

Alma Mater Studiorum Università di Bologna Archivio istituzionale della ricerca

Semantic Web Machine Reading with FRED

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version: Semantic Web Machine Reading with FRED / Gangemi A, Presutti V, Recupero DR, Nuzzolese AG, Draicchio F, Mongiovi M. - In: SEMANTIC WEB. - ISSN 1570-0844. - STAMPA. - 8:6(2017), pp. 873-893. [10.3233/SW-160240]

Availability: This version is available at: https://hdl.handle.net/11585/620509 since: 2020-02-28

Published:

DOI: http://doi.org/10.3233/SW-160240

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (https://cris.unibo.it/). When citing, please refer to the published version.

(Article begins on next page)

This is the final peer-reviewed accepted manuscript of:

Gangemi, A., Presutti, V., Reforgiato Recupero, D., Nuzzolese, A. G., Draicchio, F., & Mongiovì, M. (2017). Semantic web machine reading with FRED. Semantic Web, 8(6), 873-893, with permission from IOS Press <u>www.iospress.nl</u>

The final published version is available at IOS Press through

http://dx.doi.org/10.3233/SW-160240

Rights / License:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<u>https://cris.unibo.it/</u>)

When citing, please refer to the published version.

Semantic Web 0 (2016) 1–21 IOS Press

Semantic Web Machine Reading with FRED

Editor(s): Harith Alani, The Open University, UK

Solicited review(s): Philippe Cudré-Mauroux, Université de Fribourg, Switzerland; Antoine Zimmermann, École des Mines de Saint-Étienne, France; One anonymous reviewer

Aldo Gangemi ^{a,b,*}, Valentina Presutti ^b, Diego Reforgiato Recupero ^b, Andrea Giovanni Nuzzolese ^b, Francesco Draicchio ^b, and Misael Mongiovì ^b

^a LIPN, Université Paris 13, Sorbonne Paris Cité CNRS UMR 7030, France
 E-mail: aldo.gangemi@lipn.univ-paris13.fr ^b Semantic Technology Lab, ISTC-CNR, Rome and Catania, Italy
 E-mail: {aldo.gangemi,valentina.presutti,diego.reforgiato,andrea.nuzzolese,misael.mongiovi}@istc.cnr.it

Abstract. A machine reader is a tool able to transform natural language text to formal structured knowledge so as the latter can be interpreted by machines, according to a shared semantics. FRED is a machine reader for the semantic web: its output is a RDF/OWL graph, whose design is based on frame semantics. Nevertheless, FRED's graph are domain and task independent making the tool suitable to be used as a semantic middleware for domain- or task- specific applications. To serve this purpose, it is available both as REST service and as Python library. This paper provides details about FRED's capabilities, design issues, implementation and evaluation.

Keywords: Machine Reading, Knowledge Extraction, Semantic Web, Linked Data, Event extraction

1. Introduction

This paper describes FRED¹ [43]: a tool that automatically generates RDF/OWL ontologies and linked data from multilingual natural language text. The approach implemented by FRED falls into the *machine reading* paradigm [12], which aims to transform (part of) a natural language text into data. FRED adds to that paradigm the ability to generate knowledge graphs that can be interpreted by machines, according to a shared formal semantics, and is linked to available background knowledge.

As an example of a formal knowledge graph, the text "Valentina gave Aldo a book by Charlie Mingus." can be formalised according to OWL semantics and the OWL *n*-ary relation pattern², as sketched by Figure 1.

¹http://wit.istc.cnr.it/stlab-tools/fred ²http://ontologydesignpatterns.org/wiki/Submissions:N-Ary_Relation_Pattern_(OWL_2) The class Give and its instance give_1 represent the "gave" *n*-ary relation holding between "Valentina", "book", and "Aldo", which have their representative instances, in turn related to give_1.



Fig. 1.: Formalisation for the sentence: *Valentina gave Aldo a book by Charlie Mingus*, based on OWL semantics and the OWL *n*-ary relation pattern.

^{*}Corresponding author. E-mail: aldo.gangemi@lipn.univparis13.fr

FRED is able to produce such a formal knowledge representation, specifically for the semantic web. The tool leverages multiple natural language processing (NLP) components by integrating their outputs into a unified result, which is formalised as an RDF/OWL graph. Such a graph is enriched with links to existing semantic web knowledge, by means of ontology alignment and entity linking techniques, as well as with an RDF encoding of syntactic annotations based on the Earmark [39] and the NLP Interchange Format (NIF) [23] vocabularies.

The design of FRED's graphs is independent from any commitment on domain-specific or task-specific semantics, making the tool suitable to be used as a *semantic middleware*. In fact, FRED mainly targets developers of semantic applications, who can rely on FRED's output graphs and further manipulate them (e.g. refactor or enrich) for empowering their domainor task-specific client applications. Some examples are automatic text annotation, search engine optimisation, opinion mining, automatic summarisation. Examples of applications built on top of FRED are given in Section 4.

FRED¹ is available as RESTful API (providing RDF serialisation in many syntaxes) as well as Python API (*fredlib*³). Additionally, any user can access FRED's demo application online, which features a graphical user interface.

FRED's output graph is designed according to Frame Semantics [15] and ontology design patterns [20]. A frame is usually expressed by verbs or other linguistic constructions, hence all occurrences of frames that can be recognised in an input text are formalised as OWL n-ary relations, all being instances of some type of event or situation (e.g. the class Love in Figure 1). For example, FRED's output diagram⁴ for the sentence of Figure 1 is shown in Figure 2. The prefixes used in inline examples and diagrams are assigned as from Table 1. The reader may recognise the same formalisation shown in Figure 1. Nevertheless, the result produced by FRED is richer. The class fred:Give is modelled as equivalent to the frame vn.data:Give 13010100, defined in VerbNet [27]. The relations between the frame fred: Give and its arguments are modelled as object properties according to the semantic roles that can be

recognised, e.g. vn.role:Agent,

vn.role:Recipient, and vn.role:Theme. In case additional roles are detected but not recognised, FRED creates new (role) object properties and labels them by reusing the appropriate text from the input. Arguments are modelled as individuals and their types are induced from the input text, when available. If possible, individuals are linked to existing semantic web entities, e.g. dbpedia:Charles_Mingus, and typed with existing semantic web classes, e.g. schemaorg:Person, schemaorg:MusicGroup⁵ and dul:Event. The simple graph of Figure 2 exemplifies part of the extraction and modelling capabilities of FRED (frame detection, entity linking, type induction, etc.). In Section 2 such capabilities are discussed in detail.

In summary, the main contribution of FRED is to provide a novel form of machine reading. This is accomplished by addressing a number of challenges such as: providing a unique tool able to perform combined NLP tasks at a same time, integrating and enriching the results of diverse NLP tasks by drawing explicit relations between them, providing a unified formal representation compliant with semantic web design and standards (i.e. OWL/RDF), providing a reference cognitive semantics for its interpretation (i.e. frame semantics).

The paper is organised as follows. Section 2 presents the main capabilities of, and design issues addressed by, FRED. Section 3 describes the FRED's pipeline and provides implementation details. Section 4 shows the quality, importance and impact of FRED, reporting evaluation studies, community feedback, as well as examples and evaluation of applications that rely on FRED as a semantic middleware. Section 5 discusses relevant related work and the paper concludes with Section 6 that addresses open challenges and ongoing work.

2. Transforming Natural Language Processing output to OWL/RDF graphs

FRED leverages the results of many NLP components by reengineering and unifying them in a unique RDF/OWL graph designed by following semantic web ontology design practices. In this section, the main de-

³http://wit.istc.cnr.it/stlab-tools/fred/ fredlib

 $^{^{4}}$ FRED diagrams depict a subset of the generated triples, which cover the core semantics of the text.

⁵The reader may notice that **schemaorg:MusicGroup** is incorrect, however this data is inherited from existing Linked Data resources, which may contain some imprecisions.



Fig. 2.: FRED's output for the sentence: Valentina gave Aldo a book by Charlie Mingus..

sign issues addressed by FRED are discussed and exemplified, in order to provide an overall view of its main features. In order to improve readability of figures and examples, a set of prefixes is used in lieu of namespace URIs. They are summarised in Table 1.

2.1. From Discourse Representation Structures to RDF/OWL n-ary relations

The core of FRED takes as input Discourse Representation Structures (DRSs), based on Hans Kamp's Discourse Representation Theory (DRT) [25]. DRSs, informally called "boxes" (due to their graphical representation), represent natural language sentences, and include two parts: a set of discourse referents, and a set of conditions providing the interpretation of the discourse referents. The DRS language is within first order logic, and its discourse referents are arbitrary entities, including events modeled according to a neo-Davidsonian semantics [25].

The DRSs taken by FRED as input are produced by Boxer [5], which performs deep parsing out of Combinatory Categorial Grammar (CCG) parse trees [47]. It also makes use of both VerbNet [27] and FrameNet [4] for frame labelling and semantic role labelling, i.e. representing event types and the relations (thematic roles) between events and their participating individuals. Figure 3a exemplifies a box for the sentence "*People love movies*" as it is returned by Boxer: the box is divided into two sections, the top section contains the discourse referents, in this case x0, x1, x2; the bottom section contains the predicates that constrain their interpretation: x2 is an event described by the predicate love having two arguments, an agent x0, of type people and a patient x1, of type movie. For more details about Boxer and the syntax of its output, the reader can refer to [5].

