

A Trust Evaluation Method Based on Logic and Probability Theory

Reto Kohlas, Jacek Jonczy, and Rolf Haenni

Abstract We introduce a trust evaluation method applicable in a decentralized setting, in which no universally trusted authority exists. The method makes simultaneous use of logic and probability theory. The result of the qualitative part of the method are logical arguments for and against the reliability of an entity. The quantitative part returns the probability that the reliability of an entity can be deduced under the given assumptions and pieces of evidence, as well a corresponding probability for the counter-hypothesis. Our method is a true generalization of existing methods, in particular the Credential Networks. It relies on digital signatures for authenticating messages and accounts for many-to-many relationships between entities and public keys. Moreover, it includes eight different types of trust relations, namely the assumption or the statement that an entity is honest, competent, reliable, or malicious, and their corresponding negations.

1 Introduction

Members of global social networks and e-commerce systems regularly face the question whether they can trust other, a priori unknown entities. A rating system or a *trust evaluation* method can provide decision support in such a situation. It indicates arguments for the reliability of an entity (or a numerical value representing an entity's reliability, respectively) by taking available trust assumptions, recommendations and discredits into account.

The credibility of a statement, for example a recommendation, generally depends on the reliability of its author (i.e., the honesty and the competence of its author). In

Reto Kohlas and Jacek Jonczy
Institute of Computer Science and Applied Mathematics, University of Berne, 3012 Berne, Switzerland, e-mail: kohlas.jonczy@iam.unibe.ch

Rolf Haenni
Bern University of Applied Sciences, 2501 Biel/Bienne, Switzerland, e-mail: rolf.haenni@bfh.ch

a digital setting, messages should be authenticated, since the identity of the alleged sender of a message can typically be forged without effort. The method presented in this paper makes use of *public-key cryptography* and digital signature schemes for message authentication.

The use of public-key cryptography requires the *authentication of public keys*, i.e., the establishment to which physical entity a public key belongs. *Public-key certificates* are digitally signed statements which approve the authenticity of a public-key entity for a physical entity.¹ They contribute thus to public-key authentication and are useful for those physical entities who cannot exchange their public keys personally.

The main purpose of this paper is to introduce a novel trust evaluation method that relies on logic and probability theory. It uses digital signatures for message authentication and extends previously proposed approaches.

1.1 Existing Trust Evaluation Methods

Some authors have noted that early methods for evaluating trust or for authenticating public keys tend to return counter-intuitive results. Deficiencies in PGP’s Web of Trust for instance have been identified in [21, 18, 13], principles that such methods should ideally fulfill have been stated in [23, 18]. In search of improved techniques, a vast number of methods has been proposed in the last decade.

Some methods combine the confidence values *specifically*, in the sense that their way of combining the confidence values has been exclusively conceived for trust evaluation. Examples of such specific methods are [2, 22, 23, 26, 1, 19, 20]. Other methods treat trust evaluation as a special case of accepting or rejecting a *hypothesis* (that a public key is authentic or that an entity is reliable) under *uncertain assumptions* and *pieces of evidence* (public-key certificates, recommendations, discredits). Such methods use *formal techniques for reasoning under uncertainty*, and are often based on a *probabilistic* interpretation of the confidence values. Examples are Maurer’s Probabilistic Model [21] (based on Probabilistic Logic), Jøsang’s Certification Algebra [14] (based on Subjective Logic), Haenni’s Key Validation Method [6] and the Credential Networks [12] (both based on the Theory of Probabilistic Argumentation).

We here briefly describe Maurer’s probabilistic method (MPM), since it allows us to exemplify in Subsection 1.2 in which sense we intend to extend existing probabilistic methods. The basic idea behind MPM is the combination of *logic* and *probability theory*. MPM’s *deterministic* model consists of two so-called *inference rules*. The first inference rule asserts that if a reasoner A knows the authentic public key of X ($\text{Aut}_{A,X}$), if X is trusted by A for issuing public-key certificates ($\text{Trust}_{A,X,1}$)², and

¹ For reasons explained in Subsection 2.1, these entities are called *public-key* and *physical entities*.

² The third index is an integer and corresponds to the *trust level* (its exact meaning is irrelevant for the discussion in this paper).

if X issues a public-key certificate for Y ($\text{Cert}_{X,Y}$)³, then A can conclude to possess the authentic public key of Y ($\text{Aut}_{A,Y}$). Formally, this rule translates into

$$\forall X \forall Y : \text{Aut}_{A,X} \wedge \text{Trust}_{A,X,1} \wedge \text{Cert}_{X,Y} \vdash \text{Aut}_{A,Y}.$$

The second inference rule (which we do not print here) describes the role of recommendations for evaluating the reliability of a physical entity. Note that MPM considers positive recommendations only (i.e., there are no statements asserting that some entity is unreliable).

The *probabilistic* model of MPM lets A assign a probability to every assumption. Each probability, also called confidence value, is intended to stand for A 's degree of belief with respect to the truth of the judged assumption. MPM then defines confidence value for the hypothesis $\text{Aut}_{A,B}$ as function of the initially attributed probabilities. This confidence value corresponds to the probability that $\text{Aut}_{A,B}$ can be deduced from A 's initial view by applying consecutively the two inference rules of the deterministic model.

