

# Belief Revision and Trust

**Aaron Hunter**

British Columbia Institute of Technology  
Burnaby, Canada  
aaron.hunter@bcit.ca

## Abstract

Belief revision is the process in which an agent incorporates a new piece of information together with a pre-existing set of beliefs. When the new information comes in the form of a report from another agent, then it is clear that we must first determine whether or not that agent should be trusted. In this paper, we provide a formal approach to modeling trust as a pre-processing step before belief revision. We emphasize that trust is not simply a relation between agents; the trust that one agent has in another is often restricted to a particular domain of expertise. We demonstrate that this form of trust can be captured by associating a state-partition with each agent, then relativizing all reports to this state partition before performing belief revision. In this manner, we incorporate only the part of a report that falls under the perceived domain of expertise of the reporting agent. Unfortunately, state partitions based on expertise do not allow us to compare the relative strength of trust held with respect to different agents. To address this problem, we introduce pseudometrics over states to represent differing degrees of trust. This allows us to incorporate simultaneous reports from multiple agents in a way that ensures the most trusted reports will be believed.

## Introduction

The notion of trust must be addressed in many agent communication systems. In this paper, we consider one isolated aspect of trust: the manner in which trust impacts the process of belief revision. Some of the most influential approaches to belief revision have used the simplifying assumption that all new information must be incorporated; however, this is clearly untrue in cases where information comes from an untrusted source. In this paper, we are concerned with the manner in which an agent uses an external notion of trust in order to determine how new information should be integrated with some pre-existing set of beliefs.

Our basic approach is the following. We introduce a simple model of trust that allows an agent to determine if a source can be trusted to distinguish between different pairs of states. We use this notion of trust as a precursor to belief revision. Hence, before revising by a new formula, an agent first determines to what extent the source of the information can be trusted. In many cases, the agent will only incorporate “part” of the formula into their beliefs. We then extend our model of trust to a more general setting, by introducing quantitative measures of trust that allow us to compare the

degree to which different agents are trusted. Fundamental properties are introduced and established, and applications are considered.

## Preliminaries

### Intuition

It is important to note that an agent typically does not trust another agent universally. As such, we will not apply the label “trusted” to another agent; instead, we will say that an agent is trusted with respect to a certain domain of knowledge. This is further complicated by the fact that there are different reasons that an agent may not be trusted. For example, an agent might not be trusted due to their perceived knowledge of a domain. In other cases, an agent might not be trusted due to their perceived dishonesty, or bias. In this paper, our primary focus is on trust as a function of the perceived expertise of other agents. Towards the end, we briefly address the different formal mechanisms that would be required to deal with deceit.

### Motivating Example

We introduce a motivating example in commonsense reasoning where an agent must rely on an informal notion of trust in order to inform rational belief change; we will return to this example periodically as we introduce our formal model.

Consider an agent that visits a doctor, having difficulty breathing. Incidentally, the agent is wearing a necklace that prominently features a jewel on a pendant. During the examination, the doctor checks the patient’s throat for swelling or obstruction; at the same time, the doctor happens to look at the necklace. Following the examination, the doctor tells the patient “you have a viral infection in your throat - and by the way, you should know that the jewel in your necklace is not a diamond.”

The important part about this example is the fact that the doctor provides information about two distinct domains: human health and jewelry. In practice, a patient is very likely to trust the doctor’s diagnosis about the viral infection. On the other hand, the patient really has very little reason to trust the doctor’s evaluation of the necklace. We suggest that a rational agent should actually incorporate the doctor’s statement about the infection into their own beliefs, while essentially

ignoring the comment on the necklace. This approach is dictated by the kind of trust that the patient has in the doctor. Our aim in this paper is to formalize this kind of “localized” domain-specific trust, and then demonstrate how this form of trust is used in practice to inform belief revision.

## Trust

Trust consists of two related components. First, we can think of trust in terms of how likely an agent is to believe what another agent says. Alternatively, we can think of trust in terms of the degree to which an agent is likely to allow another to perform actions on their behalf. In this paper, we will be concerned only with the former.

