semantics may induce such FO fragment. Finally, with these considerations in mind, we observed the distribution of:

1. “Boolean-closed” fragments, viz.: $\{\exists, \wedge, \neg\}$, $\{\exists, \wedge, \neg, \forall\}$, $\{\exists, \wedge, \neg, \forall, \vee\}$, $\{\neg, \forall\}$, $\{\exists, \wedge, \forall\}$, and $\{\exists, \wedge, \forall, \vee\}$.
2. “Non-Boolean-closed” fragments, viz.: $\{\exists, \wedge\}$ and $\{\exists, \wedge, \vee\}$.

By “Boolean-closed”, we recall, we mean fragments expressive enough to encode Boolean satisfiability and which give rise to intractable semantic complexity. A “non-Boolean-closed” combination, by contrast, cannot express Boolean functions and gives rise only to tractable semantic complexity.

The pipeline of Boxer consists of the following three basic steps: (i) each part of speech in a sentence is annotated with its most likely (categorial grammar) syntactic category; (ii) the most likely of the resulting possible combinatorial categorial grammar derivations (or proofs) is computed and returned; and (iii) a neo-Davidsonian semantically weakened¹⁸ FO meaning representation is computed using discourse representation theory (DRT).

Example 5. When parsing Wh-questions from the TREC 2008 corpus such as “What is one common element of major religions?”, Boxer outputs a FO semantic representation of the form

¹⁴ <http://www.cs.utexas.edu/users/ml/nldata/geoquery.html>

¹⁵ <http://clinques.nlm.nih.gov/>

¹⁶ <http://trec.nist.gov/>

¹⁷ The sample contained only 36 questions.

¹⁸ In this settings, the semantics of verbs is represented in terms of events connected via thematic roles to verb arguments (agents, themes, etc.). In addition, the semantics of non-FO constructs such as “most” is weakened to some FO representation.

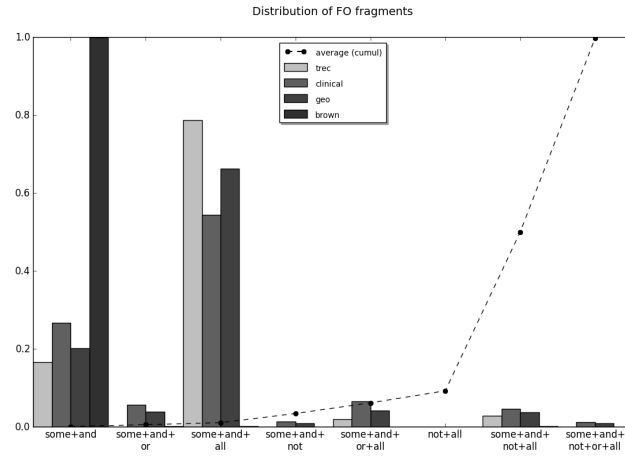


Fig. 2 Relative frequency of *co-occurring* FO operators in sample corpora. Notice the distribution of “non-Boolean-closed” sentences.

$$\exists y \exists z \exists e \exists u (\text{card}(y, u) \wedge \text{c1num}(u) \wedge \text{nnumerall}(u) \wedge \\ \text{acommon1}(y) \wedge \text{nelement1}(y) \wedge \text{amajor1}(z) \wedge \text{nreligions1}(z) \wedge \\ \text{nevent1}(e) \wedge \text{rof1}(y, z))$$

where \wedge and \exists co-occur, but not \vee , \neg , or \rightarrow . ♣

Figure 2 shows the co-occurrence distribution obtained, expressed in terms of relative frequency (i.e., number of MRs per class/total number of MRs per corpus). As the figure shows, positive existential, $\{\exists, \wedge\}$ and $\{\exists, \wedge, \vee\}$, MRs occur quite frequently. Also, it seems that the same holds for sentences expressing universal quantification whereas the opposite is true for negation (low frequency overall).

This analysis can be compared to the more linguistics-based methodology followed by (Bernardi et al., 2007), in which we analyzed the distribution in (solely) interrogative corpora of classes of *logical words* which express FO operators, e.g., “all”, “both”, “each”, “every”, “everybody”, “everyone”, “any”, “none”, “nothing”. See Figure 3.

These results suggest that, while users use negation or disjunction words as frequently as conjunction and existential words, and all these more than universal words, when combining them *within* sentences “non-Boolean-closed” combinations are preferred.

