

# DBNotes: A Post-It System for Relational Databases based on Provenance\*

Laura Chiticariu  
UC Santa Cruz  
laura@cs.ucsc.edu

Wang-Chiew Tan  
UC Santa Cruz  
wctan@cs.ucsc.edu

Gaurav Vijayvargiya  
UC Santa Cruz  
gaurav@cs.ucsc.edu

## ABSTRACT

We demonstrate DBNotes, a Post-It note system for relational data. In DBNotes, it is possible to attach zero or more notes to every value in a relation. When a relation is queried, DBNotes propagates the attached notes to the result of the query. The method by which a note is propagated is based on provenance (aka lineage). As a consequence, if every value in a relation is attached with a note that describes its address, the notes associated with a value in the result of a query show the provenance of the value.

In this demonstration, we show how one can use the query language pSQL of DBNotes to propagate and query annotations in different ways. We also demonstrate how one can use DBNotes to provide a high-level explanation of the provenance and flow of a value that may be the result of many query transformation steps. A high-level explanation shows a possible journey taken by the value through various databases and query transformation steps. DBNotes can also provide a detailed explanation at each transformation step to precisely explain why a value has moved from a source to a target database.

## 1. INTRODUCTION

We demonstrate DBNotes, a Post-It note system for relational databases where every piece of data may be associated with zero or more notes (or annotations). These annotations are transparently propagated along as data is being transformed. The method by which annotations are propagated is based on provenance: the annotations associated with a piece of data  $d$  in the result of a transformation consist of the annotations associated with each piece of data in the source where  $d$  is copied from. One immediate application of this system is to use annotations to systematically trace the provenance and flow of data. If every piece of source data is attached with an annotation that describes its address (i.e., origins), then the annotations of a piece of data in the re-

\*Supported in part by an NSF CAREER Award IIS-0347065 and an NSF grant IIS-0430994.

sult of a transformation describe its provenance. Hence, one can easily determine the provenance of data through a sequence of transformation steps simply by examining the annotations. Annotations can also be used to store additional information about data. Since a database schema is often proprietary, the ability to insert new information about data without having to change the underlying schema is a useful feature. For example, an error report could be attached to an erroneous piece of data, and this error report will be propagated to other databases along transformations, thus notifying other users of the error. Overall, the annotations on the result of a transformation can also provide an estimate on the quality of the resulting database.

**Summary of DBNotes Demonstration Features** We demonstrate four main features of DBNotes. (1) First, we demonstrate pSQL, an extension of a fragment of SQL that allows one to specify how annotations should propagate through an SQL query. We demonstrate three types of propagation schemes that are supported by pSQL. Annotations can be propagated according to where data is copied from (the default scheme) or according to where data is copied from in *all* equivalent queries (the default-all scheme) or according to the user specification (the custom scheme). (2) The second feature of DBNotes that we demonstrate is the ability to use pSQL to query both data and its annotations. With this feature, one can pose queries to retrieve, for example, all tuples that are derived from a particular source or all tuples with an attached error report. (3) The third feature of DBNotes that we demonstrate is the ability to provide a detailed explanation on the provenance of an annotation, an attribute value or a tuple in the result of a query transformation. For example, we explain the provenance of a tuple by demonstrating a proof that shows how the query transformation produces the output tuple. We show a binding for each variable in the SQL query and demonstrate that under these bindings, the conditions in the **WHERE** clause are satisfied and the desired output tuple is produced according to the **SELECT** clause of the SQL query. (4) The last feature of DBNotes we demonstrate is its capability to visually describe the journey (i.e., the provenance and flow) taken by a piece of data through various databases. For example, when a user asks about the journey of a piece of data  $d$  in a database  $D$ , DBNotes displays the sequence of databases and transformation steps that derive  $d$  in  $D$ . It also displays the sequence of databases and transformation steps that depend on  $d$  in  $D$ . At the visual interface, the user also has the option of “zooming in” on each transformation step to see a detailed explanation of each transformation step.