A great deal of existing work on trust focuses on the manner in which an agent develops a *reputation* based on past behaviour. A brief survey of reputation systems is given in (Huynh, Jennings, and Shadbolt 2006). Reputation systems can be used to inform the allocation of tasks (Ramchurn et al. 2009), or to avoid deception (Salehi-Abari and White 2009). The model of trust presented in this paper is not intended to be an alternative to existing reputation systems; we are not concerned with the manner in which an agent learns to trust another. Instead, our focus is simply on developing a suitable model of trust that is expressive enough to inform the process of *belief revision*. The manner in which this model of trust is developed over time is beyond the scope of this paper.

## Belief Revision

*Belief revision* refers to the process in which an agent must integrate new information with some pre-existing beliefs about the state of the world. One of the most influential approaches to belief revision is the AGM approach, in which an agent incorporates the new information while keeping as much of the initial belief state as consistently possible (Alchourrón, Gärdenfors, and Makinson 1985).

This approach was originally defined with respect to a finite set  $P$  of propositional variables representing properties of the world. A *state* is a propositional interpretation over  $P$ , representing a possible state of the world. A *belief set* is a deductively closed set of formulas, representing the beliefs of an agent. Since  $P$  is finite, it follows that every belief set defines a corresponding *belief state*, which is the set of states that an agent considers to be possible. A revision operator is a function that takes a belief set and a formula as input, and returns a new belief set. An AGM revision operator is a revision operator that satisfies the AGM postulates, as specified in (Alchourrón, Gärdenfors, and Makinson 1985).

It turns out that every AGM revision operator is characterized by a total pre-order over possible worlds. To be more precise, a *faithful assignment* is a function that maps each belief set to a total pre-order over states in which the models of the belief set are the minimal states. When an agent is presented with a new formula  $\phi$  for revision, the revised belief state is the set of all minimal models of  $\phi$  in the total pre-order given by the faithful assignment. We refer the reader to (Katsuno and Mendelzon 1992) for a proof of this result, as well as a complete description of the implications.

For our purposes, we simply need to know that each AGM revision operator necessarily defines a faithful assignment.

## A Model of Trust

### Domain-Specific Trust

Assume we have a fixed propositional signature  $\mathbf{F}$  as well as a set of agents  $\mathbf{A}$ . For each  $A \in \mathbf{A}$ , let  $Bel_A$  denote a deductively closed set of formulas over  $\mathbf{F}$  called the *belief set* of  $A$ . For each  $A$ , let  $*_A$  denote an AGM revision operator that intuitively captures the way that the agent  $A$  revises their beliefs when presented with new information. This revision operator represents sort of an “ideal” revision situation, in which  $A$  has complete trust in the new information. We want to modify the way this operator is used, by adding a representation of the extent to which  $A$  trusts each other agent  $B \in \mathbf{A}$  over  $\mathbf{F}$ .

We assume that all new information is reported by an agent, so each formula for revision can be labelled with the name of the reporting agent.<sup>1</sup> At this point, we are not concerned with degrees of trust or with resolving conflicts between different sources of information. Instead, we start with a binary notion of trust, where  $A$  either trusts  $B$  or does not trust  $B$  with respect to a particular domain of expertise.

We encode trust by allowing each agent  $A$  to associate a partition  $\Pi_A^B$  over possible states with each agent  $B$ .

**Definition 1** A state partition  $\Pi$  is a collection of subsets of  $2^{\mathbf{F}}$  that is collectively exhaustive and mutually exclusive. For any state  $s \in 2^{\mathbf{F}}$ , let  $\Pi(s)$  denote the element of  $\Pi$  that contains  $s$ .

If  $\Pi = \{2^{\mathbf{F}}\}$  then we call  $\Pi$  the *trivial partition* with respect to  $\mathbf{F}$ . If  $\Pi = \{\{s\} \mid s \in 2^{\mathbf{F}}\}$ , then we call  $\Pi$  the *unit partition*.

**Definition 2** For each  $A \in \mathbf{A}$  the trust function  $T_A$  is a function that maps each  $B \in \mathbf{A}$  to a state partition  $\Pi_A^B$ .

