

Securing Databases from Probabilistic Inference

Marco Guarnieri

Institute of Information Security
Department of Computer Science
ETH Zurich, Switzerland
marco.guarnieri@inf.ethz.ch

Srdjan Marinovic

The Wireless Registry, Inc.
Washington DC, US
srdjan@wirelessregistry.com

David Basin

Institute of Information Security
Department of Computer Science
ETH Zurich, Switzerland
basin@inf.ethz.ch

Abstract—Databases can leak confidential information when users combine query results with probabilistic data dependencies and prior knowledge. Current research offers mechanisms that either handle a limited class of dependencies or lack tractable enforcement algorithms. We propose a foundation for Database Inference Control based on PROBLOG, a probabilistic logic programming language. We leverage this foundation to develop ANGERONA, a provably secure enforcement mechanism that prevents information leakage in the presence of probabilistic dependencies. We then provide a tractable inference algorithm for a practically relevant fragment of PROBLOG. We empirically evaluate ANGERONA’s performance showing that it scales to relevant security-critical problems.

I. INTRODUCTION

Protecting the confidentiality of sensitive data stored in databases requires protection from both *direct* and *indirect* access. The former happens when a user observes query results, and the latter happens when a user infers sensitive information by combining results with external information, such as data dependencies or prior knowledge. Controlling indirect access to data is often referred to as *Database Inference Control* [30] (DBIC). This topic has attracted considerable attention in recent years, and current research considers different sources of external information, such as the database schema [15], [37], [39], [42], [57], [63], [64], the system’s semantics [37], statistical information [3], [16], [22]–[24], exceptions [37], error messages [44], user-defined functions [44], and data dependencies [10], [12], [54], [55], [67], [68], [72].

An important and relevant class of data dependencies are probabilistic dependencies, such as those found in genomics [43], [47], [49], social networks [40], and location tracking [52]. Attackers can exploit these dependencies to infer sensitive information with high confidence. To effectively prevent probabilistic inferences, DBIC mechanisms should (1) support a large class of probabilistic dependencies, and (2) have tractable runtime performance. The former is needed to express different attacker models. The latter is necessary for mechanisms to scale to real-world databases.

Most existing DBIC mechanisms support only precise data dependencies [10], [12], [67], [68], [72] or just limited classes of probabilistic dependencies [14], [15], [39], [45], [54], [55], [71]. As a result, they cannot reason about the complex probabilistic dependencies that exist in many realistic settings. Mardziel et al.’s mechanism [51] instead supports arbitrary

probabilistic dependencies, but no complexity bounds have been established and their algorithm appears to be intractable.

Contributions. We develop a tractable and practically useful DBIC mechanism based on probabilistic logic programming.

First, we develop ATKLOG, a language for formalizing users’ beliefs and how they evolve while interacting with the system. ATKLOG builds on PROBLOG [20], [21], [31], a state-of-the-art probabilistic extension of DATALOG, and extends its semantics by building on three key ideas from [18], [46], [51]: (1) users’ beliefs can be represented as probability distributions, (2) belief revision can be performed by conditioning the probability distribution based on the users’ observations, and (3) rejecting queries as insecure may leak information. By combining DATALOG with probabilistic models and belief revision based on users’ knowledge, ATKLOG provides a natural and expressive language to model users’ beliefs and thereby serves as a foundation for DBIC in the presence of probabilistic inferences.

Second, we identify acyclic PROBLOG programs, a class of programs where probabilistic inference’s data complexity is PTIME. We precisely characterize this class and develop a dedicated inference engine. Since PROBLOG’s inference is intractable in general, we see acyclic programs as an essential building block to effectively using ATKLOG for DBIC.

Finally, we present ANGERONA¹, a novel DBIC mechanism that secures databases against probabilistic inferences. We prove that ANGERONA is secure with respect to any ATKLOG-attacker. In contrast to existing mechanisms, ANGERONA provides precise tractability and completeness guarantees for a practically relevant class of attackers. We empirically show that ANGERONA scales to relevant problems of interest.

Structure. In §II, we illustrate the security risks associated with probabilistic data dependencies. In §III, we present our system model, which we formalize in §IV. We introduce ATKLOG in §V and in §VI we present our inference engine for acyclic programs. In §VII, we present ANGERONA. We discuss related work in §VIII and draw conclusions in §IX. A prototype of our enforcement mechanism is available at [36].

II. MOTIVATING EXAMPLE

Hospitals and medical research centres store large quantities of health-related information for purposes ranging from diag-

¹Angerona is the Roman goddess of silence and secrecy, and She is the keeper of the city’s sacred, and secret, name.

nosis to research. As this information is extremely sensitive, the databases used must be carefully secured [13], [28]. This task is, however, challenging due to the dependencies between health-related data items. For instance, information about someone’s hereditary diseases or genome can be inferred from information about her relatives. Even seemingly non-sensitive information, such as someone’s job or habits, may leak sensitive health-related information such as her predisposition to diseases. Most of these dependencies can be formalized using probabilistic models developed by medical researchers.

Consider a database storing information about the smoking habits of patients and whether they have been diagnosed with lung cancer. The database contains the tables *patient*, *smokes*, *cancer*, *father*, and *mother*. The first table contains all patients, the second contains all regular smokers, the third contains all diagnosed patients, and the last two associate patients with their parents. Now consider the following probabilistic model: (a) every patient has a 5% chance of developing cancer, (b) for each parent with cancer, the likelihood that a child develops cancer increases by 15%, and (c) if a patient smokes regularly, his probability of developing cancer increases by 25%. We intentionally work with a simple model since, despite its simplicity, it illustrates the challenges of securing data with probabilistic dependencies. We refer the reader to medical research for more realistic probabilistic models [4], [69].

The database is shared between different medical researchers, each conducting a research study on a subset of the patients. All researchers have access to the *patient*, *smokes*, *father*, and *mother* tables. Each researcher, however, has access only to the subset of the *cancer* table associated with the patients that opted-in to his research study. We want to protect our database against a malicious researcher whose goal is to infer the health status of patients not participating in the study. This is challenging since restricting direct access to the *cancer* table is insufficient. Sensitive information may be leaked even by queries involving only authorized data. For instance, the attacker may know that the patient *Carl*, which has not disclosed his health status, smokes regularly. From this, he can infer that *Carl*’s probability of developing lung cancer is, at least, 30%. If, additionally, *Carl*’s parents opted-in to the research study and both have cancer, the attacker can directly infer that the probability of *Carl* developing lung cancer is 60% by accessing his parents’ information.

Security mechanisms that ignore such probabilistic dependencies allow attackers to infer sensitive information. An alternative is to use standard DBIC mechanisms and encode all dependencies as precise, non-probabilistic, dependencies. This, however, would result in an unusable system. Medical researchers, even honest ones, would be able to access the health-related status only of those patients whose relatives also opted-in to the user study, independently of the amount of leaked information, which may be negligible. Hence, to secure the database and retain usability, it is essential to reason about the probabilistic dependencies.

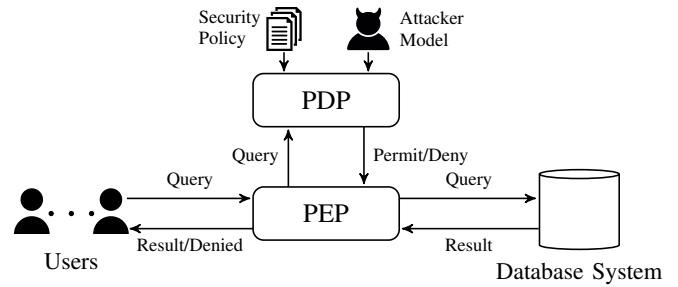


Fig. 1: System model.

III. SYSTEM MODEL

Figure 1 depicts our system model. Users interact with two components: a database system and an inference control system, which consists of a Policy Decision Point (PDP) and a Policy Enforcement Point (PEP). We assume that all communication between users and the components and between the components themselves is over secure channels.

Database System. The database system manages the system’s data. Its state is a mapping from tables to sets of tuples.

Users. Each user has a unique account used to retrieve information from the database system by issuing `SELECT` commands. Note that these commands do not change the database state. This reflects settings where users have only read-access to a database. Each command is checked by the inference control system and is executed if and only if the command is authorized by the security policy.

Security policy. The system’s security policy consists of a set of *negative permissions* specifying information to be kept secret. These permissions express bounds on users’ beliefs, formalized as probability distributions, about the actual database content. Negative permissions are formalized using commands of the form `SECRET q FOR u THRESHOLD l` , where q is a query, u is a user identifier, and l is a rational number, $0 \leq l \leq 1$. This represents the requirement that “A user u ’s belief in the result of q must be less than l .” Namely, the probability assigned by u ’s belief to q ’s result must be less than l . Requirements like “A user u is not authorized to know the result of q ” can be formalized as `SECRET q FOR u THRESHOLD 1`. The system also supports commands of the form `SECRET q FOR USERS NOT IN $\{u_1, \dots, u_n\}$ THRESHOLD l` , which represents the requirement that “For all users $u \notin \{u_1, \dots, u_n\}$, u ’s belief in the result of q must be less than l .”

Attacker. An attacker is a system user with an assigned user account, and each user is a potential attacker. An attacker’s goal is to violate the security policy, that is, to read or infer information about one of the `SECRETS` with a probability of at least the given threshold.

An attacker can interact with the system and observe its behaviour in response to his commands. Furthermore, he can reason about this information and infer information by exploiting domain-specific relationships between data items.

We assume that attackers know the database schema as well as any integrity constraints on it.

Attacker Model. An attacker model represents each user’s initial beliefs about the actual database state and how he updates his beliefs by interacting with the system and observing its behaviour in response to his commands. These beliefs may reflect the attacker’s knowledge of domain-specific relationships between the data items or prior knowledge.

Inference Control System. The inference control system protects the confidentiality of database data. It consists of a PEP and a PDP, configured with a security policy P and an attacker model ATK . For each user, the inference control system keeps track of the user’s beliefs according to ATK .

The system intercepts all commands issued by the users. When a user u issues a command c , the inference control system decides whether u is authorized to execute c . If c complies with the policy, i.e., the users’ beliefs still satisfy P even after executing c , then the system forwards the command to the database, which executes c and returns its result to u . Otherwise, it raises a *security exception* and rejects c .

IV. FORMAL MODEL

A. Database Model

We introduce here background and notation for databases and queries. Our formalization follows [2].

Let \mathcal{R} be a countably infinite set representing identifiers of relation schemas. A *database schema* D is a pair $\langle \Sigma, \mathbf{dom} \rangle$, where Σ is a first-order signature and \mathbf{dom} is a fixed domain. For simplicity, we consider just a single domain. Extensions to the many-sorted case are straightforward [2]. The signature Σ consists of a set of *relation schemas* $R \in \mathcal{R}$, each schema with arity $|R|$, and one constant symbol for each constant in \mathbf{dom} . We interpret constants by themselves in the semantics.

A *state* s of D is a finite Σ -structure with domain \mathbf{dom} that interprets each relation schema R by an $|R|$ -ary relation over \mathbf{dom} . We denote by Ω_D the set of all states. Given a schema $R \in D$, $s(R)$ denotes the tuples that belong to (the interpretation of) R in the state s . We assume that the domain \mathbf{dom} is finite, as is standard for many application areas combining databases and probabilistic reasoning [20], [33], [47], [65]. In this case, the set of all states Ω_D is finite.

A *query* q over a schema D is of the form $\{\bar{x} \mid \phi\}$, where \bar{x} is a sequence of variables, ϕ is a relational calculus formula over D , and ϕ ’s free variables are those in \bar{x} . A *boolean query* is a query $\{\mid \phi\}$, also written as ϕ , where ϕ is a sentence. The result of executing a query q on a state s , denoted by $[q]^s$, is a boolean value in $\{\top, \perp\}$, if q is a boolean query, or a set of tuples otherwise. Furthermore, given a sentence ϕ , $\llbracket \phi \rrbracket$ denotes the set $\{s \in \Omega_D \mid [\phi]^s = \top\}$. We denote by RC (respectively RC_{bool}) the set of all relational calculus queries (respectively sentences). We consider only *domain-independent queries* and we employ the standard relational calculus semantics [2].

An *integrity constraint* over D is a relational calculus sentence γ over D . Given a state s , we say that s *satisfies the constraint* γ iff $[\gamma]^s = \top$. Given a set of constraints Γ ,

	<u>patient</u>	<u>smokes</u>	<u>father</u>
	Alice	Bob	Bob Carl
	Bob	Carl	
	Carl		<u>mother</u>
			Alice Carl

Fig. 2: The template for all database states, where the content of the *cancer* table is left unspecified.

Ω_D^Γ denotes the set of all states satisfying the constraints in Γ , i.e., $\Omega_D^\Gamma = \{s \in \Omega_D \mid \bigwedge_{\gamma \in \Gamma} [\gamma]^s = \top\}$.

Example IV.1. The database associated with the example in §II consists of five relational schemas *patient*, *smokes*, *cancer*, *father*, and *mother*, where the first three schemas have arity 1 and the last two have arity 2. We assume that there are only three patients Alice, Bob, and Carl, so the domain \mathbf{dom} is $\{\text{Alice}, \text{Bob}, \text{Carl}\}$. The integrity constraints are as follows:

- Alice, Bob, and Carl are patients.

$$patient(\text{Alice}) \wedge patient(\text{Bob}) \wedge patient(\text{Carl})$$

- Alice and Bob are Carl’s parents.

$$\forall x, y. (father(x, y) \leftrightarrow (x = \text{Bob} \wedge y = \text{Carl})) \wedge \\ \forall x, y. (mother(x, y) \leftrightarrow (x = \text{Alice} \wedge y = \text{Carl}))$$

- Alice does not smoke, whereas Bob and Carl do.

$$\neg smokes(\text{Alice}) \wedge smokes(\text{Bob}) \wedge smokes(\text{Carl})$$

Given these constraints, there are just 8 possible database states in Ω_D^Γ , which differ only in their *cancer* relation. The content of the *cancer* relation is a subset of $\{\text{Alice}, \text{Bob}, \text{Carl}\}$, whereas the content of the other tables is shown in Figure 2. We denote each possible world as s_C , where the set $C \subseteq \{\text{Alice}, \text{Bob}, \text{Carl}\}$ denotes the users having cancer. ■

B. Security Policies

Existing access control models for databases are inadequate to formalize security requirements capturing probabilistic dependencies. For example, SQL cannot express statements like “A user u ’s belief that ϕ holds must be less than l .” We present a simple framework, inspired by knowledge-based policies [51], for expressing such requirements.

A *D-secret* is a tuple $\langle U, \phi, l \rangle$, where U is either a finite set of users in \mathcal{U} or a co-finite set of users, i.e., $U = \mathcal{U} \setminus U'$ for some finite $U' \subset \mathcal{U}$, ϕ is a relational calculus sentence over D , and l is rational number $0 \leq l \leq 1$ specifying the uncertainty threshold. Abusing notation, when U consists of a single user u , we write u instead of $\{u\}$. Informally, $\langle U, \phi, l \rangle$ represents that for each user $u \in U$, u ’s belief that ϕ holds in the actual database state must be less than l . Therefore, a command of the form SECRET q FOR u THRESHOLD l can be represented as $\langle u, q, l \rangle$, whereas a command SECRET q FOR USERS NOT IN $\{u_1, \dots, u_n\}$ THRESHOLD l can be represented as $\langle \mathcal{U} \setminus \{u_1, \dots, u_n\}, q, l \rangle$. Finally, a *D-security policy* is a finite set of *D-secrets*. Given a *D-security policy* P , we denote by $secrets(P, u)$ the

set of D -secrets associated with the user u , i.e., $\text{secrets}(P, u) = \{\langle u, \phi, l \rangle \mid \langle U, \phi, l \rangle \in P \wedge u \in U\}$. Note that the function secrets is computable since the set U is always either finite or co-finite.

Our framework also allows the specification of lower bounds. Requirements of the form “A user u ’s belief that ϕ holds must be greater than l ” can be formalized as $\langle u, \neg\phi, 1-l \rangle$ (since the probability of $\neg\phi$ is $1-P(\phi)$, where $P(\phi)$ is ϕ ’s probability). Security policies can be extended to support secrets over non-boolean queries. A secret $\langle u, \{\bar{x} \mid \phi(\bar{x})\}, l \rangle$ can be seen as a shorthand for the set $\{\langle u, \phi[\bar{x} \mapsto \bar{t}], l \rangle \mid \bar{t} \in \bigcup_{s \in \Omega_D^\Gamma} [\{\bar{x} \mid \phi(\bar{x})\}]^s\}$, i.e., u ’s belief in any tuple \bar{t} being in the query’s result must be less than l .

Example IV.2. Let *Mallory* denote the malicious researcher from §II and D be the schema from Example IV.1. Consider the requirement from §II: *Mallory*’s belief in a patient having cancer must be less than 50%. This can be formalized as $\langle \text{Mallory}, \text{cancer}(\text{Alice}), 1/2 \rangle$, $\langle \text{Mallory}, \text{cancer}(\text{Bob}), 1/2 \rangle$, and $\langle \text{Mallory}, \text{cancer}(\text{Carl}), 1/2 \rangle$, or equivalently as $\langle \text{Mallory}, \{p \mid \text{cancer}(p)\}, 1/2 \rangle$. In contrast, the requirement “For all users u that are not *Carl*, u ’s belief in *Carl* having cancer must be less than 50%” can be formalized as $\langle U \setminus \{\text{Carl}\}, \text{cancer}(\text{Carl}), 1/2 \rangle$, where *Carl* denotes the user identifier associated with Carl. ■

C. Formalized System Model

We now formalize our system model. We first define a system configuration, which describes the database schema and the integrity constraints. Afterwards, we define the system’s state. Finally, we define a system run, which represents a possible interaction of users with the system.

A *system configuration* is a tuple $\langle D, \Gamma \rangle$, where D is a database schema and Γ is a set of D -integrity constraints. Let $C = \langle D, \Gamma \rangle$ be a system configuration. A *C-system state* is a tuple $\langle db, U, P \rangle$, where $db \in \Omega_D^\Gamma$ is a database state, $U \subset \mathcal{U}$ is a finite set of users, and P is a D -security policy. A *C-query* is a pair $\langle u, \phi \rangle$ where $u \in \mathcal{U}$ is a user and ϕ is a relational calculus sentence over D .² We denote by Ω_C the set of all system states and by \mathcal{Q}_C the set of all queries.

A *C-event* is a triple $\langle q, a, res \rangle$, where q is a C -query in \mathcal{Q}_C , $a \in \{\top, \perp\}$ is a security decision, where \top stands for “authorized query” and \perp stands for “unauthorized query”, and $res \in \{\top, \perp, \dagger\}$ is the query’s result, where \top and \perp represent the usual boolean values and \dagger represents that the query was not executed as access was denied. Given a C -event $e = \langle q, a, res \rangle$, we denote by $q(e)$ (respectively $a(e)$ and $res(e)$) the query q (respectively the decision a and the result res). A *C-history* is a finite sequence of C -events. We denote by \mathcal{H}_C the set of all possible C -histories. Moreover, given a sequence h , $|h|$ denotes its length, $h(i)$ its i -th element, and h^i the sequence containing the first i elements of h . We

²Without loss of generality, we focus only on boolean queries [2]. We can support non-boolean queries as follows. Given a database state s and a query $q := \{\bar{x} \mid \phi\}$, if the inference control mechanism authorizes the boolean query $\bigwedge_{\bar{t} \in [q]^s} \phi[\bar{x} \mapsto \bar{t}] \wedge (\forall \bar{x}. \phi \rightarrow \bigvee_{\bar{t} \in [q]^s} \bar{x} = \bar{t})$, then we return q ’s result, and otherwise we reject q as unauthorized.

also denote by h^0 the empty sequence ϵ , and \cdot denotes the concatenation operator.

We now formalize Policy Decision Points. A *C-PDP* is a function $f : \Omega_C \times \mathcal{Q}_C \times \mathcal{H}_C \rightarrow \{\top, \perp\}$ taking as input a system state, a query, and a history and returning the security decision, accept (\top) or deny (\perp).

Let C be a system configuration, $s = \langle db, U, P \rangle$ be a C -state, and f be a C -PDP. A C -history h is *compatible with s and f* iff for each $1 \leq i \leq |h|$, (1) $f(s, q(h(i)), h^{i-1}) = a(h(i))$, (2) if $a(h(i)) = \perp$, then $res(h(i)) = \dagger$, and (3) if $a(h(i)) = \top$, then $res(h(i)) = [\phi]^{db}$, where $q(h(i)) = \langle u, \phi \rangle$. In other words, h is compatible with s and f iff it was generated by the PDP f starting in state s .

A (C, f) -run is a pair $\langle s, h \rangle$, where s is a system state in Ω_C and h is a history in \mathcal{H}_C compatible with s and f . Since all queries are SELECT queries, the system state does not change along the run. Hence, our runs consist of a state and a history instead of e.g., an alternating sequence of states and actions (as is standard for runs). We denote by $\text{runs}(C, f)$ the set of all (C, f) -runs. Furthermore, given a run $r = \langle \langle db, U, P \rangle, h \rangle$, we denote by r^i the run $\langle \langle db, U, P \rangle, h^i \rangle$, and we use dot notation to access to r ’s components. For instance, $r.db$ denotes the database state db and $r.h$ denotes the history.

Example IV.3. Consider the run $r = \langle \langle db, U, P \rangle, h \rangle$, where the database state db is the state $s_{\{\text{A,B,C}\}}$, where *Alice*, *Bob*, and *Carl* have cancer, the policy P is defined in Example IV.2, the set of users U contains only *Mallory*, and the history h is as follows (here we assume that all queries are authorized):

- 1) *Mallory* checks whether *Carl* smokes. Thus, $h(1) = \langle \langle \text{Mallory}, \text{smokes}(\text{Carl}) \rangle, \top, \top \rangle$.
- 2) *Mallory* checks whether *Carl* is *Alice*’s and *Bob*’s son. Therefore, $h(2)$ is $\langle \langle \text{Mallory}, \text{father}(\text{Bob}, \text{Carl}) \wedge \text{mother}(\text{Alice}, \text{Carl}) \rangle, \top, \top \rangle$.
- 3) *Mallory* checks whether *Alice* has cancer. Thus, $h(3) = \langle \langle \text{Mallory}, \text{cancer}(\text{Alice}) \rangle, \top, \top \rangle$.
- 4) *Mallory* checks whether *Bob* has cancer. Thus, $h(4) = \langle \langle \text{Mallory}, \text{cancer}(\text{Bob}) \rangle, \top, \top \rangle$. ■

D. Attacker Model

To reason about DBIC, it is essential to precisely define (1) how users interact with the system, (2) how they reason about the system’s behaviour, (3) their initial beliefs about the database state, and (4) how these beliefs change by observing the system’s behaviour. We formalize this in an attacker model.

Each user has an initial belief about the database state. Following [18], [19], [29], [51], we represent a user’s beliefs as a probability distribution over all database states. Furthermore, users observe the system’s behaviour and derive information about the database content. We formalize a user’s observations as an equivalence relation over runs, where two runs are equivalent iff the user’s observations are the same in both runs, as is standard in information-flow [7], [8]. A user’s knowledge is the set of all database states that he considers possible given his observations. Finally, we use Bayesian conditioning to update a user’s beliefs given his knowledge.

Let $C = \langle D, \Gamma \rangle$ be a system configuration and f be a C -PDP. A C -probability distribution is a discrete probability distribution given by a function $P : \Omega_D^\Gamma \rightarrow [0, 1]$ such that $\sum_{db \in \Omega_D^\Gamma} P(db) = 1$. Given a set $E \subseteq \Omega_D^\Gamma$, $P(E)$ denotes $\sum_{s \in E} P(s)$. Furthermore, given two sets $E', E'' \subseteq \Omega_D^\Gamma$ such that $P(E') \neq 0$, $P(E'' \mid E')$ denotes $P(E'' \cap E')/P(E')$ as is standard. We denote by \mathcal{P}_C the set of all possible C -probability distributions. Abusing notation, we extend probability distributions to formulae: $P(\psi) = P(\llbracket \psi \rrbracket)$, where $\llbracket \psi \rrbracket = \{db \in \Omega_D^\Gamma \mid [\psi]^{db} = \top\}$.

We now introduce *indistinguishability*, an equivalence relation used in information-flow control [41]. Let C be a system configuration and f be a C -PDP. Given a history h and a user $u \in \mathcal{U}$, $h|_u$ denotes the history obtained from h by removing all C -events from users other than u , namely $\epsilon|_u = \epsilon$, and if $h = \langle \langle u', q \rangle, a, res \rangle \cdot h'$, then $h|_u = h'|_u$ in case $u \neq u'$, and $h|_u = \langle \langle u, q \rangle, a, res \rangle \cdot h'|_u$ if $u = u'$. Given two runs $r = \langle \langle db, U, P \rangle, h \rangle$ and $r' = \langle \langle db', U', P' \rangle, h' \rangle$ in $runs(C, f)$ and a user $u \in \mathcal{U}$, we say that r and r' are *indistinguishable for u* , written $r \sim_u r'$, iff $h|_u = h'|_u$. This means that r and r' are indistinguishable for a user u iff the system's behaviour in response to u 's commands is the same in both runs. Note that \sim_u depends on both C and f , which we generally leave implicit. Given a run r , $[r]_{\sim_u}$ is the equivalence class of r with respect to \sim_u , i.e., $[r]_{\sim_u} = \{r' \in runs(C, f) \mid r' \sim_u r\}$, whereas $\llbracket r \rrbracket_{\sim_u}$ is set of all databases associated to the runs in $[r]_{\sim_u}$, i.e., $\llbracket r \rrbracket_{\sim_u} = \{db \mid \exists U, P, h. \langle \langle db, U, P \rangle, h \rangle \in [r]_{\sim_u}\}$.

Definition IV.1. Let $C = \langle D, \Gamma \rangle$ be a configuration and f be a C -PDP. A (C, f) -attacker model is a function $ATK : \mathcal{U} \rightarrow \mathcal{P}_C$ associating to each user $u \in \mathcal{U}$ a C -probability distribution representing u 's initial beliefs. Additionally, for all users $u \in \mathcal{U}$ and all states $s \in \Omega_D^\Gamma$, we require that $ATK(u)(s) > 0$. The *semantics of ATK* is $\llbracket ATK \rrbracket(u, r) = \lambda s \in \Omega_D^\Gamma. ATK(u)(s \mid \llbracket r \rrbracket_{\sim_u})$, where $u \in \mathcal{U}$ and $r \in runs(C, f)$. \square

The semantics of an attacker model ATK associates to each user u and each run r the probability distribution obtained by updating u 's initial beliefs given his knowledge with respect to the run r . We informally refer to $\llbracket ATK \rrbracket(u, r)(\llbracket \phi \rrbracket)$ as u 's beliefs in a sentence ϕ (given a run r).

Example IV.4. The attacker model for the example from §II is as follows. Let X_{Alice} , X_{Bob} , and X_{Car1} be three boolean random variables, representing the probability that the corresponding patient has cancer. They define the following joint probability distribution, which represents a user's initial beliefs about the actual database state: $P(X_{\text{Alice}}, X_{\text{Bob}}, X_{\text{Car1}}) = P(X_{\text{Alice}}) \cdot P(X_{\text{Bob}}) \cdot P(X_{\text{Car1}} \mid X_{\text{Alice}}, X_{\text{Bob}})$. The probability distributions of these variables are given in Figure 3 and they are derived from the probabilistic model in §II. We associate each outcome (x, y, z) of $X_{\text{Alice}}, X_{\text{Bob}}, X_{\text{Car1}}$ with the corresponding database state s_C , where C is the set of patients such that the outcome of the corresponding variable is \top . For each user $u \in \mathcal{U}$, the distribution P_u is defined as $P_u(s_C) = P(X_{\text{Alice}} = x, X_{\text{Bob}} = y, X_{\text{Car1}} = z)$,

	X_{Alice}		X_{Bob}		X_{Car1}
\top	$1/20$	\top	\top	$12/20$	$8/20$
\perp	$19/20$	\top	\perp	$9/20$	$11/20$
	X_{Bob}		X_{Car1}		
\top	$6/20$	\perp	\top	$9/20$	$11/20$
\perp	$14/20$	\perp	\perp	$6/20$	$14/20$

Fig. 3: Probability distribution for the random variables X_{Alice} , X_{Bob} , and X_{Car1} from Example IV.4.

State	Probability	State	Probability
s_\emptyset	0.4655	$s_{\{A,B\}}$	0.006
$s_{\{A\}}$	0.01925	$s_{\{A,C\}}$	0.01575
$s_{\{B\}}$	0.15675	$s_{\{B,C\}}$	0.12825
$s_{\{C\}}$	0.1995	$s_{\{A,B,C\}}$	0.009

Fig. 4: Probability distribution over all database states. Each state is denoted as s_C , where C is the content of the *cancer* table. Here we denote the patients' names with their initials.

where x (respectively y and z) is \top if Alice (respectively Bob and Car1) is in C and \perp otherwise. Figure 4 shows the probabilities associated with each state in Ω_D^Γ , i.e., a user's initial beliefs. Finally, the attacker model is $ATK = \lambda u \in \mathcal{U}. P_u$. \blacksquare

E. Confidentiality

We first define the notion of a secrecy-preserving run for a secret $\langle u, \phi, l \rangle$ and an attacker model ATK . Informally, a run r is secrecy-preserving for $\langle u, \phi, l \rangle$ iff whenever an attacker's belief in the secret ϕ is below the threshold l , then there is no way for the attacker to increase his belief in ϕ above the threshold. Our notion of secrecy-preserving runs is inspired by existing security notions for query auditing [29].

Definition IV.2. Let $C = \langle D, \Gamma \rangle$ be a configuration, f be a C -PDP, and ATK be a (C, f) -attacker model. A run r is *secrecy-preserving* for a secret $\langle u, \phi, l \rangle$ and ATK iff for all $0 \leq i < |r|$, $\llbracket ATK \rrbracket(u, r^i)(\phi) < l$ implies $\llbracket ATK \rrbracket(u, r^{i+1})(\phi) < l$. \square

We now formalize our confidentiality notion. A PDP provides data confidentiality for an attacker model ATK iff all runs are secrecy-preserving for ATK . Note that our security notion can be seen as a probabilistic generalization of opacity [60] for the database setting. Our notion is also inspired by the semantics of knowledge-based policies [51].

Definition IV.3. Let $C = \langle D, \Gamma \rangle$ be a system configuration, f be a C -PDP, and ATK be a (C, f) -attacker model. We say that the PDP f *provides data confidentiality with respect to C and ATK* iff for all runs $r = \langle \langle db, U, P \rangle, h \rangle$ in $runs(C, f)$, for all users $u \in \mathcal{U}$, for all secrets $s \in secrets(P, u)$, r is secrecy-preserving for s and ATK . \square

A PDP providing confidentiality ensures that if an attacker's initial belief in a secret ϕ is below the corresponding threshold, then there is no way for the attacker to increase his belief in ϕ above the threshold by interacting with the system. This guarantee does not however apply to *trivial non-secrets*, i.e., those secrets an attacker knows with a probability at least the threshold even before interacting with the system. No PDP can

prevent their disclosure since the disclosure does not depend on the attacker’s interaction with the database.

Example IV.5. Let r be the run given in Example IV.3, ATK be the attacker model in Example IV.4, and u be the user *Mallory*. In the following, ϕ_1 , ϕ_2 , and ϕ_3 denote $cancer(Carl)$, $cancer(Bob)$, and $cancer(Alice)$ respectively, i.e., the three secrets from Example IV.2. Furthermore, we assume that the policy contains an additional secret $\langle Mallory, \phi_4, 1/2 \rangle$, where $\phi_4 := \neg cancer(Alice)$.

Figure 5 illustrates *Mallory*’s beliefs about ϕ_1, \dots, ϕ_4 and whether the run is secrecy-preserving for the secrets ϕ_1, \dots, ϕ_4 . The probabilities in the tables can be obtained by combining the states in $\llbracket r^i \rrbracket_{\sim u}$, for $0 \leq i \leq 4$, and $\llbracket \phi_j \rrbracket$, for $1 \leq j \leq 4$, with the probabilities from Figure 4. As shown in Figure 5, the run is not secrecy-preserving for the secrets ϕ_1 and ϕ_2 as it completely discloses that *Alice* and *Bob* have cancer, in the third and fourth steps respectively. Secrecy-preservation is also violated for the secret ϕ_1 , even though r does not directly disclose any information about *Carl*’s health status. Indeed, in the last step of the run, *Mallory*’s belief in ϕ_1 is 0.6, which is higher than the threshold $1/2$, even though his belief in ϕ_1 before learning that *Bob* had cancer was below the threshold. Note that ϕ_4 is a trivial non-secret: even before interacting with the system, *Mallory*’s belief in ϕ_4 is 0.95. ■

F. Discussion

Our approach assumes that the attacker’s capabilities are well-defined. While this, in general, is a strong assumption, there are many domains where such information is known. There are, however, domains where this information is lacking. In these cases, security engineers must (1) determine the appropriate beliefs capturing the desired attacker models, and (2) formalize them. The latter can be done, for instance, using $ATKLOG$ (see §V). Note however that precisely eliciting the attackers’ capabilities is still an open problem in $DBIC$.

V. $ATKLOG$

A. Probabilistic Logic Programming

$PROBLOG$ [20], [21], [31] is a probabilistic logic programming language with associated tool support. An exact inference engine for $PROBLOG$ is available at [1].

Conventional logic programs are constructed from terms, atoms, literals, and rules. In the following, we consider only function-free logic programs, also called $DATALOG$ programs. In this setting, terms are either variable identifiers or constants.

Let Σ be a first-order signature, \mathbf{dom} be a finite domain, and Var be a countably infinite set of variable identifiers. A (Σ, \mathbf{dom}) -atom $R(v_1, \dots, v_n)$ consists of a predicate symbol $R \in \Sigma$ and arguments v_1, \dots, v_n such that n is the arity of R , and each v_i , for $1 \leq i \leq n$, is either a variable identifier in Var or a constant in \mathbf{dom} . We denote by $\mathcal{A}_{\Sigma, \mathbf{dom}}$ the set $\{R(v_1, \dots, v_{|R|}) \mid R \in \Sigma \wedge v_1, \dots, v_{|R|} \in \mathbf{dom} \cup Var\}$ of all (Σ, \mathbf{dom}) -atoms. A (Σ, \mathbf{dom}) -literal l is either a (Σ, \mathbf{dom}) -atom a or its negation $\neg a$, where $a \in \mathcal{A}_{\Sigma, \mathbf{dom}}$. We denote by $\mathcal{L}_{\Sigma, \mathbf{dom}}$ the set $\mathcal{A}_{\Sigma, \mathbf{dom}} \cup \{\neg a \mid a \in \mathcal{A}_{\Sigma, \mathbf{dom}}\}$ of

all (Σ, \mathbf{dom}) -literals. Given a literal l , $vars(l)$ denotes the set of its variables, $args(l)$ the list of its arguments, and $pred(l)$ the predicate symbol used in l . As is standard, we say that a literal l is *positive* if it is an atom in $\mathcal{A}_{\Sigma, \mathbf{dom}}$ and *negative* if it is the negation of an atom. Furthermore, we say that a literal l is *ground* iff $vars(l) = \emptyset$.

A (Σ, \mathbf{dom}) -rule is of the form $h \leftarrow l_1, \dots, l_n, e_1, \dots, e_m$, where $h \in \mathcal{A}_{\Sigma, \mathbf{dom}}$ is a (Σ, \mathbf{dom}) -atom, $l_1, \dots, l_n \in \mathcal{L}_{\Sigma, \mathbf{dom}}$ are (Σ, \mathbf{dom}) -literals, and e_1, \dots, e_m are equality and inequality constraints over the variables in h, l_1, \dots, l_n .³ Given a rule r , we denote by $head(r)$ the atom h , by $body(r)$ the literals l_1, \dots, l_n , by $ctr(r)$ the constraints e_1, \dots, e_m , and by $body(r, i)$ the i -th literal in r ’s body, i.e., $body(r, i) = l_i$. Furthermore, we denote by $body^+(r)$ (respectively $body^-(r)$) all positive (respectively negative) literals in $body(r)$. As is standard, we assume that the free variables in a rule’s head are a subset of the free variables of the positive literals in the rule’s body, i.e., $vars(head(r)) \subseteq \bigcup_{l \in body^+(r)} vars(l) \cup \bigcup_{(x=c) \in ctr(r) \wedge c \in \mathbf{dom}} \{x\}$. Finally, a (Σ, \mathbf{dom}) -logic program is a set of (Σ, \mathbf{dom}) -ground atoms and (Σ, \mathbf{dom}) -rules. We consider only programs p that do not contain negative cycles in the rules as is standard for stratified $DATALOG$ [2].

To reason about probabilities, $PROBLOG$ extends logic programming with probabilistic atoms. A (Σ, \mathbf{dom}) -probabilistic atom is a (Σ, \mathbf{dom}) -atom a annotated with a value $0 \leq v \leq 1$, denoted $v::a$. $PROBLOG$ supports both probabilistic ground atoms and rules having probabilistic atoms in their heads. $PROBLOG$ also supports *annotated disjunctions* $v_1::a_1; \dots; v_n::a_n$, where a_1, \dots, a_n are ground atoms and $\left(\sum_{1 \leq i \leq n} v_i\right) \leq 1$, which denote that a_1, \dots, a_n are mutually exclusive probabilistic events happening with probabilities v_1, \dots, v_n . Annotated disjunctions can either be used as ground atoms or as heads in rules. Both annotated disjunctions and probabilistic rules are just syntactic sugar and can be expressed using ground probabilistic atoms and standard rules [20], [21], [31]; see Appendix A.

A (Σ, \mathbf{dom}) - $PROBLOG$ program p defines a probability distribution over all possible (Σ, \mathbf{dom}) -structures, denoted $\llbracket p \rrbracket$. Note that we consider only function-free $PROBLOG$ programs. Hence, in our setting, $PROBLOG$ is a probabilistic extension of $DATALOG$. Appendix A contains a formal account of $PROBLOG$ ’s semantics.

Medical Data. We formalize the probability distribution from Example IV.4 as a $PROBLOG$ program. We reuse the database schema and the domain from Example IV.1 as the first-order signature and the domain for the $PROBLOG$ program. First, we encode the template shown in Figure 2 using ground atoms: $patient(Alice)$, $patient(Bob)$, $patient(Carl)$, $smokes(Bob)$, $smokes(Carl)$, $father(Bob, Carl)$, and $mother(Alice, Carl)$. Second, we encode the probability distribution associated with the possible values of the *cancer* table using the following $PROBLOG$ rules, which

³Without loss of generality, we assume that equality constraints involving a variable v and a constant c are of the form $v = c$.

i	$\llbracket \text{ATK} \rrbracket(u, r^i)(\llbracket \phi \rrbracket)$				$\llbracket \text{ATK} \rrbracket(u, r^{i+1})(\llbracket \phi \rrbracket)$				Secrecy			
	ϕ_1	ϕ_2	ϕ_3	ϕ_4	ϕ_1	ϕ_2	ϕ_3	ϕ_4	ϕ_1	ϕ_2	ϕ_3	ϕ_4
0	0.3525	0.3	0.05	0.95	0.3525	0.3	0.05	0.95	✓	✓	✓	*
1	0.3525	0.3	0.05	0.95	0.3525	0.3	0.05	0.95	✓	✓	✓	*
2	0.3525	0.3	0.05	0.95	0.495	0.3	1	0	✓	✓	✗	*
3	0.495	0.3	1	0	0.6	1	1	0	✗	✗	✗	*
4	0.6	1	1	0	–	–	–	–	–	–	–	–

Fig. 5: Evolution of *Mallory*'s beliefs in the secrets ϕ_1, \dots, ϕ_4 for the run r and the attacker model *ATK* from Example IV.5. In the table, \mathcal{X} and \checkmark denote that secrecy-preservation is violated and satisfied respectively, whereas $*$ denotes trivial secrets.

have probabilistic atoms in their heads:

$$\begin{aligned}
1/20::\text{cancer}(x) &\leftarrow \text{patient}(x) \\
5/19::\text{cancer}(x) &\leftarrow \text{smokes}(x) \\
3/14::\text{cancer}(y) &\leftarrow \text{father}(x, y), \text{cancer}(x), \\
&\quad \text{mother}(z, y), \neg \text{cancer}(z) \\
3/14::\text{cancer}(y) &\leftarrow \text{father}(x, y), \neg \text{cancer}(x), \\
&\quad \text{mother}(z, y), \text{cancer}(z) \\
3/7::\text{cancer}(y) &\leftarrow \text{father}(x, y), \text{cancer}(x), \\
&\quad \text{mother}(z, y), \text{cancer}(z)
\end{aligned}$$

The coefficients in the above example are derived from §II. For instance, the probability that a smoking patient x whose parents are not in the *cancer* relation has cancer is 30%. The coefficient in the first rule is $1/20$ since each patient has a 5% probability of having cancer. The coefficient in the second rule is $5/19$, which is $(6/20 - 1/20) \cdot (1 - 1/20)^{-1}$, i.e., the probability that *cancer*(x) is derived from the second rule given that it has not been derived from the first rule. This ensures that the overall probability of deriving *cancer*(x) is $6/20$, i.e., 30%. The coefficients for the last two rules are derived analogously.