**Related Work** A study on the problem of tracing the provenance of data that is the result of a query applied on a relational database was first approached by [5]. The type of provenance studied by [5] is called *why-provenance* according to [3]. This type of provenance explains why a tuple is in the result of a query transformation. In [3], a notion of *where-provenance* is also characterized. The where-provenance of an attribute value tells one exactly where that value is copied from. DBNotes propagates annotations based on where-provenance. In both [3, 5], a “reverse” query is generated in order to answer provenance. This approach is *lazy* in the sense that it requires a reverse query to be generated and evaluated only when the provenance of an output tuple is sought for. In contrast, DBNotes computes provenance *eagerly* by forwarding annotations as data is transformed. A high-level description of how a piece of data has moved along transformations can be constructed by examining the annotations in each database. A reverse query is also generated and evaluated by DBNotes when a detailed explanation on the provenance of an annotation, attribute value or tuple in a database is sought for. A lineage tracing system for data warehouses [4] has been implemented. In this system, the lineage of a tuple in the result of a complex relational view is explained based on relational operators. STAG [9] is a system that allows for querying of annotations over digital documents using an SQL-like query language. However, unlike DBNotes, STAG does not propagate annotations. The idea of propagating annotations along transformations is in fact not new and has been proposed in various forms in existing literature. We refer the interested reader to [2] for more discussion on related work.

To the best of our knowledge, DBNotes is the first Post-It notes system for relational databases that allows a user to specify how annotations should propagate, query annotations and analyze the provenance and flow of data through databases generated by query transformations.

## 2. DEMONSTRATION OVERVIEW

We make use of the following example restaurant database for our demonstration.

### 2.1 Our Example Database

Our example restaurant database contains information about restaurants owned by Hollywood celebrities in the greater Los Angeles area. We have collected information from four different sources [1, 6, 7, 8] and the corresponding schema for each source is shown below:

```
Blue.Pages(Name,Address,Owner,Cuisine,Cost)
LATimes(Res_Name,Hours,Location,Owner,Cost,Type_of.food)
Restaurant_Reviews(Name,Owned_by,Food,Address,Expensive)
Seeing_Stars(Name,Address,Getting_there,Owner,Cost,Cuisine)
```

Every attribute value of every tuple may be associated with zero or more annotations. For example, the tuple  $t_1$

$t_1$ : (Madre’s {LATimes:Good}, Jennifer Lopez, Cuban {May be the worst Cuban food in CA}, Pasadena, Expensive)

in relation LATimes has one annotation attached to the value “Madre’s” and one annotation attached to the value “Cuban”. The annotations are shown in braces. The rest of the attribute values do not have any annotations. Annotations are not part of the database schema and the method by which annotations are stored and propagated is transparent to the user.

### 2.2 Propagating Annotations

We first demonstrate a feature of pSQL, an extension of a fragment of SQL, that allows one to specify how annotations should propagate through an SQL query. For example, the pSQL query below integrates information about restaurants from the first two sources and propagates annotations in the default way (i.e., according to where data is copied from). For  $Q_1$ , this means that every tuple in Blue.Pages and LATimes as well as the associated annotations are emitted to the result.

```
Q1:
SELECT DISTINCT Name, Address, Owner, Cuisine, Cost
FROM Blue.Pages
PROPAGATE DEFAULT
UNION
SELECT DISTINCT Res_Name AS Name, Location AS Address,
Owner, Type_of.Food AS Cuisine, Cost
FROM LATimes
PROPAGATE DEFAULT
```

If a tuple occurs in both Blue.Pages and LATimes, only one of the tuples is emitted and the corresponding annotations are unioned together. For example, if  $t_1$  is a tuple in LATimes and  $t_2$  (shown below) is a tuple in Blue.Pages, then the tuple  $t$  shown below appears in the output of  $Q_1$ .  
 $t_2$ : (Madre’s {Blue.Pages:Bad}, Pasadena, Jennifer Lopez, Cuban {Bad Food}, Expensive {Expensive Restaurant})  
 $t$ : (Madre’s {Blue.Pages:Bad, LATimes:Good}, Pasadena, Jennifer Lopez, Cuban {Bad Food, May be the worst Cuban food in CA}, Expensive {Expensive Restaurant})

In addition to the default propagation scheme, we shall also demonstrate the default-all and custom propagation schemes supported by pSQL.

### 2.3 Querying Annotations

We also demonstrate a feature of pSQL which allows one to query over data and annotations through the following use cases. Suppose we have an integrated view, called Restaurants, of all four sources that is given by executing a query written against  $Q_1$  and the two sources Restaurant\_Reviews and Seeing\_Stars and the resulting schema is the same as that of  $Q_1$ . An important feature of pSQL that allows annotations to be queried is through the use of the keyword ANNOT(attribute-name) which elevates annotations of attribute-name to first-class citizens. Some examples of how one could pose queries against data and annotations are shown below in queries  $Q_2$  to  $Q_4$ .

```
Q2:
SELECT DISTINCT *
FROM Restaurants r, ANNOT(r.Name) a
WHERE r.cost = 'Expensive' AND a LIKE '%Bad%'
PROPAGATE DEFAULT
```