The partition  $\Pi_A^B$  represents the trust that  $A$  has in  $B$  over different aspects of knowledge. Informally, the partition encodes states that  $A$  will trust  $B$  to distinguish. If  $\Pi_A^B(s_1) \neq \Pi_A^B(s_2)$ , then  $A$  will trust that  $B$  can distinguish between states  $s_1$  and  $s_2$ . Conversely, if  $\Pi_A^B(s_1) = \Pi_A^B(s_2)$ , then  $A$  does not see  $B$  as an authority capable of distinguishing between  $s_1$  and  $s_2$ . We clarify by returning to our motivating example.

**Example** Let  $\mathbf{A} = \{A, D, J\}$  and let  $\mathbf{F} = \{sick, diam\}$ . Informally, the fluent *sick* is true if  $A$  has an illness and the fluent *diam* is true if a certain piece of jewelry that  $A$  is wearing contains a real diamond. If we imagine that  $D$  represents a doctor and  $J$  represents a jeweler, then we can use state partitions to represent the trust that  $A$  has in  $D$  and  $J$  with respect to different domains. Following standard shorthand notation, we represent a state  $s$  by the set of fluent symbols that are *true* in  $s$ . In order to make the descriptions of a partition more readable, we use a  $|$  symbol to visually

<sup>1</sup>This is not a significant restriction. In domains involving sensing or other forms of discovery, we could simply allow an agent  $A$  to self-report information with complete trust.

separate different cells. The following partitions are then intuitively plausible in this example:

$$\begin{aligned}\Pi_A^D &:= \{sick, diam\}, \{sick\}|\{diam\}, \emptyset \\ \Pi_A^J &:= \{sick, diam\}, \{diamond\}|\{sick\}, \emptyset\end{aligned}$$

Hence,  $A$  trusts the doctor  $D$  to distinguish between states where  $A$  is sick as opposed to states where  $A$  is not sick. However,  $A$  does not trust  $D$  to distinguish between worlds that are differentiated by the authenticity of a diamond. The formula  $sick \wedge \neg diamond$  encodes the doctor's statement that the agent is sick, and the necklace they are wearing has a fake diamond.

Although the preceding example is simple, it illustrates how a partition can be used to encode the perceived expertise of agents. In the doctor-jeweler example, we could equivalently have defined trust with respect to the set of fluents. In other words, we could have simply said that  $D$  is trusted over the fluent  $sick$ . However, there are many practical cases where this is not sufficient; we do not want to rely on the fluent vocabulary to determine what is a valid feature with respect to trust. For example, a doctor may have specific expertise over lung infections for those working in factories, but not for lung infections for those working in a space shuttle. By using state partitions to encode trust, we are able to capture a very flexible class of distinct areas of trust.

### Incorporating Trust in Belief Revision

As indicated previously, we assume each agent  $A$  has an AGM belief revision operator  $*_A$  for incorporating new information. In this section, we describe how the revision operator  $*_A$  is combined with the trust function  $T_A$  to define a new, trust-incorporating revision operator  $*_A^B$ . In many cases, the operator  $*_A^B$  will not be an AGM operator because it will fail to satisfy the AGM postulates. In particular,  $A$  will not necessarily believe a new formula when it is reported by an untrusted source. This is a desirable feature.

Our approach is to define revision as a two-step process. First, the agent considers the source and the relevant state partition to determine how much of the new information to incorporate. Second, the agent performs standard AGM revision using the faithful assignment corresponding to the belief revision operator.

**Definition 3** Let  $\phi$  be a formula and let  $T_A(B) = \Pi_A^B$ . Define:

$$\Pi_A^B[\phi] = \bigcup \{\Pi_A^B(s) \mid s \models \phi\}.$$

Hence  $\Pi_A^B[\phi]$  is the union of all cells that contain a model of  $\phi$ .

If  $A$  does not trust  $B$  to distinguish between states  $s$  and  $t$ , then any report from  $B$  that provides evidence that  $s$  is the actual state is also evidence that  $t$  is the actual state. When  $A$  performs belief revision, it should be with respect to the distinctions that  $B$  can be trusted to make. It follows that  $A$  need not believe  $\phi$  after revision; instead  $A$  should interpret  $\phi$  to be evidence of any state  $s$  that is  $B$ -indistinguishable from a model of  $\phi$ . Formally, this means that the formula  $\phi$  is construed to be evidence for each state in  $\Pi_A^B[\phi]$ .