Although the example in Figure 3a is very simple, DRSs can be very complex, depending on the input text. Therefore, the problem of transforming DRSs to RDF/OWL models requires non-trivial design decisions on how to represent discourse referents and their interpreting predicates. In the case of Figure 3a, there is one box, which encapsulates a single event controlling all discourse referents. Therefore, the box content can be represented as an RDF/OWL n-ary relation modelling the identified event (frame), with its arguments modelled as typed individuals, and thematic roles modelled as object properties: the corresponding FRED output is shown in Figure 3c. Events are modelled as subtypes of the class dul: Event. The reader may notice that FRED also performs skolemization of first-order predicate variables, e.g. x0 becomes fred:people_1, as well as alignment of predicates to existing semantic web ontologies, e.g with owl:equivalentClass triples. The skolemization of the event occurrence, its event type, and the labels for its semantic roles are given in two varieties following either VerbNet or FrameNet. In the first case, the event type fred:Love is

Default local (customisable) namespace	fred:	http://www.ontologydesignpatterns.org/ont/fred/domain.o				
VerbNet thematic roles	vn.role:	http://www.ontologydesignpatterns.org/ont/vn/abox/role/				
VerbNet verb classes	vn.data:	http://www.ontologydesignpatterns.org/ont/vn/data/				
FrameNet frame vocabulary	ff:	http://www.ontologydesignpatterns.org/ont/framenet/abox/frame/				
FrameNet frame element vocabulary	fe:	http://www.ontologydesignpatterns.org/ont/framenet/abox/fe				
DOLCE+DnS Ultra Light	dul:	http://www.ontologydesignpatterns.org/ont/dul/DUL.owl#				
WordNet	wn30:	http://www.w3.org/2006/03/wn/wn30/instances/				
a vocabulary for Boxer primitive classes	boxer:	http://ontologydesignpatterns.org/ont/boxer/boxer.owl#				
a vocabulary for some box types	boxing:	http://ontologydesignpatterns.org/ont/boxer/boxing.owl#				
DBpedia resources	dbpedia:	http://dbpedia.org/resource/				
schema.org vocabularies	schemaorg: http://schema.org/					

Table 1: List of prefixes used in FRED diagrams and inline examples.



(c) FRED output with FrameNet labels.

Fig. 3.: Boxer versus FRED result (in two varieties labelling with either VerbNet and FrameNet frames and roles) for the sentence *"People love movies"*, which exemplifies the case of a box that does not add any semantics to its content, which in turn maps to a specific frame.

aligned to vn.data:Love_31020100, and therefore roles are labeled as vn.role:Experiencer and vn.role:Theme respectively.

In the second case, the event type fred:Love is aligned to ff:Experiencer_focus, and there-fore roles are labeled as

fe:Experiencer:experiencer_focus and fe:Content:experiencer_focus respectively.

Natural language text on the web is rarely such simple as "*People love movies*", hence the possibile configurations of boxes and their contents have a wide range of complexity. In this example, the box itself does not add any special semantics to its content, but this is not always the case. One of the design issues addressed by FRED is to assess whether a box has its own semantics or not.

There are two main basic types of boxes that FRED needs to distinguish: (1) boxes only have a syntactic role in Boxer's result, meaning that FRED only needs to focus on representing their content and linking them to the rest; (2) boxes provide a unified relation to a complex state of affair (usually expressed in the text by a *copula*), meaning that they have their own semantics to be represented. Figure 3 is an example of (1), where the content of the box is represented as an *n*-ary relation expressing an event love_1 that keeps to-

gether the entities participating in the event. Hence, the box does not add any significant information: the love frame is the main informative and unifying element. Other examples of (1) are the sentences "Valentina is a researcher" and "Valentina is happy". In these cases, the output of Boxer has always the same structure, as depicted in Figure 4, which includes three boxes, one predicate (researcher) applied to a discourse referent (a named entity in the examples). The boxes here would indicate special states of affairs, e.g. the typing of an entity, and the relation of an entity to its quality, respectively. The boxes here are redundant as all discourse referents are controlled by the same and single predicate characterising their content. FRED detects such special situations and maps them to specific semantic relations after getting rid of the boxes: as for the examples in Figure 4 these are rdf:type in the first case (Figure 4c), and dul: hasQuality in the second case (Figure 4d). The reader may notice that Boxer tries to assign a type to the named entities that it detects, in the case of Figure 4a "Valentina" is assigned with the type loc, which stands for "Location". These assignments have a typically low accuracy in a sample-based validation, hence FRED bypasses them, and handles named entity typing by exploiting external semantic web resources when possible, or by simply leaving the individual untyped, otherwise.

The case (2) refers to those state of affairs that are expressed neither by means of events nor as typing/quality assertions (or similar). For example, let us consider the sentence: Valentina is Gianni's fifth daughter, from his second marriage. Figure 5 shows the results of Boxer and FRED for this sentence. In this case, the box declaring the x5 discourse referent has a specific semantic role: it frames the state of affairs (situation) expressed by the sentence. In fact, this box conveys the semantics of a unifying relation between all discourse referents and predicates contained in the first box by defining a co-referencing statement for x0and x^2 (by using a redundant box) and the from predicate. Therefore, FRED represents this box in the graph as an individual fred:situation_1 from the class boxing: Situation, which models a unified relational context for the individuals and relations corresponding to boxed content.

These basic patterns for boxes and their contents can be composed by means of formal relations such as *and*, *or*, *entails* between boxes. Furthermore, boxes can be nested, negated, or can include more than one event, or no events. These combinations and variety of content give rise to complex configurations that may cause the emergence of additional patterns to be handled by appropriate heuristic rules.

RDF vocabularies for DRSs. The development of FRED involved also the creation of two OWL vocabularies: one, referred by the prefix boxer: for representing a taxonomy of types defined and used by Boxer as first class objects, such taxonomy includes types such as "person", "organisation", "location", etc.; the other, referred by the prefix boxing:, for representing boxes (boxing:Situation), relations among boxes (i.e. boxing:entails, boxing:union), connecting boxes with their participating entities (e.g. boxing:involves), as well as other properties such as modality, negation, etc.

2.2. Representing tense, modality and negation

FRED also represents modality, tense and negation in its unified OWL/RDF graph, by identifying the corresponding patterns in Boxer output. Let us consider the example sentence "*Rahm Emanuel says he won't resign over police shooting*." and its corresponding FRED graph depicted in Figure 6.

Tense representation is addressed by modelling a time interval (e.g.fred:now_1 in Figure 6) that refers to what in the text is (possibly implicitly) expressed with the linguistic present tense: hence such "present" may not refer to the time at which the sentence was written or published, or the time of actual happening of the mentioned events. All events and situations expressed in the sentence are related to fred:now_1) by means of a relation from a set of object properties inspired by Allen's interval algebra (before, after, included, etc.), depending on the tense of the verb that expresses them. For example, in Figure 6, the event fred:say_1 is included in the fred:now_1 time interval as it is expressed with present tense, while fred:now_1 is "before" fred:resign_1, as the latter is expressed with future tense.

Modality and negation are represented with a lightweight, RDF-oriented semantics because the underlying natural language semantics is unpredictable, and poses controversial problems from both linguistic and philosophic perspectives [8].

As far as negation is concerned, given two sentences with a similar syntactic structure, the negated scope is often ambiguous to interpret. For example, the sentence "John did not go to school by car" can be formally represented in different ways. For example (us-



(a) Boxer output for a sentence expressing the type of an entity. (b) Boxer output for a sentence expressing the quality of an entity.



(c) FRED output for DRSs with one typing predicate.



(d) FRED output for DRSs with one quality predicate.

Fig. 4.: Boxer versus FRED results for examples of sentences expressing states of affairs (mostly) by means of the copula.

ing a neo-Davidsonian first order logic style), it can be represented as:

$$\neg \exists e(go(e, John, s, c) \land Event(e) \land School(s) \land Car(c))$$

meaning that "there is no event in which John went to school by car". The negation in this case is applied to the event *e*, whose arguments are *John*, the *school* and the *car*. There is no assumption as to what argument, if any, has a special role in the negated situation. Another possible formal representation of the same sentence is the following:

$$\exists e(go(e, John, s, c) \land Event(e) \land School(s) \land \neg Car(c))$$

meaning that "there is an event in which John went to school, but not by car". In this case the negation is applied to one of the arguments, i.e. the *car*, and an explicit assessment that the event *John went to school* happened, is made.⁶ To the best of our knowledge, in such cases it is impossible to automatically establish what is the correct scope without knowing about the extralinguistic context. We remark that even in simple cases, e.g. "John is not a doctor", the apparently obvious complement-based semantics is not reliable, since the linguistic or extralinguistic context can create alternative interpretations, e.g. "John is not a doctor, he is a donkey!"; let alone cases where negation actually provides a graded quality, such as "she is not unhappy", which means "she is not fully unhappy or happy".

