

# Provenance as Dependency Analysis

James Cheney<sup>1</sup>, Amal Ahmed<sup>2</sup>, and Umut A. Acar<sup>2</sup> \*

<sup>1</sup> University of Edinburgh

<sup>2</sup> Toyota Technological Institute at Chicago

**Abstract.** Provenance is information recording the source, derivation, or history of some information. Provenance tracking has been studied in a variety of settings; however, although many design points have been explored, the mathematical or semantic foundations of data provenance have received comparatively little attention. In this paper, we argue that dependency analysis techniques familiar from program analysis and program slicing provide a formal foundation for forms of provenance that are intended to show how (part of) the output of a query depends on (parts of) its input. We introduce a semantic characterization of such *dependency provenance*, show that this form of provenance is not computable, and provide dynamic and static approximation techniques.

## 1 Introduction

*Provenance* is information about the origin, ownership, influences upon, or other *historical* or *contextual* information about an object. Such information has many applications, including evaluating integrity or authenticity claims, establishing the chain of custody of or responsibility for an object, detecting and repairing errors, and memoization and caching of the results of computations. Provenance is particularly important in scientific computation and recordkeeping, since it is considered essential for ensuring the repeatability of experiments and judging the scientific value of their results.

Provenance tracking has been studied in a variety of settings, including databases, file systems, and scientific workflows; indeed, many familiar systems provide simple forms of provenance, such as the timestamp and ownership metadata in file systems, system logs, and version control systems. Although a wide variety of design points have been explored [8, 17, 21], there is relatively little understanding of the relationships among techniques or of the design considerations that should be taken into account when developing or evaluating an approach to provenance. The mathematical or semantic foundations of data provenance have received comparatively little attention, with a few relatively recent exceptions [10, 11, 15, 19].

Most prior approaches have invoked intuitive concepts such as *contribution*, *influence*, and *relevance* as motivation for their definitions of provenance. These intuitions seem adequate when considering monotone relational queries, but tend to break down when negation, grouping, or aggregation are considered. These intuitions have also motivated rigorous approaches to seemingly quite different problems, such as aiding debugging via *program slicing* [7, 16, 24], supporting efficient memoization and

---

\* Cheney was supported by EPSRC grant R37476. Acar was supported by an Intel gift.

(a)

Protein				EnzReact		Reaction		
ID	Name	MW	...	PID	RID	ID	Name	...
<b>p<sub>1</sub></b>	thioredoxin	<b>11.8</b>	...	<b>p<sub>1</sub></b>	<b>r<sub>1</sub></b>	<b>r<sub>1</sub></b>	thia-phos + ATP = thi diphos + ADP	...
<b>p<sub>2</sub></b>	flavodoxin	19.7	...	<b>p<sub>2</sub></b>	<b>r<sub>1</sub></b>	<b>r<sub>2</sub></b>	H <sub>2</sub> O + an acyl phos → phos + a carboxylate	...
<b>p<sub>3</sub></b>	ferredoxin	12.3	...	<b>p<sub>1</sub></b>	<b>r<sub>2</sub></b>	<b>r<sub>3</sub></b>	D-ribose-5-phosp = D-ribulose-5-phos	...
<b>p<sub>4</sub></b>	ArgR	<b>-700</b>	...	<b>p<sub>4</sub></b>	<b>r<sub>2</sub></b>	<b>r<sub>4</sub></b>	<i>β</i> -D-gluc-6-phos = fruct-6-phos	...
<b>p<sub>5</sub></b>	CheW	18.1	...	<b>p<sub>5</sub></b>	<b>r<sub>3</sub></b>	<b>r<sub>5</sub></b>	panteth 4'-phos + ATP = dephos-CoA + diphos	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

(b)

```

SELECT R.Name as Name, AVERAGE(P.MW) as AvgMW
  FROM Protein P, EnzymaticReaction ER, Reaction R
 WHERE P.ID = ER.ProteinID, ER.ReactionID = R.ID
 GROUP BY R.Name

```

(c)

Name	AvgMW
thia-phos + ATP = thi diphos + ADP	15.75
H <sub>2</sub> O + an acyl phos → phos + a carboxylate	<b>-338.2</b>
D-ribose-5-phosp = D-ribulose-5-phos	18.1
⋮	⋮

**Fig. 1.** Example (a) input, (b) query, and (b) output data; input field values relevant to *italicized* erroneous output value are highlighted in **bold**

caching [2, 4], and improving program security using information flow analysis [20]. As Abadi et al. have argued [1], slicing, information flow, and several other program analysis techniques can all be understood in terms of dependence.

In this paper, we argue that these *dependency analysis* and *slicing* techniques familiar from programming languages provide a suitable foundation for an interesting class of provenance techniques. To illustrate our approach, consider the (realistic, but biologically inaccurate) input data shown in Figure 1(a) and the query in Figure 1(b) which calculates the average molecular weights of proteins involved in each reaction. The result of this query is shown in Figure 1(c).

Since the MW field contains the molecular weight of a protein, it is clearly an error for the italicized value in the result to be negative. To track down the source of the error, it would be helpful to know which parts of the input *contributed to*, or were *relevant to*, the erroneous part of the output. We can formalize this intuition by saying that a part of the output *depends on* a part of the input if a change to the input part may result in a change to the output part. This is analogous to program slicing [24], a debugging aid that identifies the parts of a *program* on which a program output depends.

In this example, the input field values that the erroneous output AvgMW-value depends on are highlighted in bold. The dependences include the two summed MW values

and the ID fields which are compared by the selection and grouping query. These ID fields must be included because a change to any one of them could result in a change to the italicized output value—for example, changing the occurrence of  $p_4$  in table `EnzymaticReaction`. On the other hand, the names of the proteins and reactions are *irrelevant* to the output `AvgMW`.

This example is simplistic, but the ability to concisely explain which parts of the input influence each part of the output is much more important if we consider a realistic database with perhaps tens or hundreds of columns per table and thousands or millions of rows. Moreover, dependence information can also be useful for a variety of other applications, including estimating the *quality* or *freshness* of data in a query result by aggregating timestamps or quality annotations on the relevant inputs.

In this paper, we argue that data dependence provides a solid semantic foundation for generalizing why-provenance, lineage, and similar techniques. We consider the full nested relational calculus with set difference, equality, grouping and aggregation operations, and functions on basic types. We consider annotation-propagating semantics for such queries and define a property called *dependency-correctness*, which, intuitively, means that the annotations produced by a query correctly show all parts of the output that may change if a part of the input is changed.

The structure of the rest of this paper is as follows. We briefly review the nested relational calculus in Section 2. We then introduce (in Section 3) the annotation-propagation model and define dependency-correctness. In Section 3.1 we describe a dynamic *provenance-tracking* semantics that is dependency-correct. We also (Section 3.2) introduce a type-based *provenance analysis* which is less accurate than provenance tracking, but can be performed statically; we also prove its correctness relative to dynamic provenance tracking. We discuss our results and a preliminary implementation in Section 3.3 and discuss related and future work and conclude in Sections 4–5. Full proofs of our main results are provided in a companion technical report [14].

## 2 Background

We assume some familiarity with the *nested relational calculus* (NRC) [13], which is closely related to *monad algebra* [22]. We consider multiset (or bag) collections rather than sets. The types and expressions of our variant of NRC are as follows:

$$\begin{aligned} \tau ::= & \text{bool} \mid \text{int} \mid \tau_1 \times \tau_2 \mid \{\tau\} \\ e ::= & x \mid \text{let } x = e_1 \text{ in } e_2 \mid (e_1, e_2) \mid \pi_i(e) \mid b \mid i \mid \neg e \mid e_1 \wedge e_2 \mid e_1 + e_2 \mid \text{sum}(e) \\ & \mid e_1 \approx e_2 \mid \text{if } e_0 \text{ then } e_1 \text{ else } e_2 \mid \emptyset \mid \{e\} \mid e_1 \cup e_2 \mid e_1 - e_2 \mid \{e_2 \mid x \in e_1\} \mid \bigcup e \end{aligned}$$

Here,  $i \in \mathbb{Z} = \{\dots, -1, 0, 1, \dots\}$  represents integer constants and  $b \in \mathbb{B} = \{\text{true}, \text{false}\}$  denotes Boolean constants. The type  $\{\tau\}$  describes *collections* of elements of type  $\tau$ ; in this paper, we consider collections to be *bags* (multisets). The bag operations include  $\emptyset$ , the constant empty bag; singletons  $\{e\}$ ; bag union and difference; comprehension  $\{e_2 \mid x \in e_1\}$ ; and flattening  $\bigcup e$ . By convention, we write  $\{e_1, \dots, e_n\}$  as syntactic sugar for  $\{e_1\} \cup \dots \cup \{e_n\}$ . Finally, we include `sum`, a typical aggregation operation, which adds together all of the elements of a bag and produces a value; e.g.