**Definition 4** Let  $A, B \in \mathbf{A}$  with  $T_A(B) = \Pi_A^B$ , and let  $*_A$  be an AGM revision operator for  $A$ . For any belief set  $K$  with corresponding ordering  $\prec_K$  given by the underlying faithful assignment, the trust-sensitive revision  $K *_A^B \phi$  is the set of formulas true in

$$\min_{\prec_K}(\{s \mid s \in \Pi_A^B[\phi]\}).$$

So rather than taking the minimal models of  $\phi$ , we take all minimal states that  $B$  can not be trusted to distinguish from the minimal models of  $\phi$ .

It is worth remarking that this notion can be formulated syntactically as well. Since  $\mathbf{F}$  is finite, each state  $s$  is defined by a unique, maximal conjunction over literals in  $\mathbf{F}$ ; we simply take the conjunction of all the atomic formulas that are true in  $s$  together with the negation of all the atomic formulas that are false in  $s$ .

**Definition 5** For any state  $s$ , let  $prop(s)$  denote the unique, maximal conjunction of literals true in  $s$ .

This definition can be extended for a cell in a state partition.

**Definition 6** Let  $\Pi$  be a state partition. For any state  $s$ ,

$$prop(\Pi(s)) = \bigvee \{prop(s') \mid s' \in \Pi(s)\}.$$

Note that  $prop(\Pi(s))$  is a well-defined formula in disjunctive normal form, due to the finiteness of  $\mathbf{F}$ . Intuitively,  $prop(\Pi(s))$  is the formula that defines the partition  $\Pi(s)$ . In the case of a trust partition  $\Pi_A^B$ , we can use this idea to define the *trust expansion* of a formula.

**Definition 7** Let  $A, B \in \mathbf{A}$  with the corresponding state partition  $\Pi_A^B$ , and let  $\phi$  be a formula. The trust expansion of  $\phi$  for  $A$  with respect to  $B$  is the formula

$$\phi_A^B := \bigvee \{prop(\Pi_A^B(s)) \mid s \models \phi\}.$$

Note that this is a finite disjunction of disjunctions, which is again a well defined formula. We refer to  $\phi_A^B$  as the trust expansion of  $\phi$  because it is true in all states that are consistent with  $\phi$  with respect to distinctions that  $A$  trusts  $B$  to be able to make. It is an expansion because the set of models of  $\phi_A^B$  is normally larger than the set of models of  $\phi$ . The trust sensitive revision operator could equivalently be defined as the normal revision, following translation of  $\phi$  to the corresponding trust expansion.

**Example** Returning to our example, we consider a few different formulas for revision:

1.  $\phi_1 = sick$
2.  $\phi_2 = \neg diam$
3.  $\phi_3 = sick \wedge \neg diam$ .

Suppose that the agent initially believes that they are not sick, and that the diamond they have is real, so  $K = \neg sick \wedge diam$ . For simplicity, we will assume that the underlying pre-order  $\prec_K$  has only two levels: those states where  $K$  is true are minimal, and those where  $K$  is false are not. We have the following results for revision

1.  $K *_A^D \phi_1 = sick \wedge diam$
2.  $K *_A^D \phi_2 = \neg sick \wedge diam$
3.  $K *_A^D \phi_3 = sick \wedge diam.$

The first result indicates that  $A$  believes the doctor when the doctor reports that they are sick. The second result indicates that  $A$  essentially ignores a report from the doctor on the subject of jewelry. The third result is perhaps the most interesting. It demonstrates that our approach allows an agent to just incorporate a part of a formula. Hence, even though  $\phi_3$  is given as a single piece of information, the agent  $A$  only incorporates the part of the formula over which the doctor is trusted.

## Formal Properties

### Basic Results

We first consider extreme cases for trust-sensitive revision operators. Intuitively, if  $T_A(B)$  is the trivial partition, then  $A$  does not trust  $B$  to be able to distinguish between any states. Therefore,  $A$  should not incorporate any new information obtained from  $B$ . The following proposition makes this observation explicit.