1.2 Motivation

In accordance with other contributions, we propose to use a probabilistic framework as the basis of our method. However, we suggest to revisit existing probabilistic methods with respect to the type of assumptions and certificates (credentials) they take into account. The following list discusses some important and often neglected modeling aspects:

- *Physical entities may use multiple public-key entities.* Most methods assume that each physical entity uses *at most* one public-key entity. In MPM, for example, the supposed public key of X is precisely for this reason not included in the statement $\text{Aut}_{A,X}$. As a consequence, statements signed by different public-key entities are usually considered independent. However, it is often impossible in a decentralized system to limit the number of keys used, since each entity can generate as many public-key entities and distribute as many public keys as desired. If some physical entity controls two public-key entities, then statements signed by these two public-key entities are by no means independent.
- *Two physical entities can share the use of a public-key entity.* It is usually impossible to assure that one public-key entity is controlled by only *one* physical entity. A key holder can for instance disclose the passphrase for accessing the private key to another physical entity, and thereby share control of the public-key entity. Such sharing of public-key entities can be problematic. If both physical entities control the public-key entity, it is not possible to uniquely assign state-

³ Note that X uses her public-key entity to issue a certificate for Y 's public key. But neither X nor Y 's public key are parameters in the statement $\text{Cert}_{X,Y}$, because all physical entities are assumed to control *exactly one* public-key entity.

ments signed by the public-key entity to either of the physical entities. As a consequence, if the two physical entities are not equally trusted, it is impossible to determine the credibility of the signed statement in a unique way.

- *The opposite of trust is twofold.* Trust is often modeled as positive assumption only allowing to conclude what trusted introducers say. If an introducer is *not* trusted (e.g., in MPM the statement $\text{Trust}_{A,X,1}$ would not be valid), no conclusions are drawn within these methods. But it is possible that malicious entities lie. In the context of public-key authentication, persuading someone to use the “wrong” public key allows to decrypt messages and make statements in someone else’s name; providing “false” statements about reliability could convince somebody to enter a deal to the cheating entity’s advantage. In Subsection 2.3 we shall therefore differentiate between two opposites of trust: first, as an entity’s belief that a given introducer is *incompetent*, and second as the stronger assumption that the introducer is *malicious*, in which case the contrary of what the introducer says can be deduced.
- *Negative statements.* Many existing methods are monotonic. In MPM for instance, if A adds a public certificate to her view, the confidence value $\text{conf}(\text{Aut}_{A,B})$ remains the same or increases (but it does not decrease). There is never evidence for the hypothesis that a public key is *not* authentic. The reason for the monotonicity of the methods lies in the fact that only positive statements are taken into account. However, negative and positive statements are equally important. If a honest introducer observes that someone else is spreading false information, or that a public key is not authentic, this honest introducer should have a means at hand to warn other participants. We intend therefore to include different types of negative statements in our model.

1.3 Goal and Outline

The goal of this paper is to propose a trust evaluation method that considers the modeling aspects mentioned in the previous subsection. We base our method on the Theory of Probabilistic Argumentation [16, 9, 7] (TPA), which allows us to cope with conflicting assumptions and evidence. Moreover, hypotheses can be evaluated qualitatively and quantitatively. The qualitative part of the method provides logical arguments for and against a hypothesis. The results of the quantitative evaluation are two corresponding probabilities of derivability.

The emphasis of this paper lies primarily in the preciseness and not in the practicability of the proposed method. By suggesting a more accurate model we hope to understand the mechanisms behind trust evaluation better. Aspects of efficiency and usability will be part of future work.⁴

This paper is organized as follows. In Section 2 we introduce our model. Section 3 describes the logical and probabilistic evaluation of hypotheses concerning

⁴ We are confident that practicable implementations are possible, as recent experiences in the context of the Credential Networks have shown [11].

reliability and public-key authenticity. We conclude with Section 4 by discussing the contributions of our paper and directions for future research.

2 A Model for Reliability and Public-Key Authenticity

We start by recapitulating an existing entity-relationship model [17], which has been conceived for the public-key authentication problem (Subsection 2.1). The relationships defined within this model are used as predicates in the evidential language \mathcal{L} (described in Subsection 2.2), which allows to formalize the terms of public-key authenticity and trust in Subsection 2.3. Finally, we introduce the concept of a trust and authenticity network in Subsection 2.4.

2.1 Entities and Relationships

The model we consider consists of two types of *entities*. A *physical world entity* (*physical entity* for short) is someone that exists in the reality of the “physical” world. Examples are natural persons (human beings) or legal persons (companies, governmental agencies, sports clubs, etc.). A *public-key entity* consists of a public key and a private key, as well as a signature generation and a signature verification algorithm. Access to the private key is needed for generating signatures in the public-key entity’s name.

Finding unique and adequate names for physical entities can be difficult, especially for natural persons. Here we put the naming problem aside and assume that each entity is known under exactly one, unique identifier. We use p_1, \dots, p_m to denote the physical entities of our model, and k_1, \dots, k_n the n public-key entities. The symbol b represents the entity whose reliability is evaluated. Corresponding capital letters refer to variables, i.e., the indexed variable P_i to a physical entity and K_i to a public-key entity. We use \mathcal{P} to denote the set of physical entities, and \mathcal{K} stands for the set of public-key entities. Entities of our model can stand in the following *relationships*:

- A physical entity *controls* a public-key entity whenever she has access to its private key. Access occurs through knowledge of a password or passphrase, through possession of a physical device such as a smartcard, or through a biometric attribute. The same public-key entity can be controlled by more than one physical entity. A physical entity can control more than one public-key entity.
- The relationship *signs* involves a public-key entity and a statement. It holds if there exists a digital signature under the statement, which has been generated by using the private key of the public-key entity. Note that *signs* does not indicate which physical entity is using or controlling a public-key entity.
- The relationship *authors* stands for the fact that it was in a physical entity’s *intention* to be the author of a statement. Authoring a statement can mean to say

it, to write it on a piece of paper, to type it into a computer, or to create any other representation of the statement. The digital signature provides evidence that a physical entity has authored the signed statement.