$\text{sum}\{1, 2, 3\} = 6$  (by convention,  $\text{sum}(\emptyset) = 0$ ). We syntactically distinguish between NRC’s equality operation  $\approx$  and mathematical equality  $=$ .

**Types** Query language expressions can be typechecked using standard techniques. Contexts  $\Gamma$  are lists of pairs of variables and types  $x_1 : \tau_1, \dots, x_n : \tau_n$ , where  $x_1, \dots, x_n$  are distinct. The rules for typechecking expressions are shown in Figure 2.

**Semantics** We write  $\mathcal{M}_{\text{fin}}(X)$  for the set of all finite bags with elements drawn from  $X$ . The (standard) interpretation of base types as sets of values is as follows:

$$\begin{aligned} \mathcal{T}[\text{bool}] &= \mathbb{B} = \{\text{true}, \text{false}\} & \mathcal{T}[\tau_1 \times \tau_2] &= \mathcal{T}[\tau_1] \times \mathcal{T}[\tau_2] \\ \mathcal{T}[\text{int}] &= \mathbb{Z} = \{\dots, -1, 0, 1, \dots\} & \mathcal{T}[\{\tau\}] &= \mathcal{M}_{\text{fin}}(\mathcal{T}[\tau]) \end{aligned}$$

An environment  $\gamma$  is a function from variables to values. We define the set of environments matching context  $\Gamma$  as  $\mathcal{T}[\Gamma] = \{\gamma \mid \forall x \in \text{dom}(\Gamma). \gamma(x) \in \mathcal{T}[\Gamma(x)]\}$ .

Figure 3 gives the semantics of queries. Note that we overload notation for pair projection  $\pi_i$  and bag operations such as  $\cup$  and  $\sqcup$ ; also, if  $S$  is a bag of integers, then  $\sum S$  is the sum of their values (taking  $\sum \emptyset = 0$ ). It is straightforward to show that

**Lemma 1.** *If  $\Gamma \vdash e : \tau$  then  $\mathcal{E}[e] : \mathcal{T}[\Gamma] \rightarrow \mathcal{T}[\tau]$ .*

As discussed in previous work [13], the NRC can express a wide variety of queries including ordinary relational queries as well as grouping and aggregation. We do not consider incomplete information (NULL values). Additional primitive functions and relations such as average can also be added without difficulty (in fact, average is definable in our language). For example, using more readable named record, comprehensions, and patterns, the SQL query from the Figure 1(b) can be defined as

$$\text{let } X = \{(r.\text{Name}, p.\text{MW}) \mid r \in R, er \in ER, p \in P, er.\text{RID} = r.\text{ID}, p.\text{ID} = er.\text{PID}\} \text{ in} \\ \{(n, \text{average}\{mw \mid (n', mw) \in X, n = n'\}) \mid (n, \_) \in X\}$$

Additional examples are shown in Figure 4.

### 3 Annotations, Provenance and Dependence

We wish to define *dependency provenance* as information relating each part of the output of a query to a set of parts of the input on which the output part depends. Collection types such as sets and bags are unordered and lack a natural way to “address” parts of values, so we must introduce one. One technique (familiar from many program analyses as well as other work on provenance [10, 23]) is to enrich the data model with *annotations* that can be used to refer to parts of the value. We can then infer provenance information from functions on annotated values by observing how such functions propagate annotations; conversely, we can define provenance-tracking semantics by enriching ordinary functions with annotation-propagation behavior. In this section, we show how to define dependency provenance in this way. In the next sections, we will show how to compute dynamic and static dependency provenance for NRC queries using annotations.

We define *annotated values* (*a-values*)  $v$ , *raw values* (*r-values*)  $w$ , and *multisets of annotated values*  $V$  as follows:

$$v ::= w^\Phi \quad w ::= i \mid b \mid (v_1, v_2) \mid V \quad V ::= \{v_1, \dots, v_n\}$$

$$\boxed{\Gamma \vdash e : \tau}$$

$$\frac{x:\tau \in \Gamma \quad \Gamma \vdash e_1 : \tau_1 \quad \Gamma, x:\tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash x : \tau} \quad \frac{i \in \mathbb{Z}}{\Gamma \vdash i : \text{int}} \quad \frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1 + e_2 : \text{int}}$$

$$\frac{\Gamma \vdash e : \{\text{int}\}}{\Gamma \vdash \text{sum}(e) : \text{int}} \quad \frac{b \in \mathbb{B}}{\Gamma \vdash b : \text{bool}} \quad \frac{\Gamma \vdash e_0 : \text{bool} \quad \Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } e_0 \text{ then } e_1 \text{ else } e_2 : \tau} \quad \frac{\Gamma \vdash e : \text{bool}}{\Gamma \vdash \neg e : \text{bool}}$$

$$\frac{\Gamma \vdash e_1 : \text{bool} \quad \Gamma \vdash e_2 : \text{bool}}{\Gamma \vdash e_1 \wedge e_2 : \text{bool}} \quad \frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash (e_1, e_2) : \tau_1 \times \tau_2} \quad \frac{\Gamma \vdash e : \tau_1 \times \tau_2}{\Gamma \vdash \pi_i(e) : \tau_i}$$

$$\frac{\Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash e_1 \approx e_2 : \text{bool}} \quad \frac{\Gamma \vdash e : \tau}{\Gamma \vdash \emptyset : \{\tau\}} \quad \frac{\Gamma \vdash e_1 : \{\tau\} \quad \Gamma \vdash e_2 : \{\tau\}}{\Gamma \vdash e_1 \cup e_2 : \{\tau\}}$$

$$\frac{\Gamma \vdash e_1 : \{\tau\} \quad \Gamma \vdash e_2 : \{\tau\}}{\Gamma \vdash e_1 - e_2 : \{\tau\}} \quad \frac{\Gamma \vdash e_1 : \{\tau_1\} \quad \Gamma, x:\tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash \{e_2 \mid x \in e_1\} : \{\tau_2\}} \quad \frac{\Gamma \vdash e : \{\tau\}}{\Gamma \vdash \bigcup e : \{\tau\}}$$

**Fig. 2.** Well-formed query expressions

$$\begin{array}{ll} \mathcal{E}[\![x]\!] \gamma = \gamma(x) & \mathcal{E}[\![\text{let } x = e_1 \text{ in } e_2]\!] \gamma = \mathcal{E}[\![e_2]\!] \gamma[x \mapsto \mathcal{E}[\![e_1]\!] \gamma] \\ \mathcal{E}[\![i]\!] \gamma = i & \mathcal{E}[\![e_1 + e_2]\!] \gamma = \mathcal{E}[\![e_1]\!] \gamma + \mathcal{E}[\![e_2]\!] \gamma \\ \mathcal{E}[\![\text{sum}(e)]\!] \gamma = \sum \mathcal{E}[\![e]\!] \gamma & \mathcal{E}[\![b]\!] \gamma = b \\ \mathcal{E}[\![\neg e]\!] \gamma = \neg \mathcal{E}[\![e]\!] \gamma & \mathcal{E}[\![e_1 \wedge e_2]\!] \gamma = \mathcal{E}[\![e_1]\!] \gamma \wedge \mathcal{E}[\![e_2]\!] \gamma \\ \mathcal{E}[\![(e_1, e_2)]\!] \gamma = (\mathcal{E}[\![e_1]\!] \gamma, \mathcal{E}[\![e_2]\!] \gamma) & \mathcal{E}[\![\pi_i(e)]\!] \gamma = \pi_i(\mathcal{E}[\![e]\!] \gamma) \\ \mathcal{E}[\![\emptyset]\!] \gamma = \emptyset & \mathcal{E}[\![\{e}\!] \!] \gamma = \{\mathcal{E}[\![e]\!] \gamma\} \\ \mathcal{E}[\![e_1 \cup e_2]\!] \gamma = \mathcal{E}[\![e_1]\!] \gamma \cup \mathcal{E}[\![e_2]\!] \gamma & \mathcal{E}[\![e_1 - e_2]\!] \gamma = \mathcal{E}[\![e_1]\!] \gamma - \mathcal{E}[\![e_2]\!] \gamma \\ \mathcal{E}[\![\bigcup e]\!] \gamma = \bigcup \mathcal{E}[\![e]\!] \gamma & \mathcal{E}[\![\{e \mid x \in e_0\}]\!] \gamma = \{\mathcal{E}[\![e]\!] \gamma[x \mapsto v] \mid v \in \mathcal{E}[\![e_0]\!] \gamma\} \end{array}$$

$$\mathcal{E}[\![\text{if } e_0 \text{ then } e_1 \text{ else } e_2]\!] \gamma = \begin{cases} \mathcal{E}[\![e_1]\!] \gamma & \text{if } \mathcal{E}[\![e_0]\!] \gamma = \text{true} \\ \mathcal{E}[\![e_2]\!] \gamma & \text{if } \mathcal{E}[\![e_0]\!] \gamma = \text{false} \end{cases}$$

$$\mathcal{E}[\![e_1 \approx e_2]\!] \gamma = \begin{cases} \text{true} & \text{if } \mathcal{E}[\![e_1]\!] \gamma = \mathcal{E}[\![e_2]\!] \gamma \\ \text{false} & \text{if } \mathcal{E}[\![e_1]\!] \gamma \neq \mathcal{E}[\![e_2]\!] \gamma \end{cases}$$