**Proposition 1** *If  $T_A(B)$  is the trivial partition, then  $K *_A^B \phi = K$  for all  $K$  and  $\phi$ .*

The other extreme situation occurs when  $T_A(B)$  is the unit partition, which consists of all singleton sets. In this case,  $A$  trusts  $B$  to be able to distinguish between every possible pair of states. It follows from this result that trust sensitive revision operators are not AGM revision operators.

**Proposition 2** *If  $T_A(B)$  is the unit partition, then  $*_A^B = *_A$ .*

Hence, if  $B$  is universally trusted, then the corresponding trust sensitive revision operator is just the a priori revision operator for  $A$ .

### Refinements

There is a partial ordering on partitions based on the notion of *refinement*. We say that  $\Pi_1$  is a refinement of  $\Pi_2$  just in case, for each  $S_1 \in \Pi_1$ , there exists  $S_2 \in \Pi_2$  such that  $S_1 \subseteq S_2$ . We also say that  $\Pi_1$  is *finer* than  $\Pi_2$ . In terms of trust-partitions, refinement has a natural interpretation in terms of “breadth of trust.” If the partition corresponding to  $B$  is finer than that corresponding to  $C$ , it means that  $B$  is trusted more broadly than  $C$ . To be more precise, it means that  $B$  is trusted to distinguish between all of the states that  $C$  can distinguish, and possibly more. If  $B$  is trusted more broadly than  $C$ , it follows that a report from  $B$  should give give  $A$  more information. This idea is formalized in the following proposition.

**Proposition 3** *For any formula  $\phi$ , if  $\Pi_A^B$  is a refinement of  $\Pi_A^C$ , then  $|K *_A^B \phi| \subseteq |K *_A^C \phi|$ .*

This is a desirable property; if  $B$  is trusted over a greater range of states, then fewer states are possible after a report from  $B$ .

## Multiple Reports

One natural question that arises is how to deal with multiple reports of information from different agents, with different trust partitions. In our example, for instance, we might get a conflicting report from a jeweler with respect to the status of the necklace. In order to facilitate the discussion, we introduce a precise notion of a *report*.

**Definition 8** *A report is a pair  $(B, \phi)$ , where  $B \in \mathbf{A}$  and  $\phi$  is a formula.*

We can now extend the definition of trust sensitive revision to reports in the obvious manner. In fact, if the revising agent  $A$  is clear from the context, we can use the short hand notation:

$$K * (\phi, B) = K *_A^B \phi.$$

The following definition extends the notion of revision to incorporate multiple reports.

**Definition 9** *Let  $\{A\} \cup B \subseteq \mathbf{A}$ , and let  $\Phi = \{(\phi_i, B_i) \mid i < n\}$  be a finite set of reports. Given  $K$ ,  $*$  and  $\prec_K$ , the trust-sensitive revision  $K *_A \Phi$  is the set of formulas true in*

$$\min_{\prec_K}(\{s \mid s \in \Pi_A^{B_i}[\phi_i]\}).$$

So the trust sensitive revision for a finite set of reports from different agents is essentially the normal, single-shot revision by the conjunction of formulas. The only difference is that we expand each formula with respect to the trust partition for a particular reporting agent.

**Example** In the doctor and jeweler domain, we can consider how an agent might incorporate a set of reports from  $D$  and  $J$ . We start with the same initial belief set as before:  $K = \neg sick \wedge diam$ . Consider the following reports:

1.  $\Phi_1 = \{(sick, D), (\neg diam, D)\}$
2.  $\Phi_2 = \{(sick, J), (\neg diam, J)\}$
3.  $\Phi_3 = \{(sick, D), (\neg diam, J)\}$
4.  $\Phi_4 = \{(sick, J), (\neg diam, D)\}$

We have the following results following revision:

1.  $K *_A \Phi_1 = sick \wedge diam$
2.  $K *_A \Phi_2 = \neg sick \wedge \neg diam$
3.  $K *_A \Phi_3 = sick \wedge \neg diam$
4.  $K *_A \Phi_4 = \neg sick \wedge diam.$

These results demonstrate how the agent  $A$  essentially incorporates information from  $D$  and  $J$  in domains where they are trusted, and ignores information when they are not trusted. Note that, in this case,  $D$  and  $J$  are trusted over disjoint sets of states. As a result, it is not possible to have contradictory reports that are equally trusted.

