

N-ary Relation Extraction for Joint T-Box and A-Box Knowledge Base Augmentation

Marco Fossati ^{a,*}, Emilio Dorigatti ^b, and Claudio Giuliano ^c

^a *Data and Knowledge Management Unit, Fondazione Bruno Kessler, via Sommarive 18, 38123 Trento, Italy*

E-mail: fossati@fbk.eu

^b *Department of Computer Science, University of Trento, via Sommarive 9, 38123 Trento, Italy*

E-mail: emilio.dorigatti@unitn.it

^c *Future Media Unit, Fondazione Bruno Kessler, via Sommarive 18, 38123 Trento, Italy*

E-mail: giuliano@fbk.eu

Abstract. The Web has evolved into a huge mine of knowledge carved in different forms, the predominant one still being the free-text document. This motivates the need for Intelligent Web-reading Agents: hypothetically, they would skim through disparate Web sources corpora and generate meaningful structured assertions to fuel Knowledge Bases (KBs). Ultimately, comprehensive KBs, like Wikidata and DBpedia, play a fundamental role to cope with the issue of information overload. On account of such vision, this paper depicts the FACT EXTRACTOR, a complete Natural Language Processing (NLP) pipeline which reads an input textual corpus and produces machine-readable statements. Each statement is supplied with a confidence score and undergoes a disambiguation step via entity linking, thus allowing the assignment of KB-compliant URIs. The system implements four research contributions: it (1) executes N-ary relation extraction by applying the Frame Semantics linguistic theory, as opposed to binary techniques; it (2) jointly populates both the T-Box and the A-Box of the target KB; it (3) relies on a lightweight NLP machinery, namely part-of-speech tagging only; it (4) enables a completely supervised yet reasonably priced machine learning environment through a crowdsourcing strategy. We assess our approach by setting the target KB to DBpedia and by considering a use case of 52,000 Italian Wikipedia soccer player articles. Out of those, we yield a dataset of more than 213,000 triples with a 78.5% F_1 . We corroborate the evaluation via (i) a performance comparison with a baseline system, as well as (ii) an analysis of the T-Box and A-Box augmentation capabilities. The outcomes are incorporated into the Italian DBpedia chapter, can be queried through its SPARQL endpoint, and/or downloaded as standalone data dumps. The codebase is released as free software and is publicly available in the DBpedia Association repository.

Keywords: Information Extraction, Natural Language Processing, Frame Semantics, Crowdsourcing, Machine Learning

1. Introduction

The World Wide Web is nowadays one of the most prominent sources of information and knowledge. Despite the constantly increasing availability of semi-structured or structured data, a major portion of its content is still represented in an unstructured form, namely free text: deciphering its meaning is a complex task for machines and yet relies on subjective human interpretations. Hence, there is an ever growing need for

Intelligent Web-reading Agents, i.e., Artificial Intelligence systems that can read and understand human language in documents across the Web. Ideally, these agents should be robust enough to interchange between heterogeneous sources with agility, while maintaining equivalent reading capabilities. More specifically, given a set of input corpora (where an item corresponds to the textual content of a Web source), they should be able to navigate from corpus to corpus and to extract comparable structured assertions out of each one. Ultimately, the collected data would feed a target *Knowledge Base* (KB), namely a repository that encodes areas of human

* Corresponding author. E-mail: fossati@fbk.eu

intelligence into a richly shaped representation. Typically, KBs are composed of graphs, where real-world and abstract entities are bound together through relationships, and classified according to a formal description of the world, i.e., an ontology.

In this scenario, the encyclopedia Wikipedia contains a huge amount of data, which may represent the best digital approximation of human knowledge. Recent efforts, most notably DBPEDIA [23], FREEBASE [8], YAGO [21], and WIKIDATA [31], attempt to extract semi-structured data from Wikipedia in order to build KBs that are proven useful for a variety of applications, such as question answering, entity summarization and entity linking (EL), just to name a few. The idea has not only attracted a continuously rising commitment of research communities, but has also become a substantial focus of the largest Web companies. As an anecdotal yet remarkable proof, Google acquired Freebase in 2010,¹ embedded it in its KNOWLEDGE GRAPH,² and has lately opted to shut it down to the public.³ Currently, it is foreseen that Freebase data will eventually migrate to Wikidata⁴ via the *primary sources* tool,⁵ which aims at standardizing the flow for data donations.

However, the reliability of a general-purpose KB like Wikidata is an essential requirement to ensure credible (thus high-quality) content: as a support for their trustworthiness, data should be validated against third-party resources. Even though the Wikidata community strongly agrees on the concern,⁶ few efforts have been approached towards this direction. The addition of references to external (i.e., non-Wikimedia), authoritative Web sources can be viewed as a form of validation. Consequently, such real-world setting further consolidates the need for an intelligent agent that harvests structured data from raw text and produces e.g., Wikidata statements with reference URLs. Besides the prospective impact on the KB augmentation and quality, the agent would also dramatically shift the burden of manual data

addition and curation, by pushing the (intended) fully human-driven flow towards an assisted paradigm, where automatic suggestions of pre-packaged statements just require to be approved or rejected. Figure 1 depicts the current state of the primary sources tool interface for Wikidata editors, which is in active development yet illustrates such future technological directions. Our system already takes part in the process, as it feeds the tool back-end.

On the other hand, the DBpedia EXTRACTION FRAMEWORK⁷ is pretty much mature when dealing with Wikipedia semi-structured content like infoboxes, links and categories. Nevertheless, unstructured content (typically text) plays the most crucial role, due to the potential amount of extra knowledge it can deliver: to the best of our understanding, no efforts have been carried out to integrate an unstructured data extractor into the framework. For instance, given the Germany football team article,⁸ we aim at extracting a set of meaningful facts and structure them in machine-readable statements. The sentence In Euro 1992, Germany reached the final, but lost 0–2 to Denmark would produce a list of *triples*, such as:

(Germany, defeat, Denmark)
 (defeat, score, 0–2)
 (defeat, winner, Denmark)
 (defeat, competition, Euro 1992)

To fulfill both Wikidata and DBpedia duties, we aim at investigating in what extent can the *Frame Semantics* theory [16,17] be leveraged to perform Information Extraction over Web documents. The main purpose of Information Extraction is to gather structured data from free text via Natural Language Processing (NLP), while Frame Semantics originates from linguistic research in Artificial Intelligence. A *frame* can be informally defined as an event triggered by some term in a text and embedding a set of participants, or *Frame Elements* (FEs). Hence, the aforementioned sentence would induce the DEFEAT frame (triggered by lost) together with the WINNER, COMPETITION, and SCORE participants. Such theory has led to the creation of FRAMENET [5,6], namely a lexical database with manually annotated examples of frame usage in English. FrameNet currently adheres to a rigorous protocol for data annotation and quality control. The activity

¹<https://googleblog.blogspot.it/2010/07/deeper-understanding-with-metaweb.html>

²https://www.google.com/intl/en_us/insidesearch/features/search/knowledge.html

³<https://plus.google.com/109936836907132434202/posts/bu3z2wVqcQc>

⁴https://www.wikidata.org/wiki/Wikidata:WikiProject_Freebase

⁵https://www.wikidata.org/wiki/Wikidata:Primary_sources_tool

⁶https://www.wikidata.org/wiki/Wikidata:Referencing_improvements_input, <http://blog.wikimedia.de/2015/01/03/scaling-wikidata-success-means-making-the-pie-bigger/>

⁷<https://github.com/dbpedia/extraction-framework>

⁸http://en.wikipedia.org/wiki/Germany_national_football_team



Figure 1. Screenshot of the Wikidata primary sources gadget activated in ROBERTO BAGGIO’s page. Automatic suggestions are displayed with a blue background and can be either approved or rejected by editors

is known to be expensive with respect to time and cost, thus constituting an encumbrance for the extension of the resource [4], both in terms of additional labeled sentences and of languages. To alleviate this, crowdsourcing the annotation task is proven to dramatically reduce the financial and temporal expenses. Consequently, we foresee to exploit the novel annotation approach described in [18], which provides full frame annotation in a single step and in a bottom-up fashion, thus being also more compliant with the definition of frames as per [17].

In this paper, we focus on Wikipedia as the source corpus and on DBpedia as the target KB. We propose to apply NLP techniques to Wikipedia text in order to harvest structured facts that can be used to automatically add novel statements to DBpedia. Our FACT EXTRACTOR is set apart from related state of the art thanks to the combination of the following contributions:

1. **N-ary relation extraction**, as opposed to binary standard approaches, e.g., [15,3,2,30,14,9];
2. **joint T-Box and A-Box population** of the target KB, in contrast to e.g., [12];
3. relatively **cheap NLP machinery**, requiring grammatical analysis (i.e., part-of-speech tagging) only,

with no need for syntactic parsing (e.g., [25]) nor semantic role labeling (e.g., [22]);

4. **low-cost yet supervised machine learning** paradigm, via training set crowdsourcing as in [18], which ensures full supervision without the need for expert annotators.