This is why FRED only annotates (using the property boxing:hasTruthValue) the referenced event with the information that its truth value is false, without making any assumption on the impact of nega-

⁶More interpretations can be generated by swapping the negated argument.



(b) FRED output.

Fig. 5.: Boxer versus FRED result for the sentence "Valentina is Gianni's fifth daughter, from his second marriage", which exemplifies the case of a box carrying the semantics of a state of affairs.

tion over logical quantification and scope. The main limit of this approach is that the resulting representation does not trigger any automatic reasoning. The main benefit is that client applications have a hook to the negated event, and if needed, they can refactor the graph (e.g. with a SPARQL CONSTRUCT query) in order to express their desired interpretation, including rule-based or statistical approaches to decide what formal semantics to apply in each case.

As for modality, OWL lacks formal constructs to allow the required expressivity. The approach is then similar to the one for negation. Modality in FRED can be twofold: boxing:Necessary (corresponding to forms such as "will", "must", etc.) and boxing:Possible (for forms such as "may", "might", etc.). Both are individuals of the nominal class boxing:Modality. Note that the form "should" is represented with boxing:Necessary, by FRED. The English dictionary indicates two possible interpretations/uses of "should", one for indicating "obligation", the other for indicating something "probable". However, it also indicates that "in modern English uses of should are dominated by the senses relating to obligation". This is why FRED represents occurrences of should as expressions of "necessary" modality.

For the sentence of Figure 6 the triples:

:resign_1 boxing:hasTruthValue boxing:False ;
 boxing:hasModality boxing:Necessary .

formalise the "will not resign" fragment, including modality and negation.



Fig. 6.: FRED graph for the sentence: Rahm Emanuel says he won't resign over police shooting.

2.3. Compositional semantics, taxonomy induction and quality representation

As the output of Boxer does not tag compound terms explicitly, FRED extracts them by recognising two main patterns in the input DRSs. Given a compound term "term1 term2" the two predicates forming it ("term1" and "term2") can be represented by Boxer either as a dependency relation (i.e. term1(x0), term2(x1), nn(x0, x1)), or as two co-referent predicates (i.e. term1(x0), term2(x0)). FRED implements a compositional semantics on identified compound terms, exemplified in Figure 6 for the expression *police shooting*.

FRED creates a class representing the compound term, e.g., fred:PoliceShooting, then it builds its corresponding taxonomy, e.g.,

```
fred:PoliceShooting
  rdfs:subClassOf fred:Shooting .
```

based on the rationale that the modifier of the main concept ("police") provides a distinguishing feature to the more specific class. This feature used to be called in ancient logic *differentia specifica*, and can be either a quality (typically expressed by an adjective or an adverb), or another concept (typically expressed by a noun). In the latter case, a new class is created for the concept, e.g. fred:Police, together with a triple associating such a class with the original class derived from the compound noun, e.g. fred:PoliceShooting. To this aim, the DOLCE property dul:associatedWith is used. The following triples (represented in Turtle notation) show this compositional FRED design pattern applied to the example sentence in Figure 6:

```
fred:PoliceShooting
  rdfs:subClassOf fred:Shooting ;
  dul:associatedWith fred:Police .
```

The case of adjectives and adverbs expressing a differentia specifica is different from nouns. Adjectives in particular have an unpredictable and unstable semantics as nicely explained in [35]. FRED represents such modifiers (i.e. adjectives) as qualities of the modified term by means of the DOLCE property dul:hasQuality. But the quality can alternatively modify the individual denoted by the term, or the class expressed by that term. In a recent work, [19] reinterprets and augments Morzycki's distinctions, and provides an ontology of adjectives with a set of associated knowledge representation patterns, based on two main aspects of adjectives: sectivity and framality. Sectivity impacts on the semantics of classes and individuals that are modified by adjectives, while framality is the ability of an adjective to activate a conceptual frame (in Fillmore's sense [15]), and can be used for explaining sectivity phenomena. FRED implements an algorithm that leverages the theory and resources developed in [19] for applying the most appropriate representation pattern in presence of adjective occurrences. For example, FRED is able to distinguish the two main types of adjectives:

- intersective: the adjective that modifies a noun can be independently predicated of the individual, via entailment. For example, in a graph formalising the sentence: "Roberto Bolle is an Italian dancer.", we would find the following triples:

```
fred:Roberto_Bolle
  rdf:type fred:ItalianDancer,
  dul:hasQuality fred:Italian .
```

subsective: the adjective that modifies a noun cannot be independently predicated of an individual via entailment. For example, referring to Figure 5, we find the following triples:

```
fred:FifthDaughter
  rdfs:subClassOf fred:Daughter,
  dul:hasQuality fred:Fifth .
```

where the quality does not modify the individual, but the class: it is represented as an intensional quality, exploiting the "punning" pattern available in OWL2.

In the special case of Figure 5, the fragment *fifth daughter* is part of a "periphrastic expression", which leads to a special representation in FRED, as explained in the next section:

```
fred:Valentina
    fred:fifthDaughterOf fred:Gianni .
```

2.4. Generating periphrastic relations

There are many cases of relations that are expressed (and annotated by NLP tools) by means of prepositions, e.g. of, with, for, in, etc. Naming an object property in a semantic web graph with one of those terms, e.g. of, results meaningless and potentially identical to many other relations that actually may mean a completely different concept. For example, "survivor of" has a completely different meaning from "sister of", although they include the same preposition. In those cases FRED performs a paraphrasing task by identifying the noun to be associated with the preposition, and putting it before the preposition in order to form the label of the resulting relation. For example, the sentence "He was the only survivor of the expedition." would be formalised by generating the following triples:

```
fred:survivor_1
    fred:survivorOf fred:expedition_1;
    rdf:type fred:Survivor .
fred:expedition_1 rdf:type fred:Expedition .
```

2.5. Named Entity Recognition, Entity Linking, and Coreference Resolution

Named Entity Recognition (NER) is used for identifying elements in a text that should have a corresponding OWL individual in the graph. FRED also integrates the results of Entity Linking (EL) performed on the input text for enriching its output graph with owl:sameAs axioms. For example, in Figure 6, the reader may notice the following triple:

```
fred:Rahm_emanuel
    owl:sameAs dbpedia:Rahm_Emanuel .
```

Co-reference resolution and role propagation output is used for merging nodes, for example in Figure 6 the reader may notice that the individual fred:Rahm_emanuel is both the agent of the event say_1 and of the event resign_1, while in the text

say_1 and of the event resign_1, while in the text he was referred first by his name, and then by the pronoun "he".

2.6. Word Sense Disambiguation for ontology alignment

FRED produces RDF/OWL ontologies having classes (and related taxonomies) depending on the lexicon used in the text. FRED exploits word sense disambiguation (WSD) in order to provide a public identity to these classes by identifying equivalent or more general concepts into WordNet and BabelNet [36], and by creating alignments, where appropriate. WSD also enables FRED to generate alignments to two top-level ontologies: WordNet "supersenses" and a subset of DOLCE+DnS Ultra Lite (DUL) classes. For example, the term *programming language* is formalised by the following alignment axioms:

```
fred:ProgrammingLanguage
  owl:equivalentClass
    wn30:synset-programming-language-noun-1 ;
  rdfs:subClassOf dul:InformationEntity ;
  rdfs:subClassOf
    wn30:supersense-noun_communication .
```

Since Wikipedia is rich in "conceptual" entities, EL is also used for disambiguating the sense of words after a process of contextualisation and formal interpretation: if a text segment annotated by EL was originally annotated by FRED as an owl:Class, the resolved DBpedia entity is also given the semantics of an owl:Class. Therefore, if an individual in the graph is typed by the original FRED class, it will be typed by the DBpedia entity as well (FRED here applies an inheritance pattern). This ends out also to be a way of coercing the semantics of DBpedia entities where applicable.

2.7. Other FRED's capabilities

Multilingualism. FRED takes as input a text in one of 48 different languages, it translates it in English and then processes it for producing its corresponding graph. Therefore, the resulting graph always has English labels. As for the translation, FRED relies on Bing Translation APIs⁷. If the input language is different from English, the input text must be preceded by the tag $\langle BING_LANG: lang \rangle$, where *lang* is the code for the language⁸.

Textual annotation grounding. As part of its OWL/RDF output, FRED provides annotations that link text fragments to their corresponding graph elements. These annotations are expressed by means of the Earmark vocabulary [39] and the NLP Interchange Format (NIF) [23].

3. FRED pipeline and implementation

FRED pipeline is depicted in Figure 7. In this section this pipeline is described by providing details about the actual components that realise its implementation. FRED pipeline is characterised by three main phases: text processing, heuristic-based triplification, and RDF graph enrichment.

