

# Reasoning about Uncertain Information and Conflict Resolution through Trust Revision \*

Murat Şensoy<sup>1,4</sup>, Achille Fokoue<sup>2</sup>, Jeff Z. Pan<sup>1</sup>, Timothy J. Norman<sup>1</sup>,  
Yuqing Tang<sup>3</sup>, Nir Oren<sup>1</sup>, and Katia Sycara<sup>3</sup>

<sup>1</sup>Computing Science, University of Aberdeen, UK

<sup>2</sup>IBM T. J. Watson Research Center, NY, US

<sup>3</sup>Carnegie Mellon University, Pittsburgh, US

<sup>4</sup>Computer Science, Ozyegin University, Istanbul, Turkey

{m.sensoy, jeff.z.pan, t.j.norman, n.oren}@abdn.ac.uk,  
achille@us.ibm.com, {yuqing.tang, katia}@cs.cmu.edu

## ABSTRACT

In information driven MAS, information consumers collect information about their environment from various sources such as sensors. However, there is no guarantee that a source will provide the requested information truthfully and correctly. Even if information is provided only by trustworthy sources, it can contain conflicts that hamper its usability. In this paper, we propose to exploit such conflicts to revise trust in information. This requires a reasoning mechanism that can accommodate domain constraints, uncertainty, and trust. Our formalism — *SDL-Lite* — is an extension of a tractable subset of Description Logics with Dempster-Shafer theory of evidence. *SDL-Lite* allows reasoning about uncertain information and enables conflict detection. Then, we propose methods for conflict resolution through trust revision and analyse them through simulations. We show that the proposed methods allow reasonably accurate estimations of trust in information in realistic settings.

## Categories and Subject Descriptors

I.2.11 [Distributed Artificial Intelligence]: Multiagent Systems

## Keywords

Information Fusion, Trust, Uncertainty, Description Logics

## 1. INTRODUCTION

Uncertainty is a core feature of many domains in which agents are expected to operate. In such environments, information sources

\*This research was sponsored by the U.S. Army Research Laboratory and the U.K. Ministry of Defence and was accomplished under Agreement Number W911NF-06-3-0001. The views and conclusions contained in this document are those of the author(s) and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Army Research Laboratory, the U.S. Government, the U.K. Ministry of Defence or the U.K. Government. The U.S. and U.K. Governments are authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation hereon. This work is partially supported by the EU K-Drive project. Dr. Şensoy thanks to The Scientific and Technological Research Council of Turkey (TUBITAK) for its support under grant 111K476.

**Appears in:** *Proceedings of the 12th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2013)*, Ito, Jonker, Gini, and Shehory (eds.), May, 6–10, 2013, Saint Paul, Minnesota, USA.

Copyright © 2013, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

such as *sensors* provide agents with state information. However, in the context of a multi-agent system, different sets of sensors can be controlled by different agents, each with their own capabilities. In order to reason about the state of the environment, an agent must therefore request potentially noisy, incomplete, or misleading information from potentially untrustworthy agents. By obtaining information from multiple sources, the agent can build up a more accurate view of its environment than by utilising its own sensors alone. To achieve such a task, the agent needs to reason about the noisy, incomplete, and misleading information. This paper addresses such a need with a mechanism that can reason about uncertain information and conflict resolution through trust revision.

Various researchers have examined aspects of this problem — work in computational trust (e.g. [10]) is intended to allow an agent to determine which other agents should be asked for information, while work on information fusion [13] examines how incomplete and noisy information from different sources should be combined in order to obtain a true picture of the environment. However, neither work in isolation considers both sources of uncertainty in the domain, namely trustworthiness of information sources as well as incompleteness and vagueness of the provided information.

This paper makes two core contributions: Firstly, we combine *Description Logics* (DLs) [1] and *Dempster-Shafer theory* (DST) [15] to create a computationally tractable framework for reasoning about uncertain information obtained from different agents. This framework enables us to detect conflicts in uncertain information due to constraints imposed by the domain. It also enables us to resolve such conflicts. Secondly, we show how trust in uncertain information can be revised when additional information is received. We model this problem as an optimisation problem, and propose heuristics that allow us to identify high quality solutions.

The remainder of the paper is structured as follows. We begin by introducing Description Logics and Dempster-Shafer theory. Section 3 describes *SDL-Lite*, our extension of *DL-Lite* with subjective opinions. In Section 4, we concentrate on the problem of how much trust should be placed in a source of information, and we describe how conflicting information sources can be made consistent. We evaluate our approach in Section 5 and discuss it with respect to the existing work and future research directions in Section 6.

## 2. PRELIMINARIES

### 2.1 DL-based Ontologies

Due to limited space, we do not provide a full overview of Description Logics (DLs), but rather point the reader to [1]. We note,

however, that even for the smallest propositionally closed DL,  $\mathcal{ALC}$  (which only provides class constructors  $\neg C, C \sqcap D, C \sqcup D, \exists R.C$  and  $\forall R.C$ ), the complexity of logical entailment is EXPTIME. Recently, Calvanese *et al.* [5] proposed DL-Lite, which can express most features in UML class diagrams with a low reasoning overhead (with data complexity  $\text{AC}_0$ ). It is for this reason that we base our model on DL-Lite<sub>core</sub> (referred to here as DL-Lite, although there are extensions [3]), and hence provide a brief formalisation to ground the subsequent presentation of our model.

A DL-Lite knowledge base  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  consists of a TBox  $\mathcal{T}$  and an ABox  $\mathcal{A}$ . Axioms of the following forms compose  $\mathcal{K}$ :

1. class inclusion axioms:  $B \sqsubseteq C \in \mathcal{T}$  where  $B$  is a basic class  $B := A \mid \exists R \mid \exists R^-$  and  $C$  is a general class  $C := B \mid \neg B \mid C_1 \sqcap C_2$  (where  $A$  denotes an named class,  $R$  denotes a named property, and  $R^-$  is the inverse of  $R$ );
2. individual axioms:  $B(a), R(a, b) \in \mathcal{A}$  where  $a$  and  $b$  are named individuals.

Description Logics have a well-defined model-theoretic semantics, which are provided in terms of interpretations. An interpretation  $\mathcal{I}$  is a pair  $(\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ , where  $\Delta^{\mathcal{I}}$  is a non-empty set of objects and  $\cdot^{\mathcal{I}}$  is an interpretation function, which maps each class  $C$  to a subset  $C^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$  and each property  $R$  to a subset  $R^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$ .

Using a trivial normalisation, it is possible to convert class inclusion axioms of the form  $B_1 \sqsubseteq C_1 \sqcap C_2$  into a set of simpler class inclusions of the form  $B_1 \sqsubseteq B_i$  or  $B_1 \sqsubseteq \neg B_j$ , where  $B_1, B_i$ , and  $B_j$  are basic concepts [5]. For instance, during normalisation,  $B_1 \sqsubseteq B_2 \sqcap \neg B_3$  is replaced with  $B_1 \sqsubseteq B_2$  and  $B_1 \sqsubseteq \neg B_3$ . In Table 3, we define semantics over a normalised TBox for our variant of DL-Lite — SDL-Lite.

## 2.2 Subjective Opinions

Dempster-Shafer Theory (DST) offers a means to characterise an agent's view of the state of world by assigning *basic probability masses* to subsets of truth assignments of propositions in the logic. In DST, a *binomial opinion* about a proposition  $x$  is represented by a triple  $w_x = (b_x, d_x, u_x)$  which is derived from the basic probability masses assigned to subsets of truth assignments of the language. In the opinion  $w_x$ ,  $b_x$ , also denoted by  $b(w_x)$ , is the belief about  $x$  — the summation of the probability masses that entail  $x$ ;  $d_x$ , also denoted by  $d(w_x)$ , is the disbelief about  $x$  — the summation of the probability masses that entail  $\neg x$ ; and  $u_x$ , also denoted by  $u(w_x)$ , is the uncertainty about  $x$  — the summation of the probability masses that neither entail  $x$  nor entail  $\neg x$ . The constraints over the probability mass assignment function require that  $b_x + d_x + u_x = 1$  and  $b_x, d_x, u_x \in [0, 1]$ . When a more concise notation is necessary, we use  $(b_x, d_x)$  instead of  $(b_x, d_x, u_x)$ , since  $u_x = 1 - b_x - d_x$ . The negation over an opinion  $w_x$  is defined as  $\neg(b_x, d_x, u_x) = (d_x, b_x, u_x) = (b_{\neg x}, d_{\neg x}, u_{\neg x})$  [9].