**Fig. 3.** Semantics of query expressions

$$\begin{aligned} \Pi_A(R) &= \{x.A \mid x \in R\} \\ \sigma_{A=B}(R) &= \bigcup \{\text{if } x.A = x.B \text{ then } \{x\} \text{ else } \emptyset \mid x \in R\} \\ R \times S &= \{(A : x.A, B : x.B, C : y.C, D : y.D, E : y.E) \mid x \in R, y \in S\} \\ \Pi_{BE}(\sigma_{A=D}(R \times S)) &= \{\text{if } x.A = y.D \text{ then } \{(B : x.B, E : y.E)\} \text{ else } \emptyset \mid x \in R, y \in S\} \\ R \cup \rho_{A/C, B/D}(\Pi_{CD}(S)) &= R \cup \{(A : y, C, B : y.D) \mid y \in S\} \\ R - \rho_{A/D, B/E}(\Pi_{DE}(S)) &= R - \{(A : y, D, B : y.E) \mid y \in S\} \\ \text{sum}(\Pi_A(R)) &= \text{sum}\{x.A \mid x \in R\} \\ \text{count}(R) &= \text{sum}\{1 \mid x \in R\} \end{aligned}$$

**Fig. 4.** Example queries

For us, annotations are *sets*  $\Phi \subseteq \text{Color}$  of values from some atomic data type  $\text{Color}$  of *colors*. We often omit set brackets in the annotations, for example writing  $w^{a,b,c}$  instead of  $w^{\{a,b,c\}}$  and  $w$  instead of  $w^\emptyset$ . An a-value  $v$  is said to be *distinctly colored* if every part of it is colored with a singleton set  $\{a\}$  and no color  $c$  is repeated anywhere in  $v$ .

For each type  $\tau$ , we define the set  $\mathcal{A}_0[\tau]$  of annotated values of type  $\tau$  as follows:

$$\begin{aligned} \mathcal{A}_0[\text{bool}] &= \{b^\Phi \mid b \in \mathbb{B}\} & \mathcal{A}_0[\tau_1 \times \tau_2] &= \{(v_1, v_2)^\Phi \mid v_1 \in \mathcal{A}_0[\tau_1], v_2 \in \mathcal{A}_0[\tau_2]\} \\ \mathcal{A}_0[\text{int}] &= \{i^\Phi \mid i \in \mathbb{Z}\} & \mathcal{A}_0[\{\tau\}] &= \{V^\Phi \mid \forall v \in V. v \in \mathcal{A}_0[\tau]\} \end{aligned}$$

Annotated environments  $\hat{\gamma}$  map variables to annotated values. We define the set of annotated environments matching context  $\Gamma$  as  $\mathcal{A}_0[\Gamma] = \{\hat{\gamma} \mid \forall x \in \text{dom}(\Gamma). \hat{\gamma}(x) \in \mathcal{A}_0[\Gamma(x)]\}$ .

We define an erasure function  $|\cdot|$ , mapping a-values to ordinary values (and, abusing notation, also mapping r-values to ordinary values), as follows:

$$|i| = i \quad |b| = b \quad |(v_1, v_2)| = (|v_1|, |v_2|) \quad |\{V\}| = \{|v| \mid v \in V\} \quad |w^\Phi| = |w|$$

and an annotation extraction function  $\|\cdot\|$  which extracts the set of all colors mentioned anywhere in an a-value or r-value, defined by taking  $\|w^\Phi\| = \Phi \cup \|w\|$  and

$$\|i\| = \emptyset \quad \|b\| = \emptyset \quad \|(v_1, v_2)\| = \|v_1\| \cup \|v_2\| \quad \|\{V\}\| = \bigcup \{|v| \mid v \in V\}$$

Two a-values are said to be *compatible* (written  $v \cong v'$ ) if  $|v| = |v'|$ .

We now consider annotated functions (a-functions)  $F : \mathcal{A}_0[\Gamma] \rightarrow \mathcal{A}_0[\tau]$  on a-values. Recall that we plan to define provenance for functions  $f : \mathcal{T}[\Gamma] \rightarrow \mathcal{T}[\tau]$  by observing how a-functions transform annotations. For this to make sense, we first need to restrict attention to a-functions  $F$  whose behavior is consistent with that of some function  $f$ ; that is, such that  $\forall v \in \mathcal{A}_0[\Gamma]. f(|v|) = |F(v)|$ . If this is the case, then we say that the a-function  $F$  is an *enrichment* of  $f$ ; there can be at most one such  $f$ , so we sometimes write  $|F|$  for  $f$ . Of course, many a-functions are not enrichments of any ordinary function: for example, suppose  $F_0(1^a) = 1^a$  while  $F_0(1^b) = 2^b$ . It may be of interest to semantically characterize the a-functions that are enrichments of ordinary functions, by analogy with generic queries in relational databases and color-invariance in [10]; while this would be important for studying expressiveness, in this paper we simply restrict attention to a-functions  $F$  that are enrichments of ordinary functions, that is, for which  $|F|$  exists.

**Dependency-correctness** Intuitively, an a-function  $F$  is *dependency-correct* if its output annotations tell us how changes to parts of the input may affect parts of the output. First, we need to capture the intuitive notion of changing a value at a particular location:

**Definition 1 (Equal except at  $c$ ).** *Two a-values  $v_1, v_2$  are equal except at  $c$  ( $v_1 \equiv_c v_2$ ) provided that they have the same structure except possibly at subterms labeled with  $c$ ; this relation is defined as follows:*

$$\frac{d \in \mathbb{B} \cup \mathbb{Z}}{d \equiv_c d} \quad \frac{v_1 \equiv_c v'_1 \quad v_2 \equiv_c v'_2}{(v_1, v_2) \equiv_c (v'_1, v'_2)} \quad \frac{v_1 \equiv_c v'_1 \quad \dots \quad v_n \equiv_c v'_n}{\{v_1, \dots, v_n\} \equiv_c \{v'_1, \dots, v'_n\}} \quad \frac{w_1 \equiv_c w_2}{w_1^\Phi \equiv_c w_2^\Phi} \quad \frac{c \in \Phi_1 \cap \Phi_2}{w_1^{\Phi_1} \equiv_c w_2^{\Phi_2}}$$

*Example 1.* Consider the two a-environments:

$$\begin{aligned}\hat{\gamma} &= (\mathbf{R} : \{(1^{c_1}, 3^{c_2}, 5^{c_3})^{b_1}, \dots\}^a, \mathbf{S} : \dots) \\ \hat{\gamma}' &= (\mathbf{R} : \{(2^{c_1}, 3^{c_2}, 5^{c_3})^{b_1}, \dots\}^a, \mathbf{S} : \dots)\end{aligned}$$

We have  $\hat{\gamma} \equiv_a \hat{\gamma}'$ ,  $\hat{\gamma} \equiv_{b_1} \hat{\gamma}'$ , and  $\hat{\gamma} \equiv_{c_1} \hat{\gamma}'$ , assuming that the elided portions are identical. For distinctly-colored values, a color serves as an address uniquely identifying a subterm. Thus,  $\equiv_c$  relates a distinctly-colored value to a value which can be obtained by modifying the subterm “at  $c$ ”; that is, if we write  $v_1$  as  $C[v'_1]$  where  $C$  is a context and  $v'_1$  is the subterm labeled with  $c$  in  $v_1$ , and  $v_1 \equiv_c v_2$ , then  $v_2 = C[v'_2]$  for some subterm  $v'_2$  labeled with  $c$ . Note that  $v'_2$  and  $v_2$  need not be distinctly colored, and that  $\equiv_c$  makes sense for arbitrary a-values, not just distinctly colored ones.