3.1. Text processing.

FRED takes as input a text in natural language (NL). It can be a short text (to be processed at once) or a corpus of NL documents; at the moment, the latter is available only by using the *fredlib* Python API. The input text is processed and transformed into DRSs. This processing includes also frame detection and semantic role labelling (based on Verb-Net [27] and FrameNet [4]), identification of relations between frames, named entity recognition (NER) and coreference resolution (CRR). Boxer [5] is in charge of producing DRSs (see Section 2 for more details

⁷http://www.microsoft.com/web/post/ using-the-free-bing-translation-apis

⁸The website http://msdn.microsoft.com/en-us/ library/hh456380.aspx provides the list of language codes about Boxer) nevertheless, after observing a significant number of cases during FRED implementation and experience of usage, it emerged that Boxer's pronoun CRR capability is limited. For this reason FRED also integrates CoreNLP9 as an additional component for this specific task. A similar approach is used for NER, in fact FRED integrates TAGME [14], which uses Wikipedia content as context to disambiguate named entities. For both tasks, i.e. NER and CCR, FRED takes as input the union of Boxer, TAGME and CoreNLP outputs, where TAGME overrules Boxer for NER (when both provide results) and CoreNLP overrules Boxer for CRR, in an analogous situation. The output of TAGME is reused also later (in the RDF graph enrichment phase) for performing Entity Linking. The result of this processing is transformed into an intermediate representation and passed to the second phase.

Implementation. This phase is implemented as a Python software component that uses system calls for getting Boxer, TAGME and CoreNLP outputs, integrate them in an internal representation and pass the result to the next component.

3.2. Heuristic triplification.

In this phase a first version of the graph is created. The component in charge of this task is implemented as a manager of ~100 heuristics that remove redundancies, associate a OWL/RDF representation to each identified pattern in the input DRSs (described in Section 2), and combines them. At this stage, FRED performs taxonomy induction, variable reification, role propagation, periphrastic relation extraction, frame and situation modelling, semantic role labelling, lightweight representation of negation and modality, tense modeling, and representation of both individual and intensional qualities. Figure 8 shows the result of this phase for the sentence "Rahm Emanuel says he won't resign over police shooting.". The reader may notice that the type assigned by Boxer is kept in this preliminary model and represented by means of a boxer:possibleType triple, which will be evaluated and possibly discarded in the next phase. Two events are modeled as n-ary relations based on the detected VerbNet frames "Say" and "Resign" and the result of semantic role labeling. Negation and modality

⁹http://nlp.stanford.edu/software/corenlp. shtml



Fig. 7.: FRED pipeline.

are also represented for the event instance resign_1. Finally, a taxonomy is induced and added to the graph due to the recognition of the compound term "police shooting".

Implementation. This phase is realised by a set of Python modules that implement the heuristics described in Section 2. The resulting component is deployed as a REST service, meaning that it is possible to programmatically access the output of the heuristic-based triplification. Additional triples are generated in order to annotate the fragments from a text with the semantic entities that are extracted from the text, jointly with their syntactic part-of-speech annotations. The Earmark [39], NIF [23], and semiotics.owl¹⁰ vocabularies are used in order to provide both a *semiotic* and a processing interoperability to FRED graphs. Text fragment are represented as offsets (e.g. fred:offset_28_36_daughter), and annotated as summarized in Table 2:

3.3. Graph enrichment.

The result of the previous phase is the input of the final phase, which has the goal to further enrich the

¹⁰http://www.ontologydesignpatterns.org/cp/
owl/semiotics.owl

Gangemi et al. / Semantic Web Machine Reading with FRED

Annotation	Earmark	NIF	Labels, syntax, semiotics	Example
Offset type	earmark:PointerRange	nif-core:OffsetBasedString		offset a earmark:PointerRange
Context	earmark:refersTo	nif-core:referenceContext		offset nif-core:referenceContext fred:docuverse
Segment start	earmark:begins	nif-core:beginIndex		offset earmark:begins "28"^^ xsd:nonNegativeInteger
Segment end	earmark:ends	nif-core:endIndex		offset nif-core:endIndex "36"^^ xsd:nonNegativeInteger
Labeling			rdfs:label	<pre>offset rdfs:label "daughter"^^xsd:string</pre>
Part-of-speech			pos:pennpos	offset pos:pennpos pos.owl:NN
Denotation			semiotics.owl:denotes	<pre>offset semiotics.owl:denotes fred:daughter_1</pre>
Interpretant			semiotics.owl:hasInterpretant	offset semiotics.owl:hasInterpretant fred:Daughter

Table 2: Annotation properties in FRED. The fred:offset_28_36_daughter offset is shortened as "offset".



Fig. 8.: Output of the second phase of FRED pipeline (heuristic triplification) for the sentence: "*Rahm Emanuel says he won't resign over police shooting.*"

graph with links to existing semantic web resources and textual annotations as well as to enrich compositional associations (dul:associatedWith triples), and to perform validation of, and possible corrections on, the RDF format. The results of TAGME are reused in order to produce owl:sameAs triples for recognised entities that have a correspondence in e.g. DBpedia. Word sense disambiguation (WSD) is used for producing alignments with classes from external ontologies (expressed by means of rdfs:subClassOf and owl:equivalentClass axioms). This task is addressed by integrating an external component: at the moment FRED reuses UKB [1]; an integration with Babelfy [34] is under development as well. The result of this phase is the final graph, which is provided in RDF format: FRED supports most RDF serialisation formats.

Implementation. The third phase is implemented as an independent Java software framework, named K~ore. K~ore is a modular set of Java components, each accessible via its own RESTful interface. All components are implemented as OSGi [38] bundles, components and services based on Apache Felix¹¹. For example, the WSD service has been implemented as an OSGi wrapper for UKB [1]. Pragmatically, K~ore

¹¹http://felix.apache.org/

wraps FRED's Python components, and enriches their output by means of additional services implemented as OSGi bundles. The whole wrapper is deployed as a REST service.

Additional implementation and deployment details. FRED is also released as a Python API named *fredlib*, which relies on the K~ore REST services, and allows to query FRED with a user-specified corpus, also enabling the manipulation of the resulting graph. More specifically, a Python function hides details related to the communication with the FRED service, and returns the user a FRED graph object that is easily manageable. FRED graph objects expose methods for extracting qualified parts (motifs) from the graph. These methods include functions for extracting the set of individual and class nodes, equivalence axioms, typing axioms, categories of FRED nodes (i.e. events, situations, qualities, general concepts) and categories of edges (roles and non-roles). In addition, fredlib supports rdflib¹² (for managing RDF graphs) and net $workx^{13}$ (for managing complex networks) libraries.

4. Quality, Importance, Impact

FRED's quality and impact can be assessed from different perspectives: (i) its ability to perform specific knowledge extraction tasks such as event detection, taxonomy induction, etc. (ii) its popularity in terms of reuse and citations by a wider community of developers; (iii) its performance as a semantic middleware, based on the evaluation of applications built on top of it.

4.1. Rigorous evaluation

As far as (i) is concerned: [43] reports FRED performance on the *frame detection* task showing that it is one order of magnitude faster than Semafor [11] (see Figure 9), i.e. the best state-of-art tool at the time. Table 3 summarises the accuracy performance of the two tools: FRED and Semafor have comparable precision values for frame detection, while the value of recall is lower for FRED. Nevertheless, it has to be noted that Semafor was trained on the FrameNet reference corpus, which put it in strong advantage with respect to FRED. As additional functionality compared to Semafor, FRED provides a formal representation of the identified frame occurrences.



Fig. 9.: Time to provide answers in function of the number of sentences per document as reported in [43].

A further study [16] reports a comparison between information extraction tools, including FRED. The author defines a number of *basic* semantic tasks by providing a correspondence between NLP tasks and semantic web terminology. Table 4 reports the list of tasks with a brief explanation of such correspondences: each NLP task is informally associated with a corresponding OWL-based semantics by indicating the type of triples that may be produced starting from its output.

As many of the analysed tools provide a non-RDF output, in order to allow their comparison, a manual conversion to RDF graphs was performed, based on a reference translation table. The study compares fifteen tools (including FRED) by assessing their coverage of, and performance on, the listed tasks against a gold standard of 524 triples produced from a news text. The result shows that FRED has the largest coverage of tasks and best accuracy performance for some of them. Table 5 summarises this evaluation analysis, the reader can refer to [16] for a detailed description of each tool (here we only provide an external reference for each of them), and to online updated data¹⁴ for details about performance measures.

⁴ http://stlab.istc.cnr.it/stlab/KnowledgeExtractionToolEv						
Tool	Precision	Recall	F-Score			
FRED	75.320	57.519	65.227			
Semafor	75.325	74.797	75.060			

Table 3: Frame detection task: performance comparison between FRED and Semafor as reported in [43].

¹²http://code.google.com/p/rdflib/

¹³https://networkx.github.io/

Task	NLP terms	Semantic Web triples			
ТорЕ	Topic Extraction	dc:subject			
NER	Named Entity Recognition	owl:NamedIndividual			
NEReS	Named Entity Resolution	owl:sameAs			
TE	Terminology extraction	owl:Class owl:ObjectProperty owl:DatatypeProperty			
TReS	Terminology resolution	owl:equivalentTo∥rdfs:subClassOf∥rdfs:subPropertyOf			
Senses	Sense tagging	rdf:type			
Tax	Taxonomy induction	rdfs:subClassOf			
RE	Relation Extraction	owl:ObjectProperty owl:DatatypeProperty			
Events/Roles	Event detection and SRL	<event> rdf:type <event.type> . <event> <semrole_i> <entity_j></entity_j></semrole_i></event></event.type></event>			
Frames	Frame detection	<event.type> rdfs:subClassOf <frame/></event.type>			

Table 4: Summary of basic semantic tasks.

