

Towards a Theory of Accountability and Audit

Radha Jagadeesan^{1,*}, Alan Jeffrey², Corin Pitcher¹, and James Riely^{1,*}

¹ School of Computing, DePaul University

² Bell Labs, Alcatel–Lucent

Abstract. Accountability mechanisms, which rely on after-the-fact verification, are an attractive means to enforce authorization policies. In this paper, we describe an operational model of accountability-based distributed systems. We describe analyses which support both the design of accountability systems and the validation of auditors for finitary accountability systems. Our study provides formal foundations to explore the tradeoffs underlying the design of accountability systems including: the power of the auditor, the efficiency of the audit protocol, the requirements placed on the agents, and the requirements placed on the communication infrastructure.

1 Introduction

The context of our paper is authorization in distributed systems. The attackers that we consider are untrustworthy principals running arbitrary programs on the network. Attackers may not respect the policies of a system; for example, attackers may create authorization objects without actually having the rights to create them, aiming to subvert the global authorization policy. Traditionally, authorization policies are enforced by controls imposed before shared resources are accessed.

Recently, there has been great interest in accountability mechanisms that rely on after-the-fact verification (Weitzner et al. 2007). In this approach, audit logs record vital systems information and an auditor uses these logs to identify dishonest principals and to assign blame when there has been a violation of security policy. The fear of being “caught” helps to achieve security by deterrence, in the spirit of traditional law enforcement and organizational security. Accountability plays a critical role in the development of trust during human interaction (Friedman and Grudin 1998). Thus, accountability is viewed both as a tool to achieve practical security (Lampson 2004) and as a first-class design goal of services in federated distributed systems (Yumerefendi and Chase 2004).

While designing for accountability is subtle in general (Eriksén 2002), mechanisms to instrument systems to support accountability have been explored in several specific applications: determinate distributed systems (Haeberlen et al. 2007), network storage (Yumerefendi and Chase 2007), validating ISP quality of service claims (Argyrazi et al. 2007), internet protocol (Andersen et al. 2008) and policy enforcement on shared documents (Etalle and Winsborough 2007).

In comparison to *a priori* approaches such as access-control, however, the accountability approach to security lacks general foundations for models and programming.

* Supported by NSF Career 0347542.

Citing a small sample of references, access control has (a) *operational* models in the form of automata (Schneider 2000), with associated algebraic models based on regular expressions (Abadi et al. 2005); (b) *logic-based declarative* approaches in a fragment of many-sorted first-order predicate logic (Halpern and Weissman 2003; Li and Mitchell 2003); and (c) *static analysis* to validate the access-control properties of interfaces, e.g., types for authorization (Fournet et al. 2005; Cirillo et al. 2008).

In this paper we make two contributions toward bringing such formal foundations to the study of accountability. First, we describe an operational model of accountability based systems. Honest and dishonest principals are described as agents in a distributed system where the communication model guarantees point-to-point integrity and authenticity. Auditors and other trusted agents (such as trusted third parties) are also modeled internally as agents. Behaviors of all agents are described as processes in a process algebra with discrete time. Auditor implementability is ensured by forcing auditor behavior to be completely determined by the messages that it receives.

Second, we describe analyses to support the design of accountability systems and the validation of auditors for finitary systems (those with finitely many principals running finite state processes with finitely many message kinds). We compile finitary systems to (turn-based) games and use alternating temporal logic to specify the properties of interest. This permits us to adapt existing model-checking algorithms for verification.

Our results provide the foundations necessary to explore tradeoffs in the design of mechanisms that ensure accountability. The potentially conflicting design parameters include the efficiency of the audit, the amount of logging, and the required use of message signing, watermarking, or trusted third parties. Design choices place constraints on the auditor, the agents of the system and the underlying communication infrastructure.

The paper is organized as follows. We motivate our approach in Section 2. Section 3 describes the model and Section 4 describes the analysis framework. The ideas are illustrated using examples in Section 5. We survey related work in Section 6. In this extended abstract, we elide all proofs.

2 Overview of Our Approach

In this section, we illustrate the motivations behind the design of our framework using variants of a motivating example from (Barth et al. 2007).

In Section 5, we analyze an abstract variant of the example that permits message forwarding amongst health professionals. Our analysis yields a variety of auditors for the example, even in general distributed settings, and shows that powerful mechanisms, such as trusted third parties, are not necessary for all audit protocols.

Example 1 (My Health). The MyHealth patient portal at Vanderbilt University Hospital allows patients to interact with healthcare professionals through a web based system. There are three possible roles that can be assumed by principals: health professionals (doctors and nurses), non-health professionals (secretaries), and patients. The possible messages include health questions from patients and health answers from doctors. We focus on the two privacy policies in Barth et al. (2007): (a) a health question can only be directed to a health professional, and (b) a health answer about a patient can only be directed to the same patient or to a health professional. These policies permit health

professionals to forward health information amongst themselves. In the discussion below, we will consider the case where patient Charlie contacts the auditor because he has received a health answer from doctor Bob that was intended for a different patient. The motivation for such an audit is to aid in the detection and discovery of the source of the leak. \square

We now describe our model and its relation to the following properties. The discussion is intended to establish intuitions, with formalities deferred to later sections.

- *Upper bound*: Every agent guilty of a dishonest action is blamed by the auditor.
- *Lower bound*: Everyone blamed by the auditor is guilty.
- *Overlap*: At least one of the agents blamed by the auditor is guilty.
- *Liveness*: The auditor is always successful in blaming a non-empty subset of agents.
- *Blamelessness*: Honest agents have a strategy to avoid being pronounced a *possible* offender by an auditor.

Agents. We model the behavior of principals (both honest and dishonest) as agents in a distributed system. Auditors are also modeled internally as honest agents. We use processes to specify an upper bound on honest behavior: a principal is behaving honestly in a run whenever their contribution to the run is a trace of an honest process. A dishonest agent is unconstrained. A run of an agent reveals its dishonesty if it is not a permitted trace for an honest agent.

The communication model captures point-to-point communication over an underlying secure communication mechanism which provides integrity and authenticity guarantees, but provides no additional mechanisms for non-repudiation or end-to-end security. This model is realizable using transport mechanisms such as TLS.

Dishonest agents may collaborate arbitrarily. This means that the auditor has to achieve its objectives independent of potential cartels of dishonest agents. Honest agents may also collaborate, depending upon the specification of honest agents.