The remainder of this paper is structured as follows: Section 2 defines the specific problem we aim at tackling and illustrates the proposed solution. We introduce a use case in Section 3, which will drive the implementation of our system. Its high-level architecture is then described in Section 4, and devises the core modules, which we detail in Section 5, 6, 7, 8, and 9. A baseline system is reported in Section 10: this enables the comparative evaluation presented in Section 11, among with an assessment of the T-Box and A-Box enrichment capabilities. In Section 12, we gather a list of research and technical considerations to pave the way for future work. The state of the art is reviewed in Section 13, before our conclusions are drawn in Section 14.

Table 1
Fact extraction examples on the Germany national football team article

Sentence	Extracted statements
The first manager of the Germany national team was Otto Nerz	(Germany, team manager, Otto Nerz)
Germany has won the World Cup four times	(Germany, victory, World Cup), (victory, count, 4)
In the 70s, Germany wore Erima kits	(Germany, wearing, Erima), (wearing, period, 1970)

2. Problem and Solution

The main research challenge is formulated as a KB population problem: specifically, we tackle how to automatically enrich DBpedia resources with novel statements extracted from the text of Wikipedia articles. We conceive the solution as a machine learning task implementing the Frame Semantics linguistic theory [16,17]: we investigate how to recognize meaningful factual parts given a natural language sentence as input. We cast this as a classification activity falling into the supervised learning paradigm. Specifically, we focus on the construction of a new extractor, to be integrated into the current DBpedia infrastructure. Frame Semantics will enable the discovery of relations that hold between entities in raw text. Its implementation takes as input a collection of documents from Wikipedia (i.e., the corpus) and outputs a structured dataset composed of machine-readable statements.

3. Use Case

Soccer is a widely attested domain in Wikipedia: according to DBpedia,⁹ the English Wikipedia counts a total of 223,050 articles describing soccer-related entities, which is a significant portion (around 5%) of the whole chapter. Moreover, infoboxes on those articles are generally very rich (cf. for instance the Germany national football team article). On account of these observations, the soccer domain properly fits the main challenge of this effort. Table 1 displays three examples of candidate statements from the Germany national football team article text, which do not exist in the corresponding DBpedia resource.

⁹As per the 2014 release, based on the English Wikipedia dumps from May 2014.

4. System Description

The implementation workflow is intended as follows, depicted in Figure 2, and applied to the use case in Italian language:

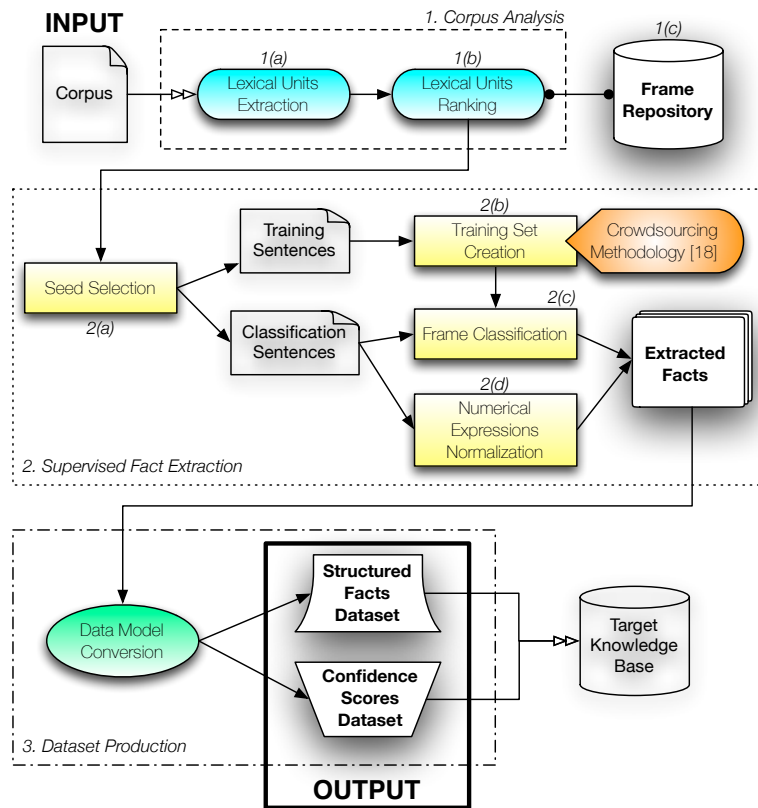
1. Corpus Analysis

- (a) **Lexical Units (LUs) extraction** via text tokenization, lemmatization, and part-of-speech (POS) tagging. LUs serve as the frame triggers;
- (b) **LUs Ranking** through lexicographical and statistical analysis of the input corpus. The selection of top-N meaningful LUs is produced via a combination of term weighting measures (i.e., TF-IDF) and purely statistical ones (i.e., standard deviation);
- (c) each selected LU will trigger one or more frames together with their FEs, depending on the definitions contained in a given **frame repository**. The repository also holds the input labels for two automatic classifiers (the former handling FEs, the latter frames) based on Support Vector Machines (SVM).

2. Supervised Fact Extraction

- (a) **Seed Selection**: two sets of sentences are gathered upon the candidate LUs, one for training examples and the other for the actual classification;
- (b) **Training Set Creation**: construction of a fully annotated training set, leveraging the crowdsourcing methodology proposed in [18];
- (c) **Frame Classification**: massive frame and FEs extraction on the input corpus seed sentences, via the classifiers trained with the result of the previous step.

3. **Dataset Production**: structuring the extraction results to fit the target KB (i.e., DBpedia) **data model** (i.e., RDF). A frame would map to a property, while participants would either map to subjects or to objects, depending on their role.

Figure 2. High level overview of the *Fact Extractor* system

We proceed with a simplification of the original Frame Semantics theory with respect to two aspects: (a) LUs may be evoked by additional POS (e.g., nouns), but we focus on verbs, since we assume that they are more likely to trigger factual information; (b) depending on the frame repository, full lexical coverage may not be guaranteed (i.e., some LUs may not trigger any frames), but we expect that ours will, otherwise LU candidates would not generate any fact.

5. Corpus Analysis

Wikipedia dumps¹⁰ are packaged as XML documents and contain text formatted according to the Mediawiki markup syntax,¹¹ with templates to be transcluded.¹² Hence, a pre-processing step is required to obtain a raw

text representation of the dump. To achieve this, we leverage the WIKIEXTRACTOR,¹³ a third-party tool that retains the text and expands templates of a Wikipedia XML dump, while discarding other data such as tables, references, images, etc. We note that the tool is not completely robust with respect to templates expansion. Such drawback is expected for two reasons: first, new templates are constantly defined, thus requiring regular maintenance of the tool; second, Wikipedia editors do not always comply to the specifications of the templates they include. Therefore, we could not obtain a fully cleaned Wikipedia plain text corpus, and noticed gaps in its content, probably due to template expansion failures. Nevertheless, we argue that the loss of information is not significant and can be neglected despite the recall cost. From the entire Italian Wikipedia corpus, we slice the use case subset by querying the ITALIAN DBPEDIA CHAPTER¹⁴ for the Wikipedia article IDs of relevant entities.

¹⁰<http://dumps.wikimedia.org/>

¹¹<https://www.mediawiki.org/wiki/Help:Formatting>

¹²<https://www.mediawiki.org/wiki/Help:Templates>

¹³<https://github.com/attardi/wikiextractor>

¹⁴<http://it.dbpedia.org/>

5.1. Lexical Units Extraction

Given the use case corpus, we first extract the complete set of verbs through a standard NLP pipeline: tokenization, lemmatization and POS tagging. POS information is required to identify verbs, while lemmas are needed to build the ranking. TREETAGGER¹⁵ is exploited to fulfill these tasks. Although our input has a relatively low dimension (i.e., 7.25 million tokens circa), we observe that the tool is not able to handle it as a whole, since it crashes with a segmentation fault even on a powerful machine (i.e., 24 cores CPU at 2.53 GHz, 64 GB RAM). Consequently, we had to run it over each document, thus impacting on the processing time. However, we believe that further investigation will lead to the optimization of such issue.

5.2. Lexical Units Ranking

The unordered set of extracted verbs is the subject of a further analysis, which aims at discovering the most representative verbs with respect to the corpus. Two measures are combined to generate a score for each verb lemma, thus enabling the creation of a rank. We first compute the term frequency–inverse document frequency (TF-IDF) of each verb lexicalization (i.e., the occurring tokens) over each document in the corpus: this weighting measure is intended to capture the *lexical* relevance of a given verb, namely how important it is with respect to other terms in the whole corpus. Then, we determine the standard deviation value out of the TF-IDF scores list: this *statistical* measure is meant to catch heterogeneously distributed verbs, in the sense that the higher the standard deviation is, the more variably the verb is used, thus helping to understand its overall usage signal over the corpus. Ultimately, we produce the final score and assign it to a verb lemma by averaging all its lexicalizations scores. The top-N lemmas serve as candidate LUs, each evoking one or more frames according to the definitions of a given frame repository.