The table reports a value for each basic task indicating that a certain tool addresses that task (coverage) with a certain accuracy performance. Accuracy is computed as the sum of true (positive and negative) results over the sum of all (true and false) results. If a value is not present it means that it was either not computed or not computable (i.e. task not addressed), yet a "–" sign indicates "not addressed" while a "+" sign indicates "addressed".

4.2. Community feedback

Besides rigorous evaluation, FRED's quality is supported by evidence of its impact in the community (ii). A public forum in the software engineering community, stackoverflow.com, contains independent discussions about extending FRED's usage to other platforms (Python, C#)¹⁵. Another forum, answers.semanticweb.com, contains independent analyses and discussion of related works, revealing FRED's uniqueness in providing a solution for producing rich RDF datasets from text¹⁶. The Wikipedia entry for Knowledge Extraction¹⁷ contains a list of tools for knowledge extraction from text, and even looking only at the column for "extracted entities", FRED has the largest set. The only comparable entry is for a proprietary tool that does not provide any demo on the web.

More literature shows the impact of FRED beyond semantic technology circles: a 2014 study about deal-

ing with big data for statistics made by the United Nations Economic Commission for Europe¹⁸ says that

"The knowledge extraction from unstructured text is still a hard task, though some preliminary automated tools are already available. For instance, the tool FRED (http://wit.istc.cnr. it/stlab-tools/fred) permits to extract an ontology from sentences in natural language." A 2013 book [2] on big data computing says in the introduction: "The technologies for this [knowledge extraction] are under intensive development currently, for example wit.istc.cnr.it/stlab-tools/ fred (accessed October 8, 2012)".

4.3. Evaluation as a semantic middleware

Since FRED is a semantic middleware (iii) its quality is primarily assessed by evaluating the performance of applications that depend on it. This is a hard task to address because a rigorous methodology would require to assess the performance of each such application both with and without using the middleware. Nevertheless, we argue that FRED's impact as a middleware can be demonstrated by showing successful results of applications depending on it that address a broad range of tasks, hence reducing the possibility that their individual success is only due to other factors. In support to this claim, it has to be noted that all tools discussed in the next sections strongly depends on FRED's output: the core of their logic relies on heuristics defined based on FRED's graph design. In other words, these tools are specialisations of FRED for specific domains or tasks. The reminder of this sec-

¹⁵http://tinyurl.com/qa2dyfj,

http://tinyurl.com/o993scy

¹⁶http://tinyurl.com/n6pzpot,

http://tinyurl.com/kb8564w

¹⁷http://en.wikipedia.org/wiki/Knowledge_ extraction

¹⁸http://tinyurl.com/ml6ystn

Tool	TopE	NER	NEReS	TE	TReS	Senses	Tax	RE	Events	Roles and Frames
AIDA ^a	-	.89	.80	-	_	.64	-	-	_	-
Alchemy ^b	.52	.89	_	.20	_	.64	-	.30	_	-
Apache Stanbol ^c	-	.77	.25	-	-	.50	-	-	_	-
CiceroLite ^d	-	.89	.75	.21	.07	.64	-	.25	.18	.22
DB Spotlight ^e	-	.79	.55	-	-	.42	-	-	-	-
FOX ^f	-	.86	.75	.33	.65	.57	-	-	-	-
FRED	-	.84	.60	.90	.07	.48	+	.82	.87	.69
NERD ^g	-	.88	.60	-	-	.69	-	-	-	-
Open Calais ^h	.48	.82	-	-	-	.57	-	-	.04	-
PoolParty KD ⁱ	.28	_	-	-	-	-	-	-	-	-
ReVerb ^j	-	_	-	-	-	-	-	.27	-	-
Semiosearch ^k	-	-	.60	-	.46	-	-	-	-	-
Wikimeta ^l	-	.86	.75	.04	.07	.80	-	-	_	-
Zemanta ^m	-	.93	_	-	—	.27	-	_	_	_

Table 5: Summary of evaluation results for basic tasks indicating accuracy values in the interval [0,1] with 1 expressing the best possible accuracy.

```
ahttp://www.mpi-inf.mpg.de/yago-naga/aida/
bhttp://www.alchemyapi.com/api/demo.html
chttp://dev.iks-project.eu:8081/enhancer
dhttp://demo.languagecomputer.com/cicerolite
ehttp://dbpedia-spotlight.github.com/demo
fhttp://aksw.org/Projects/FOX.html
ghttp://nerd.eurecom.fr
hhttp://viewer.opencalais.com/
ihttp://poolparty.biz/demozone/general
jhttp://reverb.cs.washington.edu
khttp://wik.istc.cnr.it/stlab-tools/wikifier
lhttp://www.wikimeta.com/wapi/semtag.pl
"http://www.zemanta.com/demo/
```

tion is dedicated to briefly describe three such tools addressing different tasks and their performance. Notably, FRED has been used also in [24] to automatically extract the meaning of citations in scientific research articles, and in [32] for supporting a semantic web approach to textual knowledge reconciliation.

4.4. Semantic sentiment analysis

Sentilo¹⁹ [21][44] is a semantic sentiment analysis tool built on top of FRED. The result of Sentilo is an enriched FRED graph with annotations based on an opinion-specific vocabulary. Given an opinion sentence, it runs a set of heuristics on both the input text and the FRED graph for identifying and classifying the graph entities that represent the holders of the expressed opinion (if any), the opinion topics, and the opinion-expressing words. Furthermore, it assigns sentiment scores to the expressed opinions. Sentilo also relies on additional external resources such as Sentic-Net [6], SentiWordNet [3] and Levin++²⁰ (a resource based on Levin classification of verbs [29]).

Figure 10 shows the pipeline of Sentilo. Its extension to FRED includes

- an ontology for opinion sentences for enriching FRED graphs;
- SentiloNet, a novel (frame-based) lexical resource that enables the evaluation of opinions expressed by means of events;
- a novel scoring algorithm for opinion sentences.

¹⁹http://wit.istc.cnr.it/stlab-tools/ sentilo

²⁰http://www.stlab.istc.cnr.it/documents/ sentilo/levin-opinion.zip

Sentilo is accessible as both REST service and web application featuring a graphical user interface.

An evaluation conducted on a corpus of open-rating reviews about hotels, reported in the cited papers, shows high accuracy of the system for the different addressed tasks: holder detection (F1 = .95), topic detection (F1 = .66), subtopic detection (F1 = .80), and sentiment scoring (Pearson $\rho = .81$, the correlation between Sentilo scores and open-rating review scores).



Fig. 10.: Pipeline of Sentilo (taken from [44]).

4.5. Extraction of link semantics

 $Legalo^{21}[42]$ is a novel approach that relies on FRED to automatically generate OWL properties that express the semantics of hyperlinks, and possibly align them to existing ontology properties. Given a sentence including hyperlinks, Legalo is able (i) to identify pair of entities that are relevantly related according to the meaning of the sentence (as well as to discard entity pairs that are not relevantly related), (ii) to generate a label summarising the meaning of such relation, (iii) to formally define an OWL property representing such relation, annotated with the generated label, and (iv) to align the learned OWL property to existing OWL ontologies (if possible). Legalo enriches FRED graphs with binary relations between entities that are only indirectly connected in the original graph, but whose relation is relevant in the context of the sentence. In order to address these tasks, Legalo defines a set of heuristics based on recurrent graph-patterns observed in FRED output.

Figure 11 depicts the pipeline implemented by Legalo. The system also relies on external resources





Fig. 11.: Pipeline of Legalo. Numbers indicate the order of execution of a component in the pipeline.

such as Watson²² [10], LOV²³, and NELL²⁴ [7], for addressing the property alignment task. The performance of Legalo with respect to each addressed task was evaluated with the help of crowdsourcing showing very promising results: relation relevance assessment (F1 = .87), usability of generated label (F1 = .78), property alignment (*Precision* = .84). Furthermore, the automatically generated labels were compared to human generated labels by computing a similarity score (Sim = 0.80) with the SimLibrary framework [41]²⁵ (the interval value of Sim is [0, 1], the higher the score, the more similar the two phrases).

4.6. Automatic entity typing

*Tipalo*²⁶ [18] relies on FRED to automatically type DBpedia entities based on their natural language definitions. An entity definition is extracted from its corresponding Wikipedia page and transformed by FRED in a graph representation. Tipalo defines a set of heuristics for extracting taxonomies of types that classify an entity, from FRED output. It also disambiguates the sense of learned concepts and align them to existing ontologies, hence providing an alternative to DBpedia and YAGO [48]. Figure 12 shows the pipeline of Tipalo, which includes external resources such as UKB [1] for word sense disambiguation, OntoWord-

²²http://watson.kmi.open.ac.uk/WatsonWUI/ ²³http://lov.okfn.org/dataset/lov/

²⁴http://rtw.ml.cmu.edu/rtw/

²⁵http://simlibrary.wordpress.com/

²⁶http://wit.istc.cnr.it/stlab-tools/

tipalo/



Fig. 12.: Pipeline of Tipalo. Numbers indicate the order of execution of a component in the pipeline.