7 Related Work

The work described in this chapter has been complemented by related results obtained by the authors and published elsewhere. In particular, we have applied and

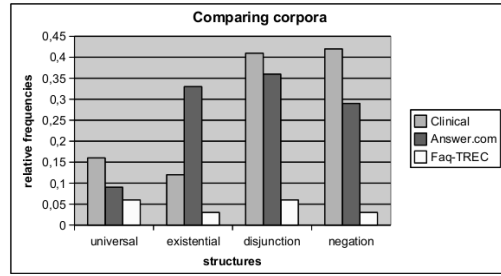


Fig. 3 Relative frequency of FO operators in question corpora (Bernardi et al., 2007).

Table 4 Non-“Boolean-closed” and “Boolean-closed” controlled English constructs.

| Data Complexity | Constructs |
|-----------------|---|
| Tractable | Negation (“not”) in a predicate VERB_PHRASE Relatives (“who”, “which”) everywhere Conjunction (“and”) everywhere Transitive verbs (“loves”) everywhere Existential quantification (“some”) everywhere |
| Intractable | Negation (“not”) in a subject NOUN Universal quantification (“only”) in a subject NOUN Disjunction (“or”) in predicate VERB_PHRASE |

generalized the methodology defined in this chapter to determine which are the fragments of ACE-OWL that are tractable (i.e., at most **PTime**) and those that are intractable (i.e., **coNP-hard**) in data complexity (see Thorne and Calvanese, 2012). Table 4 summarizes what these results mean in terms of language coverage, viz., which *maximal* combinations of (English) function and content words give rise to *tractable* (“non-Boolean closed”) controlled fragments, and which *minimal* combinations give rise to *intractable* (“Boolean closed”) controlled fragments.

Intractability arises with any combination capable of expressing full Boolean negation (“not”) and full Boolean conjunction (conjunction, relative pronouns). Note that the good computational properties of Lite-English depend ultimately on the fact that, while expressing Boolean conjunction, it cannot express full Boolean negation, but rather a very limited form of it. We observe that a similar analysis carried out on Pratt and Third’s fragments (and extending their own) by Thorne (2010) yielded similar results.

We have also applied the methodology used in this chapter to study controlled fragments of English *questions* for which the data complexity of reasoning (or evaluation) against an ontology, authored using any of their fragments, is tractable (see Thorne, 2010, Ch. 5). This work shows that positive questions (questions built with “some”, relative pronouns, conjunction and eventually, disjunction) and restricted to proper and common nouns, and intransitive and transitive verbs as content words, give rise to tractability. It also shows that they can be enriched with so-called *aggre-*

gate determiners, i.e., English constructs such as “the total number of”, “the number of”, “the average of”, etc., that express *aggregate functions*¹⁹, in formal query languages without negatively impacting on the semantic complexity of the controlled fragments.

As we hinted in the introduction, several CLs, most of which are equipped with a compositional semantics, have been proposed to provide NLI to ontologies and ontology-based systems. In particular, to provide English front ends to (i) ontology authoring systems, specifically, semantic web ontologies in the form of OWL DL ontologies (for which its fragment ACE-OWL was engineered) and (ii) controlled English querying to such ontologies. Table 5 provides an overview of the best known and used, viz., PENG (Schwitter et al., 2003), Rabbit (Schwitter et al., 2008) and OWL CNL (Schwitter and Tilbrook, 2006), which are to a big extent siblings and/or children of the main two: Attempto Controlled English (ACE) and its fragment ACE-OWL (Fuchs et al., 2006; Kaljurand, 2007).

While the coverage of ACE, ACE-OWL and its relatives is way greater than of any of the CLs defined in this chapter, they suffer from our perspective from the fact of being too expressive. Query evaluation over OWL DL (viz., *SHOIN*) ontologies and a fortiori in ACE-OWL NLI is **coNP**-hard in the size of the data, and hence intractable and unsuited for managing large data repositories. In more expressive CLs such as Rabbit or (full) ACE, reasoning is undecidable.²⁰ However ACE-OWL and kindred CLs contain (grammatically correct) fragments which may exhibit better computational properties, which we believe can be defined using our methodology.

In addition to NLI to OWL ontologies (Schwitter and Tilbrook, 2006; Kaljurand and Fuchs, 2006), systems have been proposed that, e.g., guide the user to formulate his/her natural language (NL) question via an ontology that incrementally shows the possible concepts that could be involved in the question (Franconi et al., 2010; Dongilli and Franconi, 2006). Others guide the user via an incremental parser (Bernstein et al., 2006; Damljanić, 2010), or engage the user in clarification dialogs (Gunning et al., 2010).