**Definition 2 (Dependency-correctness).** An a-function  $F : \mathcal{A}_0[\Gamma] \rightarrow \mathcal{A}_0[\tau]$  is dependency-correct if for any  $c \in \text{Color}$  and  $\hat{\gamma}, \hat{\gamma}' \in \mathcal{A}_0[\Gamma]$  satisfying  $\hat{\gamma} \equiv_c \hat{\gamma}'$ , we have  $F(\hat{\gamma}) \equiv_c F(\hat{\gamma}')$ .

*Example 2.* Recall  $\hat{\gamma}, \hat{\gamma}'$  as in the previous example. Suppose

$$F(\hat{\gamma}) = \{(1^{c_1}, 3^{c_2}, 5^{c_3})^{b_1}\}^a.$$

Since  $\hat{\gamma} \equiv_{c_1} \hat{\gamma}'$ , we know that  $F(\hat{\gamma}) \equiv_{c_1} F(\hat{\gamma}')$  so we can see that  $F(\hat{\gamma}')$  must be of the form

$$\{(s^{c_1}, 3^{c_2}, 5^{c_3})^{b_1}\}^a$$

for some  $n \in \mathbb{Z}$ . We do *not* necessarily know that  $n$  must be 2; this is not captured by dependency-correctness.

More generally, dependency-correctness tells us that for any  $c$ , we must have  $F(\hat{\gamma}) = C[v_1, \dots, v_n]$  and  $F(\hat{\gamma}') = C[v'_1, \dots, v'_n]$ , where  $C[-, \dots, -]$  is a context not mentioning  $c$  and  $v_1, \dots, v_n, v'_1, \dots, v'_n$  are labeled with  $c$ . Thus,  $F$ 's annotations tell us which parts of the output (i.e.,  $v_1, \dots, v_n$ ) *may* change if the input is changed at  $c$ . Dually, they also tell us what part of the output (i.e.,  $C[-, \dots, -]$ ) *cannot* be changed by changing the input at  $c$ .

Of course, dependency-correctness does not uniquely characterize the annotation behavior of a given  $F$ . It is possible for the annotations to be dependency-correct but *inaccurate*. For example we can always trivially annotate each part of the output with *every* color appearing in the input. This, of course, tells us nothing about the function's behavior. In general, the fewer the annotations present in the output of a dependency-correct  $F$ , the more they tell us about  $F$ 's behavior. We therefore consider a function  $F$  to be *minimally annotated* if no annotations can be removed from  $F$ 's output for any  $v$  without damaging correctness.

We say that a query  $e$  is *constant* if  $\llbracket e \rrbracket \gamma = v$  for some  $v$  and every suitable  $\gamma$ .

**Proposition 1.** *It is undecidable whether a Boolean NRC query is constant.*

*Proof.* Recall that query equivalence is undecidable for the (nested) relational calculus [3]; this holds even for queries  $e(x), e'(x)$  over a single variable  $x$ . Given two such queries, consider the expression  $\hat{e} = e(x) \approx e'(x) \vee y$  (definable as  $\neg(\neg(e(x) \approx e'(x)) \wedge \neg y)$ ). The result of this expression cannot be false everywhere since the disjunction is true for  $y = \text{true}$ , so is  $\hat{e}$  constant iff  $\llbracket \hat{e} \rrbracket v = \text{true}$  for every  $v$  iff  $e \equiv e'$ .

$$\begin{aligned}
(w^\Phi)^{+\Phi_0} &= w^{\Phi \cup \Phi_0} & (i_1^{\Phi_1}) \hat{+} (i_2^{\Phi_2}) &= (i_1 + i_2)^{\Phi_1 \cup \Phi_2} \\
\hat{\neg}(b^\Phi) &= (\neg b)^\Phi & (b_1^{\Phi_1}) \hat{\wedge} (b_2^{\Phi_2}) &= (b_1 \wedge b_2)^{\Phi_1 \cup \Phi_2} \\
\hat{\pi}_i((v_1, v_2)^\Phi) &= v_i^{+\Phi} & (w_1^{\Phi_1}) \hat{\cup} (w_2^{\Phi_2}) &= (w_1 \cup w_2)^{\Phi_1 \cup \Phi_2} \\
\widehat{\text{cond}}(\text{true}^\Phi, v_1, v_2) &= v_1^{+\Phi} & \widehat{\text{cond}}(\text{false}^\Phi, v_1, v_2) &= v_2^{+\Phi}
\end{aligned}$$

$$\begin{aligned}
\widehat{\sum}\{v_1, \dots, v_n\}^\Phi &= (v_1 \hat{+} \dots \hat{+} v_n)^{+\Phi} \\
\widehat{\cup}\{v_1, \dots, v_n\}^\Phi &= (v_1 \hat{\cup} \dots \hat{\cup} v_n)^{+\Phi} \\
\{v(x) \mid x \hat{\in} w^\Phi\} &= \{v(x) \mid x \in w\}^\Phi \\
(w_1^{\Phi_1}) \hat{\supset} (w_2^{\Phi_2}) &= \{v \in w_1 \mid |v| \notin |w_2|\}^{\Phi_1 \cup \|w_1\| \cup \Phi_2 \cup \|w_2\|} \\
v_1 \hat{\approx} v_2 &= \begin{cases} \text{true}^{\|v_1\| \cup \|v_2\|} & |v_1| = |v_2| \\ \text{false}^{\|v_1\| \cup \|v_2\|} & |v_1| \neq |v_2| \end{cases}
\end{aligned}$$

**Fig. 5.** Auxiliary annotation-propagating operations

$$\begin{aligned}
\mathcal{P}[x] \hat{\gamma} &= \hat{\gamma}(x) & \mathcal{P}[\text{let } x = e_1 \text{ in } e_2] \hat{\gamma} &= \mathcal{P}[e_2] \hat{\gamma}[x \mapsto \mathcal{P}[e_1] \hat{\gamma}] \\
\mathcal{P}[i] \hat{\gamma} &= i^\emptyset & \mathcal{P}[e_1 + e_2] \hat{\gamma} &= \mathcal{P}[e_1] \hat{\gamma} \hat{+} \mathcal{P}[e_2] \hat{\gamma} \\
\mathcal{P}[\text{sum}(e)] \hat{\gamma} &= \widehat{\sum} \mathcal{P}[e] \hat{\gamma} & \mathcal{P}[b] \hat{\gamma} &= b^\emptyset \\
\mathcal{P}[\neg e] \hat{\gamma} &= \hat{\neg} \mathcal{P}[e] \hat{\gamma} & \mathcal{P}[e_1 \wedge e_2] \hat{\gamma} &= \mathcal{P}[e_1] \hat{\gamma} \hat{\wedge} \mathcal{P}[e_2] \hat{\gamma} \\
\mathcal{P}[(e_1, e_2)] \hat{\gamma} &= (\mathcal{P}[e_1] \hat{\gamma}, \mathcal{P}[e_2] \hat{\gamma})^\emptyset & \mathcal{P}[\pi_i(e)] \hat{\gamma} &= \hat{\pi}_i(\mathcal{P}[e] \hat{\gamma}) \\
\mathcal{P}[\emptyset] \hat{\gamma} &= \emptyset^\emptyset & \mathcal{P}[\{e\}] \hat{\gamma} &= \{\mathcal{P}[e] \hat{\gamma}\}^\emptyset \\
\mathcal{P}[e_1 \cup e_2] \hat{\gamma} &= \mathcal{P}[e_1] \hat{\gamma} \hat{\cup} \mathcal{P}[e_2] \hat{\gamma} & \mathcal{P}[e_1 - e_2] \hat{\gamma} &= \mathcal{P}[e_1] \hat{\gamma} \hat{\supset} \mathcal{P}[e_2] \hat{\gamma} \\
\mathcal{P}[\cup e] \hat{\gamma} &= \widehat{\cup} \mathcal{P}[e] \hat{\gamma} & \mathcal{P}[\{e \mid x \in e_0\}] \hat{\gamma} &= \{\widehat{\mathcal{P}[e] \hat{\gamma}[x \mapsto v] \mid v \hat{\in} \mathcal{P}[e_0] \hat{\gamma}}\} \\
\mathcal{P}[e_1 \approx e_2] \hat{\gamma} &= \mathcal{P}[e_1] \hat{\gamma} \hat{\approx} \mathcal{P}[e_2] \hat{\gamma} & \mathcal{P}[\text{if } e_0 \text{ then } e_1 \text{ else } e_2] \hat{\gamma} &= \widehat{\text{cond}}(\mathcal{P}[e_0] \hat{\gamma}, \mathcal{P}[e_1] \hat{\gamma}, \mathcal{P}[e_2] \hat{\gamma})
\end{aligned}$$

**Fig. 6.** Provenance-tracking semantics

Clearly, an annotation is needed on the result of a Boolean query if and only if the query is not a constant, so finding minimal annotations is undecidable. As a result, we cannot expect to be able to compute dependency-correct annotations with perfect accuracy. It is important to note, though, that dependency-tracking is hard even if we leave out the  $e_1 - e_2$  or  $e_1 \approx e_2$  operators. For example, we can reduce (coNP-hard) propositional validity problems to dependency-tracking for ordinary Boolean expressions or conditionals.