Informally, a probabilistic ground atom $1/2::\text{cancer}(\text{Bob})$ expresses that *cancer*(Bob) holds with a probability $1/2$. Similarly, the rule $1/20::\text{cancer}(x) \leftarrow \text{patient}(x)$ states that, for any x such that *patient*(x) holds, then *cancer*(x) can be derived with probability $1/20$. This program yields the probability distribution shown in Figure 4.

B. ATKLOG's Foundations

We first introduce belief programs, which formalize an attacker's initial beliefs. Afterwards, we formalize ATKLOG.

Belief Programs. A belief program formalizes an attacker's beliefs as a probability distribution over the database states.

A database schema $D' = \langle \Sigma', \text{dom} \rangle$ extends a schema $D = \langle \Sigma, \text{dom} \rangle$ iff Σ' contains all relation schemas in Σ . The extension D' may extend Σ with additional predicate symbols necessary to encode probabilistic dependencies. Given an extension D' , a D' -state s' agrees with a D -state s iff $s'(R) = s(R)$ for all R in D . Given a D -state s , we denote by $\text{EXT}(s, D, D')$ the set of all D' -states that agree with s .

A (Σ', dom) -PROBLOG program p , where $D' = \langle \Sigma', \text{dom} \rangle$ extends D , is a *belief program over D*. The *D-semantics of p* is $\llbracket p \rrbracket_D = \lambda s \in \Omega_D. \sum_{s' \in \text{EXT}(s, D, D')} \llbracket p \rrbracket(s')$. Given a system configuration $C = \langle D, \Gamma \rangle$, a belief program

p over D complies with C iff $\llbracket p \rrbracket_D$ is a C -probability distribution. With a slight abuse of notation, we lift the semantics of belief programs to sentences: $\llbracket p \rrbracket_D = \lambda \phi \in \mathcal{RC}_{\text{bool}}. \sum_{s' \in \{s \in \Omega_D \mid \llbracket \phi \rrbracket^s = \top\}} \llbracket p \rrbracket_D(s')$.

ATKLOG. An ATKLOG model specifies the initial beliefs of all users in \mathcal{U} using belief programs.

Let D be a database schema and $C = \langle D, \Gamma \rangle$ be a system configuration. A C -ATKLOG model *ATK* is a function associating to each user $u \in U$, where $U \subset \mathcal{U}$ is a finite set of users, a belief program p_u and to all users $u \in \mathcal{U} \setminus U$ a belief program p_0 , such that for all users $u \in \mathcal{U}$, $\llbracket \text{ATK}(u) \rrbracket_D$ complies with C and for all database states $s \in \Omega_D^\Gamma$, $\llbracket \text{ATK}(u) \rrbracket_D(s) > 0$, i.e., all database states satisfying the integrity constraints are possible. Informally, a C -ATKLOG model associates a distinct belief program to each user in U , and it associates to each user in $\mathcal{U} \setminus U$ the same belief program p_0 .

Given a C -PDP f , a C -ATKLOG model *ATK* defines the (C, f) -attacker model $\lambda u \in \mathcal{U}. \llbracket \text{ATK}(u) \rrbracket_D$ that associates to each user $u \in \mathcal{U}$ the probability distribution defined by the belief program *ATK*(u). The semantics of this (C, f) -attacker model is: $\lambda u \in \mathcal{U}. \lambda r \in \text{runs}(C, f). \lambda s \in \Omega_D^\Gamma. \llbracket \text{ATK}(u) \rrbracket_D(s \mid \llbracket r \rrbracket_{\sim_u})$. Informally, given a C -ATKLOG model *ATK*, a C -PDP f , and a user u , u 's belief in a database state s , given a run r , is obtained by conditioning the probability distribution defined by the belief program *ATK*(u) given the set of database states corresponding to all runs $r' \sim_u r$.

VI. TRACTABLE INFERENCE FOR PROBLOG PROGRAMS

Probabilistic inference in PROBLOG is intractable in general. Its data complexity, i.e., the complexity of inference when only the programs' probabilistic ground atoms are part of the input and the rules are considered fixed and not part of the input, is $\#P$ -hard; see Appendix B. This limits the practical applicability of PROBLOG (and ATKLOG) for DBIC. To address this, we define acyclic PROBLOG programs, a class of programs where the data complexity of inference is PTIME.

Given a PROBLOG program p , our inference algorithm consists of three steps: (1) we compute all of p 's derivations, (2) we compile these derivations into a Bayesian Network (BN) bn , and (3) we perform the inference over bn . To ensure tractability, we leverage two key insights. First, we exploit guarded negation [9] to develop a sound over-approximation, called the relaxed grounding, of all derivations of a program that is independent of the presence (or absence) of the probabilistic atoms. This ensures that whenever a ground atom can be derived from a program (for a possible

assignment to the probabilistic atoms), the atom is also part of this program's relaxed grounding. This avoids grounding p for each possible assignment to the probabilistic atoms. Second, we introduce syntactic constraints (acyclicity) that ensure that bn is a forest of poly-trees. This ensures tractability since inference for poly-tree BNs can be performed in polynomial time in the network's size [47].

We also precisely characterize the expressiveness of acyclic PROBLOG programs. In this respect, we prove that acyclic programs are as expressive as forests of poly-tree BNs, one of the few classes of BNs with tractable inference.

As mentioned in §V, probabilistic rules and annotated disjunctions are just syntactic sugar. Hence, in the following we consider PROBLOG programs consisting just of probabilistic ground atoms and non-probabilistic rules. Note also that we treat ground atoms as rules with an empty body.

A. Preliminaries

Negation-guarded Programs. A rule r is *negation-guarded* [9] iff all the variables occurring in negative literals also occur in positive literals, namely for all negative literals l in $body^-(r)$, $vars(l) \subseteq \bigcup_{l' \in body^+(r)} vars(l')$. To illustrate, the rule $C(x) \leftarrow A(x), \neg B(x)$ is negation-guarded, whereas $C(x) \leftarrow A(x), \neg B(x, y)$ is not since the variable y does not occur in any positive literal. We say that a program p is *negation-guarded* if all rules $r \in p$ are.

Relaxed Grounding. The relaxed grounding of a program p is obtained by considering all probabilistic atoms as certain and by grounding all positive literals. For all negation-guarded programs, the relaxed grounding of p is a sound over-approximation of all possible derivations in p . Given a program p and a rule $r \in p$, $rg(p)$ denotes p 's relaxed grounding and $rg(p, r)$ denotes the set of r 's ground instances. We formalize relaxed groundings in Appendix B.

Example VI.1. Let p be the program consisting of the facts $1/2::A(1), A(2), A(3), D(1), E(2), F(1), O(1, 2)$, and $2/3::O(2, 3)$, and the rules $r_a = B(x) \leftarrow A(x), D(x)$, $r_b = B(x) \leftarrow A(x), E(x)$, and $r_c = B(y) \leftarrow B(x), \neg F(x), O(x, y)$. The relaxed grounding of p consists of the initial facts together with $B(1), B(2)$, and $B(3)$, whereas $rg(p, r_c)$ consists of $B(2) \leftarrow B(1), \neg F(1), O(1, 2)$ and $B(3) \leftarrow B(2), \neg F(2), O(2, 3)$. ■

Dependency and Ground Graphs. The *dependency graph* of a program p , denoted $graph(p)$, is the directed labelled graph having as nodes all the predicate symbols in p and having an edge $a \xrightarrow{r, i} b$ iff there is a rule r such that a occurs in i -th literal in r 's body and b occurs in r 's head. Figure 6 depicts the dependency graph from Example VI.1. The *ground graph* of a program p is the graph obtained from its relaxed grounding. Hence, there is an edge $a \xrightarrow{r, gr, i} b$ from the ground atom a to the ground atom b iff there is a rule r and a ground rule $gr \in rg(p, r)$ such that $body(gr, i) \in \{a, \neg a\}$ and $head(gr) = b$. Figure 7 depicts the ground graph from Example VI.1. Note that there are no incoming or

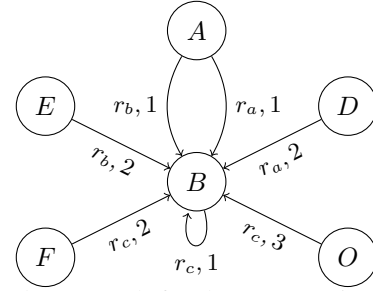


Fig. 6: Dependency graph for the program in Example VI.1.

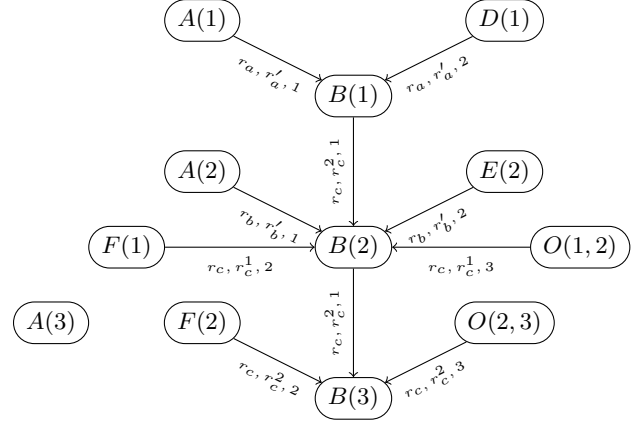


Fig. 7: Ground graph for the program in Example VI.1. The ground rules r'_a , r'_b , r_c^1 , and r_c^2 are as follows: $r'_a = B(1) \leftarrow A(1), D(1)$, $r'_b = B(2) \leftarrow A(2), E(2)$, $r_c^1 = B(2) \leftarrow B(1), \neg F(1), O(1, 2)$, and $r_c^2 = B(3) \leftarrow B(2), \neg F(2), O(2, 3)$.

outgoing edges from $A(3)$ because the node is not involved in any derivation.

Propagation Maps. We use propagation maps to track how information flows inside rules. Given a rule r and a literal $l \in body(r)$, the (r, l) -vertical map is the mapping μ from $\{1, \dots, |l|\}$ to $\{1, \dots, |head(r)|\}$ such that $\mu(i) = j$ iff $args(l)(i) = args(head(r))(j)$ and $args(l)(i) \in Var$. Given a rule r and literals l and l' in r 's body, the (r, l, l') -horizontal map is the mapping μ from $\{1, \dots, |l|\}$ to $\{1, \dots, |l'|\}$ such that $\mu(i) = j$ iff $args(l)(i) = args(l')(j)$ and $args(l)(i) \in Var$.

We say that a path links to a literal l if information flows along the rules to l . This can be formalized by posing constraints on the mapping obtained by combining horizontal and vertical maps along the path. Formally, given a literal l and a mapping $\nu : \mathbb{N} \rightarrow \mathbb{N}$, a directed path $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ ν -downward links to l iff there is a $0 \leq j < n - 1$ such that the function $\mu := \mu' \circ \mu_j \circ \dots \circ \mu_1$ satisfies $\mu(k) = \nu(k)$ for all k for which $\nu(k)$ is defined, where for $1 \leq h \leq j$, μ_h is the $(r_h, body(r_h, i_h))$ -vertical map, and μ' is the horizontal map connecting $body(r_{j+1}, i_{j+1})$ with l . Similarly, a directed path $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ ν -upward links to l iff there is a $1 \leq j \leq n - 1$ such that the function $\mu := \mu'^{-1} \circ \mu_{j+1}^{-1} \circ \dots \circ \mu_n^{-1}$ satisfies $\mu(k) = \nu(k)$ for all k for which $\nu(k)$ is defined, where μ_h is the $(r_h, body(r_h, i_h))$ -vertical map, for $j < h \leq n - 1$, and μ' is the (r_j, l) -vertical map. A path P links to a predicate symbol a iff there

is an atom $a(\bar{x})$ such that P links to $a(\bar{x})$.

Example VI.2. The horizontal map connecting $A(x)$ and $D(x)$ in r_a , i.e., the $(r_a, A(x), D(x))$ -horizontal map, is $\{1 \mapsto 1\}$. The horizontal map connecting $A(x)$ and $E(x)$ in r_b is $\{1 \mapsto 1\}$ as well. Hence, the path $A \xrightarrow{r_a, 1} B$ downward links to D and the path $A \xrightarrow{r_b, 1} B$ downward links to E for the mapping $\{1 \mapsto 1\}$. Furthermore, the path $B \xrightarrow{r_c, 1} B$ downward links to O for $\{1 \mapsto 1\}$ since the $(r_c, B(x), O(x, y))$ -horizontal map is $\{1 \mapsto 1\}$. Finally, the path $B \xrightarrow{r_c, 1} B$ upward links to O for $\{2 \mapsto 1\}$ since the $(r_c, O(x, y))$ -vertical map is $\{2 \mapsto 1\}$. ■

B. Acyclic PROLOG programs

A sufficient condition for tractable inference is that p 's ground graph is a forest of poly-trees. This requires that p 's ground graph neither contains directed nor undirected cycles, or, equivalently, the undirected version of p 's ground graph is acyclic. To illustrate, the ground graph in Figure 7 is a poly-tree. The key insight here is that a cycle among ground atoms is caused by a (directed or undirected) cycle among p 's predicate symbols. In a nutshell, acyclicity requires that all possible cycles in $graph(p)$ are guarded. This ensures that cycles in $graph(p)$ do not lead to cycles in the ground graph. Additionally, acyclicity requires that programs are negation-guarded. This ensures that the relaxed grounding and the ground graph are well-defined. In the following, let p be a (Σ, \mathbf{dom}) -PROLOG program.

Annotations. Annotations represent properties of the relations induced by the program p , and they are syntactically derived by analysing p 's ground atoms and rules.

Let $a, a' \in \Sigma$ be two predicate symbols such that $|a| = |a'|$. A *disjointness annotation* $DIS(a, a')$ represents that the relations induced by a and a' (given p 's relaxed grounding) are disjoint. We say that $DIS(a, a')$ can be derived from p iff no rules in p contain a or a' in their heads, and there is no $\bar{v} \in \mathbf{dom}^{|a|}$ where both $a(\bar{v})$ and $a'(\bar{v})$ appear as (possibly probabilistic) ground atoms in p . Hence, the relations induced by a and a' are disjoint.

Let $n \in \mathbb{N}$ and $A \subseteq \Sigma$ be a set of predicate symbols such that $|a| = 2n$ for all $a \in A$. An *ordering annotation* $ORD(A)$ represents that the transitive closure of the union of the relations induced by predicates in A given p 's relaxed grounding is a strict partial order over \mathbf{dom}^n . The annotation $ORD(A)$ can be derived from the program p iff there is no rule $r \in p$ that contains any of the predicates in A in its head and the transitive closure of $\bigcup_{a \in A} \{((v_1, \dots, v_n), (v_{n+1}, \dots, v_{2n})) \mid \exists k. k::a(v_1, \dots, v_{2n}) \in p\}$ is a strict partial order over \mathbf{dom}^n . Hence, the closure of the relation $\bigcup_{a \in A} \{((v_1, \dots, v_n), (v_{n+1}, \dots, v_{2n})) \mid a(v_1, \dots, v_{2n}) \in rg(p)\}$ induced by the relaxed grounding is a strict partial order.

Let $a \in \Sigma$ be a predicate symbol and $K \subseteq \{1, \dots, |a|\}$. A *uniqueness annotation* $UNQ(a, K)$ represents that the attributes in K are a primary key for the relation induced by a given the relaxed grounding. We say that $UNQ(a, K)$ can be derived from a program p iff no rule contains a in its head

and for all $\bar{v}, \bar{v}' \in \mathbf{dom}^{|a|}$, if (1) $\bar{v}(i) = \bar{v}'(i)$ for all $i \in K$, and (2) there are k and k' such $k::a(\bar{v}) \in p$ and $k'::a(\bar{v}') \in p$, then $\bar{v} = \bar{v}'$. This ensures that whenever $a(\bar{v}), a(\bar{v}') \in rg(p)$ and $\bar{v}(i) = \bar{v}'(i)$ for all $i \in K$, then $\bar{v} = \bar{v}'$.

A Σ -template \mathcal{T} is a set of annotations. In Appendix B, we relax our syntactic rules for deriving annotations.

Example VI.3. We can derive $DIS(D, E)$ from the program in Example VI.1 since no rule generates facts for D and E and the relations defined by the ground atoms are $\{1\}$ and $\{2\}$. We can also derive $ORD(\{O\})$ since the relation defined by O 's ground atoms is $\{(1, 2), (2, 3)\}$, whose transitive closure is a strict partial order. Finally, we can derive $UNQ(O, \{1\})$, $UNQ(O, \{2\})$, and $UNQ(O, \{1, 2\})$ since both arguments of O uniquely identify the tuples in the relation induced by O . ■

Unsafe structures. An unsafe structure models a part of the dependency graph that may introduce cycles in the ground graph. We define directed and undirected unsafe structures which may respectively introduce directed and undirected cycles in the ground graph.

A *directed unsafe structure* in $graph(p)$ is a directed cycle C in $graph(p)$. We say that a directed unsafe structure C covers a directed cycle C' iff C is equivalent to C' .

An *undirected unsafe structure* in $graph(p)$ is quadruple $\langle D_1, D_2, D_3, U \rangle$ such that (1) D_1, D_2 , and D_3 are directed paths whereas U is an undirected path, (2) D_1 and D_2 start from the same node, (3) D_2 and D_3 end in the same node, and (4) $D_1 \cdot U \cdot D_3 \cdot D_2$ is an undirected cycle in $graph(p)$. We say that an unsafe structure $\langle D_1, D_2, D_3, U \rangle$ covers an undirected cycle U' in $graph(p)$ iff $D_1 \cdot U \cdot D_3 \cdot D_2$ is equivalent to U' .

Example VI.4. The cycle introduced by the rule r_c is captured by the directed unsafe structure $B \xrightarrow{r_c, 1} B$, while the cycle introduced by r_a and r_b is captured by the structure $\langle A \xrightarrow{r_a, 1} B, A \xrightarrow{r_b, 1} B, \epsilon, \epsilon \rangle$, where ϵ denotes the empty path. ■

Connected Rules. A connected rule r ensures that a grounding of r is fully determined either by the assignment to the head's variables or to the variables of any literal in r 's body. Formally, a strongly connected rule r guarantees that for any two groundings gr', gr'' of r , if $head(gr') = head(gr'')$, then $gr' = gr''$. In contrast, a weakly connected rule r guarantees that for any two groundings gr', gr'' of r , if $body(gr', i) = body(gr'', i)$ for some i , then $gr' = gr''$. This is done by exploiting uniqueness annotations and the rule's structure.

Before formalizing connected rules, we introduce join trees. A join tree represents how multiple predicate symbols in a rule share variables. In the following, let r be a rule and \mathcal{T} be a template. A *join tree for a rule r* is a rooted labelled tree $(N, E, root, \lambda)$, where $N \subseteq body(r)$, E is a set of edges (i.e., unordered pairs over N^2), $root \in N$ is the tree's root, and λ is the labelling function. Moreover, we require that for all $n, n' \in N$, if $n \neq n'$ and $(n, n') \in E$, then $\lambda(n, n') = vars(n) \cap vars(n')$ and $\lambda(n, n') \neq \emptyset$. A join tree $(N, E, root, \lambda)$ covers a literal l iff $l \in N$. Given a join tree $J = (N, E, root, \lambda)$ and a node $n \in N$, the *support of n* , denoted $support(n)$, is the

set $\text{vars}(\text{head}(r)) \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\} \cup \{\text{vars}(n') \mid n' \in \text{anc}(J, n)\}$, where $\text{anc}(J, n)$ is the set of n 's ancestors in J , i.e., the set of all nodes (different from n) on the path from root to n . A join tree $J = (N, E, \text{root}, \lambda)$ is \mathcal{T} -strongly connected iff for all positive literals $l \in N$, there is a set $K \subseteq \{i \mid \bar{x} = \text{args}(l) \wedge \bar{x}(i) \in \text{support}(l)\}$ such that $\text{UNQ}(\text{pred}(l), K) \in \mathcal{T}$ and for all negative literals $l \in N$, $\text{vars}(l) \subseteq \text{support}(l)$. In contrast, a join tree $(N, E, \text{root}, \lambda)$ is \mathcal{T} -weakly connected iff for all $(a(\bar{x}), a'(\bar{x}')) \in E$, there are $K \subseteq \{i \mid \bar{x}(i) \in L\}$ and $K' \subseteq \{i \mid \bar{x}'(i) \in L\}$ such that $\text{UNQ}(a, K), \text{UNQ}(a', K') \in \mathcal{T}$, where $L = \lambda(a(\bar{x}), a'(\bar{x}'))$.

We now formalize strongly and weakly connected rules. A rule r is \mathcal{T} -strongly connected iff there exist \mathcal{T} -strongly connected join trees J_1, \dots, J_n that cover all literals in r 's body. This guarantees that for any two groundings gr', gr'' of r , if $\text{head}(gr') = \text{head}(gr'')$, then $gr' = gr''$.

Given a rule r , a set of literals L , and a template \mathcal{T} , a literal $l \in \text{body}(r)$ is (r, \mathcal{T}, L) -strictly guarded iff (1) $\text{vars}(l) \subseteq \bigcup_{l' \in L \cap \text{body}^+(r)} \text{vars}(l') \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\}$, and (2) there is a positive literal $a(\bar{x}) \in L$ and an annotation $\text{UNQ}(a, K) \in \mathcal{T}$ such that $\{\bar{x}(i) \mid i \in K\} \subseteq \text{vars}(l)$. A rule r is weakly connected for \mathcal{T} iff there exists a \mathcal{T} -weakly connected join tree $J = (N, E, \text{root}, \lambda)$ such that $N \subseteq \text{body}^+(r)$, and all literals in $\text{body}(r) \setminus N$ are (r, \mathcal{T}, N) -strictly guarded. This guarantees that for any two groundings gr', gr'' of r , if $\text{body}(gr', i) = \text{body}(gr'', i)$ for some i , then $gr' = gr''$.

Example VI.5. Let \mathcal{T} be the template from Example VI.3. The rule $r_c := B(y) \leftarrow B(x), \neg F(x), O(x, y)$ is \mathcal{T} -strongly connected. Indeed, the join tree having $O(x, y)$ as root and $B(x)$ and $\neg F(x)$ as leaves is such that (1) there is a uniqueness annotation $\text{UNQ}(O, \{2\})$ in \mathcal{T} such that the second variable in $O(x, y)$ is included in those of r_c 's head, (2) the variables in $B(x)$ and $\neg F(x)$ are a subset of those of their ancestors, and (3) the tree covers all literals in r_c 's body. The rule is also \mathcal{T} -weakly connected: the join tree consisting only of $O(x, y)$ is \mathcal{T} -weakly connected and the literals $B(x)$ and $\neg F(x)$ are strictly guarded. Note that the rules r_a and r_b are trivially both strongly and weakly connected. ■

Guarded undirected structures. Guarded undirected structures ensure that undirected cycles in the dependency graph do not correspond to undirected cycles in the ground graph by exploiting disjointness annotations. Formally, an undirected unsafe structure $\langle D_1, D_2, D_3, U \rangle$ is guarded by a template \mathcal{T} iff either (D_1, D_2) is \mathcal{T} -head-guarded or (D_2, D_3) is \mathcal{T} -tail-guarded.

A pair of non-empty paths (P_1, P_2) sharing the same initial node a is \mathcal{T} -head guarded iff (1) if $P_1 = P_2$, all rules in P_1 are weakly connected for \mathcal{T} , or (2) if $P_1 \neq P_2$, there is an annotation $\text{DIS}(pr, pr') \in \mathcal{T}$, a set $K \subseteq \{1, \dots, |a|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |pr|\}$ such that P_1 ν -downward links to pr and P_2 ν -downward links to pr' . Given two ground paths P'_1 and P'_2 corresponding to P_1 and P_2 , the first condition ensures that $P'_1 = P'_2$ whereas the second ensures that P'_1 or P'_2 are not in the ground graph.

Similarly, a pair of non-empty paths (P_1, P_2) sharing the same final node a is \mathcal{T} -tail guarded iff (1) if $P_1 = P_2$, all rules in P_1 are strongly connected for \mathcal{T} , or (2) if $P_1 \neq P_2$, there is an annotation $\text{DIS}(pr, pr') \in \mathcal{T}$, a set $K \subseteq \{1, \dots, |a|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |pr|\}$, such that P_1 ν -upward links to pr and P_2 ν -upward links to pr' .

Example VI.6. The only non-trivially guarded undirected cycle in the graph from Figure 6 is the one represented by the undirected unsafe structure $\langle A \xrightarrow{r_a, 1} B, A \xrightarrow{r_b, 1} B, \epsilon, \epsilon \rangle$. The structure is guarded since the paths $A \xrightarrow{r_a, 1} B$ and $A \xrightarrow{r_b, 1} B$ are head guarded by $\text{DIS}(D, E)$. Indeed, for the same ground atom $A(v)$, for some $v \in \{1, 2, 3\}$, only one of r_a and r_b can be applied since D and E are disjoint. ■

Guarded directed structures. Guarded directed structures exploit ordering annotations to ensure that directed cycles in the dependency graph do not correspond to directed cycles among ground atoms. A directed unsafe structure $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_n, i_n} pr_n$ is guarded by a template \mathcal{T} iff there is an annotation $\text{ORD}(O) \in \mathcal{T}$, integers $1 \leq y_1 < y_2 < \dots < y_e = n$, literals $o_1(\bar{x}_1), \dots, o_e(\bar{x}_e)$ (where $o_1, \dots, o_e \in O$), a non-empty set $K \subseteq \{1, \dots, |pr_1|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |o|/2\}$ such that for each $0 \leq k < e$, (1) $pr_{y_k} \xrightarrow{r_{y_k}, i_{y_k}} \dots \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν -downward connects to $o_{k+1}(\bar{x}_{k+1})$, and (2) $pr_{y_{k+1}-1} \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν' -upward connects to $o_{k+1}(\bar{x}_{k+1})$, where $\nu'(i) = \nu(x) + |o_1|/2$ for all $1 \leq i \leq |o_1|/2$, and $y_0 = 1$.

Example VI.7. The directed unsafe structure $B \xrightarrow{r_c, 1} B$ is guarded by $\text{ORD}(\{O\})$ in the template from Example VI.3. Indeed, the strict partial order induced by O breaks the cycle among ground atoms belonging to B . In particular, the path $B \xrightarrow{r_c, 1} B$ both downward links and upward links to $O(x, y)$; see Example VI.2. ■

Acyclic Programs. Let p be a negation-guarded program and \mathcal{T} be the template containing all annotations that can be derived from p . We say that p is acyclic iff (a) for all undirected cycles U in $\text{graph}(p)$ that are not directed cycles, there is a \mathcal{T} -guarded undirected unsafe structure that covers U , and (b) for all directed cycles C in $\text{graph}(p)$, there is a \mathcal{T} -guarded directed unsafe structure that covers C . This ensures the absence of cycles in the ground graph.

Proposition VI.1. *Let p be a PROLOG program. If p is acyclic, then the ground graph of p is a forest of poly-trees.*

Example VI.8. The program p from Example VI.1 is acyclic. This is reflected in the ground graph in Figure 7. The program $q = p \cup \{E(1)\}$, however, is not acyclic: we cannot derive $\text{DIS}(D, E)$ from q and the undirected unsafe structure $\langle A \xrightarrow{r_a, 1} B, A \xrightarrow{r_b, 1} B, \epsilon, \epsilon \rangle$ is not guarded. As expected, q 's ground graph contains an undirected cycle between $A(1)$ and $B(1)$, as shown in Figure 8. ■

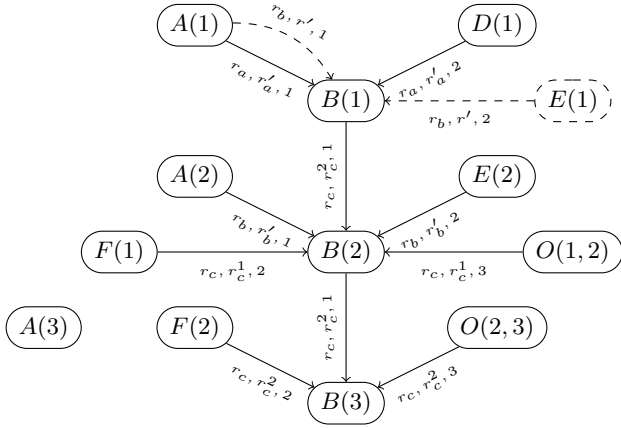
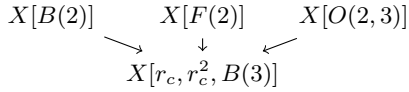


Fig. 8: Ground graph for the program in Example VI.1 extended with the atom $E(1)$. The additional edges and nodes are represented using dashed lines. The ground rules r'_a , r'_b , r_c^1 , and r_c^2 are as in Figure 7, and $r' = B(1) \leftarrow A(1), E(1)$.



Selected CPTs.

			$X[r_c, r'_c, B(3)]$	
$X[B(2)]$	$X[F(2)]$	$X[O(2,3)]$	\top	\perp
\top	\top	\top	0	1
\top	\top	\perp	0	1
\top	\perp	\top	1	0
\top	\perp	\perp	0	1
\perp	\top	\top	0	1
\perp	\top	\perp	0	1
\perp	\perp	\top	0	1
\perp	\perp	\perp	0	1

Fig. 9: Portion of the resulting BN for the atoms $B(2)$, $F(2)$, and $O(2,3)$, the rule $r_c = B(y) \leftarrow B(x), \neg F(x), O(x, y)$, and the ground rule $r_c^2 = B(3) \leftarrow B(2), \neg F(2), O(2,3)$, together with the CPT encoding r_c^2 's semantics.

Expressiveness. Acyclicity trades off the programs expressible in PROBLOG for a tractable inference procedure. Acyclic programs, nevertheless, can encode many relevant probabilistic models.

Proposition VI.2. *Any forest of poly-tree BNs can be represented as an acyclic PROBLOG program.*

To clarify this proposition's scope, observe that poly-tree BNs are one of the few classes of BNs with tractable inference procedures. From Proposition VI.2, it follows that a large class of probabilistic models with tractable inference can be represented as acyclic programs. This supports our thesis that our syntactic constraints are not overly restrictive. In Appendix B, we relax acyclicity to support a limited form of annotated disjunctions and rules sharing a part of their bodies, which are needed to encode the example from §V and for Proposition VI.2. We also provide all the proofs.

C. Inference Engine

Our inference algorithm for acyclic PROBLOG programs consists of three steps: (1) we compute the relaxed grounding of p (cf. §VI-A), (2) we compile the relaxed grounding into a Bayesian Network (BN), and (3) we perform the inference using standard algorithms for poly-tree Bayesian Networks [47].

Encoding as BNs. We compile the relaxed grounding $rg(p)$ into the BN $bn(p)$. The boolean random variables in $bn(p)$ are as follows: (a) for each atom a in $rg(p)$ and ground literal a or $\neg a$ occurring in any $gr \in \bigcup_{r \in p} rg(p, r)$, there is a random variable $X[a]$, (b) for each rule $r \in p$ and each ground atom a such that there is $gr \in rg(p, r)$ satisfying $a = head(gr)$, there is a random variable $X[r, a]$, and (c) for each rule $r \in p$, each ground atom a , and each ground rule $gr \in rg(p, r)$ such that $a = head(gr)$, there is a random variable $X[r, gr, a]$.

The edges in $bn(p)$ are as follows: (a) for each ground atom a , rule r , and ground rule gr , there is an edge from $X[r, gr, a]$ to $X[r, a]$ and an edge from $X[r, a]$ to $X[a]$, and (b) for each ground atoms a and b , rule r , and ground rule gr , there is an edge from $X[b]$ to $X[r, gr, a]$ if b occurs in gr 's body.

Finally, the Conditional Probability Tables (CPTs) of the variables in $bn(p)$ are as follows. The CPT of variables of the form $X[a]$ and $X[r, a]$ is just the OR of the values of their parents, i.e., the value is \top with probability 1 iff at least one of the parents has value \top . For variables of the form $X[r, gr, a]$ such that $body(r) \neq \emptyset$, then the variable's CPT encode the semantics of the rule r , i.e., the value of $X[r, gr, a]$ is \top with probability 1 iff all positive literals have value \top and all negative literals have value \perp . In contrast, for variables of the form $X[r, gr, a]$ such that $body(r) = \emptyset$, then the variable has value \top with probability v and \perp with probability $1 - v$, where r is of the form $v::a$ (if $r = a$ then $v = 1$).

To ensure that the size of the CPT of variables of the form $X[r, a]$ is independent of the size of the relaxed grounding, instead of directly connecting variables of the form $X[r, gr, a]$ with $X[r, a]$, we construct a binary tree of auxiliary variables where the leaves are all variables of the form $X[r, gr, a]$ and the root is the variable $X[r, a]$. Figure 9 depicts a portion of the BN for the program in Example VI.1.

Complexity. We now introduce the main result of this section.

Theorem VI.1. *The data complexity of inference for acyclic PROBLOG programs is PTIME.*

This follows from (1) the relaxed grounding and the encoding can be computed in PTIME in terms of data complexity, (2) the encoding ensures that, for acyclic programs, the resulting Bayesian Network is a forest of poly-trees, and (3) inference algorithms for poly-tree BNs [47] run in polynomial time in the BN's size. In Appendix B, we extend our encoding to handle additional features such as annotated disjunctions, and we prove its correctness and complexity.

VII. ANGERONA

ANGERONA is a DBIC mechanism that provably secures databases against probabilistic inferences. ANGERONA is

Algorithm 1: ANGERONA Enforcement Algorithm.

Input: A system state $s = \langle db, U, P \rangle$, a history h , an action $\langle u, q \rangle$, a system configuration C , and a C -ATKLOG model ATK .

Output: The security decision in $\{\top, \perp\}$.

```
begin
  for  $\langle u, \psi, l \rangle \in secrets(P, u)$  do
    if  $secure(C, ATK, h, \langle u, \psi, l \rangle)$ 
      if  $pox(C, ATK, h, \langle u, q \rangle)$ 
         $h' := h \cdot \langle \langle u, q \rangle, \top, \top \rangle$ 
        if  $\neg secure(C, ATK, h', \langle u, \psi, l \rangle)$ 
          return  $\perp$ 
        if  $pox(C, ATK, h, \langle u, \neg q \rangle)$ 
           $h' := h \cdot \langle \langle u, q \rangle, \top, \perp \rangle$ 
          if  $\neg secure(C, ATK, h', \langle u, \psi, l \rangle)$ 
            return  $\perp$ 
      return  $\top$ 

function  $secure(\langle D, \Gamma \rangle, ATK, h, \langle u, \psi, l \rangle)$ 
   $p := ATK(u)$ 
  for  $\phi \in knowledge(h, u)$  do
     $p := p \cup PL(\phi) \cup \{evidence(head(\phi), true)\}$ 
   $p := p \cup PL(\psi)$ 
  return  $\llbracket p \rrbracket_D(head(\psi)) < l$ 

function  $pox(\langle D, \Gamma \rangle, ATK, h, \langle u, \psi \rangle)$ 
   $p := ATK(u)$ 
  for  $\phi \in knowledge(h, u)$  do
     $p := p \cup PL(\phi) \cup \{evidence(head(\phi), true)\}$ 
   $p := p \cup PL(\psi)$ 
  return  $\llbracket p \rrbracket_D(head(\psi)) > 0$ 
```

parametrized by an ATKLOG model representing the attacker's capabilities and it leverages PROBLOG's inference capabilities.

A. Checking Query Security

Algorithm 1 presents ANGERONA. It takes as input a system state $s = \langle db, U, P \rangle$, a history h , the current query q issued by the user u , a system configuration C , and an ATKLOG model ATK formalizing the users' beliefs. ANGERONA checks whether disclosing the result of the current query q may violate any secrets in $secrets(P, u)$. If this is the case, the algorithm concludes that q 's execution would be insecure and returns \perp . Otherwise, it returns \top and authorizes q 's execution. Note that once we fix a configuration C and an ATKLOG model ATK , ANGERONA is a C -PDP as defined in §IV-C.

To check whether a query q may violate a secret $\langle u, \psi, l \rangle \in secrets(P, u)$, ANGERONA first checks whether the secret has been already violated. If this is not the case, ANGERONA checks whether disclosing q violates any secret. This requires checking that u 's belief about the secret ψ stays below the threshold independently of the result of the query q ; hence, we must ensure that u 's belief is below the threshold both in case the query q holds in the actual database and in case q does not hold (this ensures that the access control decision itself does not leak information). ANGERONA, therefore, first checks whether there exists at least one possible database state where q is satisfied given h , using the procedure pox . If this is the case, the algorithm extends the current history h with the new event recording that the query q is authorized and its

result is \top and it checks whether u 's belief about ψ is still below the corresponding threshold once q 's result is disclosed, using the $secure$ procedure. Afterwards, ANGERONA checks whether there exists at least a possible database state where q is not satisfied given h , it extends the current history h with another event representing that the query q does not hold, and it checks again whether disclosing that q does not hold in the current database state violates the secret. Note that checking whether there is a database state where q is (or is not) satisfied is essential to ensure that the conditioning that happens in the $secure$ procedure is well-defined, i.e., the set of states we condition on has non-zero probability.

ANGERONA uses the $secure$ subroutine to determine whether a secret's confidentiality is violated. This subroutine takes as input a system configuration, an ATKLOG model ATK , a history h , and a secret $\langle u, \psi, l \rangle$. It first computes the set $knowledge(h, u)$ containing all the authorized queries in the u -projection of h , i.e., $knowledge(h, u) = \{\phi \mid \exists i. h|_u(i) = \langle \langle u, \phi \rangle, \top, \top \rangle \cup \{\neg\phi \mid \exists i. h|_u(i) = \langle \langle u, \phi \rangle, \top, \perp \rangle\}$. Afterwards, it generates a PROBLOG program p by extending $ATK(u)$ with additional rules. In more detail, it translates each relational calculus sentence $\phi \in knowledge(h, u)$ to an equivalent set of PROBLOG rules $PL(\phi)$. The translation $PL(\phi)$ is standard [2]. For example, given a query $\phi = (A(1) \wedge B(2)) \vee \neg C(3)$, the translation $PL(\phi)$ consists of the rules $\{(h_1 \leftarrow A(1)), (h_2 \leftarrow B(2)), (h_3 \leftarrow \neg C(3)), (h_4 \leftarrow h_1, h_2), (h_5 \leftarrow h_4), (h_5 \leftarrow h_3)\}$, where h_1, \dots, h_5 are fresh predicate symbols. We denote by $head(\phi)$ the unique predicate symbol associated with the sentence ϕ by the translation $PL(\phi)$. In our example, $head(\phi)$ is the fresh predicate symbol h_5 . The algorithm then conditions the initial probability distribution $ATK(u)$ based on the sentences in $knowledge(h, u)$. This is done using *evidence statements*, which are special PROBLOG statements of the form $evidence(a, v)$, where a is a ground atom and v is either `true` or `false`; see Appendix A. For each sentence $\phi \in knowledge(h, u)$, the program p contains a statement $evidence(head(\phi), true)$. Finally, the algorithm translates ψ to a set of logic programming rules and checks whether ψ 's probability is below the threshold l .

The pox subroutine takes as input a system configuration, an ATKLOG model ATK , a history h , and a query $\langle u, \psi \rangle$. It determines whether there is a database db' that satisfies ψ and complies with the history $h|_u$. Internally, the routine again constructs a PROBLOG program starting from ATK , $knowledge(h, u)$, and ψ . Afterwards, it uses the program to check whether the probability of ψ given $h|_u$ is greater than 0.