8 Conclusions

In this chapter we have outlined a methodology for defining controlled fragments (CLs) of English for NLI to ontology-based systems, which scale to very large ontologies. In addition to their scalability, such CLs can express key ontology language constructs via a symbolic translation formally underpinned by formal semantics in the Montagovian tradition.

We have argued that this can be achieved as follows: (i) On the one hand, by focusing on semantic complexity, viz., the computational complexity of logical rea-

¹⁹ That is, second-order functions such as, resp., $sum(\cdot)$, $\#(\cdot)$, $avg(\cdot)$, defined over *sets* of individuals or data values.

²⁰ Reasoning on OWL Full or FO is undecidable (cf., Baader et al., 2003).

Table 5 An overview of some CLs; DCG stands for “definite clause grammar”, the other acronyms for known English parsers or parser APIs such as GATE, and “comp.” for “compositional”.

| CL (English) | Comp. | Maps to | Parser | Goal |
|---|-------|----------|------------|-----------------|
| ACE (Fuchs et al., 2006) | yes | FO | APE | Knowledge repr. |
| ACE-OWL (Kaljurand, 2007) | yes | OWL DL | APE | Ontology mgmt |
| PENG (Schwitter et al., 2003) | yes | OWL DL | ECOLE | Ontology mgmt |
| OWL CNL (Schwitter and Tilbrook, 2006) | yes | OWL DL | DCG parser | Ontology mgmt |
| Rabbit (Schwitter et al., 2008) | no | OWL Full | GATE | Ontology mgmt |

soning in such CLs, which can be studied via the FO fragment induced by their formal, compositional semantics. We have stressed that a key requirement is for semantic complexity to be at most polynomial in the size of the ontology (or ontology-based system), and in \mathbf{AC}^0 in the size of the data stored therein, that is, to have efficient semantic data complexity. (ii) On the other hand, by considering English constructs that express ontology languages with efficient data complexity. (iii) Finally, by putting together those English constructs via CGs to build a CL that expresses such low complexity ontology languages and that possesses appropriate semantic data complexity while expressing key ontology language constructs.

Following our methodology, we have identified the fragment of English that corresponds to an ontology language suitable for specifying and querying ontologies with optimal data complexity, namely *DL-Lite*; and based on this we have defined an efficient CL, Lite English, using CGs (via the CG-lite grammar).

We have also performed a preliminary corpus analysis regarding the distribution of relevant English constructs. We believe that this methodology could, if further developed, help the CL community in identifying suitable CLs that provide good trade-offs between coverage and tractability.

Appendix

In this appendix we sketch how CG-lite formally captures *DL-Lite*_{R,□} (and a fortiori *DL-Lite*_{core}). That is, we show that for every *DL-Lite*_{R,□} TBox assertion $Cl \sqsubseteq Cr$, there exists a CG-lite derivation D rooted in $s_{tb} \vdash s_{tb} : \forall x(\varphi(Cl, x) \rightarrow \varphi(Cr, x))$.

Remark 1 (Cl and Cr vs. λ -FO). Recall that Cl and Cr concepts are defined as below.

$$Cl \longrightarrow B \mid \exists R \mid Cl_1 \sqcap Cl_2 \quad \text{and} \quad Cr \longrightarrow B \mid \neg B \mid \exists R \mid \neg \exists R \mid \exists R.B.$$

Left concepts correspond to: (i) B , i.e., $\lambda x.(\mathbb{B}(x))$ (in λ -FO), (ii) $\exists R$, i.e., $\lambda x.\exists y.R(x,y)$ (in λ -FO), and (iii) $Cl_1 \sqcap Cl_2$, i.e., $\lambda x.(\varphi(Cl_1,x) \wedge \varphi(Cl_2,x))$ (in λ -FO). Regarding right concepts, the new concepts that are not Cl s are: (i') $\neg B$, i.e., $\lambda x.\neg \mathbb{B}(x)$ (in λ -FO), (ii') $\neg \exists R$, i.e., $\lambda x.\neg \exists y.R(x,y)$ (in λ -FO), and (iii') $\exists R.B$, i.e., $\lambda x.\exists y.(R(x,y) \wedge \mathbb{B}(y))$ (in λ -FO). †

Remark 2. In a CG-lite derivation of $\Gamma \vdash A : \alpha$, the resulting category A will match a subcategory A' occurring in a positive position within the categories occurring in Γ . This means that, when expressing left and right concepts we are interested in derivations where $Cl_1) A = n$, $Cl_2) A = np \setminus s_{cl}$, $Cl_3) A = np \setminus s_{who}$ and $Cr) A = np \setminus s_{cr}$. †

Lemma 1 (Left Cl concepts). For every $DL\text{-Lite}_{R,\sqcap}$ left concept Cl , there exists a CG-lite derivation D satisfying Remarks 1 and 2 that expresses it.

