

## **Tuurbine: A Generic CBR Engine over RDFS**

Emmanuelle Gaillard, Laura Infante-Blanco, Jean Lieber, Emmanuel Nauer

► **To cite this version:**

Emmanuelle Gaillard, Laura Infante-Blanco, Jean Lieber, Emmanuel Nauer. Tuurbine: A Generic CBR Engine over RDFS. Case-Based Reasoning Research and Development, Sep 2014, Cork, Ireland. 8765, pp.140 - 154, 2014, <<http://link.springer.com/chapter/10.1007>

**HAL Id: hal-01082372**

**<https://hal.inria.fr/hal-01082372>**

Submitted on 13 Nov 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Tuurbine: A Generic CBR Engine Over RDFS\*

Emmanuelle Gaillard, Laura Infante-Blanco,  
Jean Lieber, and Emmanuel Nauer

Université de Lorraine, LORIA — 54506 Vandœuvre-lès-Nancy, France  
CNRS — 54506 Vandœuvre-lès-Nancy, France  
Inria — 54602 Villers-lès-Nancy, France  
`emmanuelle.gaillard@loria.fr`, `laura.infanteblanco@loria.fr`,  
`jean.lieber@loria.fr`, `emmanuel.nauer@loria.fr`

**Abstract.** This paper presents TUURBINE, a case-based reasoning (CBR) system for the Semantic Web. TUURBINE is built as a generic CBR system able to reason on knowledge stored in RDF format; it uses Semantic Web technologies like RDF/RDFS, RDF stores, SPARQL, and optionally Semantic Wikis. TUURBINE implements a generic case-based inference mechanism in which adaptation consists in retrieving similar cases and in replacing some features of these cases in order to obtain one or more solutions for a given query. The search for similar cases is based on a generalization/specialization method performed by means of generalization costs and adaptation rules. The whole knowledge (cases, domain knowledge, costs, adaptation rules) is stored in an RDF store.

**Keywords:** generic CBR engine, Semantic Web, RDFS, triple store

## 1 Introduction and motivations

This paper presents TUURBINE (<http://tuurbine.loria.fr/>), a new generic case-based reasoning system (CBR), the reasoning procedure of which is based on a domain ontology. TUURBINE was created on the basis of the experience acquired during five years with the TAAABLE system [1].

This research work is motivated by the will to develop a generic CBR system. TUURBINE implements a generic case-based inference. Searching for similar cases is based on a generalization/specialization process performed by means of generalization costs and adaptation rules. The generic adaptation approach is inspired from the TAAABLE's adaptation.

The second major motivation is the development of a system able to exploit the huge and growing amount of knowledge available on the Web, especially the knowledge coming from The Linked Data or contained in Semantic Wikis.

---

\* The development of TUURBINE was supported by an Inria ADT funding from October 2011 to October 2013. The authors would like to thank the reviewers who have helped improving the quality of this paper: it was not possible to take into account all their remarks in the paper, but these remaining remarks point out interesting issues that the authors plan to address as future work.

This is why TUUURBINE is built according to Semantic Web standards (RDF, SPARQL, RDFS), to facilitate the interoperability with Semantic Web knowledge. With TUUURBINE, the knowledge (cases, domain knowledge, costs, adaptation rules) is encoded in a triple store, which may be installed in the same machine as the reasoner or in a remote one, and could be additionally be interfaced with a Semantic Wiki, in order to benefit from collaborative web edition and knowledge management involved in the reasoning process [2].

Storing the knowledge in a triple store provides an important advantage: the knowledge is not managed by the CBR system anymore, but by an external tool that is efficient and based on standards. A clear separation is made between the reasoning inference engine sub-system and the knowledge handling and storage sub-system of the CBR scheme: the knowledge used by the CBR system is not loaded in cache once and for all, but on demand, only the necessary knowledge is retrieved to solve a new problem, using SPARQL language. This approach also ensures the exploitation of up-to-date knowledge to best solve a problem in environments where the knowledge evolves continuously. The choice of such an architecture is the result of the experience acquired with TAAABLE. Indeed, in the first version of TAAABLE, the knowledge was completely loaded in memory from several files encoded in various languages [3]. This caused the problem of imperistence of updated knowledge cache, which appears frequently when tuning the system. The second version of TAAABLE improved the knowledge evolution using a Semantic Wiki, called WIKITAAABLE [2] enabling collaborative edition of the knowledge, but still the TAAABLE CBR system required, at that time, loading a new dump of knowledge in memory, after every single knowledge modification, in order to use up-to-date knowledge. The last version of TAAABLE, which uses now the TUUURBINE CBR engine, benefits from a complete synchronization with the knowledge base thanks to a dynamic access to the triple store.

The paper is organized as follows. Section 2 describes briefly the Semantic Web technologies used in TUUURBINE. Section 3 describes the generic CBR approach of TUUURBINE and details the knowledge the approach is based on and the knowledge representation choices that were made. Section 4 details the architecture of the system. Section 5 presents two use-cases. Section 6 discusses this work and relates it to other works. Section 7 concludes the paper.

## 2 Semantic Web technologies used in Tuurbine

**RDF, RDFS and the Triple Stores.** RDF<sup>1</sup>, SPARQL<sup>2</sup> and RDFS<sup>3</sup> are three standards recommendations of the W3C (World Wide Web Consortium) for the Semantic Web.