**DEFINITION 1.** Let  $w_1 = (b_1, d_1, u_1)$  and  $w_2 = (b_2, d_2, u_2)$  be two opinions about the same proposition. We call  $w_1$  a specialisation of  $w_2$  ( $w_1 \preceq w_2$ ) iff  $b_2 \leq b_1$  and  $d_2 \leq d_1$  (implies  $u_1 \leq u_2$ ). Similarly, we call  $w_1$  a generalisation of  $w_2$  ( $w_2 \preceq w_1$ ) iff  $b_1 \leq b_2$  and  $d_1 \leq d_2$  (implies  $u_2 \leq u_1$ ). ■

An agent  $i$ 's opinion about a proposition  $x$  is denoted by  $w_x^i = (b_x^i, d_x^i, u_x^i)$ . This opinion  $w_x^i$  may not be directly used by another agent  $j$ . Agent  $j$  could have a view of the reliability or competence of  $i$  with respect to  $x$ . Shafer [15] proposed a discounting operator  $\otimes$  to normalise the belief and disbelief in  $w_x^j$  based on the degree of trust  $j$  has of  $i$  with respect to  $x$ :  $t_j^i$ . The normalised opinion,  $w_x^j$ , is computed as  $(b_x^i \times t_j^i, d_x^i \times t_j^i, u_x^i)$ .

The trustworthiness of information sources can be modelled using Beta probability density functions [10]. A Beta distribution has

**Table 1: Information sources and their trustworthiness**

Source	Definition	Evidence	Degree of trust
C	Local civilian sources	(4, 0)	0.83
P	Local police sources	(10, 3)	0.786
M <sub>2</sub>	Collaborating military forces	(50, 5)	0.89
A	Acoustic sensors of M <sub>1</sub>	(1000, 0)	0.999

**Table 2: DL-Lite TBox and opinions about ABox assertions**

Initial TBox	Opinions about ABox assertions
$BombedRoad \sqsubseteq SabotagedRoad$	$Blocked(R): (0.71, 0.09, 0.2)$
$\exists road BombedBy \sqsubseteq BombedRoad$	$Safe(R): (0.63, 0.066, 0.304)$
$SabotagedRoad \sqsubseteq \neg Safe \sqcap Blocked$	$BombedRoad(R): (0.2, 0.3, 0.5)$
Normalised TBox	
$BombedRoad \sqsubseteq SabotagedRoad$	
$\exists road BombedBy \sqsubseteq BombedRoad$	
$SabotagedRoad \sqsubseteq \neg Safe$	
$SabotagedRoad \sqsubseteq Blocked$	

two parameters  $(r + 1, s + 1)$ , where  $r$  is the amount of positive evidence and  $s$  is the amount of negative evidence for the trustworthiness agent  $i$  agent has for agent  $j$ . The degree of trust  $t_j^i$  is then computed as the expectation value of the Beta distribution:  $t_j^i = (r + 1)/(r + s + 2)$ .

We now provide a running example used throughout the remainder of the paper. Following this, we will present the semantics and reasoning mechanisms for a language that combines DL-Lite and DST with the view to offering a model for tractable reasoning with uncertain information and trust.

## 2.3 Example Scenario

Consider a region in which insurgents are active and where civilian groups are in need of support. An NGO operating in the region has identified a safe zone  $Z$  and aims to bring relief to the injured in village  $V$  by transporting them to the safe zone. There is only one road  $R$  between  $Z$  and  $V$ , but there is conflict between groups  $G_1$  and  $G_2$  within the region. As part of a multi-national peace effort,  $M_1$  and  $M_2$  operate within the region, and part of their remit is to protect and support NGOs. The resources available to  $M_2$  include Unmanned Aerial Vehicles (UAVs).  $M_1$  acts as liaison to the NGO, and has intelligence from the information sources listed in Table 1 along with models of the trustworthiness of these sources.  $M_1$  collects the following pieces of information from the sources in the area: 1)  $M_2$  informs  $M_1$  that  $R$  is blocked with opinion  $(0.8, 0.1, 0.1)$ ; 2)  $P$  reports that  $R$  is safe with opinion  $(0.8, 0.1, 0.1)$ ; and 3)  $A$  reports an explosion on route  $R$  with opinion  $(0.2, 0.3, 0.5)$ . This intelligence is interpreted by  $M_1$  given its degree of trust in the sources into Table 1. The opinions from information sources are discounted by the trustworthiness of the sources. For instance,  $M_2$ 's original opinion for 'R is blocked' is  $(0.8, 0.1, 0.1)$ , but it is discounted to  $(0.71, 0.09, 0.2)$  using  $M_2$ 's trustworthiness, i.e., 0.89 (see Table 2 Abox for other discounted opinions).

Having provided an overview of DLs and DST, we now turn to describing our core contribution — a description logic able to represent uncertainty in the manner of DST.

## 3. SUBJECTIVE DL-LITE

We propose *Subjective DL-Lite* (or *SDL-Lite* for short), which extends DL-Lite<sub>core</sub> with subjective opinion assertions of the form  $\mathcal{B}:w$ , where  $w$  is an opinion and  $\mathcal{B}$  is an ABox axiom (i.e., assertion). Each ABox axiom is associated with one opinion. ABox axioms have the form  $B(a)$  or  $R(a, b)$ , where  $B$  is basic class,  $R$  is a property, and  $a$  and  $b$  are individuals.

### 3.1 SDL-Lite Semantics

In common with DL-Lite ontologies, the semantics of an ontology in *SDL-Lite* is defined in terms of *subjective interpretations*. Let  $\mathcal{W}$  be the set of all possible subjective binary opinions. A subjective interpretation is a pair  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$  where the domain  $\Delta^{\mathcal{I}}$

**Table 3: Semantics of Subjective DL-Lite**

Syntax	Semantics
$\top$	$\top^{\mathcal{I}}(o) = (1, 0, 0)$
$\perp$	$\perp^{\mathcal{I}}(o) = (0, 1, 0)$
$\exists R$	$b((\exists R)^{\mathcal{I}}(o_1)) \geq \max_{\forall o_2} \{b(R^{\mathcal{I}}(o_1, o_2))\}$ and $d((\exists R)^{\mathcal{I}}(o_1)) \leq \min_{\forall o_2} \{d(R^{\mathcal{I}}(o_1, o_2))\}$
$\neg B$	$(\neg B)^{\mathcal{I}}(o) = \neg B^{\mathcal{I}}(o)$
$R^-$	$(R^-)^{\mathcal{I}}(o_2, o_1) = R^{\mathcal{I}}(o_1, o_2)$
$B_1 \sqsubseteq B_2$	$\forall o \in \Delta^{\mathcal{I}}, b(B_1^{\mathcal{I}}(o)) \leq b(B_2^{\mathcal{I}}(o))$ and $d(B_1^{\mathcal{I}}(o)) \leq d(B_2^{\mathcal{I}}(o))$
$B_1 \sqsubseteq \neg B_2$	$\forall o \in \Delta^{\mathcal{I}}, b(B_1^{\mathcal{I}}(o)) \leq d(B_2^{\mathcal{I}}(o))$ and $b(B_2^{\mathcal{I}}(o)) \leq d(B_1^{\mathcal{I}}(o))$
$B(a):w$	$b(w) \leq b(B^{\mathcal{I}}(a^{\mathcal{I}}))$ and $d(w) \leq d(B^{\mathcal{I}}(a^{\mathcal{I}}))$
$R(a, b):w$	$b(w) \leq b(R^{\mathcal{I}}(a^{\mathcal{I}}, b^{\mathcal{I}}))$ and $d(w) \leq d(R^{\mathcal{I}}(a^{\mathcal{I}}, b^{\mathcal{I}}))$

is a non-empty set of objects and  $\mathcal{I}$  is a subjective interpretation function, which maps:

- an individual  $a$  to an element of  $a^{\mathcal{I}} \in \Delta^{\mathcal{I}}$ ,
- a named class  $A$  to a function  $A^{\mathcal{I}} : \Delta^{\mathcal{I}} \rightarrow \mathcal{W}$ ,
- a named property  $R$  to a function  $R^{\mathcal{I}} : \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}} \rightarrow \mathcal{W}$ .

To provide a semantics for *SDL-Lite*, we extend interpretations of DL-Lite class and property descriptions, and of axioms under unique name assumption. The semantics are presented in Table 3.

The semantics of  $\exists R$  is derived from the rule  $R(a^{\mathcal{I}}, b^{\mathcal{I}}) \rightarrow \exists R(a^{\mathcal{I}}), \forall b^{\mathcal{I}} \in \Delta^{\mathcal{I}}$ . This rule constrains the minimum belief and the maximum disbelief that  $\exists R(a^{\mathcal{I}})$  can have. For any individuals  $a$  and  $b$ , the belief in  $a$  having a property  $R$  (i.e.,  $\exists R(a)$ ), is not less than belief in  $a$  having the property  $R$  with  $b$  (i.e.,  $R(a, b)$ ), and disbelief in  $\exists R(a)$  is not more than disbelief in  $R(a, b)$ .