Proof. (Sketch) We show, by (structural) induction on left concepts Cl , that there exists a CG-lite derivation D rooted in either of the following three Lambek sequents: (1) $n \vdash n : \lambda x.\varphi(Cl,x)$ or (2) $np \setminus s_{cl} \vdash np \setminus s_{cl} : \lambda x.\varphi(Cl,x)$ or (3) $np \setminus s_{who} \vdash np \setminus s_{who} : \lambda x.\varphi(Cl,x)$, with categories found in or derived from CG-lite's lexicon CAT_{lex} .

- Base cases: Cl is an atomic concept B or a qualified existential $\exists R$.
 1. Consider the lexicon entry $n \vdash n : \lambda x.Student(x)$; (1) holds.
 2. Consider the lexicon entry $np \setminus s \vdash np \setminus s : \lambda x.Left(x)$; (2) holds.
 3. Consider the two entries $((np \setminus s_{\exists})/np) \setminus (np \setminus s) \vdash ((np \setminus s_{\exists})/np) \setminus (np \setminus s) : \lambda z.\lambda y.\exists x.Z(y,x)$ and $(np \setminus s_{\exists})/np \vdash (np \setminus s_{\exists})/np : \lambda x.\lambda z.Reads(z,x)$. By applying one to each other, we derive $np \setminus s \vdash np \setminus s : \lambda x.\exists y.Reads(x,y)$. Since $s \Rightarrow s_{cl}$, (1) holds.
- Inductive cases: Cl is a complex concept $Cl_1 \sqcap Cl_2$. By I.H. the property holds for Cl_1 and Cl_2 . There are several cases. As they are similar, we deal only with one.
 1. Consider the lexicon entry $((np \setminus s_{cl}) \setminus (np \setminus s_{cl})) / (np \setminus s_{cl}) : \lambda X.\lambda Y.\lambda z.(X(z) \wedge Y(z))$, expressing conjunction. By I.H., we may combine it in turn with the (derived) sequents $np \setminus s_{cl} \vdash np \setminus s_{cl} : \lambda x.\varphi(Cl_1,x)$ and $np \setminus s_{cl} \vdash np \setminus s_{cl} : \lambda x.\varphi(Cl_2,x)$ (i.e., verifying (2)). This results in a derivation rooted in $np \setminus s_{cl} \vdash np \setminus s_{cl} : \lambda x.(\varphi(Cl_1,x) \wedge \varphi(Cl_2,x))$, which satisfies (2). □

Lemma 2 (Right Cr concepts). For every $DL\text{-Lite}_{R,\sqcap}$ right concept Cr , there exists a CG-lite derivation D satisfying Remarks 1 and 2 that expresses it.

Proof. (Sketch) The claim can be proven by case analysis on Cr as in the preceding lemma (there is no inductive clause in the Cr definition), viz., by showing that a derivation D rooted in (4) $np \setminus s_{cr} \vdash np \setminus s_{cr} : \lambda x.\varphi(Cr,x)$ exists. □

Theorem 1. For every $DL\text{-Lite}_{R,\sqcap}$ TBox assertion $Cl \sqsubseteq Cr$, there exists a CG-lite derivation D satisfying Remarks 1 and 2 that expresses it.

Proof. The proof follows from the two Lemmas above and by the fact that the only two lexical entries with s_{tb} in a positive position are those for:

1. “every”, i.e., $(s_{tb}/(np \setminus s_{cr}))/n_{cl} \vdash (s_{tb}/(np \setminus s_{cr}))/n_{cl} : \lambda X. \lambda Y. \forall x. (X(x) \rightarrow Y(x))$; and, on the other hand,
2. “everyone”, i.e., $(s_{tb}/(np \setminus s_{cr}))/n_{cl} \vdash (s_{tb}/(np \setminus s_{cr}))/n_{cl} : \lambda X. \lambda Y. \forall x. (X(x) \rightarrow Y(x))$.