The problem with Definition 9 is that the set of states in the minimization may be empty. This occurs when multiple agents give conflicting reports, and we trust each agent on the domain. In order to resolve this kind of conflict, we need a more expressive form of trust that allows some agents to be trusted more than others. We introduce such a representation in the next section.

## Trust Pseudometrics

### Measuring Trust

In the previous section, we were concerned with a binary notion of trust that did not include any measure of the strength of trust held in a particular agent or domain. Such an approach is appropriate in cases where we only receive new information from a single source, or from a set of sources that are equally reliable. However, it is not sufficient if we consider cases where several different sources may provide conflicting information. In such cases, we need to determine which information source is the most trust worthy with respect to the domain currently under consideration.

In the binary approach, we associated a partition of the state space with each agent. In order to capture different levels of trust, we would like to introduce a measure of the distance between two states from the perspective of a particular agent. In other words, an agent  $A$  would like to associate a distance function  $d_B$  over states with each other agent  $B$ . If  $d_B(s, t) = 0$ , then  $B$  can not be trusted to distinguish between the states  $s$  and  $t$ . On the other hand, if  $d_B(s, t)$  is very large, then  $A$  has a high level of trust in  $B$ 's ability to distinguish between  $s$  and  $t$ . The notion of distance that we introduce will be a *pseudometric* on the state space. A pseudometric is a function  $d$  that satisfies the following properties for all  $x, y, z$  in the domain  $X$ :

1.  $d(x, x) = 0$
2.  $d(x, y) = d(y, x)$
3.  $d(x, z) \leq d(x, y) + d(y, z)$

The difference between a *metric* and a pseudometric is that we do not require that  $d(x, y) = 0$  implies  $x = y$  (the so-called law of indiscernables). This would be undesirable in our setting, because we want to use the distance 0 to represent states that are indistinguishable rather than identical. The first two properties are clearly desirable for a measure of our trust in another agent's ability to discern states. The third property is the triangle inequality, and it is required to guarantee that our trust in other agents is transitive across different domains.

**Definition 10** For each  $A \in \mathbf{A}$ , a pseudometric trust function  $\mathcal{T}_A$  is a function that maps each  $B \in \mathbf{A}$  to a pseudometric  $d_B$  over  $2^{\mathbf{F}}$ .

The pair  $(2^{\mathbf{F}}, \mathcal{T}_A)$  is called a pseudometric trust space. We would like to model the situation where a sequence of formulas  $\Phi = \phi_1, \dots, \phi_n$  is received from the agents  $\mathbf{B} = B_1, \dots, B_n$ , respectively. Note that the order does not matter, we think of the formulas as arriving at the same instant with no preference between them other than the preference induced by the pseudometric trust space.

We associate a sequence of state partitions with each pseudometric trust space.

**Proposition 4** Let  $(2^{\mathbf{F}}, \mathcal{T}_A)$  be a pseudometric trust space, let  $B \in \mathbf{A} - A$ , and let  $i$  be a natural number. For each state  $s$ , define the set  $\Pi_B^A(i)(s)$  as follows:

$$\Pi_B^A(i)(s) = \{t \mid d_B(s, t) \leq i\}.$$

The collection of sets  $\{\Pi_B^A(i)(s) \mid s \in 2^{\mathbf{F}}\}$  is a state partition.

We let  $\Pi_B^A(i)$  denote the state partition obtained from this proposition. The cells of the partition  $\Pi_B^A(i)$  consist of all states are separated by a distance of no more than  $i$ . The following proposition is immediate.