An ontology provides us with domain constraints in the form of TBox axioms. For instance, the axiom  $B_1 \sqsubseteq B_2$  means that every instance of class  $B_1$  is also an instance of class  $B_2$ . This trivially implies  $\neg B_2 \sqsubseteq \neg B_1$ , i.e., an individual that is not an instance of  $B_2$  cannot be an instance of  $B_1$ . Therefore, given an individual  $a$ , the axiom  $B_1 \sqsubseteq B_2$  implies that our belief in  $B_2(a)$  cannot be less than our belief in  $B_1(a)$  and our disbelief in  $B_2(a)$  cannot be more than our disbelief in  $B_1(a)$ . That is,  $b(B_1^{\mathcal{I}}(a^{\mathcal{I}})) \leq b(B_2^{\mathcal{I}}(a^{\mathcal{I}}))$  and  $d(B_2^{\mathcal{I}}(a^{\mathcal{I}})) \leq d(B_1^{\mathcal{I}}(a^{\mathcal{I}}))$  must hold. Similar constraints also exist in Table 3 for  $B_1 \sqsubseteq \neg B_2$ .

**DEFINITION 2.** An *SDL-Lite* knowledge base  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$  is consistent if and only if it has a model. A model of  $\mathcal{K}$  is an interpretation of  $\mathcal{K}$  that satisfies the constraints in Table 3.  $\blacksquare$

If  $\mathcal{K}$  is consistent, it can have many models, but one of them is the most general model with respect to the partial ordering on opinions by Definition 1. In the next section, we describe how to detect consistency, and how to compute the most general model of a consistent *SDL-Lite* knowledge base.

### 3.2 Reasoning with *SDL-Lite*

The aims of this section are twofold. First, we detail the reasoning mechanisms necessary to infer opinions about the world given an *SDL-Lite* ontology from the information available to an agent. Second, we define the conditions under which an *SDL-Lite* knowledge base becomes inconsistent, and prove that our reasoning mechanism ensures the maintenance of a consistent knowledge base under certain conditions.

The TBox  $\mathcal{T}$  of an *SDL-Lite* knowledge base contains: (i) positive inclusions (PIs) in the form  $B_1 \sqsubseteq B_2$ , where  $B_1$  and  $B_2$  are

**Table 4: *SDL-Lite* knowledge base derived from Table 2**

Normalised and Extended TBox	Computed Interpretations
$t_1 : BombedRoad \sqsubseteq SabotagedRoad$	$Blocked(R):(0.71, 0.09, 0.2)$
$t_2 : BombedRoad \sqsubseteq \neg Safe$	$Safe(R):(0.63, 0.2, 0.17)$
$t_3 : BombedRoad \sqsubseteq Blocked$	$BombedRoad(R):(0.2, 0.63, 0.17)$
$t_4 : SabotagedRoad \sqsubseteq \neg Safe$	$SabotagedRoad(R):(0.2, 0.63, 0.17)$
$t_5 : SabotagedRoad \sqsubseteq Blocked$	$\exists roadBombedBy(R):(0, 0.63, 0.37)$
$t_6 : \exists roadBombedBy \sqsubseteq BombedRoad$	$\exists roadBombedBy^-(G_1):(0, 0, 1)$
$t_7 : \exists roadBombedBy \sqsubseteq SabotagedRoad$	$\exists roadBombedBy^-(G_2):(0, 0, 1)$
$t_8 : \exists roadBombedBy \sqsubseteq Blocked$	$roadBombedBy(R, G_1):(0, 0.63, 0.37)$
$t_9 : \exists roadBombedBy \sqsubseteq \neg Safe$	$roadBombedBy(R, G_2):(0, 0.63, 0.37)$
$t_{10} : Safe \sqsubseteq \neg SabotagedRoad$	$roadBombedBy^-(G_1, R):(0, 0.63, 0.37)$
$t_{11} : Safe \sqsubseteq \neg BombedRoad$	$roadBombedBy^-(G_2, R):(0, 0.63, 0.37)$
$t_{12} : Safe \sqsubseteq \neg \exists roadBombedBy$	

basic concepts (being either an atomic or an existential concept); and (ii) negative inclusions (NIs) of the form  $B_1 \sqsubseteq \neg B_2$ . In order to reason over an *SDL-Lite* knowledge base  $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ , we compute the extended knowledge base  $\mathcal{K}^* = (\mathcal{T}^*, \mathcal{A}^*)$ .

The *extended TBox*  $\mathcal{T}^*$  is initialised as equivalent to  $\mathcal{T}$  and expanded by computing all (nontrivial) NIs and PIs between basic concepts implied by  $\mathcal{T}^*$ ; i.e.  $\mathcal{T}^*$  is closed with respect to the rules:

- if  $B_1 \sqsubseteq B_2$  occurs in  $\mathcal{T}^*$  and either  $B_2 \sqsubseteq \neg B_3$  or  $B_3 \sqsubseteq \neg B_2$  occurs in  $\mathcal{T}^*$ , add  $B_1 \sqsubseteq \neg B_3$  to  $\mathcal{T}^*$
  - if  $B_1 \sqsubseteq B_2$  and  $B_2 \sqsubseteq B_3$  occur in  $\mathcal{T}^*$ , add  $B_1 \sqsubseteq B_3$  to  $\mathcal{T}^*$
- We, therefore, extend  $\mathcal{T}$  in  $O(n^2)$  steps in the size of  $\mathcal{T}$  and form  $\mathcal{T}^*$  that contains all possible inclusions derived from  $\mathcal{T}$ .

The *extended ABox*  $\mathcal{A}^*$  is initialised as equivalent to  $\mathcal{A}$  and expanded in  $O(n)$  steps in the size of  $\mathcal{A}$  using the rules:

- if  $R(a, b):w$  occurs in  $\mathcal{A}^*$ , add  $R^-(b, a):w$
- for each individual  $a$  where  $R(a, b_1):w_1, \dots, R(a, b_n):w_n \in \mathcal{A}^*$ , add a new assertion  $\exists R(a):w$  where  $d(w) = 0$  and  $b(w) = \max(\{b(w_1), \dots, b(w_n)\})$ .

We may now specify how *interpretations* of classes and roles may be computed based on the computed  $\mathcal{T}^*$  and  $\mathcal{A}^*$ . Given a class  $B_n$ ,  $\mathcal{T}^*$  may contain the following axioms that constrain the interpretation of an ABox assertion  $B_n(a)$ :

- $B_0 \sqsubseteq B_n, B_1 \sqsubseteq B_n, \dots, B_l \sqsubseteq B_n$
- $B_{l+1} \sqsubseteq \neg B_n, B_{l+2} \sqsubseteq \neg B_n, \dots, B_i \sqsubseteq \neg B_n$
- $B_n \sqsubseteq B_{i+1}, B_n \sqsubseteq B_{i+2}, \dots, B_n \sqsubseteq B_j$
- $B_n \sqsubseteq \neg B_{j+1}, B_n \sqsubseteq \neg B_{j+2}, \dots, B_n \sqsubseteq \neg B_k$

Given these kinds of axioms, Our aim here is to define appropriate inferences regarding the opinion of some assertion in  $\mathcal{A}^*$ . Let  $w_n$  refer to the opinion related to the ABox assertion  $B_n(a)$ , i.e.,  $B_n(a):w_n \in \mathcal{A}^*$ . Clearly, if  $B_n(a)$  does not appear in  $\mathcal{A}^*$  — representing the case that we have no evidence regarding  $B_n(a)$  — then  $w_n = (0, 0, 1)$ , representing maximal uncertainty regarding  $B_i(a)$ . We use  $w'_n$  to refer to the interpretation of  $B_n(a)$ , i.e.  $w' = B_n^{\mathcal{I}}(a)$ . The most general opinion  $w''_n$  that satisfy the constraints for  $w'_n$  is then computed as follows:

$$w''_n = (\max(S_b), \max(S_d)) \text{ where}$$

$$S_b = \{b(w_n), b(w_0), b(w_1), \dots, b(w_l)\}$$

$$S_d = \{d(w_n), b(w_{l+1}), b(w_{l+2}), \dots, b(w_i),$$

$$d(w_{i+1}), d(w_{i+2}), \dots, d(w_j),$$

$$b(w_{j+1}), b(w_{j+2}), \dots, b(w_k)\}$$

If  $b(w''_n) + d(w''_n) > 1$ , then there is no opinion satisfying the constraints defined by the semantics for  $B_n(a)$ ; otherwise, we take  $w''_n$  as an interpretation of  $B_n(a)$ .

Let us explain the interpretations for class assertions through the scenario in Section 2.3. Table 4 shows the extended TBox for the scenario. The interpretation of *SabotagedRoad*(R) is constrained by TBox axioms  $t_1, t_4, t_5, t_7$ , and  $t_{10}$  of Table 4. Therefore, given the ABox in Table 2 and based on Equation 1, the interpretation is computed as  $(\max(\{0, 0.2\}), \max(\{0, 0.63, 0.09\})) = (0.2, 0.63)$ .