Now, by Lemmas 1 and 2, we know that left concepts Cl and right concepts Cr are CG-lite-expressible, i.e., that there exist derivations for them rooted in $np \setminus s_{cl} \vdash np \setminus s_{cl} : \lambda x. \varphi(Cl, x)$ and $np \setminus s_{cr} \vdash np \setminus s_{cr} : \lambda x. \varphi(Cr, x)$, resp.

When we combine such sequents with the entry for “every”, we obtain immediately $s_{tb} \vdash s_{tb} : \forall x (\varphi(Cl, x) \rightarrow \varphi(Cr, x))$. In the case of “everyone”, we need to combine them with the entry for “who”, viz., $(np \setminus s_{who})/(np \setminus s_{cl}) \vdash np \setminus s_{who}/(np \setminus s_{cl}) : \lambda P. \lambda z. P(z)$, and we again derive $s_{tb} \vdash s_{tb} : \forall x (\varphi(Cl, x) \rightarrow \varphi(Cr, x))$. \square

For reasons of space, we omit the proof of the converse, viz., that every (complete) sentence w in Lite-English expresses a $DL\text{-}Lite_{R, \square}$ assertion. It can be constructed in a manner similar to Theorem 1, by induction on CG-lite derivations, i.e., by showing how every CG-lite constituent of category n or $np \setminus s_{cr}$ (resp. n or $n \setminus s_{cl}$) gives rise to a right (resp. left) concept. Such constituents are then combined together into a sentence expressing an assertion via the function words “every” or by “everyone who”. The sentential levels and the derivability relations that ensue (see Sections 5.2 and 5.3) prevent over-generation.

Acknowledgements This research has been partially supported by the EU under the large-scale integrating project (IP) Optique (Scalable End-user Access to Big Data), grant agreement n. FP7-318338.

References

- Androutsopoulos, I., G. D. Ritchie, and P. Thanish (1995). Natural language interfaces to databases – An introduction. *J. of Natural Language Engineering* 1, 29–81.
- Artale, A., D. Calvanese, R. Kontchakov, and M. Zakharyashev (2009). The *DL-Lite* family and relations. *J. of Artificial Intelligence Research* 36, 1–69.
- Baader, F., D. Calvanese, D. McGuinness, D. Nardi, and P. F. Patel-Schneider (Eds.) (2003). *The Description Logic Handbook: Theory, Implementation and Applications*. Cambridge University Press.
- Barwise, J. and R. Cooper (1980). Generalized quantifiers and natural language. *Linguistics and Philosophy* 4(2), 159–219.
- Bernardi, R. (2002). *Reasoning with Polarity in Categorical Type Logic*. Ph. D. thesis, UiL, OTS, Utrecht University.