2.2 An Evidential Language

We use a formal language \mathcal{L} to model assumptions, pieces of evidence, and their logical relationship in the context of trust evaluation. Due to limited space, we cannot provide the exact formal definitions of \mathcal{L} here, but we give at least the basic idea behind \mathcal{L} .

The relationships introduced in the previous subsection are used as predicates in \mathcal{L} , the elements of \mathcal{L} are called \mathcal{L} -formulas. \mathcal{L} is a *many-sorted logic* without function symbols.⁵ “Many-sorted” means that all variables and the arguments of the predicates are of a specified sort. We consider three sorts, namely the physical entities, the public-key entities, and the statements. The atoms of \mathcal{L} are the relationships introduced in the previous subsection; a distinguished predicate symbol is the equality sign. An atomic \mathcal{L} -formula is a predicate symbol together with arguments of appropriate sort. An argument of predicate is either a constant symbol or a variable. Examples of atomic \mathcal{L} -formulas are `controls(p1, k3)`, `controls(P1, k3)`, or `P1 = p2` (note that in the two latter formulas P₁ stands for a *variable* of sort physical entity, and not for a constant). \mathcal{L} contains the usual logical connectives: \wedge (logical and), \vee (logical or), \neg (not), \rightarrow (material implication), \leftrightarrow (bidirectional material implication), and \forall (universal quantifier). In the sequel, let L and the indexed variable L_i stand each for an \mathcal{L} -formula.

2.3 Formalizing Reliability and Public-Key Authenticity

2.3.1 Reliability and Maliciousness

We differentiate among three types of introducers in our model. A physical entity is *reliable* if she is competent and honest; statements authored by a reliable principal are believable. The second type of introducers are those who are *incompetent*. If a statement L is authored by an incompetent introducer, it is impossible to decide whether L is true or false; L can be true by chance, independently of the introducer’s honesty. Therefore statements made by incompetent entities should be simply ignored. The third type of physical entities are the *malicious* introducers. A malicious entity is competent but dishonest, and tries to deceive other physical entities by spreading credentials that contain a false statement. Under the assumption that someone is malicious one can conclude the contrary of what the suspected in-

⁵ Except for constant symbols such as p₁, which can be seen as 0-ary function symbols.

roducer says. We therefore define the reliability (rel) and maliciousness (mal) of a physical entity depending on her competence (comp) and honesty (hon) as follows:⁶

Rule 1: *Reliable physical entity.*

$$\forall.P : (\text{rel}(P) \leftrightarrow (\text{comp}(P) \wedge \text{hon}(P))) \quad (1)$$

Rule 2: *Malicious physical entity.*

$$\forall.P : (\text{mal}(P) \leftrightarrow (\text{comp}(P) \wedge \neg \text{hon}(P))) \quad (2)$$

The logical relationship between a physical entity P's honesty and her competence, as well as the truth of a statement L authored by P are captured by the following two rules. On the one hand, if P is believed to be reliable, the statement L authored by P can be believed. On the other hand, if P is assumed to be malicious, we invert the truth of the uttered statement L:

Rule 3: *Statement authored by a reliable physical entity.*

$$\forall.P \forall.L : ((\text{authors}(P,L) \wedge \text{rel}(P)) \rightarrow L) \quad (3)$$

Rule 4: *Statement authored by a malicious physical entity.*

$$\forall.P \forall.L : ((\text{authors}(P,L) \wedge \text{mal}(P)) \rightarrow \neg L) \quad (4)$$

2.3.2 Public-Key Authenticity

Public-key authenticity of K_1 for P_1 means that P_1 , but no other entity P_2 , controls K . This formally translates into Rule (5).

Rule 5: *Definition of public-key authenticity.*

$$\forall.P_1 \forall.P_2 \forall.K : \text{aut}(P_1,K) \leftrightarrow (\text{controls}(P_1,K) \wedge ((P_1 \neq P_2) \rightarrow \neg \text{controls}(P_2,K))) \quad (5)$$

Because the variables P_1 , P_2 , and K are universally quantified, Rule (5) is valid for all physical entities P_1 and P_2 , as well as all public-key entities K .

Rule (6) formalizes a simplified view of the security of a digital signature scheme: If only P has access to K (i.e., $\text{aut}(P,K)$ holds), and if there is a digital signature under the statement L by K, then P authored the statement L:

Rule 6: *Ascribing digital signatures to physical entities.*

$$\forall.P \forall.K \forall.L : ((\text{aut}(P,K) \wedge \text{signs}(K,L)) \rightarrow \text{authors}(P,L)) \quad (6)$$

⁶ We do not have to define incompetent introducers at this point, since the assumptions that someone is incompetent allows no conclusion.