6. Use Case Frame Repository

Among the top 50 LUs that emerged from the corpus analysis phase, we manually selected a subset of 5 items to facilitate the full implementation of our pipeline.

Once the approach has been tested and evaluated, it can scale up to the whole ranking (cf. Section 12 for more observations). The selected LUs comply to two criteria: first, they are picked from both the best and the worst ranked ones, with the purpose of assessing the validity of the corpus analysis as a whole; second, they fit the use case domain, instead of being generic. Consequently, we proceed with the following LUs: *esordire* (to start out), *giocare* (to play), *perdere* (to lose), *rimanere* (to stay, remain), and *vincere* (to win).

The next step consists of finding a language resource (i.e., frame repository) to suitably represent the use case domain. Given a resource, we first need to define a relevant subset, then verify that both its frame and FEs definitions are a relevant fit. After an investigation of FrameNet and KICKTIONARY [29], we notice that:

- to the best of our knowledge, no suitable domain-specific Italian FrameNet or Kicktionary are publicly available, in the sense that neither LU sets nor annotated sentences for the Italian language match our purposes;
- FrameNet is too coarse-grained to encode our domain knowledge. For instance, the FINISH_COMPETITION frame may seem a relevant candidate at a first glimpse, but does not make the distinction between a victory and a defeat (as it can be triggered by both *to win* and *to lose* LUs), thus rather fitting as a superframe (but no subframes exist);
- Kicktionary is too specific, since it is built to model the speech transcriptions of football matches. While it indeed contains some in-scope frames such as VICTORY (evoked by *to win*), most LUs are linked to frames that are not likely to appear in our input corpus, e.g., *to play* with PASS (occurring in sentences like *Ronaldinho played the ball in for Deco*).

Therefore, we adopted a custom frame repository, maximizing the reuse of the available ones as much as possible, thus serving as a hybrid between FrameNet and Kicktionary. Moreover, we tried to provide a challenging model for the classification task, prioritizing FEs overlap among frames and LU ambiguity (i.e., focusing on very fine-grained semantics with subtle sense differences). We believe this does not only apply to machines, but also to humans: we can view it as a stress test both for the machine learning and the crowdsourcing parts. A total of 6 frames and 15 FEs are modeled with Italian labels as follows:

¹⁵<http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

- ATTIVITÀ (activity), FEs AGENTE (agent), COMPETIZIONE (competition), DURATA (duration), LUOGO (place), SQUADRA (team), TEMPO (time). Evoked by *esordire* (to start out), *giocare* (to play), *rimanere* (to stay, remain), as in Roberto Baggio played with Juventus in Serie A between 1990 and 1995. Frame label translated from FrameNet ACTIVITY, FEs from a subset of FrameNet ACTIVITY;
- PARTITA (match), FEs SQUADRA_1 (team 1), SQUADRA_2 (team 2), COMPETIZIONE, LUOGO, TEMPO, PUNTEGGIO (score), CLASSIFICA (ranking). Evoked by *giocare*, as in Juventus played Milan at the UEFA cup final (2-0). Frame label translated from Kicktionary MATCH, FEs from a subset of FrameNet COMPETITION, LU shared by both;
- SCONFITTA (defeat), FEs PERDENTE, VINCITORE, COMPETIZIONE, LUOGO, TEMPO, PUNTEGGIO, CLASSIFICA. Evoked by *perdere* (to lose), as in Milan lost 0-2 against Juventus at the UEFA cup final. Frame label translated from Kicktionary DEFEAT, FEs from a subset of FrameNet BEAT_OPPONENT, LU from Kicktionary;
- STATO (status), FEs ENTITÀ (entity), STATO (status), DURATA, LUOGO, SQUADRA, TEMPO. Evoked by *rimanere*, as in Roberto Baggio remained faithful to Juventus until 1995. Custom frame and FEs derived from corpus evidence, to augment the *rimanere* LU ambiguity;
- TROFEO (trophy), FEs AGENTE, COMPETIZIONE, SQUADRA, PREMIO (prize), LUOGO, TEMPO, PUNTEGGIO, CLASSIFICA. Evoked by *vincere* (to win), as in Roberto Baggio won a UEFA cup with Juventus in 1992. Custom frame label, FEs from a subset of FrameNet WIN_PRIZE, LU from FrameNet;
- VITTORIA (victory), FEs VINCITORE, PERDENTE, COMPETIZIONE, LUOGO, TEMPO, PUNTEGGIO, CLASSIFICA. Evoked by *vincere*, as in Juventus won 2-0 against Milan at the UEFA cup final. Frame label translated from Kicktionary VICTORY, FEs from a subset of FrameNet BEAT_OPPONENT, LU from Kicktionary.

7. Supervised Fact Extraction

The first step involves the creation of the training set: we leverage the crowdsourcing platform CROWD-

FLOWER¹⁶ and the method described in [18], which requires users to detect the *core* FEs: these are the fundamental items to distinguish between frames, as opposed to *extra* ones, thus allowing to automatically induce the correct frame. The training set has a double outcome, as it will feed two classifiers: one will identify FEs, and the other is responsible for frames.

Both frame and FEs recognition are cast to a multi-class classification task: while the former can be related to text categorization, the latter should answer questions such as *can this entity be this FE?* or *is this entity this FE in this context?*. Such activity boils down to semantic role labeling (cf. [24] for an introduction), and usually requires a more fine-grained text analysis. Previous work in the area exploits deeper NLP layers, such as syntactic parsing (e.g., [25]). We alleviate this through EL techniques, which perform word sense disambiguation by linking relevant parts of a source sentence to URIs of a target KB. We leverage THE WIKI MACHINE¹⁷ [19], a state-of-the-art [26] approach conceived for connecting text to Wikipedia URLs, thus inherently entailing DBpedia URIs. EL results are part of the FE classifier feature set. We claim that EL enables the automatic addition of features based on existing entity attributes within the target KB (notably, the class of an entity, which represents its semantic type).

Given as input an unknown sentence, the full frame classification workflow involves the following tasks: tokenization, POS tagging, EL, FE classification, and frame classification.

7.1. Seed Selection

The seed selection procedure allows to harvest meaningful sentences from the input corpus, and to feed the classifier. Therefore, its outcome is two-fold: to build a representative training set and to extract relevant sentences for classification. We experimented multiple strategies as follows. They all share the same base constraint, i.e., each seed must contain a LU lexicalization.

- *Baseline*: the seed must not exceed a given length in words;
- *Sentence splitter*: the seed forms a complete sentence extracted with a sentence splitter. This strategy requires training data for the splitter;
- *Syntactic*: the seed must match a pattern expressed via a context-free chunker grammar. This strat-

¹⁶<http://www.crowdfLOWER.com/>

¹⁷<http://thewikimachine.fbK.eu/>

egy requires a POS tagger and engineering effort for defining the grammar (e.g., a noun phrase, followed by a verb phrase, followed by a noun phrase);

- *Lexical*: the seed must match a pattern based on lexicalizations of candidate entities. This strategy requires querying a KB for instances of relevant classes (e.g., soccer-related ones as per the use case).

First, we note that all the strategies but the baseline necessitate a significant cost overhead in terms of language resources availability and engineering. Furthermore, given the soccer use case input corpus of 52,000 articles circa, those strategies dramatically reduce the number of seeds, while the baseline performed an extraction with a .95 article/seed ratio (despite some noise). Hence, we decided to leverage the baseline for the sake of simplicity and for the compliance to our contribution claims. We still foresee further investigation of the other strategies for scaling besides the use case.

7.2. Training Set Creation

We apply the one-step, bottom-up approach described in [18] to let the crowd perform a full frame annotation over a set of training sentences. The set is randomly sampled from the input corpus and contains 3,055 items. The outcome is the same amount of frame examples and 55,385 FE examples. The task is sent to the CrowdFlower platform.

7.2.1. Task Anatomy

We ask the crowd to (a) read the given sentence, (b) focus on the “topic” (i.e., the frame) written above it, and (c) assign the correct “label” (i.e., the FE) to each “word” (i.e., unigram) or “group of words” (i.e., n-grams) from the multiple choices provided below each n-gram. Figure 3 displays the front-end interface of a sample sentence.

During the preparation phase of the task input data, the main challenge is to automatically provide the crowd with relevant candidate FE text chunks, while minimizing the production of noisy ones. To tackle this, we experimented with the following chunking strategies:

- third-party full-stack NLP pipeline, namely TEXTPRO [27] for Italian, by extracting nominal chunks with the CHUNKPRO module;¹⁸

Attività

Dal gennaio 2010 gioca con il Legnano in Lega Pro Seconda Divisione.