The interpretation for  $R_n(a, b)$  is constrained by the interpretations of  $\exists R_n(a)$  and  $\exists R_n^-(b)$ . Let  $R_n(a, b):w \in \mathcal{A}$  and  $\exists R_n(a):w_n$

be the ABox assertion added to  $\mathcal{A}^*$  while extending  $\mathcal{A}$ . We know that  $b(w_n) \geq b(w)$  and  $d(w) \geq d(w_n)$ . Belief in  $\exists R_n^{\mathcal{I}}(\mathbf{a})$  does not constrain belief in  $R_n^{\mathcal{I}}(\mathbf{a}, \mathbf{b})$ . However, in a consistent knowledge base, disbelief in  $R_n^{\mathcal{I}}(\mathbf{a}, \mathbf{b})$  cannot be lower than disbelief in  $\exists R_n^{\mathcal{I}}(\mathbf{a})$  or  $\exists R_n^{\mathcal{I}}(\mathbf{b})$ . Hence, disbelief in  $\exists R_n^{\mathcal{I}}(\mathbf{a})$  is constrained by the following TBox axioms in  $\mathcal{T}^*$ :

- $B_1 \sqsubseteq \neg \exists R_n, B_2 \sqsubseteq \neg \exists R_n, \dots, B_{a-1} \sqsubseteq \neg \exists R_n, B_a \sqsubseteq \neg \exists R_n$
- $\exists R_n \sqsubseteq B_{a+1}, \exists R_n \sqsubseteq B_{a+2}, \dots, \exists R_n \sqsubseteq B_{b-1}, \exists R_n \sqsubseteq B_b$
- $\exists R_n \sqsubseteq \neg B_{b+1}, \exists R_n \sqsubseteq \neg B_{b+2}, \dots, \exists R_n \sqsubseteq \neg B_{c-1}, \exists R_n \sqsubseteq \neg B_c$

Disbelief in  $\exists R_n^{\mathcal{I}}(\mathbf{b})$  is constrained by the axioms in  $\mathcal{T}^*$ :

- $B_{c+1} \sqsubseteq \neg \exists R_n^-, B_{c+2} \sqsubseteq \neg \exists R_n^-, \dots, B_{d-1} \sqsubseteq \neg \exists R_n^-, B_d \sqsubseteq \neg \exists R_n^-$
- $\exists R_n^- \sqsubseteq B_{d+1}, \exists R_n^- \sqsubseteq B_{d+2}, \dots, \exists R_n^- \sqsubseteq B_{e-1}, \exists R_n^- \sqsubseteq B_e$
- $\exists R_n^- \sqsubseteq \neg B_{e+1}, \exists R_n^- \sqsubseteq \neg B_{e+2}, \dots, \exists R_n^- \sqsubseteq \neg B_{f-1}, \exists R_n^- \sqsubseteq \neg B_f$

Let  $w_i$  refer to the opinion related to the ABox assertion for  $B_i(\mathbf{a})$  if  $1 \leq i \leq c$  and for  $B_i(\mathbf{b})$  if  $c+1 \leq i \leq f$ . The most general opinion  $w''$  that satisfies the constraints for  $R_n^{\mathcal{I}}(\mathbf{a}, \mathbf{b})$  is computed as follows:

$$w'' = (b(w), \max(S_d)) \quad \text{where} \\ S_d = \{d(w), b(w_1), \dots, b(w_a), b(w_{c+1}), \dots, b(w_d), \\ d(w_{a+1}), d(w_{a+2}), \dots, d(w_b), \\ d(w_{d+1}), d(w_{d+2}), \dots, d(w_e), \\ b(w_{b+1}), b(w_{b+2}), \dots, b(w_c), \\ b(w_{e+1}), b(w_{e+2}), \dots, b(w_f)\} \quad (2)$$

If  $b(w'') + d(w'') > 1$ , there is no opinion satisfying the constraints defined by the semantics for  $R_n(\mathbf{a}, \mathbf{b})$ ; otherwise, we take  $w''$  as interpretation of  $R_n(\mathbf{a}, \mathbf{b})$ .

Let us now compute the interpretation of  $\text{roadBombedBy}(\mathbf{R}, \mathbf{G}_1)$  in our example scenario. It is constrained by the interpretations of two other assertions  $\exists \text{roadBombedBy}(\mathbf{R})$  and  $\exists \text{roadBombedBy}^-(\mathbf{G}_1)$ . The disbelief in the interpretation of  $\exists \text{roadBombedBy}(\mathbf{R})$  is constrained by TBox axioms  $t_6, t_7, t_8, t_9$ , and  $t_{12}$  of Table 4. We have opinion assertions only for the assertions  $\text{Blocked}(\mathbf{R})$ ,  $\text{Safe}(\mathbf{R})$ , and  $\text{BombedRoad}(\mathbf{R})$  in the extended ABox derived from Table 2. Therefore, based on Equation 2, we compute the interpretation as  $(0, \max(\{0, 0.63, 0.3, 0.09\})) = (0, 0.63)$ . Table 4 shows the interpretations computed for the scenario using Equations 1 and 2.

The computational complexity of these calculations is  $O(n)$  in the size of  $\mathcal{T}^*$  and  $\mathcal{A}^*$ . Now, we introduce Theorem 1, which defines the conditions necessary and sufficient for inconsistency.

**THEOREM 1.** *An extended SDL-Lite KB  $\mathcal{K}^* = (\mathcal{T}^*, \mathcal{A}^*)$  with a coherent  $\mathcal{T}^*$  is inconsistent with respect to the semantics in Table 3 if and only if one of the following conditions hold:*

1. *Given  $B_m(\mathbf{a}):w_m, B_n(\mathbf{a}):w_n \in \mathcal{A}^*$ , and  $B_m \sqsubseteq B_n \in \mathcal{T}^*$ , we have  $b(w_m) + d(w_n) > 1$*
2. *Given  $B_m(\mathbf{a}):w_m, B_n(\mathbf{a}):w_n \in \mathcal{A}^*$ , and  $B_m \sqsubseteq \neg B_n \in \mathcal{T}^*$ , we have  $b(w_m) + b(w_n) > 1$*

**Proof:** *The inconsistency arises if and only if at least one class or role does not have a valid interpretation satisfying the semantics in Table 3. Let us first analyse the inconsistencies due to the interpretations of classes. The most general interpretation for  $B_n(\mathbf{a})$  is computed as in Equation 1 and referred to as  $w_n''$ . Let  $w_n$  be the opinion for  $B_n(\mathbf{a})$  in  $\mathcal{A}^*$ . If  $b(w_n'') + d(w_n'') > 1$ , there is no opinion satisfying the constraints defined by the semantics for  $B_n(\mathbf{a})$  and  $\mathcal{K}^*$  is inconsistent. To have  $b(w_n'') + d(w_n'') > 1$ , one of the following conditions must hold based on Equation 1:*

- $b(w_n'') = b(w_n)$

- $d(w_n'') \in \{b(w_{i+1}), b(w_{i+2}), \dots, b(w_i)\}$ : *This implies that there exists a TBox axiom  $B_m \sqsubseteq \neg B_n$  with ABox assertions  $B_m(\mathbf{a}):w_m$  and  $B_n(\mathbf{a}):w_n$  such that  $b(w_n) + b(w_m) > 1$ .*
- $d(w_n'') \in \{d(w_{i+1}), d(w_{i+2}), \dots, d(w_i)\}$ : *This implies that there exists a TBox axiom  $B_n \sqsubseteq B_m$  with ABox assertions  $B_m(\mathbf{a}):w_m$  and  $B_n(\mathbf{a}):w_n$  such that  $b(w_n'') + d(w_m) > 1$ .*
- $d(w_n'') \in \{b(w_{j+1}), b(w_{j+2}), \dots, b(w_k)\}$ : *This implies that there exists a TBox axiom  $B_n \sqsubseteq \neg B_m$  with ABox assertions  $B_m(\mathbf{a}):w_m$  and  $B_n(\mathbf{a}):w_n$  such that  $b(w_n) + b(w_m) > 1$ .*
- $b(w_n'') \in \{b(w_0), b(w_1), \dots, b(w_l)\}$ 
  - $d(w_n'') = d(w_n)$ : *This implies that there exists a TBox axiom  $B_m \sqsubseteq B_n$  with ABox assertions  $B_m(\mathbf{a}):w_m$  and  $B_n(\mathbf{a}):w_n$  such that  $b(w_m) + d(w_n) > 1$ .*
  - $d(w_n'') \in \{b(w_{i+1}), b(w_{i+2}), \dots, b(w_i)\}$ : *This implies TBox axioms  $B_x \sqsubseteq B_n$  and  $B_y \sqsubseteq \neg B_n$  with ABox assertions  $B_x(\mathbf{a}):w_x$  and  $B_y(\mathbf{a}):w_y$ . These TBox axioms imply that the extended  $\mathcal{T}^*$  contains  $B_x \sqsubseteq \neg B_y$  and  $b(w_x) + b(w_y) > 1$ .*
  - $d(w_n'') \in \{d(w_{i+1}), d(w_{i+2}), \dots, d(w_j)\}$ : *This implies TBox axioms  $B_x \sqsubseteq B_n$  and  $B_n \sqsubseteq B_y$  with ABox assertions  $B_x(\mathbf{a}):w_x$  and  $B_y(\mathbf{a}):w_y$ . These TBox axioms imply that  $\mathcal{T}^*$  contains  $B_x \sqsubseteq B_y$  and  $b(w_x) + d(w_y) > 1$ .*
  - $d(w_n'') \in \{b(w_{j+1}), b(w_{j+2}), \dots, b(w_k)\}$ : *This implies TBox axioms  $B_x \sqsubseteq B_n$  and  $B_n \sqsubseteq \neg B_y$  with ABox assertions  $B_x(\mathbf{a}):w_x$  and  $B_y(\mathbf{a}):w_y$ . These TBox axioms imply that  $\mathcal{T}^*$  contains  $B_x \sqsubseteq \neg B_y$  and  $b(w_x) + b(w_y) > 1$ .*