RDF (Resource Description Framework) is a format to encode resources over the Web, existing in various syntaxes. A *resource* is any kind of entity that has been reified (i.e., associated with an identifier). A *literal* is a constant datatype

<sup>1</sup> <http://www.w3.org/RDF>

<sup>2</sup> <http://www.w3.org/2009/sparql>

<sup>3</sup> <http://www.w3.org/TR/rdf-schema>

(e.g., an integer, a float, a string, etc.). A *property* relates a resource to another resource or a literal. An RDF base is a set of *triples*, a triple being an expression of the form  $\langle s \ p \ o \rangle$ .  $s$ ,  $p$  and  $o$  are respectively called the *subject* (a resource), the *predicate* (a property) and the *object* (a resource or a literal) of a triple. For example, if `romeo` and `juliet` are two resources to be understood as the main characters of [4], if `loves` and `ofTheFamily` are the properties meaning “loves” and “is a member of the family”, and if `age` relates a resource to a literal of integer datatype indicating the age of the resource, then the RDF base  $\mathcal{B} = \{\langle \text{romeo ofTheFamily montague} \rangle, \langle \text{romeo loves juliet} \rangle, \langle \text{juliet age 13} \rangle\}$  means that Romeo, a member of the family Montague, loves Juliet who is 13.

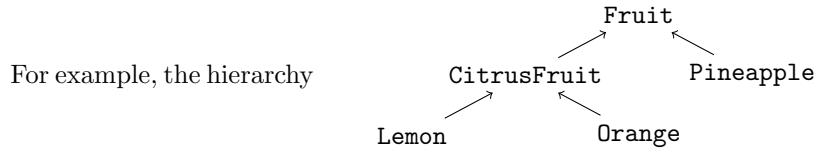
SPARQL (recursive acronym for SPARQL Protocol And RDF Query Language) is at the same time a protocol and a query language for RDF. It can be likened to SQL: SPARQL is to RDF as SQL is to relational databases. For instance, the following SPARQL query to the base  $\mathcal{B}$  above returns the resources  $x$  such that  $\langle x \ \text{loves} \ \text{juliet} \rangle \in \mathcal{B}$  and  $\langle x \ \text{ofTheFamily} \ \text{montague} \rangle \in \mathcal{B}$ :

```
SELECT ?x
WHERE { ?x loves juliet .           // “.” stands for “and”
       ?x ofTheFamily montague }
```

RDFS (RDF Schema) can be considered as a knowledge representation formalism, the syntax of which is RDF (a formula of RDFS is an RDF triple) and the semantics of which is associated with a set of resources having a predefined semantics, called the RDFS *vocabulary*. In this paper, only the properties `rdf:type` and `rdfs:subClassOf` of the RDFS vocabulary are considered and they are abbreviated by `type` and `subc`. `type` can be understood as “is an element of” ( $\in$ ) and `subc` as “is a subset of” ( $\subseteq$ ). The following inference rules are used for RDFS entailment reduced to the vocabulary  $\{\text{type}, \text{subc}\}$ :

$$\frac{\langle a \ \text{type} \ C \rangle \quad \langle C \ \text{subc} \ D \rangle}{\langle a \ \text{type} \ D \rangle} \qquad \frac{\langle C \ \text{subc} \ D \rangle \quad \langle D \ \text{subc} \ E \rangle}{\langle C \ \text{subc} \ E \rangle}$$

`subc` enables to define hierarchies of classes based on the “is more specific than” relation. Moreover, if  $C$  denotes an RDFS class, then  $\langle C \ \text{subc} \ C \rangle$  is an RDFS axiom. If  $\mathcal{B}$  is an RDFS base and  $\tau$  is a triple,  $\mathcal{B} \vdash \tau$  means that  $\tau$  is a logical consequence of  $\mathcal{B}$ . Thus,  $\mathcal{B} \not\vdash \tau$  means that  $\tau$  cannot be entailed from  $\mathcal{B}$ .



can be represented by the following RDFS base:

$$\mathcal{H} = \{\langle \text{CitrusFruit subc Fruit} \rangle, \langle \text{Pineapple subc Fruit} \rangle, \langle \text{Lemon subc CitrusFruit} \rangle, \langle \text{Orange subc CitrusFruit} \rangle\}$$

which entails, e.g., the triple  $\langle \text{Orange subc Fruit} \rangle$ , i.e. every orange is a fruit.

Let  $SQ$  be a SPARQL query and  $\mathcal{B}$  be a RDFS base. Then  $\text{Result}_{\vdash}(SQ, \mathcal{B})$  denotes (in this paper) the result of the SPARQL query on  $\mathcal{B}$  using the entailment. For example,

$$\text{Result}_{\vdash}\left(\boxed{\begin{array}{l} \text{SELECT } ?x \\ \text{WHERE } \{ ?x \text{ subc CitrusFruit} \} \end{array}}, \mathcal{H}\right) = \left\{ \begin{array}{l} \text{Lemon, Orange,} \\ \text{CitrusFruit} \end{array} \right\}.$$

A *triple store* (also called *RDF store*) engine is a database management system for RDF. Using an RDF store enables in particular to avoid loading the whole RDF knowledge base in cache. Some triple store engines support RDFS entailment. However, for performance reasons, in the current version of TUURBINE, the triple store chosen by default is 4Store<sup>4</sup> which does not draw entailments, these being implemented by TUURBINE.

**Semantic Wikis.** A *semantic wiki* is a wiki the contents of which are not only documents and links between documents (as in classical Wikis) but also machine-processible data. SMW is a Semantic Wiki engine extending MediaWiki.<sup>5</sup> The semantic data in a semantic wiki using SMW are managed via a triple store. From a knowledge engineering viewpoint, a semantic wiki engine can be considered as a cooperative knowledge management tool: the knowledge can be edited thanks to this tool, associated with (non formalized) pieces of knowledge in plain text. To exploit the data of a Semantic Wiki, one can create a dump to an RDF file or query directly the triple store, using a SPARQL endpoint.

### 3 Tuurbine reasoning principles

This section presents the principles upon which the TUURBINE engine has been implemented. First, the running example is introduced. This example is then developed in the subsequent sections: representation of pieces of knowledge (cases, domain knowledge, similarity and adaptation knowledge), the representation of TUURBINE queries, and the retrieval and adaptation procedures.

#### 3.1 Introduction of the running example

Let us consider an application of CBR where a case is a recipe (and a case base represents a recipe book). Such application has been developed for the Computer Cooking Contest at ICCBR for the past years.