2.4 Trust and Authenticity Networks

For evaluating a hypothesis concerning the reliability of a physical entity or the authenticity of a public key, a reasoner A takes certain *assumptions* and collects a set of *credentials*. Assumptions and credentials are either with respect to the authenticity of public keys or the reliability of entities. Assumptions are *subjective*; A decides which assumptions are acceptable for her. A credential is a statement which is either digitally signed by a public-key entity or authored by a physical entity.

A's assumptions and credentials form what we call her *Trust and Authenticity Network (TAN)*. A TAN can be depicted by a *multigraph*. We use drawn-through arrows for authenticity assumptions and credentials, similarly to [21, 6, 12]. The graph in Fig. 1 (a) shows A's assumption that k_1 is authentic for p_1 , the graph in Fig. 1 (b) represents the statement that k_2 is authentic for p_2 , and is digitally signed by k_1 . An example of a *negative* authenticity credential (i.e., a statement that a public key is not authentic) is depicted in Figure 1 (c); negative statements are indicated by the negation sign \neg . For the moment, we consider only assumptions and credentials about the `aut` predicate, but it is conceivable to incorporate `controls` statements in a future method. Dashed arrows represent trust assumptions and credentials. Whereas

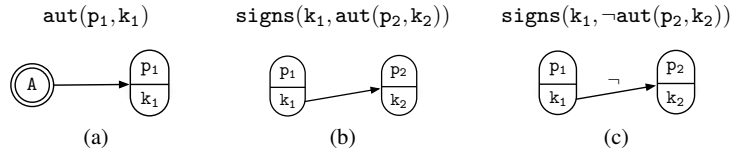


Fig. 1 An authenticity assumption and two authenticity credentials.

A's assumption that p_1 is reliable constitutes a positive trust assumption, her belief that p_1 is incompetent or cheating is negative. We use the following abbreviations: R for a `rel` assumption and credential, I for `incomp`, and M for `mal`. Figure 2 shows

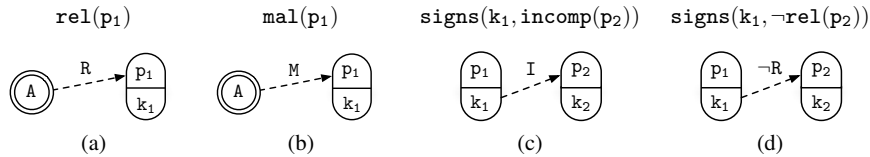


Fig. 2 Two trust assumptions (of entity A) and two trust credentials (digitally signed by the public-key entity k_1).

examples of trust assumptions and credentials. In the graph of Fig. 2 (a), A believes that p_1 is reliable, in Fig. 2 (b) A assumes that p_1 is malicious. The graph in Fig. 2 (c) shows a statement, digitally signed by k_1 , which asserts that p_2 is incompetent. Finally, the graph in Figure 2 (d) provides an example of a negated trust statement: the key owner of k_1 claims that p_2 is not reliable (which is not equal to the statement that p_2 is malicious). TAN-assumptions and credentials can be connected, which results in a multigraph as depicted in Fig. 3.

3 Reasoning about Reliability and Public-Key Authenticity

Our reasoner A is possibly *uncertain* about several of her assumptions; A might doubt the reliability of an introducer; authenticity assumptions can be uncertain if the identification process of an alleged public-key entity owner is error-prone. In analogous manner, credentials can also be uncertain; an introducer can express her uncertainty about an assertion contained within a credential by assigning a *weight* to it. The logical and probabilistic reasoning allows A to evaluate her hypotheses under uncertain assumptions and credentials.

3.1 Logical Reasoning

In this subsection we explain the basic ideas behind *scenarios*, *assumptions* and *arguments*, which are TPA's building blocks for reasoning logically about hypotheses. The definitions coincide to some extent with those provided in [7].

3.1.1 Assumptions

An *assumption* (in the sense of TPA) is a basic unit of concern which is uncertain from the point of view of an entity A. With respect to a TAN, all edges in the multigraph are assumptions, as discussed in the introduction of this section. From a syntactic point of view, an assumption is an \mathcal{L} -formula which consists of a predicate symbol and constant symbols of appropriate sort. In the example of Fig. 3 (a), $\text{aut}(p_1, k_1)$, $\text{aut}(p_2, k_2)$, $\text{rel}(p_1) \text{rel}(p_2)$, $\text{signs}(k_1, \text{aut}(b, k_3))$, and $\text{signs}(k_2, \text{aut}(p_3, k_3))$ are assumptions.

3.1.2 Scenarios

A *scenario* is specified by a truth value assigned to each of the assumptions. Given n assumptions, 2^n scenarios exist. A scenario is denoted by the symbol S . If an assumption A is true (false) in S , we write $S(A) = 1$ ($S(A) = 0$). It is assumed that there

is exactly one scenario which represents the real state of the world. Unfortunately, A does not know which scenario meets this condition.

With respect to a knowledge base (and hence with respect to a given TAN), the set of scenarios can be divided into *conflicting* and *consistent* scenarios. A conflicting scenario stands in contradiction with the knowledge base, a consistent scenario on the other hand is non-conflicting. In Fig. 3 (a), the scenario in which all assumptions hold is conflicting. An informal explanation is the following: from the two signs and aut assumptions we can conclude - by applying Rule (6) - that p_1 authored $\text{aut}(b, k_3)$ and p_2 authored $\text{aut}(p_3, k_3)$. Since p_1 and p_2 are trusted, by using Rule (2) $\text{aut}(b, k_3)$ and $\text{aut}(p_3, k_3)$ can be derived. But Rule (5) asserts that $\text{aut}(b, k_3)$ and $\text{aut}(p_3, k_3)$ cannot hold both at same time. Hence the scenario is conflicting. All other scenarios are consistent with respect to the TAN of Fig. 3.