### 3.1 Dynamic Provenance Tracking

We now consider a *provenance tracking* approach in which we interpret each expressions  $e$  as dependency-correct a-functions  $\mathcal{P}[e]$ . The definition of the provenance-tracking semantics is shown in Figure 6. Auxiliary operations are used to define  $\mathcal{P}[-]$ ; these are shown in Figure 5. In particular, note that we define  $(w^\Phi)^{+\Psi} = w^{\Phi \cup \Psi}$ .

Many cases involving ordinary programming constructs are self-explanatory. Constants always have empty annotations: nothing in the input can affect them. Built-in



functions such as  $+$ ,  $\wedge$ ,  $\neg$  propagate all annotations on their arguments to the result. For a conditional if  $e_0$  then  $e_1$  else  $e_2$ , the result is the result of evaluating  $e_1$  or  $e_2$ , combined with the annotations of  $e_0$ . A constructed pair has an empty top-level annotation; in a projection, the top-level annotation of the pair is merged with that of the returned value.

The cases involving collection types deserve some explanation. The empty set is a constant, so has an empty top-level annotation. Similarly, a singleton set constructor has an empty annotation. For union, we take the union of the underlying bags (of annotated values) and fuse the top-level annotations. For comprehension, we leave the top-level annotation alone. For flattening  $\bigcup e$ , we take the lifted union ( $\widehat{\cup}$ ) of the elements of  $e$  and add the top-level annotation of  $e$ . Similarly,  $\widehat{\text{sum}}(e)$  uses  $\widehat{+}$  to add together the elements of  $e$ , fusing their annotations with that of  $e$ . For set difference, to ensure dependency correctness, we must conservatively include all of the colors present on either side in the annotation of the top-level expression. Similarly, for equality tests, we must include all of the colors present in either value in the result annotation.

*Remark 1.* Our approach to handling negation and equality is somewhat awkward since it may result in very large annotations. For example, consider  $\{1^a, 2^b\}^c - \{1^d, 3^e\}^f$ : clearly, changing any of the input locations  $a, b, c, d, e, f$  can cause the output to change. In contrast, other approaches such as lineage associate tuple  $t \in R - S$  only with  $t \in R$  and all tuples of  $S$ . This may seem more “accurate”, but it is not dependency-correct. On the other hand, our approach can also be more “accurate” than lineage in the presence of negation; for example, in  $\{1\} - \{\pi_1(x) \mid x \in S\}$ , our approach will indicate that the output does not depend on the second components of elements of  $S$ , whereas the lineage of this query include all the records in  $S$ .

However, although our approach to negation and equality has pathological behavior in some cases, it does seem to provide useful provenance for typical queries. Developing more sophisticated forms of dependence that are better-behaved in the presence of negation or equality is an interesting area for future work.

*Example 3.* Consider an annotated input environment  $\widehat{\gamma}$ , shown in Figure 7(a), of schema  $R : (A : \text{int}, B : \text{int}), S(C : \text{int}, D : \text{int}, E : \text{int})$  (we use more compact relational schema notation with field names for readability). Figure 7(b) shows the provenance tracking semantics of the example queries from Figure 4. We write  $a_{123}$  as an abbreviation for  $a_1, a_2, a_3$ , etc. Note that in the count example query, the output depends on *no* individual of the input; we cannot change the number of elements of a multiset by changing field values. However, in a query such as  $\text{count}(\sigma_{A=B}(R))$ , the result depends on all of the  $A$  and  $B$  fields.

These annotated results can easily be used to “highlight” the parts of the input that may be relevant to a part of the output, by examining the annotations appearing above the output part of interest. This is how the example in Figure 1 was constructed.

Note that equivalent expressions  $e \equiv e'$  need not satisfy  $\mathcal{P}[e] \equiv \mathcal{P}[e']$ ; for example,  $x - x \equiv \emptyset$  but  $\mathcal{P}[x - x] \not\equiv \emptyset^\emptyset$ .

We summarize the main results concerning  $\mathcal{P}[-]$  as follows:

**Theorem 1.** *If  $\Gamma \vdash e : \tau$  then (1)  $\mathcal{P}[e] : \mathcal{A}_0[\Gamma] \rightarrow \mathcal{A}_0[\tau]$ , (2)  $|\mathcal{P}[e]| = \mathcal{E}[e]$ , and (3)  $\mathcal{P}[e]$  is dependency-correct.*

(a)

$$\hat{\gamma} = [R := \{(A : 1^{a_1}, B : 1^{b_1}), (A : 1^{a_2}, B : 2^{b_2}), (A : 2^{a_3}, B : 3^{b_3})\}, \\ S := \{(C : 1^{c_1}, D : 2^{d_1}, E : 3^{e_1}), (C : 1^{c_2}, D : 1^{d_2}, E : 4^{e_2})\}]$$

(b)

$$\begin{aligned} \mathcal{P}[\Pi_A(R)]\hat{\gamma} &= \{(A : 1^{a_1}), (A : 1^{a_2}), (A : 2^{a_3})\} \\ \mathcal{P}[\sigma_{A=B}(R)]\hat{\gamma} &= \{(A : 1^{a_1}, B : 1^{b_1})\}^{a_1, b_1} \\ \mathcal{P}[R \times S]\hat{\gamma} &= \{(A : 1^{a_1}, B : 1^{b_1}, C : 1^{c_1}, D : 2^{d_1}, E : 3^{e_1}), \\ &\quad (A : 1^{a_1}, B : 1^{b_1}, C : 1^{c_2}, D : 1^{d_2}, E : 4^{e_2}), \dots\} \\ \mathcal{P}[\Pi_{BE}(\sigma_{A=D}(R \times S))]\hat{\gamma} &= \{(B : 1^{b_1}, E : 4^{e_2}), (B : 2^{b_2}, E : 4^{e_2}), \\ &\quad (B : 3^{b_3}, E : 3^{e_1})\}^{a_{123}, d_{12}} \\ \mathcal{P}[R \cup \rho_{A/C, B/D}(\Pi_{CD}(S))]\hat{\gamma} &= \{(A : 1^{a_1}, B : 1^{b_1}), (A : 1^{a_2}, B : 2^{b_2}), (A : 2^{a_3}, B : 3^{b_3}), \\ &\quad (A : 1^{c_1}, B : 2^{d_1}), (A : 1^{c_2}, B : 1^{d_2})\} \\ \mathcal{P}[R - \rho_{A/D, B/E}(\Pi_{DE}(S))]\hat{\gamma} &= \{(A : 1^{a_1}, B : 1^{b_2}), (A : 1^{a_2}, B : 2^{b_2})\}^{a_{123}, b_{123}, d_{12}, e_{12}} \\ \mathcal{P}[\text{sum}(\Pi_A(R))]\hat{\gamma} &= 4^{a_1, a_2, a_3} \quad \mathcal{P}[\text{count}(R)]\hat{\gamma} = 3 \end{aligned}$$

Fig. 7. (a) Annotated input environment (b) Examples of provenance tracking

### 3.2 Static Provenance Analysis

Although it can often be quite accurate, dynamic provenance seems expensive to compute and nontrivial to implement in a standard relational database system. Moreover, dynamic analysis cannot tell us anything about a query without looking at (annotated) input data. In this section we consider a *static provenance analysis* which statically approximates the dynamic provenance, but can be calculated more easily and without accessing the input.