Internal auditors. Auditors are intended to be realizable agents in a distributed system without global knowledge. Thus, they are unaware of transactions that do not involve them, and their local state is only influenced by the messages that they receive. In contrast, the strategies adopted by dishonest agents can potentially depend on viewing traffic on the network between other agents. The internalization of auditors limits them. Auditors can only address dishonest behaviors using the information available on individual runs of a system; they cannot audit violations of security properties that need sets of traces for their specification (such as non-interference). Auditors cannot detect cartels of dishonest agents who conduct dishonest exchanges amongst themselves.

Thus, our auditors cannot in general satisfy *Upper bound*. To see this, consider the leakage of patient records to a dishonest non-health professional by a dishonest health professional via out-of-band mechanisms without using the MyHealth website in Example 1. Such leakage of records by dishonest agents solely to dishonest agents will not be detected at all by an auditor in our framework.

Mandatory logging and responsiveness. Even in the case that the auditor has become aware of dishonest behavior and initiates an audit, the auditor is powerless unless there are statutory and enforceable reporting requirements on the honest agents.

In Example 1, if there is no requirement for maintaining and presenting records, doctor Bob can achieve “absence of provable guilt” by maintaining no records. Such

reasoning motivates requirements on honest agents to maintain audit logs in several accountability systems. Furthermore, a guarantee that honest principals provide answers to audit queries is needed for the auditor to achieve *Liveness*.

These desiderata motivate the inclusion of time in our system specification formalism to enable systems to mandate promptness on honest agents. Thus, auditors can use tardiness as evidence for dishonesty and assign blame to such tardy agents. Our model uses *discrete time*, which is abstract and logical rather than quantitative.

Communication model. The absence of non-repudiation in our basic communication model limits the accuracy of the audit process to *Overlap*. For example, in the audit scenario of Example 1, the auditor commences by querying doctor Bob: if Bob disagrees that he sent the message to patient Charlie, the auditor can deduce that at least one of Bob or Charlie is compromised. The absence of non-repudiation prevents further disambiguation. Alternately, Bob might point to another principal, Eve, as the sender of the patient health message. In this case, the auditor proceeds to question Eve. This process either ends in one of two ways. (a) The auditor discovers two principals (perhaps one of whom is a health professional) who disagree on messages sent by one and received by the other, as sketched above; in this case, both the principals are deemed guilty. (b) The auditor discovers a cycle of non-health professionals, each claiming to have received the message from the predecessor in the cycle; in this case, the entire cycle of principals is deemed guilty. In either case, the auditor achieves *Overlap*.

This situation may be unsatisfactory to an honest agent, since it is not possible for an honest agent to achieve *Blamelessness*. In addition the auditor cannot achieve stronger properties for the auditor, such as *Lower bound*. Such properties require more detailed and secure logging of messages.

We do not limit attention to strong models of communication—such as those enabling non-repudiation—because weaker models are often more realistic. For example, in the IETF Session Initiation Protocol (SIP), a typical SIP proxy is expected to handle large volume of calls; thus, it is difficult to successfully mandate computationally expensive signature based methods on each point-to-point communication link.

As evidence for the flexibility of our modeling, we show that our model can indeed encode notaries as trusted third parties. This permits us to address the stronger communication guarantees required to accurately capture examples such as the MyHealth website of Example 1.

3 Formalizing the Model

Based on a notion of process, defined below, we will define an arena $\langle \mathcal{A}, \mathcal{M}, \mathcal{H} \rangle$ to be a set \mathcal{A} of principals, a set \mathcal{M} of messages, and a set \mathcal{H} of processes, which define the honest behaviors of agents. Later, we shall give example arenas, and then give desirable properties of auditors in such an arena.

Our formal model is based on Communicating Sequential Processes (Brookes et al. 1984), I/O automata (Lynch 2003), and discrete timed process algebra (Hennessy and Regan 1995). Our processes are *input-enabled*, to prevent a (perhaps dishonest) agent from blocking the output of other agents. We use *discrete time*, and the *timeouts* that it engenders, to specify conditions on prompt response. Our communication model provides integrity and authenticity guarantees but provides no additional

mechanisms for non-repudiation or end-to-end security. We use processes as a *safety* specification of honest behavior: a principal is behaving honestly in a run whenever their contribution to the run is a trace of an honest process.

Actions. Fix a countable set \mathcal{A} of *principals* and a countable set \mathcal{M} of *messages*. Let a, b, c, d, h range over elements of \mathcal{A} ; A, B, C, D, H over subsets of \mathcal{A} ; and m over elements of \mathcal{M} .

The set of *actions* \mathcal{H} over $\langle \mathcal{A}, \mathcal{M} \rangle$ is then generated by the grammar

$$k, \ell ::= a \rightarrow b : m \mid \sigma$$

where $a \rightarrow b : m$ represents a message m sent from a to b , and σ represents a timeout.

Relative to a set of principals A , an action may be output, input, internal, disjoint or timeout. The action $(a \rightarrow b : m)$ is *output* from A if $a \in A$ and $b \notin A$, *input* to A if $a \notin A$ and $b \in A$, *internal* to A if $a \in A$ and $b \in A$, and *disjoint* from A if $a \notin A$ and $b \notin A$. The action σ is timeout from A , for any A .

We often describe actions from the point of view of a particular principal, using $?$ for inputs and $!$ for outputs. Thus, when giving the example of a process for a , we will write $a \rightarrow b : m$ as $a \rightarrow b ! m$ and $b \rightarrow a : m$ as $b \rightarrow a ? m$.

Processes. A *process* over $\langle \mathcal{A}, \mathcal{M} \rangle$ is a quadruple $P = \langle A, S, s_0, \rightarrow \rangle$ where (a) $A \subseteq \mathcal{A}$ is a subset of principals, (b) S is a set of states, ranged over by s and t , (c) $s_0 \in S$ is a distinguished start state, (d) $\rightarrow \subseteq S \times \mathcal{H} \times S$ is a labeled transition relation in which labels are actions over $\langle \mathcal{A}, \mathcal{M} \rangle$. We call A the *principals of P* , written $\pi(P)$.

We say that s allows k whenever there exists a t such that $s \xrightarrow{k} t$. We also require that no label in \rightarrow is disjoint from A , every state in S allows every input for A (input-enabling), every state in S allows at least one timeout or output for A (timeout-enabling).