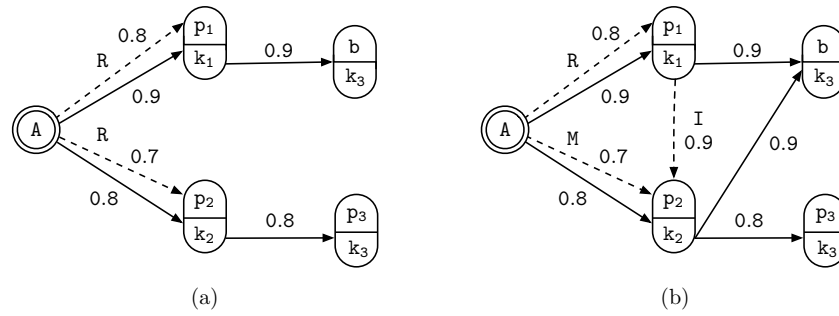


Fig. 3 Two simple TANs.

With respect to a hypothesis h , the set of consistent scenarios can be divided into *supporting*, *refuting*, and *neutral* scenarios [7]. A supporting scenario is a consistent scenario that allows the deduction of h . A refuting scenario is a scenario supporting the counter-hypothesis $\neg h$. A neutral scenario with respect to h is a consistent scenario which is neither supporting nor refuting h . An example of a supporting scenario for $\text{aut}(b, k_3)$ is

$$\begin{aligned} S(\text{aut}(p_1, k_1)) &= 1, & S(\text{rel}(p_1)) &= 1, & S(\text{signs}(k_1, \text{aut}(b, k_3))) &= 1, \\ S(\text{aut}(p_2, k_2)) &= 1, & S(\text{rel}(p_2)) &= 0, & S(\text{signs}(k_2, \text{aut}(p_3, k_3))) &= 1. \end{aligned}$$

The assumptions $\text{aut}(p_1, k_1)$, $\text{signs}(k_1, \text{aut}(b, k_3))$, and $\text{rel}(p_1)$ allow to conclude $\text{aut}(b, k_3)$ (by Rule (5) and Rule (2)). The scenario is not conflicting, since by the assumed falsity of $\text{rel}(p_2)$ the lower “certification path” in Fig. 3 (a) is broken; hence A cannot conclude $\text{aut}(p_3, k_3)$ (otherwise this would lead to a contradiction). An example of a refuting scenario for $\text{aut}(b, k_3)$ is

$$\begin{aligned} S(\text{aut}(p_1, k_1)) &= 1, & S(\text{rel}(p_1)) &= 0, & S(\text{signs}(k_1, \text{aut}(b, k_3))) &= 1, \\ S(\text{aut}(p_2, k_2)) &= 1, & S(\text{rel}(p_2)) &= 1, & S(\text{signs}(k_2, \text{aut}(p_3, k_3))) &= 1. \end{aligned}$$

The scenario is supporting $\text{aut}(p_3, k_3)$, and since we do not accept a public key as authentic for two physical entities, $\neg\text{aut}(b, k_3)$ follows.

$qs(\perp)$
(1) $\text{aut}(p_1, k_1) \wedge \text{aut}(p_2, k_2) \wedge \text{rel}(p_1) \wedge \text{rel}(p_2) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \text{signs}(k_2, \text{aut}(p_3, k_3))$
$sp(\text{aut}(b, k_3))$
(1) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \neg\text{aut}(p_2, k_2)$
(2) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \neg\text{rel}(p_2)$
(3) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \neg\text{signs}(k_2, \text{aut}(p_3, k_3))$
$sp(\neg\text{aut}(b, k_3))$
(1) $\text{aut}(p_2, k_2) \wedge \text{rel}(p_2) \wedge \text{signs}(k_2, \text{aut}(p_3, k_3)) \wedge \neg\text{aut}(p_1, k_1)$
(2) $\text{aut}(p_2, k_2) \wedge \text{rel}(p_2) \wedge \text{signs}(k_2, \text{aut}(p_3, k_3)) \wedge \neg\text{rel}(p_1)$
(3) $\text{aut}(p_2, k_2) \wedge \text{rel}(p_2) \wedge \text{signs}(k_2, \text{aut}(p_3, k_3)) \wedge \neg\text{signs}(k_1, \text{aut}(b, k_3))$

Table 1 $qs(\perp)$, $sp(\text{aut}(b, k_3))$, and $sp(\neg\text{aut}(b, k_3))$ for the TAN of Fig. 3 (a).

3.1.3 Arguments

A compact logical representation of scenarios is achieved by means of *arguments*. Technically, an *argument* is a conjunction of assumption literals. There are conflicting, supporting and refuting arguments, analogously to the different types of scenarios. The expression $qs(\perp)$ represents the set of minimal conflicting assumptions; $sp(h)$ and $sp(\neg h)$ stand for the sets of minimal arguments supporting and refuting h , respectively.