Net [17] and WordNet 3.0 Supersenses RDF²⁷ for concept alignment. The tool was evaluated against a sample of DBpedia entities and showed high accuracy for the entity typing tasks [18]: type selection (F1 = .92), type (concept) sense disambiguation (F1 = .75). Concepts are extracted from the original definitions used in Wikipedia, hence they provide an alternative to DB-pedia and YAGO [48].

5. Related work

The work on Open Information Extraction (OIE, [13]) is the foundation of *machine reading* and relies on an open domain and unsupervised paradigm. The main antecedent to OIE is probably the 1999 Open Mind Common Sense project [46], which adopted an ante-litteram crowdsourcing and games-with-purpose approach to populate a large informal knowledge base of facts expressed in triplet-based natural language. Another OIE project is Never Ending Language Learning (NELL) [31], a learning tool that since 2010 processes the web for building an evolving knowledge base of facts, categories and relations. In this case there is a (shallow) attempt to build a structured ontology of recognised entities and predicates from the facts learnt by NELL. Ontology learning [9] aims to address

a similar task as machine reading: it uncovers statistical regularities in linguistic features from large corpora, which could justify e.g. a subsumption or disjointness relations, in order to generate logical axioms. Mostly, the results are sparse and rely on shallow parsing. Works such as [5,33] assume an axiomatic form and make the extraction process converge to that form. In these cases, although the output is formalised its transformation to semantic web languages is far from being straightforward.

FRED provides a means to perform machine reading where the result is formally represented according to semantic web standards and design, carrying the semantics of cognitive linguistics frame and enabling automatic reasoning and reuse from other software agents. We call this variety of OIE *Open Knowledge Extraction* [42].

Most research work in the area of natural language processing (NLP) is relevantly related to machine reading. NLP research is characterised by the fact that the developed methods focus on specific tasks such as relation extraction, named entity recognition, frame detection, semantic role labeling, etc. Furthermore, the formal representation of NLP methods' results is mostly overlooked in their development. Therefore, two among the challenges for advancing the state of the art in machine reading are: (i) to address combined NLP tasks at the same time and (ii) to identify a unified formal representation for the output. Among FRED's merits is the ability to address these challenges: it leverages different existing NLP methods and unifies their results into a formal representation.

The integration between NLP and semantic web (often referred to as "semantic technologies") is progressing fast. Most work has been opportunistic: on one hand exploiting NLP algorithms and applications, such as named-entity recognisers and sense taggers, to populate semantic web datasets or ontologies, or for creating NL query interfaces; on the other hand exploiting large semantic web datasets and ontologies (DBpedia, YAGO, Freebase) to improve NLP algorithms. For example, Alchemy API²⁸, Apache Stanbol²⁹, NERD [45], and FOX³⁰ perform grounding of extracted entities in publicly available identities such as Wikipedia, DBpedia and Freebase. Nevertheless, their output lacks information about how the identified

27http://semanticweb.cs.vu.nl/lod/wn30/

²⁸http://www.alchemyapi.com

²⁹http://stanbol.apache.org

³⁰http://aksw.org/Projects/FOX.html

entities are related to each other, which is one of the features provided by FRED.

Relation extraction and question answering are mostly domain-dependent and supervised: they use a finite vocabulary of predicates (e.g. from the DBpedia ontology), and rely on their extensional interpretation in data (e.g. DBpedia) to either link two entities recognised in some text (as in [28,26]), or to find an answer to a question, from which some entities have been recognised (as in [30]). FRED performs a different form of relation extraction: it identifies what terms convey the semantics of a relation and creates a formal representation for them in the form of an OWL object property. Currently it is partly independent from any existing source of properties, which leads its behaviour to generate redundancies. However, its goal is not to detect existing relations, but to formalise the ones it recognises as such. An extension of its capabilities to cope with redundancies may be desirable as future work.

As the integration between NLP and semantic web is becoming tighter, clearer practices about how to represent linguistic data are strongly desirable. Some work propose linked data versions of linguistic resources such as WordNet³¹ [17,49] and FrameNet [37], and the recent proposal of Ontolex-Lemon by the Ontolex W3C Community Group³² aims to improve linguistic resource reuse. Some existing tools such as Apache Stanbol³³, NERD [45], and FOX³⁴ provide RDF descriptions of their outputs, making it easier to reuse them. The NLP Interchange Format (NIF) [23] and EARMARK [40] provide RDF vocabularies for annotating text fragments, hence supporting the integration of results from different NLP tools. FRED graphs include both NIF and EARMARK annotations.

6. Discussion and conclusions

This paper describes FRED: a tool able to provide a formal representation of natural language text based on semantic web design principles and technologies. After providing a detailed description of FRED's architecture and capabilities, the paper reports on its impact by showing both rigorous evaluations and in-

³³http://stanbol.apache.org

dependent quotes from a wider community of researchers, adopters and practitioners. These references demonstrate that FRED currently stands as the noncommercial tool having the largest coverage of formally defined tasks, and, to our knowledge, the largest coverage for the semantic web specifically. The final aim motivating FRED development is to support natural language understanding.

FRED currently responds to a very important challenge: to leverage existing natural language methods and tools in order to obtain a unified, formalised representation of both facts and concepts expressed by a natural language text. Most natural language tasks are addressed by specialised tools separately, and the semantic assumptions behind their output is neither harmonised nor formalised in most (if not all) cases. The main issues implied by developing FRED were to identify such a unified formal interpretation and to minimise the heuristic rules needed for producing a sound formal result, when combining the diverse discourse patterns.

There are still important open issues in dealing with discourse phenomena: they are more diverse and broader than what can be currently extracted and represented. Some of the challenges depend on implicit knowledge, others on ambiguity, and some more depend on higher order modal and conceptual structures.

Ongoing work for extending FRED is dealing with some of them, for example certain kinds of coercion, adjective semantics, polarity, sentiment, frame composition, a subset of presuppositions and paraphrases, etc. For some (e.g. polarity, sentiment) very promising results have been achieved [21][44].

The current state of the art for machine reading, either grounded in the semantic web or not, is still at an early stage when compared to the grand vision of natural language understanding. Two relevant examples of difficult tasks are: (1) the accurate extraction of implicit discourse relations and conventional implicatures, which do not only require background knowledge, but also reasoning on that knowledge in a way close to the appropriate discourse level [22]; (2) the recognition of cultural framing out of real world facts, as in political discourse, which requires the extraction, representation, and reasoning over high-level frames (attitudes, values, metaphors), which tend to control the factual frames that are currently the most complex grasp offered by automated discourse representation.

A current focus of our work is creating large repositories of FRED graphs, using typed named graphs and reconciliation techniques [32] for the cases when the

³¹http://semanticweb.cs.vu.nl/lod/wn30/ ³²http://www.w3.org/community/ontolex/wiki/ Main_Page

³⁴http://aksw.org/Projects/FOX.html

source texts are related for some reason, e.g. with news series, large texts, abstracts of categorised scientific articles, etc. The final goal is to produce a large repository of knowledge graphs that can be used to perform deep and formal annotation of large archives of documents, and to automatically produce formal relations between them. Another ongoing evolution of FRED is in the area of robot-human interaction, where FRED graphs extracted from natural language dialogues need to be grounded to physical environments.

Acknowledgements

The research leading to these results has received funding from the European Union Horizons 2020 the Framework Programme for Research and Innovation (2014-2020) under grant agreement 643808 Project MARIO Managing active and healthy aging with use of caring service robots.

References

- [1] E. Agirre and A. Soroa. Personalizing PageRank for word sense disambiguation. In A. Lascarides, C. Gardent, and J. Nivre, editors, EACL 2009, 12th Conference of the European Chapter of the Association for Computational Linguistics, Proceedings of the Conference, Athens, Greece, March 30 - April 3, 2009, pages 33–41, Athens, Greece, 2009. The Association for Computer Linguistics. http://www.aclweb. org/anthology/E09-1005.
- [2] R. Akerkar. *Big Data Computing*. CRC Press, 2013. DOI:10.1201/b16014.
- [3] S. Baccianella, A. Esuli, and F. Sebastiani. SentiWord-Net 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In N. Calzolari, K. Choukri, B. Maegaard, J. Mariani, J. Odijk, S. Piperidis, M. Rosner, and D. Tapias, editors, *Proceedings of the International Conference on Language Resources and Evaluation*, *LREC 2010, 17-23 May 2010, Valletta, Malta*, pages 2200– 2204, Valletta, Malta, 2010. European Language Resources Association (ELRA). http://www.lrec-conf.org/ proceedings/lrec2010/summaries/769.html.
- [4] C. F. Baker, C. J. Fillmore, and J. B. Lowe. The Berkeley FrameNet Project. In C. Boitet and P. Whitelock, editors, 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, COLING-ACL '98, August 10-14, 1998, Université de Montréal, Montréal, Quebec, Canada. Proceedings of the Conference., pages 86–90, Montreal, Québec, Canada, 1998. Association for Computational Linguistics. http://aclweb. org/anthology/P/P98/P98-1013.pdf.
- [5] J. Bos. Wide-Coverage Semantic Analysis with Boxer. In J. Bos and R. Delmonte, editors, *Conference on Semantics in Text Processing (STEP)*, pages 277–286, Venice, Italy, 2008. College Publications. DOI:10.3115/1626481.1626503.