The query  $Q_2$  retrieves all expensive restaurants with a bad rating. The *annotation variable* “a” ranges over elements in the set of annotations associated with each Name value of each tuple in Restaurants. This query demonstrates how one can query over data and annotations.

```
Q3:
SELECT DISTINCT r.Name AS Name
FROM Restaurants r, ANNOT(r.Name) a1,
ANNOT(r.Name) a2
WHERE a1 LIKE '%LATimes:Good%' AND
a2 LIKE '%Blue.Pages:Bad%'
PROPAGATE DEFAULT
```

The query  $Q_3$  retrieves all restaurants that have a good rating from LATimes and a bad rating from Blue.Pages.

This query demonstrates how one can query about source information provided that annotations carry information about the sources.

```

Q4:
SELECT DISTINCT r.Name AS Name, COUNT(a) AS Bad_Ratings
FROM Restaurant r, ANNOT(r.Name) a
WHERE a LIKE '%Bad%'
PROPAGATE a TO Name
GROUP BY Name
HAVING COUNT(a) > 2

```

The query  $Q_4$  (with a custom propagation scheme) retrieves bad restaurants where a restaurant is considered as bad if it has more than two bad ratings. This query demonstrates how one can count the number of annotations.

## 2.4 Explaining Data Provenance

The third feature of DBNotes that we demonstrate is its ability to provide detailed explanations on the provenance of an element (i.e., an annotation, an attribute value or a tuple) in the result of a query transformation. We achieve this in a declarative fashion, by demonstrating a proof that shows how the transformation produced that element.

The query  $Q_5$  below retrieves all celebrities that frequent expensive restaurants which received a bad rating from Blue\_Pages. (We assume our example database also contains a relation Frequent(Name, Restaurant) that records information about famous stars and the restaurants they visit. This information is obtained from [8].)

```

Q5:
SELECT DISTINCT f.Name AS Celebrity, r.Name AS Restaurant
FROM Frequent f, Blue_Pages r,
      ANNOT(r.Name) a
WHERE f.Restaurant = r.Name AND
      r.Cost = 'Expensive' AND
      a LIKE '%Blue_Pages:Bad%'
PROPAGATE DEFAULT

```

Let  $t_3$  be a tuple (shown below) in Frequent which states that Ben Affleck is a regular of Madre’s. Suppose  $t_2$  (from Section 2.2) is a tuple in Blue\_Pages, then the following tuple  $t'$  appears in the output of  $Q_5$ :

```

t3: (Ben Affleck {Dated Jennifer Lopez}, Madre’s)
t': (Ben Affleck {Dated Jennifer Lopez}, Madre’s {Blue_Pages:Bad})

```

When asked to provide a detailed explanation of the provenance of the attribute value “Madre’s” in the result of  $Q_5$ , DBNotes will first generate a reverse query. The purpose of this query is to extract the relevant source tuples that produced the attribute value “Madre’s”. DBNotes will show, for each tuple and annotation variable of  $Q_5$ , a binding to source tuples and respectively, source annotations, and demonstrate that under these bindings, the conditions in the WHERE clause are satisfied and the desired output element is produced according to the SELECT clause of  $Q_5$ . In this case, DBNotes shows the bindings for  $r$ ,  $f$  and  $a$  and detailed explanation as shown in Figure 1.

## 2.5 Tracing the Provenance and Flow of Data

The last feature of DBNotes that we demonstrate is its capability to provide a visual of the “journey” taken by a piece of data through various databases and transformation steps, where the term “journey” refers to both the provenance and flow of that piece of data. When asked to show the provenance and flow of a piece of data  $d$  in a database  $D$ , DBNotes shows the sequence of databases and transformation steps that derived  $d$  in  $D$  as well as the sequence of

Blue_Pages:					
	Name	Address	Owner	Cuisine	Cost
$r$	Madre’s	Pasadena	Jennifer Lopez	Cuban	Expensive

Frequent:		ANNOT(r.Name)
	Name	Restaurant
$f$	Ben Affleck	Madre’s

	$a$
	Blue_Pages:Bad

The conditions in the WHERE clause of  $Q_5$  are satisfied as follows:

1. The condition  $f.Restaurant = r.Name$  is satisfied because  $f.Restaurant = \text{“Madre’s”} = r.Name$ .
2. The condition  $r.Cost = \text{‘Expensive’}$  is satisfied because  $r.Cost = \text{“Expensive”}$ ;
3. The condition  $a LIKE \text{‘%Blue_Pages:Bad%’}$  is satisfied because the annotation “Blue\_Pages:Bad” currently bound to the annotation variable  $a$  matches the pattern ‘%Blue\_Pages:Bad%’.