In this application, let us consider the following query:

$$Q = \boxed{\text{a cocktail recipe with mint, gin, orange juice but no wine.}} \quad (1)$$

If at least one recipe matches exactly  $Q$ , the application returns it. Otherwise, a recipe matching approximately  $Q$  is searched (retrieval step) and then is modified in order to answer  $Q$  (adaptation step).

<sup>4</sup> <http://4store.org>

<sup>5</sup> <http://semantic-mediawiki.org/>

Let us assume that no recipe matching exactly  $Q$  can be found but that the following source case similar to  $Q$  is retrieved:

Source =	<b>Recipe</b> “Mexican cocktail in my way”
	<b>Dish type:</b> cocktail
	<b>Ingredients:</b> 40 cl tequila, 1 l guava juice, 1 l pineapple juice, 30 cl apple juice, 1 pkt vanilla sugar, 3 mint leaves
	<b>Preparation:</b> Combine guava juice and pineapple. Add tequila and apple juice. [...]

Source matches approximately  $Q$  since:

- There is an exact match on the dish type (cocktail), the ingredient mint and the absence of the ingredient wine.
- The ingredients **Tequila** (Source) and **Gin** ( $Q$ ) are subclasses of **Liquor**, and the ingredients **GuavaJuice**, **PineappleJuice**, **AppleJuice** (Source) and **OrangeJuice** ( $Q$ ) are subclasses of **FruitJuice**.

Finally, adaptation modifies Source so that the modified recipe matches exactly  $Q$ . The adaptation is usually based on the approximate matching between Source and  $Q$ . In this example, it consists in applying the following modification:

Replace tequila with gin and guava juice, pineapple juice, or apple juice with orange juice.	(2)
--	-----

Furthermore, there could be available rules which can be applied to adapt the case. In the example, let us consider the following adaptation rule:

$AR_1 =$ In the context of a cocktail dish without anise, guava juice and vanilla sugar can be substituted with orange juice and sugar cane syrup.
--

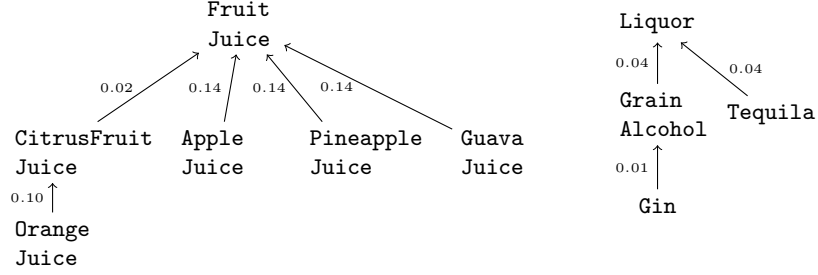
This adaptation rule enables the retrieval of another (hopefully, better) adaptation of Source, given by the following modification:

Replace tequila with gin, and guava juice and vanilla sugar with orange juice and sugar cane syrup.	(3)
---	-----

### 3.2 Knowledge containers of a Tuurbine application

A CBR knowledge base can be split into four knowledge containers [5]. This knowledge is represented by an RDFS base  $KB$  and managed by a triple store.

**Representation of cases.** In the running example, the reasoning takes into account only the dish type and the ingredients of the recipe (neither its title nor its preparation). Moreover, the ingredient quantities are not taken into account.



**Fig. 1.** The hierarchy forming the domain knowledge used in the running example with the generalization costs as retrieval knowledge.

Therefore, this recipe can be represented by the following RDFS base:

$$\text{Source} = \left\{ \begin{array}{l} \langle s \text{ type SourceCase} \rangle, \langle s \text{ dishType CocktailDish} \rangle, \\ \langle s \text{ ingredient Tequila} \rangle, \langle s \text{ ingredient GuavaJuice} \rangle, \\ \langle s \text{ ingredient PineappleJuice} \rangle, \langle s \text{ ingredient AppleJuice} \rangle, \\ \langle s \text{ ingredient VanillaSugar} \rangle, \langle s \text{ ingredient Mint} \rangle \end{array} \right\}$$

More generally, a case of the case base is identified by a resource  $s$  and is defined by an RDFS base  $\text{Source} \subseteq \text{KB}$  containing the triples of the form  $\langle s \text{ prop val} \rangle$ : this is a simple feature-value representation (the attribute being  $\text{prop}$  and the value being  $\text{val}$ ). The triple  $\langle s \text{ type SourceCase} \rangle$  is mandatory to indicate that this case belongs to  $\text{CaseBase}$ , the case base: it enables the obtention of the whole

case base via the following SPARQL query:

```
SELECT ?s
WHERE { ?s type SourceCase }
```

**Representation of domain knowledge.** The domain knowledge is represented by an RDFS base  $\text{DK}$  ( $\text{DK} \subseteq \text{KB}$ ) consisting in a set of triples in the form  $\langle C \text{ subc } D \rangle$ . Fig. 1 represents the domain knowledge for the running example by a hierarchy the edges  $C \xrightarrow{x} D$  of which represent the triples  $\langle C \text{ subc } D \rangle$  (the meaning of  $x$  is explained hereafter).

**Representation of similarity (retrieval knowledge).** The retrieval knowledge is encoded by a cost function, associating a triple  $\langle C \text{ subc } D \rangle \in \text{DK}$  to a positive real number  $\text{cost}(\langle C \text{ subc } D \rangle)$  for a given property. If  $C \xrightarrow{x} D$  is an edge of the Fig. 1 hierarchy then  $\text{cost}(\langle C \text{ subc } D \rangle) = x$ . This cost can be understood intuitively as the measure of “the generalization effort” from  $C$  to  $D$ . The cost function is assumed to be additive:

$$\text{cost}(\langle C \text{ subc } E \rangle) = \text{cost}(\langle C \text{ subc } D \rangle) + \text{cost}(\langle D \text{ subc } E \rangle)$$

Therefore  $\text{cost}(\langle \text{OrangeJuice subc FruitJuice} \rangle) = 0.10 + 0.02 = 0.12$  can be deduced from Fig. 1 and, for any RDFS class  $C$ ,  $\text{cost}(\langle C \text{ subc } C \rangle) = 0$ .