Now, we look into the interpretations for role axioms to test consistency. The most general interpretation for  $R_n(\mathbf{a}, \mathbf{b})$  is computed as in Equation 2 and referred to as  $w''$ . Let  $w$  be the opinion for  $R_n(\mathbf{a}, \mathbf{b})$  in  $\mathcal{A}^*$ . If  $b(w'') + d(w'') > 1$ , there is no opinion satisfying the constraints defined by the semantics for  $R_n(\mathbf{a}, \mathbf{b})$  and  $\mathcal{K}^*$  is inconsistent. To have  $b(w'') + d(w'') > 1$ , one of the following conditions must hold based on Equation 2:

- $B_m \sqsubseteq \neg \exists R_n$  or  $B_m \sqsubseteq \neg \exists R_n^-$  and  $b(w_n) + b(w_m) > 1$  since  $b(w) + b(w_m) > 1$
- $\exists R_n \sqsubseteq B_m$  or  $\exists R_n^- \sqsubseteq B_m$  and  $b(w_n) + d(w_m) > 1$  since  $b(w) + d(w_m) > 1$
- $\exists R_n \sqsubseteq \neg B_m$  or  $\exists R_n^- \sqsubseteq \neg B_m$  and  $b(w_n) + b(w_m) > 1$  since  $b(w) + b(w_m) > 1$

As shown, in the case of inconsistency, one of the conditions defined in Theorem 1 must hold. Furthermore, if none of these conditions holds in  $\mathcal{K}^*$ , we guarantee that  $\mathcal{K}^*$  is consistent. ■

In an inconsistent extended SDL-Lite knowledge base  $\mathcal{K}^* = (\mathcal{T}^*, \mathcal{A}^*)$ , the inconsistencies exist only because of conflicting opinions. Two opinions  $w_m$  and  $w_n$ , which are about  $B_m(\mathbf{a})$  and  $B_n(\mathbf{a})$  respectively, are in conflict if they satisfy one of the conditions in Theorem 1. We label the portion of  $w_m$  which conflicts with  $w_n$  as  $c_{mn}$ , and refer to it as the *conflicting portion*. If the conflict is due to the axiom  $B_m \sqsubseteq B_n \in \mathcal{T}^*$ , then the conflict arises because  $b(w_m) + d(w_n) > 1$ ; hence  $c_{mn} = b(w_m)$  and  $c_{nm} = d(w_n)$ . On the other hand, if the conflict is due to the axiom  $B_m \sqsubseteq \neg B_n \in \mathcal{T}^*$ , we have conflict because  $b(w_m) + b(w_n) > 1$ ; hence  $c_{mn} = b(w_m)$  and  $c_{nm} = b(w_n)$ . If all conflicts in  $\mathcal{K}^*$  are resolved, then the knowledge base becomes consistent.

In the rest of the paper, we assume that the opinion about a specific ABox assertions is provided by a single source. When there is more than one source for an assertion, only one of them is chosen (e.g. based on their trustworthiness). This will be relaxed in future.

**Table 5: Extended ABox for case II**

$a_1$	: $roadBombedBy(R, G_1): (0.67, 0.083, 0.247)$
$a_2$	: $roadBombedBy^-(G_1, R): (0.67, 0.083, 0.247)$
$a_3$	: $\exists roadBombedBy(R): (0.67, 0.0, 0.33)$
$a_4$	: $\exists roadBombedBy^-(G_1): (0.67, 0.0, 0.33)$
$a_5$	: $Blocked(R): (0.71, 0.09, 0.2)$
$a_6$	: $Safe(R): (0.63, 0.066, 0.304)$
$a_7$	: $BombedRoad(R): (0.2, 0.3, 0.5)$

Having described *SDL-Lite* we now examine a novel application of the system, describing how evidence from multiple sources can be reasoned about based on the trust placed in these sources.

## 4. TRUST-BASED EVIDENCE ANALYSIS

Here we get to the crux of the problem being addressed in this paper: how can we draw reliable conclusions regarding the state of the world, given evidence acquired from disparate sources (agents), about whom we have variable trust? We refer to this process as trust-based evidence analysis. Our aim is not to offer a new mechanism for assessing the trustworthiness of information sources; in fact, we exploit a widely-studied model [10] for this purpose based on Beta distributions as described in Section 2.2. The novelty of this work lies in the use of such models to guide evidence analysis.

### 4.1 Handling Inconsistencies

*SDL-Lite* presented in the previous section provides a tractable means to capture and interpret evidence acquired from other agents. The fact that we have evidence from multiple agents, however, means that there are likely to be inconsistencies in the evidence received. Thus, given evidence (i.e., opinions) from various sources, our knowledge-base may not be consistent. This is despite the use of *discounting* through DST. Discounting provides us with a “best-guess” of the reliability of agents based on an aggregation of our prior experiences with, and other knowledge of them as evidence sources. As with any computational model of trust, the trust assessments that drive discounting are vulnerable to: lack of evidence about other agents and the effects of whitewashing [2]; a conflation of the probability of malicious behaviour and lack competence/expertise in the evidence-provider; strategic liars; and collusion among evidence-providers. In our running example, for instance, local police and civilian sources have relatively low trustworthiness, not because of any perceived malicious intent but due to a belief that they lack experience in providing precise information. With more evidence, trustworthiness of information sources may be modelled more accurately, but our challenge is to support the analysis of evidence given the status quo.

To illustrate this challenge, consider an adaptation of our example (case II) in which additional evidence is received from a third source, agent C, about R: C reports that R was bombed by  $G_1$  with opinion  $(0.8, 0.1, 0.1)$ . With this additional report, our ABox contains  $roadBombedBy(R, G_1): (0.67, 0.083, 0.247)$  after discounting the opinion with C’s trustworthiness 0.83 listed in Table 1 ( $a_1$  in Table 5 where the resulting extended ABox is presented). The extended ABox will now have a conflict between  $a_1$  and  $a_6$ , because  $0.67 + 0.63 > 1$  and the extended TBox contains  $\exists roadBombedBy \sqsubseteq \neg Safe$ . Let  $w_1 = (0.63, 0.066, 0.304)$  and  $w_2 = (0.67, 0.0, 0.33)$ . The conflicting portions of  $w_1$  and  $w_2$  are  $c_{12} = 0.63$  and  $c_{21} = 0.67$ . Let us refer to the trustworthiness of the sources of  $w_1$  and  $w_2$  as  $t_1$  and  $t_2$  respectively. In our example, from Table 1,  $t_1 = 0.83$  and  $t_2 = 0.786$ . In order for us to transform our inconsistent knowledge-base into a *consistent* knowledge-base, from which we can draw valid conclusions given our semantics, we need to determine additional discounting factors  $x_1$  and  $x_2$  for opinions  $w_1$  and  $w_2$  such that  $0 \leq c_{12}.x_1 + c_{21}.x_2 \leq 1$ .

In this paper, we specify this problem as that of finding *additional* discounting factors for the belief-mass distributions of pieces

of evidence to make our knowledge-base consistent. In general, our conflict resolution problem is a tuple  $\langle \mathcal{C}, \mathcal{X} \rangle$  where  $\mathcal{C}$  is the set of conflicting portions that appear in the extended knowledge base  $\mathcal{K}^*$ , and  $\mathcal{X}$  is a set of additional discounting factors corresponding to  $\mathcal{C}$ . We require that, in  $\langle \mathcal{C}, \mathcal{X} \rangle$ ,  $\forall c_{ij} \in \mathcal{C}$ ,  $\exists c_{ji} \in \mathcal{C}$  and  $\exists x_i, x_j \in \mathcal{X}$ . Then, a solution to this problem is an assignment of values to each  $x_i \in \mathcal{X}$  such that

$$\forall c_{ij}, c_{ji} \in \mathcal{C}, \forall x_i, x_j \in \mathcal{X} \quad 0 \leq c_{ij}.x_i + c_{ji}.x_j \leq 1$$

There are many heuristic approaches to solving this problem, among them being to consider only consistent knowledge to draw conclusions from the evidence received; i.e.  $\forall x_i \in \mathcal{X}, x_i = 0$ . This, however, could lead to a significant loss of evidence. Here, we explore a number of increasingly refined approaches that guarantee the generation of a consistent knowledge-base: *trust-based deleting*, *trust-based discounting* and *evidence-based discounting*.