According to the SELECT clause of  $Q_5$ :

1. Since  $f.Name$  is “Ben Affleck” and  $f.Name$  is copied to Celebrity in the output, the Celebrity value of the tuple produced by the current binding is “Ben Affleck”.
2. Since  $r.Name$  is “Madre’s” and  $r.Name$  is copied to Restaurant in the output, the Restaurant value of the tuple produced by the current binding is “Madre’s”.

**Figure 1: The bindings for variables  $r$ ,  $f$  and  $a$  in  $Q_5$  and an explanation for why “Madre’s” is in the result according to these bindings.**

databases and transformation steps that rely on  $d$  in  $D$ .

If every value is annotated with its address, then the annotations associated with a value in the result of a query describe the provenance of that value since annotations are, by default, propagated based on provenance. We call such annotations, *provenance annotations*. In fact, DBNotes automatically maintains provenance annotations. Whenever a user writes a query  $Q$  over a database  $D$ , DBNotes creates a new database based on  $Q(D)$  that modifies each annotation by adding the address of each value to each existing annotation associated with that value. To exemplify, suppose the result of query  $Q_1$  is stored in a relation called Partial\_Restaurants (“PR” in short). DBNotes automatically attaches the address of “Madre’s” in Partial\_Restaurants to the associated annotations. We obtain the following annotations for “Madre’s” in the relation Partial\_Restaurants:

```

{(PR, t, Name).(Q1)(LATimes, t1, Res_Name):“LATimes:Good”,
 (PR, t, Name).(Q1)(Blue_Pages, t2, Name):“Blue_Pages:Bad”}.

```

The first annotation states that the value at the address  $(PR, t, Name)$  (i.e., the Name value of tuple  $t$  in Partial\_Restaurants) is copied from the value at the address  $(LATimes, t_1, Res_Name)$  through query  $Q_1$  and this value is the source value with an associated annotation “LATimes:Good”. On the whole, the two annotations tell us that the value “Madre’s” at address  $(PR, t, Name)$  was copied from the corresponding values at addresses in LATimes and Blue\_Pages via the transformation  $Q_1$ . It is from these provenance annotations that DBNotes could determine the sequence of databases and query transformations that derived a value. For the provenance of the value “Madre’s” in the integrated view Restaurants, DBNotes displays the diagram shown in Figure 2(a) (assuming query  $Q$  integrates information from Partial\_Restaurants, Restaurant\_Reviews and Seeing\_Stars and the latter two also contain information about the restaurant “Madre’s”).

In addition to provenance, DBNotes is also capable of dis-

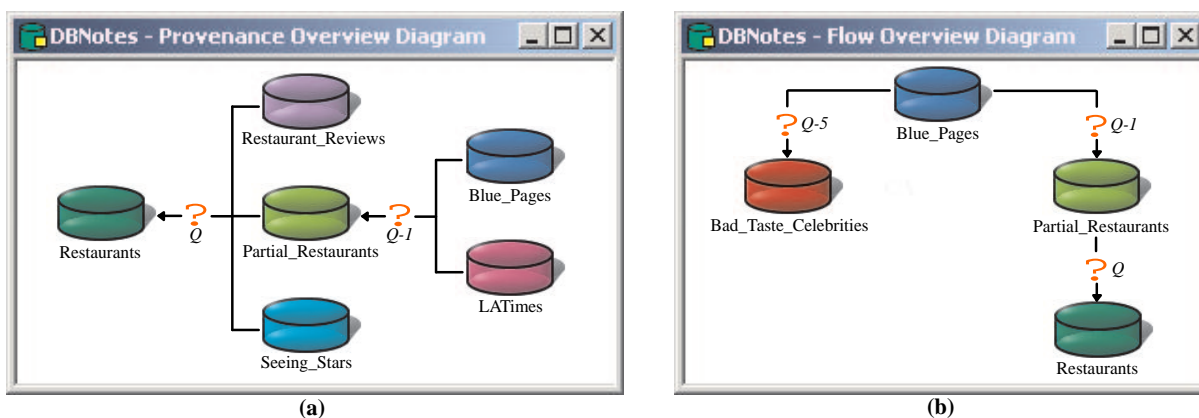


Figure 2: (a) A visual of the provenance of the value “Madre’s” in the relation Restaurants, and (b) a visual of the flow of the value “Madre’s” from the relation Blue\_Pages.