Since it is a tedious work to set manually all these costs, some default values are computed by TUURBINE according to the following formula:

$$\text{cost}(\langle C \text{ subc } D \rangle) = K \cdot \frac{\#\text{CasesWith}(D) - \#\text{CasesWith}(C)}{|\text{CaseBase}|}$$

$$\text{where } \#\text{CasesWith}(X) = \left| \text{Result}_{\perp} \left( \begin{array}{l} \text{SELECT } ?x \\ \text{WHERE } \{ ?x \text{ type SourceCase } ., \text{KB} \\ \quad ?x ?p X \} \end{array} \right) \right|,$$

$K$  is a coefficient cost factor depending of  $?p$

and  $|A|$  is the cardinal of the set  $A$

A similar formula is used by the TAAABLE system [1].

**Representation of adaptation knowledge.** The adaptation of a case **Source** to answer a query  $Q$  is performed using domain knowledge, retrieval knowledge and, when available, adaptation knowledge in the form of a finite set **AK** of adaptation rules. Syntactically, an adaptation rule is defined as follows:

$$\text{AR} = \begin{array}{ll} \mathbf{p}\text{-context } p_1, p_2, \dots & // \text{ positive context} \\ \mathbf{n}\text{-context } n_1, n_2, \dots & // \text{ negative context} \\ \mathbf{replace } r_1, r_2, \dots & // \text{ left part} \\ \mathbf{with } w_1, w_2, \dots & // \text{ right part} \end{array}$$

where  $p_1, p_2, \dots, n_1, n_2, \dots, r_1, r_2, \dots, w_1, w_2, \dots$  are expressions in the form  $\text{prop} : C$ ,  $\text{prop}$  being an RDFS property and  $C$  being an RDFS class. For example, the adaptation rule introduced in the running example (§3.1) can be formalized by

$$\text{AR}_1 = \begin{array}{l} \mathbf{p}\text{-context } \text{dishType} : \text{CocktailDish} \\ \mathbf{n}\text{-context } \text{ingredient} : \text{Anise} \\ \mathbf{replace } \text{ingredient} : \text{GuavaJuice}, \text{ingredient} : \text{VanillaSugar} \\ \mathbf{with } \text{ingredient} : \text{OrangeJuice}, \text{ingredient} : \text{SugarCaneSyrup} \end{array}$$

Let **Source** be a case identified by  $s$ . **AR** is applicable on **Source** if:

- For each term  $t = \text{prop} : C \in \{p_1, p_2, \dots, r_1, r_2, \dots\}$ ,  $\text{KB} \vdash \langle s \text{ prop } C \rangle$  (recall that  $\text{Source} \subseteq \text{KB}$ );
- For each term  $t = \text{prop} : C \in \{n_1, n_2, \dots\}$ ,  $\text{KB} \not\vdash \langle s \text{ prop } C \rangle$ .<sup>6</sup>

If **AR** is applicable on **Source** then the application of **AR** on **Source** gives a case **AR(Source)** obtained by:

<sup>6</sup> This amounts to a closed world assumption (CWA) defined by the inference rule  $\frac{\text{KB} \not\vdash \langle s \text{ prop } C \rangle}{\neg \langle s \text{ prop } C \rangle}$ . CWA is justified by the fact that there is no negation in RDFS. From a CBR viewpoint, this means that a case represents a specific situation (it is not a generalized case [6]): every fact  $\tau$  that is expressible in the KB vocabulary that is not entailed by **Source** is assumed not to hold for **Source**.



- Deleting all the triples of `Source` matching  $r_1, r_2, \dots$ ;
- Adding all the triples  $\langle s \text{ prop } C \rangle$  for each  $\text{prop}:C \in \{w_1, w_2, \dots\}$ ;
- Substituting all the occurrences of  $s$  by  $s'$  (a new case identifier).

For example, the rule  $\text{AR}_1$  is applicable on the recipe `Source` of the running example and  $\text{AR}_1(\text{Source}) = \text{AS}$  (adapted `Source`) with:

$$\text{AS} = \left\{ \begin{array}{l} \langle s' \text{ type } \text{SourceCase} \rangle, \langle s' \text{ dishType } \text{CocktailDish} \rangle, \\ \langle s' \text{ ingredient } \text{Tequila} \rangle, \langle s' \text{ ingredient } \text{PineappleJuice} \rangle, \\ \langle s' \text{ ingredient } \text{AppleJuice} \rangle, \langle s' \text{ ingredient } \text{Mint} \rangle, \\ \langle s' \text{ ingredient } \text{OrangeJuice} \rangle, \langle s' \text{ ingredient } \text{SugarCaneSyrup} \rangle \end{array} \right\}$$

In other words, the application of  $\text{AR}_1$  consists in applying the following substitution to the ingredients:

$$\Sigma = \text{GuavaJuice} \wedge \text{VanillaSugar} \rightsquigarrow \text{OrangeJuice} \wedge \text{SugarCaneSyrup}$$

Any adaptation rule  $\text{AR} \in \text{AK}$  has an associated value  $\text{cost}(\text{AR}) > 0$ , used during the adaptation process (see Section 3.5).