Whenever A and B are disjoint we define the *composition* of processes $P = \langle A, S, s_0, \rightarrow_1 \rangle$ and $Q = \langle B, T, t_0, \rightarrow_2 \rangle$ to be $P \parallel Q = \langle A \cup B, S \parallel T, s_0 \parallel t_0, \rightarrow \rangle$ where $S \parallel T = \{(s \parallel t) \mid s \in S \text{ and } t \in T\}$ and \rightarrow is defined as follows:

$$\begin{array}{c} \frac{s \xrightarrow{k_1} s'}{s \parallel t \xrightarrow{k} s' \parallel t} \quad k \text{ is disjoint from } B \qquad \frac{t \xrightarrow{k_2} t'}{s \parallel t \xrightarrow{k} s \parallel t'} \quad k \text{ is disjoint from } A \\ \frac{s \xrightarrow{k_1} s' \quad t \xrightarrow{k_2} t'}{s \parallel t \xrightarrow{k} s' \parallel t'} \quad \begin{array}{l} k \text{ is input to } A \text{ and output from } B, \text{ or} \\ k \text{ is input to } B \text{ and output from } A, \text{ or} \\ k \text{ is } \sigma \end{array} \end{array}$$

We write $\prod_{i \in I} P_i$ for the composition of processes P_i .

A *trace*, v, w , is a finite sequence of actions. Write ε for the empty trace, and $v.w$ for trace composition. Write $s \xRightarrow{v} s'$ to indicate that there exists a sequence of transitions from s to s' labeled by v . A trace has principals A whenever it contains no action disjoint from A . A *run* of a process P with start state s_0 is a trace v such $s_0 \xRightarrow{v} s'$ for some s' . Note that any run of a process with principals A must have principals A .

Write $v \upharpoonright A$ for the projection of v onto actions relating to A :

$$\begin{array}{lll} \varepsilon \upharpoonright A = \varepsilon & v.(b \rightarrow c : m) \upharpoonright A = v \upharpoonright A & \text{if } A \cap \{b, c\} = \emptyset \\ v.\sigma \upharpoonright A = v \upharpoonright A.\sigma & v.(b \rightarrow c : m) \upharpoonright A = v \upharpoonright A.(b \rightarrow c : m) & \text{if } A \cap \{b, c\} \neq \emptyset \end{array}$$

Note that for any P with principals A and Q with principals B , ν is a run of $P \parallel Q$ whenever ν has principals $A \cup B$, $\nu|A$ is a run of P , and $\nu|B$ is a run of Q . This is the usual trace semantics of parallel composition in CSP (Brookes et al. 1984).

Arenas. An arena $\langle \mathcal{A}, \mathcal{M}, \mathcal{H} \rangle$ comprises a countable set \mathcal{A} of principals, a countable set \mathcal{M} of messages, and a countable set \mathcal{H} of process over $\langle \mathcal{A}, \mathcal{M} \rangle$.

Given a principal set H and a process P , we say that H is *honest* in P if $P = \prod_{i \in I} P_i$ and for all $h \in H$ there exists $P_i \in \mathcal{H}$ such that $\pi(P_i) = \{h\}$, that is, if every honest principal must be represented by an honest process.

We note that honesty for processes is down-closed (if $H \supseteq H'$ and H is honest in P then H' is honest in P) and union-closed (if H and H' are honest in P then so is $H \cup H'$); so any process has a maximum honest set of principals.

Given a principal set H and a trace ν , we say that H is *honest* in ν whenever there exists a process P with run ν such that H is honest in P . Honesty for traces is down-closed and union-closed; so any trace has a maximum honest set of principals.

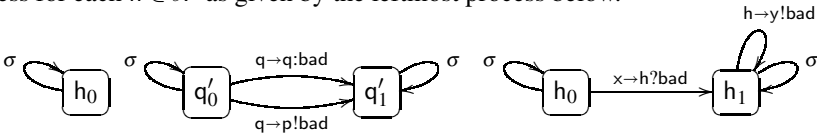
Honesty is a global property of traces, that is for any trace ν capturing all behavior of the system, we can determine the principals H who have behaved honestly in that trace. The problem of audit is that the auditor is not provided with the trace ν but only a local fragment of ν .

Example Arenas. In the following examples, we write principal names without subscripts or superscripts (p, h, b) and write states belonging to a principal with numeric subscripts and optional primes (p_0, h''_5, b'_2) with the convention that p_i is a state of an honest process with principals $\{p\}$, and that p_0 is the start state of the process.

Honesty and dishonesty are properties of principals with respect to a trace, rather than an intrinsic property of a principal. Nonetheless, we find it helpful to use suggestive names in examples to indicate principals that are intended to be honest or dishonest. We use p, q, r for general principals; h, g, f for principals with honest behaviors; and d, c, b for principals with dishonest behaviors. We also use x, y, z for parameters and a for auditors, as discussed below.

We elide transitions required solely for input enabling, assuming an implicit transition $p_i \xrightarrow{k} p_i$ for any input action of p that is not explicitly given.

Example 2. Consider an arena with $\mathcal{A} = \{p, q\}$ and $\mathcal{M} = \{\text{bad}\}$, and define an honest process for each $h \in \mathcal{A}$ as given by the leftmost process below.



The initial state $p_0 \parallel q_0$ shows that both p and q are honest in any trace containing only timeouts.

p is also honest in the traces $q \rightarrow q:\text{bad}$ (because the action is disjoint with p) and $q \rightarrow p:\text{bad}$ (because p is input enabled), as witnessed by the initial state $p_0 \parallel q'_0$ where q'_0 is defined by the center process above. q is dishonest in these traces, since there is no process that is honest for q which allows them. Moreover, q is dishonest in any trace containing $q \rightarrow q:\text{bad}$ or $q \rightarrow p:\text{bad}$. Symmetrically, q is honest in $p \rightarrow p:\text{bad}$ and $p \rightarrow q:\text{bad}$, whereas p is dishonest in these traces.

Auditing in this arena is trivial, since the sender of a bad message is guaranteed to be dishonest. This example corresponds to the case in Example 1 when no principal is allowed to forward health answers to anyone except the patient in question.

There is a problem with initiating audit, however, in that honest agents have no mechanism for reporting dishonest behavior, for example p cannot report the receipt of the message $q \rightarrow p:\text{bad}$ to an auditor. \square

Next, we model the message forwarding capabilities of principals in Example 1.