We formulate the analysis as a type-based analysis; annotated types (a-types)  $\hat{\tau}$  and raw types (r-types)  $\omega$  are defined as follows:

$$\hat{\tau} ::= \omega^\Phi \quad \omega ::= \text{int} \mid \text{bool} \mid \hat{\tau} \times \hat{\tau}' \mid \{\hat{\tau}\}$$

We write  $\hat{\Gamma}$  for a typing context mapping variables to a-types. We lift the auxiliary a-value operations of erasure ( $|\hat{\tau}|$ ), annotation extraction ( $\|\hat{\tau}\|$ ), and compatibility ( $\hat{\tau}_1 \cong \hat{\tau}_2$ ) to a-types in the obvious way.

We also define a “union” operation  $\sqcup$  on compatible types as follows:

$$\begin{aligned} \omega_1^{\Phi_1} \sqcup \omega_2^{\Phi_2} &= (\omega_1 \sqcup \omega_2)^{\Phi_1 \cup \Phi_2} \quad \text{int} \sqcup \text{int} = \text{int} \quad \text{bool} \sqcup \text{bool} = \text{bool} \\ (\hat{\tau}_1 \times \hat{\tau}_2) \sqcup (\hat{\tau}'_1 \times \hat{\tau}'_2) &= (\hat{\tau}_1 \sqcup \hat{\tau}'_1) \times (\hat{\tau}_2 \sqcup \hat{\tau}'_2) \quad \{\hat{\tau}\} \sqcup \{\hat{\tau}'\} = \{\hat{\tau} \sqcup \hat{\tau}'\} \end{aligned}$$

Finally, we write  $\hat{\tau} \sqsubseteq \hat{\tau}'$  if  $\hat{\tau}' = \hat{\tau} \sqcup \hat{\tau}'$ ; this is essentially a subtyping relation.

$$\begin{array}{c}
\frac{}{\widehat{\Gamma} \vdash i : \text{int} \ \& \ \emptyset} \quad \frac{\widehat{\Gamma} \vdash e_1 : \text{int} \ \& \ \Phi_1 \quad \widehat{\Gamma} \vdash e_2 : \text{int} \ \& \ \Phi_2}{\widehat{\Gamma} \vdash e_1 + e_2 : \text{int} \ \& \ \Phi_1 \cup \Phi_2} \quad \frac{\widehat{\Gamma} \vdash e : \{\text{int}^{\Phi_0}\} \ \& \ \Phi}{\widehat{\Gamma} \vdash \text{sum}(e) : \text{int} \ \& \ \Phi_0 \cup \Phi} \\
\frac{}{\widehat{\Gamma} \vdash b : \text{bool} \ \& \ \emptyset} \quad \frac{\widehat{\Gamma} \vdash e : \text{bool} \ \& \ \Phi}{\widehat{\Gamma} \vdash \neg e : \text{bool} \ \& \ \Phi} \quad \frac{\widehat{\Gamma} \vdash e_1 : \text{bool} \ \& \ \Phi_2 \quad \widehat{\Gamma} \vdash e_2 : \text{bool} \ \& \ \Phi_2}{\widehat{\Gamma} \vdash e_1 \wedge e_2 : \text{bool} \ \& \ \Phi_1 \cup \Phi_2} \\
\frac{\widehat{\Gamma} \vdash e_1 : \widehat{\tau}_1 \quad \widehat{\Gamma} \vdash e_2 : \widehat{\tau}_2}{\widehat{\Gamma} \vdash (e_1, e_2) : (\widehat{\tau}_1 \times \widehat{\tau}_2) \ \& \ \emptyset} \quad \frac{\widehat{\Gamma} \vdash e : \omega_1^{\Phi_1} \times \omega_2^{\Phi_2} \ \& \ \Phi}{\widehat{\Gamma} \vdash \pi_i(e) : \omega_i \ \& \ \Phi_i \cup \Phi} \\
\frac{\widehat{\Gamma} \vdash e_1 : \widehat{\tau}_1 \quad \widehat{\Gamma} \vdash e_2 : \widehat{\tau}_2 \quad \widehat{\tau}_1 \cong \widehat{\tau}_2}{\widehat{\Gamma} \vdash e_1 \approx e_2 : \text{bool} \ \& \ \|\widehat{\tau}_1\| \cup \|\widehat{\tau}_2\|} \quad \frac{\widehat{\Gamma} \vdash e_0 : \text{bool} \ \& \ \Phi_0 \quad \widehat{\Gamma} \vdash e_1 : \widehat{\tau}_1 \quad \widehat{\Gamma} \vdash e_2 : \widehat{\tau}_2 \quad \omega_1 \cong \omega_2}{\widehat{\Gamma} \vdash \text{if } e_0 \text{ then } e_1 \text{ else } e_2 : (\widehat{\tau}_1 \sqcup \widehat{\tau}_2)^{+\Phi_0}} \\
\frac{}{\widehat{\Gamma} \vdash \emptyset : \{\widehat{\tau}\} \ \& \ \emptyset} \quad \frac{\widehat{\Gamma} \vdash e : \widehat{\tau}}{\widehat{\Gamma} \vdash \{e\} : \{\widehat{\tau}\} \ \& \ \emptyset} \quad \frac{\widehat{\Gamma} \vdash e_1 : \{\widehat{\tau}_1\} \ \& \ \Phi_1 \quad \widehat{\Gamma} \vdash e_2 : \{\widehat{\tau}_2\} \ \& \ \Phi_2 \quad \widehat{\tau}_1 \cong \widehat{\tau}_2}{\widehat{\Gamma} \vdash e_1 \cup e_2 : \{\widehat{\tau}_1 \sqcup \widehat{\tau}_2\} \ \& \ \Phi_1 \cup \Phi_2} \\
\frac{\widehat{\Gamma} \vdash e_1 : \{\widehat{\tau}_1\} \ \& \ \Phi_1 \quad \widehat{\Gamma} \vdash e_2 : \{\widehat{\tau}_2\} \ \& \ \Phi_2 \quad \widehat{\tau}_1 \cong \widehat{\tau}_2}{\widehat{\Gamma} \vdash e_1 - e_2 : \{\widehat{\tau}_1\} \ \& \ \|\{\widehat{\tau}_1\}^{\Phi_1}\| \cup \|\{\widehat{\tau}_2\}^{\Phi_2}\|} \quad \frac{\widehat{\Gamma} \vdash e : \{\{\widehat{\tau}\}^{\Phi_2}\} \ \& \ \Phi_1}{\widehat{\Gamma} \vdash \bigcup e : \{\widehat{\tau}\} \ \& \ \Phi_1 \cup \Phi_2} \\
\frac{\widehat{\Gamma} \vdash e_1 : \{\widehat{\tau}_1\} \ \& \ \Phi_1 \quad \widehat{\Gamma}, x : \widehat{\tau}_1 \vdash e_2 : \omega \ \& \ \Phi_2}{\widehat{\Gamma} \vdash \{e_2 \mid x \in e_1\} : \{\omega_2^{\Phi}\} \ \& \ \Phi_1} \quad \frac{x : \widehat{\tau} \in \widehat{\Gamma}}{\widehat{\Gamma} \vdash x : \widehat{\tau}} \quad \frac{\widehat{\Gamma} \vdash e_1 : \widehat{\tau}_1 \quad \widehat{\Gamma}, x : \widehat{\tau}_1 \vdash e_2 : \widehat{\tau}_2}{\widehat{\Gamma} \vdash \text{let } x = e_1 \text{ in } e_2 : \widehat{\tau}_2}
\end{array}$$

**Fig. 8.** Type-based provenance analysis

We interpret a-types  $\widehat{\tau}$  as sets  $\mathcal{A}[\widehat{\tau}]$  of a-values. We interpret the annotations in a-types as upper bounds on the annotations in the corresponding a-values:

$$\begin{aligned}
\mathcal{A}[\text{int}] &= \{i \mid i \in \mathbb{Z}\} & \mathcal{A}[\widehat{\tau}_1 \times \widehat{\tau}_2] &= \{(v_1, v_2) \mid v_1 \in \mathcal{A}[\widehat{\tau}_1], v_2 \in \mathcal{A}[\widehat{\tau}_2]\} \\
\mathcal{A}[\text{bool}] &= \{b \mid b \in \mathbb{B}\} & \mathcal{A}[\{\widehat{\tau}\}] &= \{V \mid \forall v \in V. v \in \mathcal{A}[\widehat{\tau}]\} \\
\mathcal{A}[\omega^{\Phi}] &= \{w^{\Psi} \mid \Psi \subseteq \Phi, w \in \mathcal{A}[\omega]\}
\end{aligned}$$