Finally, TUURBINE proposes another kind of adaptation rules called “specific adaptation rules” (SAR) with  $\text{cost}(\text{SAR}) = 0$ . Such a rule is associated to a source case which constitutes its context. For the recipe application, this rule can be seen as a way of encoding variants of the recipe. For example, the piece of information “Basil can be used instead of mint in the recipe of the running example” could be represented by the following specific adaptation rule:

$$\text{SAR}_1 = \boxed{\begin{array}{l} \text{p-context "Mexican cocktail in my way" recipe} \\ \text{replace ingredient: Mint} \\ \text{with ingredient: Basil} \end{array}}$$

### 3.3 Representation of Tuurbine queries

Syntactically, a TUURBINE query is a conjunction of expressions in the form  $\text{sign prop:val}$  where  $\text{sign} \in \{\epsilon, +, !, -\}$ ,  $\text{prop}$  is an RDF property and  $\text{val}$  is either a resource representing a class or a literal. For example, the following query is a TUURBINE translation of the query (1):<sup>7</sup>

$$\begin{aligned} \text{Q} = & +\text{dishType:CocktailDish} \wedge \text{ingredient:Mint} \\ & \wedge \text{ingredient:Gin} \wedge \text{ingredient:OrangeJuice} \wedge !\text{ingredient:Wine} \end{aligned} \quad (4)$$

The signs  $\epsilon$  and  $+$  are “positive signs”: they prefix features that the requested case must have.  $+$  indicates that this feature must also occur in the source case whereas  $\epsilon$  indicates that the source case may not have this feature, thus the adaptation phase has to make it appear in the final case.

<sup>7</sup>  $\epsilon$  denotes the empty word. Thus,  $\epsilon\text{prop:val}$  is simply written  $\text{prop:val}$ . The  $+$  in front of  $\text{dishType:CocktailDish}$  means that the part of the case base searched corresponds to cocktail recipes.

The signs ! and – are “negative signs”: they prefix features that the requested case must not have. – indicates that this feature must not occur in the source case whereas ! indicates that the source case may have this feature, and that the adaptation phase has to remove it.

### 3.4 Case retrieval in Tuurbine

Let  $Q$  be a TUURBINE query. The goal of retrieval is to find the cases  $\text{Source} \in \text{CaseBase}$  that best match  $Q$ . If no source case exactly match  $Q$ , then the query is relaxed and an approximate matching is searched.

**Exact matching search.**  $Q$  can be written  $Q = \bigwedge_i \text{sign}_i \text{prop}_i : \text{val}_i$ . For each

$i$ , let us consider the SPARQL query  $\text{SQ}_i =$

<pre> SELECT ?s WHERE {   ?s type SourceCase .   ?s prop_i ?x .   ?x subc val_i }                 </pre>
--

For example, the element of query  $\text{+dishType:CocktailDish}$  gives a SPARQL query that can be read “get the source cases with a cocktail dish type” (i.e., the  $?s$  such that  $?s$  is a instance of  $\text{SourceCase}$  and such that  $?s$  has a dish type  $?x$  which is a subclass of  $\text{CocktailDish}$ ).

The exact matching of the query  $Q$  can be done by executing the SPARQL queries  $\text{SQ}_i$  and combining them as follows (EMS stands for “Exact Matching Search”):

$$\text{EMS}(Q) = \bigcap_{i, \text{sign}_i \in \{\epsilon, +\}} \text{Result}_+(\text{SQ}_i) \setminus \bigcup_{i, \text{sign}_i \in \{!, -\}} \text{Result}_+(\text{SQ}_i)$$

In other words, the result of the exact matching search are the source cases that match the SPARQL queries  $\text{SQ}_i$  such that  $\text{sign}_i$  is a positive sign and that does not match the SPARQL queries  $\text{SQ}_i$  such that  $\text{sign}_i$  is a negative sign.<sup>8</sup>

**Approximate search.** The principle of this search is to find a generalization function  $\Gamma$  with minimal cost such that the execution of the query  $Q$  modified by  $\Gamma$  returns at least one source case:  $\text{EMS}(\Gamma(Q)) \neq \emptyset$ .

Let  $Q = \bigwedge_i \text{sign}_i \text{prop}_i : \text{val}_i$  be a query. A one step-generalization  $\gamma(Q)$  of  $Q$  consists in generalizing a term  $\text{sign}_i \text{prop}_i : \text{val}_i$  such that  $\text{sign}_i \notin \{+, -\}$ :

- If  $\text{sign}_i = \epsilon$  and  $\text{val}_i$  is an RDFS class, then the generalizations of this term are the terms  $\text{sign}_i \text{prop}_i : \text{val}$  such that  $\langle \text{val}_i \text{ subc val} \rangle \in \text{DK}$ . This one step generalization is written  $\gamma = \text{prop}_i : \text{val}_i \rightsquigarrow \text{prop}_i : \text{val}$  or, simply,  $\gamma = \text{val}_i \rightsquigarrow \text{val}$ . The cost of such a generalization is  $\text{cost}(\langle \text{val}_i \text{ subc val} \rangle)$ .
- If  $\text{sign}_i = !$  or if  $\text{sign}_i = \epsilon$  and  $\text{val}_i$  is a literal, then  $\text{val}_i$  is directly generalized to the ontology top  $\top$ . This one step generalization is written  $\gamma = \text{val}_i \rightsquigarrow \top$ . The cost of such a generalization is 0.

<sup>8</sup> In practice,  $\text{EMS}(Q)$  could be computed thanks to the execution of fewer SPARQL queries thus giving the same result with a lower computational cost.