Given a run $\langle s, h \rangle$ and a user u , the $secure$ and pox subroutines condition u 's initial beliefs based on the sentences in $knowledge(h, u)$, instead of using the set $\llbracket r \rrbracket_{\sim_u}$ as in the ATKLOG semantics. The key insight is that, as we prove in Appendix C, the set of possible database states defined by the sentences in $knowledge(h, u)$ is equivalent to $\llbracket r \rrbracket_{\sim_u}$, which contains all database states derivable from the runs $r' \sim_u r$. This allows us to use PROBLOG to implement ATKLOG's semantics without explicitly computing $\llbracket r \rrbracket_{\sim_u}$.

Example VII.1. Let ATK be the attacker model in Example IV.4, u be the user *Mallory*, the database state be $s_{\{A,B,C\}}$, where Alice, Bob, and Carl have cancer, and the policy P be the one from Example IV.2. Furthermore, let q_1, \dots, q_4 be the queries issued by *Mallory* in Example IV.3. ANGERONA permits the execution of the first two queries since they do not violate the policy. In contrast, it denies the execution of the last two queries as they leak sensitive information. ■

Confidentiality. As we prove in Appendix C, ANGERONA provides the desired security guarantees for any ATKLOG-attacker. Namely, it authorizes only those queries whose disclosure does not increase an attacker’s beliefs in the secrets above the corresponding thresholds. ANGERONA also provides precise completeness guarantees: it authorizes all secrecy-preserving queries. Informally, a query $\langle u, q \rangle$ is *secrecy-preserving* given a run r and a secret $\langle u, \psi, l \rangle$ iff disclosing the result of $\langle u, q \rangle$ in any run $r' \sim_u r$ does not violate the secret.

Theorem VII.1. *Let a system configuration C and a C -ATKLOG model ATK be given, and let ANGERONA be the C -PDP f . ANGERONA provides data confidentiality with respect to C and $\lambda u \in \mathcal{U}$. $\llbracket ATK(u) \rrbracket_D$, and it authorizes all secrecy-preserving queries.*

Complexity. ANGERONA’s complexity is dominated by the complexity of inference. We focus our analysis only on data complexity, i.e., the complexity when only the ground atoms in the PROBLOG programs are part of the input while everything else is fixed. A *literal query* is a query consisting either of a ground atom $a(\bar{c})$ or its negation $\neg a(\bar{c})$. We call an ATKLOG model *acyclic* if all belief programs in it are acyclic. Furthermore, a *literal secret* is a secret $\langle U, \phi, l \rangle$ such that ϕ is a literal query. We prove in Appendix C that for acyclic ATKLOG models, literal queries, and literal secrets, the PROBLOG programs produced by the *secure* and *pox* subroutines are acyclic. We can therefore use our dedicated inference engine from §VI to reason about them. Hence, ANGERONA can be used to protect databases in PTIME in terms of data complexity. Literal queries are expressive enough to formulate queries about the database content such as “does Alice have cancer?”.

Theorem VII.2. *For all acyclic ATKLOG attackers, for all literal queries q , for all runs r whose histories contain only literal queries and contain only secrets expressed using literal queries, ANGERONA’s data complexity is PTIME.*

Discussion. Our tractability guarantees apply only to acyclic ATKLOG models, literal queries, and literal secrets. Nevertheless, ANGERONA can still handle relevant problems of interest. As stated in §VI, acyclic models are as expressive as poly-tree Bayesian Networks, one of the few classes of Bayesian Networks with tractable inference. Hence, for many probabilistic models that cannot be represented as acyclic ATKLOG models, exact probabilistic inference is intractable.

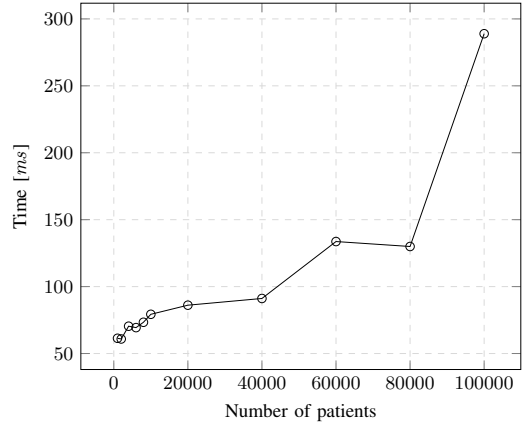


Fig. 10: ANGERONA execution time in seconds.

Literal queries are expressive enough to state simple facts about the database content. More complex (non-literal) queries can be simulated using (possibly large) sequences of literal queries. Similarly, policies with non-literal secrets can be implemented as sets of literal secrets, and the Boole–Fréchet inequalities [38] can be used to derive the desired thresholds. In both cases, however, our completeness guarantees hold only for the resulting literal queries, not for the original ones.

Finally, whenever our tractability constraints are violated, ANGERONA can still be used by directly using PROBLOG’s inference capabilities. In this case, one would retain the security and completeness guarantees (Theorem VII.1) but lose the tractability guarantees (Theorem VII.2).

B. Implementation and Empirical Evaluation

To evaluate the feasibility of our approach in practice, we implemented a prototype of ANGERONA, available at [36]. The prototype implements our dedicated inference algorithm for acyclic PROBLOG programs (§VI), which computes the relaxed grounding of the input program p , constructs the Bayesian Networks BN , and performs the inference over BN using belief propagation [47]. For inference over BN , we rely on the GRRM library [66]. Observe that evidence statements in PROBLOG are encoded by fixing the values of the corresponding random variables in the BN. Note also that computing the relaxed grounding of p takes polynomial time in terms of data complexity, where the exponent is determined by p ’s rules. A key optimization is to pre-compute the relaxed grounding and construct BN off-line. This avoids grounding p and constructing the (same) Bayesian Network for each query. In our experiments we measure this time separately.

We use our prototype to study ANGERONA’s efficiency and scalability. We run our experiments on a PC with an Intel i7 processor and 32GB of RAM. For our experiments, we consider the database schema from §IV. For the belief programs, we use the PROBLOG program given in §V, which can be encoded as an acyclic program when the parent-child relation is a poly-tree. We evaluate ANGERONA’s efficiency and scalability in terms of the number of ground atoms in the belief programs. We generate synthetic belief programs

containing 1,000 to 100,000 patients and for each of these instances, we generate 100 random queries of the form $R(\bar{t})$, where R is a predicate symbol and \bar{t} is a tuple. For each instance and sequence of queries, we check the security of each query with our prototype, against a policy containing 100 randomly generated secrets specified as literal queries.

Figure 10 reports the execution times for our case study. Once the BN is generated, ANGERONA takes under 300 milliseconds, even for our larger examples, to check a query’s security. During the initialization phase of our dedicated inference engine, we ground the original PROLOG program and translate it into a BN. Most of the time is spent in the grounding process, whose data complexity is polynomial, where the polynomial’s degree is determined by the number of free variables in the belief program. Our prototype uses a naive bottom-up grounding technique, and for our larger examples the initialization times is less than 2.5 minutes. We remark, however, that the initialization is performed just once per belief program. Furthermore, it can be done offline and its performance can be greatly improved by naive parallelization.

VIII. RELATED WORK

Database Inference Control. Existing DBIC approaches protect databases (either at design time [54], [55] or at runtime [14], [15], [39]) only against restricted classes of probabilistic dependencies, e.g., those arising from functional and multi-valued dependencies. ATKLOG, instead, supports arbitrary probabilistic dependencies, and even our acyclic fragment can express probabilistic dependencies that are not supported by [14], [15], [39], [54], [55]. Weise [71] proposes a DBIC framework, based on possibilistic logic, that formalizes secrets as sentences and expresses policies by associating bounds to secrets. Possibility theory differs from probability theory, which results in subtle differences. For instance, there is no widely accepted definition of conditioning for possibility distributions, cf. [11]. Thus, the probabilistic model from §II cannot be encoded in Weise’s framework [71].

Statistical databases store information associated to different individuals and support queries that return statistical information about an entire population [16]. DBIC solutions for statistical databases [3], [16], [22]–[24] prevent leakages of information about single individuals while allowing the execution of statistical queries. These approaches rely on various techniques, such as perturbing the original data, synthetically generating data, or restricting the data on which the queries are executed. Instead, we protect specific secrets in the presence of probabilistic data dependencies and we return the original query result, without modifications, if it is secure.

Differential Privacy. Differential Privacy [25], [26] is widely used for privacy-preserving data analysis. Systems such as ProPer [27] or PINQ [53] provide users with automated ways to perform differentially private computations. A differentially private computation guarantees that the presence (or absence) of an individual’s data in the input data set affects the probability of the computation’s result only in limited way, i.e., by at

most a factor e^ϵ where ϵ is a parameter controlling the privacy-utility trade-off. While differential privacy does not make any assumption about the attacker’s beliefs, we assume that the attacker’s belief is known and we guarantee that for all secrets in the policy, no user can increase his beliefs, as specified in the attacker model, over the corresponding thresholds by interacting with the system.

Information Flow Control. Quantified Information Flow [6], [17], [48], [50] aims at quantifying the amount of information leaked by a program. Instead of measuring the amount of leaked information, we focus on restricting the information that an attacker may obtain about a set of given secrets.

Non-interference has been extended to consider probabilities [5], [58], [70] for concurrent programs. Our security notion, instead, allows those leaks that do not increase an attacker’s beliefs in a secret above the threshold, and it can be seen as a probabilistic extension of *opacity* [60], which allows any leak except leaking whether the secret holds.

Mardziel et al. [51] present a general DBIC architecture, where users’ beliefs are expressed as probabilistic programs, security requirements as threshold on these beliefs, and the beliefs are updated in response to the system’s behaviour. Our work directly builds on top of this architecture. However, instead of using an imperative probabilistic language, we formalize beliefs using probabilistic logic programming, which provides a natural and expressive language for formalizing dependencies arising in the database setting, e.g., functional and multi-valued dependencies, as well as common probabilistic models, like Bayesian Networks.

Mardziel et al. [51] also propose a DBIC mechanism based on abstract interpretation. They do not provide any precise complexity bound for their mechanism. Their algorithm’s complexity class, however, appears to be intractable, since they use a probabilistic extension of the polyhedra abstract domain, whose asymptotic complexity is exponential in the number of program variables [62]. In contrast, ANGERONA exploits our inference engine for acyclic programs to secure databases against a practically relevant class of probabilistic inferences, and it provides precise tractability and completeness guarantees.

We now compare (unrestricted) ATKLOG with the imperative probabilistic language used in [51]. ATKLOG allows one to concisely encode probabilistic properties specifying relations between tuples in the database. For instance, a property like “the probability of $A(x)$ is $1/2^n$, where n is the number of tuples (x, y) in B ” can be encoded as $1/2::A(x) \leftarrow B(x, y)$. Encoding this property as an imperative program is more complex; it requires a **for** statement to iterate over all variables representing tuples in B and an **if** statement to filter the tuples. In contrast to [51], ATKLOG provides limited support for numerical constraints (as we support only finite domains). Mardziel et al. [51] formalize queries as imperative probabilistic programs. They can, therefore, also model probabilistic queries or the use of randomization to limit disclosure. While all these features are supported by ATKLOG, our goal is to protect databases from attackers that use standard query lan-

guages like SQL. Hence, we formalize queries using relational calculus and ignore probabilistic queries. Similarly to [51], our approach can be extended to handle some uncertainty on the attackers' capabilities. In particular, we can associate to each user a finite number of possible beliefs, instead of a single one. However, like [51], we cannot handle infinitely many alternative beliefs.

Probabilistic Programming. Probabilistic programming is an active area of research [35]. Here, we position PROBLOG with respect to expressiveness and inference. Similarly to [33], [51], PROBLOG can express only discrete probability distributions, and it is less expressive than languages supporting continuous distributions [32], [34], [61]. Current exact inference algorithms for probabilistic programs are based on program analysis techniques, such as symbolic execution [32], [61] or abstract interpretation [51]. In this respect, we present syntactic criteria that *ensure* tractable inference for PROBLOG. Sampson et al. [59] symbolically execute probabilistic programs and translate them to BNs to verify probabilistic assertions. In contrast, we translate PROBLOG programs to BNs to perform exact inference and our translation is tailored to work together with our acyclicity constraints to allow tractable inference.

IX. CONCLUSION

Effectively securing databases that store data with probabilistic dependencies requires an expressive language to capture the dependencies and a tractable enforcement mechanism. To address these requirements, we developed ATKLOG, a formal language providing an expressive and concise way to represent attackers' beliefs while interacting with the system. We leveraged this to design ANGERONA, a provably secure DBIC mechanism that prevents the leakage of sensitive information in the presence of probabilistic dependencies. ANGERONA is based on a dedicated inference engine for a fragment of PROBLOG where exact inference is tractable.

We see these results as providing a foundation for building practical protection mechanisms, which include probabilistic dependencies, as part of real-world database systems. As future work, we plan to extend our framework to dynamic settings where the database and the policy change. We also intend to explore different fragments of PROBLOG and relational calculus for which exact inference is practical

Acknowledgments. We thank Ognjen Maric, Dmitriy Traytel, Der-Yuean Yu, and the anonymous reviewers for their comments.

REFERENCES

- [1] "ProbLog – Probabilistic Programming," Online at <http://dtai.cs.kuleuven.be/problog/index.html>.
- [2] S. Abiteboul, R. Hull, and V. Vianu, *Foundations of databases*. Addison-Wesley Reading, 1995, vol. 8.
- [3] N. R. Adam and J. C. Worthmann, "Security-control methods for statistical databases: a comparative study," *ACM Computing Surveys (CSUR)*, vol. 21, no. 4, pp. 515–556, 1989.
- [4] K. Ahmed, A. A. Emran, T. Jesmin, R. F. Mukti, M. Z. Rahman, and F. Ahmed, "Early detection of lung cancer risk using data mining," *Asian Pacific Journal of Cancer Prevention*, vol. 14, no. 1, pp. 595–598, 2013.
- [5] A. Aldini, "Probabilistic information flow in a process algebra," in *Proceedings of the 12th International Conference on Concurrency Theory*. Springer, 2001, pp. 152–168.
- [6] M. Alvim, M. Andrés, and C. Palamidessi, "Probabilistic information flow," in *Proceedings of the 25th Annual IEEE Symposium on Logic in Computer Science*. IEEE, 2010, pp. 314–321.
- [7] A. Askarov and S. Chong, "Learning is change in knowledge: Knowledge-based security for dynamic policies," in *Proceedings of the 25th IEEE Computer Security Foundations Symposium*. IEEE, 2012, pp. 308–322.
- [8] A. Askarov and A. Sabelfeld, "Gradual release: Unifying declassification, encryption and key release policies," in *Proceedings of the 28th IEEE Symposium on Security and Privacy*. IEEE, 2007, pp. 207–221.
- [9] V. Bárány, B. Ten Cate, and M. Otto, "Queries with guarded negation," *Proceedings of the VLDB Endowment*, vol. 5, no. 11, pp. 1328–1339, 2012.
- [10] P. A. Bonatti, S. Kraus, and V. Subrahmanian, "Foundations of secure deductive databases," *IEEE Transactions on Knowledge and Data Engineering*, vol. 7, no. 3, pp. 406–422, 1995.
- [11] B. Bouchon-Meunier, G. Coletti, and C. Marsala, "Independence and possibilistic conditioning," *Annals of Mathematics and Artificial Intelligence*, vol. 35, no. 1, pp. 107–123, 2002.
- [12] A. Brodsky, C. Farkas, and S. Jajodia, "Secure databases: Constraints, inference channels, and monitoring disclosures," *IEEE Transactions on Knowledge and Data Engineering*, vol. 12, no. 6, pp. 900–919, 2000.
- [13] Centers for Medicare & Medicaid Services, "The Health Insurance Portability and Accountability Act of 1996 (HIPAA)," Online at <http://www.cms.hhs.gov/hipaahp>, 1996.
- [14] Y. Chen and W. W. Chu, "Database security protection via inference detection," in *International Conference on Intelligence and Security Informatics*. Springer, 2006, pp. 452–458.
- [15] —, "Protection of database security via collaborative inference detection," *IEEE Transactions on Knowledge and Data Engineering*, vol. 20, no. 8, pp. 1013–1027, 2008.
- [16] F. Y. Chin and G. Ozsoyoglu, "Auditing and inference control in statistical databases," *IEEE Transactions on Software Engineering*, vol. 8, no. 6, pp. 574–582, 1982.
- [17] D. Clark, S. Hunt, and P. Malacaria, "A static analysis for quantifying information flow in a simple imperative language," *Journal of Computer Security*, vol. 15, no. 3, pp. 321–371, 2007.
- [18] M. R. Clarkson, A. C. Myers, and F. B. Schneider, "Belief in information flow," in *Proceedings of the 18th IEEE Workshop on Computer Security Foundations*. IEEE, 2005, pp. 31–45.
- [19] —, "Quantifying information flow with beliefs," *Journal of Computer Security*, vol. 17, no. 5, pp. 655–701, 2009.
- [20] L. De Raedt and A. Kimmig, "Probabilistic (logic) programming concepts," *Machine Learning*, vol. 100, no. 1, pp. 5–47, 2015.
- [21] L. De Raedt, A. Kimmig, and H. Toivonen, "Problog: A probabilistic prolog and its application in link discovery," in *Proceedings of the 20th International Joint Conference on Artificial Intelligence*. Morgan Kaufmann, 2007, pp. 2468–2473.
- [22] D. E. Denning, "Secure statistical databases with random sample queries," *ACM Transactions on Database Systems (TODS)*, vol. 5, no. 3, pp. 291–315, 1980.
- [23] D. Dobkin, A. K. Jones, and R. J. Lipton, "Secure databases: protection against user influence," *ACM Transactions on Database Systems*, vol. 4, no. 1, pp. 97–106, 1979.
- [24] J. Domingo-Ferrer, *Inference control in statistical databases: From theory to practice*. Springer, 2002, vol. 2316.
- [25] C. Dwork, "Differential privacy," in *Proceedings of the 33rd International Colloquium on Automata, Languages and Programming*. Springer, 2006, pp. 1–12.
- [26] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3-4, pp. 211–407, 2014.
- [27] H. Ebadi, D. Sands, and G. Schneider, "Differential privacy: Now it's getting personal," in *Proceedings of the 42nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 2015, pp. 69–81.
- [28] European Parliament, "General Data Protection Regulation (2016/679)," Online at <http://eur-lex.europa.eu/eli/reg/2016/679/oj>, 2016.
- [29] A. Evfimievski, R. Fagin, and D. Woodruff, "Epistemic privacy," *Journal of ACM*, vol. 58, no. 1, pp. 2:1–2:45, 2010.
- [30] C. Farkas and S. Jajodia, "The inference problem: a survey," *ACM SIGKDD Explorations Newsletter*, vol. 4, no. 2, pp. 6–11, 2002.
- [31] D. Fierens, G. Van den Broeck, J. Renkens, D. Shterionov, B. Gutmann, I. Thon, G. Janssens, and L. De Raedt, "Inference and learning in

- probabilistic logic programs using weighted boolean formulas,” *Theory and Practice of Logic Programming*, vol. 15, no. 03, pp. 358–401, 2015.
- [32] T. Gehr, S. Misailovic, and M. Vechev, “Psi: Exact symbolic inference for probabilistic programs,” in *Proceedings of the 28th International Conference on Computer Aided Verification*. Springer, 2016, pp. 62–83.
- [33] L. Getoor, *Introduction to statistical relational learning*. MIT press, 2007.
- [34] A. D. Gordon, T. Graepel, N. Rolland, C. Russo, J. Borgstrom, and J. Guiver, “Tabular: A schema-driven probabilistic programming language,” in *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 2014, pp. 321–334.
- [35] A. D. Gordon, T. A. Henzinger, A. V. Nori, and S. K. Rajamani, “Probabilistic programming,” in *Proceedings of the Conference on The Future of Software Engineering*. ACM, 2014, pp. 167–181.
- [36] M. Guarnieri, S. Marinovic, and D. Basin, “Securing databases from probabilistic inference — prototype,” Online at <http://www.infsec.ethz.ch/research/projects/FDAC.html>.
- [37] —, “Strong and provably secure database access control,” in *Proceedings of the 1st European Symposium on Security and Privacy*. ACM, 2005, pp. 118–127.
- [38] T. Hailperin, “Probability logic,” *Notre Dame Journal of Formal Logic*, vol. 25, no. 3, pp. 198–212, 1984.
- [39] J. Hale and S. Shenoi, “Catalytic inference analysis: Detecting inference threats due to knowledge discovery,” in *Proceedings of the 18th IEEE Symposium on Security and Privacy*. IEEE, 1997, pp. 188–199.
- [40] J. He, W. W. Chu, and Z. V. Liu, “Inferring privacy information from social networks,” in *International Conference on Intelligence and Security Informatics*. Springer, 2006, pp. 154–165.
- [41] D. Hedin and A. Sabelfeld, “A perspective on information-flow control,” in *Proceedings of the 2011 Marktobendorf Summer School*.
- [42] T. H. Hinke, H. S. Delugach, and R. P. Wolf, “Protecting databases from inference attacks,” *Computers & Security*, vol. 16, no. 8, pp. 687–708, 1997.
- [43] M. Humbert, E. Ayday, J.-P. Hubaux, and A. Telenti, “Addressing the concerns of the lacks family: Quantification of kin genomic privacy,” in *Proceedings of the 20th ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2013, pp. 1141–1152.
- [44] G. Kabra, R. Ramamurthy, and S. Sudarshan, “Redundancy and information leakage in fine-grained access control,” in *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data*. ACM, 2006.
- [45] V. Katos, D. Vrakas, and P. Katsaros, “A framework for access control with inference constraints,” in *Proceedings of 35th IEEE Annual Conference on Computer Software and Applications*. IEEE, 2011, pp. 289–297.
- [46] K. Kenthapadi, N. Mishra, and K. Nissim, “Simulatable auditing,” in *Proceedings of the 24th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. ACM, 2005, pp. 118–127.
- [47] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [48] B. Köpf and D. Basin, “An information-theoretic model for adaptive side-channel attacks,” in *Proceedings of the 14th ACM Conference on Computer and Communications Security*. ACM, 2007, pp. 286–296.
- [49] S. L. Lauritzen and N. A. Sheehan, “Graphical models for genetic analyses,” *Statistical Science*, vol. 18, no. 4, pp. 489–514, 2003.
- [50] G. Lowe, “Quantifying information flow,” in *Proceedings of the 15th IEEE Workshop on Computer Security Foundations*. IEEE, 2002, pp. 18–31.
- [51] P. Mardziel, S. Magill, M. Hicks, and M. Srivatsa, “Dynamic enforcement of knowledge-based security policies using probabilistic abstract interpretation,” *Journal of Computer Security*, vol. 21, no. 4, pp. 463–532, 2013.
- [52] W. Mathew, R. Raposo, and B. Martins, “Predicting future locations with hidden markov models,” in *Proceedings of the 2012 ACM Conference on Ubiquitous Computing*. ACM, 2012, pp. 911–918.
- [53] F. D. McSherry, “Privacy integrated queries: an extensible platform for privacy-preserving data analysis,” in *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*. ACM, 2009, pp. 19–30.
- [54] M. Morgenstern, “Security and inference in multilevel database and knowledge-base systems,” in *Proceedings of the 1987 ACM SIGMOD International Conference on Management of data*. ACM, 1987, pp. 357–373.
- [55] —, “Controlling logical inference in multilevel database systems,” in *Proceedings of the 9th IEEE Symposium on Security and Privacy*. IEEE, 1988, pp. 245–255.
- [56] A. Pagourtzis and S. Zachos, “The complexity of counting functions with easy decision version,” in *International Symposium on Mathematical Foundations of Computer Science*. Springer, 2006, pp. 741–752.
- [57] X. Qian, M. E. Stickel, P. D. Karp, T. F. Lunt, and T. D. Garvey, “Detection and elimination of inference channels in multilevel relational database systems,” in *Proceedings of the 14th IEEE Symposium on Security and Privacy*. IEEE, 1993, pp. 196–205.
- [58] A. Sabelfeld and D. Sands, “Probabilistic noninterference for multi-threaded programs,” in *Proceedings of the 13th IEEE Workshop on Computer Security Foundations*. IEEE, 2000, pp. 200–214.
- [59] A. Sampson, P. Panckheka, T. Mytkowicz, K. S. McKinley, D. Grossman, and L. Ceze, “Expressing and verifying probabilistic assertions,” in *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, 2014, pp. 112–122.
- [60] D. Schoepe and A. Sabelfeld, “Understanding and enforcing opacity,” in *Proceedings of the 28th IEEE Computer Security Foundations Symposium*. IEEE, 2015, pp. 539–553.
- [61] C.-c. Shan and N. Ramsey, “Exact bayesian inference by symbolic disintegration,” in *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages*. ACM, 2017, pp. 130–144.
- [62] G. Singh, M. Püschel, and M. Vechev, “Fast polyhedra abstract domain,” in *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages*. ACM, 2017, pp. 46–59.
- [63] T.-A. Su and G. Ozsoyoglu, “Controlling FD and MVD inferences in multilevel relational database systems,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 3, no. 4, pp. 474–485, 1991.
- [64] —, “Data dependencies and inference control in multilevel relational database systems,” in *Proceedings of the 7th IEEE Symposium on Security and Privacy*. IEEE, 1987, p. 202.
- [65] D. Suciu, D. Olteanu, C. Ré, and C. Koch, “Probabilistic databases,” *Synthesis Lectures on Data Management*, vol. 3, no. 2, pp. 1–180, 2011.
- [66] C. Sutton, “GRMM: Graphical Models in Mallet,” Online at <http://mallet.cs.umass.edu/grmm/>.
- [67] M. Thuraisingham, “Security checking in relational database management systems augmented with inference engines,” *Computers & Security*, vol. 6, no. 6, pp. 479–492, 1987.
- [68] T. S. Toland, C. Farkas, and C. M. Eastman, “The inference problem: Maintaining maximal availability in the presence of database updates,” *Computers & Security*, vol. 29, no. 1, pp. 88–103, 2010.
- [69] P. J. Villeneuve and Y. Mao, “Lifetime probability of developing lung cancer, by smoking status, Canada,” *Cancer Journal of Public Health*, vol. 85, no. 6, pp. 385–388, 1994.
- [70] D. Volpano and G. Smith, “Probabilistic noninterference in a concurrent language,” *Journal of Computer Security*, vol. 7, no. 2-3, pp. 231–253, 1999.
- [71] L. Wiese, “Keeping secrets in possibilistic knowledge bases with necessity-valued privacy policies,” in *Computational Intelligence for Knowledge-Based Systems Design*. Springer, 2010, pp. 655–664.
- [72] R. W. Yip and E. Levitt, “Data level inference detection in database systems,” in *Proceedings of the 11th IEEE Computer Security Foundations Workshop*. IEEE, 1998, pp. 179–189.

APPENDICES

Here, we provide additional details about numerous aspects of the paper, such as acyclic PROBLOG programs and ANGERONA. We also provide complete proofs of all our results.

Structure. Appendix A reviews PROBLOG’s syntax and semantics. In Appendix B, we formalize acyclicity for PROBLOG programs, we present our encoding from acyclic programs to poly-tree Bayesian Networks, we prove the correctness of this encoding, and we give complexity results for inference for acyclic PROBLOG programs. Finally, Appendix C contains additional details about ANGERONA as well as all the security and complexity proofs.

Here, we provide a formal account of PROBLOG, which follows [20], [21], [31]. In addition to PROBLOG's semantics, we present here again some (revised) aspects of PROBLOG's syntax, which we introduced in §V. As mentioned in §V, we restrict ourselves to the function-free fragment of PROBLOG.

Syntax. As introduced in §V, a (Σ, \mathbf{dom}) -probabilistic atom is an atom $a \in \mathcal{A}_{\Sigma, \mathbf{dom}}$ annotated with a value $0 \leq v \leq 1$, denoted $v::a$. If $v = 1$, then we write $a(\bar{c})$ instead of $1::a(\bar{c})$. A (Σ, \mathbf{dom}) -PROBLOG program is a finite set of ground probabilistic (Σ, \mathbf{dom}) -atoms and (Σ, \mathbf{dom}) -rules. Note that ground atoms $a \in \mathcal{A}_{\Sigma, \mathbf{dom}}$ are represented as $1::a$. Observe also that rules do not involve probabilistic atoms (as formalized in §V). This is without loss of generality: as we show below probabilistic rules and annotated disjunctions can be represented using only probabilistic atoms and non-probabilistic rules. We denote by $prob(p)$ the set of all probabilistic ground atoms $v::a$ in p , i.e., $prob(p) := \{v::a \in p \mid 0 \leq v \leq 1 \wedge a \in \mathcal{A}_{\Sigma, \mathbf{dom}}\}$, and by $rules(p)$ the non-probabilistic rules in p , i.e., $rules(p) := p \setminus prob(p)$. As already stated in §V, we consider only programs p that do not contain negative cycles in the rules. Finally, we say that a PROBLOG program p is a *logic program* iff $v = 1$ for all $v::a \in prob(p)$, i.e., p does not contain probabilistic atoms.

Semantics. Given a (Σ, \mathbf{dom}) -PROBLOG program p , a p -grounded instance is a PROBLOG program $A \cup R$, where the set of ground atoms A is a subset of $\{a \mid v::a \in prob(p)\}$ and $R = rules(p)$. Informally, a grounded instance of p is one of the logic programs that can be obtained by selecting some of the probabilistic atoms in p and keeping all rules in p . A p -probabilistic assignment is a total function associating to each probabilistic atom $v::a$ in $prob(p)$ a value in $\{\top, \perp\}$. We denote by $A(p)$ the set of all p -probabilistic assignments. The *probability* of a p -probabilistic assignment f is $prob(f) = \prod_{v::a \in prob(p)} (\prod_{f(v::a)=\top} v \cdot \prod_{f(v::a)=\perp} (1-v))$. Given a p -probabilistic assignment f , $instance(p, f)$ denotes the p -grounded instance $\{a \mid \exists v. f(v::a) = \top\} \cup rules(p)$. Finally, given a p -grounded instance p' , $WFM(p')$ denotes the well-founded model associated with the logic program p' , as defined by standard logic programming semantics [2].

The semantics of a (Σ, \mathbf{dom}) -PROBLOG program p is defined as a probability distribution over all possible p -grounded instances. Note that PROBLOG's semantics relies on the *closed world assumption*, namely every fact that is not in a given model is considered false. The semantics of p , denoted by $\llbracket p \rrbracket$, is as follows: $\llbracket p \rrbracket(p') = \sum_{f \in F_{p, p'}} prob(f)$, where $F_{p, p'} = \{f \in A(p) \mid p' = instance(p, f)\}$. We remark that a (Σ, \mathbf{dom}) -PROBLOG program p implicitly defines a probability distribution over (Σ, \mathbf{dom}) -structures. Indeed, the probability of a given (Σ, \mathbf{dom}) -structure s is the sum of the probabilities of all p -grounded instances whose well-founded model is s . With a slight abuse of notation, we extend the semantics of p to (Σ, \mathbf{dom}) -structures and ground atoms as follows: $\llbracket p \rrbracket(s) = \sum_{f \in \mathcal{M}(p, s)} prob(f)$, where s is a $(\Sigma,$

$\mathbf{dom})$ -structure and $\mathcal{M}(p, s)$ is the set of all assignments f such that $WFM(instance(p, f)) = s$. Finally, p 's semantics can be lifted to sentences as follows: $\llbracket p \rrbracket(\phi) = \sum_{s \in \llbracket \phi \rrbracket} \llbracket p \rrbracket(s)$, where $\llbracket \phi \rrbracket = \{s \in \Omega_D^\Gamma \mid \llbracket \phi \rrbracket^s = \top\}$.

Evidence. PROBLOG supports expressing *evidence* inside programs [20]. To express evidence, i.e., to condition a distribution on some event, we use statements of the form $evidence(a, v)$, where a is a ground atom and $v \in \{\mathbf{true}, \mathbf{false}\}$. Let p be a (Σ, \mathbf{dom}) -PROBLOG program p with evidence $evidence(a_1, v_1), \dots, evidence(a_n, v_n)$, and p' be the program without the evidence statements. Furthermore, let $POX(p')$ be the set of all (Σ, \mathbf{dom}) -structures s complying with the evidence, i.e., the set of all states s such that a_i holds in s iff $v_i = \mathbf{true}$. Then, $\llbracket p \rrbracket(s)$, for a (Σ, \mathbf{dom}) -structure $s \in POX(p)$, is $\llbracket p' \rrbracket(s) \cdot \left(\sum_{s' \in POX(p)} \llbracket p' \rrbracket(s') \right)^{-1}$.

Syntactic Sugar. Following [20], [21], [31], we extend PROBLOG programs with two additional constructs: annotated disjunctions and probabilistic rules. As shown in [20], these constructs are just syntactic sugar. A *probabilistic rule* is a PROBLOG rule where the head is a probabilistic atom. The probabilistic rule $v::h \leftarrow l_1, \dots, l_n$ can be encoded using the additional probabilistic atoms $v::sw(_)$ and the rule $h \leftarrow l_1, \dots, l_n, sw(\bar{x})$, where sw is a fresh predicate symbol, \bar{x} is the tuple containing the variables in $vars(h) \cup \bigcup_{1 \leq i \leq n} vars(l_i)$, and $v::sw(_)$ is a shorthand representing the fact that there is a probabilistic atom $v::sw(\bar{t})$ for each tuple $\bar{t} \in \mathbf{dom}^{|\bar{x}|}$.

An *annotated disjunction* $v_1::a_1; \dots; v_n::a_n$, where a_1, \dots, a_n are ground atoms and $\left(\sum_{1 \leq i \leq n} v_i \right) \leq 1$, denotes that a_1, \dots, a_n are mutually exclusive probabilistic events happening with probabilities v_1, \dots, v_n . It can be encoded as:

$$\begin{aligned} p_1::sw_1(_) \\ \vdots \\ p_n::sw_n(_) \\ a_1(\bar{t}_1) \leftarrow sw_1(\bar{t}_1) \\ \vdots \\ a_n(\bar{t}_n) \leftarrow \neg sw_1(\bar{t}_1), \dots, \neg sw_{n-1}(\bar{t}_{n-1}), sw_n(\bar{t}_n), \end{aligned}$$

where each p_i , for $1 \leq i \leq n$, is $v_i \cdot \left(1 - \sum_{1 \leq j < i} v_j \right)^{-1}$. Probabilistic rules can be easily extended to support annotated disjunctions in their heads.

Example A.1. Let Σ be a first-order signature with two predicate symbols V and W , both with arity 1, \mathbf{dom} be the domain $\{a, b\}$, and p be the program consisting of the facts $^{1/4}::T(a)$ and $^{1/2}::T(b)$, the annotated disjunction $^{1/4}::W(a); ^{1/2}::W(b)$, and the rule $^{1/2}::T(x) \leftarrow W(x)$.

The probability associated to each (Σ, \mathbf{dom}) -structure by the program p is shown in the following table.

		W			
		\emptyset	$\{a\}$	$\{b\}$	$\{a, b\}$
T	\emptyset	$3/32$	$3/64$	$3/32$	0
	$\{a\}$	$1/32$	$5/64$	$1/32$	0
	$\{b\}$	$3/32$	$3/64$	$9/32$	0
	$\{a, b\}$	$1/32$	$5/64$	$3/32$	0

The empty structure has probability $3/32$. The only grounded instance whose well-formed model is the empty database is the instance i_1 that does not contain grounded atoms. Its probability is $3/32$ because the probability that $T(a)$ is not in i_1 is $3/4$, the probability that $T(b)$ is not in i_1 is $1/2$, and the probability that neither $W(a)$ nor $W(b)$ are in i_1 is $1/4$ and all these events are independent.

The probability of some structures is determined by more than one grounded instance. For example, the probability of the structure s where $s(T) = \{a, b\}$ and $s(W) = \{a\}$ is $5/64$. There are two grounded instances i_2 and i_3 whose well-founded model is s . The instance i_2 has probability $1/16$ and it consists of the atoms $T(b), W(a), sw(a)$ and the rule $T(x) \leftarrow W(x), sw(x)$, whereas the instance i_3 has probability $1/64$ and it consists of the atoms $T(a), T(b), W(a)$ and the rule $T(x) \leftarrow W(x), sw(x)$. Note that before computing the ground instances, we translated probabilistic rules and annotated disjunctions into standard PROBLOG rules. ■

The belief programs in ATKLOG are based on PROBLOG, which is a very expressive probabilistic language and it can be used to encode a wide range of probabilistic models. As a result, performing exact inference in PROBLOG is intractable in general, i.e., it is $\#P$ -hard in terms of data complexity.

We identify acyclic PROBLOG programs, a restricted class of PROBLOG programs where exact inference can be done efficiently, in linear time in terms of data complexity.

We first introduce acyclic programs. Afterwards, we introduce relaxed acyclic programs, a generalization of acyclic programs, and we provide an encoding of relaxed acyclic programs as a Bayesian Network. Observe that here we state again (sometimes with some generalizations) some of the concepts we introduced in the main paper. Note also that in the main paper we do not distinguish between acyclic programs and relaxed acyclic programs.

A. Acyclic Programs

Here, we define acyclic PROBLOG programs. Intuitively, an acyclic program is a PROBLOG program whose probability distribution can be represented using a poly-tree Bayesian Network. While inference for general Bayesian Networks is intractable, i.e., NP -hard, inference for poly-tree Bayesian Networks can be done in linear time in the size of the corresponding Bayesian Network.

1) *Acyclic PROBLOG Programs:* In the following, let Σ be a first-order signature, \mathbf{dom} be a finite domain, and p be a (Σ, \mathbf{dom}) -program. Without loss of generality, in the following we assume that (1) inside rules (not ground atoms), constants are used just inside equality and inequality atoms, i.e., instead of $B(x, 7) \leftarrow D(x, y)$ we consider only the equivalent rule $B(x, k) \leftarrow D(x, y), k = 7$, and (2) there are no distinct v_1 and v_2 such that $v_1::a(\bar{c}) \in p$ and $v_2::a(\bar{c}) \in p$ (this is without loss of generality since multiple occurrences of the same ground atom with different probabilities can be represented as a single ground atom an adjusted probability). Observe that in the following we re-define, extend, or make more precise the definitions from §VI.

Dependency Graph. The p -dependency graph, denoted $graph(p)$, is the directed labelled graph (N, E) where $N = \Sigma$ and $E = \{pred(body(r, i)) \xrightarrow{r, i} pred(head(r)) \mid r \in p \wedge i \in \mathbb{N}\}$.

Directed Paths and Cycles. A directed path in $graph(p)$ is a path $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ in $graph(p)$ such that for all $1 \leq j < n$, $p_j = p_{j+1}$. Two directed paths $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ and $pr'_1 \xrightarrow{r'_1, i'_1} \dots \xrightarrow{r'_{m-1}, i'_{m-1}} pr'_m$ are head-connected iff $pr_1 = pr'_1$. Similarly, two directed paths $pr'_1 \xrightarrow{r'_1, i'_1} \dots \xrightarrow{r'_{n-1}, i'_{n-1}} pr'_n$ and $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{m-1}, i_{m-1}} pr'_m$ are tail-connected iff $pr_n = pr'_m$. Note that an empty path consists of just one node in N , i.e., one predicate symbol. A directed cycle in $graph(p)$ is a directed path $pr_1 \xrightarrow{r_1, i_1}$

$\dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ in $graph(p)$ such that $pr_n = pr_1$. Given a directed path $P = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ in $graph(p)$, we denote by $start(P)$ the predicate symbol pr_1 and by $end(P)$ the predicate symbol pr_n . We say that a directed cycle C is simple iff (1) each edge occurs at most once in C , and (2) there is a predicate symbol pr such that pr occurs twice in C and for all $pr' \neq pr$, pr' occurs at most once in C .

Undirected Paths and Cycles. A directed path P_1 in $graph(p)$ is connected with P_2 iff $end(P_1) = end(P_2)$, $start(P_1) = start(P_2)$, or $end(P_1) = start(P_2)$. An undirected path in $graph(p)$ is a sequence of directed paths P_1, \dots, P_n such that, for all $1 \leq i < n$, P_i is connected with P_{i+1} . An undirected cycle in $graph(p)$ is a sequence of directed paths P_1, \dots, P_n such that (1) P_1, \dots, P_n is an undirected path and (2) P_n is connected with P_1 . Note that a directed path (respectively cycle) is also an undirected path (respectively cycle). We say that two (directed or undirected) cycles are equivalent iff they are the same cycle.