playing a visual of the flow of a value in a database. The flow of a value, however, cannot be eagerly computed, as we are unaware of what future transformations might depend on the value. Instead, DBNotes computes the flow of a value through a Transformations catalog which records information about source and target databases along with the transformations between them. As an example, when DBNotes is asked to compute the flow of the value “Madre’s” (at address (Blue.Pages,  $t_2$ , Name)) in Blue.Pages of Figure 2(a), it begins by examining the Transformations catalog in search for databases that depend on Blue.Pages. Assuming that the result of  $Q_5$  is stored in a relation called Bad.Taste.Celebrities, the two databases that are discovered in this phase are Partial.Restaurants and Bad.Taste.Celebrities. In the next phase, DBNotes scans both these relations and, by examining the annotations, discovers that values “Madre’s” at addresses (Partial.Restaurants,  $t$ , Name) and (Bad.Taste.Celebrities,  $t'$ , Restaurant) were copied from the value “Madre’s” at the address (Blue.Pages,  $t_2$ , Name) via the transformations  $Q_1$  and  $Q_5$  respectively. In another iteration, the Transformations catalog is re-examined in search for relations that are derived from Partial.Restaurants or Bad.Taste.Celebrities. In this iteration, DBNotes discovers that a value in Restaurants is copied from (Partial.Restaurants,  $t$ , Name). The value “Madre’s” in Bad.Taste.Celebrities and Restaurants did not flow elsewhere and hence the resulting figure in Figure 2(b) is displayed.

DBNotes also offers the option of “zooming in” on each transformation step, in order to see a detailed explanation (similar to what was described in Section 2.4) of how a value was produced in the output of that transformation step.

### 3. SYSTEM ARCHITECTURE

The architecture of our system is illustrated in Figure 3. The explainer module is described in Section 2.4. The provenance and the flow tracer modules are described in Section 2.5. The purpose of the Translator and Postprocessor modules is to compute the result of pSQL queries. The translator module translates a pSQL query into an SQL query against DBNotes’s underlying storage scheme for annotations. The postprocessor module merges together annotations associated with identical addresses. We refer the interested reader to [2] for more details about these two mod-

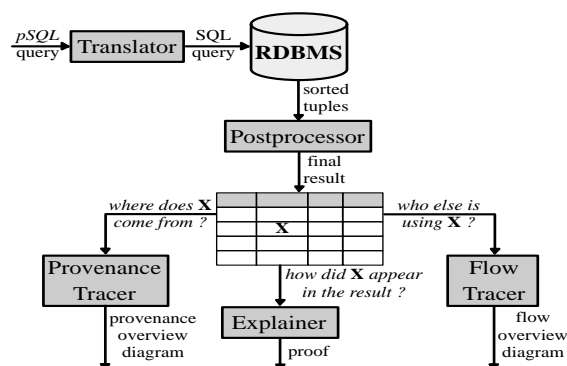


Figure 3: Architecture of DBNotes.

ules. DBNotes is currently implemented with Java v1.4.2 on top of Oracle 9i.

### 4. REFERENCES

- [1] Restaurant Reviews - Los Angeles Southern California. <http://losangeles.about.com/cs/restaurantreviews/>.
- [2] D. Bhagwat, L. Chiticariu, W. Tan, and G. Vijayvargiya. An Annotation Management System for Relational Databases. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, pages 900–911, Toronto, Canada, 2004.
- [3] P. Buneman, S. Khanna, and W. Tan. Why and Where: A Characterization of Data Provenance. In *Proceedings of the International Conference on Database Theory (ICDT)*, pages 316–330, London, United Kingdom, 2001.
- [4] Y. Cui and J. Widom. Lineage Tracing in a Data Warehousing System. In *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 683–684, San Diego, California, 2000.
- [5] Y. Cui, J. Widom, and J. Wiener. Tracing the Lineage of View Data in a Warehousing Environment. *ACM Transactions on Database Systems (TODS)*, 25(2):179–227, 2000.
- [6] Los Angeles Times. <http://www.latimes.com>.
- [7] Online Blue Pages for Restaurants. <http://www.restaurants.com>.
- [8] Seeing Stars in Hollywood. <http://www.seeing-stars.com>.
- [9] S. Yamini and A. Gupta. Spatiotemporal Annotation Graph (STAG): A Data Model for Composite Digital Objects. In *Proceedings of the International Conference on Data Engineering (ICDE) (To appear)*, Tokyo, Japan, 2005.