A generalization function  $\Gamma$  is a composition of one-step generalizations  $\gamma_1, \gamma_2, \dots, \gamma_n$ :  $\Gamma = \gamma_n \circ \dots \circ \gamma_2 \circ \gamma_1$ . Its cost is the sum of the costs of  $\gamma_i$ . For example, the generalization function  $\Gamma$  can be applied on  $\mathbb{Q}$  defined by equation (4):

$$\begin{aligned}
\Gamma &= \text{Gin} \rightsquigarrow \text{Liquor} \circ \text{OrangeJuice} \rightsquigarrow \text{FruitJuice} \\
\text{and } \Gamma(\mathbb{Q}) &= +\text{dishType}:\text{CocktailDish} \wedge \text{ingredient}:\text{Mint} \\
&\quad \wedge \text{ingredient}:\text{Liquor} \wedge \text{ingredient}:\text{FruitJuice} \\
&\quad \wedge !\text{ingredient}:\text{Wine} \\
\text{cost}(\Gamma) &= \text{cost}(\langle \text{Gin} \text{ subc } \text{Liquor} \rangle) \\
&\quad + \text{cost}(\langle \text{OrangeJuice} \text{ subc } \text{FruitJuice} \rangle)
\end{aligned} \tag{5}$$

The execution of the query  $\Gamma(\mathbb{Q})$  returns the recipe of the example:  $\text{Source} \in \text{EMS}(\Gamma(\mathbb{Q}))$  so, provided that  $\text{cost}(\Gamma)$  is the minimum of the costs of the generalization functions  $\Lambda$  such that  $\text{EMS}(\Lambda(\mathbb{Q})) \neq \emptyset$ ,  $\text{Source}$  is a retrieved case. Technically,  $\Gamma$  is searched by increasing cost in a generalization function space.

**Searching for less similar cases.** It may occur that a user of a TUURBINE application wants to find other cases than the ones returned in a first launch. TUURBINE offers the possibility to do so: it simply consists in resuming the search after  $\Gamma$  has been found. This way, a second generalization function  $\Gamma'$  can be found. For example:

$$\Gamma' = \text{Gin} \rightsquigarrow \text{Alcohol} \circ \text{Mint} \rightsquigarrow \text{Herb}$$

### 3.5 Case adaptation in Tuurbine

There are two adaptation processes in TUURBINE.

The first adaptation process consists in obtaining the matching between  $\text{Source}$  and  $\mathbb{Q}$  that has permitted the retrieval of  $\text{Source}$ , this matching being composed of the matching between  $\text{Source}$  and  $\Gamma(\mathbb{Q})$  (given by the fact that  $\text{Source} \in \text{EMS}(\Gamma(\mathbb{Q}))$ ) and the matching between  $\Gamma(\mathbb{Q})$  and  $\mathbb{Q}$  (given by the generalization function  $\Gamma$ ). This first kind of adaptation, for the running example, works as follows.  $\Gamma(\mathbb{Q})$  defined by equation (5) matches exactly  $\text{Source}$  (wrt to DK). In fact, there are several matchings between  $\text{Source}$  and  $\Gamma(\mathbb{Q})$ :

Tequila matches Liquor

each  $ing \in \{\text{GuavaJuice}, \text{PineappleJuice}, \text{AppleJuice}\}$  matches FruitJuice

By composing this matching between  $\text{Source}$  and  $\Gamma(\mathbb{Q})$  and the matching between  $\Gamma(\mathbb{Q})$  and  $\mathbb{Q}$  given by  $\Gamma$ , it comes the following adaptation:

In the ‘‘Mexican cocktail in my way’’ recipe, substitute Tequila by Gin

and substitute  $\left. \begin{array}{l} \text{GuavaJuice and/or} \\ \text{PineappleJuice and/or} \\ \text{AppleJuice} \end{array} \right|$  by OrangeJuice.

which represents the expected adaptation (2).

The second kind of adaptation uses the adaptation rules  $AR \in AK$  and the query generalization procedure used for case retrieval. More precisely it searches a pair  $(\Sigma, \Gamma)$  such that:

- $\Sigma(\text{Source})$  matches exactly  $\Gamma(Q)$  (wrt DK);
- $\Sigma$  is a composition of adaptation rules that is applicable on **Source** ( $\Sigma = AR_p \circ \dots \circ AR_2 \circ AR_1$  such that  $AR_1$  is applicable on **Source**,  $AR_2$  is applicable on  $AR_1(\text{Source})$ , etc.);
- $\Gamma$  is a generalization function;
- $\text{cost}(\Sigma) + \text{cost}(\Gamma)$  is minimal (where  $\text{cost}(\Sigma) = \text{cost}(AR_p) + \dots + \text{cost}(AR_2) + \text{cost}(AR_1)$ ).

Technically, this kind of adaptation relies on a best-first search: the states are pairs  $(\Sigma, \Gamma)$ ; a final state is such that  $\Sigma(\text{Source})$  matches exactly  $\Gamma(Q)$ ; the state space is searched by increasing  $\text{cost}(\Sigma) + \text{cost}(\Gamma)$ . It must be noticed that:

- $\text{cost}(\Sigma) + \text{cost}(\Gamma) \leq \text{cost}(\Gamma_{\text{retrieval}})$  where  $\Gamma_{\text{retrieval}}$  is the generalization function generated during retrieval (and used by the first kind of adaptation);
- If  $AK = \emptyset$  then  $\Sigma$  is the identity function and  $\Gamma = \Gamma_{\text{retrieval}}$ .

For the example, if  $AK = \{AR_1\}$  and  $\text{cost}(AR_1) + \text{cost}(\langle \text{Gin subc Liquor} \rangle) < \text{cost}(\Gamma_{\text{retrieval}})$  then the adaptation coincides with (3).

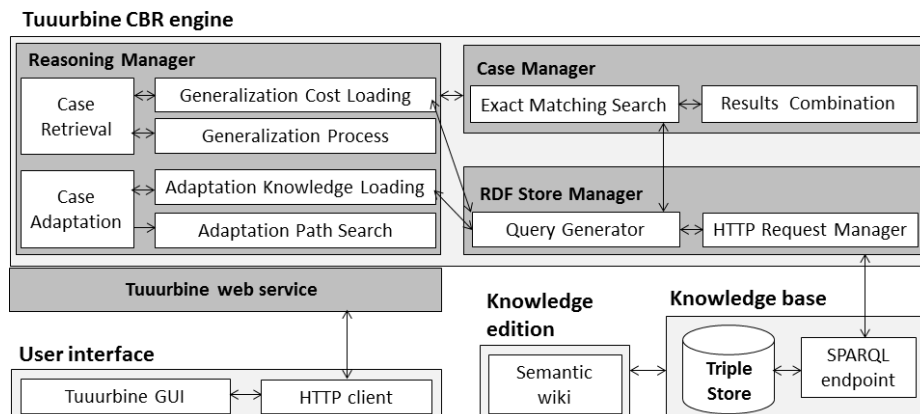
## 4 Implementation

The general architecture of TUURBINE is presented in Fig. 2. Its CBR engine is composed of three main modules: the RDF store manager, the case manager and the reasoning manager. The case manager implements EMS (the exact matching search). The case retrieval triggers the RDF store manager, in charge of generating and executing a set of SPARQL queries, using a SPARQL endpoint. The results are returned to the case manager which combines them to build the cases.