Example 3. In the variant given by the rightmost process above, honest processes are allowed to forward the first bad message that they receive; for example, reporting dishonest behavior to an auditor. If an auditor knows that a message $p \rightarrow q:\text{bad}$ has been sent, then there must be a dishonest principal, but does not know who is dishonest – there are traces containing $p \rightarrow q:\text{bad}$ in which p is honest, or q is honest, or both. The goal of the auditor should be to determine the agent that *initiated* the bad message. \square

Auditors. An *arena with audit* is an arena with a distinguished honest principal a and a set of distinguished messages blame B for every $B \subseteq \mathcal{A}$, indicating the blame set B . For simplicity, we treat the blame action as internal to the auditor, and thus we abbreviate the action “ $a \rightarrow a:\text{blame } B$ ” as “ $a:\text{blame } B$ ”.

We now consider various notions of correctness for auditors. Many of these notions, while appealing, have serious technical problems, and so we will not consider them further in this paper.

In these definitions, we will discuss a trace with dishonest principals D , defined to be $\mathcal{A} \setminus H$ where H is the largest honest set.

Candidate 1 (Upper bound). An *arena with audit* provides an upper bound on dishonesty if, for any trace v with dishonest $D \not\subseteq a$ containing $a:\text{blame } B$, we have $D \subseteq B$.

Unfortunately, the only auditor capable of providing an upper bound on dishonesty is one which blames all principals who are capable of dishonesty, regardless of whether they acted dishonestly or not.

$A \subseteq \mathcal{A}$ are said to be *capable of dishonesty* in an arena whenever there is a trace v internal to A (i.e., all messages from $a \in A$ are sent to some $b \in A$) with dishonest A .

Proposition 4. In any arena where audit provides an upper bound on dishonesty, and where $A \not\subseteq a$ are capable of dishonesty, we have that any trace containing $a:\text{blame } B$ must have $A \subseteq B$.

We do not consider this notion of correctness further.

Candidate 2 (Lower bound). An *arena with audit* provides a lower bound on dishonesty if, for any trace v with dishonest $D \not\subseteq a$ containing $a:\text{blame } B$, we have $B \subseteq D$.

Unfortunately, auditors are only capable of blaming dishonest principals who confessed their own dishonesty. Dishonest principals who do not confess will never be blamed.

In a trace v with dishonest $D \ni d$, we say that d *confessed* whenever, for any w such that $v \upharpoonright \{d, a\} = w \upharpoonright \{d, a\}$ we have that w has dishonest $D' \ni d$.

Proposition 5. *In any arena where audit provides a lower bound on dishonesty, and any trace containing $a:\text{blame } B$ with $d \in B$, we have that d confessed.*

Trust mechanisms (such as trusted third parties) are required to establish the non-repudiation implied in the above proposition. We discuss these in Section 5.

Candidate 3 (Overlap). *An arena with audit provides overlap with dishonesty if, for any trace v with dishonest $D \not\equiv a$ containing $a:\text{blame } B$, we have $B \cap D = \emptyset$ implies $B = \emptyset$.*

Overlap is a more general property than providing a finite lower bound, since any lower bound $\{d_1, \dots, d_n\}$ can be replaced by a series of singleton overlaps $\{d_1\}, \dots, \{d_n\}$.

We do, however, note one problem with this definition, which is that although it is up-closed, it is not intersection-closed, that is there may be $v.a:\text{blame } B.w$ and $v.a:\text{blame } C.w$ which overlap with dishonesty, but $v.a:\text{blame } (B \cap C).w$ does not. This may arise in cases of separation of duty (Ferraiolo et al. 2003), if p and q must dishonestly collude to cause some action, then an auditor might choose to blame either $\{p\}$ or $\{q\}$, but not \emptyset . We leave this problem for future work.

Candidate 4 (Liveness). *An arena with audit is n -live if for any run $v.k.w$ such that a is honest, k is an input to a , and w contains at least n timeout actions, there is an action $a:\text{blame } B$ in w . An arena with audit is live whenever it is n -live for some n .*

As is common with correctness criteria, we distinguish between safety properties and liveness properties. In this case, liveness is quite simple to specify and verify (since an arena is n -live precisely when the honest processes for a are n -live).

4 Analysis Using Turn-Based Games

This section describes the use of game-based methods to automate the analysis of the properties described in the prior section. We refer the reader to (Alur et al. 2002) for background motivation and detailed examples.

Definition 6. A *turn-based game graph* over n -players player 1 to player n is $\mathcal{G} = (q, \mathcal{S} = \mathcal{S}_1 \uplus \dots \uplus \mathcal{S}_n, \mathcal{E}, \Pi, \pi)$ where:

- $(\mathcal{S}, \mathcal{E})$ is a directed graph with a total transition relation \mathcal{E} over the finite stateset \mathcal{S} .
- $\mathcal{S}_1, \dots, \mathcal{S}_n$ is a partition of \mathcal{S} and $q \in \mathcal{S}$ is the start state
- Π is a set of propositions; $\pi : \mathcal{S} \rightarrow \Pi$ yields the propositions true at each state. \square

An evolution proceeds as follows. States in \mathcal{S}_i are *player- i* states, where player i decides the successor state. A path in the game graph is a finite or infinite sequence of states. By totality, every finite path extends to a *play*, an infinite path of states.

Strategies. A (pure) strategy for a player is a recipe to extend a play, i.e., given a finite sequence of states, representing the history of the play, a strategy for a player chooses a unique successor state to extend the play.

Let mem_i be a set called *memory* that encodes the information about the history of the play. A player i strategy can be described as a pair of functions: a *memory-update*

function $\zeta^U: 2^{\Pi_i} \times mem_i \rightarrow mem_i$ to update the memory with the current state and a *next-move* function ζ^M that yields a new player i move for every element of $\mathcal{S}_i \times mem_i$. A strategy must prescribe only available moves, i.e., for all $s \in \mathcal{S}_i$, for all $m \in mem$, we have $(s, \zeta^M(s, m)) \in \mathcal{E}$.

Let Σ_i stand for the set of valid player i strategies under consideration. Strategies interact as follows. Player i follows the strategy ζ_i if in each player i move, she chooses the next state according to ζ_i^M . Once a starting state $s \in \mathcal{S}$ and strategies $\zeta_i \in \Sigma_i$ of the players are fixed, it is clear that the resulting outcome is a play of the game.