### 4.2 Trust-based deleting

If two opinions  $w_1$  and  $w_2$  are in conflict, the opinion from the less trustworthy source is deleted, and if both sources are equally trustworthy both opinions are deleted. Thus, if the trust we have in the source of opinion  $w_1$  is greater than that of the source of  $w_2$  ( $t_1 > t_2$ ) then  $x_2 = 0$  and  $x_1 = 1$ , and in the event that  $t_1 = t_2$  we assign  $x_1 = x_2 = 0$ . In our example, the local police sources P are slightly less trustworthy than the local civilian sources C. Hence, the opinion about  $Safe(R)$  is changed to  $(0, 0, 1)$  and the conflict is resolved. This approach, however, neglects the amount of evidence used to calculate trust, and it does not consider the difference between trust values ( $t_C = 0.83 \approx t_P = 0.786$ ).

### 4.3 Trust-based discounting

If two opinions  $w_1$  and  $w_2$  are in conflict, they are discounted in proportion to the trustworthiness of their sources. That is, the additional discounting factor for  $w_1$  and  $w_2$  is computed using  $t_1/(c_{12}t_1 + c_{21}t_2)$  and  $t_2/(c_{12}t_1 + c_{21}t_2)$ , respectively, where  $t_1$  and  $t_2$  are the trustworthiness of the sources of the opinions. In our example, an additional discount factor of  $roadBombedBy(R, G_1)$  is 0.79 and that of  $Safe(R)$  is 0.75, since the trustworthiness of C and P are 0.83 and 0.786, respectively. Therefore, to resolve the conflict, the original opinion of C about  $roadBombedBy(R, G_1)$  is discounted by  $0.83 \times 0.79 = 0.65$  and that of P about  $Safe(R)$  is discounted by  $0.786 \times 0.75 = 0.59$ . However, this approach neglects the amount of evidence used to calculate trust in sources.

### 4.4 Evidence-based discounting

Within the evidence analysis domain, the information that we have to work with relates to past experiences with a specific agent (i.e., information source)  $\varrho_k$  where information received has proven reliable or unreliable according to some criteria (as would be captured in any trust assessment model). In other words, the amount of positive evidence we have for agent  $\varrho_k$ , namely  $r_k$ , and the amount of negative evidence for that agent, namely  $s_k$ . From this evidence, we calculate trustworthiness of  $\varrho_k$ , denoted as  $t_k$  described in Section 2.2. When we receive opinion  $w_i^k$  from  $\varrho_k$ , we discount it by  $t_k$  and add the resulting opinion  $w_i$  to our knowledge base. However, as explained before, additional discounting by factor  $x_i$  is required when  $w_i$  is in conflict with another opinion in the knowledge base. Discounting  $w_i$  by  $x_i$  implies discounting the original opinion  $w_i^k$  by  $t_k.x_i$ . This corresponds to revising the trustworthiness of  $w_i^k$  as  $t_k.x_i$  by speculating about the trustworthiness of  $\varrho_k$  regarding this single opinion. That is, even though the trustworthiness of  $\varrho_k$  is  $t_k$  based on the existing evidence  $(r_k, s_k)$ , it becomes  $t_k.x_i$  for this specific opinion  $w_i^k$ ; so,  $t_k.x_i$  effectively becomes the trust in  $w_i^k$ .

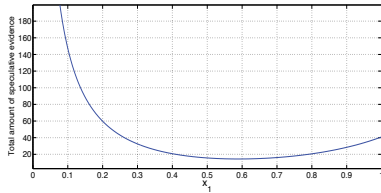


Figure 1: Speculative evidence required for case II ( $\kappa = 1$ ).

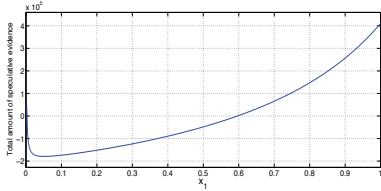


Figure 2: Speculative evidence required for case III ( $\kappa = 1$ ).

Here, we create a metric to measure how much we speculate about the trustworthiness of  $\rho_k$  regarding  $w_i^k$ .

First, to decrease trust from  $t_k$  to  $t_k \cdot x_i$ , we need additional negative evidence, which is called *speculative evidence* and denoted as  $\rho_i$ . Our intuition is that it is less likely for a trustworthy agent to present additional negative *speculative evidence* than it is for an untrustworthy agent, and thus the receipt of such evidence should be tempered by  $(\bar{t}_k)^\kappa$ . Here,  $\bar{t}_k$  represents the *distrust* we have in agent  $\rho_k$ ; i.e. the likelihood that we will receive additional negative evidence given our experiences with the source. The calibration constant  $\kappa \geq 0$  enables us to vary the influence that prior experience has on our prediction that an individual will present negative evidence in the future. If  $\kappa = 0$ , for example, we assume that all sources are equally likely to provide negative evidence. Now, using the Beta distribution formula for trust, we obtain:

$$\begin{aligned} t_k \cdot x_i &= \frac{r_k + 1}{r_k + s_k + 2} \cdot x_i = \frac{r_k + 1}{s_k + r_k + 2 + \rho_i \cdot (\bar{t}_k)^\kappa} \\ &= \frac{r_k + 1}{s_k + r_k + 2 + \rho_i \cdot (\frac{s_k + 1}{r_k + s_k + 2})^\kappa} \end{aligned}$$

Rearranging this for  $\rho_i$  yields:

$$\rho_i = \frac{\nu_i(1 - x_i)}{x_i} \quad \text{where} \quad \nu_i = \frac{(r_k + s_k + 2)^{\kappa+1}}{(s_k + 1)^\kappa} \quad (3)$$

To illustrate this, let us return to case II above in which agent C reports that road R is bombed by  $G_1$ . Using Equation 3, we can compute the total amount of *speculative evidence* necessary to discount  $w_1$  and  $w_2$  by  $x_1$  and  $x_2$ , respectively. If we assume that  $c_{12} \cdot x_1 + c_{21} \cdot x_2 = 1$ , we have  $x_2 = (1 - c_{12} \cdot x_1) / c_{21}$ . Then, the total amount of *speculative evidence* (i.e.  $\rho_1 + \rho_2$ ) can be formulated as a function of single variable  $x_1$  by Equation 4, which is plotted in Figure 1. This function has a minimum at  $x_1 = 0.5892$  in the interval  $[0, 1]$  and the corresponding  $x_2$  is 0.9607. That is, for a consistent knowledge base, trust in C's opinion about *roadBombedBy*(R,  $G_1$ ) should be reduced to 0.489 from 0.83, but the trust in the opinion of P about *Safe*(R) is reduced only slightly to 0.755 from 0.786. This reflects the relative level of positive and negative evidence we have from prior experience from both parties, and results in a consistent knowledge-based from which we can draw conclusions.

$$f(x_1) = \frac{\nu_1(1 - x_1)}{x_1} + \frac{\nu_2(1 - \frac{(1 - c_{12} \cdot x_1)}{c_{21}})}{\frac{(1 - c_{12} \cdot x_1)}{c_{21}}} \quad (4)$$

Table 6: Extended ABox for case III

<i>roadBombedBy</i> (R, $G_1$ ):(0.67, 0.083, 0.247)
<i>roadBombedBy</i> <sup>-</sup> ( $G_1$ , R):(0.67, 0.083, 0.247)
$\exists$ <i>roadBombedBy</i> (R):(0.67, 0.0, 0.33)
$\exists$ <i>roadBombedBy</i> <sup>-</sup> ( $G_1$ ):(0.67, 0.0, 0.33)
<i>Blocked</i> (R):(0.71, 0.09, 0.2)
<i>Safe</i> (R):(0.63, 0.066, 0.304)
<i>BombedRoad</i> (R):(0.2, 0.6, 0.2)
<i>SabotagedRoad</i> (R):(0.801, 0, 0.199)

Table 7: After extra discounting for case III ( $\kappa = 1$ )

Extended ABox	Computed Interpretations
<i>Blocked</i> (R):(0.71, 0.09)	<i>Blocked</i> (R):(0.71, 0.09)
<i>Safe</i> (R):(0.63, 0.066)	<i>Safe</i> (R):(0.63, 0.066)
<i>roadBombedBy</i> (R, $G_1$ ):(0.0342, 0.0043)	<i>BombedRoad</i> (R):(0.2, 0.63)
<i>roadBombedBy</i> <sup>-</sup> ( $G_1$ , R):(0.0342, 0.0043)	<i>SabotagedRoad</i> (R):(0.0342, 0.63)
$\exists$ <i>roadBombedBy</i> (R):(0.0342, 0)	$\exists$ <i>roadBombedBy</i> (R):(0.0342, 0.63)
$\exists$ <i>roadBombedBy</i> <sup>-</sup> ( $G_1$ ):(0.0342, 0)	$\exists$ <i>roadBombedBy</i> <sup>-</sup> ( $G_1$ ):(0.0342, 0)
<i>Blocked</i> (R):(0.71, 0.09)	$\exists$ <i>roadBombedBy</i> <sup>-</sup> ( $G_2$ ):(0, 0, 1)
<i>Safe</i> (R):(0.63, 0.066)	<i>roadBombedBy</i> (R, $G_1$ ):(0.0342, 0.63)
<i>BombedRoad</i> (R):(0.2, 0.6)	<i>roadBombedBy</i> (R, $G_2$ ):(0, 0.63)
<i>SabotagedRoad</i> (R):(0.0304, 0)	<i>roadBombedBy</i> <sup>-</sup> ( $G_1$ , R):(0.0342, 0.63)
	<i>roadBombedBy</i> <sup>-</sup> ( $G_2$ , R):(0, 0.63)

Until now, we considered only one conflict between two opinions. When we have multiple conflicts, they may interact in such a way that resolving one may also affect the resolution of another. To illustrate this, consider two new intelligence reports (case III):

- A reports a bomb explosion on R with opinion (0.2, 0.6, 0.2).
- $M_2$  informs  $M_1$  that R is sabotaged with opinion (0.9, 0, 0.1).