**Proposition 5**  $\Pi_B^A(i)$  is a refinement of  $\Pi_B^A(j)$ , for any  $i < j$ .

Hence, a pseudometric trust space defines a sequence of partitions for each agent. This sequence of partitions gets coarser as we increase the index; increasing the index corresponds to requiring a higher level of trust that an agent can distinguish between states. Since we can use Definition 4 to define a trust sensitive revision operator from a state partition, we can now define a trust sensitive revision operator for any fixed distance  $i$  between states. Informally, as  $i$  increases, we require  $B$  to have a greater degree of certainty in order to trust them to distinguish between states. However, it is not clear in advance exactly which  $i$  is the right threshold. Our approach will be to find the lowest possible threshold that yields a consistent result.

Note that  $\Pi_B^A(i)$  will be a trivial partition for any  $i$  that is less than the minimum distance assigned by the underlying pseudometric trust function.

**Definition 11** Let  $(2^{\mathbf{F}}, \mathcal{T}_A)$  be a pseudometric trust space, and let  $m$  be the least natural number such that  $\Pi_B^A(m)$  is non-trivial. The trust sensitive revision operator for  $A$  with respect to  $B$  is the trust sensitive revision operator given by  $\Pi_B^A(m)$ .

This is a simple extension of our approach based on state partitions. In the next section, we take advantage of the added expressive power of pseudometrics.

**Example** We modify the doctor example. In order to consider different levels of trust, it is more interesting to consider a domain involving two doctors: a general practitioner  $D$  and a specialist  $S$ . We also assume that the vocabulary includes two fluents: *ear* and *skin*. Informally, *ear* is understood to be true if the patient has an ear infection, whereas *skin* is true if the patient has skin cancer. The important point is that an ear infection is something that can easily be diagnosed by any doctor, whereas skin cancer is typically diagnosed by a specialist. In order to capture these facts, we define two pseudometrics  $d_D$  and  $d_S$ . For simplicity, we label the possible states as follows:

$$\begin{aligned} s_1 &= \{ear, skin\} \\ s_2 &= \{ear\} \\ s_3 &= \{skin\} \\ s_4 &= \emptyset \end{aligned}$$

We define the pseudometrics as follows: With these pseudo-

	$s_1, s_2$	$s_1, s_3$	$s_1, s_4$	$s_2, s_3$	$s_2, s_4$	$s_3, s_4$
$d_D$	1	2	2	2	2	1
$d_S$	2	2	2	2	2	2

metrics, it is easy to see that both  $D$  and  $S$  can distinguish all

of the states. However,  $S$  is more trusted to distinguish between states related to a skin cancer diagnosis. In our framework, we would like to ensure that this implies  $S$  will be trusted in the case of conflicting reports from  $D$  and  $S$  with respect to skin cancer.

## Multiple Reports

We view the distances in a pseudometric trust space as absolute measurements. As such, if  $d_B(s, t) > d_C(s, t)$ , then we have greater trust in  $B$  as opposed to  $C$  as far as the ability to discern the states  $s$  and  $t$  is concerned. We would like to use this intuition to resolve conflicting reports between agents.

**Proposition 6** *Let  $\{A\} \cup B \subseteq A$ , and let  $\Phi = \{(\phi_i, B_i) \mid i < n\}$  be a finite set of reports. There exists a natural number  $m$  such that*

$$\bigcap_{i < n} (\Pi_A^{B_i}[\phi_i](m)) \neq \emptyset.$$

Hence, for any set of reports, we can get a non-intersecting intersection if we take a sufficiently coarse state partition. In many cases this partition will be non-trivial. Using this proposition, we define multiple report revision as follows.

**Definition 12** *Let  $(2^F, \mathcal{T}_A)$  be a pseudometric trust space, let  $\Phi = \{(\phi_i, B_i) \mid i < n\}$  be a finite set of reports, and let  $m$  be the least natural number such that  $\bigcap_{i < n} (\Pi_A^{B_i}[\phi_i](m)) \neq \emptyset$ . Given  $K, *$  and  $\prec_K$ , the trust-sensitive revision  $K *_A^B \Phi$  is the set of formulas true in*

$$\min_{\prec_K} (\{s \mid s \in \Pi_A^{B_i}[\phi_i](m)\}).$$

Hence, trust-sensitive revision in this context involves finding the finest possible partition that provides a meaningful combination of the reports, and then revising with the corresponding state partition.