The syntactic operations  $|-$ ,  $\|\cdot\|$ ,  $\sqsubseteq$  and  $\sqcup$  correspond to appropriate semantic operations on sets of a-values. We note some useful properties of these operations:

- Lemma 2.**
1. If  $v \in \mathcal{A}[\widehat{\tau}]$  then  $|v| \in \mathcal{T}[\widehat{\tau}]$  and  $\|v\| \subseteq \|\widehat{\tau}\|$ .
  2. If  $\widehat{\tau}_1 \sqsubseteq \widehat{\tau}_2$  then  $\mathcal{A}[\widehat{\tau}_1] \subseteq \mathcal{A}[\widehat{\tau}_2]$  and  $\|\widehat{\tau}_1\| \subseteq \|\widehat{\tau}_2\|$ .
  3. If  $\widehat{\tau}_1 \sqcup \widehat{\tau}_2$  is defined then  $\mathcal{A}[\widehat{\tau}_1 \sqcup \widehat{\tau}_2] = \mathcal{A}[\widehat{\tau}_1] \cup \mathcal{A}[\widehat{\tau}_2]$  and  $\|\widehat{\tau}_1 \sqcup \widehat{\tau}_2\| = \|\widehat{\tau}_1\| \cup \|\widehat{\tau}_2\|$ .

The annotated typing judgment  $\widehat{\Gamma} \vdash e : \widehat{\tau}$  (sometimes written  $\widehat{\Gamma} \vdash e : \omega \ \& \ \Phi$  for readability, provided  $\widehat{\tau} = \omega^{\Phi}$ ) extends the plain typing judgment shown in Figure 2.

**Proposition 2.** *The judgment  $\Gamma \vdash e : \tau$  is derivable if and only if for any  $\widehat{\Gamma}$  enriching  $\Gamma$ , there exists a  $\widehat{\tau}$  enriching  $\tau$  such that  $\widehat{\Gamma} \vdash e : \widehat{\tau}$ . Moreover, given  $\Gamma \vdash e : \tau$ , and  $\widehat{\Gamma}$  enriching  $\Gamma$ , we can compute  $\widehat{\tau}$  in polynomial time (by a simple syntax-directed algorithm).*

The correctness of the analysis is proved with respect to the provenance-tracking semantics. This property takes the form of a type-soundness theorem: we simply need to show that if the input environment  $\widehat{\gamma}$  is well-formed at annotated context  $\widehat{\Gamma}$  then the result  $\mathcal{P}[e]\widehat{\gamma}$  is well-formed at type  $\widehat{\tau}$ :

$$\begin{aligned}
& \text{(a)} && \widehat{\Gamma} = [R : \{(A : \text{int}^a, B : \text{int}^b)\}, S : \{(C : \text{int}^c, D : \text{int}^d, E : \text{int}^e)\}] \\
& \text{(b)} && \widehat{\Gamma} \vdash \Pi_A(R) && : \{(A : \text{int}^a)\} \\
& && \widehat{\Gamma} \vdash \sigma_{A=B}(R) && : \{(A : \text{int}^a, B : \text{int}^b)\}^{a,b} \\
& && \widehat{\Gamma} \vdash R \times S && : \{(A : \text{int}^a, B : \text{int}^b, C : \text{int}^c, D : \text{int}^d, E : \text{int}^e)\} \\
& && \widehat{\Gamma} \vdash \Pi_{BE}(\sigma_{A=D}(R \times S)) && : \{(B : \text{int}^b, E : \text{int}^e)\}^{a,d} \\
& && \widehat{\Gamma} \vdash R \cup \rho_{A/C, B/D}(\Pi_{CD}(S)) && : \{(A : \text{int}^{a,c}, B : \text{int}^{b,d})\} \\
& && \widehat{\Gamma} \vdash R - \rho_{A/D, B/E}(\Pi_{DE}(S)) && : \{(A : \text{int}^a, B : \text{int}^b)\}^{a,b,d,e} \\
& && \widehat{\Gamma} \vdash \text{sum}(\Pi_A(R)) && : \text{int}^a \\
& && \widehat{\Gamma} \vdash \text{count}(R) && : \text{int}
\end{aligned}$$

**Fig. 9.** (a) Annotated input context (b) Examples of provenance analysis

**Theorem 2.** *If  $\widehat{\Gamma} \vdash e : \widehat{\tau}$  then  $\mathcal{P}[[e]] : \mathcal{A}[[\widehat{\Gamma}]] \rightarrow \mathcal{A}[[\widehat{\tau}]]$ .*

This theorem tells us that the annotations we obtain (statically) by provenance analysis over-approximate those obtained (dynamically) by provenance tracking provided the initial value  $\widehat{\gamma}$  matches  $\mathcal{A}[[\widehat{\Gamma}]]$ .

*Example 4.* Consider an annotated type context  $\widehat{\Gamma}$ , shown in Figure 9(a), where we have annotated field values  $A, B, C, D, E$  with colors  $a, b, c, d, e$  respectively. Figure 9(b) shows the results of static analysis for the queries in Figure 7. In some cases, the type information simply reflects the field names which are present in the output. However, the colors are not affected by renamings, as in  $\rho_{A/C, B/D}$ . Furthermore, note that (if we replace the colors  $a, b, c, d, e$  with color sets  $\{a_1, a_2, a_3\}$ , etc.) in each case the type-level colors over-approximate the value-level colors calculated in Figure 7.

*Example 5.* To illustrate how the analysis works in practice, we consider an extended example for a query that performs grouping and aggregation:

$$Q(R) = \{(\pi_1(x), \text{sum}(G(x))) \mid x \in R\}$$

where we make the following abbreviations:

$$\begin{aligned}
G(x) &:= \bigcup \{\text{if } \pi_1(y) \approx \pi_1(x) \text{ then } \{\pi_2(y)\} \text{ else } \{\} \mid y \in R\} \\
\widehat{\tau}_R &:= \text{int}^a \times \text{int}^b & \widehat{\Gamma} &:= R : \{\widehat{\tau}_R\} \\
\widehat{\Gamma}_1 &:= \widehat{\Gamma}, x : \widehat{\tau}_R & \widehat{\Gamma}_2 &:= \widehat{\Gamma}_1, y : \widehat{\tau}_R
\end{aligned}$$

We will derive  $\widehat{\Gamma} \vdash Q(R) : \{\text{int}^a \times \text{int}^{a,b}\}$ . The derivation illustrates how color  $a$  is propagated from the to both parts of the result type, while color  $b$  is only propagated to the second column.

First, we can reduce the analysis of  $Q$  to analyzing  $G$  as follows:

$$\frac{\frac{\widehat{I}_1 \vdash x : \widehat{\tau}_R}{\widehat{I}_1 \vdash \pi_1(x) : \text{int}^a} \quad \frac{\widehat{I}_1 \vdash G(x) : \{\text{int}^b\}^a}{\widehat{I}_1 \vdash \text{sum}(G(x)) : \text{int}^{a,b}}}{\widehat{I}_1 \vdash R : \{\widehat{\tau}_R\} \quad \widehat{I}_1 \vdash (\pi_1(x), \text{sum}(G(x))) : \text{int}^a \times \text{int}^{a,b}} \quad \frac{}{\widehat{I} \vdash \{(\pi_1(x), \text{sum}(G(x))) \mid x \in R\} : \{\text{int}^a \times \text{int}^{a,b}\}}$$

We next reduce the analysis of  $G$  to an analysis of the conditional inside  $G$ :

$$\frac{\frac{\widehat{I}_1 \vdash R : \{\widehat{\tau}_R\} \quad \widehat{I}_2 \vdash \text{if } \pi_1(y) \approx \pi_1(x) \text{ then } \{\pi_2(y)\} \text{ else } \{\} : \{\text{int}^b\}^a}{\widehat{I}_1 \vdash \{\text{if } \pi_1(y) \approx \pi_1(x) \text{ then } \{\pi_2(y)\} \text{ else } \{\} \mid y \in R\} : \{\{\text{int}^b\}^a\}}}{\widehat{I}_1 \vdash \bigcup \{\text{if } \pi_1(y) \approx \pi_1(x) \text{ then } \{\pi_2(y)\} \text{ else } \{\} \mid y \in R\} : \{\text{int}^b\}^a}$$

Finally, we can analyze the conditional as follows:

$$\frac{\frac{\frac{\widehat{I}_2 \vdash y : \widehat{\tau}_R}{\widehat{I}_2 \vdash \pi_1(y) : \text{int}^a} \quad \frac{\widehat{I}_2 \vdash y : \widehat{\tau}_R}{\widehat{I}_2 \vdash \pi_1(x) : \text{int}^a} \quad \frac{\widehat{I}_2 \vdash y : \widehat{\tau}_R}{\widehat{I}_2 \vdash \pi_2(y) : \text{int}^b}}{\widehat{I}_2 \vdash \pi_1(y) \approx \pi_1(x) : \text{bool}^a} \quad \frac{}{\widehat{I}_2 \vdash \{\pi_2(y)\} : \{\text{int}^b\}} \quad \frac{}{\widehat{I}_2 \vdash \{\} : \{\text{int}\}}}{\widehat{I}_2 \vdash \text{if } \pi_1(y) \approx \pi_1(x) \text{ then } \{\pi_2(y)\} \text{ else } \{\} : \{\text{int}^b\}^a}$$

### 3.3 Discussion

We have implemented a prototype interpreter for the NRC that performs ordinary type-checking and evaluation as well as provenance tracking and analysis. We used this implementation to construct the examples.