The resulting ABox is shown in Table 6 and implies three relevant conflicts:  $0.67 + 0.63 > 1$ ,  $0.67 + 0.6 > 1$ , and  $0.63 + 0.801 > 1$ . Let us refer to (0.2, 0.6, 0.2) and (0.801, 0, 0.199) as  $w_3$  and  $w_4$ , respectively. We refer to the conflicting portions as  $c_{31} = 0.6$  and  $c_{42} = 0.801$ . We also use  $x_3$  and  $x_4$  to refer to the additional discounting necessary for  $w_3$  and  $w_4$ , respectively, to resolve the conflicts. The overall amount of *speculative evidence* necessary to resolve all of these relevant conflicts is computed as in Equation 5.

$$\begin{aligned} f(x_1) &= \frac{\nu_1(1 - x_1)}{x_1} + \frac{\nu_2(1 - x_2)}{x_2} + \frac{\nu_3(1 - x_3)}{x_3} + \frac{\nu_4(1 - x_4)}{x_4} \\ \text{such that} \quad &0 \leq c_{12} \cdot x_1 + c_{21} \cdot x_2 \leq 1 \text{ and} \\ &0 \leq c_{13} \cdot x_1 + c_{31} \cdot x_3 \leq 1 \text{ and} \\ &0 \leq c_{24} \cdot x_2 + c_{42} \cdot x_4 \leq 1 \end{aligned} \quad (5)$$

Since these conflicts are relevant, we can write  $x_2$ ,  $x_3$  and  $x_4$  in terms of  $x_1$  if we set  $c_{12}x_1 + c_{21}x_2 = 1$ ,  $c_{13}x_1 + c_{31}x_3 = 1$ , and  $c_{24}x_2 + c_{42}x_4 = 1$ . The resulting function is shown in Figure 2 and has a minimum at  $x_1 = 0.0514$  in the interval  $[0, 1]$ . The other discounting factors are computed as  $x_2 = 1$ ,  $x_3 = 1$ , and  $x_4 = 0.043$  in the same interval. That is, trust in the opinion of C about *roadBombedBy*(R,  $G_1$ ) is reduced to 0.0427 and trust in the opinion of M about *SabotagedRoad*(R) is reduced to 0.0338. The ABox and the computed interpretations after extra discounting is shown in Table 7.

We generalise this approach for any number of conflicts with arbitrary relations. Assume we have a set of conflicting opinions  $\{\langle w_i, w_j \rangle, \dots, \langle w_m, w_n \rangle\}$  and, derived from trust evidence about agents, coefficients  $\{\nu_i, \nu_j, \dots, \nu_m, \nu_n\}$ . To determine the optimum discounting factors  $\{x_i, x_j, \dots, x_m, x_n\}$  for these opinions, we construct the following optimisation problem with a multivariate non-linear objective function and linear constraints.

$$\begin{aligned} \arg \min_{\vec{x}} f(\vec{x}) \quad \text{where} \\ f(\langle x_i, x_j, \dots, x_m, x_n \rangle) &= \frac{\nu_i(1 - x_i)}{x_i} + \frac{\nu_j(1 - x_j)}{x_j} + \dots \\ &\quad \frac{\nu_m(1 - x_m)}{x_m} + \frac{\nu_n(1 - x_n)}{x_n} \\ \text{such that} \quad &0 \leq x_i \leq 1, 0 \leq x_j \leq 1, \dots \\ \text{and} \quad &0 \leq c_{ij}x_i + c_{ji}x_j \leq 1, \dots \end{aligned} \quad (6)$$

Existing constrained non-linear programming methods can be used to solve this problem in order to estimate the best discounting factors. There are various techniques that may be used including *Interior-Point* and *Active-Set* algorithms. In this work, we use *Interior-Point* approximation. Details of these methods are out of the scope of this paper and can be found elsewhere [14].

In this section we have formalised the problem of computing additional discounting factors for *opinions* received about the world from other agent so that we may formulate a consistent SDL-Lite



knowledge-base from which we can draw reliable conclusions. We have presented a number of approaches to the resolutions of inconsistencies between opinions including an optimisation-based approach, evidence-based discounting. Next, we evaluate these approaches with respect to their robustness in the face of liars.

## 5. EVALUATION

We have evaluated our approach through a set of simulations. In each simulation, we define the domain by randomly generating an SDL-Lite TBox that contains 100 concepts and roles, as well as axioms over those, e.g.,  $B_1 \sqsubseteq B_2$  and  $B_2 \sqsubseteq \neg \exists R_3$ . For each role or concept, there is one information source that provides opinions about its instances, e.g.,  $B_1(a):(0.8, 0, 0.2)$  and  $R_3(a, b):(0.5, 0.1, 0.4)$ . There are 10 information sources in total, each is an expert on 10 concepts and roles, and provides its opinions about those.

In our simulations, we assume there is one information consumer that uses the information from sources to make decisions. Each simulation is composed of 10 iterations. At each iteration  $t$ , the consumer needs to gather information about an individual  $a$ . We generate ground truth about  $a$ , which is composed of one assertion about  $a$  for each concept and role with an associated opinion. Each information source knows the ground truth only about the concepts and roles of their expertise. However, they may not provide the ground truth to the consumer when it is requested. Behaviours of the information sources are determined by their behavioural type, which are summarised as follows.

- **Honest:** Most of the time, this type of sources provide the ground truth about the assertion of their expertise with small Gaussian noise  $N(0, 0.01)$ . With probability  $P_b$ , honest sources behave like malicious ones and provide bogus information.
- **Malicious:** This type of sources aim at misleading the information consumer by providing bogus opinions. More specifically, given  $(b, d, \_)$  is the ground truth about an assertion, a malicious source provides the opinion  $(abs(\epsilon_1), 0.9 + \epsilon_2, \_)$  if  $b \approx d$ ; otherwise it provides the opinion  $(d + \epsilon_1, b + \epsilon_2, \_)$ , where  $\epsilon_1, \epsilon_2 \in [-0.05, 0.05]$ . There are two types of malicious sources, which are defined as follows:
  - i. **Simple liars:** they always provide bogus opinions.
  - ii. **Strategic liars:** they behave like honest sources to build trust and then provides bogus information exploiting the built trust. After providing misleading information to the consumer, they change their identity to avoid negative evidence against them.

After collecting opinions about different assertions from information sources, the information consumer uses its trust in these sources to discount these opinions and uses the proposed reasoning mechanisms for SDL-Lite to compute interpretations. Ideally, these interpretations should be close to the ground truth if all sources are accurate and their trustworthiness is modelled correctly. If there are some malicious sources, there may be conflicts in the collected information. In the case of conflicts, the consumer resolve the conflicts using *Naive Deleting* (NDL), *Trust-based Deleting* (TDL), *Trust-based Discounting* (TDC), or *Evidence-based Discounting* (EDC) with  $\kappa = 1$ . In NDL, all conflicting opinions are deleted from the knowledge base to resolve the conflicts. The consumer computes the interpretations for concept and role assertions related to  $a$ , after resolving the conflicts if any. Then, we measure the performance as the *mean absolute error* in the computed interpretations. Let  $(b, d, u)$  be the ground truth and  $(b', d', u')$  be the computed interpretation for assertion  $B(a)$ , then the *absolute error* in the interpretation is computed as  $err_{B(a)} = abs(\delta_b) + abs(\delta_d)$ , where  $\delta_b = b' - b$  and  $\delta_d = d' - d$ . For instance, if the ground

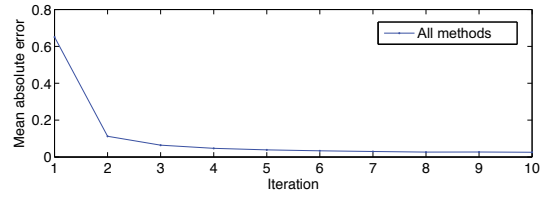


Figure 3: Simple liars ( $R_{liar} = 0.5$  and  $P_b = 0$ )

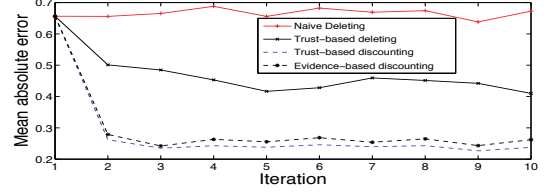


Figure 4: Simple liars ( $R_{liar} = 0.5$  and  $P_b = 0.1$ )

truth about  $B(a)$  is  $(0.9, 0.05, 0.05)$ , but the computed interpretation is  $(0.05, 0.9, 0.05)$ , then the error would be 1.7.