<p>gennaio</p> <p><input type="radio"/> Nessuno</p> <p><input type="radio"/> Agente</p> <p><input type="radio"/> Competizione</p> <p><input type="radio"/> Luogo</p> <p><input type="radio"/> Squadra</p>	<p>Legnano</p> <p><input type="radio"/> Agente</p> <p><input type="radio"/> Competizione</p> <p><input type="radio"/> Luogo</p> <p><input type="radio"/> Nessuno</p> <p><input type="radio"/> Squadra</p>	<p>Lega Pro Seconda Divisione</p> <p><input type="radio"/> Agente</p> <p><input type="radio"/> Luogo</p> <p><input type="radio"/> Competizione</p> <p><input type="radio"/> Nessuno</p> <p><input type="radio"/> Squadra</p>
--	--	---

Figure 3. Worker interface example

- custom noun phrase chunker via a context-free grammar;
- EL surface forms;

We surprisingly observed that the full-stack pipeline outputs a significant amount of noisy chunks, besides being the slowest strategy. On the other hand, the custom chunker was the fastest one, but still too noisy to be crowdsourced. EL resulted in the best trade-off, and we adopted it for the final task. Obviously, noise cannot be automatically eliminated: we cast such validation to the crowd by allowing the None answer along with the candidate FE labels.

The task parameters are as follows:

- we set 3 judgments per sentence to enable the computation of an agreement based on majority vote;
- the pay sums to 5 \$ cents per page, where one page contains 5 sentences;
- we limit the task to Italian native speakers only by targeting the Italian country and setting the required language skills to Italian;
- the minimum worker accuracy is set to 70% in quiz mode (i.e., the warm-up phase where workers are only shown gold units¹⁹ and are recruited according to their accuracy) and relaxed to 65% in work mode (i.e., the actual annotation phase) to avoid extra cost in terms of time and expenses to collect judgments;
- on account of a personal calibration, the minimum time per page threshold is set to 30 seconds, which allows to automatically discard a contributor when triggered;
- we set the maximum number of judgments per contributor to 280, in order to prevent each contributor from answering more than once on a given

¹⁸<http://textpro.fbk.eu/>

¹⁹cf. [18] for details

Table 2
Training set crowdsourcing task outcomes

Sentences	3,111
Gold units	56
Trusted judgments	9,198
Untrusted judgments	972
Total cost	152.46 \$

sentence, while avoiding to remove proficient contributors from the task.

The outcomes are resumed in Table 2.

Finally, the crowdsourced annotation results are processed and translated into a suitable format to serve as input training data for the classifier.

7.3. Frame Classification: Features

We train our classifiers with the following linguistic features, in the form of bag-of-features vectors:

1. *both classifiers*: for each input word token, both the token itself (bag of terms) and the lemma (bag of lemmas);
2. *FE classifier*: contextual sliding window of width 5 (i.e., 5-gram, for each token, consider the 2 previous and the 2 following ones);
3. *frame classifier*: we implement our bottom-up frame annotation approach, thus including the set of FE labels (bag of roles) to help this classifier induce the frame;
4. *gazetteer*: defined as a map of key-value pairs, where each key is a feature and its value is a list of n-grams, we automatically build a wide-coverage gazetteer with relevant DBpedia ontology (DBPO) classes as keys (e.g., SoccerClub) and instances as values (e.g., Juventus), by way of a query to the target KB.

8. Numerical Expressions Normalization

During our pilot crowdsourcing annotation experiments, we noticed a low agreement on numerical FEs. Moreover, asking the crowd to label such frequently occurring FEs would represent a considerable overhead, resulting in a higher temporal cost (i.e., more annotations per sentence) and lower overall annotation accuracy. Hence, we opted for the implementation of a rule-based system to detect and normalize numerical expressions. The normalization process takes as input

a numerical expression such as a date, a duration, or a score, and outputs a transformation into a standard format suitable for later inclusion into the target KB.

The task is not formulated as a classification one, but we argue it is relevant for the completeness of the extracted facts: rather, it is carried out via matching and transformation rule pairs. Given for instance the input expression *tra il 1920 e il 1925* (between 1920 and 1925), our normalizer first matches it through a regular expression rule, then applies a transformation rule complying to the XML Schema Datatypes²⁰ (typically dates and times) standard, and finally produces the following output:²¹

```
duration: "P5Y"^^xsd:duration
start: "1920"^^xsd:gYear
end: "1925"^^xsd:gYear
```

All rule pairs are defined with the programming language-agnostic YAML²² syntax. The pair for the above example is as follows:

```
Regular Expression:
tra il (?P<y1>\ d{{2,4}}) e il (?P<y2>\ d{{2,4}})

Transformation:
{
  'duration':
  ``P{}Y"^^<{}>' .format (
  int (match.group('y2')) - int (match.group('y1')),
  schema['duration']
  ),
  'start':
  ``{}"^^<{}>' .format (
  abs_year (match.group('y1')), schema['year']
  ),
  'end':
  ``{}"^^<{}>' .format (
  abs_year (match.group('y2')), schema['year']
  )
}
```

9. Dataset Production

The integration of the extraction results into DBpedia requires their conversion to a suitable data model, i.e., RDF. Frames intrinsically bear N-ary relations through FEs, while RDF naturally represents binary relations. Hence, we need a method to express FEs relations in RDF, which we call *reification*. This can be achieved in multiple ways:

²⁰<http://www.w3.org/TR/xmlschema-2/>

²¹We use the *xsd* prefix as a short form for the full URI <http://www.w3.org/2001/XMLSchema#>

²²<http://www.yaml.org/spec/1.2/spec.html>

- standard reification;²³
- N-ary relations,²⁴ and similarly [13];
- named graphs.²⁵

A recent overview [20] highlighted that all the mentioned strategies are similar with respect to query performance. Given as input n frames and m FEs, we argue that:

- standard reification is too verbose, since it would require $3(n + m)$ triples;
- applying Pattern 1 of the aforementioned W3C Working Group note to N-ary relations would allow us to build $n + m$ triples;
- named graphs can be used to encode provenance or context metadata, e.g., the article URI from where a fact was extracted. In our case however, the fourth element of the quad would be the frame (which represents the context), thus boiling down to minting $n + m$ quads instead of triples;

We opted for the less verbose strategy, namely N-ary relations. Given the running example sentence *In Euro 1992, Germany reached the final, but lost 0–2 to Denmark*, classified as a DEFEAT frame and embedding the FEs WINNER, LOSER, COMPETITION, SCORE, we generate RDF as per the following Turtle serialization:

```
:Germany :defeat :Defeat_01 .
:Defeat_01
  :winner :Denmark ;
  :loser :Germany ;
  :competition :Euro_1992 ;
  :score "0-2" .
```

We add an extra instance type triple to assign an ontology class to the reified frame, as well as a provenance triple to indicate the original sentence:

```
:Defeat_01
  a :Defeat ;
  :extractedFrom "In Euro 1992,
    Germany reached the final,
    but lost 0-2 to Denmark"@it .
```

In this way, the generated statements amount to $n + m + 2$.

It is not trivial to decide on the subject of the main frame statement, since not all frames are meant to have

exactly one core FE that would serve as a plausible logical subject candidate: most have many, e.g., FINISH_COMPETITION has COMPETITION, COMPETITOR and OPPONENT as core FEs in FrameNet. Therefore, we tackle this as per the following assumption: given the encyclopedic nature of our input corpus, both the logical and the topical subjects correspond in each document. Hence, each candidate sentence inherits the document subject. We acknowledge that such assumption strongly depends on the corpus: it applies to entity-centric documents, but will not perform well for general-purpose ones such as news articles. However, we believe it is still a valid in-scope solution fitting our scenario.

9.1. Confidence Scores

Besides the fact datasets, we also keep track of confidence scores and generate additional datasets accordingly. Therefore, it is possible to filter facts that are not considered as confident by setting a suitable threshold. When processing a sentence, our pipeline outputs two different scores for each FE, stemming from the entity linker and the supervised classifier. We merge both signals by calculating the F-score between them, as if they were representing precision and recall, in a fashion similar to the standard classification metrics. The final score can be then produced via an aggregation of the single FE scores in multiple ways, namely: (a) arithmetic mean; (b) weighted mean based on core FEs (i.e., they have a higher weight than extra ones); (c) harmonic mean, weighted on core FEs as well.

10. Baseline Classifier

To enable a performance evaluation comparison with the supervised method, we developed a rule-based algorithm that handles the full frame and FEs annotation. The main intuition is to map FEs defined in the frame repository to ontology classes of the target KB: such mapping serves as a set of rule pairs ($FE, class$), e.g., (WINNER, SoccerClub). In the FrameNet terminology, this is homologous to the assignment of *semantic types* to FEs: for instance, in the ACTIVITY frame, the AGENT is typed with the generic class Sentient. The idea would allow the implementation of the bottom-up one-step annotation flow described in [18]: to achieve this, we run EL over the input sentences and check whether the attached ontology class metadata appear in the frame repository, thus fulfilling the FE classification task.