- [6] E. Cambria, D. Olsher, and D. Rajagopal. SenticNet 3: A common and common-sense knowledge base for cognitiondriven sentiment analysis. In C. E. Brodley and P. Stone, editors, Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec, Canada., pages 1515–1521, Québec, Canada, 2014. AAAI Press. http://www.aaai.org/ocs/index. php/AAAI/AAAI14/paper/view/8479.
- [7] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. H. Jr., and T. M. Mitchell. Toward an architecture for neverending language learning. In M. Fox and D. Poole, editors, Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010, pages 1306–1313, Georgia, USA, 2010. AAAI Press. http://www.aaai.org/ocs/index. php/AAAI/AAAI10/paper/view/1879.
- [8] L. Champollion. Quantification and negation in event semantics. *The Baltic International Yearbook of Cognition, Logic and Communication*, 6:1–23, October 2011. DOI:10.4148/biyclc.v6i0.1563.
- [9] P. Cimiano and J. Völker. Text2onto a framework for ontology learning and data-driven change discovery. In A. Montoyo, R. Muñoz, and E. Métais, editors, Natural Language Processing and Information Systems, 10th International Conference on Applications of Natural Language to Information Systems, NLDB 2005, Alicante, Spain, June 15-17, 2005, Proceedings, volume 3513 of Lecture Notes in Computer Science, pages 227–238, Alicante, Spain, 2005. Springer. DOI:10.1007/11428817_21.
- [10] M. d'Aquin, E. Motta, M. Sabou, S. Angeletou, L. Gridinoc, V. Lopez, and D. Guidi. Toward a new generation of semantic web applications. *IEEE Intelligent Systems*, 23(3):20–28, 2008. DOI:10.1109/MIS.2008.54.
- [11] D. Das, N. Schneider, D. Chen, and N. A. Smith. Probabilistic frame-semantic parsing. In R. Kaplan, J. Burstein, M. Harper, and G. Penn, editors, *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, June 2-4, 2010, Los Angeles, California, USA*, pages 948–956, Los Angeles, California, 2010. The Association for Computational Linguistics. http://www.aclweb.org/anthology/N10-1138.
- [12] O. Etzioni, M. Banko, and M. J. Cafarella. Machine reading. In Y. Gil and R. J. Mooney, editors, *Proceedings, The Twenty-First National Conference on Artificial Intelligence and the Eighteenth Innovative Applications of Artificial Intelligence Conference, July 16-20, 2006, Boston, Massachusetts, USA, pages 1517–1519, Boston, Massachusetts, 2006. AAAI Press. http://www.aaai.org/Library/ AAAI/2006/aaai06-239.php.*
- [13] O. Etzioni, A. Fader, J. Christensen, S. Soderland, and Mausam. Open information extraction: The second generation. In T. Walsh, editor, *IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011*, pages 3–10, Barcelona, Spain, 2011. IJCAI/AAAI. DOI:10.5591/978-1-57735-516-8/IJCAI11-012.
- [14] P. Ferragina and U. Scaiella. TAGME: on-the-fly annotation of short text fragments (by Wikipedia entities). In J. Huang, N. Koudas, G. J. F. Jones, X. Wu, K. Collins-Thompson, and A. An, editors, *Proceedings of the 19th* ACM Conference on Information and Knowledge Manage-

ment, CIKM 2010, Toronto, Ontario, Canada, October 26-30, 2010, pages 1625–1628, Toronto, Canada, 2010. ACM. DOI:10.1145/1871437.1871689.

- [15] C. J. Fillmore. Frame semantics. In L. S. of Korea, editor, *Linguistics in the Morning Calm*, pages 111–137. Hanshin Publishing Co., 1982. DOI:10.1016/B0-08-044854-2/00424-7.
- [16] A. Gangemi. A comparison of knowledge extraction tools for the Semantic Web. In P. Cimiano, Ó. Corcho, V. Presutti, L. Hollink, and S. Rudolph, editors, *The Semantic Web: Semantics and Big Data, 10th International Conference, ESWC* 2013, Montpellier, France, May 26-30, 2013. Proceedings, volume 7882 of Lecture Notes in Computer Science, pages 351– 366, Montpellier, France, 2013. Springer. DOI:10.1007/978-3-642-38288-8_24.
- [17] A. Gangemi, R. Navigli, and P. Velardi. The OntoWordNet Project: Extension and axiomatization of conceptual relations in WordNet. In R. Meersman, Z. Tari, and D. C. Schmidt, editors, On The Move to Meaningful Internet Systems 2003: CoopIS, DOA, and ODBASE - OTM Confederated International Conferences, CoopIS, DOA, and ODBASE 2003, Catania, Sicily, Italy, November 3-7, 2003, volume 2888 of Lecture Notes in Computer Science, pages 820–838, Catania, Italy, 2003. Springer. DOI:10.1007/978-3-540-39964-3_52.
- [18] A. Gangemi, A. G. Nuzzolese, V. Presutti, F. Draicchio, A. Musetti, and P. Ciancarini. Automatic typing of dbpedia entities. In P. Cudré-Mauroux, J. Heflin, E. Sirin, T. Tudorache, J. Euzenat, M. Hauswirth, J. X. Parreira, J. Hendler, G. Schreiber, A. Bernstein, and E. Blomqvist, editors, *The Semantic Web - ISWC 2012 - 11th International Semantic Web Conference, Boston, MA, USA, November 11-15, 2012, Proceedings, Part I, volume 7649 of Lecture Notes in Computer Science, pages 65–81, Boston, MA, USA, 2012. Springer. DOI:10.1007/978-3-642-35176-1_5.*
- [19] A. Gangemi, A. G. Nuzzolese, V. Presutti, and D. R. Recupero. Adjective semantics in open knowledge extraction. In R. Ferrario and W. Kuhn, editors, *Formal Ontology in Information Systems Proceedings of the 9th International Conference, FOIS 2016, Annecy, France, July 6-9, 2016, volume 283 of Frontiers in Artificial Intelligence and Applications*, pages 167–180, Annecy, France, 2016. IOS Press. DOI:10.3233/978-1-61499-660-6-167.
- [20] A. Gangemi and V. Presutti. Ontology Design Patterns. In S. Staab and R. Studer, editors, *Handbook on Ontologies*, International Handbooks on Information Systems, pages 221–243. Springer, 2nd edition, 2009. DOI:10.1007/978-3-540-92673-3_10.
- [21] A. Gangemi, V. Presutti, and D. Reforgiato Recupero. Framebased detection of opinion holders and topics: A model and a tool. *IEEE Computational Intelligence*, 9(1):20–30, 2014. DOI:10.1109/MCI.2013.2291688.
- [22] N. D. Goodman and A. Stuhlmüller. Knowledge and implicature: Modeling language understanding as social cognition. *Topics in Cognitive Science*, 5(1):173–184, 2013. DOI:10.1111/tops.12007.
- [23] S. Hellmann, J. Lehmann, S. Auer, and M. Brümmer. Integrating NLP using Linked Data. In H. Alani, L. Kagal, A. Fokoue, P. T. Groth, C. Biemann, J. X. Parreira, L. Aroyo, N. F. Noy, C. Welty, and K. Janowicz, editors, *The Semantic Web - ISWC* 2013 - 12th International Semantic Web Conference, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part II, volume 8219 of Lecture Notes in Computer Science, pages 98–

113, Sydney, Australia, 2013. Springer. DOI:10.1007/978-3-642-41338-4_7.