We chose to study provenance via the NRC because it is a clean and system-independent model; we believe our results can be specialized to common database implementations and physical operators without much difficulty. We have not yet investigated scaling this approach to large datasets. There are several apparent obstacles to implementing annotation-based provenance tracking in standard database systems that do not natively support annotation. Recent research has begun to address this problem [6, 18] and we plan to investigate whether these techniques can be used to implement our approach.

Static provenance analysis is also more expensive than ordinary typechecking, but since the overhead is proportional only to the size of the *query*, not the (usually much larger) data, this seems acceptable. Moreover, static analysis may be useful in optimizing provenance tracking, for example by using the results of static analysis to avoid tracking annotations that are statically irrelevant to the output.

## 4 Related and Future Work

**Slicing and other dependence analyses** Dependence tracking and analysis have been shown to be useful in many contexts [1] such as program slicing [7, 16, 24], memoization and caching [2, 4], and information-flow security [20]. In program slicing [24],

the goal is to identify a (small) set of program points whose execution contributes to the value of an output variable (or other observable behavior). This is analogous to our approach to provenance, except that provenance identifies relevant parts of the *input database*, not the *program* (i.e. query). Our approach is inspired by, and in some cases could be viewed as an adaptation of, these techniques to a database setting with collection types.

**Provenance in databases** Most work on provenance in databases [23, 15, 11, 12, 10] has focused on identifying information that explains *why* some data is present in the output of a query (or view) or *where* some data in the output was copied from in the input. However, semantic characterizations of these intuitions have been elusive and difficult to generalize beyond monotone relational queries. Our work generalizes some of these techniques and provides clear semantic guarantees and qualitatively useful provenance information in the presence of grouping and aggregation.

**Updates** Some recent work has generalized where-provenance to database updates [9, 10], motivated by “curated” scientific databases that are updated frequently, often by (manual) copying from other sources. Our approach addresses an orthogonal issue; we plan to investigate dependency provenance for updates.

**Workflow provenance** Provenance has also been studied in geospatial and scientific “grid” computation [8, 17, 21], particularly for *workflows* (visual programs written by scientists). At present, formal correctness criteria have not been identified for most of these approaches, but we believe it to be worthwhile to seek a foundation for workflow provenance using dependence analysis.

**Annotations** Recent research on annotations, uncertainty, and incomplete information [6, 5, 18, 19] is closely related to provenance. Green et al. [19] showed that relations with semiring-valued annotations generalize several variations of the relational model, including set, bag, probabilistic, and incomplete information models, and identified a relationship between free semiring-valued relations and why-provenance.

## 5 Conclusions

Provenance information that relates parts of the result of a query to “relevant” parts of the input is useful for many purposes, including judging the reliability of information based on the relevant sources and identifying parts of the database that may be responsible for an error in the output of a query. We have argued that the notion of *dependence*, familiar from program slicing, information flow security, and other analyses, provides a solid semantic foundation for understanding provenance for complex database queries. In this paper we introduced a semantic characterization of *dependency provenance*, showed that minimal dependency provenance is not computable, and presented approximate tracking and analysis techniques. We believe there are many promising directions for future work, including implementing efficient practical techniques, identifying more sophisticated and useful dependency properties, and studying dependency provenance in other settings such as update languages and workflows.

*Acknowledgments* We wish to thank Peter Buneman and Stijn Vansummeren for helpful discussions on this work.

## References

1. Martín Abadi, Anindya Banerjee, Nevin Heintze, and Jon G. Riecke. A core calculus of dependency. In *POPL*, pages 147–160. ACM Press, 1999.
2. Martín Abadi, Butler Lampson, and Jean-Jacques Lévy. Analysis and caching of dependencies. In *ICFP*, pages 83–91, New York, NY, USA, 1996. ACM Press.
3. Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
4. Umut A. Acar, Guy E. Blelloch, and Robert Harper. Selective memoization. In *Proceedings of the 30th Annual ACM Symposium on Principles of Programming Languages*, 2003.
5. Omar Benjelloun, Anish Das Sarma, Alon Y. Halevy, and Jennifer Widom. ULDBs: Databases with uncertainty and lineage. In *VLDB*, pages 953–964, 2006.
6. Deepavali Bhagwat, Laura Chiticariu, Wang-Chiew Tan, and Gaurav Vijayvargiya. An annotation management system for relational databases. *VLDB Journal*, 14(4):373–396, 2005.
7. S. Biswas. *Dynamic Slicing in Higher-Order Programming Languages*. PhD thesis, University of Pennsylvania, 1997.
8. Rajendra Bose and James Frew. Lineage retrieval for scientific data processing: a survey. *ACM Comput. Surv.*, 37(1):1–28, 2005.
9. Peter Buneman, Adriane Chapman, and James Cheney. Provenance management in curated databases. In *SIGMOD 2006*, pages 539–550, 2006.
10. Peter Buneman, James Cheney, and Stijn Vansummeren. On the expressiveness of implicit provenance in query and update languages. In *ICDT 2007*, number 4353 in Lecture Notes in Computer Science, pages 209–223. Springer, 2007.
11. Peter Buneman, Sanjeev Khanna, and Wang-Chiew Tan. Why and where: A characterization of data provenance. In *ICDT 2001*, number 1973 in LNCS, pages 316–330. Springer, 2001.
12. Peter Buneman, Sanjeev Khanna, and Wang-Chiew Tan. On propagation of deletions and annotations through views. In *PODS*, pages 150–158, 2002.
13. Peter Buneman, Shamim A. Naqvi, Val Tannen, and Limsoon Wong. Principles of programming with complex objects and collection types. *Theor. Comp. Sci.*, 149(1):3–48, 1995.
14. James Cheney, Amal Ahmed, and Umut Acar. Provenance as dependency analysis. Technical Report arXiv:0708.2173v1, arXiv.org e-Print archive, 2007.
15. Yingwei Cui, Jennifer Widom, and Janet L. Wiener. Tracing the lineage of view data in a warehousing environment. *ACM Trans. Database Syst.*, 25(2):179–227, 2000.
16. John Field and Frank Tip. Dynamic dependence in term rewriting systems and its application to program slicing. *Information and Software Technology*, 40(11–12):609–636, November/December 1998.
17. Ian Foster and Luc Moreau, editors. *Proceedings of the 2006 International Provenance and Annotation Workshop (IPAW 2006)*. Number 4145 in LNCS. Springer-Verlag, 2006.
18. Floris Geerts, Anastasios Kementsietsidis, and Diego Milano. Mondrian: Annotating and querying databases through colors and blocks. In *ICDE 2006*, page 82, 2006.
19. Todd J. Green, Grigoris Karvounarakis, and Val Tannen. Provenance semirings. In *PODS*, pages 31–40, New York, NY, USA, 2007. ACM Press.
20. Andrei Sabelfeld and Andrew Myers. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications*, 21(1):5–19, 2003.
21. Yogesh Simmhan, Beth Plale, and Dennis Gannon. A survey of data provenance in e-science. *SIGMOD Record*, 34(3):31–36, 2005.
22. P. Wadler. Comprehending monads. *Mathematical Structures in Computer Science*, 2:461–493, 1992.
23. Y. Richard Wang and Stuart E. Madnick. A polygen model for heterogeneous database systems: The source tagging perspective. In *VLDB*, pages 519–538, 1990.
24. Mark Weiser. Program slicing. In *ICSE*, pages 439–449, Piscataway, NJ, USA, 1981. IEEE Press.