²³<http://www.w3.org/TR/2004/REC-rdf-primer-20040210/#reification>

²⁴<http://www.w3.org/TR/swbp-n-aryRelations/>

²⁵<http://www.w3.org/TR/rdf11-concepts/>

Algorithm 1 Rule-based baseline classifier**Input:** S ; F ; L **Output:** C

```

1:  $C \leftarrow \emptyset$ 
2: for all  $s \in S$  do
3:    $E \leftarrow \text{entityLinking}(s)$ 
4:    $T \leftarrow \text{tokenize}(s)$ 
5:   for all  $t \in T$  do
6:     if  $t \in L$  then #Check whether a sentence
       token matches a LU token
7:       for all  $f \in F$  do
8:          $\text{core} \leftarrow \text{false}$ 
9:          $O \leftarrow \text{getLinkedEntityClasses}(E)$ 
10:        for all  $o \in O$  do
11:           $fe \leftarrow \text{lookup}(f)$  #Get the FE that
            maps to the current linked entity class
12:           $\text{core} \leftarrow \text{checkIsCore}(fe)$ 
13:        end for
14:        if  $\text{core}$  then #Relaxed classification
15:           $c \leftarrow [s, f, fe]$ 
16:           $C \leftarrow C \cup \{c\}$ 
17:        else
18:          continue #Skip to the next frame
19:        end if
20:      end for
21:    end if
22:  end for
23: end for
24: return  $C$ 

```

Besides that, we exploit the notion of core FEs: this would cater for the frame disambiguation part. Since a frame may contain at least one core FE, we proceed with a *relaxed* assignment, namely we set the frame if a given input sentence contains at least one entity whose ontology class maps to a core FE of that frame. The implementation workflow is illustrated in Algorithm 1: it takes as input the set S of sentences, the frame repository F embedding frame and FEs labels, core/non-core annotations and rule pairs, and the set L of trigger LU tokens.

It is expected that the relaxed assignment strategy will not handle the overlap of FEs across competing frames that are evoked by a single LU. Therefore, if at least one core FE is detected in multiple frames, the baseline makes a random assignment for the frame. Furthermore, the method is not able to perform FE classification in case different FEs share the ontology class (e.g., both WINNER and LOSER map to SoccerClub): we opt for a FE random guess as well.

11. Evaluation

We assess our main research contributions through the analysis of the following aspects:

- Classification performance;
- T-Box property coverage extension;
- A-Box statements addition.

11.1. Classification Performance

Table 3 describes the overall performance of the baseline and the supervised system over a gold standard dataset. We randomly sampled 500 sentences containing at least one occurrence of our use case LU set from the input corpus. We first outsourced the annotation to the crowd as per the training set construction and the results were further manually validated twice by the authors. Measures are computed as follows: (1) a true positive is triggered if the predicted label is correct and the predicted text chunk at least partially matches the expected one; (2) chunks that should not be labeled are marked with a “O” and not explicitly counted as true or false positives.

11.1.1. Supervised Classification Performance Breakdown

Figure 4 and Figure 5 respectively display the FE and frame classification confusion matrices. First, since the “O” markers were discarded from the evaluation, all the respective performance measures amount to 0.

Concerning FEs, we observe that COMPETIZIONE is frequently mistaken for PREMIO and ENTITÀ, while rarely for TEMPO and DURATA, or just missed. On the other hand, TEMPO is mistaken for COMPETIZIONE: our hypothesis is that competition mentions, such as World Cup 2014, are disambiguated as a whole entity by the linker, since a specific target Wikipedia article exists. However, it overlaps with a temporal expression, thus confusing the classifier. AGENTE is often mistaken for ENTITÀ, due to their equivalent semantic type, which is always a person. Finally, we notice that

Table 3

Classification performance standard evaluation over a gold standard of 500 random sentences from the Italian Wikipedia corpus

Approach	Task	Precision	Recall	F1
Baseline	Frames	.73	.61	.66
Baseline	FES	.66	.64	.65
Supervised	Frames	.80	.78	.79
Supervised	FES	.82	.75	.78

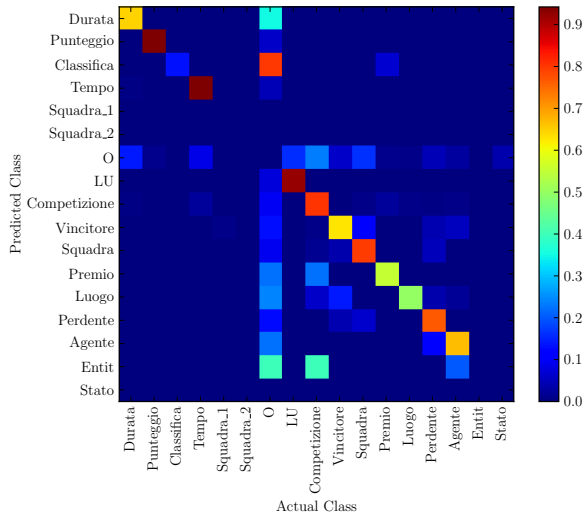


Figure 4. Supervised FE classification confusion matrix

some FEs (i.e., SQUADRA_1, SQUADRA_2 and CLASSIFICA) did not appear in the training set, but did in the gold standard, and viceversa (i.e., ENTITÀ).

With respect to frames, we note that ATTIVITÀ is often mistaken for STATO or not classified at all: in fact, the difference between these two frames is quite subtle with respect to their sense. The former is more generic and could also be labeled as CAREER: if we viewed it in a frame hierarchy, it would serve as a superframe of the latter. The latter instead encodes the development modality of a soccer player’s career, e.g., when he remains unbound from some team due to contracting issues. Hence, we may conclude that distinguishing between these frames is a challenge even for

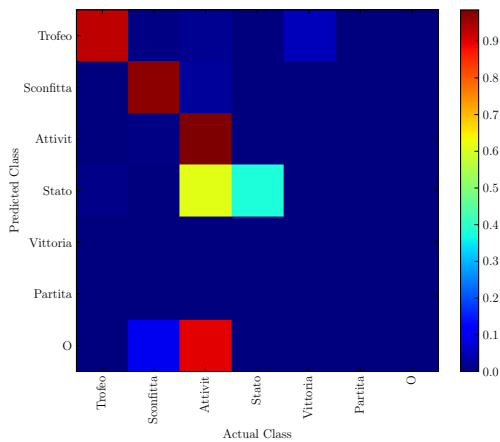


Figure 5. Supervised frame classification confusion matrix

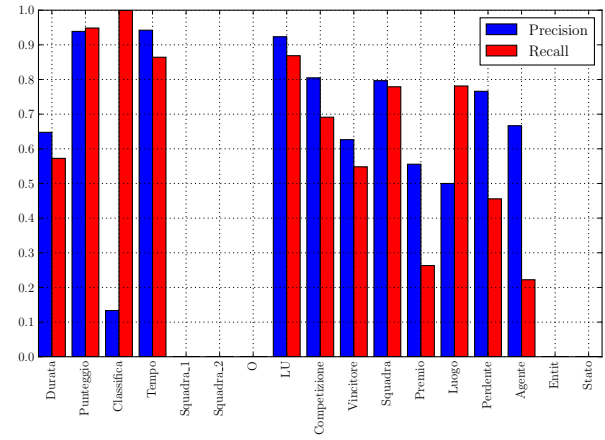


Figure 6. Supervised FE classification precision and recall breakdown

humans. Furthermore, frames with no FEs are classified as “O”, thus considered wrong despite the correct prediction. VITTORIA is almost never mistaken for TROFEO: this is positively surprising, since the FE COMPETIZIONE (frame VITTORIA) is often mistaken for PREMIO (frame TROFEO), but those FEs do not seem to affect the frame classification. Again, such FE distinction must take into account a delicate sense nuance, which is hard for humans as well. Due to an error in the training set crowdsourcing step, we lack of VITTORIA and PARTITA samples.

Figure 6 and Figure 7 respectively plot the FE and frame classification performance, broken down to each label.

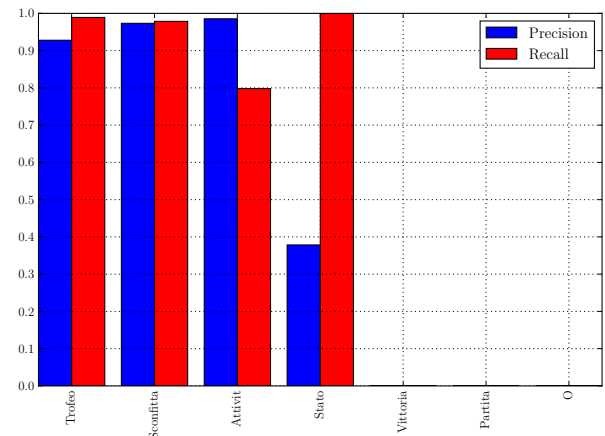


Figure 7. Supervised frame classification precision and recall breakdown

Table 4

Lexicographical analysis of the Italian Wikipedia soccer player subcorpus

Stems (frequency %)	Candidate frames (FrameNet)
gioc (47), partit (39), campionat (34), stagion (36), presen (30), disput (20), serie (14), nazional (13), titolar (13), competizion (5), scend (5), torne (5)	COMPETITION
pass (24), trasfer (19), prest (15), contratt (11)	ACTIVITY_START, EMPLOYMENT_START
termin (12), contratt, ced (10), lasc (6), vend (2)	ACTIVITY_FINISH, EMPLOYMENT_END
gioc, disput (20), scend	FINISH_GAME
campionat, stagion, serie, nazional, competizion, torne	FINISH_COMPETITION
vins/vinc (18), pers/perd (11), sconfi (8)	BEAT_OPONENT, FINISH_GAME
vins/vinc, conquis (8), otten (7), raggiun (6), aggiud (2)	WIN_PRIZE, PERSONAL_SUCCESS

11.2. T-Box Enrichment

One of our main objectives is to extend the target KB ontology with new properties on existing classes. We focus on the use case and argue that our approach will have a significant impact if we manage to identify non-existing properties. This would serve as a proof of concept which can ideally scale up to all kinds of input. In order to assess such potential impact in discovering new relations, we need to address the following question: *which extractable relations are not already mapped in DBPO or do not even exist in the raw infobox properties datasets?*. Table 4 illustrates an empirical lexicographical study gathered from the Italian Wikipedia soccer player subcorpus (circa 52,000 articles). It contains absolute occurrence frequencies of word stems (in descending order) that are likely to trigger domain-relevant frames, thus providing a rough overview of the extraction potential.