Compilation. We compile a finite collection of finite state processes (with a finite universe of messages) into a turn-based game. The translation uses new propositions $guilt_p$ for every $p \in \mathcal{A}$. If $guilt_p$ is satisfied by a state on a path, then p is dishonest on that path.

Logic. We use a fragment of the logic ATL^* (Alur et al. 2002). The usable propositions are restricted to the ones of interest. The path formulas exclude the next modality.

As a result, the properties are insensitive to the extra transitions introduced by the above compilation of arenas into turn-based games.

We refer to (Alur et al. 2002) for precise semantics. The state (ϕ) and path (ψ) are given by the following grammar: (A is any subset of principals)

$$\begin{aligned} \phi &::= true \mid guilt_p \mid \sigma \mid a:blame B \mid p \rightarrow q:m \mid \neg\phi \mid \phi \vee \phi \mid \langle\langle A \rangle\rangle\psi \\ \psi &::= true \mid \phi \mid \neg\psi \mid \psi \vee \psi \mid \Box\psi \mid \Diamond\psi \mid \psi \mathcal{U} \psi \end{aligned}$$

The formula $\langle\langle A \rangle\rangle\psi$ is true at a state if there exist strategies for the players in set A such that no matter what strategies the other players (in the complement of A) choose, the resulting play satisfies the path formula ψ .

We use existential and universal quantification over finite sets instead of finite disjunction and conjunction; e.g., $(\exists p)\psi$ is shorthand for $(\vee_p)\psi$. Also, we define

- $NonZero \triangleq \Box\Diamond\sigma$, to identify live traces with infinitely many σ actions.
- $AInit \triangleq \exists p, m. \Diamond(p \rightarrow a:m)$, to identify traces where the auditor a has been initialized by being sent some message.
- $Succ(B) \triangleq NonZero \wedge Ainit \wedge \Diamond(a:blame B)$, to identify Non-Zeno traces where the auditor has been contacted and the auditor has assigned blame to B .

and

<i>Overlap</i>	$\langle\langle \emptyset \rangle\rangle Succ(B) \Rightarrow (\exists p \in B) guilt_p$
<i>Lower bound</i>	$\langle\langle \emptyset \rangle\rangle Succ(B) \Rightarrow (\forall p \in B) guilt_p$

Since the auditor is fixed, these are LTL properties. Thus, when expressed in ATL^* , they have the strategy quantifier with the empty set to capture universal quantification over all traces reflecting other player choices. The soundness of the logical encoding above w.r.t. the trace based definitions of the earlier section follows from the soundness of the compilation w.r.t. the trace semantics of an arena.

Blamelessness of p for a fixed audit protocol is true at a state only if the agent a has a strategy to ensure that p never ends up in the blame set assigned by the auditor, independent of the given fixed auditor strategy and independent of any choice of strategies for the scheduler and the other players. Formally, we define

<i>Blamelessness</i>	$\langle\langle p \rangle\rangle \neg (\exists B \ni p) \Diamond (a:blame B)$
----------------------	---

The model-checking problem for ATL* is 2EXPTIME in the size of the formula and PTIME-hard for bounded-size formulas (Alur et al. 2002). So, we have:

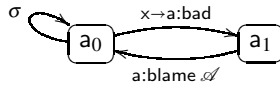
Proposition 7. *The model-checking problem for Overlap, Lower bound and Blamelessness for an arena $\langle \mathcal{A}, \mathcal{M}, \mathcal{H} \rangle$ with a fixed audit protocol is solvable in EXPTIME in the size of the arena and 2EXPTIME in the formula size.*

The formulas of interest are small. The bottleneck is the EXPTIME dependence on arenas caused by the determinization of the honest processes in the compilation process.

5 Example Auditors

We present a series of examples in which the auditor aims to detect the origin of a special bad message. At the end of this section, we relate the discussion to Example 1.

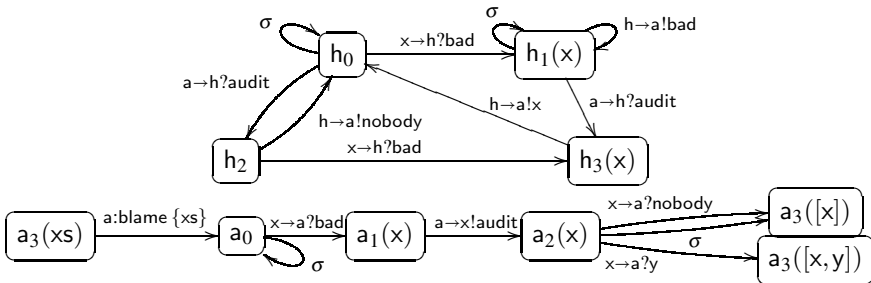
Example 8. Consider auditing the arena in Example 3. When the auditor receives a bad message, they know that there is a dishonest principal. However, since the arena does not permit them to query principals for further information, they have no way to discover the guilty parties. So the best they can do is blame everyone:



This auditor provides liveness and overlap with dishonesty, albeit trivially. This auditor does not provide lower bound. □

We now consider audit protocols where honest principals are required to respond to requests for information from the auditor. From a principal h , the auditor will request the identity of the sender of a bad message to h . We analyze the variations that arise depending on whether: (a) honest principals may forward bad messages to principals other than the auditor; (b) honest principals are *required* to report bad messages (all are *allowed* to report bad messages); (c) senders report to whom they have forwarded bad.

Example 9. Extend the arena of Example 2 to accommodate audit by setting $\mathcal{M} = \mathcal{A} \cup \{\text{bad, audit, nobody}\}$. The honest processes are described below, where we describe the potentially infinite state transition systems using state variables x, y, xs ranging over \mathcal{A} . States h_1, a_1 and a_2 are parameterized by principal x , which sent the bad message. State a_3 is parameterized by the principal list xs , which are blamed; we use ML notation for lists ($[]$ for the empty list, $::$ for prefixing, $@$ for concatenation).