Undirected Unsafe Structures. An undirected unsafe structure in $graph(p)$ is quadruple $\langle D_1, D_2, D_3, U \rangle$ such that (1) D_1, D_2, D_3 are directed paths in $graph(p)$, (2) U is an undirected path in $graph(p)$, (3) D_1 and D_2 are non-empty or D_2 and D_3 are non-empty, (4) D_1 and D_2 are head-connected, (5) D_2 and D_3 are tail-connected, and (6) D_1, U, D_3, D_2 is an undirected cycle in $graph(p)$. An undirected unsafe structure $\langle D_1, D_2, D_3, U \rangle$ covers an undirected cycle U' iff D_1, U, D_3, D_2 is equivalent to U' .

Directed Unsafe Structures. A directed unsafe structure in $graph(p)$ is a directed cycle C in $graph(p)$. A directed unsafe structure C covers a directed cycle C' iff C is equivalent to C' .

Deriving Ordering Annotations. A Σ -ordering annotation is $ORD(A)$, where $A \subseteq \Sigma$ and there is a $k \in \mathbb{N}$ such that for all $a \in A$, $|a| = 2k$. Let \preceq be a well-founded partial order over 2^Σ . We say that the ordering annotation $ORD(A)$ can be derived from the program p given \preceq , written $p, \preceq \models ORD(A)$, iff one of the following conditions hold:

- 1) For all $pr \in A$, there is no rule r in p such that $pred(head(r)) = pr$ and $body(r) \neq \emptyset$, there is no $pr' \in \Sigma$ such that $pr' \prec pr$ and $p, \preceq \models ORD(pr')$, and the transitive closure of the relation $R = \bigcup_{pr \in A} \{(\bar{c}, \bar{v}) \mid \exists r \in p. ((head(r) = pr(\bar{c}, \bar{v}) \vee \exists v'. head(r) = v'::pr(\bar{c}, \bar{v})) \wedge body(r) = \emptyset) \wedge |\bar{c}| = |\bar{v}| = |pr|/2\}$ is strict partial order.
- 2) There is a set $A' \subseteq \Sigma$ such that (1) $A' \preceq A$, (2) $p, \preceq \models ORD(A')$, (3) there is a $pr \in \Sigma$ and a $pr' \in A'$ such that (a) $|pr| = |pr'|$, (b) $A = (A' \setminus \{pr'\}) \cup \{pr\}$, and (c) for all rules r in p such that $pred(head(r)) = pr$, there are sequences of variables \bar{x} and \bar{y} and an i , such that $head(r) = pr(\bar{x}, \bar{y})$, $body(r, i) = pr'(\bar{x}, \bar{y})$, and $|\bar{x}| = |\bar{y}|$.

A Σ -ordering template \mathcal{O} is a set of ordering annotations. We say that p complies with \mathcal{O} with respect to \preceq iff for all ordering annotations $a \in \mathcal{O}$, $p, \preceq \models a$ holds. We say that p complies

with \mathcal{O} iff there exists a well-founded partial order \preceq over Σ such that p complies with \mathcal{O} with respect to \preceq . Intuitively, a program p complies with an ordering template \mathcal{O} iff all sets of predicates in \mathcal{O} represent strict partial orders in all possible models of p .

Deriving Disjointness Annotations. A Σ -disjointness annotation is $DIS(a, a')$, where $a, a' \in \Sigma$ and $|a| = |a'|$. Let \preceq be a well-founded partial order over Σ^2 . We say that the disjointness annotation $DIS(pr, pr')$ can be derived from the program p given \preceq , written $p, \preceq \models DIS(pr, pr')$, iff one of the following conditions hold:

- 1) There are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$ or $pred(head(r)) = pr'$, and there is no $(pr_1, pr_2) \in \Sigma^2$ such that $(pr_1, pr_2) \prec (pr, pr')$ and $p, \preceq \models DIS(pr_1, pr_2)$, and the sets $\{\bar{c} \mid pr(\bar{c}) \in p \vee \exists v. v :: pr(\bar{c}) \in p\}$ and $\{\bar{c} \mid \exists r \in p. head(r) = pr'(\bar{c}) \vee \exists v. v :: pr'(\bar{c}) \in p\}$ are disjoint.
- 2) There are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$, and there is $(pr, pr'_1) \in \Sigma^2$ such that (a) $p, \preceq \models DIS(pr, pr'_1)$, (b) $(pr, pr'_1) \prec (pr, pr')$, and (c) for all rules $r \in p$ such that $head(r) = pr'(\bar{x})$ there is $1 \leq i \leq |body(r)|$ such that $body(r, i) = pr'_1(\bar{x})$.
- 3) There are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr'$, there is $(pr_1, pr') \in \Sigma^2$ such that (a) $p, \preceq \models DIS(pr_1, pr')$, (b) $(pr_1, pr') \prec (pr, pr')$, and (c) for all rules $r \in p$ such that $head(r) = pr(\bar{x})$ there is $1 \leq i \leq |body(r)|$ such that $body(r, i) = pr_1(\bar{x})$.
- 4) There is $(pr_1, pr'_1) \in \Sigma^2$ such that (a) $p, \preceq \models DIS(pr_1, pr'_1)$, (b) $(pr_1, pr'_1) \prec (pr, pr')$, and (c) for all rules $r, r' \in p$ such that $head(r) = pr(\bar{x})$ and $head(r') = pr'(\bar{x}')$, there are $1 \leq i \leq |body(r)|$, $1 \leq i' \leq |body(r')|$, and $(pr_1, pr'_1) \in \mathcal{D}$ such that $body(r, i) = pr_1(\bar{x})$, $body(r', i') = pr'_1(\bar{x}')$.

A Σ -disjointness template \mathcal{D} is a set of disjointness annotations. We say that p complies with \mathcal{D} iff there exists a well-founded partial order \preceq over Σ^2 such that for all disjointness annotations $a \in \mathcal{D}$, $p, \preceq \models a$ holds. Intuitively, a program p complies with a disjointness template \mathcal{D} iff all pairs of predicates (pr, pr') involved in a disjointness annotation in \mathcal{D} are pairwise disjoint in all possible models of p .

Deriving Uniqueness Annotations. A Σ -uniqueness annotation is $UNQ(a, K)$, where $a \in \Sigma$ and K is a non-empty subset of $\{1, \dots, |a|\}$. Let \preceq be a well-founded order over Σ . We say that the uniqueness annotation $UNQ(pr, K)$ can be derived from p given \preceq , written $p, \preceq \models UNQ(pr, K)$, iff one of the following conditions hold:

- 1) There are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$, and there is no pr' such that $pr' \prec pr$, and for all \bar{t}, \bar{v} in $\{\bar{c} \mid pr(\bar{c}) \in p \vee \exists v. v :: pr(\bar{c}) \in p\}$, if $\bar{t}(i) = \bar{v}(i)$ for all $i \in K$, then $\bar{t} = \bar{v}$.
- 2) Both the following conditions hold:
 - a) For each rule $r \in p$ such that $head(r) = pr(\bar{x})$ and $body(r) \neq \emptyset$, there is a mapping μ' associating to each predicate $pr' \prec pr$ a set K such that (a) $p,$

$\preceq \models UNQ(pr', \mu(pr'))$, and (b) the following hold:

$$V \subseteq \bigcup_{l \in bound(head(r), K, body^+(r), \mu')} vars(l),$$

where $V = \{\bar{x}(i) \mid i \notin K \wedge \bar{x}(i) \in Var\}$, $u(\bar{y}, K) = \{\bar{y}(i) \mid i \in K\}$, and

$$bound(h, K, L, \mu') = \bigcup_{\substack{b(\bar{y}) \in L \wedge b \prec pr \wedge \\ u(\bar{y}, \mu'(b)) \subseteq u(args(h), K)}} \{b(\bar{y})\} \cup \bigcup_{\substack{b(\bar{y}) \in L \wedge b \prec pr \wedge \\ \exists l' \in bound(h, K, L, \mu'). (u(\bar{y}, \mu'(b)) \subseteq vars(l'))}} \{b(\bar{y})\}.$$

- b) For all rules $r_1, r_2 \in p$, if $r_1 \neq r_2$, $pred(head(r_1)) = pred(head(r_2)) = pr$, and $K \neq \{1, \dots, |pr|\}$, then there is a value $i \in K$ such that $args(head(r_1)) \in \mathbf{dom}$, $args(head(r_2)) \in \mathbf{dom}$, and $args(head(r_1)) \neq args(head(r_2))$.

A Σ -uniqueness template \mathcal{U} is a set of uniqueness annotations. We say that p complies with \mathcal{U} with respect to \preceq iff for all uniqueness annotations $a \in \mathcal{U}$, $p, \preceq \models a$ holds. We say that p complies with \mathcal{U} iff there exists a well-founded partial order \preceq over Σ such that p complies with \mathcal{U} with respect to \preceq . Intuitively, a program complies with a uniqueness template \mathcal{U} iff for all predicates pr each set K such that $UNQ(pr, K) \in \mathcal{U}$ represents a primary key for pr .

Propagation Maps. We use propagation maps to track how information flows inside rules. Given a rule r and a literal $l \in body(r)$, the (r, l) -vertical map is the mapping μ from $\{1, \dots, |l|\}$ to $\{1, \dots, |head(r)|\}$ such that $\mu(i) = j$ iff $args(l)(i) = args(head(r))(j)$ and $args(l)(i) \in Var$. Given a rule r and literals l, l' in r 's body, the (r, l, l') -horizontal map is the mapping μ from $\{1, \dots, |l|\}$ to $\{1, \dots, |l'|\}$ such that $\mu(i) = j$ iff $args(l)(i) = args(l')(j)$ and $args(l)(i) \in Var$.

We say that a path links to a literal l if information flows along the rules to l . This can be formalized by posing constraints on the mapping obtained by combining horizontal and vertical maps along the path. Formally, given a literal l and a mapping $\nu : \mathbb{N} \rightarrow \mathbb{N}$, a directed path $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ ν -downward links to l iff there is a $0 \leq j < n - 1$ such that the function $\mu := \mu' \circ \mu_j \circ \dots \circ \mu_1$ satisfies $\mu(k) = \nu(k)$ for all k for which $\nu(k)$ is defined, where for $1 \leq h \leq j$, μ_h is the vertical map connecting $body(r_h, i_h)$ and r_h , and μ' is the horizontal map connecting $body(r_{j+1}, i_{j+1})$ with l . Similarly, a directed path $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ ν -upward links to l iff there is a $1 \leq j \leq n - 1$ such that the function $\mu := \mu'^{-1} \circ \mu_{j+1}^{-1} \circ \dots \circ \mu_{n-1}^{-1}$ satisfies $\mu(k) = \nu(k)$ for all k for which $\nu(k)$ is defined, where μ_h is the $(r_h, body(r_h, i_h))$ -vertical map, for $j < h \leq n - 1$, and μ' is the (r_j, l) -vertical map. A path P links to a predicate symbol a iff there is an atom $a(\bar{x})$ such that P links to $a(\bar{x})$.

Negation-guarded Programs. A rule r is *negation-guarded* [9] iff for all negative literals l in r , $vars(l) \subseteq \bigcup_{l' \in body^+(r)} vars(l')$. We say that a program p is *negation-guarded* if all rules $r \in p$ are negation-guarded.

Join Trees. A join tree represents how multiple predicate symbols in a rule share variables. A *join tree for a rule r* is a rooted labelled tree $(N, E, \text{root}, \lambda)$, where $N \subseteq \text{body}(r)$, E is a set of edges (i.e., unordered pairs over N^2), $\text{root} \in N$ is the tree's root, and λ is the labelling function. Moreover, we require that for all $n, n' \in N$, if $n \neq n'$ and $(n, n') \in E$, then $\lambda(n, n') = \text{vars}(n) \cap \text{vars}(n')$ and $\lambda(n, n') \neq \emptyset$. A join tree $(N, E, \text{root}, \lambda)$ *covers* a literal l iff $l \in N$. Given a join tree $J = (N, E, \text{root}, \lambda)$ and a node $n \in N$, the *support of n* , denoted $\text{support}(n)$, is the set $\text{vars}(\text{head}(r)) \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\} \cup \{\text{vars}(n') \mid n' \in \text{anc}(J, n)\}$, where $\text{anc}(J, n)$ is the set of n 's ancestors in J , i.e., the set of all nodes on the path from root to n (if such a path exists).

Let \mathcal{U} be a uniqueness template. A join tree $J = (N, E, \text{root}, \lambda)$ is \mathcal{U} -*strongly connected* iff for all positive literals $l \in N$, there is a set $K \subseteq \{i \mid \bar{x} = \text{args}(l) \wedge \bar{x}(i) \in \text{support}(l)\}$ such that $\text{UNQ}(\text{pred}(l), K) \in \mathcal{U}$ and for all negative literals $l \in N$, $\text{vars}(l) \subseteq \text{support}(l)$. In contrast, a join tree $(N, E, \text{root}, \lambda)$ is \mathcal{U} -*weakly connected* iff for all $(a(\bar{x}), a'(\bar{x}')) \in E$, there are $K \subseteq \{i \mid \bar{x}(i) \in L\}$ and $K' \subseteq \{i \mid \bar{x}'(i) \in L'\}$ such that $\text{UNQ}(a, K), \text{UNQ}(a', K') \in \mathcal{U}$, where $L = \lambda(a(\bar{x}), a'(\bar{x}'))$.

Connected rules. A connected rule r ensure that a grounding of r is fully determined either by the assignment to the head's variables (strongly connected rule) or to the variable of any literal in r 's body (weakly connected rule). This is done by exploiting uniqueness annotations and the rule's structure. In the following, let \mathcal{U} be a uniqueness template.

We say that a rule r is *strongly connected for \mathcal{U}* iff there exist join trees J_1, \dots, J_n such that (a) the trees cover all literals in $\text{body}(r)$, and (b) for each $1 \leq i \leq n$, J_i is strongly connected for \mathcal{U} . We call J_1, \dots, J_n a *witness for the strong connectivity of r with respect to \mathcal{U}* . A strongly connected rule r guarantees that for any two groundings r', r'' of r , if $\text{head}(r') = \text{head}(r'')$, then $r' = r''$.

Given a rule r , a set of literals L , and a uniqueness template \mathcal{U} , we say that a literal $l \in \text{body}(r)$ is (r, \mathcal{U}, L) -*strictly guarded* iff (1) $\text{vars}(l) \subseteq \bigcup_{l' \in L \cap \text{body}^+(r)} \text{vars}(l') \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\}$, and (2) there is a positive literal $a(\bar{x}) \in L \cap \text{body}^+(r)$ and an annotation $\text{UNQ}(a, K) \in \mathcal{U}$ such that $\{\bar{x}(i) \mid i \in K\} \subseteq \text{vars}(l)$. We say that a rule r is *weakly connected for \mathcal{U}* iff there exists a join tree $J = (N, E, \text{root}, \lambda)$ such that (a) J is weakly connected for \mathcal{U} , (b) $N \subseteq \text{body}^+(r)$, and (c) all literals in $\text{body}(r) \setminus N$ are (r, \mathcal{U}, N) -strictly guarded. We call J a *witness for the weak connectivity of r with respect to \mathcal{U}* . A weakly connected rule r guarantees that for any two groundings r', r'' of r , if $\text{body}(r', i) = \text{body}(r'', i)$ for some i , then $r' = r''$.

Guarded Cycles. A set of predicates \mathcal{O} *guards a directed cycle* $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n \xrightarrow{r_n, i_n} pr_1$ iff there are integers $1 \leq y_1 < y_2 < \dots < y_e = n$, literals $o_1(\bar{x}_1), \dots, o_e(\bar{x}_e)$ (where $o_j \in \mathcal{O}$ and $|\bar{x}_j| = |o_j|$), a non-empty set $K \subseteq \{1, \dots, |pr_1|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |o_1|/2\}$ such that for each $0 \leq k < e$, (1) $pr_{y_k} \xrightarrow{r_{y_k}, i_{y_k}} \dots \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν -downward

connects to $o_{k+1}(\bar{x}_{k+1})$, and (2) $pr_{y_{k+1}-1} \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν' -upward connects to $o_{k+1}(\bar{x}_{k+1})$, where $\nu'(i) = \nu(x) + |o_1|/2$ for all $1 \leq i \leq |o_1|/2$, and $y_0 = 1$. A directed cycle C is guarded by an ordering template \mathcal{O} iff there is an annotation $\text{ORD}(A) \in \mathcal{O}$ such that A guards C .

Head-Guarded Paths. A pair of head-connected non-empty directed paths (P_1, P_2) is $(\mathcal{D}, \mathcal{U})$ -*head guarded* iff one of the following conditions hold:

- if $P_1 = P_2$, then all rules in P_1 are weakly connected for \mathcal{U} , or
- if $P_1 \neq P_2$, then there is an annotation $\text{DIS}(pr, pr') \in \mathcal{D}$, a set $K \subseteq \{1, \dots, |a|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |pr|\}$ such that P_1 ν -downward links to pr and P_2 ν -downward links to pr' .

Given two ground paths P'_1 and P'_2 corresponding to P_1 and P_2 , the first condition ensures that $P'_1 = P'_2$ whereas the second ensures that P'_1 or P'_2 are not in the ground graph.

Tail-Guarded Paths. A pair of tail-connected non-empty directed paths (P_1, P_2) are $(\mathcal{D}, \mathcal{U})$ -*tail guarded* iff one of the following conditions hold:

- if $P_1 = P_2$, then all rules in P_1 are strongly connected for \mathcal{U} , or
- if $P_1 \neq P_2$, there is an annotation $\text{DIS}(pr, pr') \in \mathcal{D}$, a set $K \subseteq \{1, \dots, |a|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |pr|\}$, such that P_1 ν -upward links to pr and P_2 ν -upward links to pr' .

Guarded Unsafe Structures. An undirected unsafe structure $\langle D_1, D_2, D_3, U \rangle$ is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -*guarded* iff either (D_1, D_2) is $(\mathcal{D}, \mathcal{U})$ -head-guarded or (D_2, D_3) is $(\mathcal{D}, \mathcal{U})$ -tail-guarded. In contrast, a directed unsafe structure C is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -*guarded* iff C is guarded by \mathcal{O} .

Acyclic Programs. A (Σ, \mathbf{dom}) -program p is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -*acyclic*, where \mathcal{O} is a Σ -ordering template, \mathcal{D} is a Σ -disjointness template, and \mathcal{U} is a Σ -uniqueness template, iff

- 1) p complies with \mathcal{O} , \mathcal{D} , and \mathcal{U} ,
- 2) p is a negation-guarded program, and
- 3) for all undirected cycles U in $\text{graph}(p)$ that are not simple directed cycles, there is an undirected unsafe structure S in $\text{graph}(p)$ that covers U , and is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded
- 4) for all directed cycles C in $\text{graph}(p)$, there is a directed unsafe structure S in $\text{graph}(p)$ that covers C , and is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded

A (Σ, \mathbf{dom}) -program p is *acyclic* iff there are a Σ -ordering template \mathcal{O} , a Σ -disjointness template \mathcal{D} , and a Σ -uniqueness template \mathcal{U} such that p is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -acyclic.

2) *Exact Grounding:* We now introduce the exact grounding of a PROBLOG program p given a p -probabilistic assignment. The exact grounding encodes the PROBLOG semantics.

Let p be a PROBLOG program and f be a p -probabilistic assignment. Furthermore, let μ be the mapping from predicate symbols in p to \mathbb{N} : $\mu(a) = \max_{t \in \text{paths}(p, a)} (w(t))$, where $\text{paths}(p, a)$ is the set of all directed paths that ends in a in p 's

dependency graph and the weight of each path is

$$w(t) = \begin{cases} 0 & \text{if } t = \epsilon \\ w(t') & \text{if } t = a_1 \xrightarrow{r_1, i_1} a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_n, i_n} a_{n+1} \\ & \wedge t' = a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_n, i_n} a_{n+1} \\ & \wedge \text{body}(r_1, i_1) \in \mathcal{A}_{\Sigma, \text{dom}}^+ \\ 1 + w(t') & \text{if } t = a_1 \xrightarrow{r_1, i_1} a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_n, i_n} a_{n+1} \\ & \wedge t' = a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_n, i_n} a_{n+1} \\ & \wedge \text{body}(r_1, i_1) \notin \mathcal{A}_{\Sigma, \text{dom}}^+ \end{cases}$$

We remark that μ is a valid stratification for the program p . From more details about logic programming with stratified negation we refer the reader to [2].

The functions $g_f(p, j, i)$, $g_f(p, j)$, and $g_f(p)$ are defined as follows:

$$\begin{aligned} g_f(p, 0, 0) &= \{a(\bar{c}) \mid a(\bar{c}) \in p \wedge \mu(a) = 0\} \cup \\ &\quad \{a(\bar{c}) \mid \exists v. v::a(\bar{c}) \in p \wedge f(v::a(\bar{c})) = \top \\ &\quad \wedge \mu(a) = 0\} \\ g_f(p, 0, i) &= g_f(p, 0, i-1) \cup \\ &\quad \{\text{head}(r)\Theta \mid \Theta \in \text{ASS}(r) \wedge r \in p \wedge \\ &\quad \forall l \in \text{body}^+(r). l\Theta \in g_f(p, 0, i-1) \wedge \\ &\quad \text{body}^-(r) = \emptyset \wedge \mu(\text{pred}(\text{head}(r))) = 0\} \\ g_f(p, j, 0) &= g_f(p, j-1) \cup \\ &\quad \{a(\bar{c}) \mid a(\bar{c}) \in p \wedge \mu(a) = j\} \cup \\ &\quad \{a(\bar{c}) \mid \exists v. v::a(\bar{c}) \in p \wedge f(v::a(\bar{c})) = \top \\ &\quad \wedge \mu(a) = j\} \\ g_f(p, j, i) &= g_f(p, j, i-1) \cup \\ &\quad \{\text{head}(r)\Theta \mid \Theta \in \text{ASS}(r) \wedge r \in p \wedge \\ &\quad \forall l \in \text{body}^+(r). l\Theta \in g_f(p, j, i-1) \wedge \\ &\quad \forall l \in \text{body}^-(r). l\Theta \notin g_f(p, j-1) \wedge \\ &\quad \wedge \mu(\text{pred}(\text{head}(r))) = j\} \\ g_f(p, j) &= \bigcup_{i \in \mathbb{N}} g_f(p, j, i) \\ g_f(p) &= \bigcup_{j \in \mathbb{N}} g_f(p, j) \end{aligned}$$

Furthermore, the functions $g_f(p, r, j, i)$, $g_f(p, r, j)$, and $g_f(p, r)$ are defined as follows:

$$\begin{aligned} g_f(p, a(\bar{c}), j, i) &= \{a(\bar{c}) \mid a(\bar{c}) \in p \wedge \mu(a) = j\} \\ g_f(p, v::a(\bar{c}), j, i) &= \{a(\bar{c}) \mid v::a(\bar{c}) \in p \wedge f(v::a(\bar{c})) = \top \wedge \\ &\quad \mu(a) = j\} \\ g_f(p, r, 0, 0) &= \emptyset \end{aligned}$$

$$\begin{aligned} g_f(p, r, 0, i) &= g_f(p, r, 0, i-1) \cup \\ &\quad \{h\Theta \leftarrow l_1\Theta, \dots, l_n\Theta \mid \\ &\quad r = h \leftarrow l_1, \dots, l_n \wedge \\ &\quad \Theta \in \text{ASS}(r) \wedge \\ &\quad \forall l \in \text{body}^+(r). l\Theta \in g_f(p, 0, i-1) \wedge \\ &\quad \text{body}^-(r) = \emptyset \wedge \mu(\text{pred}(h)) = 0\} \\ g_f(p, r, j, 0) &= g_f(p, r, j-1) \\ g_f(p, r, j, i) &= g_f(p, r, j, i-1) \cup \\ &\quad \{h\Theta \leftarrow l_1\Theta, \dots, l_n\Theta \mid \\ &\quad r = h \leftarrow l_1, \dots, l_n \wedge \\ &\quad \Theta \in \text{ASS}(r) \wedge \\ &\quad \forall l \in \text{body}^+(r). l\Theta \in g_f(p, j, i-1) \wedge \\ &\quad \forall l \in \text{body}^-(r). l\Theta \notin g_f(p, j-1) \wedge \\ &\quad \mu(\text{pred}(h)) = j\} \\ g_f(p, r, j) &= \bigcup_{i \in \mathbb{N}} g_f(p, r, j, i) \\ g_f(p, r) &= \bigcup_{j \in \mathbb{N}} g_f(p, r, j) \end{aligned}$$

The fact below follows directly from g_f 's definition and the notion of well-founded models for stratified logic programs [2].

Fact B.1. For any PROBLOG program p and any p -probabilistic assignment f , $g_f(p) = \text{WMF}(\text{instance}(p, f))$.

3) *Relaxed Grounding:* Here we introduce the notion of relaxed grounding. This is a key notion for our acyclicity and correctness proofs and it is also a key component of our compilation from PROBLOG programs to poly-tree Bayesian Networks.

Given a PROBLOG acyclic program p , the *relaxed grounding of p* is a set M of ground atoms and rules containing all possible ground atoms and rules that can be derived in any ground instance of p . We remark that a relaxed grounding is an over-approximation of the ground atoms and/or rules that could be derived by the ground instances of p .

Formally, let Σ be a first-order signature, dom be a finite domain, p be a (Σ, dom) -acyclic PROBLOG, and $i \in \mathbb{N}$ be an natural number. The function $\text{ground}(p, i)$ is defined as follows: $\text{ground}(p, 0) = \{a(\bar{c}) \mid a(\bar{c}) \in p \vee \exists v. v::a(\bar{c}) \in p\}$, and $\text{ground}(p, i) = \text{ground}(p, i-1) \cup \{h\Theta \mid \Theta \in \text{ASS}(r) \wedge \exists r. r \in p \wedge h = \text{head}(r) \wedge \forall l \in \text{body}^+(r). l\Theta \in \text{ground}(p, i-1) \wedge \forall i, j. \text{consistent}(\text{body}(r, i)\Theta, \text{body}(r, j)\Theta)\}$, where $\text{ASS}(r)$ is the set of all mappings $\text{Var} \rightarrow \text{dom}$ that satisfy the equality and inequality constraints in r while $\text{consistent}(a(\bar{v}), \neg a(\bar{v})) = \text{consistent}(\neg a(\bar{v}), a(\bar{v})) = \perp$ and $\text{consistent}(l, l') = \top$ otherwise. The function $\text{ground}(p, r, i)$, where $r \in p$, is defined as follows: $\text{ground}(p, a(\bar{c}), i) = \{a(\bar{c})\}$, $\text{ground}(p, v::a(\bar{c}), i) = \{a(\bar{c})\}$, $\text{ground}(p, r, 0) = \emptyset$, and $\text{ground}(p, r, i)$ is the set $\{h\Theta \leftarrow l_1\Theta, \dots, l_n\Theta \mid \Theta \in \text{ASS}(r) \wedge \forall l_i \in \text{body}^+(r). l_i\Theta \in \text{ground}(p, i-1) \wedge \forall i, j. \text{consistent}(l_i\Theta, l_j\Theta)\}$, where $r = h \leftarrow l_1, \dots, l_n$.

Finally, the relaxed grounding of p , denoted by $ground(p)$, is the set $ground(p) = \bigcup_{i \in \mathbb{N}} ground(p, i)$, whereas $ground(p, r) = \bigcup_{i \in \mathbb{N}} ground(p, r, i)$.

The fact below follows directly from $ground$'s definition and the definition of acyclic PROBLOG programs.

Fact B.2. *For any acyclic PROBLOG program p and any p -probabilistic assignment f , $g_f(p) \subseteq ground(p)$ and $g_f(p, r) \subseteq ground(p, r)$.*

4) *Ground Graphs:* Let Σ be a first-order signature, \mathbf{dom} be a finite domain, and p be a (Σ, \mathbf{dom}) -acyclic PROBLOG program. The ground graph of p , denoted $gg(p)$, is the directed graph (N, E) defined as follows:

- There is one node $a(\bar{c})$ in N per literal $a(\bar{c})$ or $\neg a(\bar{c})$ used in a ground rule $r' \in ground(p, r)$ for some $r \in p$. For each rule $r \in p$, ground rule $s \in ground(p, r)$, and position $1 \leq j \leq |body(r')|$, there is a node (r, s, j) .
- There is an edge from $b(\bar{v})$ to (r, s, j) if $body(s, j) = b(\bar{v})$ or $body(s, j) = \neg b(\bar{v})$. There is an edge from (r, s, j) to $a(\bar{c})$ if $head(s) = a(\bar{c})$.

5) *Proofs about Templates:* Here we prove preliminary results about the semantic relationships enforced by the ordering, disjointness, and uniqueness templates.

Proposition B.1. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, \mathcal{O} be a Σ -ordering template, and \preceq be a well-founded partial order over 2^Σ . If p complies with \mathcal{O} with respect to \preceq and $ORD(A) \in \mathcal{O}$, then if $pr_1(\bar{a}_1, \bar{a}_2), pr_2(\bar{a}_2, \bar{a}_3), \dots, pr_{n-1}(\bar{a}_{n-1}, \bar{a}_n)$ are in $ground(p)$ and $\{pr_1, \dots, pr_{n-1}\} \subseteq A$, then $\bar{a}_1 \neq \bar{a}_n$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, \mathcal{O} be a Σ -ordering template, and \preceq be a well-founded partial order over 2^Σ . We assume that p complies with \mathcal{O} with respect to \preceq , $ORD(A) \in \mathcal{O}$, $pr_1(\bar{a}_1, \bar{a}_2), \dots, pr_{n-1}(\bar{a}_{n-1}, \bar{a}_n)$ are in $ground(p)$, and $\{pr_1, \dots, pr_{n-1}\} \subseteq A$. We now prove, by induction over \preceq , that the transitive closure of the unions of the relations induced by pr_1, \dots, pr_{n-1} over $ground(p)$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$. From this, it follows that if $pr_1(\bar{a}_1, \bar{a}_2), pr_2(\bar{a}_2, \bar{a}_3), \dots, pr_{n-1}(\bar{a}_{n-1}, \bar{a}_n)$ are in $ground(p)$, then $\bar{a}_1 \neq \bar{a}_n$.

Base Case. Let A be a set such that there is no $A' \preceq A$ such that $p, \preceq \models ORD(A')$. From this, $ORD(A) \in \mathcal{O}$, and p complies with \mathcal{O} with respect to \preceq , it follows that (1) for all $pr \in A$, there is no rule r in p such that $pred(head(r)) = pr$ and $body(r) \neq \emptyset$, and (2) the transitive closure of the relation $R = \bigcup_{pr \in A} \{(\bar{c}, \bar{v}) \mid \exists r \in p. ((head(r) = pr(\bar{c}, \bar{v}) \vee \exists v'. head(r) = v'::pr(\bar{c}, \bar{v})) \wedge body(r) = \emptyset) \wedge |\bar{c}| = |\bar{v}| = |pr|/2\}$ is strict partial order. From this and $\{(\bar{c}, \bar{v}) \mid \exists pr \in A. pr(\bar{c}, \bar{v}) \in ground(p) \wedge |\bar{c}| = |\bar{v}| = |pr|/2\} \subseteq R$, it follows that the transitive closure of the union of the relations induced by pr over $ground(p)$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$.

Induction Step. Assume that the relation induced by any $A' \prec A$ over $ground(p)$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$

if $p, \preceq \models ORD(A')$, where pr' is any predicate symbol in A' . Without loss of generality, we assume that there is at least one $A' \prec A$ such that $p, \preceq \models ORD(A')$ (otherwise the proof is the same as in the base case). We now prove that also the relation induced by A over $ground(p)$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$. Let pr, A', pr_1 be such that $A = (A' \setminus \{pr_1\}) \cup \{pr\}$, $A' \prec A$, $p, \preceq \models ORD(A')$, and $pr_1 \in A'$. Note that the previous values always exist. We now prove that also the relation induced by A over $ground(p)$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$. Since $pr \in A$, $A = (A' \setminus \{pr_1\}) \cup \{pr\}$, $pr_1 \in A'$, and $ORD(A) \in \mathcal{O}$, it follows that for all rules r in p such that $pred(head(r)) = pr$, there are sequences of variables \bar{x} and \bar{y} and an i , $head(r) = pr(\bar{x}, \bar{y})$, $body(r, i) = pr_1(\bar{x}, \bar{y})$, and $|\bar{x}| = |\bar{y}|$. From this, it follows that the $\{(\bar{v}, \bar{w}) \mid pr(\bar{v}, \bar{w}) \in ground(p)\} \subseteq \{(\bar{v}, \bar{w}) \mid pr_1(\bar{v}, \bar{w}) \in ground(p)\}$. From the induction hypothesis, it follows that the transitive closure of $\bigcup_{pr' \in A'} \{(\bar{v}, \bar{w}) \mid pr'(\bar{v}, \bar{w}) \in ground(p)\}$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$. From this, $pr_1 \in A'$, $A = (A' \setminus \{pr_1\}) \cup \{pr\}$, $pr_1 \in A'$, and $\{(\bar{v}, \bar{w}) \mid pr(\bar{v}, \bar{w}) \in ground(p)\} \subseteq \{(\bar{v}, \bar{w}) \mid pr_1(\bar{v}, \bar{w}) \in ground(p)\}$, it follows that the transitive closure of $\bigcup_{pr' \in A} \{(\bar{v}, \bar{w}) \mid pr'(\bar{v}, \bar{w}) \in ground(p)\}$ is a strict partial order over $\mathbf{dom}^{|pr|/2}$. This completes the proof of our claim. \square

Proposition B.2. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, \mathcal{D} be a Σ -disjointness template, and \preceq be a well-founded partial order over Σ^2 . If p complies with \mathcal{D} with respect to \preceq and $DIS(pr, pr') \in \mathcal{D}$, then there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, \mathcal{D} be a Σ -disjointness template, and \preceq be a well-founded partial order over Σ^2 . Furthermore, we assume that p complies with \mathcal{D} with respect to \preceq and $DIS(pr, pr') \in \mathcal{D}$. We prove, by induction over \preceq , that for any $DIS(pr, pr') \in \mathcal{D}$, there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$.

Base Case. Let (pr, pr') be a pair in Σ^2 such that there is no $(pr_1, pr'_1) \prec (pr, pr')$ and $p, \preceq \models DIS(pr, pr')$. From this, $DIS(pr, pr') \in \mathcal{D}$, and p complies with \mathcal{D} , it follows that (1) there are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$ or $pred(head(r)) = pr'$, and (2) the sets $A = \{\bar{c} \mid pr(\bar{c}) \in p \vee \exists v'. v'::pr(\bar{c}) \in p\}$ and $A' = \{\bar{c} \mid \exists r \in p. head(r) = pr'(\bar{c}) \vee \exists v'. v'::pr'(\bar{c}) \in p\}$ are disjoint. From this, it follows that there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$.

Induction Step. Assume that there is no tuple \bar{v} such that $pr_1(\bar{v}) \in ground(p)$ and $pr'_1(\bar{v}) \in ground(p)$ for any $(pr_1, pr'_1) \prec (pr, pr')$ such that $p, \preceq \models DIS(pr_1, pr'_1)$. We now prove that there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$. There are four cases:

- There are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$, and there is a $(pr, pr'_1) \in \Sigma^2$ such that (a) $DIS(pr, pr'_1) \in \mathcal{D}$, (b) $(pr, pr'_1) \prec (pr, pr')$, and (c) for all rules $r \in p$ such that $head(r) = pr'(\bar{x})$

there is $1 \leq i \leq |body(r)|$ such that $body(r, i) = pr'_1(\bar{x})$. From $(pr, pr'_1) \prec (pr, pr')$, $DIS(pr, pr') \in \mathcal{D}$, and the induction's hypothesis, it follows that there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'_1(\bar{v}) \in ground(p)$. Furthermore, the relation associated with pr' is a subset of the relation associated to pr'_1 . Therefore, there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$.

- There are no rules $r \in p$ such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr'$, and there is a $(pr_1, pr') \in \Sigma^2$ such that (a) $DIS(pr_1, pr') \in \mathcal{D}$, (b) $(pr_1, pr') \prec (pr, pr')$, and (c) for all rules $r \in p$ such that $head(r) = pr(\bar{x})$ there is $1 \leq i \leq |body(r)|$ such that $body(r, i) = pr_1(\bar{x})$. From $(pr_1, pr') \prec (pr, pr')$, $DIS(pr_1, pr') \in \mathcal{D}$, and the induction's hypothesis, it follows that there is no tuple \bar{v} such that $pr_1(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$. Furthermore, the relation associated with pr is a subset of the relation associated to pr_1 . Therefore, there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$.
- There is a $(pr_1, pr'_1) \in \Sigma^2$ such that (a) $DIS(pr_1, pr'_1) \in \mathcal{D}$, (b) $(pr_1, pr'_1) \prec (pr, pr')$, and (c) for all rules $r, r' \in p$ such that $head(r) = pr(\bar{x})$ and $head(r') = pr'(\bar{x}')$, there are $1 \leq i \leq |body(r)|$, $1 \leq i' \leq |body(r')|$ such that $body(r, i) = pr_1(\bar{x})$, $body(r', i') = pr'_1(\bar{x}')$. From $(pr_1, pr'_1) \prec (pr, pr')$, $DIS(pr_1, pr'_1) \in \mathcal{D}$, and the induction's hypothesis, it follows that there is no tuple \bar{v} such that $pr_1(\bar{v}) \in ground(p)$ and $pr'_1(\bar{v}) \in ground(p)$. Furthermore, the relation associated with pr is a subset of the relation associated with pr_1 and the relation associated with pr' is a subset of the relation associated with pr'_1 . Therefore, there is no tuple \bar{v} such that $pr(\bar{v}) \in ground(p)$ and $pr'(\bar{v}) \in ground(p)$.

This completes the proof of our claim. \square

Proposition B.3. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, \mathcal{U} be a Σ -uniqueness template, and \preceq be a well-founded order over Σ . If p complies with \mathcal{U} with respect to \preceq and $UNQ(pr, K) \in \mathcal{U}$, then for all tuples $pr(\bar{v}), pr(\bar{w}) \in ground(p)$, if $\bar{v}(i) = \bar{w}(i)$ for all $i \in K$, then $\bar{v} = \bar{w}$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, \mathcal{U} be a Σ -uniqueness template, and \preceq be a well-founded partial order over \mathcal{U} . Furthermore, we assume that p complies with \mathcal{U} with respect to \preceq and $UNQ(pr, K) \in \mathcal{U}$. We prove, by induction over \preceq , that for all $UNQ(pr, K) \in \mathcal{U}$, and all tuples $pr(\bar{v}), pr(\bar{w}) \in ground(p)$, if $\bar{v}(i) = \bar{w}(i)$ for all $i \in K$, then $\bar{v} = \bar{w}$. Without loss of generality, in the following we assume that $K \neq \{1, \dots, |pr|\}$. If this is not the case, then our claim holds trivially.

Base Case. Let pr be a predicate such that there is no $pr' \prec pr$. From this and $UNQ(pr, K) \in \mathcal{U}$, it follows that there are no rules r in p such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$, and for all \bar{w}, \bar{v} in $A = \{\bar{c} \mid pr(\bar{c}) \in p \vee \exists v.v::pr(\bar{c}) \in p\}$, if $\bar{v}(i) = \bar{w}(i)$ for all $i \in K$, then $\bar{v} = \bar{w}$. From this and

$\{\bar{c} \mid pr(\bar{c}) \in ground(p)\} \subseteq A$, it follows that for all tuples $pr(\bar{v}), pr(\bar{w}) \in ground(p)$, if $\bar{v}(i) = \bar{w}(i)$ for all $i \in K$, then $\bar{v} = \bar{w}$.