The corpus analysis phase (cf. Section 5) yielded a ranking of LUs evoking the frames ACTIVITY, DEFEAT, MATCH, TROPHY, STATUS, and VICTORY: these frames would serve as ontology property candidates, together with their embedded FEs. DBPO already has most of the classes that are needed to represent the main entities involved in the use case: SoccerPlayer, SoccerClub, SoccerManager, SoccerLeague, SoccerTournament, SoccerClubSeason, SoccerLeagueSeason, although some of them lack an exhaustive description (cf. SoccerClubSeason²⁶ and SoccerLeagueSeason).²⁷

²⁶<http://mappings.dbpedia.org/server/ontology/classes/SoccerClubSeason>

²⁷<http://mappings.dbpedia.org/server/ontology/classes/SoccerLeagueSeason>

For each of the aforementioned classes, we computed the amount and frequency of ontology and raw infobox properties via a query, with results in ascending order of frequency: Figure 8 illustrates their distribution. The horizontal axis stands for the normalized (log scale) frequency, encoding the current usage of properties in the target KB; the vertical axis represents the ratio (which we call coverage) between the position of the property in the ordered result set of the query and the total amount of distinct properties (i.e., the size of the result set). Properties with a null frequency are ignored.

First, we observe a lack of ontology property usage in 4 out of 7 classes, probably due to missing mappings between Wikipedia template attributes and DBPO. On the other hand, the ontology properties have a more homogenous distribution compared to the raw ones: this serves as an expected proof of concept, since the main purpose of DBPO and the ontology mappings is to merge heterogenous and multilingual Wikipedia template attributes into a unique representation. In av-

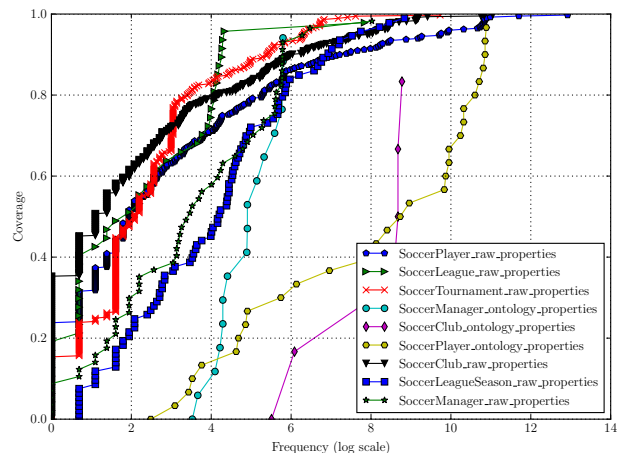


Figure 8. Italian DBpedia soccer property statistics

Table 5

Relative A-Box population gain compared to pre-existing T-Box property assertions in the Italian DBpedia chapter

Property	Dataset	Assertions (#)	Gain (%)
careerStation	DBpedia	2,073	N.A.
	Baseline all	20,430	89.8
	Supervised all	26,316	92.12
award	DBpedia	7,755	N.A.
	Baseline all	4,953	-56.57
	Supervised all	10,433	25.66
playerStatus	DBpedia	0	N.A.
	Baseline all	0	0
	Supervised all	26	100

erage, most raw properties are concentrated below coverage and frequency threshold values of 0.8 and 4 respectively: this means that roughly 80% of them have a significantly low usage, further highlighted by the log scale. While ontology properties are better distributed, most still do not reach a high coverage/frequency trade-off, except for SoccerPlayer, which benefits from both rich data (cf. Section 3) and mappings.²⁸

On the light of the two analyses discussed above, it is clear that our approach would result in a larger variety and finer granularity of facts than those encoded into Wikipedia infoboxes and DBPO classes. Moreover, we believe the lack of dependence on infoboxes would enable more flexibility for future generalization to sources beyond Wikipedia.

Subsequent to the use case implementation, we manually identified the following mappings from frames and FEs to DBPO properties:

- Frames: (ACTIVITY, careerStation), (AWARD, award), (STATUS, playerStatus);
- FEs: (TEAM, team), (SCORE, score), (DURATION, [duration, startYear, endYear]).

In conclusion, we claim that 3 out of 6 frames and 12 out of 15 FEs represent novel T-Box properties.

11.3. A-Box Population

Our methodology enables a joint T-Box and A-Box augmentation: while frames and FEs serve as T-Box properties, the extracted facts feed the A-Box part. Out of 49,063 input sentences, we generated a total of

213,479 and 216,451 triples (i.e., with a 4.35 and 4.41 ratio per sentence) from the supervised and the baseline classifiers respectively. 52% and 55% circa are considered *confident*, namely facts with confidence scores (cf. Section 9.1) above the dataset average threshold.

To assess the domain coverage gain, we can exploit two signals: (a) the amount of produced novel data with respect to pre-existing T-Box properties and (b) the overlap with already extracted assertions, regardless of their origin (i.e., whether they stem from the raw infobox or the ontology-based extractors). Given the same Italian Wikipedia dump input dating 21 January 2015, we ran both the baseline and the supervised fact extraction, as well as the DBpedia extraction framework to produce an Italian DBpedia chapter release, thus enabling the coverage comparison.

Table 5 describes the analysis of signal (a) over the 3 frames that are mapped to DBPO properties. For each property and dataset, we computed the amount of available assertions and reported the gain relative to the fact extraction datasets. Although we considered the whole Italian DBpedia KB in these calculations, we observe that it has a generally low coverage with respect to the analysed properties, probably due to missing ontology mappings. For instance, the amount of assertions is always zero if we analyse the use case subset only, as no specific relevant mappings (e.g., *Carriera_sportivo*²⁹ to *careerStation*) currently exist. We view this as a major achievement, since our automatic approach also serves as a substitute for the manual mapping procedure.

Table 6

Overlap with pre-existing assertions in the Italian DBpedia chapter and relative gain in A-Box population

	Overlap (#)	Gain (%)
Baseline all	3,341	98.2
Supervised all	4,546	97.4
Baseline confident	2,387	97.6
Supervised confident	2,841	96.8

Table 6 shows the results for signal (b). To obtain them, we proceeded as follows.

1. slice the use case DBpedia subset;
2. gather the subject-object patterns from all datasets. Properties are not included, as they are not comparable;

²⁸http://mappings.dbpedia.org/index.php/Mapping_it:Sportivo

²⁹https://it.wikipedia.org/wiki/Template:Carriera_sportivo

3. compute the patterns overlap between DBpedia and each of the fact extraction datasets (including the confident subsets);
4. compute the gain in terms of novel assertions relative to the fact extraction datasets.

The A-Box enrichment is clearly visible from the results, given the low overlap and high gain in all approaches, despite the rather large size of the DBpedia use case subset, namely 6, 167, 678 assertions.

12. Observations

First, we acknowledge that the use case frame repository is not exhaustive: LUs may have a higher ambiguity. For instance, *giocare* (to play) may trigger an additional frame depending on the context (as in the sentence to play as a defender); *esordire* (to start out) may also trigger the frame *PARTITA* (match). Second, if a sentence is not in the gold standard, the supervised classifier should discard it (abstention). Third, the baseline approach may contain rules that are more harmful than beneficial, depending on the target KB reliability: for instance, the *SportsEvent* *DBPO* class leads to wrongly typed instances, due to the misuse of the template by Wikipedia editors. Finally, both the input corpus and the target KB originate from a relatively small Wikipedia chapter (i.e., Italian, with 1.23 million articles) if compared to the largest one (i.e., English, with almost 5 million articles). Therefore, we recognize that the T-Box and A-Box evaluation results may be proportionally different if obtained with English data.

12.1. Scaling Up

Our approach has been tested on the Italian language, a specific domain, and with a small frame repository. Hence, we may consider the use case implementation as a monolingual closed-domain information extraction system. We outline below the points that need to be addressed for scaling up to multilingual open information extraction:

1. *Language*: training data availability for POS tagging and lemmatization. The LUs automatically extracted through the corpus analysis phase should be projected to a suitable frame repository;
2. *Domain*:
 - Baseline: mapping between FEs and target KB ontology classes;
 - Supervised:

- * financial resources for the crowdsourced training set construction, in average 4.79 \$ cents per annotated sentence;
- * adapt the query to generate the gazetteer.