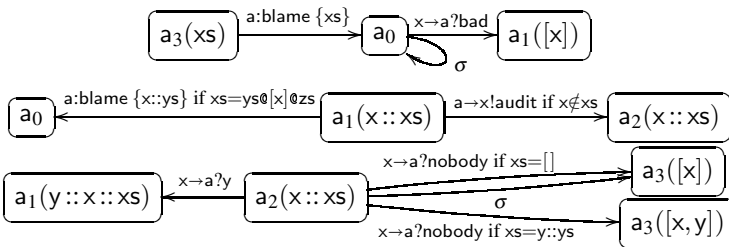
Honest principals may initiate an audit by reporting the receipt of a bad message to the auditor (the transition $h \rightarrow a!bad$ at h_1). The auditor responds to audit with a request for the sender of the bad message. If the auditor's request (at a_2) times out, or the response is nobody, then the principal initiating the audit is dishonest and is blamed. If the response $x \rightarrow a?y$ is received, then the auditor blames $\{x, y\}$ because it is unable to detect whether y initiated or forwarded bad, or whether x is lying about receipt of bad from y . This auditor provides liveness and overlap with dishonesty. The algorithms of Section 4 verify this for the case when the number of principals is finite. \square

We now analyze the consequences of allowing honest principals to forward bad messages to principals other than the auditor.

Example 10. Allowing honest principals to forward bad messages to all principals, as in Example 3, necessitates changing the auditor to track down the original source of the bad message. To see this, first modify the arena from Example 9 by replacing the $h \rightarrow a!bad$ self-loop on $h_1(x)$ with $h \rightarrow y!bad$.

The auditor from Example 9 does not provide an overlap with dishonesty for this new arena, because a trace of the form $d \rightarrow h!bad, h \rightarrow g!bad, g \rightarrow a!bad, \dots$ would result in g and perhaps h being blamed incorrectly (their forwarding behavior is honest) when only d has been dishonest (initially sending bad).

To identify an originator of a bad message (there may be several), the auditor below follows a chain of forwarders until it receives: (a) no response (a timeout); (b) the answer nobody; (c) the answer a ; or (d) it finds a cycle of forwarders. The auditor then blames: (a) the principal that did not respond to an audit request (it is dishonest to ignore the auditor); (b) the principal p that responded with nobody and, if there is one, the principal q that claimed p forwarded bad to q (either q is lying about receiving a forwarded bad or p is unable to identify a principal that forwarded bad to them); (c) the principal that claimed a forwarded bad (that principal is lying because the auditor does not send bad); (d) all principals in the cycle (one of them is lying about the source). States a_1, a_2 and a_3 are parameterized by the list of suspected principals.



This auditor provides liveness and overlap with dishonesty. The algorithms of Section 4 verify this for the case when the number of principals is finite. \square

It is important in the above example that the arena requires honest principals to record the *initial* sender of bad rather than the most recent sender to make the auditor overlap with dishonesty. (If instead the *most recent* sender of bad was reported to the auditor, and we saw a trace ending with a cycle of the form $d \rightarrow h!bad, h \rightarrow g!bad, g \rightarrow h!bad, h \rightarrow a!bad, \dots$, then the auditor would find and blame the cycle h to g to h . Neither g nor

h are dishonest, so the auditor above would not overlap with dishonesty for the modified arena.)

In both Example 9 and Example 10, an honest agent is unable to achieve *Blamelessness*. We address this issue next by encoding the use of notaries as trusted third parties to permit honest agents to establish blamelessness.

Notaries. The presence of notaries provides a non-repudiation function and disables the ability of a dishonest principal d to get an honest principal h blamed (by simply claiming that h sent bad to d). The notary principals are assumed to be honest. For this reason, we refer to the notary principals as Trusted Third Parties (TTPs). Here we consider a single non-auditor principal for the sake of simplicity, but it is not essential that there be only one TTP.

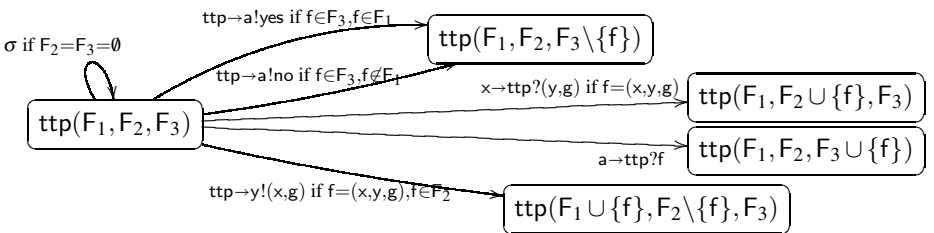
We assume a collection of messages $\mathcal{G} \ni g$ that pass uninterpreted through the TTP. We define \mathcal{G} to include the bad message with different provenance chains indicating the path of the bad message. With \mathcal{G} fixed, we then define the messages of the arena by:

$$\begin{aligned} \mathcal{M} &\triangleq (\mathcal{A} \times \mathcal{G}) && \text{(Messages to and from TTP)} \\ &\cup (\mathcal{A} \times \mathcal{A} \times \mathcal{G}) && \text{(Message query by auditor to TTP)} \\ &\cup \{\text{yes, no, unknown}\} && \text{(Response to auditor)} \end{aligned}$$

We use f to range over *forwarding records* of the form (x, y, g) indicating that the TTP forwarded g from x to y . We use F to range over sets of forwarding records.

The TTP interacts with principals by forwarding messages on their behalf. A principal x sends a forwarding request of the form (y, g) to the TTP (indicating the target). Subsequently, the TTP forwards the message (x, g) (indicating the source) to y , and adds the forwarding record (x, y, g) to its store. The TTP also respond to queries from the auditor that ask whether $f = (x, y, g)$ has been forwarded in the past. It can only respond honestly with yes (resp. no) if its store contains the forwarding record f (resp. does not contain the forwarding record f).