Induction Step. Assume that the claim holds for any $pr' \prec pr$ and $K \subseteq \{1, \dots, |pr|\}$ such that $UNQ(pr, K) \in \mathcal{U}$ (we denote this induction hypothesis as \clubsuit). We now prove that the claim holds for pr as well. Without loss of generality, we assume that there is at least one rule r such that $body(r) \neq \emptyset$ and $pred(head(r)) = pr$ (otherwise the proof is identical to the base case). Assume, for contradiction's sake, that there are two ground atoms $pr(\bar{v}), pr(\bar{w}) \in ground(p)$ such that $\bar{v}(i) = \bar{w}(i)$ for all $i \in K$ but $\bar{v} \neq \bar{w}$. There are two cases:

- $pr(\bar{v})$ and $pr(\bar{w})$ are generated by two ground instances of the same rule r . From $\bar{v} \neq \bar{w}$, it follows that there is $j \notin K$ such that $\bar{v}(j) \neq \bar{w}(j)$. From this and the fact that both atoms are generated by the same rule r , it follows that there is a variable x such that $x \in \{\bar{x}(i) \mid i \notin K\}$ and $x = \bar{x}(j)$, where $\bar{x} = args(head(r))$. From this and $pr \in \mathcal{U}$, it follows that $x \in \bigcup_{l \in bound(head(r), K, body^+(r), \mu')} vars(l)$. We claim that the value of x is determined by the values of the variables whose positions are in K . From this and the fact that \bar{v} and \bar{w} agree on all values whose positions are in K , it follows that $\bar{v} = \bar{w}$, leading to a contradiction.

We now prove our claim that the value of x is determined by the values of the variables whose positions are in K . The value of the variable x is determined by the grounding of one of the atoms in $bound(head(r), K, body^+(r), \mu')$. We first slightly modify the definition of $bound$. We denote by $bd^0(h, K, L, \mu')$ the function $\bigcup_{\substack{b(\bar{y}) \in L \wedge b \prec pr \wedge \\ u(\bar{y}, \mu'(b)) \subseteq u(args(h), K)}} \{b(\bar{y})\}$ and by $bd^i(h,$

$K, L, \mu')$, where $i > 0$, the function $bd^{i-1}(h, K, L, \mu') \cup \bigcup_{\substack{b(\bar{y}) \in L \wedge b \prec pr \wedge \\ \exists l' \in bound^{i-1}(h, K, L, \mu'). (u(\bar{y}, \mu'(b)) \subseteq vars(l'))}} \{b(\bar{y})\}$. It

is easy to see that (1) $bound(head(r), K, body^+(r), \mu') = \bigcup_{i \in \mathbb{N}} bd^i(h, K, L, \mu')$, and (2) the fix-point is always reached in a finite amount of steps (bounded by $|body^+(r)|$). We now prove by induction on i that the groundings of the literals in $bd^i(h, K, L, \mu')$ is always determined by the values of the variables in $u(args(head(r)), K)$. From this, our claim trivially follows. For the base case, let $i = 0$. Then, for any literal $b(\bar{y}) \in bd^0(h, K, L, \mu')$, it follows that (a) $b \prec pr$, and (b) $u(\bar{y}, \mu'(b)) \subseteq u(args(head(r)), K)$. From this, p complies with \mathcal{U} , and the induction's hypothesis (\clubsuit), it follows that the variables in $UNQ(b, \mu'(b))$ uniquely determine the grounding of $b(\bar{y})$. From this and $u(\bar{y}, \mu'(b)) \subseteq u(args(head(r)), K)$, it follows that the variables in $u(args(head(r)), K)$ uniquely determine the grounding of $b(\bar{y})$. For the induction's step, assume that the claim hold for all $j < i$ (we denote this induction hypothesis as \spadesuit). Let $b(\bar{y}) \in bd^i(h, K, L, \mu')$ (without loss of generality, $b(\bar{y}) \notin bd^{i-1}(h, K, L, \mu')$). From this, it follows that $b \preceq pr$ and there is a $l' \in bound^{i-1}(h, K, L, \mu')$ such that $u(\bar{y}, \mu'(b)) \subseteq vars(l')$. From the

induction hypothesis (\spadesuit), it follows that the grounding of l' is directly determined by the grounding of the values of the variables in $u(\text{args}(\text{head}(r)), K)$. Furthermore, from $b \preceq pr$, $p, \preceq \models \text{UNQ}(b, \mu'(b))$, and the induction hypothesis (\clubsuit), it follows that the values of the variables in $u(\bar{y}, \mu'(b))$ uniquely determine the grounding of $b(\bar{y})$. Therefore, the values of the variables in $u(\text{args}(\text{head}(r)), K)$ uniquely determine the grounding of $b(\bar{y})$. This completes the proof of the claim.

- $pr(\bar{v})$ and $pr(\bar{w})$ are generated by ground instances of two different rules r_1 and r_2 . From this, $\text{UNQ}(pr, K) \in \mathcal{U}$, and $K \neq \{1, \dots, |pr|\}$, it follows that it follows that there is a value $i \in K$ such that $\text{args}(\text{head}(r_1)) \in \mathbf{dom}$, $\text{args}(\text{head}(r_2)) \in \mathbf{dom}$, and $\text{args}(\text{head}(r_1)) \neq \text{args}(\text{head}(r_2))$. Therefore, there is an $i \in K$ such that $\bar{v}(i) \neq \bar{w}(i)$. This contradicts the fact that $\bar{v}(i) = \bar{w}(i)$ for all $i \in K$.

This completes the proof of our claim. \square

6) *Proofs about Propagation Maps:* Here we prove some results about propagation maps, which we afterwards use in proving the main result.

Proposition B.4. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , and l be the i -th literal in $\text{body}(r)$. Furthermore, let μ be the (r, l) -vertical map. Given a rule $r' \in \text{ground}(p, r)$, then $\bar{b}(j) = \bar{h}(\mu(j))$ for any j such that $\mu(j)$ is defined, where $\bar{h} = \text{args}(\text{head}(r'))$ and $\bar{b} = \text{args}(\text{body}(r', i))$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , and l be the i -th literal in $\text{body}(r)$. Furthermore, let μ be the (r, l) -vertical map and $r' \in \text{ground}(p, r)$. From the definition of $\text{ground}(p, r)$, it follows that there is an assignment Θ from variables to elements in \mathbf{dom} such that $r' = r\Theta$. From this, $\bar{b}(j) = \Theta(\text{args}(l)(j))$ and $\bar{h}(\mu(j)) = \Theta(\text{args}(\text{head}(r))(\mu(j)))$. Note that from the definition of (r, l) -vertical map it follows that $\text{args}(l)(j) = \text{args}(\text{head}(r))(\mu(j))$. From this, $\bar{b}(j) = \bar{h}(\mu(j))$ whenever $\mu(j)$ is defined. \square

Proposition B.5. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , l be the i -th literal in $\text{body}(r)$, and l' be the j -th literal in $\text{body}(r)$. Furthermore, let μ be the (r, l, l') -horizontal map. Given a rule $r' \in \text{ground}(p, r)$, then $\bar{b}_1(k) = \bar{b}_2(\mu(k))$ for any k such that $\mu(k)$ is defined, where $\bar{b}_1 = \text{args}(\text{body}(r', i))$ and $\bar{b}_2 = \text{args}(\text{body}(r', j))$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , l be the i -th literal in $\text{body}(r)$, and l' be the j -th literal in $\text{body}(r)$. Furthermore, let μ be the (r, l, l') -horizontal map and $r' \in \text{ground}(p, r)$. From the definition of $\text{ground}(p, r)$, it follows that there is an assignment Θ from variables to elements in \mathbf{dom} such that $r' = r\Theta$. From this, $\bar{b}_1(j) = \Theta(\text{args}(l)(j))$ and $\bar{b}_2(\mu(j)) = \Theta(\text{args}(l')(\mu(j)))$.

Note that from the definition of (r, l, l') -vertical map it follows that $\text{args}(l)(j) = \text{args}(l')(\mu(j))$. From this, $\bar{b}_1(j) = \bar{b}_2(\mu(j))$ whenever $\mu(j)$ is defined. \square

Proposition B.6. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, $P = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ be a directed path in $\text{graph}(p)$, and $\nu : \mathbb{N} \rightarrow \mathbb{N}$ be a mapping. Furthermore, let $P' = a_1 \rightarrow (r_1, s_1, i_1) \rightarrow a_2 \rightarrow (r_2, s_2, i_2) \rightarrow \dots \rightarrow (r_{n-1}, s_{n-1}, i_{n-1}) \rightarrow a_n$ be a directed path in $\text{gg}(p)$ corresponding to P . If P ν -downward links to l , where l is the k -th literal in r_j , then $\bar{b}_1(m) = \bar{v}_2(\nu(m))$ whenever $\nu(m)$ is defined, where $\bar{b}_1 = \text{args}(\text{body}(s_1, i_1))$ and $\bar{v}_2 = \text{args}(\text{body}(s_j, k))$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, $P = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ be a directed path in $\text{graph}(p)$, and $\nu : \mathbb{N} \rightarrow \mathbb{N}$ be a mapping. Furthermore, let $P' = a_1 \rightarrow (r_1, s_1, i_1) \rightarrow a_2 \rightarrow (r_2, s_2, i_2) \rightarrow \dots \rightarrow (r_{n-1}, s_{n-1}, i_{n-1}) \rightarrow a_n$ be a directed path in $\text{gg}(p)$ corresponding to P . Assume that P ν -downward links to l , where l is the k -th literal in r_j . From this, it follows that the function $\mu := \mu' \circ \mu_j \circ \dots \circ \mu_1$ satisfies $\mu(m) = \nu(m)$ for all m for which $\nu(k)$ is defined, where for $1 \leq h \leq j$, μ_h is the vertical map connecting $\text{body}(r_h, i_h)$ and r_h , and μ' is the horizontal map connecting $\text{body}(r_{j+1}, i_{j+1})$ with l . By repeatedly applying Proposition B.4 to the rules in P' , we have that $\bar{b}_1(m) = \bar{b}_j(\phi(m))$ whenever $\phi(m)$ is defined, where $\bar{b}_1 = \text{args}(\text{body}(s_1, i_1))$, $\bar{b}_j = \text{args}(\text{body}(s_j, i_j))$, and $\phi = \mu_j \circ \dots \circ \mu_1$. Moreover, by applying Proposition B.5, we have that $\bar{b}_j(m) = \bar{v}_2(\mu'(m))$ whenever $\mu'(m)$ is defined, where $\bar{b}_j = \text{args}(\text{body}(s_j, i_j))$ and $\bar{v}_2 = \text{args}(\text{body}(s_j, k))$. Therefore, we have that $\bar{b}_1(m) = \bar{v}_2(\mu(m))$ whenever $\mu'(m)$ is defined (by composing the previous results). From this and $\mu(m) = \nu(m)$ for all m for which $\nu(m)$ is defined, it follows that $\bar{b}_1(m) = \bar{v}_2(\nu(m))$ whenever $\mu'(m)$ is defined. \square

Proposition B.7. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, $P = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ be a directed path in $\text{graph}(p)$, and $\nu : \mathbb{N} \rightarrow \mathbb{N}$ be a mapping. Furthermore, let $P' = a_1 \rightarrow (r_1, s_1, i_1) \rightarrow a_2 \rightarrow (r_2, s_2, i_2) \rightarrow \dots \rightarrow (r_{n-1}, s_{n-1}, i_{n-1}) \rightarrow a_n$ be a directed path in $\text{gg}(p)$ corresponding to P . If P ν -upward links to l , where l is the k -th literal in r_j , then $\bar{b}_n(m) = \bar{v}_2(\nu(m))$ whenever $\nu(m)$ is defined, where $\bar{b}_n = \text{args}(\text{head}(r_{n-1}))$ and $\bar{v}_2 = \text{args}(\text{body}(s_j, k))$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, $P = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n$ be a directed path in $\text{graph}(p)$, and $\nu : \mathbb{N} \rightarrow \mathbb{N}$ be a mapping. Furthermore, let $P' = a_1 \rightarrow (r_1, s_1, i_1) \rightarrow a_2 \rightarrow (r_2, s_2, i_2) \rightarrow \dots \rightarrow (r_{n-1}, s_{n-1}, i_{n-1}) \rightarrow a_n$ be a directed path in $\text{gg}(p)$ corresponding to P . Assume that P ν -upward links to l , where l is the k -th literal in r_j . From this, it follows that the function $\mu := \mu'^{-1} \circ \mu_{j+1}^{-1} \circ \dots \circ \mu_{n-1}^{-1}$ satisfies $\mu(k) = \nu(k)$ for all k for which $\nu(k)$ is defined,

where μ_h is the $(r_h, \text{body}(r_h, i_h))$ -vertical map, for $j < h \leq n-1$, and μ' is the (r_j, l) -vertical map. By repeatedly applying Proposition B.4 to the rules in P' , we obtain that $\overline{b_{j+1}}(\phi^{-1}(m)) = \overline{b_n}(m)$ whenever $\phi^{-1}(m)$ is defined, where $\overline{b_{j+1}} = \text{args}(\text{body}(s_{j+1}, i_{j+1}))$, $\overline{b_n} = \text{args}(\text{head}(s_{n-1}))$, and $\phi^{-1} = \mu_{j+1}^{-1} \circ \dots \circ \mu_{n-1}^{-1}$. Furthermore, by applying Proposition B.4 to the (r_j, l) -vertical map, we have that $\overline{v_2}(\mu'^{-1}(m)) = \overline{v_1}(m)$ whenever $\mu'^{-1}(m)$ is defined, where $\overline{v_1} = \text{args}(\text{head}(s_j))$ and $\overline{v_2} = \text{args}(\text{body}(s_j, k))$. From this and $\mu(m) = \nu(\overline{m})$ for all m for which $\nu(m)$ is defined, it follows that $\overline{b_n}(m) = \overline{v_2}(\nu(m))$ whenever $\nu(m)$ is defined. \square

7) *Proofs about Connected Rules:* We now prove that, for strongly connected rules, the grounding of a rule's head uniquely determines the grounding of the rule's body (Proposition B.8), whereas for weakly connected rules the grounding of one of the atoms in the body determines the head's grounding (Proposition B.9).

Proposition B.8. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , and U be a Σ -uniqueness template. If p complies with \mathcal{U} and r is strongly connected for \mathcal{U} , then for all $r_1, r_2 \in \text{ground}(p, r)$, if $\text{head}(r_1) = \text{head}(r_2)$, then $r_1 = r_2$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , and U be a Σ -uniqueness template. Furthermore, we assume that (1) p complies with \mathcal{U} , and (2) there are join trees J_1, \dots, J_n such that (a) the trees cover all literals in $\text{body}(r)$, and (b) for each $1 \leq i \leq n$, J_i is strongly connected for \mathcal{U} . Finally, let $r_1, r_2 \in \text{ground}(p, r)$ be two ground rules such that $\text{head}(r_1) = \text{head}(r_2)$. Assume, for contradiction's sake, that $r_1 \neq r_2$. Therefore, there is a position $1 \leq i \leq |\text{body}(r)|$ such that $\text{body}(r_1, i) \neq \text{body}(r_2, i)$. Let $J = (N, E, \text{root}, \lambda)$ be one of the join trees that cover $l = \text{body}(r, i)$. We claim that the grounding of $\text{head}(r)$ determines the grounding of all literals in J . From this, it follows that $\text{body}(r_1, i) = \text{body}(r_2, i)$ leading to a contradiction.

We now prove our claim that the grounding of $\text{head}(r)$ determines the grounding of all literals in J . Let $V(J, i)$ be the set of all the nodes in J at distance at most i from the root root . Furthermore, we denote by $\text{ground}(r, r', J, i)$ the set $\{\text{body}(r', j) \mid \text{body}(r, j) \in V(J, i)\}$. We prove, by induction on i , that for all i and all ground rules r_1 and r_2 instances of r , if $\text{head}(r_1) = \text{head}(r_2)$, then $\text{ground}(r, r_1, J, i) = \text{ground}(r, r_2, J, i)$. From this, it follows that the grounding of $\text{head}(r)$ determines the grounding of all literals in J .

Base case. For $i = 0$, there is a j such that $V(J, 0) = \{\text{body}(r, j)\}$. From this, it follows that $\text{ground}(r, r_1, J, 0) = \{\text{body}(r_1, j)\}$ and $\text{ground}(r, r_2, J, 0) = \{\text{body}(r_2, j)\}$. Furthermore, $\text{anc}(J, \text{body}(r, j)) = \emptyset$. There are two cases depending on whether $\text{body}(r, j)$ is a positive literal or not:

- 1) If $\text{body}(r, j)$ is a positive literal of the form $a(\overline{x})$, then $\text{body}(r_1, j) = a(\overline{c_1})$ and $\text{body}(r_2, j) = a(\overline{c_2})$. From the fact that J is strongly connected for \mathcal{U} , it

follows that there is a set of variables $K \subseteq \{i \mid \overline{x}(i) \in \text{support}(a(\overline{x}))\}$ such that $UNQ(a, K) \in \mathcal{U}$. From $\text{support}(a(\overline{x})) = \text{vars}(\text{head}(r)) \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\}$ and $\text{head}(r_1) = \text{head}(r_2)$, it follows that the values assigned to the variables associated to the indexes in K are the same in r_1 and r_2 . From this, $UNQ(a, K) \in \mathcal{U}$, p complies with \mathcal{U} , $\{a(\overline{c_1}), a(\overline{c_2})\} \subseteq \text{ground}(p)$, and Proposition B.3, it follows that $\overline{c_1} = \overline{c_2}$ and $\text{ground}(r, r_1, J, 0) = \text{ground}(r, r_2, J, 0)$.

- 2) If $\text{body}(r, j)$ is a negative literal of the form $\neg a(\overline{x})$, then $\text{body}(r_1, j) = \neg a(\overline{c_1})$ and $\text{body}(r_2, j) = \neg a(\overline{c_2})$. From the fact that J is strongly connected for \mathcal{U} , it follows that $\text{vars}(\neg a(\overline{x})) \subseteq \text{support}(\neg a(\overline{x}))$. From this, $\text{support}(\neg a(\overline{x})) = \text{vars}(\text{head}(r)) \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\}$, and $\text{head}(r_1) = \text{head}(r_2)$, it follows that the values of the variables in $\text{support}(\neg a(\overline{x}))$ are the same in r_1 and r_2 . From this, it follows that $\overline{c_1} = \overline{c_2}$ and $\text{ground}(r, r_1, J, 0) = \text{ground}(r, r_2, J, 0)$.

Induction Step. Assume that for all $j < i$ and all ground rules $r_1, r_2 \in \text{ground}(p, r)$, if $\text{head}(r_1) = \text{head}(r_2)$, then $\text{ground}(r, r_1, J, j) = \text{ground}(r, r_2, J, j)$. We now show that $\text{ground}(r, r_1, J, i) = \text{ground}(r, r_2, J, i)$. Assume, for contradiction's sake, that this is not the case, namely $\text{ground}(r, r_1, J, i) \neq \text{ground}(r, r_2, J, i)$. From the definition of $\text{ground}(r, r', J, i)$, it follows that $\text{ground}(r, r_1, J, i) = \text{ground}(r, r_1, J, i-1) \cup \{\text{body}(r_1, j) \mid \text{body}(r, j) \in V(J, i) \setminus V(J, i-1)\}$ and $\text{ground}(r, r_2, J, i) = \text{ground}(r, r_2, J, i-1) \cup \{\text{body}(r_2, j) \mid \text{body}(r, j) \in V(J, i) \setminus V(J, i-1)\}$. From this, the induction's hypothesis, and $\text{ground}(r, r_1, J, i) \neq \text{ground}(r, r_2, J, i)$, it follows that $\{\text{body}(r_1, j) \mid \text{body}(r, j) \in V(J, i) \setminus V(J, i-1)\} \neq \{\text{body}(r_2, j) \mid \text{body}(r, j) \in V(J, i) \setminus V(J, i-1)\}$. Therefore, there is a j such that $\text{body}(r, j) \in V(J, i) \setminus V(J, i-1)$ and $\text{body}(r_1, j) \neq \text{body}(r_2, j)$. There are two cases, depending on whether $\text{body}(r, j)$ is a positive literal:

- 1) If $\text{body}(r, j) = a(\overline{x})$, then $\text{body}(r_1, j) = a(\overline{c_1})$, $\text{body}(r_2, j) = a(\overline{c_2})$, and $\overline{c_1} \neq \overline{c_2}$. From the fact that J is strongly connected for \mathcal{U} , it follows that there is a set of variables $K \subseteq \{i \mid \overline{x}(i) \in \text{support}(a(\overline{x}))\}$ such that $UNQ(a, K) \in \mathcal{U}$. From $\text{support}(a(\overline{x})) = \text{vars}(\text{head}(r)) \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\} \cup \bigcup_{l \in V(J, i-1)} \text{vars}(l)$, $\text{ground}(r, r_1, J, i-1) = \text{ground}(r, r_2, J, i-1)$ (from the induction's hypothesis), and $\text{head}(r_1) = \text{head}(r_2)$, it follows that the values assigned to the variables associated to the indexes in K are the same in r_1 and r_2 . From this, $UNQ(a, K) \in \mathcal{U}$, p complies with \mathcal{U} , $\{a(\overline{c_1}), a(\overline{c_2})\} \subseteq \text{ground}(p)$, and Proposition B.3, it follows that $\overline{c_1} = \overline{c_2}$ leading to a contradiction.
- 2) If $\text{body}(r, j) = \neg a(\overline{x})$, then $\text{body}(r_1, j) = \neg a(\overline{c_1})$, $\text{body}(r_2, j) = \neg a(\overline{c_2})$, and $\overline{c_1} \neq \overline{c_2}$. From the fact that J is strongly connected for \mathcal{U} , it follows that $\text{vars}(\neg a(\overline{x})) \subseteq \text{support}(\neg a(\overline{x}))$. From $\text{support}(\neg a(\overline{x})) = \text{vars}(\text{head}(r)) \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\} \cup \bigcup_{l \in V(J, i-1)} \text{vars}(l)$, $\text{ground}(r, r_1, J, i-1) = \text{ground}(r, r_2, J, i-1)$ (from the induction's hypothesis), it follows that the values assigned to the variables in $\text{support}(\neg a(\overline{x}))$

are the same in r_1 and r_2 . From this, it follows that $\bar{c}_1 = \bar{c}_2$ leading to a contradiction.

Since both cases lead to a contradiction, $\text{ground}(r, r_1, J, i) = \text{ground}(r, r_2, J, i)$. \square

Proposition B.9. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , and U be a Σ -uniqueness template. If p complies with \mathcal{U} and r is weakly connected for \mathcal{U} , then for all $r_1, r_2 \in \text{ground}(p, r)$, if there is an $1 \leq i \leq |\text{body}(r)|$ such that $\text{body}(r_1, i) = \text{body}(r_2, i)$, then $r_1 = r_2$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -PROBLOG program, r be a rule in p , and U be a Σ -uniqueness template. Furthermore, we assume that p complies with \mathcal{U} and that r is weakly connected for \mathcal{U} , namely there is a join tree $J = (N, E, \text{root}, \lambda)$ such that (a) J is weakly connected for \mathcal{U} , (b) $N \subseteq \text{body}^+(r)$, and (c) all literals in $\text{body}(r) \setminus N$ are (r, \mathcal{U}, N) -strictly guarded. Let $r_1, r_2 \in \text{ground}(p, r)$ be two ground rules such that there is an $1 \leq i \leq |\text{body}(r)|$ such that $\text{body}(r_1, i) = \text{body}(r_2, i)$, and let l be the literal $\text{body}(r, i)$. There are two cases:

- 1) $\text{body}(r, i) \in N$. Given a node n in a join tree J , we denote by $\text{adjacent}(n, i)$, where $i \in \mathbb{N}$, the sub-tree obtained by considering only the nodes reachable from n using at most i edges. We claim that for all nodes in $\text{adjacent}(l, j)$, the corresponding ground atoms in r_1 and r_2 are the same. From this and the fact that there is a j such that $\text{adjacent}(l, j) = J$, it follows that for all $i \in \{j \mid \text{body}(r, j) \in N\}$, $\text{body}(r_1, i) = \text{body}(r_2, i)$. From this and the fact that the literals in $\text{body}(r) \setminus N$ are (r, \mathcal{U}, N) -strictly guarded, it follows that for all $i \in \{j \mid \text{body}(r, j) \in \text{body}(r) \setminus N\}$, $\text{body}(r_1, i) = \text{body}(r_2, i)$ (since $\text{vars}(l) \subseteq \bigcup_{l' \in N \cap \text{body}^+(r)} \text{vars}(l') \cup \{x \mid (x = c) \in \text{cstr}(r) \wedge c \in \mathbf{dom}\}$ for any literal l in $\text{body}(r) \setminus N$). Therefore, $\text{body}(r_1) = \text{body}(r_2)$.
- 2) $\text{body}(r, i) \notin N$. From this and the fact that r is weakly connected, it follows that there a j such that $\text{body}(r, j) = a(\bar{x})$, $\text{body}(r, j) \in N$, and a $\text{UNQ}(a, K) \in \mathcal{U}$ such that $\{\bar{x}(i) \mid i \in K\} \subseteq \text{vars}(\text{body}(r, i))$. From this, p complies with \mathcal{U} , $\text{body}(r_1, i) = \text{body}(r_2, i)$, $\{\text{body}(r_1, j), \text{body}(r_2, j)\} \subseteq \text{ground}(p)$, and Proposition B.3, it follows that $\text{body}(r_1, j) = \text{body}(r_2, j)$. We proved above that if $\text{body}(r_1, j) = \text{body}(r_2, j)$ and $\text{body}(r, j) \in N$, then $\text{body}(r_1) = \text{body}(r_2)$. Therefore, $\text{body}(r_1) = \text{body}(r_2)$.

From $\text{vars}(\text{head}(r)) \subseteq \bigcup_{l \in \text{body}^+(r)} \text{vars}(l)$ and $\text{body}(r_1) = \text{body}(r_2)$, it follows that $\text{head}(r_1) = \text{head}(r_2)$. Therefore, $r_1 = r_2$.

We now prove, by induction on j , that for all nodes in $\text{adjacent}(l, j)$, the corresponding ground atoms in r_1 and r_2 are the same.

Base Case. The sub-tree $\text{adjacent}(l, 0)$ contains only the node $l = \text{body}(r, i)$. Since $\text{body}(r_1, i) = \text{body}(r_2, i)$, the claim holds for the base case.

Induction step. Assume now that the claim holds for all $j' < j$. We now prove that the claim holds also for j . The sub-tree

$\text{adjacent}(l, j)$ is obtained by extending $\text{adjacent}(l, j-1)$ with either edges of the form $n_1 \xrightarrow{L_1} n'_1$, where $n'_1 \in \text{adjacent}(l, j-1)$, or $n'_2 \xrightarrow{L_2} n_2$, where $n'_2 \in \text{adjacent}(l, j-1)$. In both cases, from the induction hypothesis, the ground literals corresponding to n'_1 and n'_2 are l'_1 and l'_2 and they are the same in r_1 and r_2 . From the definition of weakly connected join tree, there are variables $K_1 \subseteq L_1$ and $K_2 \subseteq L_2$ such that $\text{UNQ}(\text{pred}(n_1), K_1) \in \mathcal{U}$ and $\text{UNQ}(\text{pred}(n_2), K_2) \in \mathcal{U}$. From this, p complies with \mathcal{U} , the fact that the value of K_1 and K_2 are fixed by l'_1 and l'_2 , and Proposition B.3, it follows that the ground atoms corresponding to n_1 and n_2 are the same in r_1 and r_2 (because the values in K_1 and K_2 determines all values in n_1 and n_2). \square

8) *Acyclicity Proof:* We are now ready to prove our first key result, namely that the ground graph associated to an acyclic PROBLOG program is a forest of *poly-trees*, i.e., its undirected version does not contain simple cycles (which are cycles without repetitions of edges and vertices other than the starting and ending vertices).

Proposition B.10. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, and p be a (Σ, \mathbf{dom}) -acyclic PROBLOG program. The graph $gg(p)$ is a forest of poly-trees, i.e., its undirected version does not contain simple cycles.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, and p be a (Σ, \mathbf{dom}) -acyclic PROBLOG program. We prove that the undirected version of $gg(p)$ is acyclic. Assume, for contradiction's sake, that this is not the case, namely there is a simple cycle $C := n_1 \rightarrow n_2 \rightarrow \dots \rightarrow n_k \rightarrow n_1$ in the undirected version of $gg(p)$. There are two cases: (1) there is a directed simple cycle in $gg(p)$ that directly corresponds to C , or (2) there are n directed paths P_1, \dots, P_n in $gg(p)$ that induce a simple cycle C in the undirected version of $gg(p)$ (and P_1, \dots, P_n does not correspond to any directed simple cycle in $gg(p)$). From the first case, it follows that there is a directed cycle in $\text{graph}(p)$, whereas from the second case, it follows that there is an undirected cycle in $\text{graph}(p)$.

Directed Cycle. Assume that $C := n_1 \rightarrow n_2 \rightarrow \dots \rightarrow n_k \rightarrow n_1$ directly corresponds to a directed cycle in $gg(p)$. Without loss of generality, we assume that n_1 is a node of the form $a(\bar{c})$. Therefore, the cycle has form $a_1(\bar{c}_1) \rightarrow (r_1, s_1, i_1) \rightarrow a_2(\bar{c}_2) \rightarrow \dots \rightarrow a_n(\bar{c}_n) \rightarrow (r_n, s_n, i_n) \rightarrow a_{n+1}(\bar{c}_{n+1})$, where $n = k/2$ and $a_{n+1}(\bar{c}_{n+1}) = a_1(\bar{c}_1)$. From this, it follows that there are rules r_1, \dots, r_n and ground rules s_1, \dots, s_n such that: (1) there is a directed cycle $a_1 \xrightarrow{r_1, i_1} a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_{n-1}, i_{n-1}} a_n \xrightarrow{r_n, i_n} a_1$ in p 's dependency graph $\text{graph}(p)$, and (2) for all $1 \leq i \leq n$, $s_i \in \text{ground}(p, r_i)$ and $\text{head}(s_i) = a_{i+1}(c_{i+1})$. From this, it follows that p 's dependency graph $\text{graph}(p)$ contains a directed cycle $C' = a_1 \xrightarrow{r_1, i_1} a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_{n-1}, i_{n-1}} a_n \xrightarrow{r_n, i_n} a_1$. Note that C' may contain loops (i.e., it is not simple). Since p is acyclic and C' is a directed cycle in $\text{graph}(p)$, it follows that there is an ordering template \mathcal{O} , a disjointness template \mathcal{D} , and a uniqueness template \mathcal{U} such that (1) p complies with \mathcal{O}, \mathcal{D} ,

and \mathcal{U} and (2) there is a directed unsafe structure S that covers C' and is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded. Without loss of generality, we assume that $S = C'$. From S covers C' and S is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded, it follows that there is an ordering annotation $ORD(O) \in \mathcal{O}$ such that there are integers $1 \leq y_1 < y_2 < \dots < y_e = n$, literals $o_1(\bar{x}_1), \dots, o_e(\bar{x}_e)$ (where $o_j \in \mathcal{O}$ and $|\bar{x}_j| = |o_j|$), a non-empty set $K \subseteq \{1, \dots, |pr_1|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |o_1|/2\}$ such that for each $0 \leq k < e$,

(1) $pr_{y_k} \xrightarrow{r_{y_k}, i_{y_k}} \dots \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν -downward

connects to $o_{k+1}(\bar{x}_{k+1})$, and (2) $pr_{y_{k+1}-1} \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν' -upward connects to $o_{k+1}(\bar{x}_{k+1})$, where $\nu'(i) = \nu(x) + |o_1|/2$ for all $1 \leq i \leq |o_1|/2$, and $y_0 = 1$. By applying Proposition B.6 and Proposition B.7 to the

paths $pr_{y_k} \xrightarrow{r_{y_k}, i_{y_k}} \dots \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν -downward connects to $o_{k+1}(\bar{x}_{k+1})$ for each $0 \leq k < e$, it follows that (1) ground atoms $o_1(\bar{b}_1, \bar{b}_2), \dots, o_e(\bar{b}_{|K|}, \bar{b}_{|K|+1})$ are in $ground(p)$, and (2) for all $1 \leq w \leq |pr_1|/2$, both $\bar{b}_1(w) = args(body(s_1, i_1))(\nu(w))$ and $\bar{b}_{|K|+1}(w) = args(head(s_n))(\nu'(w))$ hold. From this, $ORD(O) \in \mathcal{O}$, p complies with \mathcal{O} , and Proposition B.1, it follows that $\bar{b}_1 \neq \bar{b}_{|K|+1}$. From this $\bar{b}_1(w) = args(body(s_1, i_1))(\nu(w))$ and $\bar{b}_{|K|+1}(w) = args(head(s_n))(\nu'(w))$ for all $1 \leq w \leq |pr_1|/2$, it follows that $body(s_1, i_1) \neq head(s_n)$. This contradicts $head(s_n) = a(\bar{c}_1)$ and $body(s_1, i_1) = a(\bar{c}_1)$ (which directly follows from the existence of the cycle C).

Undirected Cycle. Assume that C does not directly correspond to any directed simple cycle in $gg(p)$. From this, it follows that there are n directed paths P_1, \dots, P_n in $gg(p)$ such that P_1, \dots, P_n correspond to the simple cycle C in the undirected version of $gg(p)$. From this, it follows that P_1, \dots, P_n form an undirected cycle in $gg(p)$. Therefore:

- 1) For $1 \leq j \leq n$, there is a directed path $D_j := a_1 \xrightarrow{r_1, i_1} a_2 \xrightarrow{r_2, i_2} \dots \xrightarrow{r_{n_j-1}, i_{n_j-1}} a_{n_j}$ in p 's dependency graph $graph(p)$, where P_j is $a_1(\bar{c}_1) \rightarrow (r_1, s_1, i_1) \rightarrow a_2(\bar{c}_2) \rightarrow \dots \rightarrow a_{n_j-1}(\bar{c}_{n_j-1}) \rightarrow (r_{n_j-1}, s_{n_j-1}, i_{n_j-1}) \rightarrow a_{n_j}(\bar{c}_{n_j})$.
- 2) For all $1 \leq j \leq n$, $1 \leq h < n_j$, $s_h \in ground(p, r_h)$ and $head(s_h) = a_{h+1}(\bar{c}_{h+1})$, where $P_j := a_1(\bar{c}_1) \rightarrow (r_1, s_1, i_1) \rightarrow a_2(\bar{c}_2) \rightarrow \dots \rightarrow a_{n_j-1}(\bar{c}_{n_j-1}) \rightarrow (r_{n_j-1}, s_{n_j-1}, i_{n_j-1}) \rightarrow a_{n_j}(\bar{c}_{n_j})$.
- 3) D_1, \dots, D_n form an undirected cycle in $graph(p)$.
- 4) D_1, \dots, D_n is not a directed simple cycle. Indeed, if D_1, \dots, D_n is a directed simple cycle, then C would contain loops (contradicting our assumption that C is a simple cycle).

Since p is acyclic and D_1, \dots, D_n form an undirected cycle in $graph(p)$ that is not a directed simple cycle, it follows that there is an ordering template \mathcal{O} , a disjointness template \mathcal{D} , and a uniqueness template \mathcal{U} such that (1) p complies with \mathcal{O} , \mathcal{D} , and \mathcal{U} and (2) there is an unsafe structure S that covers D_1, \dots, D_n and is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded. We assume that there are values i, a, b, c such that

- 1) $S = \langle D_a, D_b, D_c, U \rangle$,
- 2) $a = i$, $b = (i + 1)\%n$, and $c = (i + 2)\%n$ (i.e., $P_a, P_b,$

and P_c are adjacent in the cycle), and

- 3) U is the undirected path containing all D_i 's that are different from $D_a, D_b,$ and D_c .

Note that the previous assumption is without loss of generality: we can always pick the directed paths in $gg(p)$ inducing C in such a way that they match the guarded structure S .

Let P_a be $a_1(\bar{a}_1) \rightarrow (r_1, s_1, i_1) \rightarrow a_2(\bar{a}_2) \rightarrow (r_2, s_2, i_2) \rightarrow \dots \rightarrow (r_{n_a-1}, s_{n_a-1}, i_{n_a-1}) \rightarrow a_{n_a}(\bar{a}_{n_a})$, P_b be $b_1(\bar{b}_1) \rightarrow (r'_1, s'_1, i'_1) \rightarrow b_2(\bar{b}_2) \rightarrow (r'_2, s'_2, i'_2) \rightarrow \dots \rightarrow (r'_{n_b-1}, s'_{n_b-1}, i'_{n_b-1}) \rightarrow b_{n_b}(\bar{b}_{n_b})$, and P_c be $c_1(\bar{c}_1) \rightarrow (r''_1, s''_1, i''_1) \rightarrow c_2(\bar{c}_2) \rightarrow (r''_2, s''_2, i''_2) \rightarrow \dots \rightarrow (r''_{n_c-1}, s''_{n_c-1}, i''_{n_c-1}) \rightarrow c_{n_c}(\bar{c}_{n_c})$. From $S = \langle D_a, D_b, D_c, U \rangle$, it follows that $a_1(\bar{a}_1) = b_1(\bar{b}_1)$ and $b_{n_b}(\bar{b}_{n_b}) = c_{n_c}(\bar{c}_{n_c})$. From $S = \langle D_a, D_b, D_c, U \rangle$, S covers D_1, \dots, D_n , and S is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded, there are two cases:

- 1) (D_a, D_b) is $(\mathcal{D}, \mathcal{U})$ -head guarded. Therefore, D_a and D_b are non-empty. There are two cases:
 - a) $D_a = D_b$. From this, it follows that (1) P_a and P_b are directed paths in $gg(p)$, (2) $a_1(\bar{a}_1) = b_1(\bar{b}_1)$, (3) $n_a = n_b$, and (4) D_a and D_b are head-connected. From this and (D_a, D_b) is head-guarded, it follows that the rules in D_a and D_b are weakly connected for \mathcal{U} . We claim that $s_h = s'_h$ for any $1 \leq h \leq n_a$. From this, it follows that $P_a = P_b$. This contradicts the fact that P_1, \dots, P_n induces a simple cycle in the undirected version of $gg(p)$ (because the cycle is not simple). We now prove, by induction on d , that $s_d = s'_d$ and $a_d(\bar{a}_d) = b_d(\bar{b}_d)$.

Base Case. Assume that $d = 1$. From $r_1 = r'_1$, r_1 is weakly connected for \mathcal{U} , $body(s_1, i_1) = body(s'_1, i'_1) = a_1(\bar{c}_1)$, $s_d \in ground(p, r_d)$, $s'_d \in ground(p, r'_d)$, and Proposition B.9, it follows that $s_1 = s'_1$. From this, $a_1(\bar{a}_1) = head(s_1)$, and $b_1(\bar{b}_1) = head(s'_1)$, it follows that $a_1(\bar{a}_1) = b_1(\bar{b}_1)$.

Induction Step. Assume that the claim holds for all $d' < d$. We now prove that $s_d = s'_d$ and $a_d(\bar{a}_d) = b_d(\bar{b}_d)$ hold as well. From the induction's hypothesis, it follows that $head(s_{d-1}) = head(s'_{d-1})$. From this and gg 's definition, it follows that $body(s_d, i_d) = body(s'_d, i'_d)$. From this, $r_d = r'_d$, $i_d = i'_d$, r_d is weakly connected for \mathcal{U} , $s_d \in ground(p, r_d)$, $s'_d \in ground(p, r'_d)$, and Proposition B.9, it follows that $s_d = s'_d$. From this, $a_d(\bar{a}_d) = head(s_d)$, and $b_d(\bar{b}_d) = head(s'_d)$, it follows that $a_d(\bar{a}_d) = b_d(\bar{b}_d)$.

- b) $D_a \neq D_b$. From this and (D_a, D_b) are head-guarded, it follows that there is an annotation $DIS(pr, pr') \in \mathcal{D}$ a set $K \subseteq \{1, \dots, |a|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |pr|\}$ such that D_a ν -downward links to pr and D_a ν -downward links to pr' . From D_a ν -downward links to pr , P_a is a ground instance of D_a , and Proposition B.6, it follows that there is a positive literal $pr(\bar{v})$ in the body of one of the ground rules such that $a_1(\bar{a}_1)(k) = \bar{v}(\nu(K))$ for any $k \in K$. From this and the definition of $ground$, it follows that $pr(\bar{v}) \in ground(p)$. From D_b ν -downward links to pr' , P_b is a

ground instance of D_b , and Proposition B.6, it follows that there is a positive literal $pr'(\bar{v}')$ in the body of one of the ground rules such that $b_1(\bar{b}_1)(k) = \bar{v}'(\nu(K))$ for any $k \in K$. From this and the definition of *ground*, it follows that $pr'(\bar{v}') \in \text{ground}(p)$. Finally, from the fact that P_a and P_b are head-connected, it follows that $a_1(\bar{a}_1) = b_1(\bar{b}_1)$. From this, $a_1(\bar{a}_1)(k) = \bar{v}(\nu(K))$ for any $k \in K$, and $b_1(\bar{b}_1)(k) = \bar{v}'(\nu(K))$ for any $k \in K$, it follows that $\bar{v} = \bar{v}'$. From this, $pr(\bar{v}) \in \text{ground}(p)$ and $pr'(\bar{v}) \in \text{ground}(p)$. This contradicts Proposition B.2, $DIS(pr, pr') \in \mathcal{D}$, and p complies with \mathcal{D} .