12.2. Technical Future Work

We report below a list of technical improvements left for planned implementation:

- LUs are handled as unigrams, but n-grams should be considered too;
- tagging n-grams with ontology classes retrieved at the EL step may be an impactful additional feature;
- the gazetteer is currently being matched at the token level, but it may be more useful if run over the whole input (sentence);
- in order to reduce the noise in the training set, we foresee to leverage a sentence splitter and extract 1-sentence examples only;
- further evaluation experiments will also count EL surface forms instead of links;
- the inclusion of the frame confidence would further refine the final confidence score.

13. Related Work

We locate our effort at the intersection of the following research areas:

- Information Extraction;
- KB Construction;
- Open Information Semantification.

13.1. Information Extraction

Although the borders are blurred, nowadays we can distinguish two principal fields in Information Extraction, namely Relation Extraction (RE) and Open Information Extraction (OIE). While both aim at structuring information in the form of relations between items, their difference relies in the relations set size, either fixed or potentially infinite. It is commonly argued that the main OIE drawback is the generation of noisy data [11,32], while RE is usually more accurate, but requires expensive supervision in terms of language resources [2,30,32].

13.1.1. Relation Extraction

RE traditionally takes as input a finite set R of relations and a document d , and induces assertions in the form $rel(subj, obj)$, where rel represent binary relations between a subject entity $subj$ and an object entity obj mentioned in d . Hence, it may be viewed as a closed-domain paradigm. Recent efforts [3,2,30] have focused on alleviating the cost of full supervision via distant supervision. Distant supervision leverages available KBs to automatically annotate training data in the input documents. This is in contrast to our work, since we aim at enriching the target KB with external data, rather than using it as a source. Furthermore, our relatively cheap crowdsourcing technique serves as a substitute to distant supervision, while ensuring full supervision. Other approaches such as [7,33] instead leverage text that is not covered by the target KB, like we do.

13.1.2. Open Information Extraction

OIE is defined as a function $f(d)$ over a document d , yielding a set of triples (np_1, rel, np_2) , where nps are noun phrases and rel is a relation between them. Known complete systems include OLLIE [25], REVERB [14], and NELL [9]. Recently, it has been discussed that cross-utterance processing can improve the performance through logical entailments [1]. This paradigm is called “open” since it is not constrained by any schemata, but rather attempts to learn them from unstructured data. In addition, it takes as input heterogeneous sources of information, typically from the Web.

In general, most efforts have focused on English, due to the high availability of language resources. Approaches such as [15] explore multilingual directions, by leveraging English as a source and applying statistical machine translation (SMT) for scaling up to target languages. Although the authors claim that their approach do not directly depend on language resources, we argue that SMT still heavily relies on them. Furthermore, all the above efforts concentrate on binary relations, while we generate n-ary ones: under this perspective, EXEMPLAR [10] is a rule-based system which is closely related to ours.

13.2. Knowledge Base Construction

DBPEDIA [23], FREEBASE [8] and YAGO [21] represent the most mature approaches for automatically building KBs from Wikipedia. Despite its crowd-sourced nature (i.e., fully manual), WIKIDATA [31] benefits from a rapidly growing community of active users,

who have developed several robots for automatic imports of Wikipedia and third-party data. The KNOWLEDGE VAULT [11] is an example of KB construction combining Web-scale textual corpora, as well as additional semi-structured Web data such as HTML tables. Although our system may potentially create a KB from scratch from an input corpus, we prefer to improve the quality of existing resources and integrate into them, rather than developing a standalone one.

13.3. Open Information Semantification

OIE output can indeed be considered structured data compared to free text, but it still lacks of a disambiguation facility: extracted facts generally do not employ unique identifiers (i.e., URIs), thus suffering from intrinsic natural language polysemy (e.g., Jaguar may correspond to the animal or a known car brand). To tackle the issue, [12] propose a framework that clusters OIE facts and maps them to elements of a target KB. Similarly to us, they leverage EL techniques for disambiguation and choose DBpedia as the target KB. Nevertheless, the authors focus on A-Box population, while we also cater for the T-Box part. Moreover, OIE systems are used as a black boxes, in contrast to our full implementation of the extraction pipeline. Finally, relations are still binary, instead of our n-ary ones. Taking as input Wikipedia articles, LEGALO [28] exploits page links manually inserted by editors and attempts to induce the relations between them via NLP. Again, the extracted relations are binary and are not mapped to a target KB for enrichment purposes.

14. Conclusion

In a Web where the profusion of unstructured data limits its automatic interpretation, the necessity of *Intelligent Web-reading Agents* turns more and more evident. These agents should preferably be conceived to browse an extensive and variegated amount of Web sources corpora, harvest structured assertions out of them, and finally cater for target KBs, which can attenuate the problem of information overload. As a support to such vision, we have outlined two real-world scenarios involving general-purpose KBs:

- (a) WIKIDATA would benefit from a system that reads reliable third-party resources, extracts statements complying to the KB data model, and leverages them to validate existing data with reference URLs,

or to recommend new items for inclusion. This would both improve the overall data quality and, most importantly, underpin the costly manual data insertion and curation flow;

- (b) DBPEDIA would naturally evolve towards the extraction of unstructured Wikipedia content. Since Wikidata is designed to be the hub for serving structured data across Wikimedia projects, it will let DBpedia focus on content besides infoboxes, categories and links.

In this article, we presented a system that puts into practice our fourfold research contribution: first, we perform (1) *N-ary relation extraction* thanks to the implementation of Frame Semantics, in contrast to traditional binary approaches; second, we (2) *jointly enrich both the T-Box and the A-Box* parts of our target KB, through the discovery of candidate relations and the extraction of facts respectively. We achieve this with a (3) *shallow layer of NLP* technology only, namely grammatical analysis, instead of more sophisticated ones, such as syntactic parsing. Finally, we ensure a (4) *fully supervised* learning paradigm via an affordable *crowdsourcing* methodology. Our work concurrently bears the advantages and leaves out the weaknesses of RE and OIE: although we assess it in a closed-domain fashion via a use case (Section 3), the corpus analysis module (Section 5) allows to discover an exhaustive set of relations in an open-domain way. In addition, we overcome the supervision cost bottleneck through crowdsourcing. Therefore, we believe our approach can represent a trade-off between open-domain high noise and closed-domain high cost.

The FACT EXTRACTOR is a full-fledged Information Extraction NLP pipeline that analyses a natural language textual corpus and generates structured machine-readable assertions. Such assertions are disambiguated by linking text fragments to entity URIs of the target KB, namely DBpedia, and are assigned a confidence score. For instance, given the sentence Buffon plays for Serie A club Juventus since 2001, our system produces the following dataset:

```
@prefix dbpedia: <http://it.dbpedia.org/resource/> .
@prefix dbpo: <http://dbpedia.org/ontology/> .
@prefix fact: <http://fact.extraction.org/> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .

dbpedia:Gianluigi_Buffon
  dbpo:careerStation dbpedia:CareerStation_01 .

dbpedia:CareerStation_01
  dbpo:team dbpedia:Juventus_Football_Club ;
  fact:competition dbpedia:Serie_A ;
  dbpo:startYear "2001"^^xsd:gYear ;
```

```
fact:confidence "0.906549"^^xsd:float .
```

We estimate the validity of our approach by means of a use case in a specific domain and language, i.e., soccer and Italian. Out of roughly 52,000 Italian Wikipedia articles describing soccer players, we output more than 213,000 triples with an average 78.5% F_1 . Since our focus is the improvement of existing resources rather than the development of a standalone one, we integrated these results into the ITALIAN DBPEDIA CHAPTER³⁰ and made them accessible through its SPARQL endpoint. Moreover, the codebase is publicly available as part of the DBPEDIA ASSOCIATION repository.³¹

We have started to expand our approach under the Wikidata umbrella, where we feed the *primary sources* tool. The community is currently concerned by the trustworthiness of Wikidata assertions: in order to authenticate them, they should be validated against references to external Web sources. Under this perspective, the FACT EXTRACTOR can serve as a reference suggestion mechanism for statement validation. To achieve this, we have successfully managed to switch the input corpus from Wikipedia to third-party corpora and translated our output to fit the Wikidata data model. The soccer use case has already been partially implemented: we have ran the baseline classifier and generated a small demonstrative dataset, named FBK-STREPHIT-SOCCER, which has been uploaded to the primary sources tool back-end. We invite the reader to play with it, by following the instructions in the project page.³²

For future work, we foresee to scale up the implementation towards multilingual open information extraction, thus paving the way to (a) its full deployment into the DBpedia Extraction Framework, and to (b) a thorough referencing system for Wikidata.

Acknowledgments. The FACT EXTRACTOR has been developed within the DBPEDIA ASSOCIATION and was partially funded by GOOGLE under the SUMMER OF CODE 2015 program.