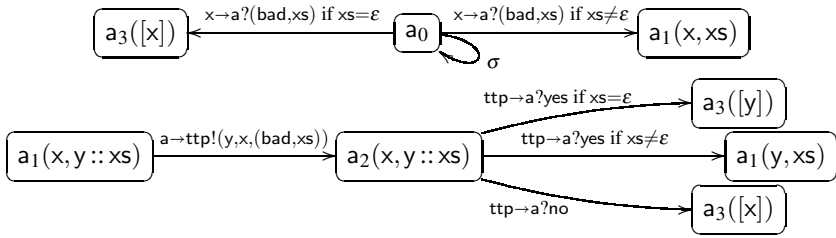
The TTP state $\text{ttp}(F_1, F_2, F_3)$ is parameterized by three sets of forwarding records. The set F_1 stores which messages have been forwarded. The set F_2 maintains the forwarding requests received but not yet acted upon. The set F_3 maintains the auditor requests received but not yet acted upon. The TTP may only timeout when there are no actions to complete, i.e., $F_2 = F_3 = \emptyset$. The sets of actions not yet completed are present to ensure that the TTP is input enabled. The behavior of the TTP is specified as:



Provenance. We add *provenance* information to the messages. In this context, provenance is a sequence (possibly empty) of principal names indicating the path of a forwarded message. An empty provenance sequence indicates that the message was not

forwarded, i.e., in $x \rightarrow y!(\text{bad}, \varepsilon)$, x is confessing to sending bad directly. In contrast, a non-empty provenance sequence of the form $(z :: xs)$ indicates that the message was forwarded with z being the most recent forwarder, i.e., in $x \rightarrow y!(\text{bad}, (z :: xs))$, x is claiming that they received the forwarded message from z as $z \rightarrow x?(\text{bad}, xs)$. We now demand that honest communication between principals occurs via the TTP (operating without knowledge of the provenance structure), and so we define $\mathcal{G} \triangleq \{\text{bad}\} \times \mathcal{A}^*$.

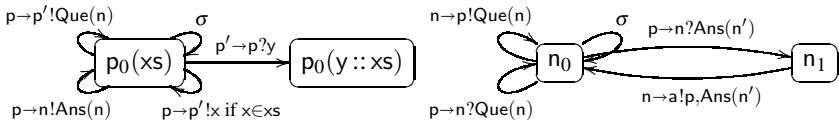
Auditor. When a principal x initiates an audit by sending bad paired with a provenance sequence to the auditor, the auditor can verify the entire provenance sequence step-by-step, using the TTP to determine whether each forward indicated in the provenance sequence is genuine. If the empty provenance sequence is ultimately found, the initial sender is blamed. If the TTP responds with no at any point, then a principal has claimed that it forwarded a message but is unable to prove its claim, that principal is blamed. The auditor is formalized as:



Honest Agents. Honest agent h are constrained as follows. (a) h is required to report bad to the auditor, and (b) h can only forward bad messages that are received via the TTP, after honestly updating the provenance and using the TTP. We elide the straightforward formalization.

Since the auditor is able to verify evidence: (a) Honest agents have *Blamelessness*, and (b) the auditor has *Lower bound*. The algorithms from Section 4 verify these statements when there are finitely many principals and messages and the length of the provenance chain is bounded.

Example 1 revisited. We discuss briefly the implications for Example 1. Let p, p' range over health professionals (doctors and nurses) and n, n' over patients and non-health professionals (secretaries). Let $\text{Que}(n)$ and $\text{Ans}(n)$ be messages representing question and answers concerning patient n . In the following processes, xs represents messages that have been sent to p that may be forwarded.



Honest health care professionals have unrestricted exchange of messages amongst themselves. However their answers and questions to patients are constrained to be about the receiver. A patient who has received an answer about another patient is allowed to initiate an audit via a message to the auditor.

The model closest to the original example Example 1 is the one from Example 5 since the MyHealth website is effectively a TTP. The techniques of Example 5 permit

an auditor permit an auditor to achieve *Lower bound*, and the honest agents to have *Blamelessness*.

Perhaps of greater interest, our analysis in Section 5 shows that even without TTPs, auditors can achieve *Overlap* in a distributed setting with only integrity assumptions on communication. This demonstrates that powerful (and expensive) mechanisms such as notaries are not necessary for all audit protocols.

The algorithms from Section 4 verify these statements for the special case when there are finitely many principals, the length of the provenance chain is bounded, and the internal state of the honest health care professional is bounded (i.e., they remember only a bounded number of messages). The extension of our methods to symbolic methods that permit handling infinite state spaces is left for future work.

6 Related Work

The security of the audit trail has built on advances in authenticated data structures (e.g., secure histories (Maniatis and Baker 2002), Persistent Authenticated Dictionaries (Anagnostopoulos et al. 2001) and Undeniable Attestations (Buldas et al. 2000)). This research has been used in specific applications. For example, PeerReview (Haeberlen et al. 2007) creates a per-node secure log, which records the messages a node has sent and received, and the inputs and outputs of the application. Node failures are detected by replaying such a trace against a reference implementation that is assumed to be determinate. CATS (Yumerefendi and Chase 2007) validates the integrity of storage hosted by a service provider. The clients are provided with the means to verify that all (and only) updates from authorized users are applied and seen. AudIt (Argyrazi et al. 2007) is an explicit accountability interface for ISPs to supply feedback to traffic sources on QoS considerations. Accountability for the Internet protocol has also been investigated (Andersen et al. 2008). The APPLE system (Etalle and Winsborough 2007) suggests an architecture for a posteriori policy enforcement on documents: documents are always associated with policies, all clients operations on documents are logged, and distributed auditors occasionally verify the compliance with policies. In value-commitment protocols, a principal commits to a hidden value. Other principals cannot read this value, but can detect unlawful updates after the commitment. Fournet et al. (2008) studies such protocols using an applied pi-calculus.

These papers focus on efficient and expressive audit mechanisms to realize specific accountability policies. We study general models and limitations of accountability, aiming to provide a foundational analysis that can be incorporated as a component in the design of such systems. For example, the design goals of PeerReview include *Blamelessness* for honest agents and *Lower bound* for auditors. Our analysis provides a justification for the need to use secure ACKS to achieve these goals of PeerReview.

Cederquist et al. (2005) describe a policy language for data ownership and administrative issues. Cederquist et al. (2007) describe a system that uses audits to enforce compliance to such policies. Proof-carrying-authorization forces the requestors of access to provide proofs validating their request. The AURA project (Vaughan et al. 2008) reuses these proofs for accountability via the “proofs as log entries” approach. These papers focus on the design of logical methods to specify policies and enforce them via

accountability. We study the design of the policies themselves, exploring the tradeoffs between the requirements that system policies place on honest agents and the power of audit protocols.

Our analysis methods are based on game-based logics for multiagent systems with *perfect* information, such as Alternating Temporal Logic (Alur et al. 2002).

7 Conclusions

We aim to develop foundations for distributed accountability systems. We have suggested an operational model and developed analysis methods using translations into games. Our running example suggests that our framework permits the designer of audit based accountability systems to explore the tradeoffs between the requirements on (a) the honest principals, (b) the guarantees provided by the communication network, and (c) the precision demanded of the audit protocol.