- [24] A. D. Iorio, A. G. Nuzzolese, and S. Peroni. Towards the automatic identification of the nature of citations. In A. G. Castro, C. Lange, P. W. Lord, and R. Stevens, editors, *Proceedings of the 3rd Workshop on Semantic Publishing, Montpellier, France, May 26th, 2013.*, volume 994 of *CEUR Workshop Proceedings*, pages 63–74, Montpellier, France, 2013. CEUR-WS.org. http://ceur-ws.org/ Vol-994/paper-06.pdf.
- [25] H. Kamp. A theory of truth and semantic representation. In J. A. G. Groenendijk, T. M. V. Janssen, and M. B. J. Stokhof, editors, *Formal Methods in the Study of Language, Part I*, pages 277–322. Mathematisch Centrum, 1981.
- [26] A. Khalili, S. Auer, and A. N. Ngomo. conTEXT Lightweight text analytics using Linked Data. In V. Presutti, C. d'Amato, F. Gandon, M. d'Aquin, S. Staab, and A. Tordai, editors, *The Semantic Web: Trends and Challenges - 11th International Conference, ESWC 2014, Anissaras, Crete, Greece, May 25-*29, 2014. Proceedings, volume 8465 of Lecture Notes in Computer Science, pages 628–643, Crete, Greece, 2014. Springer. DOI:10.1007/978-3-319-07443-6_42.
- [27] K. Kipper Schuler. VerbNet: A Broad-Coverage, Comprehensive Verb Lexicon. PhD thesis, University of Pennsylvania, 2006.
- [28] S. Krause, L. Hennig, A. Gabryszak, F. Xu, and H. Uszkoreit. Sar-graphs: A Linked Linguistic Knowledge Resource Connecting Facts with Language. In C. Chiarcos, P. Cimiano, N. I. Vassar, J. P. McCrae, and P. Osenova, editors, *Fourth Workshop on Linked Data in Linguistics: Resources and Applications (LDL)*, pages 30–38. Association for Computational Linguistics, 2015. DOI:10.18653/v1/W15-4204.
- [29] B. Levin. English Verb Classes and Alternations A Preliminary Investigation. University of Chicago Press, Chicago, USA, 1993.
- [30] V. Lopez, A. Nikolov, M. Sabou, V. S. Uren, E. Motta, and M. d'Aquin. Scaling up question-answering to linked data. In P. Cimiano and H. S. Pinto, editors, *Knowledge Engineering and Management by the Masses - 17th International Conference, EKAW 2010, Lisbon, Portugal, October 11-15, 2010. Proceedings*, volume 6317 of *Lecture Notes in Computer Science*, pages 193–210, Lisbon, Portugal, 2010. Springer. DOI:10.1007/978-3-642-16438-5_14.
- [31] T. M. Mitchell, W. W. Cohen, E. R. H. Jr., P. P. Talukdar, J. Betteridge, A. Carlson, B. D. Mishra, M. Gardner, B. Kisiel, J. Krishnamurthy, N. Lao, K. Mazaitis, T. Mohamed, N. Nakashole, E. A. Platanios, A. Ritter, M. Samadi, B. Settles, R. C. Wang, D. T. Wijaya, A. Gupta, X. Chen, A. Saparov, M. Greaves, and J. Welling. Never-ending learning. In B. Bonet and S. Koenig, editors, *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*, pages 2302–2310, Austin, Texas, USA, 2015. AAAI Press. http://www.aaai.org/ocs/ index.php/AAAI/AAAI15/paper/view/10049.
- [32] M. Mongiovì, D. Reforgiato Recupero, A. Gangemi, V. Presutti, A. G. Nuzzolese, and S. Consoli. Semantic reconciliation of knowledge extracted from text through a novel machine reader. In K. Barker and J. M. Gómez-Pérez, editors, *Proceedings of the 8th International Conference on Knowledge Capture, K-CAP 2015, Palisades, NY, USA, October 7-*10, 2015, pages 25:1–25:4, Palisades, NY, USA, 2015. ACM.

DOI:10.1145/2815833.2816945.

- [33] R. Moot and C. Retoré. *The Logic of Categorial Grammars: A deductive account of natural language syntax and semantics*. Springer, Berlin Heidelberg, 2012. DOI:10.1007/978-3-642-31555-8.
- [34] A. Moro, A. Raganato, and R. Navigli. Entity linking meets word sense disambiguation: A unified approach. *Transactions of the Association for Computational Linguistics (TACL)*, 2:231–244, 2014. https://tacl2013.cs.columbia. edu/ojs/index.php/tacl/article/view/291.
- [35] M. Morzycki. *Modification*. Key Topics in Semantics and Pragmatics. Cambridge University Press, 2015.
- [36] R. Navigli and S. P. Ponzetto. BabelNet: The Automatic Construction, Evaluation and Application of a Wide-Coverage Multilingual Semantic Network. *Artificial Intelli*gence, 193:217–250, 2012. DOI:10.1016/j.artint.2012.07.001.
- [37] A. G. Nuzzolese, A. Gangemi, and V. Presutti. Gathering lexical Linked Data and knowledge patterns from FrameNet. In M. A. Musen and Ó. Corcho, editors, *Proceedings of the 6th International Conference on Knowledge Capture (K-CAP 2011), June 26-29, 2011, Banff, Alberta, Canada*, pages 41–48, Banff, AB, Canada, 2011. ACM. DOI:10.1145/1999676.1999685.
- [38] OSGi Service Platform Release 4 Version 4.2, Core Specification, September 2009.
- [39] S. Peroni, A. Gangemi, and F. Vitali. Dealing with markup semantics. In C. Ghidini, A. N. Ngomo, S. N. Lindstaedt, and T. Pellegrini, editors, *Proceedings of the 7th International Conference on Semantic Systems, I-SEMANTICS 2011, Graz, Austria, September 7-9, 2011, ACM International Conference* Proceeding Series, pages 111–118, Graz, Austria, 2011. ACM. DOI:10.1145/2063518.2063533.
- [40] S. Peroni and F. Vitali. Annotations with EARMARK for arbitrary, overlapping and out-of order markup. In U. M. Borghoff and B. Chidlovskii, editors, *Proceedings of the 2009* ACM Symposium on Document Engineering, Munich, Germany, September 16-18, 2009, pages 171–180, Munich, Germany, 2009. ACM. DOI:10.1145/1600193.1600232.
- [41] G. Pirrò and J. Euzenat. A feature and information theoretic framework for semantic similarity and relatedness. In P. F. Patel-Schneider, Y. Pan, P. Hitzler, P. Mika, L. Zhang, J. Z. Pan, I. Horrocks, and B. Glimm, editors, *The Semantic Web - ISWC 2010 - 9th International Semantic Web Conference, ISWC 2010, Shanghai, China, November 7-11, 2010, Revised Selected Papers, Part I*, volume 6496 of *Lecture Notes in Computer Science*, pages 615–630, Shanghai, China, 2010. Springer. DOI:10.1007/978-3-642-17746-0_39.
- [42] V. Presutti, S. Consoli, A. G. Nuzzolese, D. Reforgiato Recupero, A. Gangemi, I. Bannour, and H. Zargayouna. Uncovering the semantics of Wikipedia pagelinks. In K. Janowicz, S. Schlobach, P. Lambrix, and E. Hyvönen, editors, *Knowl*edge Engineering and Knowledge Management - 19th Interna-

tional Conference, EKAW 2014, Linköping, Sweden, November 24-28, 2014. Proceedings, volume 8876 of Lecture Notes in Computer Science, pages 413–428, Linköping, Sweden, 2014. Springer. DOI:10.1007/978-3-319-13704-9_32.

- [43] V. Presutti, F. Draicchio, and A. Gangemi. Knowledge extraction based on discourse representation theory and linguistic frames. In A. ten Teije, J. Völker, S. Handschuh, H. Stuckenschmidt, M. d'Aquin, A. Nikolov, N. Aussenac-Gilles, and N. Hernandez, editors, *Knowledge Engineering and Knowledge Management - 18th International Conference, EKAW* 2012, Galway City, Ireland, October 8-12, 2012. Proceedings, volume 7603 of Lecture Notes in Computer Science, pages 114–129, Galway, Ireland, 2012. Springer. DOI:10.1007/978-3-642-33876-2 12.
- [44] D. Reforgiato Recupero, V. Presutti, S. Consoli, A. Gangemi, and A. Nuzzolese. Sentilo: Frame-based sentiment analysis. *Cognitive Computation*, 7(2):211–225, 2015. DOI:10.1007/s12559-014-9302-z.
- [45] G. Rizzo, R. Troncy, S. Hellmann, and M. Brümmer. NERD meets NIF: lifting NLP extraction results to the linked data cloud. In C. Bizer, T. Heath, T. Berners-Lee, and M. Hausenblas, editors, WWW2012 Workshop on Linked Data on the Web, Lyon, France, 16 April, 2012, volume 937 of CEUR Workshop Proceedings, Lyon, France, 2012. CEUR-WS.org. http://ceur-ws.org/ Vol-937/ldow2012-paper-02.pdf.
- [46] P. Singh. The public acquisition of commonsense kknowledge. In J. Karlgren, B. Gambäck, and P. Kanerva, editors, Acquiring (and Using) Linguistic (and World) Knowledge for Information Access - Papers from the 2002 AAAI Spring Symposium, pages 47–52, Palo Alto, CA, USA, 2002. AAAI Press. http://www.aaai.org/Library/Symposia/ Spring/2002/ss02-09-011.php.
- [47] M. Steedman. *The Syntactic Process*. MIT Press, Cambridge, MA, USA, 2000.
- [48] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: A core of semantic knowledge. In C. L. Williamson, M. E. Zurko, P. F. Patel-Schneider, and P. J. Shenoy, editors, *Proceedings of the 16th International Conference on World Wide Web, WWW 2007, Banff, Alberta, Canada, May 8-12, 2007*, pages 697–706, Banff, AB, Canada, 2007. ACM. DOI:10.1145/1242572.1242667.
- [49] M. van Assem, A. Gangemi, and G. Schreiber. Conversion of WordNet to a standard RDF/OWL representation. In N. Calzolari, K. Choukri, A. Gangemi, B. Maegaard, J. Mariani, J. Odijk, and D. Tapias, editors, *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC)*, pages 237–242, Genova, Italy, 2006. European Language Resources Association (ELRA). http://www. lrec-conf.org/proceedings/lrec2006/.