2) (D_b, D_c) is $(\mathcal{D}, \mathcal{U})$ -tail guarded. Therefore, D_b and D_c are non-empty. There are two cases:

a) $D_b = D_c$. From this, it follows that (1) P_c and P_b are directed paths in $gg(p)$, (2) $c_{n_c}(\bar{c}_{n_c}) = b_{n_b}(\bar{b}_{n_b})$, (3) $n_b = n_c$, and (4) D_b and D_c are tail-connected. From this and (D_b, D_c) is tail-guarded, it follows that (1) D_c and D_b are tail-connected, and (2) the rules in D_c and D_b are strongly connected for \mathcal{U} . We claim that $s''_h = s'_h$ for any $1 \leq h \leq n_c$. From this, it follows that $P_c = P_b$. This contradicts the fact that P_1, \dots, P_n induces a simple cycle in the undirected version of $gg(p)$ (because the cycle is not simple).

We now prove, by induction on d , that $s''_{n_b-d} = s'_{n_b-d}$ and $\text{head}(s''_{n_b-d}) = \text{head}(s'_{n_b-d})$.

Base Case. Assume $d = 0$. From $r''_{n_b} = r'_{n_b}$, r'_{n_b} is strongly connected for \mathcal{U} , $\text{head}(s''_{n_b}) = \text{head}(s'_{n_b}) = b_{n_b}(\bar{b}_{n_b})$, $s''_{n_b} \in \text{ground}(p, r''_{n_b})$, $s'_{n_b} \in \text{ground}(p, r'_{n_b})$, and Proposition B.8, it follows that $s_1 = s'_1$. From this, $a_1(\bar{a}_1) = \text{head}(s_1)$, and $b_1(\bar{b}_1) = \text{head}(s'_1)$, it follows that $b_{n_b}(\bar{b}_{n_b}) = c_{n_b}(\bar{c}_{n_b})$.

Induction Step. Assume that the claim holds for all $d' < d$. We now prove that $s''_{n_b-d} = s'_{n_b-d}$ and $c_{n_b-d}(\bar{c}_{n_b-d}) = b_{n_b-d}(\bar{b}_{n_b-d})$ hold as well. From the induction's hypothesis, it follows that $s''_{n_b-d+1} = s'_{n_b-d+1}$. From this and gg 's definition, it follows that $\text{head}(s''_{n_b-d}, i''_{n_b-d}) = \text{head}(s'_{n_b-d}, i'_{n_b-d})$. From this, $r''_{n_b-d} = r'_{n_b-d}$, $i''_{n_b-d} = i'_{n_b-d}$, r''_{n_b-d} is strongly connected for \mathcal{U} , $s''_{n_b-d} \in \text{ground}(p, r''_{n_b-d})$, $s'_{n_b-d} \in \text{ground}(p, r'_{n_b-d})$, and Proposition B.8, it follows that $s_{n_b-d} = s'_{n_b-d}$. From this, $c_{n_b-d}(\bar{c}_{n_b-d}) = \text{head}(s''_{n_b-d})$, and $b_{n_b-d}(\bar{b}_{n_b-d}) = \text{head}(s'_{n_b-d})$, it follows that $c_{n_b-d}(\bar{c}_{n_b-d}) = b_{n_b-d}(\bar{b}_{n_b-d})$.

b) $D_b \neq D_c$. From this and (D_b, D_c) is tail-guarded, it follows that there is an annotation $DIS(pr, pr') \in \mathcal{D}$, a set $K \subseteq \{1, \dots, |a|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |pr|\}$, such that D_b ν -upward links to pr and D_c ν -upward links to pr' . From D_b ν -upward links to pr , P_b is a path in the ground graph corresponding to D_b , and Proposition B.7, it follows that a positive literal $pr(\bar{v})$ in the body of one of the ground rules such that $b_{n_b}(\bar{b}_{n_b})(k) = \bar{v}(\nu(K))$ for any $k \in K$. From this and the definition of *ground*, it follows that $pr(\bar{v}) \in \text{ground}(p)$. From D_c ν -upward links to

pr' , P_c is a path in the ground graph corresponding to D_c , and Proposition B.7, it follows that a positive literal $pr'(\bar{v}')$ in the body of one of the ground rules such that $c_{n_c}(\bar{c}_{n_c})(k) = \bar{v}'(\nu(K))$ for any $k \in K$. From this and the definition of *ground*, it follows that $pr'(\bar{v}') \in \text{ground}(p)$. Finally, from the fact that P_b and P_c are tail-connected, it follows that $b_{n_b}(\bar{b}_{n_b}) = c_{n_c}(\bar{c}_{n_c})$. From this, $b_{n_b}(\bar{b}_{n_b})(k) = \bar{v}(\nu(K))$ for any $k \in K$, and $c_{n_c}(\bar{c}_{n_c})(k) = \bar{v}'(\nu(K))$ for any $k \in K$, it follows that $\bar{v} = \bar{v}'$. Therefore, $pr(\bar{v}) \in \text{ground}(p)$ and $pr'(\bar{v}) \in \text{ground}(p)$. This contradicts Proposition B.2, $DIS(pr, pr') \in \mathcal{D}$, and p complies with \mathcal{D} .

This completes the proof of our claim. \square

9) *Auxiliary Results:* Here we prove two auxiliary results that help in establishing that programs are acyclic. In particular, Proposition B.11 states pre-conditions that allows reducing the guardedness of a complex undirected structure to the guardedness of simpler structures. Similarly, Proposition B.12 states pre-conditions that allows reducing the guardedness of a complex directed structure to the guardedness of a sequence of simpler structures. Note that Proposition B.12 can be easily extended to support (1) different forms of cycle combination, and (2) combinations of non-self-loop cycles.

Proposition B.11. *Let p be a PROBLOG program, \mathcal{O} be an ordering template, \mathcal{D} be a disjointness template, \mathcal{U} be a uniqueness template, C be an undirected cycle in $graph(p)$, and $S = \langle D_1, D_2, D_3, U \rangle$ be an undirected unsafe structure. If (1) p complies with \mathcal{O} , \mathcal{D} , and \mathcal{U} , (2) C is equivalent to $D_1, U_1, U, U_2, D_3, D_2$, (3) S is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded, then there is an undirected unsafe structure that covers C and is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded.*

Proof. Let p be a PROBLOG program, \mathcal{O} be an ordering template, \mathcal{D} be a disjointness template, \mathcal{U} be a uniqueness template, C be an undirected cycle in $graph(p)$, and $S = \langle D_1, D_2, D_3, U \rangle$ be an undirected structure. We assume that (1) p complies with \mathcal{O} , \mathcal{D} , and \mathcal{U} , (2) C is equivalent to $D_1, U_1, U, U_2, D_3, D_2$, (3) S is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded. We define the undirected unsafe structure $S' = \langle D_1, D_2, D_3, U_1, U, U_2 \rangle$. From our assumption C is equivalent to $D_1, U_1, U, U_2, D_3, D_2$. Thus, S' covers C . Furthermore, since S and S' agree on all directed paths, which are the only ones that determine whether an undirected structure is guarded, it follows that S' is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -guarded \square

In the following, we say that a directed cycle is guarded by a set of predicates O and a bijection ν iff the definition of guarded cycle is satisfied for the specific bijection ν .

Proposition B.12. *Let p be a PROBLOG program, $C_1 = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n \xrightarrow{r_n, i_n} pr_1$ be a directed cycle in $graph(p)$, O be a set of predicate symbols, $pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$, and o be a predicate symbol such that for all $o' \in O$, $|o| = |o'|$, K be a non-empty set $K \subseteq \{1, \dots, |pr_1|\}$, and a bijection $\nu : K \rightarrow \{1, \dots, |o|/2\}$. If (1) C_1 is guarded for O_1 and*

ν , and (2) $pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ is guarded for o and ν , then $C_1, pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ is guarded for $O_1 \cup \{o\}$.

Proof. Let p be a PROBLOG program, $C_1 = pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n \xrightarrow{r_n, i_n} pr_1$ be a directed cycle in $graph(p)$, O be a set of predicate symbols, $pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$, and o be a predicate symbol such that for all $o' \in O$, $|o| = |o'|$, K be a non-empty set $K \subseteq \{1, \dots, |pr_1|\}$, ν be a bijection $\nu : K \rightarrow \{1, \dots, |o|/2\}$ and ν' be the bijection $\nu'(i) = \nu(x) + |o|/2$ for all $1 \leq i \leq |o|/2$. Furthermore, we assume that (1) C_1 is guarded for O_1 and ν , and (2) $pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ is guarded for o and ν .

First, we rewrite $C_1, pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ as $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n \xrightarrow{r_n, i_n} pr_{n+1} \xrightarrow{r_{n+1}, i_{n+1}} pr_1$. From C_1 is guarded for O_1 , it follows that there are integers $1 \leq y_1 < y_2 < \dots < y_e \leq n$ such that $y_e = n$ and literals $o_1(\bar{x}_1), \dots, o_e(\bar{x}_e)$ (where $o_j \in O$ and $|\bar{x}_j| = |o_j|$), such that for each $0 \leq k < e$, (1) $pr_{y_k} \xrightarrow{r_{y_k}, i_{y_k}} \dots \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν -downward connects to $o_{k+1}(\bar{x}_{k+1})$, and (2) $pr_{y_{k+1}-1} \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν' -upward connects to $o_{k+1}(\bar{x}_{k+1})$, where $y_0 = 1$. Furthermore, from $pr_1 \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ is guarded for o and ν , it follows that there is a literal $o(\bar{x})$ such that (3) $pr_{n+1} \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ ν -downward connects to $o(\bar{x})$, and (4) $pr_{y_{n+1}} \xrightarrow{r_{n+1}, i_{n+1}} pr_{y_1}$ ν' -upward connects to $o_{k+1}(\bar{x}_{k+1})$. From (1)–(4), it therefore follows that $pr_1 \xrightarrow{r_1, i_1} \dots \xrightarrow{r_{n-1}, i_{n-1}} pr_n \xrightarrow{r_n, i_n} pr_{n+1} \xrightarrow{r_{n+1}, i_{n+1}} pr_1$ is guarded for ν and $O \cup \{o\}$. Indeed, there are integers $1 \leq y_1 < y_2 < \dots < y_e < y_{e+1} \leq n+1$ such that $y_{e+1} = n+1$ and literals $o_1(\bar{x}_1), \dots, o_e(\bar{x}_e), o_{e+1}(\bar{x}_{e+1})$ (where $o_j \in O \cup \{o\}$ and $|\bar{x}_j| = |o_j|$) such that for each $0 \leq k < e+1$, (1) $pr_{y_k} \xrightarrow{r_{y_k}, i_{y_k}} \dots \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν -downward connects to $o_{k+1}(\bar{x}_{k+1})$, and (2) $pr_{y_{k+1}-1} \xrightarrow{r_{y_{k+1}-1}, i_{y_{k+1}-1}} pr_{y_{k+1}}$ ν' -upward connects to $o_{k+1}(\bar{x}_{k+1})$. \square

B. Relaxed Acyclic Programs

Acyclic PROBLOG programs as defined above are very restrictive. For instance, they cannot encode rules like those in Section V. We now design *relaxed acyclic programs*, a larger fragment of PROBLOG where these examples can be encoded.

The key components of this new fragment are two syntactic transformations over PROBLOG programs that, while not preserving the program's semantics, preserve certain key aspects of the program's structure. Intuitively, a PROBLOG program p is a *relaxed acyclic program* if the program obtained from p by applying the transformations is an acyclic PROBLOG program.

Finally, we develop a procedure for compiling any relaxed acyclic PROBLOG program into a poly-tree Bayesian Network. Since any acyclic program is a relaxed acyclic program as well, this procedure can be applied also to acyclic programs.

1) *Rule Domination*: Rule r_1 is *dominated* by rule r_2 , written $r_1 \sqsubseteq r_2$, iff:

- $head(r_1) = head(r_2)$,
- $cstr(r_1) = cstr(r_2)$,
- for all $1 \leq i \leq |body(r_1)|$, then $pred(body(r_1, i)) = pred(body(r_2, i))$ and $args(body(r_1, i)) = args(body(r_2, i))$.

We extend the domination relation also to atoms and probabilistic atoms as follows: $a_1(\bar{c}_1) \sqsubseteq a_2(\bar{c}_1)$ iff $a_1 = a_2$ and $\bar{c}_1 = \bar{c}_2$, and $v_1::a_1(\bar{c}_1) \sqsubseteq v_2::a_2(\bar{c}_1)$ iff $v_1 = v_2$, $a_1 = a_2$, and $\bar{c}_1 = \bar{c}_2$.

Given a program p and a rule $r \in p$, $[r]_{\sqsubseteq p}$ denotes the set $\{r' \in p \mid r' \sqsubseteq r\}$. We say that a rule r is *maximal* in a program p iff there is no rule $r' \in p$ such that $r \sqsubseteq r'$. The *kernel* of $[r]_{\sqsubseteq p}$, denoted $k([r]_{\sqsubseteq p})$, is the rule r' defined as follows: $head(r') = head(r)$, $|body(r')| = |body(r)|$, $cstr(r') = cstr(r)$, and for all $1 \leq i \leq |body(r')|$, then $body(r', i) = a_i(\bar{c}_i)$ if for all rules $r'' \in [r]_{\sqsubseteq p}$ (1) $|body(r'')| \geq i$ and (2) $body(r'', i)$ is a positive literal and $body(r', i) = -a_i(\bar{c}_i)$ otherwise, where $a_i = pred(body(r, i))$ and $\bar{c}_i = args(body(r, i))$.

The *maximal projection* of p , denoted $p \downarrow_{\sqsubseteq}$, is $\{k([r]_{\sqsubseteq p}) \mid \neg \exists r' \in p. r \sqsubseteq r'\}$. A *maximal partition* PP of p is a subset of $\mathbb{P}(p)$ such that: (1) for each $R \in PP$, $R = [r]_{\sqsubseteq p}$ for some maximal rule $r \in p$, and (2) for each $R_1, R_2 \in PP$, $R_1 \cap R_2 = \emptyset$. A maximal partition PP induces a *maximal assignment* $\nu_{PP} : p \rightarrow PP$ such that $\nu_{PP}(r) = R$, where $R \in PP$ and $r \in R$ if r is maximal, and $\nu_{PP}(r) = \emptyset$ otherwise.

The syntactic transformation α , which removes all non-maximal rules, takes as input a PROBLOG program p and returns as its maximal projection $p \downarrow_{\sqsubseteq}$.

2) *CPT-like predicates*: We now present CPT-like predicates, a special kind of probabilistic structure that can be formalized using annotated disjunctions. While general annotated disjunctions cannot be used in our encoding, since they introduce, in general, cycles in the underlying Bayesian Network, CPT-like predicates can be still encoded as poly-trees.

Let Σ be a first-order signature, \mathbf{dom} be a finite domain, pr be a predicate symbol in Σ of arity $|pr|$, p be a (Σ, \mathbf{dom}) -PROBLOG program, and K be a set of distinct integer values such that for all $i \in K$, $1 \leq i \leq |pr|$. Given a tuple \bar{t} , we denote by $\bar{t} \downarrow_{\{i_1, \dots, i_n\}}$ the tuple $(\bar{t}(i_1), \dots, \bar{t}(i_n))$.

We say that pr is *K-fixed domain with respect to p* iff for all rules $r \in p$, the following conditions hold:

- if $pred(head(r)) = pr$, then $args(head(r)) \downarrow_K$ is a tuple in $\mathbf{dom}^{|K|}$, and
- if $pred(body(r, j)) = pr$, for some $1 \leq j \leq |body(r)|$, then $args(body(r, j)) \downarrow_K$ is a tuple in $\mathbf{dom}^{|K|}$.

We denote by $pdom(pr, p, K)$ the set of $|K|$ -tuples representing all constant values associated with the positions in K . Namely, $pdom(pr, p, K) = \{args(head(r)) \downarrow_K \mid r \in p\} \cup \{args(body(r, j)) \downarrow_K \mid r \in p \wedge 1 \leq j \leq |body(r)|\}$. Finally, given an atom $a(\bar{x})$ and a K , we denote by $\beta(a(\bar{x}), K)$ the atom $a(\bar{x} \downarrow_K)$.

A set of rules r_1, \dots, r_n in p is an *annotated disjunction-set* iff they are the translation of an annotated disjunction to plain PROBLOG. Formally, a set of rules r_1, \dots, r_n in p

is an *annotated disjunction-set* iff there are a set of predicate symbols sw_1, \dots, sw_n and a sequence of literals L such that (1) for $1 \leq i \leq n$, $body(r_i) = L, \neg sw_1(\bar{x}_1), \dots, \neg sw_{i-1}(\bar{x}_{i-1}), sw_i(\bar{x}_i)$, where $\bar{x}_j = args(head(r_j))$ for $1 \leq j \leq n$, and (2) for $1 \leq i \leq n$, sw_i is used only to define probabilistic ground atoms such that if $v_1::sw_i(\bar{c}_1) \in p$ and $v_2::sw_i(\bar{c}_2) \in p$, then $v_1 = v_2$, and (3) $\sum_{1 \leq i \leq n} p(i) \leq 1$, where $p(i) = v_i \cdot (1 - \sum_{1 \leq j < i} p(j))$, where v_i is the probability associated to the predicate symbol sw_i in the program p . We denote by (1) $common(R)$, where R is an annotated disjunction-set, the list of literals that are common to all the rules in R , namely $common(R) = \bigcap_{r \in R} body(r)$ (with a slight abuse of notation, we use intersection for lists), by (2) $heads(R)$ the set containing the various heads of the rules in R , namely $heads(R) = \bigcup_{r \in R} \{head(r)\}$, and by (3) $common(R, pr)$ the list of all literals in $common(R)$ whose predicate symbol is pr .

We say that an annotated disjunction set R is an *pr-annotated disjunction-set* iff for all rules $r \in R$, $pred(head(r)) = pr$. Let R_1 and R_2 be two pr -annotated disjunction-sets. We say that R_1 and R_2 are (pr_1, pr_2) -disjoint iff $pr_1(\bar{x}_1 \downarrow_{\{1, \dots, |pr_1|\} \setminus K}) \in common(R_1)$ and $pr_2(\bar{x}_2 \downarrow_{\{1, \dots, |pr_2|\} \setminus K}) \in common(R_2)$, $pr(\bar{x}_1) \in heads(R_1)$, and $pr(\bar{x}_2) \in heads(R_2)$. We say that R_1 and R_2 are (pr, K) -row-distinct iff (1) all literals in $common(R_1, pr)$ are positive, (2) all literals in $common(R_2, pr)$ are positive, (3) $|common(R_1, pr)| = |common(R_2, pr)|$, and (4) $\{\bar{x}_1 \downarrow_K \mid \exists l \in common(R_1, pr). \bar{x}_1 = args(l)\} \neq \{\bar{x}_2 \downarrow_K \mid \exists l \in common(R_2, pr). \bar{x}_2 = args(l)\}$.

We say that pr is $(K, \mathcal{D}, \mathcal{U})$ -CPT-like in a program p iff

- for all rules $r \in p$ such that $pred(head(r)) = pr$, (1) r is strongly connected for \mathcal{U} , and (2) $body(r) \neq \emptyset$,
- pr is K -fixed domain with respect to p ,
- for all maximal pr -annotated disjunction-sets $R \subseteq p$, for all $pr(\bar{x}) \in heads(R)$, $\bar{x} \downarrow_K \in pdom(pr, p, K)$ for all $1 \leq i \leq n$,
- for all maximal pr -annotated disjunction-sets $R \subseteq p$, $|heads(R)| = |pdom(pr, p, K)|$,
- for all maximal pr -annotated disjunction-sets $R \subseteq p$, for all distinct $pr(\bar{x}_1), pr(\bar{x}_2) \in heads(R)$, $\bar{x}_1 \downarrow_K \neq \bar{x}_2 \downarrow_K$ and $\bar{x}_1 \downarrow_{\{1, \dots, |pr|\} \setminus K} = \bar{x}_2 \downarrow_{\{1, \dots, |pr|\} \setminus K}$,
- for any two distinct and maximal pr -annotated disjunction-sets $R_1, R_2 \subseteq p$, one of the following conditions hold:
 - there is pair of predicate symbols $(pr_1, pr_2) \in \mathcal{D}$ such that R_1 and R_2 are (pr_1, pr_2) -disjoint, or
 - R_1 and R_2 are (pr, K) -row-distinct.
- There is no rule $r \in p$ such that (1) $pred(head(r)) = pr$, and (2) r is not in a maximal pr -annotated disjunction set.

Let pr be a predicate symbol and K be a set of natural numbers. The syntactic transformation $\Downarrow_{pr, K}$ is defined as follows:

- $(pr(\bar{x})) \Downarrow_{pr, K}$ is $pr(\bar{x} \downarrow_{\{1, \dots, |\bar{x}|\} \setminus K})$,
- $(pr'(\bar{x})) \Downarrow_{pr, K}$ is $pr'(\bar{x})$, where $pr \neq pr'$,
- $(v::a(\bar{x})) \Downarrow_{pr, K}$ is $v::(a(\bar{x}) \Downarrow_{pr, K})$,

- $(\neg a(\bar{x})) \Downarrow_{pr, K}$ is $\neg(a(\bar{x}) \Downarrow_{pr, K})$, and
- $(h \leftarrow l_1, \dots, l_n, c_1, \dots, c_n) \Downarrow_{pr, K}$ is $h \Downarrow_{pr, K} \leftarrow l_1 \Downarrow_{pr, K}, \dots, l_n \Downarrow_{pr, K}, c_1, \dots, c_n$, and
- $p \Downarrow_{pr, K}$, where p is a program, is $\{r \Downarrow_{pr, K} \mid r \in p\}$.

Finally, we define the transformation $\Downarrow_{(pr_1, K_1), \dots, (pr_n, K_n)}$ that is just $\Downarrow_{pr_1, K_1} \circ \dots \circ \Downarrow_{pr_n, K_n}$.

Given a program p , the function $\mu_{p, \mathcal{D}, \mathcal{U}}$ associates to each predicate symbol pr in Σ the maximal set K such that pr is $(K, \mathcal{D}, \mathcal{U})$ -CPT-like with respect to p . The syntactic transformation $\beta_{p, \mathcal{D}, \mathcal{U}}$ takes as input a PROBLOG program p' and it returns as output the program $p' \Downarrow_{(pr_1, \mu_{p, \mathcal{D}, \mathcal{U}}(pr_1)), \dots, (pr_n, \mu_{p, \mathcal{D}, \mathcal{U}}(pr_n))}$, where pr_1, \dots, pr_n are all predicate symbols in Σ .

Note that p' is defined on a different vocabulary than p . The program p' is defined over the “reduced vocabulary” obtained by applying β to the atoms in p .

3) *Relaxed Acyclic PROBLOG programs*: A *relaxed acyclic* (Σ, \mathbf{dom}) -PROBLOG program is a (Σ, \mathbf{dom}) -PROBLOG program p such that there are a Σ -ordering template \mathcal{O} , a Σ -disjointness template \mathcal{D} , and a Σ -uniqueness template \mathcal{U} such that $\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p))$ is $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -acyclic. An *acyclicity witness for a program* p is a tuple $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ such that $\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p))$ is a $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ -acyclic PROBLOG program.

Given a relaxed acyclic (Σ, \mathbf{dom}) -PROBLOG program p , we associate to each rule r in $\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p))$, the set $[r]_p$ of rules in p defined as follows:

$$[r]_p := \bigcup_{r' \in [r] \square_{\beta_{p, \mathcal{D}, \mathcal{U}}(p)}} \{r'' \in p \mid \beta_{p, \mathcal{D}, \mathcal{U}}(r'') = r'\}$$

The set $[r]_p$ contains all rules in p that are represented by the rule r in $\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p))$.

4) *Compilation to Bayesian Networks*: Given a relaxed acyclic PROBLOG program p , its encoding as a Bayesian Network $BN = (N, E, cpt)$, denoted $bn(p)$, is defined by Algorithm 2.

5) *Encoding’s acyclicity*: Here, we prove that Algorithm 2 produces a Bayesian Network that is a forest of poly-trees. Note that Algorithm 2 produces a forest of poly-trees and not just a single poly-tree because some predicates symbols may be independent. For instance, the program consisting of the rules $a(x) \leftarrow b(x)$ and $c(x) \leftarrow d(x)$ corresponds to two poly-trees (one per rule).

Proposition B.13. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program, and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p ’s acyclicity. The Bayesian Network (N, E, cpt) produced by Algorithm 2 on input p is a forest of poly-trees.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program, and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p ’s acyclicity. Furthermore, let $BN = (N, E, cpt)$ be the Bayesian Network produced by Algorithm 2 on input p . There is a one-to-one mapping from paths in the ground graph $gg(\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p)))$ and paths in the Bayesian Network produced by Algorithm 2. Namely, there is a path from $a(\bar{c})$ to $a'(\bar{c}')$ in $gg(\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p)))$ iff

Algorithm 2: Constructing the Bayesian Network. The sub-routines *tree*, *CPT*, and CPT_{\oplus} are shown in Algorithm 3.

Input: A relaxed acyclic (Σ, dom) -PROBLOG program p_0 and a witness $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ for p .

Output: A Bayesian Network (N, E, cpt) .

▷ Transform the original program.
 $p = \alpha(\beta_{p_0, \mathcal{D}, \mathcal{U}}(p_0))$
 ▷ Initialize the domain of each predicate symbol.
 $D = \emptyset$
for $pr \in \Sigma$ **do**
 if $|\mu_{p_0, \mathcal{D}, \mathcal{U}}(pr)| > 0$ **then**
 $D(pr) = \text{pdom}(pr, p_0, \mu_{p_0, \mathcal{D}, \mathcal{U}}(pr)) \cup \{\perp\}$
 else
 $D(pr) = \{\top, \perp\}$
 ▷ Initialize the BN.
 $N = \emptyset$
 $E = \emptyset$
 $\text{cpt} = \emptyset$
 ▷ Creates the nodes.
for $r \in p$ **do**
 for $r' \in \text{ground}(p, r)$ **do**
 $N = N \cup \{X[\text{head}(r')]\}$
 for $a(\bar{c}) \in \text{body}^+(r')$ **do**
 $N = N \cup \{X[a(\bar{c})]\}$
 for $\neg a(\bar{c}) \in \text{body}^-(r')$ **do**
 $N = N \cup \{X[a(\bar{c})]\}$
 for $r \in p$ **do**
 $\mu = \emptyset$
 for $s \in \text{ground}(p, r)$ **do**
 $\mu(\text{head}(s)) = \mu(\text{head}(s)) \cup \{\text{pos}(s)\}$
 for $a(\bar{c}) \in \text{domain}(\mu)$ **do**
 $N = N \cup \{X[r, a(\bar{c})]\}$
 $E = E \cup \{X[r, a(\bar{c})] \rightarrow X[a(\bar{c})]\}$
 $K = \emptyset$
 for $I \in \mu(a(\bar{c}))$ **do**
 $K = K \cup \{X[r, I, a(\bar{c})]\}$
 if $I = \emptyset$ **then**
 ▷ Here, $[r]_{p_0} = \{r\}$ and $D(a) = \{\top, \perp\}$.
 if $\exists v. \text{head}(r) = v::a(\bar{c})$ **then**
 $p = v$
 else
 $p = 1$
 $\text{cpt}(X[r, \emptyset, a(\bar{c})]) = (\top \mapsto p,$
 $\perp \mapsto (1 - p))$
 else
 for $b(\bar{v}) \in I$ **do**
 $E = E \cup \{X[b(\bar{v})] \rightarrow X[r, I, a(\bar{c})]\}$
 for $\neg b(\bar{v}) \in I$ **do**
 $E = E \cup \{X[b(\bar{v})] \rightarrow X[r, I, a(\bar{c})]\}$
 $\text{cpt}(X[r, I, a(\bar{c})]) = \text{CPT}(I, r, D, \mu_{p_0, \mathcal{D}, \mathcal{U}}, p_0)$
 $(N', E', \text{cpt}') = \text{tree}(K, X[r, a(\bar{c})], D(a))$
 $N = N \cup N'$
 $E = E \cup E'$
 $\text{cpt} = \text{cpt} \cup \text{cpt}'$
 ▷ Set the CPT for the variables associated to the atoms.
 for $X[a(\bar{c})] \in N$ **do**
 $\text{cpt}(X[a(\bar{c})]) = \text{CPT}_{\oplus}(|\{X[r, a(\bar{c})] \in N\}| + 1, D(a))$
 return (N, E, cpt)

there is a path from $X[a(\bar{c})]$ to $X[a'(\bar{c}')$ in BN . Assume that there is a cycle in the undirected version of BN . From this, it follows that there is an undirected cycle in $gg(\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p)))$. This, however, contradicts Proposition B.10 (since p is relaxed acyclic, then $\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p))$ is acyclic and there are no undirected cycles in $gg(\alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p)))$). \square

6) *Bayesian Networks:* A Bayesian Network BN is a tuple (N, E, cpt) , where N is the set of nodes, E is the set of edges, cpt is a function associating to each node $n \in N$ its Conditional Probability Table. For each node $n \in N$, we denote by $D(n)$ its domain, i.e., the possible values it can have. Note that $D(n)$ can be immediately derived from $\text{cpt}(n)$. We denote by $p(n)$ the set of n 's parents, i.e., $p(n) = \{n' \mid n' \rightarrow n \in E\}$. Furthermore, we denote by $\text{ancestors}(n)$ the set of n 's ancestors, i.e., $\text{ancestors}(n) = p(n) \cup \bigcup_{n' \in p(n)} \text{ancestors}(n')$.

A *BN-total assignment* is a total function ν that associates to each $n \in N$ a value in $n \in D(n)$, whereas a *BN-partial assignment* is a partial function ν that associates to each $n \in N'$, where $N' \subseteq N$, a value in $n \in D(n)$. The probability defined by BN given ν , written $\llbracket BN \rrbracket(\nu)$ is $[\sum_{y_1} \dots \sum_{y_m} (\prod_{n \in N} \text{cpt}(n)) (Y_1 = y_1, \dots, Y_m = y_m)] (X_1 = v_1, \dots, X_n = v_n)$, where $\text{dom}(\nu) = \{X_1, \dots, X_n\}$, $N \setminus \text{dom}(\nu) = \{Y_1, \dots, Y_m\}$, and $v_i = \nu(X_i)$ for $1 \leq i \leq n$.

7) *Correctness Proof:* Before presenting our correctness proof, we introduce some machinery. Let p be a relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , and $BN = (N, E, \text{cpt})$ be the Bayesian Network produced by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as inputs. For each rule atom $a(\bar{c}) \in \text{ground}(p)$, we denote by $a(\bar{c}) \downarrow_p$ the atom obtained by applying the transformations α and β , whereas we denote by $a(\bar{c}) \uparrow_p$ the constants not used in the transformed program, namely $a(\bar{c}) \uparrow_p$ is $\bar{c} \downarrow_{\{1, \dots, |\bar{c}|\} \setminus \mu_{p, \mathcal{D}, \mathcal{U}}(a)}$. Furthermore, given a rule $r \in p$ (respectively a ground rule $s \in \text{ground}(p, r)$), we denote by $r \downarrow_p$ (respectively $s \downarrow_p$) the rule obtained by applying the transformations α and β . If $\mu_{p, \mathcal{D}, \mathcal{U}}(a) = \emptyset$, then $a(\bar{c}) \uparrow_p = \top$. We say that the *kernel of BN*, denoted $K(BN)$, is the set of variables $\{X[r, \emptyset, a(\bar{c})] \in N \mid \exists v. r = v::a(\bar{c})\}$. Note that all nodes $n \in K(BN)$ are boolean random variables by construction, namely $D(n) = \{\top, \perp\}$. Given a BN -total assignment ν , we say that ν is *consistent* iff for all variables $n \in N \setminus K(BN)$, $\text{cpt}(n)(\nu(P_1), \dots, \nu(P_m), \nu(n)) = 1$, where $p(n) = \{P_1, \dots, P_m\}$. Furthermore, we say that ν is a *model for a state s*, written $\nu \models s$, iff $a(\bar{c}) \in s \Leftrightarrow \nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$.

We first prove a simple results that connect the original program p and its transformed version $\alpha(\beta(p))$.

Proposition B.14. *Let Σ be a first-order signature, dom be*

Algorithm 3: Auxiliary functions used in Algorithm 2.

```

function  $CPT((a_1(\bar{c}_1), \dots, a_n(\bar{c}_n)), r, D, \mu, p_0)$ 
   $A = D(a_1) \times \dots \times D(a_n) \times D(\text{pred}(\text{head}(r)))$ 
   $cpt = \emptyset$ 
  for  $\bar{a} \in A$  do
     $\bar{a}' = \text{removeDuplicates}((a_1(\bar{c}_1), \dots, a_n(\bar{c}_n)), \bar{a})$ 
    if  $\text{satisfiable}(r, \bar{a}, D, \mu, p_0) \wedge \forall 1 \leq i \leq n. \forall 1 \leq j \leq n. (i \neq j \wedge a_i(\bar{c}_i) = a_j(\bar{c}_j) \Rightarrow \bar{a}(i) = \bar{a}(j))$  then
       $cpt = cpt \cup \{\bar{a}' \mapsto 1\}$ 
    else
       $cpt = cpt \cup \{\bar{a}' \mapsto 0\}$ 
  return  $cpt$ 

```

```

function  $\text{satisfiable}(r, (v_1, \dots, v_n), D, \mu, p_0)$ 
  if  $v_n \neq \perp$  then
     $\triangleright$ At least one satisfiable assignment.
     $res = \perp$ 
    for  $r' \in [r]_{p_0}$  do
      if  $\text{filter}(\text{head}(r'), D, \mu) = v_n$  then
         $sat = \top$ 
        for  $1 \leq i \leq |\text{body}(r')|$  do
          if  $v_i \notin \text{filter}(\text{body}(r', i), D, \mu)$  then
             $sat = \perp$ 
           $res = res \vee sat$ 
        return  $res$ 
  if  $v_n = \perp$  then
     $\triangleright$ All assignments must be unsatisfiable.
     $res = \perp$ 
    for  $r' \in [r]_{p_0}$  do
       $sat = \top$ 
      for  $1 \leq i \leq |\text{body}(r')|$  do
        if  $v_i \notin \text{filter}(\text{body}(r', i), D, \mu)$  then
           $sat = \perp$ 
         $res = res \vee sat$ 
    return  $\neg res$ 

```

```

function  $\text{filter}(l, D, \mu)$ 
  if  $\exists a \in \Sigma, \bar{x} \in (\text{Var} \cup \text{dom})^{|\bar{x}|}. l = a(\bar{x})$  then
    if  $\mu(a) \neq \emptyset$  then
      return  $\{x \downarrow_{\mu(a)}\}$ 
    else
      return  $\{\top\}$ 
  if  $\exists a \in \Sigma, \bar{x} \in (\text{Var} \cup \text{dom})^{|\bar{x}|}. l = \neg a(\bar{x})$  then
    if  $\mu(a) \neq \emptyset$  then
      return  $(D(a) \setminus \{\mu(a)\}) \cup \{\perp\}$ 
    else
      return  $\{\perp\}$ 

```

```

function  $CPT_{\oplus}(n, D)$ 
   $\triangleright$ If  $D \neq \{\top, \perp\}$ , then  $n = 2$  (since the program is relaxed-acyclic).
  if  $D \neq \{\top, \perp\} \wedge n = 2$  then
     $cpt = \emptyset$ 
    for  $(v_1, v_2) \in D^2$  do
      if  $v_1 = v_2$  then
         $cpt = cpt \cup \{(v_1, v_2) \mapsto 1\}$ 
      else
         $cpt = cpt \cup \{(v_1, v_2) \mapsto 0\}$ 
    return  $cpt$ 
  if  $D = \{\top, \perp\}$  then
     $E = \{\top, \perp\}^n$ 
     $cpt = \emptyset$ 
    for  $(v_1, \dots, v_n) \in E$  do
       $K = \{v_i \mid 1 \leq i \leq n-1 \wedge v_i = \top\}$ 
      if  $(v_n = \top \wedge K \neq \emptyset) \vee (v_n = \perp \wedge K = \emptyset)$  then
         $cpt = cpt \cup \{\bar{v} \mapsto 1\}$ 
      else
         $cpt = cpt \cup \{\bar{v} \mapsto 0\}$ 
    return  $cpt$ 

```

```

function  $\text{tree}(\mathbb{X}, \text{root}, D)$ 
   $p = \text{nil}$ 
   $N = \mathbb{X}$ 
   $E = \emptyset$ 
   $cpt = \emptyset$ 
   $\mathbb{X}' = \emptyset$ 
  if  $\mathbb{X} = \{x\}$  then
     $N = N \cup \{\text{root}\}$ 
     $E = E \cup \{x \rightarrow \text{root}\}$ 
     $cpt(\text{root}) = CPT_{\oplus}(2, D)$ 
  else
    for  $v \in \mathbb{X}$  do
      if  $p = \text{nil}$  then
         $p = v$ 
      else
         $\mathbb{X}' = \mathbb{X}' \cup X[p, v]$ 
         $N = N \cup X[p, v]$ 
         $E = E \cup \{p \rightarrow X[p, v], v \rightarrow X[p, v]\}$ 
         $cpt(X[p, v]) = CPT_{\oplus}(3, D)$ 
         $p = \text{nil}$ 
      if  $p \neq \text{nil}$  then
         $\mathbb{X}' = \mathbb{X}' \cup p$ 
         $(N', E', cpt') = \text{tree}(\mathbb{X}', \text{root}, D)$ 
    return  $(N \cup N', E \cup E', cpt \cup cpt')$ 

```

a finite domain, p be a (Σ, dom) -relaxed acyclic PROBL0G program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , s be a (Σ, dom) -structure, and p' be the program $\alpha(\beta_{p_0, \mathcal{D}, \mathcal{U}}(p_0))$. Then, (1) $a \in \text{ground}(p)$ implies $a \downarrow_p \in \text{ground}(p')$, and (2) $s \in \text{ground}(p, r)$ implies $s \downarrow_p \in \text{ground}(p, r \downarrow_p)$.

Proof. Let Σ be a first-order signature, dom be a finite domain, p be a (Σ, dom) -relaxed acyclic PROBL0G program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , s be a (Σ, dom) -structure, and p' be the program $\alpha(\beta_{p_0, \mathcal{D}, \mathcal{U}}(p_0))$.

We first prove our first claim. Let $a \in \text{ground}(p)$. From this, it follows that there is an i such that $a \in \text{ground}(p, i)$. We now prove our claim by induction on i . The base case

is $a \in \text{ground}(p, 0)$. From this, it follows that a is either a ground atom or a probabilistic ground atom. From this and the definition of α and β , $a \downarrow_p = a$. From this, it follows that $a \in p'$ and, therefore, $a \in \text{ground}(p')$. For the induction step, assume that our claim holds for all $j < i$, we now show that it holds also for i . The only interesting case is $a \in \text{ground}(p, i) \setminus \text{ground}(p, i-1)$. From this, it follows that there is a rule r and a ground rule s such that all $\text{body}^+(s) \subseteq \text{ground}(p, i-1)$. From this, $r \downarrow_p \in p'$, the fact that the transformation does not introduce new positive literals, and the induction hypothesis, it follows that $\{b \downarrow_p \mid b \in \text{body}^+(s)\} \subseteq \text{ground}(p')$. From this and ground' 's definition, it follows that $a \downarrow_p \in \text{ground}(p')$.

The proof of our second claim is similar to the first one. \square

Here are now ready to prove the two key lemmas for the encoding's correctness.

Proposition B.15. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , s be a (Σ, \mathbf{dom}) -structure, and $BN = (N, E, cpt)$ be the Bayesian Network generated by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as input. There is a p -probabilistic assignment f such that $prob(f) = k$ and $g_f(p) = s$ iff there is a BN -total assignment ν such that $\llbracket BN \rrbracket(\nu) = k$, ν is consistent, and $\nu \models s$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , s be a (Σ, \mathbf{dom}) -structure, and $BN = (N, E, cpt)$ be the Bayesian Network generated by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as input. We now prove both directions of our claim.