- Bernardi, R., F. Bonin, D. Carbotta, D. Calvanese, and C. Thorne (2007). English querying over ontologies: E-QuOnto. In *Proc. of the 10th Congress of the Italian Association for Artificial Intelligence (AI*IA 2007)*.
- Bernstein, A., E. Kaufmann, C. Kaiser, and C. Kiefer (2006). Ginseng: A Guided Input Natural language Search Engine for querying ontologies. In *2006 Jena User Conference*.
- Bos, J. (2008). Wide-coverage semantic analysis with Boxer. In *Proc. of the 2008 Conf. on Semantics in Text Processing (STEP 2008)*.
- Calvanese, D., G. De Giacomo, D. Lembo, M. Lenzerini, A. Poggi, M. Rodriguez-Muro, R. Rosati, M. Ruzzi, and D. F. Savo (2011). The Mastro system for ontology-based data access. *Semantic Web J.* 2(1), 43–53.
- Calvanese, D., G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati (2007). Tractable reasoning and efficient query answering in description logics: The *DL-Lite* family. *J. of Automated Reasoning* 39(3), 385–429.
- Calvanese, D., G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati (2013). Data complexity of query answering in description logics. *Artificial Intelligence* 195, 335–360.
- Damljanovic, D. (2010). Towards portable controlled natural languages for querying ontologies. In *Proc. of the 2nd Workshop on Controlled Natural Languages (CNL 2010)*.
- Dongilli, P. and E. Franconi (2006). An intelligent query interface with natural language support. In *Proc. of the 19th Int. Florida Artificial Intelligence Research Society Conf. (FLAIRS 2006)*.
- Franconi, E., P. Guagliardo, and M. Trevisan (2010). Quello: NL-based intelligent query interface. In *Proc. of the 2nd Workshop on Controlled Natural Languages (CNL 2010)*.
- Fuchs, N. E., K. Kaljurand, and G. Schneider (2006). Attempto Controlled English meets the challenges of knowledge representation, reasoning, interoperability and user interfaces. In *Proc. of the 19th Int. Florida Artificial Intelligence Research Society Conf. (FLAIRS 2006)*.
- Gunning, D., V. K. Chaudhri, P. Clark, K. Barker, S. Chaw, M. Greaves, B. Grosz, A. Leung, D. McDonald, S. Mishra, J. Pacheco, B. Porter, A. Spaulding, D. Tecuci, and J. Tien (2010). Project Halo update – Progress toward digital Aristotle. *AI Magazine* 31(3), 33–58.
- Horrocks, I., P. F. Patel-Schneider, and F. van Harmelen (2003). From *SHIQ* and RDF to OWL: The making of a Web Ontology Language. *J. of Web Semantics* 1(1), 7–26.
- Huijsen, W. O. (1998). Controlled language – An introduction. In *Proc. of the 2nd Int. Workshop on Controlled Language Applications (CLAW 1998)*.
- Kaljurand, K. (2007). *Attempto Controlled English as a Semantic Web Language*. Ph. D. thesis, University of Tartu. Available at <http://attempto.ifi.uzh.ch/site/pubs/>.
- Kaljurand, K. and N. E. Fuchs (2006). Birectional mapping between OWL-DL and Attempto Controlled English. In *Proc. of the 4th Int. Workshop on Principles and Practice of Semantic Web Reasoning (PPSWR 2006)*.

- Keenan, E. and L. Faltz (1985). *Boolean Semantics for Natural Language*. Dordrecht: Reidel.
- Kittredge, R. I. (2003). Sublanguages and controlled languages. In R. Mitkov (Ed.), *The Oxford Handbook of Computational Linguistics*, pp. 430–447. Oxford University Press.
- Kurtonina, N. and M. Moortgat (1995). Structural control. In P. Blackburn and M. de Rijke (Eds.), *Logic, Structures and Syntax*. Dordrecht: Kluwer.
- Lambek, J. (1958). The mathematics of sentence structure. *American Mathematical Monthly* 65, 154–170.
- Moortgat, M. (1997). Categorical Type Logics. In J. van Benthem and A. ter Meulen (Eds.), *Handbook of Logic and Language*, pp. 93–178. Cambridge: MIT.
- Moot, R. (1998). Grail: An automated proof assistant for categorial grammar logics. In *Proc. of the 1998 User Interfaces for Theorem Provers Conf.*
- Pratt, I. and A. Third (2006). More fragments of language. *Notre Dame J. of Formal Logic* 47(2), 151–177.
- Schwitter, R., K. Kaljurand, A. Cregan, C. Dolbear, and G. Hart (2008). A comparison of three controlled natural languages for OWL 1.1. In *Proc. of the 4th Int. Workshop on OWL: Experiences and Directions (OWLED 2008)*.
- Schwitter, R., A. Ljungberg, and D. Hood (2003). ECOLE – A look-ahead editor for a controlled language. In *Proc. of the 8th Int. Workshop of the European Association for Machine Translation and the 4th Controlled Language Applications Workshop (EAMT/CLAW 2003)*.
- Schwitter, R. and M. Tilbrook (2006). Let’s talk in description logic via controlled natural language. In *Proc. of the 3rd Int. Workshop on Logic and Engineering of Natural Language Semantics (LENLS 2006)*.
- Thorne, C. (2010). *Query Answering over Ontologies Using Controlled Natural Languages*. Ph. D. thesis, Faculty of Computer Science, Free University of Bozen-Bolzano.
- Thorne, C. and D. Calvanese (2012). Tractability and intractability of controlled languages for data access. *Studia Logica* 100, 787–813.
- van Benthem, J. (1987). Categorical grammar and lambda calculus. In D. Skordev (Ed.), *Mathematical Logic and its Applications*, pp. 39–60. Plenum, New York.
- van Eijck, J. (1985). *Aspects of Quantification in Natural Language*. Ph. D. thesis, University of Groningen.
- Vardi, M. (1982). The complexity of relational query languages. In *Proc. of the 14th Annual ACM Symp. on Theory of Computing*, pp. 137–146.