The arguments of the TANs discussed in this paper have been determined by translating first the TAN into a *Propositional Argumentation System* (i.e., a knowledge base in which all variables have been instantiated and the universal quantifiers have been removed). The so-obtained propositional knowledge base was implemented in ABEL [10], a framework for evaluating propositional knowledge bases qualitatively and quantitatively.

3.1.4 Examples

Table 1 shows the minimal argument sets for the example depicted Fig. 3 (a). As mentioned, there is only one conflicting scenario. Hence we have only one conflicting argument containing all assumptions. The common part of the supporting arguments for $\text{aut}(b, k_3)$ are the three assumptions of the upper certification path of our example. The assumptions $\neg\text{aut}(p_2, k_2)$, $\neg\text{rel}(p_2)$, and $\neg\text{signs}(k_2, \text{aut}(p_3, k_3))$ all guarantee that the argument is not conflicting. Each argument stands for four

scenarios (because there are two missing assumptions in each argument supporting $\text{aut}(b, k_3)$). The supporting arguments for $\neg\text{aut}(b, k_3)$ are in a certain sense symmetric to the arguments for $\text{aut}(b, k_3)$. They actually correspond to the supporting arguments for $\text{aut}(p, k_3)$. Note that Table 1 lists only the *minimal* arguments. For example, the argument

$$\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \neg\text{aut}(p_2, k_2) \wedge \text{rel}(p_2)$$

supports also $\text{aut}(b, k_3)$, but is contained in argument (1) of $sp(\text{aut}(b, k_3))$ of Table 1.

$qs(\perp)$
(1) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{incomp}(p_2)) \wedge \text{mal}(p_2)$
(2) $\text{aut}(p_1, k_1) \wedge \text{aut}(p_2, k_2) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \supset$ $\text{mal}(p_2) \wedge \text{signs}(k_2, \text{aut}(b, k_3))$
$sp(\text{aut}(b, k_3))$
(1) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \neg\text{mal}(p_2)$
(2) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \supset$ $\neg\text{signs}(k_1, \text{incomp}(p_2)) \wedge \neg\text{signs}(k_2, \text{aut}(b, k_3))$
(3) $\text{aut}(p_1, k_1) \wedge \text{rel}(p_1) \wedge \text{signs}(k_1, \text{aut}(b, k_3)) \wedge \supset$ $\neg\text{signs}(k_1, \text{incomp}(p_2)) \wedge \neg\text{aut}(p_2, k_2)$
$sp(\neg\text{aut}(b, k_3))$
(1) $\text{aut}(p_2, k_2) \wedge \text{mal}(p_2) \wedge \text{signs}(k_2, \text{aut}(b, k_2)) \wedge \neg\text{aut}(p_1, k_1)$
(2) $\text{aut}(p_2, k_2) \wedge \text{mal}(p_2) \wedge \text{signs}(k_2, \text{aut}(b, k_2)) \wedge \neg\text{rel}(p_1)$
(3) $\text{aut}(p_2, k_2) \wedge \text{mal}(p_2) \wedge \text{signs}(k_2, \text{aut}(b, k_2)) \wedge \supset$ $\neg\text{signs}(k_1, \text{aut}(b, k_3)) \wedge \neg\text{signs}(k_1, \text{incomp}(p_2))$

Table 2 $qs(\perp)$, $sp(\text{aut}(b, k_3))$, and $sp(\neg\text{aut}(b, k_3))$ for the TAN of Fig. 3 (b).

Figure 3 (b) shows an example which is more complicated. In contrast to the previous example, A believes that p_2 is malicious. There is an additional digitally signed trust credential which claims incompetence for p_2 . The owner of k_2 provides conflicting information, as she claims simultaneously public-key authenticity of k_3 for b and p_3 .

The qualitative evaluation provides some interesting insights: The second conflicting argument in Table 2 is equal to the only conflicting argument of our first example. Argument (1) of $qs(\perp)$ conflicts with the given TAN because $\text{aut}(p_1, k_1)$, $\text{rel}(p_1)$, and $\text{signs}(k_1, \text{incomp}(p_2))$ allow to conclude that p_2 is incompetent. This, however, conflicts with the assumption $\text{mal}(p_2)$, which stands for the assumption that p_2 is competent (and dishonest).

All supporting arguments for $\text{aut}(b, k_3)$ in the second TAN contain the three assumptions of the upper certification path. Again, some negated assumptions have to be added to guarantee the consistency of the supporting arguments. For example, the first supporting argument contains the literal $\neg\text{mal}(p_2)$. By adding this assumption,

a contradiction can be prevented. Note that - in contrast to the previous example - there are no supporting arguments for $\text{aut}(p_3, k_3)$. If p_2 is indeed *not* malicious, she is incompetent *or* honest. From this clause $\text{aut}(p_3, k_3)$ can *not* be deduced.

3.2 Probabilistic Reasoning

The idea of the probabilistic part of TPA (and hence of our method) is that from A's point of view each scenario corresponds with a certain probability to the real state of the world. A has to choose the probabilities such that the sum of the probabilities assigned to the scenarios equals one. Given the exponential number of scenarios, it is infeasible for A to estimate the probability of each single scenario. It is often justifiable to consider the statements of a TAN as being stochastically independent. In this case, A assigns a probability to all of her authenticity and trust assumptions. The weights assigned to each credential are represented by a probability, too. Under

	$dqs(\perp)$	$dsp(h_1)$	$dsp(\neg h_1)$	$dsp(h_2)$	$dsp(\neg h_2)$
TAN of Fig. 3(a)	0.290	0.504	0.222	0.222	0.504
TAN of Fig. 3(b)	0.486	0.403	0.282	0.000	0.654

Table 3 Qualitative evaluation of TANs: $h_1 = \text{aut}(b, k_3)$, $h_2 = \text{aut}(p_3, k_3)$.

the independence assumption, the probability of a scenario can be computed as the product of the marginal probabilities of the assumptions.