When no case satisfies the EMS, the reasoning manager is triggered for an approximative search. A first step consists in loading the generalization costs from the triple store (only the costs that are potentially useful). This is done once, at the beginning of the approximative search process. The generalisation process iterates until the EMS returns at least one case for a generalized query. These cases are then adapted using possibly adaptation knowledge and generalization knowledge as explained in §3.5.

TUURBINE is implemented as a web service. This allows to use HTTP requests to query the CBR engine, but a generic user interface is also provided (an instantiation of this GUI is visible in Fig. 3). Additionally, a Semantic Wiki may be used to manage the whole knowledge in a more convenient way, as it is done for example in WIKITAAABLE, a semantic wiki in the domain of cooking.

TUURBINE required about 6000 lines of Java code. The user interface is developed in PHP, javascript, HTML/CSS. The exchange between the interface



**Fig. 2.** Architecture of TUURBINE

and the CBR engine, as well as the configuration files are based on JSON encoding language. TUURBINE is distributed under an Affero GPL Licence and is available from <http://tuurbine.loria.fr/>.

The retrieval time for the query (4) is 3156 ms on a server Intel Xeon E5520 64 bits with  $8 \times 2.27$  GHz processors, and a 32 GB RAM, running on Linux Ubuntu 12.4. 85% of this time is for disk access and data communication. It required the execution of 117 SPARQL queries. This example is realistic: the RDF base size is 277 megabytes, it contains more than  $10^6$  triples, 1641 of them being triples of the form  $\langle s \text{ type } \text{Recipe} \rangle$  (hence 1641 source cases) and 2316 triples of the form  $\langle C \text{ subc } D \rangle$  constitute the domain knowledge. A systematic performance study of TUURBINE remains to be done as future work.

## 5 Tuurbine in action

The first use case is the instantiation of TUURBINE in the cooking domain, as the new version of TAAABLE, a CBR system which retrieves and creates cooking recipes by adaptation that has been developed to participate in the Computer Cooking Contest<sup>9</sup> since 2008 [3]. Fig. 3 presents the TUURBINE interface running query 4 and is reachable at <http://tuurbine.loria.fr/taabile/>. The TAAABLE knowledge base (<http://wikitaaable.loria.fr/>) is composed of the four classical knowledge containers: (1) the domain knowledge as an ontology of the cooking domain which includes several hierarchies (about food, dish types, etc.), (2) the case base, which are recipes described by their titles, the dish type they produces, the ingredients which are required, the preparation steps, etc., (3) the adaptation knowledge takes the form of adaptation rules as

<sup>9</sup> <http://computercookingcontest.net>

**Search**

Enter a query

**Example:** For a result containing apple without cinnamon : **+apple -cinnamon**

Need some **help?**

■ Your request is: « **+dish type:cocktail dish, ingredient:gin, ingredient:orange juice, ingredient:mint, !ingredient:wine** »

#	Original result name (click to get more details)	Adaptations	Explanations
1	Mexican cocktail in my way	Guava juice and/or Apple juice and/or Pineapple juice → Orange juice Tequila → Gin	Ccc cocktail dish → Cocktail dish

Result: 1-1 on 1 [9.26 seconds]

**Fig. 3.** Generic TUURBINE interface applied to the cooking domain.

introduced in this paper, and (4) the retrieval knowledge which is stored as cost values on subclass-of relations and adaptation rules.

To show that TUURBINE is independent of the application domain, a second instantiation of TUURBINE has been implemented in the domain of music. This application use case is about searching for a piece of music you may play using a given set of instruments. The cases are pieces of music described by its composer(s), its music genre(s), origin(s) and year of creation, and the instruments required for playing them. The adaptation consists in replacing some instruments by other ones. The music ontology and the TUURBINE interface for the music CBR system can be respectively reached at <http://muuusic.loria.fr/> and <http://tuuurbine.loria.fr/muuusic/>.

## 6 Discussion and related work

Several tools address CBR in a generic way. myCBR and jCOLIBRI are probably the most famous of them. jCOLIBRI [7] is an object-oriented framework for developing CBR applications. This framework includes connectors with several kinds of data sources (database, XML, etc.) for loading the cases in cache, on which the retrieval, reuse, revise and retain tasks can be performed. Specialized modules take into account various types of CBR applications, like CBR on textual cases, data or knowledge-intensive CBR. With jCOLIBRI, the knowledge intensive CBR approach consists in retrieving cases according to an ontology based similarity measure and in replacing some case features (instance of an ontology concept) with the closest ones in the ontology (based on similarity minimization). TUURBINE also follows this substitution approach but the way the retrieval and the adaptation of cases are made is very different. Indeed, using

the structure of the ontology only in order to compute a numerical measure has a major limitation, due to the lack of semantics of the measure, which can, moreover, be computed with various similarity functions (cosine, fdeep, etc.) which take into account various criteria (depth of the two concepts  $c_1$  and  $c_2$  which are compared, depth and relations of  $c_1$  and  $c_2$  with their least common subsumer, etc.). Such similarity functions do not take into account the distribution of the instances in the set of cases, in comparison with our cost function which represents a generalization effort. Moreover, the cost function is independent of the level of structuration of the ontology. So, having many intermediate concepts (like for example `CitrusFruitJuice` between `OrangeJuice` and `FruitJuice`) do not impact the case retrieval. This point is crucial for improving case retrieval with a better structured ontology [8].

myCBR [9] is a Java open source similarity-based tool. Cases, described by attributes with different weights, may be created thanks to a graphical user interface or through a file import (e.g., CSV file). A “Linked Open Data Connector” provides an access to open data sources for building taxonomies. Many similarity functions are provided and may be combined to rank cases according to their similarity with the query (during the retrieval step). The adaptation consists in replacing some case features with others, using adaptation rules retrieved from the myCBR knowledge models (e.g. similarity tables or taxonomies). So, the myCBR retrieving/adapting procedure is rather similar to the jCOLIBRI one.