(\Rightarrow) . Let f be a p -probabilistic assignment f such that $prob(f) = k$ and $g_f(p) = s$. Furthermore, let ν be the following BN -total assignment:

$$\nu(n) = \begin{cases} f(v::a(\bar{c})) & \text{if } n = X[v::a(\bar{c}), \emptyset, a(\bar{c})] \\ \top & \text{if } n = X[a(\bar{c}), \emptyset, a(\bar{c})] \\ k & \text{if } k \in D(n) \wedge p(n) = \{P_1, \dots, P_m\} \\ & \wedge \nu(P_1) = v_1 \wedge \dots \wedge \nu(P_m) = v_m \\ & \wedge cpt(n)(v_1, \dots, v_m, k) = 1 \end{cases}$$

The assignment ν is well-defined since (1) BN is a forest of poly-trees (see Proposition B.13), and (2) for all nodes of the form $X[v::a(\bar{c}), \emptyset, a(\bar{c})]$ and $X[a(\bar{c}), \emptyset, a(\bar{c})]$, $D(n) = \{\top, \perp\}$ by construction (cf. Algorithm 2). Furthermore, the assignment ν is consistent by construction. Indeed, for all variables $n \in N \setminus K(BN)$, the corresponding entry in the CPT is 1.

We now prove that $\nu \models s$. From Proposition B.16, ν 's consistency, and $f(v::a(\bar{c})) = \nu(X[\{v::a(\bar{c})\}, \emptyset, a(\bar{c})])$ for all $v::a(\bar{c}) \in p$, it follows that $a(\bar{c}) \in g_f(p)$ iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$. From this and $g_f(p) = s$, it follows $\nu \models s$.

Finally, we show that $\llbracket BN \rrbracket(\nu) = prob(f)$. In more detail, $\llbracket BN \rrbracket(\nu) = [(\prod_{n \in N} cpt(n))](\nu)$. This can be equivalently rewritten as follows: $\llbracket BN \rrbracket(\nu) = [(\prod_{n \in K(BN)} cpt(n)) \cdot (\prod_{n \in N \setminus K(BN)} cpt(n))](\nu)$. Furthermore, since ν is consistent and the CPTs associated with the variables in $N \setminus K(BN)$ are deterministic, $\llbracket BN \rrbracket(\nu)$ can be simplified as $\llbracket BN \rrbracket(\nu) = [(\prod_{n \in K(BN)} cpt(n))](\nu)$. From this and BN 's definition, it follows that $\llbracket BN \rrbracket(\nu) = \prod_{\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = \top} v \cdot \prod_{\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = \perp} (1 - v)$. From this and ν 's definition, it follows $\llbracket BN \rrbracket(\nu) = \prod_{f(v::a(\bar{c})) = \top} v \cdot \prod_{f(v::a(\bar{c})) = \perp} (1 - v)$, which is equivalent to $prob(f)$.

(\Leftarrow) . Let ν be a BN -total assignment such that $\llbracket BN \rrbracket(\nu) = k$, ν is consistent, and $\nu \models s$, and f be the following p -probabilistic assignment: $f(v::a(\bar{c})) = \nu(X[\{v::a(\bar{c})\}, \emptyset, a(\bar{c})])$.

We now show that $\llbracket BN \rrbracket(\nu) = prob(f)$. In more detail, $\llbracket BN \rrbracket(\nu) = [(\prod_{n \in N} cpt(n))](\nu)$. Furthermore, since ν is consistent and the CPTs associated with the variables in $N \setminus K(BN)$ are deterministic, $\llbracket BN \rrbracket(\nu)$ can be simplified as $\llbracket BN \rrbracket(\nu) = [(\prod_{n \in K(BN)} cpt(n))](\nu)$. From this and BN 's definition, it follows that $\llbracket BN \rrbracket(\nu) = \prod_{\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = \top} v \cdot \prod_{\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = \perp} (1 - v)$. From this and ν 's definition, it follows that $\llbracket BN \rrbracket(\nu) = \prod_{f(v::a(\bar{c})) = \top} v \cdot \prod_{f(v::a(\bar{c})) = \perp} (1 - v)$, which is equivalent to $prob(f)$.

We still have to prove that $g_f(p) = s$. From Proposition B.16, ν 's consistency, $f(v::a(\bar{c})) = \nu(X[\{v::a(\bar{c})\}, \emptyset, a(\bar{c})])$, it follows that $a(\bar{c}) \in g_f(p)$ iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$. From this and $\nu \models s$, it follows $g_f(p) = s$. \square

Proposition B.16. *Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , $BN = (N, E, cpt)$ be the Bayesian Network generated by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as input, f be a p -probabilistic assignment, and ν be a consistent BN -variable assignment such that $\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = f(v::a(\bar{c}))$ for all $v::a(\bar{c}) \in p$. Then, $a(\bar{c}) \in g_f(p, j, i)$, for some j and i , iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$.*

Proof. Let Σ be a first-order signature, \mathbf{dom} be a finite domain, p be a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , $BN = (N, E, cpt)$ be the Bayesian Network generated by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as input, f be a p -probabilistic assignment, and ν be a consistent BN -variable assignment such that $\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = f(v::a(\bar{c}))$ for all $v::a(\bar{c}) \in p$. Furthermore, let p' be the acyclic PROBLOG program obtained after applying the α and β transformations. We claim that $a(\bar{c}) \in g_f(p, \mu(a))$ iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$. From this, $a(\bar{c}) \in g_f(p)$ iff $a(\bar{c}) \in g_f(p, j, i)$ for some j and i , and $a(\bar{c}) \in g_f(p)$ iff $a(\bar{c}) \in g_f(p, \mu(a))$, it follows that $a(\bar{c}) \in g_f(p, j, i)$, for some j and i , iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$. However, there may be some nodes in BN that do not correspond to any ground rule or atom in $g_f(p, r)$ or $g_f(p)$.

We now prove, by induction on $\mu(a)$ (define in §B-A2), our claim that $a(\bar{c}) \in g_f(p, \mu(a))$ iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$ (we denote the corresponding induction hypothesis as (\star)).

Base Case. For the base case, we assume that $\mu(a) = 0$. We prove separately the two directions, namely (1) if $a(\bar{c}) \in g_f(p, 0, i)$ and $\mu(a) = 0$, then $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$, and (2) if $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$ and $\mu(a) = 0$, then $a(\bar{c}) \in g_f(p, 0, i)$. From this, it follows that if $\mu(a) = 0$, then $a(\bar{c}) \in g_f(p, 0)$ iff $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$.

(\Rightarrow) . We prove, by induction on i , that if $a(\bar{c}) \in g_f(p, 0, i)$, then $\nu(X[a(\bar{c})\downarrow_p]) = a(\bar{c})\uparrow_p$ (we denote this induction hypothesis as (\dagger)). From $a(\bar{c}) \in g_f(p)$, $g_f(p) \subseteq ground(p)$, and Proposition B.14, then $X[a(\bar{c})\downarrow_p] \in N$. The base case is as follows. If $a(\bar{c}) \in g_f(p, 0, 0)$, then there is a rule r such that either $r = a(\bar{c})$ or $r = v::a(\bar{c})$ and $a(\bar{c}) \in g_f(p, r, 0, 0)$. If $r = a(\bar{c})$, then $\nu(X[r, \emptyset, a(\bar{c})]) = \top$ due to ν 's consistency, $g_f(p) \subseteq ground(p)$, Proposition B.14, $r\downarrow_p = r$, and $a(\bar{c})\downarrow_p = a(\bar{c})$. If $r = v::a(\bar{c})$, then $f(v::a(\bar{c})) = \top$ follows

from $a(\bar{c}) \in g_f(p, 0, 0)$. From this, $g_f(p) \subseteq \text{ground}(p)$, Proposition B.14, $r \downarrow_p = r$, $a(\bar{c}) \downarrow_p = a(\bar{c})$, and $\nu(X[v::a(\bar{c}), \emptyset, a(\bar{c})]) = f(v::a(\bar{c}))$, it follows that there is a variable $X[r, \emptyset, a(\bar{c})] \in N$ such that $\nu(X[r, \emptyset, a(\bar{c})]) = \top$. From $\nu(X[r, \emptyset, a(\bar{c})]) = \top$, BN 's definition, and nu 's definition, it follows that there are nodes $X[r, a(\bar{c})], X[a(\bar{c})] \in N$ such that $\nu(X[r, a(\bar{c})]) = \top$ and $\nu(X[a(\bar{c})]) = \top$ as well (note that $r \downarrow_p = r$, and $a(\bar{c}) \downarrow_p = a(\bar{c})$ in this case). For the induction step, we assume that our claim holds for all $i' < i$. The only interesting case is when $a(\bar{c}) \in g_f(p, 0, i) \setminus g_f(p, 0, i - 1)$. From this, it follows that there is rule $r \in p$ such that $a(\bar{c}) \leftarrow b_1, \dots, b_m \in g_f(p, r, 0, i)$. From this, it follows that $b_1, \dots, b_m \in g_f(p, 0, i - 1)$ and $\text{body}^-(r) = \emptyset$. From this, $g_f(p, 0, i - 1) \subseteq \text{ground}(p)$, Proposition B.14, and the induction's hypothesis (\dagger), it follows that there are nodes $X[b_1 \downarrow_p], \dots, X[b_m \downarrow_p] \in N$ and $\nu(X[b_1 \downarrow_p]) = b_1 \uparrow_p, \dots, \nu(X[b_m \downarrow_p]) = b_m \uparrow_p$. From this, $\text{body}^-(r) = \emptyset$, $g_f(p, r, 0, i) \subseteq g_f(p, r) \subseteq \text{ground}(p, r)$, Proposition B.14, and BN 's construction, it follows that there is a variable $X[r \downarrow_p, \{b_1 \downarrow_p, \dots, b_m \downarrow_p\}, a(\bar{c}) \downarrow_p] \in N$ such that $\text{cpt}(X[r \downarrow_p, \{b_1 \downarrow_p, \dots, b_m \downarrow_p\}, a(\bar{c}) \downarrow_p]) (\nu(X[b_1 \downarrow_p]), \dots, \nu(X[b_m \downarrow_p]), k) = 1$ iff $k = a(\bar{c}) \uparrow_p$. As a result, $\nu(X[r \downarrow_p, \{b_1 \downarrow_p, \dots, b_m \downarrow_p\}, a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ since ν is consistent. From this and BN 's construction, it follows that there are nodes $X[r \downarrow_p, a(\bar{c}) \downarrow_p], X[a(\bar{c}) \downarrow_p] \in N$ such that $\nu(X[r \downarrow_p, a(\bar{c}) \downarrow_p]) = \top$ and $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$. This completes the proof of the if direction.

(\Leftarrow). To prove that if $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, then $a(\bar{c}) \in g_f(p, 0, i)$, for some i , we prove a stronger claim. Namely, if $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, then $a(\bar{c}) \in g_f(p, 0, \text{depth}(X[a(\bar{c}) \downarrow_p]))$, where $\text{depth}(n) = 0$ if $\text{ancestors}(n)$ does not contain any variable of the form $X[b(\bar{v})]$ and $\text{depth}(n) = 1 + \max_{n' \in \text{ancestors}(n)} \text{depth}(n')$ otherwise. We prove our claim by induction on $\text{depth}(X[a(\bar{c}) \downarrow_p])$ (we denote the induction hypothesis as (Δ)). The base case is as follows. If $\text{depth}(X[a(\bar{c}) \downarrow_p]) = 0$, then from $\nu(X[a(\bar{c}) \downarrow_p]) = \top$ and ν 's consistency, it follows that there must be a rule r , of the form $r = v::a(\bar{c})$ or $r = a(\bar{c})$, such that $X[r, a(\bar{c})] \in N$ and $\nu(X[r, a(\bar{c})]) = \top$ (note that, in this case, $r \downarrow_p = r$ and $a(\bar{c}) \downarrow_p = a(\bar{c})$). From this and ν 's consistency, there is a variable $X[r, \emptyset, a(\bar{c})] \in N$ such that $\nu(X[r, \emptyset, a(\bar{c})]) = \top$. If $r = a(\bar{c})$, then $a(\bar{c}) \in g_f(p, 0, 0)$ by definition. If $r = v::a(\bar{c})$, then from $\nu(X[r, \emptyset, a(\bar{c})]) = \top$ and $\nu(X[r, \emptyset, a(\bar{c})]) = f(v::a(\bar{c}))$, it follows that $f(v::a(\bar{c})) = \top$. From this, it follows that $a(\bar{c}) \in g_f(p, 0, 0)$. For the induction step, we assume that our claim holds for all random variables of depth less than k . From $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ and ν 's consistency, it follows that there is a rule $r \in p'$ such that $\nu(X[r, a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$. From this and ν 's consistency, there is a set of positive ground literals $I = (b_1, \dots, b_m)$ such that $X[r, (b_1 \downarrow_p, \dots, b_m \downarrow_p), a(\bar{c}) \downarrow_p] \in N$ and $\nu(X[r, (b_1 \downarrow_p, \dots, b_m \downarrow_p), a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$. From this and BN 's construction (cf. the *cpt* and *satisfiable* procedures), it follows that there is a rule $r' \in p$ such that (1) $r \in [r']_p$, and (2) the values assigned by ν to the random variables associated to

the atoms in I produce a grounding s of r that satisfies all constraints and is consistent (namely, the body of s does not contain both an atom and its negation and multiple copies of the same CPT-like atom in r' are assigned to the same value in r). From this, BN 's definition, and $\mu(a) = 0$, it follows that $\nu(X[b \downarrow_p]) = b \uparrow_p$ for all $b \in I$. From this and the induction's hypothesis (Δ), it follows that $b \in g_f(p, 0, \text{depth}(X[a(\bar{c})]) - 1)$ for all $b \in I$. From this, $\nu(X[r \downarrow_p, \{b_1 \downarrow_p, \dots, b_m \downarrow_p\}, a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, and BN 's construction (cf. the *cpt* and *satisfiable* procedures), it follows that there is a rule $r' \in p$ such that (1) $r \downarrow_p = r' \downarrow_p$, and (2) r' is satisfied by $\{b_1, \dots, b_m\}$, i.e., there is a grounding s' of r' obtained by using a subset of the literals in $\{b_1, \dots, b_m\}$ and adding repeated occurrences of the literals if needed. From this, it follows that s is in $g_f(p, r, 0, \text{depth}(X[a(\bar{c}) \downarrow_p]))$. From this, it follows that $a(\bar{c}) \in g_f(p, 0, \text{depth}(X[a(\bar{c}) \downarrow_p]))$. This completes the proof of the only if direction.

Induction Step. For the induction's step, we assume that $b(\bar{d}) \in g_f(p, \mu(b))$ iff $\nu(X[b(\bar{d}) \downarrow_p]) = b(\bar{d}) \uparrow_p$ holds for all b such that $\mu(b) < \mu(a)$. We now prove that $a(\bar{c}) \in g_f(p, \mu(a))$ iff $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ as well. In the following, let k be $\mu(a)$. We prove separately the two directions, namely (1) if $a(\bar{c}) \in g_f(p, \mu(a), i)$, for some i , then $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, and (2) if $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, then $a(\bar{c}) \in g_f(p, \mu(a), i)$, for some i . From this, it follows that $a(\bar{c}) \in g_f(p, \mu(a))$ iff $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$.

(\Rightarrow). We prove, by induction on i , that if $a(\bar{c}) \in g_f(p, k, i)$, for some i , then $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ (we denote the induction hypothesis associated to this proof as (\clubsuit)). From $a(\bar{c}) \in g_f(p) \subseteq \text{ground}(p)$ and Proposition B.14, then $X[a(\bar{c}) \downarrow_p] \in N$. The base case is as follows. Assume that $a(\bar{c}) \in g_f(p, k, 0)$. From this and $a(\bar{c}) \notin g_f(p, k - 1)$ (since $\mu(a) > k - 1$), it follows that there is a rule r such that either $r = a(\bar{c})$ or $r = v::a(\bar{c})$. If $r = a(\bar{c})$, then there is a node $X[r, \emptyset, a(\bar{c})] \in N$ such that $\nu(X[r, \emptyset, a(\bar{c})]) = \top$ due to ν 's consistency, $g_f(p, r) \subseteq \text{ground}(p, r)$, $r \downarrow_p = r$, $a(\bar{c}) \downarrow_p = a(\bar{c})$, and Proposition B.14. If $r = v::a(\bar{c})$, then $f(v::a(\bar{c})) = \top$ follows from $a(\bar{c}) \in g_f(p, k, 0)$. From this, $g_f(p, r) \subseteq \text{ground}(p, r)$, Proposition B.14, and $\nu(X[\{v::a(\bar{c})\}, \emptyset, a(\bar{c})]) = f(v::a(\bar{c}))$, it follows that there is a node $X[r, \emptyset, a(\bar{c})] \in N$ such that $\nu(X[r, \emptyset, a(\bar{c})]) = \top$. From $\nu(X[r, \emptyset, a(\bar{c})]) = \top$, BN 's definition, and nu 's definition, it follows that there are nodes $X[r, a(\bar{c})], X[a(\bar{c})] \in N$ such that $\nu(X[r, a(\bar{c})]) = \top$ and $\nu(X[a(\bar{c})]) = \top$ as well (note that $r \downarrow_p = r$ and $a(\bar{c}) \downarrow_p = a(\bar{c})$). For the induction step, we assume that our claim holds for all $i' < i$. Assume that $a(\bar{c}) \in g_f(p, k, i)$. The only interesting case is when $a(\bar{c}) \in g_f(p, k, i) \setminus g_f(p, k, i - 1)$. From this, it follows that there is rule r such that $r' = a(\bar{c}) \leftarrow b_1, \dots, b_m$ and $r' \in g_f(p, r, k, i)$. Furthermore, we denote by b'_1, \dots, b'_m the atoms $\text{pos}(b_1), \dots, \text{pos}(b_m)$. From this, it follows that $b_e \in g_f(p, k, i - 1)$ for all $b_e \in \text{body}^+(r')$ and $b'_d \notin g_f(p, k - 1)$ for all $b'_d \in \text{body}^-(r')$. From $b_e \in g_f(p, k, i - 1)$ for all $b_e \in \text{body}^+(r')$, $g_f(p, r) \subseteq \text{ground}(p, r)$, Proposition B.14, and the induction's hypothesis (\clubsuit), it follows that there is a node $X[b_e \downarrow_p] \in N$ such that $\nu(X[b_e \downarrow_p]) = b_e \uparrow_p$ for all

$b_e \in \text{body}^+(r')$. From $b'_d \notin g_f(p, k-1)$ for all $b_d \in \text{body}^-(r')$, it follows that $b'_d \notin g_f(p, \mu(\text{pred}(b_d)))$ for all $b_d \in \text{body}^-(r')$. From this, $r' \in g_f(p, r, k, i)$, $g_f(p, r, k, i) \subseteq \text{ground}(p, r)$, Proposition B.14, $\mu(\text{pred}(b_d)) < \mu(a)$ for all $b_d \in \text{body}^-(r')$, and the induction's hypothesis (\star), it follows that there is a node $X[b'_d \downarrow_p] \in N$ such that $\nu(X[b'_d \downarrow_p]) \neq b'_d \uparrow_p$ for all $b_d \in \text{body}^-(r')$. From $r' \in g_f(p, r, k, i)$, $g_f(p, r, k, i) \subseteq \text{ground}(p, r)$, Proposition B.14, $\nu(X[b_e \downarrow_p]) = b_e \uparrow_p$ for all $b_e \in \text{body}^+(r')$, $\nu(X[b'_d \downarrow_p]) \neq b'_d \uparrow_p$ for all $b_d \in \text{body}^-(r')$, and BN 's construction, there is $X[r \downarrow_p, \{b'_1 \downarrow_p, \dots, b'_m \downarrow_p\}, a(\bar{c}) \downarrow_p] \in N$ such that $\text{cpt}(X[r \downarrow_p, \{b'_1 \downarrow_p, \dots, b'_m \downarrow_p\}, a(\bar{c}) \downarrow_p]) (\nu(X[b'_1 \downarrow_p]), \dots, \nu(X[b'_m \downarrow_p]), k) = 1$ iff $k = a(\bar{c}) \uparrow_p$. As a result, $\nu(X[r \downarrow_p, \{b'_1 \downarrow_p, \dots, b'_m \downarrow_p\}, a(\bar{c}) \downarrow_p]) = \top$ since ν is consistent. From this and BN 's construction, it follows that there are variables $X[r \downarrow_p, a(\bar{c}) \downarrow_p], X[a(\bar{c}) \downarrow_p] \in N$ such that $\nu(X[r \downarrow_p, a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ and $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$. This completes the proof of the if direction.

(\Leftarrow). To prove that if $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, then $a(\bar{c}) \in g_f(p, k, i)$, for some i , we prove a stronger claim. Namely, if $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, then $a(\bar{c}) \in g_f(p, k, \text{depth}(X[a(\bar{c}) \downarrow_p]))$, where $\text{depth}(n)$ is as above. We prove our claim by induction on $\text{depth}(X[a(\bar{c}) \downarrow_p])$ (we denote the induction hypothesis associated to this proof as (\spadesuit)). The base case is as follows. If $\text{depth}(X[a(\bar{c}) \downarrow_p]) = 0$, then from $\nu(X[a(\bar{c}) \downarrow_p]) = \top$ and ν 's consistency, it follows that there must be a rule r , of the form $r = v::a(\bar{c})$ or $r = a(\bar{c})$, such that $X[r, a(\bar{c})] \in N$ and $\nu(X[r, a(\bar{c})]) = \top$ (note that $r \downarrow_p = r$ and $a(\bar{c}) \downarrow_p = a(\bar{c})$). From this and ν 's consistency, there is a variable $X[r, \emptyset, a(\bar{c})] \in N$ such that $\nu(X[r, \emptyset, a(\bar{c})]) = \top$. If $r = a(\bar{c})$, then $a(\bar{c}) \in g_f(p, k, 0)$ by definition. If $r = v::a(\bar{c})$, then from $\nu(X[r, \emptyset, a(\bar{c})]) = \top$ and $\nu(X[r, \emptyset, a(\bar{c})]) = f(v::a(\bar{c}))$, it follows that $f(v::a(\bar{c})) = \top$. From this, it follows that $a(\bar{c}) \in g_f(p, k, 0)$. For the induction step, we assume that our claim holds for all random variables of depth less than k . We now show that it holds also for a variable $X[a(\bar{c}) \downarrow_p]$ such that $\text{depth}(X[a(\bar{c}) \downarrow_p]) = i$. From $\nu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ and ν 's consistency, it follows that there is a rule $r' \in p'$ such that $X[r', a(\bar{c}) \downarrow_p] \in N$ and $\nu(X[r', a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$. From this and ν 's consistency, there is a set of ground atoms $I = (b_1, \dots, b_m)$ such that $X[r', (b_1 \downarrow_p, \dots, b_m \downarrow_p), a(\bar{c}) \downarrow_p] \in N$, $\nu(X[r', (b_1 \downarrow_p, \dots, b_m \downarrow_p), a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, and $X[b \downarrow_p] \in N$ for all $b \in I$. From this and BN 's construction (cf. the *cpt* and *satisfiable* procedures), it follows that there is a rule $r' \in p$ such that (1) $r \in [r']_p$, and (2) the values assigned by ν to the random variables associated to the atoms in I produce a grounding s of r that satisfies all constraints and is consistent (namely, the body of s does not contain both an atom and its negation and multiple copies of the same CPT-like atom in r' are assigned to the same value in r). Let I^+ be the atoms in I that are assigned to positive literals in s and I^- be the atoms in I that are assigned to negative literals (note that these two sets are disjoint). From $\nu(X[r', (b_1 \downarrow_p, \dots, b_m \downarrow_p), a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$, and ν 's consistency, it follows that $\nu(X[b_e \downarrow_p]) = b_e \uparrow_p$ for all $b_e \in I^+$ and $\nu(X[b_d \downarrow_p]) \neq b_d \uparrow_p$ for all $b_d \in I^-$. From $\nu(X[b_e \downarrow_p]) = b_e \uparrow_p$

and $\text{depth}(X[b_e \downarrow_p]) < \text{depth}(X[a(\bar{c}) \downarrow_p])$ for all $b_e \in I^+$ and the induction's hypothesis (\spadesuit), it follows that $b_e \in g_f(p, k, \text{depth}(X[b_e \downarrow_p]))$ for all $b_e \in I^+$. From this, the induction's hypothesis, and $\text{depth}(X[b_e \downarrow_p]) < \text{depth}(X[a(\bar{c}) \downarrow_p])$, it follows that $b_e \in g_f(p, k, i-1)$ for all $b_e \in \text{body}^+(r')$. From $\nu(X[b_d \downarrow_p]) \neq b_d \uparrow_p$ and $\mu(\text{pred}(b_d)) < \mu(a)$ for all $b_d \in I^-$ and the induction's hypothesis (\star), it follows that $b_d \notin g_f(p, \mu(\text{pred}(b_d)))$ for all $b_d \in I^-$. From this and $\mu(\text{pred}(b_d)) < \mu(a)$, it follows that $b_d \notin g_f(p, k-1)$ for all $b_d \in I^-$. From r' definition, $b_e \in g_f(p, k, i-1)$ for all $b_e \in I^+$, and $b_d \notin g_f(p, k-1)$ for all $b_d \in I^-$, it follows that $r' \in g_f(p, r, k, i)$. From this, it follows that $r' \in g_f(p, k, i)$ and $a(\bar{c}) \in g_f(p, k, i)$. This completes the proof the only if direction. \square

Proposition B.17 shows the correctness of our encoding.

Proposition B.17. *Let Σ be a first-order signature, dom be a finite domain, p be a (Σ, dom) -relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , BN be the Bayesian Network generated by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as input, and s be a (Σ, dom) -structure. Furthermore, let μ be the BN -partial assignment such that (1) for any ground atom $a(\bar{c})$, $\mu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ iff $a(\bar{c}) \in s$, and (2) $\mu(v)$ is undefined otherwise. Then, $\llbracket p \rrbracket(s) = \llbracket BN \rrbracket(\mu)$.*

Proof. Let Σ be a first-order signature, dom be a finite domain, p be a (Σ, dom) -relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , BN be the Bayesian Network generated by Algorithm 2 having p and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ as input, and s be a (Σ, dom) -structure. Furthermore, let μ be the BN -partial assignment such that (1) for any ground atom $a(\bar{c})$, $\mu(X[a(\bar{c}) \downarrow_p]) = a(\bar{c}) \uparrow_p$ iff $a(\bar{c}) \in s$, and (2) $\mu(v)$ is undefined otherwise.

The probability $\llbracket p \rrbracket(s)$ is $\sum_{f \in \mathcal{M}(p, s)} \text{prob}(f)$, where $\mathcal{M}(p, s)$ is the set of all assignments f such that $WFM(\text{instance}(p, f)) = s$. Equivalently, $\mathcal{M}(p, s)$ is the set of all probabilistic assignments f such that $g_f(p) = s$. Let K be the set of all total assignments that agree with μ for all variables of the form $X[a(\bar{c}) \downarrow_p]$. The probability $\llbracket BN \rrbracket(\mu)$ is $\sum_{\nu \in K} \llbracket BN \rrbracket(\nu)$. Since any non-consistent assignment has probability 0, $\llbracket BN \rrbracket(\mu) = \sum_{\nu \in K'} \llbracket BN \rrbracket(\nu)$, where K' is the set of all consistent assignments in K . From this, it follows that $\llbracket p \rrbracket(s) = \llbracket BN \rrbracket(\mu)$ iff $\sum_{f \in \{f \mid g_f(p) = s\}} \text{prob}(f) = \sum_{\nu \in K'} \llbracket BN \rrbracket(\nu)$. The latter follows trivially from Proposition B.15, which establishes a one-to-one mapping from $\{f \mid g_f(p) = s\}$ and K' that preserves probabilities. \square

8) *Size of the encoding:* Given a Bayesian Network $BN = (N, E, \text{cpt})$, the size of BN , denoted $|BN|$, is $|N| + |E| + \sum_{n \in N} |\text{cpt}(n)|$, where the size of a conditional probability table is just the number of rows in the table (i.e., the number of all assignments). The size of an atom $a(\bar{c})$ is $|a|$, whereas the size of a rule $h \leftarrow l_1, \dots, l_n$ is $|h| + \sum_{1 \leq i \leq n} l_i$. Finally, the size of a program p is $\sum_{r \in p} |r|$. The *ground version* of p , denoted $gv(p)$, is $\bigcup_{r \in p} \text{ground}(p, r)$, namely the relaxed grounding of all the rules in p . Note that $\text{ground}(p) \subseteq gv(p)$. Given a PROBLOG program p , we denote by $edb(p)$ the ground

(possibly probabilistic) atoms in p , whereas we denote by $rules(p)$ the rules, i.e., $edb(p) = \{a(\bar{c}) \in p\} \cup \{v::a(\bar{c}) \in p \mid 0 \leq v \leq 1\}$ and $rules(p) = \{r \in p \mid body(r) \neq \emptyset\}$.

Let p be a relaxed acyclic PROBLOG program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , $p' = \alpha(\beta_{p, \mathcal{D}, \mathcal{U}}(p))$ be the transformed program, $g = gv(p')$ be its ground version, and $bn(p, (\mathcal{O}, \mathcal{D}, \mathcal{U})) = (\mathcal{N}, \mathcal{E}, cpt)$ be the corresponding Bayesian Network derived by Algorithm 4. The number of nodes in N is $O(|rules(p')| \cdot |g|)$. Indeed, there is a node $X[a(\bar{c})]$ for each ground atom in g . Moreover, for each rule $r \in p'$ and ground rule $r' \in g$, there are nodes $X[r, body(r'), head(r')]$ and $X[r, head(r')]$. Finally, the number of intermediate nodes generated by the *tree* procedure is twice the number of nodes of the form $X[r, body(r'), head(r')]$.

The number of edges in E is $O(|rules(p')|^2 \cdot |g|)$. Indeed, there is an edge $X[r, a(\bar{c})] \rightarrow X[a(\bar{c})]$ for each rule r and ground atom $a(\bar{c})$. Furthermore, there is an edge $X[b(\bar{v})] \rightarrow X[r, I, a(\bar{c})]$ for each rule r , ground rule $a(\bar{c}) \leftarrow I$, and atom $b(\bar{v}) \in I$. Finally, the number of edges introduced by the *tree* procedure is $O(|g|)$.

The size of *cpt* is $O(|rules(p')| \cdot |g| \cdot \max(2, |rules(p)|^{2+|rules(p)|})$. Indeed, we have a CPT for each node $n \in N$. The size of each CPT depends on (1) the number of parents for each node, and (2) the size of the domain associated to each variable. In the worst case, the size of the domain of each node is $\max(2, |rules(p)|)$ (since the size of the domain of a CPT-like predicate depends only on the rules — there are no CPT-like predicates defined by ground atoms). The number of parents for intermediate nodes is 2, and the size of each CPT is $O(\max(2, |rules(p)|)^3)$. The maximum number of parents for nodes of the form $X[r, a(\bar{c})]$ is 1, and the size of each CPT is $O(\max(2, |rules(p)|)^2)$. The maximum number of parents for nodes of the form $X[r, I, a(\bar{c})]$ is $|rules(p')|$, and the size of each CPT is $O(\max(2, |rules(p)|)^{|rules(p')|+1})$. Similarly, the maximum number of parents for nodes of the form $X[a(\bar{c})]$ is $1 + |rules(p)|$, and the size of each CPT is $O(\max(2, |rules(p)|)^{|rules(p')|+2})$.

As a result, $|bn(p, (\mathcal{O}, \mathcal{D}, \mathcal{U}))|$ is $O(|rules(p')|^2 \cdot |g| \cdot \max(2, |rules(p)|^{2+|rules(p)|})$. Furthermore, since $|g| \in O(|edb(p')|^{|rules(p')|})$, it follows that $|bn(p, (\mathcal{O}, \mathcal{D}, \mathcal{U}))|$ is $O(|rules(p')|^2 \cdot |edb(p')|^{|rules(p')|} \cdot \max(2, |rules(p)|^{2+|rules(p)|})$. Finally, since $|edb(p')| \in O(|edb(p)|)$ and $|rules(p')| \in O(|rules(p)|)$, we can simplify the result as follows: $|bn(p, (\mathcal{O}, \mathcal{D}, \mathcal{U}))|$ is $O(|rules(p)|^2 \cdot |edb(p)|^{|rules(p)|} \cdot \max(2, |rules(p)|^{2+|rules(p)|})$.

9) *Complexity Proofs*: We first define the inference problem INF. Afterwards, we analyse its complexity.

Problem 1. INF denotes the following decision problem:

Input: A first-order signature Σ , a finite domain \mathbf{dom} , a (Σ, \mathbf{dom}) -relaxed acyclic PROBLOG program p , a set of ground literals E , and a ground atom $a(\bar{c})$.

Output: The probability of $a(\bar{c})$ given evidence in E . \square

The data complexity of $\text{INF}(\Sigma, \mathbf{dom}, p, E, a)$ for relaxed

acyclic programs can be obtained by (1) fixing $rules(p)$ and varying only $edb(p)$ (and indirectly \mathbf{dom}), and (2) requiring the program to be relaxed acyclic.

The data complexity of the INF problem for relaxed-acyclic programs p is the complexity of the following decision problem:

Problem 2. Let Σ be a first-order signature Σ , R be a fixed set of PROBLOG rules over Σ , E be a set of ground literals, $a(\bar{c})$ be a ground atom, and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be an ordering template, a disjointness template, and a uniqueness template. $\text{INF}_{\Sigma, R, E, a(\bar{c}), (\mathcal{O}, \mathcal{D}, \mathcal{U})}^{\alpha}$ denotes the following problem:

Input: A set of probabilistic atoms E' such that (1) atoms in E and $a\bar{c}$ refer only to constant values in E' and R , and (2) $R \cup E'$ is a relaxed acyclic PROBLOG program and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ is a witness for the acyclicity of $R \cup E'$.

Output: The probability of $a(\bar{c})$ given evidence in E . \square

Proposition B.18. $\text{INF}_{\Sigma, R, E, a(\bar{c}), (\mathcal{O}, \mathcal{D}, \mathcal{U})}^{\alpha}$ is in PTIME.

Proof. Let Σ be a first-order signature Σ , R be a fixed set of PROBLOG rules over Σ , E be a set of ground literals, $a(\bar{c})$ be a ground atom, and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be an ordering template, a disjointness template, and a uniqueness template. We consider only inputs E' such that $E' \cup R$ is a relaxed acyclic PROBLOG program and $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ is a witness for the acyclicity of $R \cup E'$. The size of E' is the sum of the sizes of all atoms, where the size of an atom is its cardinality.

Let e be the number of distinct constants occurring in $E' \cup R$ and r be $|R|$. Computing $\text{INF}(\Sigma, E', p, E, a)$ can be done in the following steps:

- 1) Transform the original program p into the program p' (by applying the α and β transformations).
- 2) Construct the ground version g of the program p' .
- 3) Construct the Bayesian Network $bn(p, (\mathcal{O}, \mathcal{D}, \mathcal{U}))$ from g .
- 4) Perform the inference on $bn(p, (\mathcal{O}, \mathcal{D}, \mathcal{U}))$.

The first step can be performed in linear time in the size of $O(e + r)$. The second step can be performed in $O(e^r)$ (because the grounding of p' can be done in $O(e^{|rules(p')|})$ and $|rules(p')| \leq r$). The third step can be performed in $O(r^4 \cdot e^{2 \cdot r} \cdot \max(2, r)^{2+r})$ (constructing N and E can be done in $O(r^2 \cdot e^{2 \cdot r})$, whereas the *cpt* can be constructed in $O(r^4 \cdot e^r \cdot \max(2, r)^{2+r})$). The fourth step can be performed in $O(|bn(p)|)$ since $bn(p)$ is a forest of poly-trees (see Proposition B.13). In particular, inference can be performed by (1) identifying the poly-tree that contains the atom $a(\bar{x})$ (in $O(|bn(p)|)$), (2) identifying the subset $E'' \subseteq E$ of all atoms in the poly-tree of $a(\bar{x})$ (again in $O(|bn(p)|)$), and (3) performing inference using belief-propagation (in $O(|bn(p)|)$ [47]). Therefore, the fourth step can be performed in $O(r^2 \cdot e^r \cdot \max(2, r)^{2+r})$. As a result, the answer to INF can be computed in $O(r^2 \cdot e^{2 \cdot r} \cdot \max(2, r)^{2+r})$. Furthermore, since r is fixed and the number of constants in E' is $O(|E'|)$, there is a $k \in \mathbb{N}$ such that $e \in O(|E'| + k)$. From this, it follows that there is a $k \in \mathbb{N}$ such that INF can be computed in $O(|E'|^k)$. Hence, the complexity of $\text{INF}_{\Sigma, R, E, a(\bar{c}), (\mathcal{O}, \mathcal{D}, \mathcal{U})}^{\alpha}$ (and the data complexity

of INF) is PTIME. \square

INF for PROBLG is $\#P$ -hard. This follows from the $\#P$ -hardness of inference on arbitrary Bayesian networks (BNs) [47], which can be encoded in PROBLG. We now show that inference for PROBLG programs is $\#P$ -hard even in terms of data complexity. We do this using a reduction from the $\#SAT$ problem (counting the number of satisfying assignments of a propositional formula ϕ), which is $\#P$ -hard [56].

Proposition B.19. *INF is $\#P$ -hard in terms of data complexity for PROBLG programs.*

Proof. We show this by reducing the $\#SAT$ problem to inference for a PROBLG program p where (1) the formula can be encoded in $edb(p)$, and (2) the rules are fixed. Let ϕ be a propositional formula.

First-order Signature. Let Σ be the signature containing the following predicate symbols:

- e of arity 1, which is used to store the identifiers associated to all sub-expressions of ϕ ,
- $state$ of arity 1, which is used to store the propositions in the model,
- sat of arity 1, which is used to denote whether an expression is satisfiable or not,
- $conj$ of arity 3 which is used to encode conjunctions,
- $disj$ of arity 3 which is used to encode disjunctions,
- neg of arity 2 which is used to encode negations, and
- $prop$ of arity 2 used to encode propositions.

Rules. We now define a fixed set of rules encoding the semantics of propositional logic.

$$\begin{aligned} sat(x) &\leftarrow e(x), prop(x, y), state(y) \\ sat(x) &\leftarrow e(x), neg(x, y), \neg sat(y) \\ sat(x) &\leftarrow e(x), conj(x, y, z), sat(y), sat(z) \\ sat(x) &\leftarrow e(x), disj(x, y, z), sat(y) \\ sat(x) &\leftarrow e(x), disj(x, y, z), sat(z) \end{aligned}$$

Encoding a formula ϕ . Given a formula ϕ , we define the following encoding using probabilistic ground atoms. We first associate to each sub-formula ψ of ϕ a unique identifier id_ψ . We also associate a unique identifier id_p to each propositional symbol p . For each sub-formula $\psi \wedge \gamma$, the ground atoms $e(id_{\psi \wedge \gamma})$ and $conj(id_{\psi \wedge \gamma}, id_\psi, id_\gamma)$ are in E' . For each sub-formula $\psi \vee \gamma$, the ground atoms $e(id_{\psi \vee \gamma})$ and $disj(id_{\psi \vee \gamma}, id_\psi, id_\gamma)$ are in E' . For each sub-formula $\neg\psi$, the ground atoms $e(id_{\neg\psi})$ and $neg(id_{\neg\psi}, id_\psi)$ are in E' . For each sub-formula ψ such that ψ is a propositional symbol p , the ground atoms $e(id_\psi)$ and $prop(id_\psi, id_p)$ are in E' . Finally, for each propositional symbol p , there is a propositional atom $1/2::state(id_p)$ in E' .