Three important issues remain need to be addressed in future work: (a) the full integration with cryptographic primitives in the operational model, (b) quantitative models and methods such as Bloom filters are critical to achieving efficient audits of large datasets (Calandrino et al. 2007), and (c) equilibria notions provide an analysis of player intentions that is crucial to mechanism design.

References

- Abadi, M., Birrell, A., Wobber, T.: Access control in a world of software diversity. In: Proc. of the Tenth workshop on Hot Topics in Operating Systems (2005), <http://www.usenix.org/events/hotos05/>
- Alur, R., Henzinger, T., Kupferman, O.: Alternating time temporal logic. *Journal of ACM* 49, 672–713 (2002)
- Anagnostopoulos, A., Goodrich, M.T., Tamassia, R.: Persistent authenticated dictionaries and their applications. In: Davida, G.I., Frankel, Y. (eds.) *ISC 2001*. LNCS, vol. 2200, pp. 379–393. Springer, Heidelberg (2001)
- Andersen, D.G., Balakrishnan, H., Feamster, N., Koponen, T., Moon, D., Shenker, S.: Accountable Internet Protocol (AIP). In: *SIGCOMM*, pp. 339–350. ACM Press, New York (2008)
- Argyrazi, K., Maniatis, P., Irzak, O., Shenker, S.: An accountability interface for the Internet. In: *Proceedings of the 14th IEEE International Conference on Network Protocols (2007)*
- Barth, A., Mitchell, J.C., Datta, A., Sundaram, S.: Privacy and utility in business processes. In: *CSF*, pp. 279–294. IEEE Computer Society, Los Alamitos (2007)
- Brookes, S.D., Hoare, C.A.R., Roscoe, A.W.: A theory of communicating sequential processes. *J. ACM* 31(3), 560–599 (1984)
- Buldas, A., Laud, P., Lipmaa, H.: Accountable certificate management using undeniable attestations. In: *ACM Conference on Computer and Communications Security*, pp. 9–17 (2000)
- Calandrino, J.A., Halderman, J.A., Felten, E.W.: Machine-assisted election auditing. In: *EVT 2007: Proceedings of the USENIX Workshop on Accurate Electronic Voting Technology*, p. 9. USENIX Association (2007)
- Cederquist, J.G., Corin, R., Dekker, M.A.C., Etalle, S., den Hartog, J.I.: An audit logic for accountability. In: *POLICY*, pp. 34–43. IEEE Computer Society Press, Los Alamitos (2005)
- Cederquist, J.G., Corin, R., Dekker, M.A.C., Etalle, S., den Hartog, J.I., Lenzini, G.: Audit-based compliance control. *Int. J. Inf. Sec.* 6(2-3), 133–151 (2007)

- Cirillo, A., Jagadeesan, R., Pitcher, C., Riely, J.: TAPIDO: Trust and authorization via provenance and integrity in distributed objects (extended abstract). In: Drossopoulou, S. (ed.) ESOP 2008. LNCS, vol. 4960, pp. 208–223. Springer, Heidelberg (2008)
- Eriksén, S.: Designing for accountability. In: Proceedings of the second Nordic conference on Human-computer interaction, pp. 177–186 (2002)
- Etalle, S., Winsborough, W.H.: A posteriori compliance control. In: SACMAT, pp. 11–20. ACM, New York (2007)
- Ferraiolo, D.F., Kuhn, D.R., Chandramouli, R.: Role-Based Access Control. Computer Security Series. Artech House (2003)
- Fournet, C., Gordon, A.D., Maffei, S.: A type discipline for authorization policies. In: Sagiv, M. (ed.) ESOP 2005. LNCS, vol. 3444, pp. 141–156. Springer, Heidelberg (2005)
- Fournet, C., Guts, N., Nardelli, F.Z.: A formal implementation of value commitment. In: Drossopoulou, S. (ed.) ESOP 2008. LNCS, vol. 4960, pp. 383–397. Springer, Heidelberg (2008)
- Friedman, B., Grudin, J.: Trust and accountability: preserving human values in interactional experience. In: CHI 1998: CHI 1998 conference summary on Human factors in computing systems, p. 213. ACM, New York (1998)
- Haerberlen, A., Kouznetsov, P., Druschel, P.: PeerReview: practical accountability for distributed systems. In: Proceedings of 21st ACM SIGOPS symposium on Operating systems principles, pp. 175–188. ACM, New York (2007)
- Halpern, J.Y., Weissman, V.: Using first-order logic to reason about policies. In: CSFW, pp. 118–130 (2003)
- Hennesy, M., Regan, T.: A process algebra for timed systems. *Inf. Comput.* 117(2), 221–239 (1995)
- Lampson, B.W.: Computer security in the real world. *IEEE Computer* 37(6), 37–46 (2004)
- Li, N., Mitchell, J.C.: A role-based trust-management framework. In: DISCEX (1), p. 201. IEEE Computer Society Press, Los Alamitos (2003)
- Lynch, N.A.: Input/output automata: Basic, timed, hybrid, probabilistic, dynamic... In: Amadio, R.M., Lugiez, D. (eds.) CONCUR 2003. LNCS, vol. 2761, pp. 187–188. Springer, Heidelberg (2003)
- Maniatis, P., Baker, M.: Secure history preservation through timeline entanglement. In: USENIX Security Symposium, pp. 297–312. USENIX (2002)
- Schneider, F.B.: Enforceable security policies. *Information and System Security* 3(1), 30–50 (2000)
- Vaughan, J.A., Jia, L., Mazurak, K., Zdancewic, S.: Evidence-based audit. In: CSF, pp. 177–191. IEEE Computer Society, Los Alamitos (2008)
- Weitzner, D.J., Abelson, H., Berners-Lee, T., Feigenbaum, J., Hendler, J., Sussman, G.J.: Information accountability. Technical Report MIT-CSAIL-TR-2007-034, MIT (June 2007), <http://hdl.handle.net/1721.1/37600>
- Yumerefendi, A.R., Chase, J.S.: Trust but verify: accountability for network services. In: EW11: Proceedings of the 11th workshop on ACM SIGOPS European workshop, p. 37. ACM, New York (2004)
- Yumerefendi, A.R., Chase, J.S.: Strong accountability for network storage. *Trans. Storage* 3(3), 11 (2007)