## Trust and Deceit

To this point, we have only been concerned with modeling the trust that one agent holds in another due to perceived knowledge or expertise. Of course, the issue of trust also arises in cases where one agent suspects that another may be dishonest. However, the manner in which trust must be handled differs greatly in this context. If  $A$  does not trust  $B$ , then there is little reason for  $A$  to believe any part of a message sent directly from  $B$ .

## Discussion

### Related Work

We are not aware of any other work on trust that explicitly deals with the interaction between trust and formal belief revision operators. There is, however, a great deal of work on frameworks for modelling trust. As noted previously, the focus of such work is often on building reputations. One notable approach to this problem with an emphasis on knowledge representation is (Wang and Singh 2007),

in which trust is built based on evidence. This kind of approach could be used as a precursor step to build a trust metric, although one would need to account for domain expertise.

Different levels of trust are treated in (Krukow and Nielsen 2007), where a lattice structure is used to represent various levels of trust strength. This is similar to our notion of a trust pseudometric, but it permits incomparable elements. There are certainly situations where this is a reasonable advantage. However, the emphasis is still on the representation of trust in *an agent* as opposed to trust in an agent with respect to a domain.

One notable approach that is similar to ours is the semantics of trust presented in (Krukow and Nielsen 2007), which is a domain-based approach to differential trust in an agent. The emphasis there is on *trust management*, however. That is, the authors are concerned with how agents maintain some record of trust in the other agents; they are not concerned with a differential approach to belief revision.

## Conclusion

In this paper, we have developed an approach to trust sensitive belief revision in which an agent is trusted only with respect to a particular domain. This has been formally accomplished first by using state partitions to indicate which states an agent can be trusted to distinguish, and then by using distance functions to quantify the strength of trust. In both cases, the model of trust is used as sort of a precursor to belief revision. Each agent is able to perform belief revision based on a pre-order over states, but the actual formula for revision is parametrized and expanded based on the level of trust held in the reporting agent.

There are many directions for future work, in terms of both theory and applications. As noted previously, one of the subtle distinctions that must be addressed is the difference between trusted *expertise* and trusted *honesty*. The present framework does not explicitly deal with the problem of deception or *belief manipulation* (Hunter 2013); it would be useful to explore how models of trust must differ in this context. In terms of applications, our approach could be used in any domain where agents must make decisions based on beliefs formulated from multiple reports. This is the case, for example, in many networked communication systems.

## References

- [Alchourrón, Gärdenfors, and Makinson 1985] Alchourrón, C.; Gärdenfors, P.; and Makinson, D. 1985. On the logic of theory change: Partial meet functions for contraction and revision. *Journal of Symbolic Logic* 50(2):510–530.
- [Hunter 2013] Hunter, A. 2013. Belief manipulation: A formal model of deceit in message passing systems. In *Proceedings of the Pacific Asia Workshop on Security Informatics*, 1–8.
- [Huynh, Jennings, and Shadbolt 2006] Huynh, T. D.; Jennings, N. R.; and Shadbolt, N. R. 2006. An integrated trust and reputation model for open multi-agent systems. *Autonomous Agents and Multi-Agent Systems* 13(2):119–154.

- [Katsuno and Mendelzon 1992] Katsuno, H., and Mendelzon, A. 1992. Propositional knowledge base revision and minimal change. *Artificial Intelligence* 52(2):263–294.
- [Krukow and Nielsen 2007] Krukow, K., and Nielsen, M. 2007. Trust structures. *International Journal of Information Security* 6(2-3):153–181.
- [Ramchurn et al. 2009] Ramchurn, S.; Mezzetti, C.; Giovannucci, A.; Rodriguez-Aguilar, J.; Dash, J.; and Jennings, N. 2009. Trust-based mechanisms for robust and efficient task allocation in the presence of execution uncertainty. *JAIR* 35:119–159.
- [Salehi-Abari and White 2009] Salehi-Abari, A., and White, T. 2009. Towards con-resistant trust models for distributed agent systems. In *IJCAI*, 272–277.
- [Wang and Singh 2007] Wang, Y., and Singh, M. P. 2007. Formal trust model for multiagent systems. In *IJCAI*, 1551–1556.