Formally, let A_i stand for the i th assumption, and let p_i be the probability attached to A_i . Given a scenario S , let S^+ denote the assumptions which are positive in S , and S^- the assumptions that occur negatively:

$$S^+ = \{A \in S \mid S(A) = 1\}, \quad S^- = \{A \in S \mid S(A) = 0\}.$$

The probability $P(S)$ of scenario S is then defined as

$$P(S) = \prod_{A_i \in S^+} p_i \cdot \prod_{A_i \in S^-} (1 - p_i).$$

The *degree of conflict*, denoted by $dqs(\perp)$, is obtained by summing up the *probabilities of the conflicting scenarios*. It is a measure of how conflicting the assumptions are with respect to the knowledge base (i.e., the TAN). Let $dqs(h)$ stand for the sum of $dqs(\perp)$ and the probabilities of all the scenarios supporting h (i.e., the sum of the probabilities of all scenarios allowing the deduction of h , including the conflicting ones). The *degree of support* for the hypothesis h , denoted by $dsp(h)$, is the probability that h can be derived, provided that the real scenario is not conflicting.

Formally, the degree of support corresponds to

$$dsp(h) = \frac{dqs(h)}{1 - dqs(\perp)}.$$

Table 3 shows the degrees of support for our two examples of Fig. 3. In example of Fig. 3 (a), $\text{aut}(b, k_3)$ and $\neg\text{aut}(p_3, k_3)$ are quite probable. In the TAN of Fig. 3 (b), $\text{aut}(b, k_3)$ is less probable, but we have no evidence for $\text{aut}(p_3, k_3)$. In both cases, A either accepts $\text{aut}(b, k_3)$ or collects additional evidence to gain more certainty for or against the validity of $\text{aut}(b, k_3)$. A discussion of how to validate a hypothesis based on $dsp(h)$ and $dsp(\neg h)$ can be found in [8].

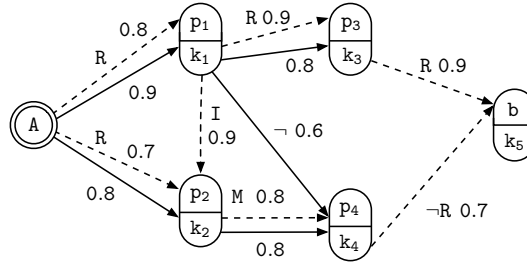


Fig. 4 A more complex TAN.

We end the presentation of our probabilistic method by analyzing the more complicated TAN depicted in Figure 4. The TAN involves most trust and authenticity statements introduced in this paper. It contains a negative authenticity statement ($\text{signs}(k_1, \neg\text{aut}(p_4, k_4))$), and a negative trust statement ($\text{signs}(k_3, \neg\text{rel}(b))$). We are interested in the evaluation of the reliability of b , i.e., the hypothesis of interest is $h = \text{rel}(b)$. Although the TAN is not that large, it has already a considerable complexity: there are twelve edges in the TAN, and hence $2^{12} = 4096$ possible scenarios. The qualitative evaluation yields two conflicting and seven supporting arguments (which we do not write down here). One of the two conflicting arguments corresponds to

$$\begin{aligned} & \text{aut}(p_1, k_1) \wedge \text{aut}(p_2, k_2) \wedge \text{rel}(p_1) \wedge \text{rel}(p_2) \wedge \triangleright \\ & \text{signs}(k_1, \neg\text{aut}(p_4, k_4)) \wedge \text{signs}(k_2, \text{aut}(p_4, k_4)). \end{aligned}$$

The above argument is conflicting, since it allows the deduction of $\text{aut}(p_4, k_4)$ and $\neg\text{aut}(p_4, k_4)$. The degree of conflict $dqs(\perp)$ is quite high and is approximately 0.473; the degree of support $dsp(\text{rel}(b))$ is roughly 0.442. Interestingly, there is no argument for the hypothesis that b is not reliable ($dsp(\text{rel}(b)) = 0$).

4 Conclusion

We have introduced a trust evaluation method, which can also be used for authenticating public keys. The used, extended model considers the possibility that a public-key entity is shared by different physical entities, and that a physical entity controls several public-key entities at the same time. Negative statements are an integral part of the method. Reliability is decomposed into honesty and competence, which allows to differentiate between incompetent and malicious physical entities. The assumptions and the available evidence from the perspective of a physical entity A can be represented by an evidential language and by a multigraph. We make use of the Theory of Probabilistic Argumentation, which allows to cope with conflicting assumptions. TPA provides logical arguments *and* probabilities of derivability for and against the hypotheses in question.

Future work consists in investigating the applicability of our method in concrete systems, and in devising specific algorithms for the evaluation of trust and authenticity networks. Possible extensions of the model are the inclusion of trust scopes and time aspects, as well as modeling the revocation of statements and public keys.

Acknowledgements This research was supported by the Swiss National Science Foundation, Project No. PP002-102652/1.