³⁰<http://it.dbpedia.org/2015/09/meno-chiacchiere-piu-fatti-una-marea-di-nuovi-dati-estratti-dal-testo-di-wikipedia/?lang=en>

³¹<https://github.com/dbpedia/fact-extractor>

³²https://www.wikidata.org/wiki/Wikidata:Primary_sources_tool#How_to_use

References

- [1] Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D. Manning. Leveraging linguistic structure for open domain information extraction. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers*, pages 344–354. The Association for Computer Linguistics, 2015.
- [2] Gabor Angeli, Julie Tibshirani, Jean Wu, and Christopher D. Manning. Combining distant and partial supervision for relation extraction. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1556–1567. ACL, 2014.
- [3] Isabelle Augenstein, Diana Maynard, and Fabio Ciravegna. Relation extraction from the web using distant supervision. In Krzysztof Janowicz, Stefan Schlobach, Patrick Lambrix, and Eero Hyvönen, editors, *Knowledge Engineering and Knowledge Management - 19th International Conference, EKAW 2014, Linköping, Sweden, November 24-28, 2014. Proceedings*, volume 8876 of *Lecture Notes in Computer Science*, pages 26–41. Springer, 2014.
- [4] Collin F. Baker. Framenet, current collaborations and future goals. *Language Resources and Evaluation*, 46(2):269–286, 2012.
- [5] Collin F. Baker. Framenet: A knowledge base for natural language processing. *ACL 2014*, 1929:1–5, 2014.
- [6] Collin F. Baker, Charles J. Fillmore, and John B. Lowe. The berkeley framenet project. In Christian Boitet and Pete White-lock, 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, Québec, Canada. Proceedings of the Conference.*, pages 86–90. Morgan Kaufmann Publishers / ACL, 1998.
- [7] Jonathan Berant and Percy Liang. Semantic parsing via paraphrasing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 1: Long Papers*, pages 1415–1425. The Association for Computer Linguistics, 2014.
- [8] Kurt D. Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In Jason Tsong-Li Wang, editor, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2008, Vancouver, BC, Canada, June 10-12, 2008*, pages 1247–1250. ACM, 2008.
- [9] Andrew Carlson, Justin Betteridge, Bryan Kisiel, Burr Settles, Estevam R. Hruschka Jr., and Tom M. Mitchell. Toward an architecture for never-ending language learning. In Maria Fox and David Poole, editors, *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010*. AAAI Press, 2010.
- [10] Filipe de Sá Mesquita, Jordan Schmidek, and Denilson Barbosa. Effectiveness and efficiency of open relation extraction. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 447–457. ACL, 2013.
- [11] Xin Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmman, Shaohua Sun, and Wei Zhang. Knowledge vault: a web-scale approach to probabilistic knowledge fusion. In Sofus A. Macskassy, Claudia Perlich, Jure Leskovec, Wei Wang, and Rayid Ghani, editors, *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014*, pages 601–610. ACM, 2014.
- [12] Arnab Dutta, Christian Meilicke, and Heiner Stuckenschmidt. Enriching structured knowledge with open information. In *Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015*, pages 267–277, 2015.
- [13] Fredo Erxleben, Michael Günther, Markus Kröttsch, Julian Mendez, and Denny Vrandečić. Introducing wikidata to the linked data web. In *The Semantic Web - ISWC 2014 - 13th International Semantic Web Conference, Riva del Garda, Italy, October 19-23, 2014. Proceedings, Part I*, pages 50–65, 2014.
- [14] Anthony Fader, Stephen Soderland, and Oren Etzioni. Identifying relations for open information extraction. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, EMNLP 2011, 27-31 July 2011, John McIntyre Conference Centre, Edinburgh, UK, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1535–1545. ACL, 2011.
- [15] Manaal Faruqi and Shankar Kumar. Multilingual open relation extraction using cross-lingual projection. In *NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015*, pages 1351–1356, 2015.
- [16] Charles Fillmore. Frame semantics. *Linguistics in the morning calm*, pages 111–137, 1982.
- [17] Charles J. Fillmore. Frame Semantics and the nature of language. In *Annals of the New York Academy of Sciences: Conference on the Origin and Development of Language*, pages 20–32. Blackwell Publishing, 1976.
- [18] Marco Fossati, Claudio Giuliano, and Sara Tonelli. Outsourcing framenet to the crowd. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, ACL 2013, 4-9 August 2013, Sofia, Bulgaria, Volume 2: Short Papers*, pages 742–747. The Association for Computer Linguistics, 2013.
- [19] Claudio Giuliano, Alfio Massimiliano Gliozzo, and Carlo Strapparava. Kernel methods for minimally supervised WSD. *Computational Linguistics*, 35(4):513–528, 2009.
- [20] Daniel Hernández, Aidan Hogan, and Markus Kröttsch. Reifying RDF: what works well with wikidata? In *Proceedings of the 11th International Workshop on Scalable Semantic Web Knowledge Base Systems co-located with 14th International Semantic Web Conference (ISWC 2015), Bethlehem, PA, USA, October 11, 2015.*, pages 32–47, 2015.
- [21] Johannes Hoffart, Fabian M. Suchanek, Klaus Berberich, and Gerhard Weikum. YAGO2: A spatially and temporally enhanced knowledge base from wikipedia. *Artif. Intell.*, 194:28–61, 2013.
- [22] Richard Johansson and Pierre Nugues. Dependency-based semantic role labeling of propbank. In *2008 Conference on Empirical Methods in Natural Language Processing, EMNLP 2008, Proceedings of the Conference, 25-27 October 2008, Honolulu,*

- Hawaii, USA, A meeting of SIGDAT, a Special Interest Group of the ACL, pages 69–78. ACL, 2008.
- [23] Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N. Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleef, Sören Auer, and Christian Bizer. DBpedia - a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web Journal*, 6(2):167–195, 2015.
- [24] Lluís Màrquez, Xavier Carreras, Kenneth C. Litkowski, and Suzanne Stevenson. Semantic role labeling: An introduction to the special issue. *Computational Linguistics*, 34(2):145–159, 2008.
- [25] Mausam, Michael Schmitz, Stephen Soderland, Robert Bart, and Oren Etzioni. Open language learning for information extraction. In Jun’ichi Tsujii, James Henderson, and Marius Pasca, editors, *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, EMNLP-CoNLL 2012, July 12-14, 2012, Jeju Island, Korea*, pages 523–534. ACL, 2012.
- [26] Pablo N. Mendes, Max Jakob, Andrés García-Silva, and Christian Bizer. Dbpedia spotlight: shedding light on the web of documents. In Chiara Ghidini, Axel-Cyrille Ngonga Ngomo, Stefanie N. Lindstaedt, and Tassilo Pellegrini, editors, *Proceedings the 7th International Conference on Semantic Systems, ISEMANTICS 2011, Graz, Austria, September 7-9, 2011*, ACM International Conference Proceeding Series, pages 1–8. ACM, 2011.
- [27] Emanuele Pianta, Christian Girardi, and Roberto Zanolì. The textpro tool suite. In *Proceedings of the International Conference on Language Resources and Evaluation, LREC 2008, 26 May - 1 June 2008, Marrakech, Morocco*, 2008.
- [28] Valentina Presutti, Sergio Consoli, Andrea Giovanni Nuzzolese, Diego Reforgiato Recupero, Aldo Gangemi, Ines Bannour, and Haïfa Zargayouna. Uncovering the semantics of wikipedia pagelinks. In Krzysztof Janowicz, Stefan Schlobach, Patrick Lambrix, and Eero Hyvönen, editors, *Knowledge Engineering and Knowledge Management - 19th International Conference, EKAW 2014, Linköping, Sweden, November 24-28, 2014. Proceedings*, volume 8876 of *Lecture Notes in Computer Science*, pages 413–428. Springer, 2014.
- [29] Thomas Schmidt. The kicktionary revisited. In Angelika Storrer, Alexander Geyken, Alexander Siebert, and Kay-Michael Würzner, editors, *Text Resources and Lexical Knowledge. Selected Papers from the 9th Conference on Natural Language Processing, KONVENS 2008, Berlin, Germany*, pages 239–251. Mouton de Gruyter, 2008.
- [30] Mihai Surdeanu, Julie Tibshirani, Ramesh Nallapati, and Christopher D. Manning. Multi-instance multi-label learning for relation extraction. In Jun’ichi Tsujii, James Henderson, and Marius Pasca, editors, *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, EMNLP-CoNLL 2012, July 12-14, 2012, Jeju Island, Korea*, pages 455–465. ACL, 2012.
- [31] Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85, 2014.
- [32] Fei Wu and Daniel S. Weld. Open information extraction using wikipedia. In Jan Hajic, Sandra Carberry, and Stephen Clark, editors, *ACL 2010, Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, July 11-16, 2010, Uppsala, Sweden*, pages 118–127. The Association for Computer Linguistics, 2010.
- [33] Xuchen Yao and Benjamin Van Durme. Information extraction over structured data: Question answering with freebase. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014, June 22-27, 2014, Baltimore, MD, USA, Volume 1: Long Papers*, pages 956–966. The Association for Computer Linguistics, 2014.