At the end of each iteration, the consumer learns the ground truth and updates the trustworthiness of the information sources with new evidence  $(r^t, s^t)$  computed as in Equation 7, which is based on the intuition that the information is still useful if it has a small amount of noise or is slightly discounted.

$$(r^t, s^t) = \begin{cases} (0, 1), & \text{if } \delta_b > 0.1 \text{ or } \delta_d > 0.1 \\ (1, 0), & \text{if } -0.1 \leq \delta_b \leq 0.01 \text{ and } -0.1 \leq \delta_d \leq 0.01 \\ (0, 0), & \text{otherwise.} \end{cases} \quad (7)$$

Each of our simulations are repeated 10 times and our results are significant based on *t-test* with a confidence interval of 0.95.

Without any evidence, the trustworthiness of sources is computed as 0.5. Thus, there are no conflict in the beginning of our simulations. If all sources have deterministic behaviours, i.e., malicious sources are simple liars and  $P_b = 0$ , then trustworthiness of sources are easily modelled over time and the opinions from liars are significantly discounted. In such settings, conflicts are totally avoided and information consumers using either of the four proposed methods have the same level of success. Figure 3 shows an example of this setting where honest sources always provides the truth ( $P_b = 0$ ) and malicious sources are simple liars. Here, the *ratio of liars* ( $R_{liar}$ ) is 0.5, i.e., half of the sources are malicious.

When honest sources provide bogus information occasionally, the conflicts may arise in the knowledge base of the consumer, because the information from these sources are not significantly discounted. Figure 4 shows our results for  $R_{liar} = 0.5$  and  $P_b = 0.1$ , where all malicious sources are simple liars. In this setting, NDL leads to significant errors in the computed interpretations. While TDL does much better than NDL, it is outperformed by discounting based approaches TDC and EDC. Both of these approaches have similarly good performance though TDC does slightly better.

Simple liars may not be enough to model malicious sources in real life. That is why we change the type of malicious sources

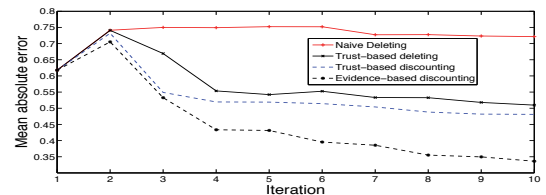


Figure 5: Strategic liars ( $R_{liar} = 0.5$  and  $P_b = 0.1$ )

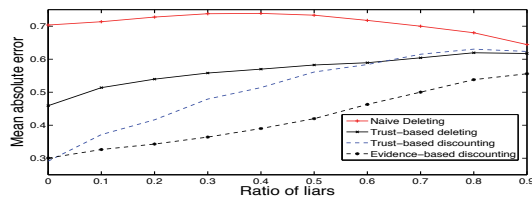


Figure 6: Strategic liars with varying  $R_{liar}$  ( $P_b = 0.1$ )

to strategic liars and repeat our simulations. Figure 5 shows our results for  $R_{liar} = 0.5$  and  $P_b = 0.1$ . In this settings, trust evaluations become misleading, since strategic liars build trust, make their impact and then change their identity to avoid any negative evidence. As a result, as shown in the figure, TDC fails significantly more than EDC after a few iterations. We repeat the simulations with strategic liars for different  $R_{liar}$  values; our results are shown in Figure 6. Our results indicate that evidence-based discounting is much more robust in the presence of realistic malicious behaviour than trust-based discounting or deletion.

## 6. DISCUSSION

DL-Lite is a tractable subset of DLs with a large number of application areas [4]. Its scalability makes it very useful especially for the settings where large amount of data should be queried. However, in a network of heterogeneous sources, any information provided by the sources could be uncertain, incomplete, and even conflicting. DL-Lite cannot accommodate such information. Pan et al. [11] proposed a framework of tractable query answering algorithms for a family of fuzzy query languages over large fuzzy DL-Lite [16] ontologies. On the other hand, DST and its extensions such as Subjective Logic explicitly takes into account *uncertainty* and *belief ownership* [9].

Gobeck and Halaschek [8] present a belief revision algorithm for OWL-DL, which is based on trust degrees to remove conflicting statements from a knowledge base. However, as the authors point out, the proposed algorithm is not guaranteed to be optimal. In our work, we embed statement retraction implicitly into the opinion revision procedure with a global optimal criteria which is grounded on a Beta distribution formalisation of trust.

Fact-finding algorithms aim to identify the *truth* given conflicting claims. Pasternack and Roth [12] propose to translate these claims to a linear program, which is solved to obtain belief scores over claims. For example, with TruthFinder [17], the belief scores obtained can be interpreted as the result of simultaneously minimising the frustration coming from the sources against the claims. These approaches do not consider semantics while reasoning about belief and trustworthiness as we do here.

Costa Pereira and Tettamanzi [6] deal with belief changes in an agent's mental state considering trust in information sources. Dong et al. [7] propose to resolve conflicts in information from multiple sources by a voting mechanism. Double counting in votes is avoided by taking into account information dependence among the sources. The dependence is derived from Bayesian analysis over data sets held by the sources with a statistical interpretation.

In this paper, we propose SDL-Lite, which is expressive enough to represent and reason about uncertain information using trust and domain knowledge. It allows us to efficiently identify conflicting information with respect to domain constraints. Then, these conflicts are resolved through the methods we propose for trust revision. Through simulations, we show that our approach can successfully handle highly misleading information in challenging settings. The simulations also show that the approach is robust in the face of

strategic liars. In this paper, mostly for clarity, we assume opinions about each assertion is provided by a single information source. In the future, we will extend our approach to handle multiple sources.

## 7. REFERENCES

- [1] F. Baader, D. L. McGuinness, D. Nardi, and P. Patel-Schneider, editors. *Description Logic Handbook: Theory, implementation and applications*. Cambridge University Press, 2002.
- [2] C. Burnett, T. J. Norman, and K. Sycara. Bootstrapping trust evaluations through stereotypes. In *Proceedings of the Ninth International Conference on Autonomous Agents and Multiagent Systems*, pages 241–248, 2010.
- [3] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Data complexity of query answering in description logics. In *Proc. of KR 2006*, pages 260–270, 2006.
- [4] D. Calvanese, G. Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Tractable reasoning and efficient query answering in description logics: The dl-lite family. *J. Autom. Reason.*, 39(3):385–429, 2007.
- [5] D. Calvanese, G. D. Giacomo, M. Lenzerini, R. Rosati, and G. Vetere. DL-Lite: Practical Reasoning for Rich DLs. In *Proc. of the DL2004 Workshop*, 2004.
- [6] C. da Costa Pereira and A. G. B. Tettamanzi. Goal generation with relevant and trusted beliefs. In *Proceedings of the 7th International Conference on Autonomous Agents and Multiagent Systems*, Estoril, Portugal, 2008.
- [7] X. L. Dong, L. Berti-Equille, and D. Srivastava. Integrating conflicting data: The role of source dependence. In *Proc. of the 35th International Conference on Very Large Databases*, Lyon, France, August 2009.
- [8] J. Golbeck and C. Halaschek-Wiener. Trust-based revision for expressive web syndication. *Journal of Logic and Computation*, 19(5):771–790, Oct. 2009.
- [9] A. Jøsang. *Subjective Logic*. Book Draft, 2011.
- [10] A. Jøsang and R. Ismail. The beta reputation system. In *Proc. of the 15th Bled Electronic Commerce Conference e-Reality: Constructing the e-Economy*, pages 48–64, 2002.
- [11] J. Z. Pan, G. Stamou, G. Stoilos, S. Taylor, and E. Thomas. Scalable Querying Services over Fuzzy Ontologies. In *the Proc. of the 17th International World Wide Web Conference (WWW2008)*, 2008.
- [12] J. Pasternack and D. Roth. Knowing what to believe (when you already know something). In *Proc. of the 23rd International Conference on Computational Linguistics*, Beijing, China, 2010.
- [13] A. Rogers, R. K. Dash, and N. R. Jennings. Computational mechanism design for information fusion within sensor networks. In *In Proceedings of The 9th International Conference on Information Fusion*, 2006.
- [14] A. Ruszczyński. *Nonlinear optimization*, volume 13. Princeton university press, 2011.
- [15] G. Shafer. *A mathematical theory of evidence*. Princeton university press, 1976.
- [16] U. Straccia. Answering vague queries in fuzzy DL-Lite. In *Proc. of the 11th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*, pages 2238–2245, 2006.
- [17] X. Yin, J. Han, and P. S. Yu. Truth discovery with multiple conflicting information providers on the web. In *Proceedings of the Conference on Knowledge and Data Discovery*, 2007.