References

1. The official ebay website. <http://www.ebay.com>, April 2006.
2. T. Beth, M. Borchering, and B. Klein. Valuation of trust in open networks. In *ESORICS'94, 3rd European Symposium on Research in Computer Security*, LNCS 875, pages 3–18. Springer, 1994.
3. M. Burrows, M. Abadi, and R. Needham. A logic of authentication. *ACM Transactions on Computer Systems*, 8(1):18–36, February 1990.
4. W. Diffie and M. E. Hellman. New directions in cryptography. *IEEE Transactions on Information Theory*, IT-22(6):644–654, 1976.
5. L. Gong, R. Needham, and R. Yahalom. Reasoning About Belief in Cryptographic Protocols. In Deborah Cooper and Teresa Lunt, editors, *Proceedings 1990 IEEE Symposium on Research in Security and Privacy*, pages 234–248. IEEE Computer Society, 1990.
6. R. Haenni. Using probabilistic argumentation for key validation in public-key cryptography. *International Journal of Approximate Reasoning*, 38(3):355–376, 2005.
7. R. Haenni. Probabilistic argumentation (submitted). *Elsevier*, 2007.
8. R. Haenni, J. Jonczy, and R. Kohlas. Two-layer models for managing authenticity and trust. In R. Song, L. Korba, and G. Yee, editors, *Trust in E-Services: Technologies, Practices and Challenges*. 2006.
9. R. Haenni, J. Kohlas, and N. Lehmann. Probabilistic argumentation systems. In D. M. Gabbay and P. Smets, editors, *Handbook of Defeasible Reasoning and Uncertainty Management Systems*, volume 5: Algorithms for Uncertainty and Defeasible Reasoning, pages 221–288. Kluwer Academic Publishers, Dordrecht, Netherlands, 2000.
10. R. Haenni and N. Lehmann. ABEL: an interactive tool for probabilistic argumentative reasoning. In *ECSQARU'03, 7th European Conference on Symbolic and Quantitative Approaches to Reasoning under Uncertainty*, pages 588–593, Aalborg, Denmark, 2003.

11. J. Jonczyk. Evaluating trust and authenticity with CAUTION. In *iTrust'06, 4rd International Conference on Trust Management*, pages 449–453, Pisa, Italy, 2006.
12. J. Jonczyk and R. Haenni. Credential networks: a general model for distributed trust and authenticity management. In A. Ghorbani and S. Marsh, editors, *PST'05: 3rd Annual Conference on Privacy, Security and Trust*, pages 101–112, St. Andrews, Canada, 2005.
13. J. Jonczyk, M. Wüthrich, and R. Haenni. A probabilistic trust model for GnuPG. In *23C3, 23rd Chaos Communication Congress*, pages 61–66, Berlin, Germany, 2006.
14. A. Jøsang. An algebra for assessing trust in certification chains. In *NDSS'99: 6th Annual Symposium on Network and Distributed System Security*, San Diego, USA, 1999.
15. A. Jøsang. A logic for uncertain probabilities. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 9(3):279–311, 2001.
16. J. Kohlas and P. A. Monney. *A Mathematical Theory of Hints. An Approach to the Dempster-Shafer Theory of Evidence*, volume 425 of *Lecture Notes in Economics and Mathematical Systems*. Springer, 1995.
17. R. Kohlas, R. Haenni, and J. Jonczyk. A new model for public-key authentication. In T. Braun, G. Carle, and B. Stiller, editors, *KiVS'07, Kommunikation in Verteilten Systemen*, pages 213–224, Berne, Switzerland, 2007.
18. R. Kohlas and U. Maurer. Confidence valuation in a public-key infrastructure based on uncertain evidence. In H. Imai and Y. Zheng, editors, *PKC'2000, Third International Workshop on Practice and Theory in Public Key Cryptography*, LNCS 1751, pages 93–112, Melbourne, Australia, 2000. Springer.
19. R. Levien and A. Aiken. Attack-resistant trust metrics for public key certification. In *7th on USENIX Security Symposium*, pages 229–242, 1998.
20. G. Mahoney, W. Myrvold, and G. C. Shojja. Generic reliability trust model. In A. Ghorbani and S. Marsh, editors, *PST'05: 3rd Annual Conference on Privacy, Security and Trust*, pages 113–120, St. Andrews, Canada, 2005.
21. U. Maurer. Modelling a public-key infrastructure. In E. Bertino, H. Kurth, G. Martella, and E. Montolivo, editors, *ESORICS, European Symposium on Research in Computer Security*, LNCS 1146, pages 324–350. Springer, 1996.
22. M. K. Reiter and S. G. Stubblebine. Path independence for authentication in large-scale systems. In *CCS'97, 4th ACM Conference on Computer and Communications Security*, pages 57–66, Zürich, Switzerland, 1997. Academic Press.
23. M. K. Reiter and S. G. Stubblebine. Toward acceptable metrics of authentication. In *SP'97: 18th IEEE Symposium on Security and Privacy*, pages 10–20, Oakland, USA, 1997.
24. R. L. Rivest, A. Shamir, and L. M. Adelman. A method for obtaining digital signatures and public-key cryptosystems. Technical Report TM-82, MIT, Cambridge, USA, 1977.
25. C. G. Zarba. Many-sorted logic. <http://theory.stanford.edu/~zarba/snow/ch01.pdf>.
26. P. R. Zimmermann. *PGP User's Guide Volume I: Essential Topics*, 1994.