Like myCBR and jCOLIBRI, TUURBINE is a generic CBR engine. The originality of TUURBINE is the implementation of a CBR based on RDFS and the exploitation of the semantics of the RDF and RDFS models. The whole knowledge (cases, domain knowledge, costs, adaptation rules) is described using RDF and RDFS and is managed outside TUURBINE using a triple store and SPARQL, two standard tools of the Semantic Web. The edition of the knowledge is also facilitated with the use of a Semantic Wiki, another Semantic Web tool.

## 7 Conclusion

This paper has presented TUURBINE, a generic CBR engine based on RDFS and using Semantic Web technologies. The main principles for representation of the four knowledge containers as well as the reasoning processes have been detailed. Information about technical use and configuration of TUURBINE is available at <http://tuurbine.loria.fr/>. TUURBINE is used in the new version of the TAAABLE system, a CBR system adapting cooking recipes, which will participate in the 2014’s Computer Cooking Contest.

The current version of the TUURBINE case-based inference engine using query generalizations is based on the `subc` property (generalization of a class by a super-class). Other constructs of RDFS could be used as well in the TUURBINE reasoning process: the subproperty relation (`subp`), domains and ranges of properties, literals (e.g., numerical values). This constitutes a future direction of research. It is planned to address first the use of `subp`. For instance, if  $\langle \text{mainIngredient subp ingredient} \rangle \in \text{DK}$  (if  $x$  is a main ingredient of  $y$  then

$x$  is an ingredient of  $y$ ), the query  $Q = \text{mainIngredient:TomatoJuice}$  can be generalized into  $I(Q) = \text{ingredient:TomatoJuice}$ .

Another direction of work related to RDFS expressiveness consists in managing “deeper” cases: in the current implementation, only values directly related to  $s$  are taken into account, but such values could be related to other values participating to the case representation. Following this direction would make TUURBINE evolve from an attribute-value representation to something similar to an object-based representation (like the ones of jCOLIBRI and myCBR).

A third direction of work is the integration of TUURBINE with other generic CBR systems such as jCOLIBRI and myCBR (cf. Section 6) or Revisor [10].

## References

1. A. Cordier, V. Dufour-Lussier, J. Lieber, E. Nauer, F. Badra, J. Cojan, E. Gaillard, L. Infante-Blanco, P. Molli, A. Napoli, and H. Skaf-Molli. Taaable: a Case-Based System for personalized Cooking. In Stefania Montani and Lakhmi C. Jain, editors, *Successful Case-based Reasoning Applications-2*, volume 494 of *Studies in Computational Intelligence*, pages 121–162. Springer, January 2014.
2. A. Cordier, J. Lieber, P. Molli, E. Nauer, H. Skaf-Molli, and Y. Toussaint. WIKITAAABLE: A semantic wiki as a blackboard for a textual case-based reasoning system. In *SemWiki 2009 - 4rd Semantic Wiki Workshop at the 6th European Semantic Web Conference - ESWC 2009*, Heraklion, Grèce, May 2009.
3. F. Badra, R. Bendaoud, R. Bentebibel, P.-A. Champin, J. Cojan, A. Cordier, S. Després, S. Jean-Daubias, J. Lieber, T. Meilender, A. Mille, E. Nauer, A. Napoli, and Y. Toussaint. TAAABLE: Text Mining, Ontology Engineering, and Hierarchical Classification for Textual Case-Based Cooking. In Martin Schaaf, editor, *9th European Conference on Case-Based Reasoning - ECCBR 2008, Workshop Proceedings*, pages 219–228, Trier, Allemagne, 2008.
4. W. Shakespeare. *Romeo and Juliet*. 1597.
5. M. M. Richter. Introduction. In M. Lenz, B. Bartsch-Spörl, H.-D. Burkhard, and S. Wess, editors, *Case-Based Reasoning Technologies. From Foundations to Applications*, LNAI 1400, chapter 1, pages 1–15. Springer, 1998.
6. K. Maximini, R. Maximini, and R. Bergmann. An investigation of generalized cases. In K. D. Ashley and D. Bridge, editors, *Proceedings of the 5th International Conference on Case Base Reasoning (ICCBR'03)*, volume 2689 of *LNAI*, pages 261–275, Trondheim, Norway, June 2003. Springer.
7. Juan A. Recio-García, Pedro A. González-Calero, and Belén Díaz-Agudo. jcolibri2: A framework for building case-based reasoning systems. *Science of Computer Programming*, (79):126–145–, 2014.
8. V. Dufour-Lussier, J. Lieber, E. Nauer, and Y. Toussaint. Improving case retrieval by enrichment of the domain ontology. In *19th International Conference on Case Based Reasoning - ICCBR'2011*, London, 2011.
9. K. Bach and K. Althoff. Developing case-based reasoning applications using mycbr 3. In Belén Díaz-Agudo and Ian Watson, editors, *ICCBR*, volume 7466 of *LNAI*, pages 17–31. Springer, 2012.
10. J. Cojan and J. Lieber. Applying belief revision to case-based reasoning. In Henri Prade and Gilles Richard, editors, *Computational Approaches to Analogical Reasoning: Current Trends*, volume 548 of *Studies in Computational Intelligence*. Springer, 2014.