Reduction. There are $2^{n(\phi)}$ grounded instances of $R \cup E'$, where $n(\phi)$ is the number of predicate symbols in ϕ , each one with probability $1/2^{n(\phi)}$. The grounded instances represent all possible assignments of \top and \perp to the proposition symbols

in ϕ . It is easy to see that $sat(id_\phi)$ can be derived in a grounded instance iff it represents a model for ϕ . Therefore, the probability $\llbracket R \cup E' \rrbracket(id_\phi)$ is $k/2^{n(\phi)}$, where k is the number of models that satisfy ϕ . \square

Note that the encoding shown in the previous proof can be tweaked to work for acyclic PROBLG programs but only for propositional formulae without repetitions of propositional symbols. We remark that the $\#SAT$ problem restricted to formulae without repetitions is no longer $\#P$ -hard. Indeed, it can be solved in PTIME as follows: given a formula ϕ without repetitions, we can construct a poly-tree boolean Bayesian Network BN encoding ϕ in $O(|\phi|)$, where the nodes associated to sub-formulae of the from p have a uniform probability distribution (i.e., p is \top with probability $1/2$ and \perp with probability $1/2$). Then, the probability associated to the root of ϕ is going to be $k/2^{n(\phi)}$, where k is the number of satisfying assignments. Since the inference on BN can be done in linear time in $|BN|$ and $|BN| \in O(|\phi|)$, then the whole problem is in PTIME.

10) Expressiveness: Here we show that any BN that consists of a forest of poly-trees can be represented as a relaxed acyclic program.

Proposition B.20. *Any BN that is a forest of poly-trees can be represented as a relaxed acyclic PROBLG program.*

Proof. Let bn be a BN that is a forest of poly-trees. We now construct the corresponding relaxed acyclic PROBLG program. In particular, for each random variable X in bn , we show how to equivalently encode it as PROBLG rules. Without loss of generality, we assume there is a unique mapping id from random variables in bn to predicate symbols identifiers. With a slight abuse of notation, we use X to refer both to the random variable and to the corresponding symbol $id(X)$.

Boolean Random Variables without parents. For each boolean random variable X such that (1) $p(X) = \emptyset$, and (2) $cpt(X) = \{\top \mapsto v, \perp \mapsto (1-v)\}$, we introduce a probabilistic atom $v::X$.

Non-Boolean Random Variables without parents. For each non-boolean random variable X with domain $\{q_1, \dots, q_n\}$ such that (1) $p(X) = \emptyset$, and (2) $cpt(X) = \{q_1 \mapsto v_1, \dots, q_n \mapsto v_n\}$, we introduce an annotated disjunction $v_1::X(q_1); \dots; v_n::X(q_n)$.

Boolean Random Variables with parents. Let X be a boolean random variable X with parents $p(X) = \{Y_1, \dots, Y_n\}$, $cpt(X)$ be the function associated with X , and $\bar{v}_1, \dots, \bar{v}_m$ be all possible assignments to the variables in $p(X)$. We introduce m fresh predicate symbols $sw_{X, \bar{v}_1}, \dots, sw_{X, \bar{v}_m}$. For each \bar{v}_i , we introduce the probabilistic atom $c_i::sw_{X, \bar{v}_i}$, where $cpt(X)(\bar{v}_i, \top) = c_i$. Finally, for each \bar{v}_i , we also introduce

the following rules:

$$\begin{aligned} X &\leftarrow Y_1(\bar{v}_i(1)), \dots, Y_n(\bar{v}_i(n)), sw(\bar{l}_1), sw_{X, \bar{v}_i} \\ &\vdots \\ X &\leftarrow Y_1(\bar{v}_i(1)), \dots, Y_n(\bar{v}_i(n)), sw(\bar{l}_{2^{i-1}}), sw_{X, \bar{v}_i} \end{aligned}$$

where $Y_k(\top) = Y_k$, $Y_k(\perp) = \neg Y_k$, and $Y_1(v) = Y_1(v)$ if $v \notin \{\top, \perp\}$, $\bar{l}_1, \dots, \bar{l}_{2^{i-1}}$ are all possible values in $\{\top, \perp\}^{i-1}$, and $sw(l_1, \dots, l_{i-1})$ is the list of literals $sw_{X, \bar{v}_1}(l_1), \dots, sw_{X, \bar{v}_{i-1}}(l_{i-1})$.

Non-Boolean Random Variables with parents. Let X be a non-boolean random variable X with domain $\{q_1, \dots, q_m\}$ and parents $p(X) = \{Y_1, \dots, Y_n\}$, $cpt(X)$ be the function associated with X , and $\bar{v}_1, \dots, \bar{v}_m$ be all possible assignments to the variables in $p(X)$. We introduce m fresh predicate symbols $sw_{X, \bar{v}_1}, \dots, sw_{X, \bar{v}_m}$. For each \bar{v}_i , we introduce the probabilistic atom $c_i :: sw_{X, \bar{v}_i}$, where $cpt(X)(\bar{v}_i, \top) = c_i$. Finally, for each \bar{v}_i , we also introduce the following rules:

$$\begin{aligned} p_1^i :: X(q_1); \dots; p_m^i :: X(q_m) &\leftarrow Y_1(\bar{v}_i(1)), \dots, Y_n(\bar{v}_i(n)), \\ &\quad sw(\bar{l}_1), sw_{X, \bar{v}_i} \\ &\vdots \\ p_1^i :: X(q_1); \dots; p_m^i :: X(q_m) &\leftarrow Y_1(\bar{v}_i(1)), \dots, Y_n(\bar{v}_i(n)), \\ &\quad sw(\bar{l}_{2^{i-1}}), sw_{X, \bar{v}_i} \end{aligned}$$

where $Y_k(\top) = Y_k$, $Y_k(\perp) = \neg Y_k$, and $Y_1(v) = Y_1(v)$ if $v \notin \{\top, \perp\}$, $\bar{l}_1, \dots, \bar{l}_{2^{i-1}}$ are all possible values in $\{\top, \perp\}^{i-1}$, $sw(l_1, \dots, l_{i-1})$ is the list of literals $sw_{X, \bar{v}_1}(l_1), \dots, sw_{X, \bar{v}_{i-1}}(l_{i-1})$, and $cpt(\bar{v}_i, q_j) = p_j^i$.

Correctness of the encoding. The encoding presented above encodes the CPTs for the random variables. The probability that a random variable X has value v is exactly the same as the probability associated with the ground literal $X(v)$ (with the notation defined above). The encoding for variables without parents directly encodes the corresponding CPTs. For variables with parents, the only non-standard part is the use of $sw(\bar{l}_{2^{i-1}}), sw_{X, \bar{v}_i}$. Note, however, that the literals in $sw(\bar{l}_{2^{i-1}})$ do not influence the derivation since there is a rule for any possible values for them. We need them only for the encoding. Therefore, also the encoding of random variables with parents directly encodes the corresponding CPTs.

Relaxed Acyclicity. Let p be the program produced by the above construction. It is easy to see that all predicates associated with non-boolean random variables are CPT-like (as we just encoded their CPTs). Therefore, the program $\beta(p)$ is obtained by removing all constant values associated with the domains of non-boolean random values. Finally, the program $\alpha(\beta(p))$ collapses the rules for random variables with parents into a single rule (this follows from our construction and the use of $sw(\bar{l}_{2^{i-1}}), sw_{X, \bar{v}_i}$ in the rules' bodies). The acyclicity of $p' = \alpha(\beta(p))$ follows from (1) p' does not contain free-variables (i.e., both the literals and the ground atoms are only propositional facts), (2) bn is a forest of poly-trees, and (3) each predicate symbol sw_{X, \bar{v}_i} occurs only in the rule of X .

From (1), it follows that all rules are both strongly and weakly connected. From (2) and (3), it follows that (a) there are no directed cycles in $graph(p)$, and (b) for all undirected cycles U in $graph(p)$, there are U', U'' such that U is equivalent to $U' \cdot U''$ and U' is $P' \xleftarrow{r, i} P \xrightarrow{r, i} P'$ for some P, P', r, i . Since all rules are both strongly and weakly connected, it directly follows that the undirected unsafe structure $\langle P \xrightarrow{r, i} P', P \xrightarrow{r, i} P', P', U'' \rangle$ is guarded. Therefore, p' is acyclic and p is a relaxed acyclic PROBLOG program. \square

Here, we present additional details about ANGERONA. We also prove its security, complexity, and completeness.

A. Checking Query Security

In the following, let $D = \langle \Sigma, \mathbf{dom} \rangle$ be a database schema. Without loss of generality, we focus only on relational calculus formulae ϕ where no distinct pair of quantifiers binds the same variable.

Normal Form. We say that a formula $\psi \wedge \neg\gamma$ is *guarded* iff $free(\gamma) \subseteq free(\psi)$. We say that a relational calculus formula ϕ is in *Normal Form* (NF) iff (1) ϕ uses only existential quantifiers, (2) negation is used only in sub-formulae of the form $\psi \wedge \neg\gamma$ and it is always guarded, (3) for any sub-formula $\psi \vee \gamma$, $free(\psi) = free(\gamma)$, (4) no distinct pair of quantifiers binds the same variable, and (5) there are no equality and inequality constraints.

Most of the time, domain-independent relational calculus formulae can be easily written in NF by just re-arranging sub-formulae. We remark that any domain-independent relational calculus formula can be written in NF by (1) extending the database schema with two relations eq and neq encoding $=$ and \neq among constants in \mathbf{dom} (this is always possible because \mathbf{dom} is finite), (2) renaming the quantified variables in a unique way, (3) replacing all universally quantified sub-formula $\forall x. \phi_e \rightarrow \psi$ with the equivalent existentially quantified version $\neg \exists x. \phi_e \wedge \neg\psi$, (4) replacing each negated sub-formula $\neg\psi$ with the equivalent sub-formula $(\bigwedge_{x \in free(\psi)} adom(x)) \wedge \neg\psi$, where $adom(x)$ is $\bigvee_{R \in D} \bigvee_{1 \leq i \leq |R|} \exists x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_{|R|}. R(x_1, \dots, x_{i-1}, x, x_{i+1}, \dots, x_{|R|})$, and (5) replacing sub-formulae of the form $\psi \vee \gamma$ with the equivalent formula $(\bigwedge_{x \in free(\gamma) \setminus free(\psi)} adom(x) \wedge \psi) \vee (\bigwedge_{x \in free(\psi) \setminus free(\gamma)} adom(x) \wedge \gamma)$. Note that the resulting NF formula is equivalent to the original one (because the rewriting does not modify the formula's semantics). Therefore, in the following we consider only NF relational calculus formulae.

From Relational Calculus to Logic Programming. Let ϕ be a NF relational calculus sentence. We denote by $sub(\phi)$ the sequence ϕ_1, \dots, ϕ_n of all ϕ 's sub-formulae ordered from the smallest to the largest, i.e., $\phi_n = \phi$, and by $DISub(\phi)$ the sequence obtained by removing from the sequence $sub(\phi)$ (1) all formulae of the form $\neg\psi$, and (2) removing all sub-formulae consisting just of equality and inequality terms. Since ϕ is in NF, all negated sub-formulae in ϕ appear only in NF in the sequence $DISub(\phi)$.

The function $PL(\phi)$ encodes any NF relational calculus sentence as a set of equivalent logic programming rules. We associate a unique predicate symbol H_i and a set of rules $r(H_i)$ to each ψ_i in $DISub(\phi)$ as follows:

- If $\psi_i := R(\bar{x}_i)$, for some $R \in \Sigma$, then $r(H_i)$ is $\{H_i(\bar{x}_i) \leftarrow R(\bar{x}_i)\}$, where \bar{x}_i are ψ_i 's free variables.
- If $\psi_i := \psi_j \wedge \neg\psi_j$, then $r(H_i)$ is $\{H_i(\bar{x}_i) \leftarrow H_j(\bar{x}_j), \neg H_k(\bar{x}_k)\}$, where \bar{x}_i are ψ_i 's free variables, \bar{x}_j are ψ_j 's

Algorithm 4: ANGERONA Enforcement Algorithm (repeated from §VII).

Input: A system configuration C , a system state $s = \langle db, U, P \rangle$, a history h , a C -ATKLOG model ATK , and an action $\langle u, q \rangle$.

Output: The security decision in $\{\top, \perp\}$.

```

begin
  for  $\langle u, \psi, l \rangle \in secrets(P, u)$  do
    if  $secure(C, ATK, h, \langle u, \psi, l \rangle)$ 
      if  $pox(C, ATK, h, \langle u, q \rangle, \top)$ 
         $h' := h \cdot \langle \langle u, q \rangle, \top, \top \rangle$ 
        if  $\neg secure(C, ATK, h', \langle u, \psi, l \rangle)$ 
          return  $\perp$ 
        if  $pox(C, ATK, h, \langle u, q \rangle, \perp)$ 
           $h' := h \cdot \langle \langle u, q \rangle, \top, \perp \rangle$ 
          if  $\neg secure(C, ATK, h', \langle u, \psi, l \rangle)$ 
            return  $\perp$ 
      return  $\top$ 

function  $secure(\langle D, \Gamma \rangle, ATK, h, \langle u, \psi, l \rangle)$ 
   $p := ATK(u)$ 
  for  $\phi \in knowledge(h, u)$  do
     $p := p \cup PL(\phi) \cup \{evidence(head(\phi), true)\}$ 
   $p := p \cup PL(\psi)$ 
  return  $\llbracket p \rrbracket_D(head(\psi)) < l$ 

function  $pox(\langle D, \Gamma \rangle, ATK, h, \langle u, \psi \rangle, v)$ 
   $p := ATK(u)$ 
  for  $\phi \in knowledge(h, u)$  do
     $p := p \cup PL(\phi) \cup \{evidence(head(\phi), true)\}$ 
  if  $v = \top$ 
     $p := p \cup PL(\psi)$ 
    return  $\llbracket p \rrbracket_D(head(\psi)) > 0$ 
  else
     $p := p \cup PL(\neg\psi)$ 
    return  $\llbracket p \rrbracket_D(head(\neg\psi)) > 0$ 

```

free variables, and \bar{x}_k are ψ_k 's free variables (note that $\bar{x}_i = \bar{x}_j$ and $\bar{x}_k \subseteq \bar{x}_j$).

- If $\psi_i = \psi_j \wedge \psi_k$, then $r(H_i)$ is $\{H_i(\bar{x}_i) \leftarrow H_j(\bar{x}_j), H_k(\bar{x}_k)\}$, where \bar{x}_i are ψ_i 's free variables, \bar{x}_j are ψ_j 's free variables, and \bar{x}_k are ψ_k 's free variables.
- If $\psi_i = \psi_j \vee \psi_k$, then $r(H_i)$ contains the rules $H_i(\bar{x}_i) \leftarrow H_j(\bar{x}_j)$ and $H_i(\bar{x}_i) \leftarrow H_k(\bar{x}_k)$ (note that $free(\psi_i) = free(\psi_j) = free(\psi_k)$).
- If $\psi_i = \exists x. \psi_j$, then $r(H_i)$ is $\{H_i(\bar{x}_i) \leftarrow H_j(\bar{x}_j), c_1, \dots, c_n\}$, where \bar{x}_i are ψ_i 's free variables and \bar{x}_j are ψ_j 's free variables (i.e., those in \bar{x}_i and x).

Furthermore, we denote $head(\phi)$ the predicate symbol associated to ψ_m .

Checking Query Security. ANGERONA's enforcement algorithm is shown in Algorithm 4 (repeated from §VII).

B. Security Proof

Let $C = \langle D, \Gamma \rangle$ be a system configuration and h be a C -history. The set $knowledge(h, u)$ is $\{\phi \mid \exists i. h(i) = \langle \langle u, \phi \rangle, \top, \top \rangle\} \cup \{\neg\phi \mid \exists i. h(i) = \langle \langle u, \phi \rangle, \top, \perp \rangle\}$.

We first prove that ANGERONA's result depends just on the queries in the history, and not on the actual database's state.

Proposition C.1. Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, and $\langle u, q \rangle$ be a C -query. $\text{ANGERONA}(C, s, h, ATK, \langle u, q \rangle) = \text{ANGERONA}(C, s', h', ATK, \langle u, q \rangle)$ whenever $\text{knowledge}(h, u) = \text{knowledge}(h', u)$.

Proof. Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be an ATKLOG model, and $\langle u, q \rangle$ be a C -query. Furthermore, let s, s' be C -states and h, h' be C -histories such that $\text{knowledge}(h, u) = \text{knowledge}(h', u)$. Assume, for contradiction's sake, that $\text{ANGERONA}(C, s, h, ATK, \langle u, q \rangle) \neq \text{ANGERONA}(C, s', h', ATK, \langle u, q \rangle)$. Without loss of generality, assume that $\text{ANGERONA}(C, s, h, ATK, \langle u, q \rangle) = \top$ and $\text{ANGERONA}(C, s', h', ATK, \langle u, q \rangle) = \perp$. This happens iff either the result of *pox* or *secure* is different in the two cases. The results of *pox* and *secure*, however, depend just on the attacker's initial beliefs, the query q , and the set of formulae derive $\text{knowledge}(h, u)$. From this and $\text{knowledge}(h, u) = \text{knowledge}(h', u)$, it follows that $\text{ANGERONA}(C, s, h, ATK, \langle u, q \rangle) = \text{ANGERONA}(C, s', h', ATK, \langle u, q \rangle)$. \square

We now prove a key result for our security proof, namely that considering only the sentences in $\text{knowledge}(h, u)$ is enough to determine secrecy-preservation.

Proposition C.2. Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $u \in \mathcal{U}$ be a user, and $r = \langle s, h \rangle$ be a (C, f) -run. The following fact holds: $\llbracket r \rrbracket_{\sim_u} = \{db \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in \text{knowledge}(h, u)} [\phi]^{db} = \top\}$.

Proof. (\Rightarrow). Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be an ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $u \in \mathcal{U}$ be a user, and $r = \langle \langle db, U, P \rangle, h \rangle$ be a (C, f) -run. Furthermore, let $r' = \langle \langle db', U', P' \rangle, h' \rangle$ be a run in $\llbracket r \rrbracket_{\sim_u}$. From this, it follows that $h|_u = h'|_u$. From this, it follows that all queries in $\text{knowledge}(h, u)$ have the same result in db and db' . Therefore, $db' \in \{db \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in \text{knowledge}(h, u)} [\phi]^{db} = \top\}$.

(\Leftarrow). Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be an ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $u \in \mathcal{U}$ be a user, and $r = \langle \langle db, U, P \rangle, h \rangle$ be a (C, f) -run. Furthermore, let db' be a database state in $\{db' \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in \text{knowledge}(h, u)} [\phi]^{db'}\}$. We now construct a C -state s and an history h' such that (1) $\langle s, h' \rangle$ is a run, (2) $s.db = db'$, and (3) $\langle s, h' \rangle \in \llbracket r \rrbracket_{\sim_u}$. The C -state s is defined as $\langle db', U, P \rangle$, whereas the history $h' = h|_u$. We claim that $\langle s, h' \rangle$ is a run. From this and $h' = h|_u$, it follows that $\langle s, h' \rangle \sim_u r$. From this, it follows that $db' \in \llbracket r \rrbracket_{\sim_u}$.

To prove our claim that $\langle \langle db', U, P \rangle, h|_u \rangle$ is a run, we prove the stronger fact that for all $0 \leq i \leq |h|$, $\langle \langle db', U, P \rangle, h^i|_u \rangle$ is a run. We prove this by induction on i .

Base case. The base case $\langle \langle db', U, P \rangle, h^0|_u \rangle$ is trivial since $db' \in \Omega_D^\Gamma$ and $h^0|_u = \epsilon$, therefore $\langle \langle db', U, P \rangle, h^0|_u \rangle$ is a run.

Induction Step. Assume that $\langle \langle db', U, P \rangle, h^{i-1}|_u \rangle$ is a run. We now show that $\langle \langle db', U, P \rangle, h^i|_u \rangle$ is a run as well. Let $\langle \langle u', q \rangle, a, res \rangle$ be the last C -event in h^i . There are two cases:

- $u' \neq u$. From this, $h^i|_u = h^{i-1}|_u$. Therefore, $\langle \langle db', U, P \rangle, h^{i-1}|_u \rangle = \langle \langle db', U, P \rangle, h^i|_u \rangle$ and our claim directly follows from the induction hypothesis.
- $u' = u$. From this, $h^i|_u = h^{i-1}|_u \cdot \langle \langle u, q \rangle, a, res \rangle$. Assume, for contradiction's sake, that $\langle \langle db', U, P \rangle, h^i|_u \rangle$ is not a run. From $\langle \langle db', U, P \rangle, h^{i-1}|_u \rangle$ is a run (which directly follows from the induction hypothesis), it follows that there are only three cases:
 - 1) $f(\langle \langle db', U, P \rangle, \langle u, q \rangle, h^{i-1}|_u \rangle) \neq a$. Since a is the security decision associated with the last event in h^i , it follows that $a = f(\langle \langle db, U, P \rangle, \langle u, q \rangle, h^{i-1} \rangle)$. From this, it follows that $f(\langle \langle db, U, P \rangle, \langle u, q \rangle, h^{i-1} \rangle) \neq f(\langle \langle db', U, P \rangle, \langle u, q \rangle, h^{i-1}|_u \rangle)$. From knowledge 's definition, it follows that $\text{knowledge}(h^{i-1}, u) = \text{knowledge}(h^{i-1}|_u, u)$. From this and Proposition C.1, it follows that $f(\langle \langle db, U, P \rangle, \langle u, q \rangle, h^{i-1} \rangle) = f(\langle \langle db', U, P \rangle, \langle u, q \rangle, h^{i-1}|_u \rangle)$, leading to a contradiction.
 - 2) $a = \perp$ but $res \neq \dagger$. This contradicts the fact that the history is derived from h , which comes from a proper run.
 - 3) $a = \top$ but $res \neq [q]^{db'}$. From $res \neq [q]^{db'}$ and $res = [q]^{db}$ (since res comes from the run r), it follows that $[q]^{db} \neq [q]^{db'}$. There are two cases:
 - $[q]^{db} = \top$. From this and the definition of knowledge , it follows that $q \in \text{knowledge}(h, u)$. Therefore, $[q]^{db'} = \top$ follows from $db' \in \{db' \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in \text{knowledge}(h, u)} [\phi]^{db'}\}$, leading to a contradiction.
 - $[q]^{db} = \perp$. From this and the definition of knowledge , it follows that $\neg q \in \text{knowledge}(h, u)$. From this and $db' \in \{db' \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in \text{knowledge}(h, u)} [\phi]^{db'}\}$, it follows that $[\neg q]^{db'} = \top$. From this, $[q]^{db'} = \perp = [q]^{db}$, leading to a contradiction.

This completes the proof of our claim. \square

We now prove that ANGERONA provides the desired security guarantees.

Proposition C.3. Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , and $ATK' = \lambda u \in \mathcal{U}. \llbracket ATK(u) \rrbracket_D$ be the (C, f) -attacker model associated to ATK . The PDP shown in Algorithm 4, parametrized with ATK , provides data confidentiality with respect to C and ATK' .

Proof. Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , and $ATK' = \lambda u \in \mathcal{U}. \llbracket ATK(u) \rrbracket_D$ be the (C, f) -attacker model associated to ATK . Furthermore, let f be the PDP shown in Algorithm 4. Assume, for contradiction's sake, that f does not provide confidentiality with respect to C and ATK' . From this, it follows that there is a run $r = \langle \langle db, U, P \rangle, h \rangle$, a user $u \in \mathcal{U}$, a secret $\langle u, \phi, l \rangle \in P$, and a $0 \leq i \leq |h| - 1$ such that $\llbracket ATK' \rrbracket(u, r^i)([\phi]) < l$ and $\llbracket ATK' \rrbracket(u, r^{i+1})([\phi]) \geq l$. As a result, the

i -th C -query is the one that leaked information. There are two cases:

- The C -query is $\langle u, q \rangle$, for some q . There are three cases:
 - The result of the PDP is \top and the query q holds in the current state. From this, it follows that there is a database state (namely, the one in r) where the query q holds and the database is consistent with $knowledge(h^i, u)$. From this, it follows that $\llbracket ATK'(u, r^i) \rrbracket(\llbracket q \rrbracket) > 0$. From this and Proposition C.5, it follows that $pox(C, ATK, h^i, \langle u, q \rangle, \top)$ returns true. From this and the fact that the query has been authorized, it follows that $secure(C, ATK, h^{i+1}, \langle u, \phi, l \rangle)$ returned true. From this, the fact that r is actually a run (since q is satisfied in $r.db$), and Proposition C.4, it follows that $\llbracket ATK'(u, r^{i+1}) \rrbracket(\llbracket \phi \rrbracket) < l$, leading to a contradiction.
 - The result of the PDP is \top and the query q does not hold in the current state. From this, it follows that there is a database state (namely, the one in r) where the query $\neg q$ holds and the database is consistent with $knowledge(h^i, u)$. From this, it follows that $\llbracket ATK'(u, r^i) \rrbracket(\llbracket \neg q \rrbracket) > 0$. From this and Proposition C.5, it follows that $pox(C, ATK, h^i, \langle u, q \rangle, \perp)$ returns true. From this and the fact that the query has been authorized, it follows that $secure(C, ATK, h^{i+1}, \langle u, \phi, l \rangle)$ returned true. From this, the fact that r is actually a run (since q is satisfied in $r.db$), and Proposition C.4, it follows that $\llbracket ATK'(u, r^{i+1}) \rrbracket(\llbracket \phi \rrbracket) < l$, leading to a contradiction.
 - The result of the PDP is \perp (namely the query is not authorized). From this and Proposition C.1, it follows that for any run $r' \sim_u r^i$ the PDP result is the same. Therefore, $\llbracket r^i \rrbracket_{\sim_u} = \llbracket r^{i+1} \rrbracket_{\sim_u}$. From this, $\llbracket ATK'(u, r^i) \rrbracket(\llbracket \phi \rrbracket) < l$, and ATKLOG semantics, it follows that $\llbracket ATK'(u, r^{i+1}) \rrbracket(\llbracket \phi \rrbracket) < l$, leading to a contradiction.
- The C -query is $\langle u', q \rangle$, where $u' \neq u$ and q is a query. From this, it follows that $\llbracket ATK'(u, r^i) \rrbracket(\llbracket \phi \rrbracket) = \llbracket ATK'(u, r^{i+1}) \rrbracket(\llbracket \phi \rrbracket)$ since u 's belief does not change in response to a query from another user (since $r^i|_u = r^{i+1}|_u$). From this and $\llbracket ATK'(u, r^i) \rrbracket(\llbracket \phi \rrbracket) < l$, it follows that both $\llbracket ATK'(u, r^{i+1}) \rrbracket(\llbracket \phi \rrbracket) < l$ and $\llbracket ATK'(u, r^{i+1}) \rrbracket(\llbracket \phi \rrbracket) \geq l$, leading to a contradiction.

Since all cases ended in contradiction, this completes the proof of our claim. \square

Proposition C.4. *Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $ATK' = \lambda u \in \mathcal{U}. \llbracket ATK(u) \rrbracket_D$ be the (C, f) -attacker model associated to ATK , r be a run in $runs(C, f)$, and $\langle u, \psi, l \rangle$ be a secret in $r.S$. Then, $secure(C, ATK, h, \langle u, \psi, l \rangle)$ returns true iff $\llbracket ATK'(u, \langle s, h \rangle) \rrbracket(\llbracket \psi \rrbracket) < l$, for any state s that is compatible with h .*

Proof. Let $r = \langle s, h \rangle$ be a run such that s is compatible with h . The $secure$ procedure returns $\llbracket ATK(u) \rrbracket_D(\llbracket \psi \rrbracket) \cap_{\phi \in knowledge(h, u)} \llbracket \phi \rrbracket < l$. From

Proposition C.2, it follows that $\llbracket r \rrbracket_{\sim_u} = \{db \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in knowledge(h, u)} \llbracket \phi \rrbracket^{db} = \top\} = \bigcap_{\phi \in knowledge(h, u)} \llbracket \phi \rrbracket$. We can therefore rewrite $secure$'s result as $\llbracket ATK(u) \rrbracket_D(\llbracket \psi \rrbracket) \cap_{\langle s, h \rangle \sim_u} < l$. This is exactly the definition of $\llbracket ATK'(u, \langle s, h \rangle) \rrbracket(\llbracket \psi \rrbracket) < l$. \square

Proposition C.5. *Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $ATK' = \lambda u \in \mathcal{U}. \llbracket ATK(u) \rrbracket_D$ be the (C, f) -attacker model associated to ATK , r be a run in $runs(C, f)$, and $\langle u, q \rangle$ be a query. Then, $pox(C, ATK, h, \langle u, q \rangle, v)$ returns true iff $\llbracket ATK'(u, \langle s, h \rangle) \rrbracket(\llbracket r(q, v) \rrbracket) > 0$ for any state s that is compatible with h , where $r(q, \top) = q$ and $r(q, \perp) = \neg q$.*

Proof. Let $r = \langle s, h \rangle$ be a run such that s is compatible with h . The pox procedure returns $\llbracket ATK(u) \rrbracket_D(\llbracket r(q, v) \rrbracket) \cap_{\phi \in knowledge(h, u)} \llbracket \phi \rrbracket > 0$. From Proposition C.2, it follows that $\llbracket r \rrbracket_{\sim_u} = \{db \in \Omega_D^\Gamma \mid \bigwedge_{\phi \in knowledge(h, u)} \llbracket \phi \rrbracket^{db} = \top\} = \bigcap_{\phi \in knowledge(h, u)} \llbracket \phi \rrbracket$. We can therefore rewrite pox 's result as $\llbracket ATK(u) \rrbracket_D(\llbracket r(q, v) \rrbracket) \cap_{\langle s, h \rangle \sim_u} < l$. This is exactly the definition of $\llbracket ATK'(u, \langle s, h \rangle) \rrbracket(\llbracket r(q, v) \rrbracket) > 0$. \square

C. Complexity Proof

First, we formalize literal queries, a fragment of relational calculus that can be composed with relaxed acyclic programs without modifying acyclicity. Afterwards, we prove that for acyclic ATKLOG models and literal queries, ANGERONA has PTIME data complexity.

A literal query is a quantifier-free relational calculus (Σ, \mathbf{dom}) -formula either of the form $R(\bar{c})$ or $\neg R(\bar{c})$, where $R \in \Sigma$ and $\bar{c} \in \mathbf{dom}^{|\mathbf{R}|}$. We say that a literal query $R(\bar{c})$ (or $\neg R(\bar{c})$) is *compatible with a relaxed acyclic program p* (with witness $(\mathcal{O}, \mathcal{D}, \mathcal{U})$) iff $\mu_{p, \mathcal{D}, \mathcal{U}}(R) \neq \emptyset$ implies $\bar{c} \downarrow_K \in pdom(R, p, \mu_{p, \mathcal{D}, \mathcal{U}}(R))$. We now show that literal queries can be composed with acyclic PROBLOG programs without introducing cycles.

Proposition C.6. *Let $D = \langle \Sigma, \mathbf{dom} \rangle$ be a database schema, p be a relaxed acyclic program, $(\mathcal{O}, \mathcal{D}, \mathcal{U})$ be a witness for p , R be a predicate symbol in Σ , $\bar{c} \in \mathbf{dom}^{|\mathbf{R}|}$ be a tuple, and ϕ be a boolean D -query. If ϕ is a literal query compatible with p , then $p \cup PL(\phi)$ is a relaxed acyclic PROBLOG program as well.*

Proof. Let $D = \langle \Sigma, \Gamma \rangle$ be a database schema, p be a relaxed acyclic program, and ϕ be a boolean D -query. Furthermore, let ϕ be a literal query. Assume, for contradiction's sake, that $p \cup PL(\phi)$ is not a relaxed acyclic program. This happens iff the program $\alpha(\beta(p \cup PL(\phi)))$ is not an acyclic program. This happens iff the program $\alpha(\beta(p \cup PL(\phi)))$ contains an unguarded unsafe structure S in the dependency graph of $\alpha(\beta(p \cup PL(\phi)))$. The only interesting case is when S contains rules from $PL(\phi)$. If $\phi = R(\bar{c})$, then $PL(\phi)$ contains a single rule $fresh \leftarrow R(\bar{c})$, where $fresh$ is a fresh predicate symbol. Similarly, if $\phi = \neg R(\bar{c})$, then $PL(\phi)$ contains two rules $fresh_1 \leftarrow fresh_2$ and $fresh_2 \leftarrow \neg R(\bar{c})$, where $fresh_1$ and $fresh_2$ are fresh predicate symbols. From this and the fact that

all rules in $PL(\phi)$ are both strongly and weakly connected for \mathcal{U} , it follows that there is always unguarded unsafe structure S' that can be obtained from S by removing the rules in $PL(\phi)$. This contradicts the fact that p is a relaxed acyclic program. \square

Here, we show that for acyclic ATKLOG models, ANGERONA has PTIME data complexity.

Proposition C.7. *Let C be a system configuration and f be the PDP shown in Figure 4. For any acyclic C -ATKLOG model ATK , any action $\langle u, q \rangle$ such that q is a literal query compatible with p , and any run $r \in \text{runs}(C, f)$ such that (1) all sentences in $\text{knowledge}(r.h, u)$, for any $u \in r.U$, are literal queries compatible with p , and (2) all secrets in $r.S$ are literal queries compatible with p , the algorithm shown in Figure 4 has PTIME data complexity.*

Proof. Let B be the PROBLG program obtained from ATK . The data complexity of f is its complexity when only the database $r.db$ and $edb(B)$ change. The algorithm in Figure 4 calls three times the *secure* procedure and twice the *pox* procedure. The set $\text{knowledge}(h, u)$ can be constructed in $O(|h|)$. From this, it follows that the program p has size $O(|edb(B)| + |\text{rules}(B)| + |h|)$. Furthermore, since all queries are literal queries, B is a relaxed acyclic PROBLG program, and Proposition C.6, it follows that p is a relaxed acyclic PROBLG program as well. From this and Proposition B.18, the inference can be performed in $O(r^2 \cdot e^r \cdot \max(2, r)^{2+r})$, where $r \in O(|\text{rules}(B)| + (|r| + |\text{edb}(K)|)^{|\text{rules}(K)|})$ and $e \in O(|edb(B)|)$. Since $\text{rules}(B)$, r , and K are fixed, then the data complexity of the *secure* procedure is $O(|edb(B)|^k)$, for some $k \in \mathbb{N}$. From this and the fact that Algorithm 4 three times the *secure* procedure and twice the *pox* procedure per secret in $r.S$, it follows that the data complexity of Algorithm 4 is PTIME. \square

D. Completeness Proof

We first introduce the notion of unconditionally secrecy-preserving query. Informally, a query $\langle u', q \rangle$ is unconditionally secrecy-preserving given a run r and a secret $\langle u, \psi, l \rangle$ iff disclosing the result of $\langle u', q \rangle$ in any run $r' \sim_u r$ does not violate the secret.

Definition C.1. Let $C = \langle D, \Gamma \rangle$ be a configuration, f be a C -PDP, and ATK be a (C, f) -attacker model. A query q is *unconditionally secrecy-preserving* for a (C, f) -run r , a user u , a secret $\langle u, \psi, l \rangle$, and ATK iff $\llbracket ATK \rrbracket(u, r)(\psi) < l$ implies that (a) if $\llbracket ATK \rrbracket(u, r)(q) > 0$, then $\llbracket ATK \rrbracket(u, r)(\psi \mid \{db \in \llbracket r \rrbracket_{\sim_u} \mid [q]^{db} = \top\}) < l$, and (b) if $\llbracket ATK \rrbracket(u, r)(\neg q) > 0$, then $\llbracket ATK \rrbracket(u, r)(\psi \mid \{db \in \llbracket r \rrbracket_{\sim_u} \mid [q]^{db} = \perp\}) < l$. \square

We say that a PDP is *complete* if it authorizes all unconditionally secrecy-preserving queries. ANGERONA is complete. This directly follows from (1) Proposition C.2, (2) the use of exact inference procedures for PROBLG programs, and (3) the fact that ANGERONA directly checks whether queries are unconditionally secrecy-preserving or not.

Proposition C.8. *Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $ATK' = \lambda u \in \mathcal{U}. \llbracket ATK(u) \rrbracket_D$ be the (C, f) -attacker model associated to ATK , $r = \langle \langle db, U, P \rangle, h \rangle$ be a run, and $\langle u, q \rangle$ be a C -query. If q is unconditionally secrecy-preserving for u , for all secrets $\langle u, \psi, l \rangle \in \text{secrets}(P, u)$, r , and f authorizes $\langle u, q \rangle$.*

Proof. Let $C = \langle D, \Gamma \rangle$ be a system configuration, ATK be a C -ATKLOG model, f be the C -PDP obtained by parametrizing Algorithm 4 with C and ATK , $ATK' = \lambda u \in \mathcal{U}. \llbracket ATK(u) \rrbracket_D$ be the (C, f) -attacker model associated to ATK , $r = \langle \langle db, U, P \rangle, h \rangle$ be a run, and $\langle u, q \rangle$ be a C -query. Assume, for contradiction's sake, that our claim does not hold. Namely, there is query $\langle u, q \rangle$ such that (1) the query is unconditionally secrecy-preserving u , for all secrets $\langle u, \psi, l \rangle \in \text{secrets}(P, u)$, r , and ATK' , and (2) ANGERONA (parametrized with ATK) does not authorize $\langle u, q \rangle$. Since ANGERONA does not authorize $\langle u, q \rangle$, it means that $\text{ANGERONA}(C, \langle db, U, P \rangle, h, ATK, \langle u, q \rangle) = \perp$. From this, it follows that there is a secret $\langle u, \psi, l \rangle \in \text{secrets}(P, u)$ such that:

- $\text{secure}(C, ATK, h, \langle u, \psi, l \rangle) = \top$, and
- one of the two cases hold:
 - $\text{pox}(C, ATK, h, \langle u, q \rangle, \top) = \top$ and $\text{secure}(C, ATK, h', \langle u, \psi, l \rangle) = \perp$, where $h' = h \cdot \langle \langle u, q \rangle, \top, \text{top} \rangle$, or
 - $\text{pox}(C, ATK, h, \langle u, q \rangle, \perp) = \top$ and $\text{secure}(C, ATK, \langle db, U, P \rangle, h'', \langle u, \psi, l \rangle) = \perp$, where $h'' = h \cdot \langle \langle u, q \rangle, \top, \perp \rangle$.

Without loss of generality, we assume that $\text{pox}(C, ATK, h, \langle u, q \rangle, \top) = \top$ and $\text{secure}(C, ATK, h' \langle u, \psi, l \rangle) = \perp$, where $h' = h \cdot \langle \langle u, q \rangle, \top, \text{top} \rangle$ (the proof for the other case is identical). From $\text{secure}(C, ATK, h, \langle u, \psi, l \rangle) = \top$, Proposition C.4, and the fact that s is compatible with h , it follows that $\llbracket ATK' \rrbracket(u, r)(\psi) < l$. From $\text{pox}(C, ATK, h, \langle u, q \rangle, \top) = \top$, Proposition C.4, and the fact that s is compatible with h , it follows that $\llbracket ATK' \rrbracket(u, r)(q) > 0$. From this and Proposition C.1, it follows that there is a system state $s' = \langle db', U, P \rangle$ such that (1) q holds in db' , (2) s' is compatible with $h|_u$, and (3) $r' = \langle s', h|_u \rangle$ is a run indistinguishable from r . From this, $\text{secure}(C, ATK, h', \langle u, \psi, l \rangle) = \perp$, $\text{secure}(C, ATK, h', \langle u, \psi, l \rangle) = \text{secure}(C, ATK, \langle db, U, P \rangle, h'|_u, \langle u, \psi, l \rangle)$, and Proposition C.4, it follows that $\llbracket ATK' \rrbracket(u, \langle s', h'|_u \rangle)(\psi) \geq l$. From this and the definition of $h'|_u$, it follows that $\llbracket ATK' \rrbracket(u, r')(\psi \mid \{db \in \llbracket r \rrbracket_{\sim_u} \mid [q]^{db} = \top\}) \geq l$. From this, $r \sim_u r'$, and $\llbracket ATK' \rrbracket(u, r) = \llbracket ATK' \rrbracket(u, r')$ if $r \sim_u r'$, it follows that $\llbracket ATK' \rrbracket(u, r)(\psi \mid \{db \in \llbracket r \rrbracket_{\sim_u} \mid [q]^{db} = \top\}) \geq l$. Therefore, we have that $\llbracket ATK' \rrbracket(u, r)(\psi) < l$, $\llbracket ATK' \rrbracket(u, r)(q) > 0$, and $\llbracket ATK' \rrbracket(u, r)(\psi \mid \{db \in \llbracket r \rrbracket_{\sim_u} \mid [q]^{db} = \top\}) \geq l$. This contradicts the fact that q is an unconditionally secrecy-preserving query and completes the proof of our claim. \square