

Panda: A System for Provenance and Data*

Robert Ikeda and Jennifer Widom
Stanford University

1 Introduction

In its most general form, *provenance* (also sometimes called *lineage*) captures where data came from, how it was derived, manipulated, and combined, and how it has been updated over time. Provenance can serve a number of important functions:

- **Explanation.** Users may be particularly interested in or wary of specific portions of a derived data set. Provenance supports “drilling down” to examine the sources and evolution of data elements of interest, enabling a deeper understanding of the data.
- **Verification.** Derived data may appear suspect—due to possible bugs in data processing and manipulation, because the data may be stale, or even due to maliciousness. Provenance enables auditing how data was produced, either for verifying its correctness, or for identifying the erroneous or outdated source data or processing nodes that are responsible for erroneous or outdated output data.
- **Recomputation.** Having found outdated or incorrect source data, or buggy processing nodes, we may want to correct the errors and propagate the corrections forward to all “downstream” data that are affected. Provenance helps us recompute only those data elements that are affected by the corrections.

There has been a large body of very interesting work in lineage and provenance over the past two decades. (Space limitations preclude detailed discussion of previous work here, but surveys are presented in, e.g., [3, 4, 5], and a growing list is maintained on our Panda project web site [1].) Nevertheless, we believe there are still many limitations and open areas. Specifically:

1. Most work has been either: *data-based*, in which fine-grained provenance of data elements is tracked based on well-defined, transparent properties of data models and query languages; or *process-based*, in which coarse-grained provenance is tracked, typically involving workflows and data at the schema level.
2. Often the primary focus is on *modeling* and *capturing* provenance: How is provenance information represented? How is it generated? There has been considerably less work on *querying* provenance: What can we do with provenance information once we’ve

captured it?

3. Many projects have focused on specific functions or application domains, rather than developing a general provenance system that can be used for different purposes and across domains.

In the nascent *Panda* (for “*provenance and data*”) project at Stanford, our goal is to fill these gaps. Specifically, we want to:

1. Seamlessly merge data-based and process-based provenance, so that the two types of provenance can be combined (e.g., workflows that combine “opaque” processing nodes with well-understood relational queries and transformations). We also want to develop a model and system that offers users a full range from fine-grained to coarse-grained provenance.
2. Define a set of useful operators for taking advantage of provenance after it has been captured, as well as a general-purpose language for querying and analyzing provenance, and for combining provenance with relevant data.
3. Develop a general-purpose open-source system that is flexible and configurable enough to be used for a wide variety of applications. The system will support its own mechanisms for provenance capture, storage, operators, and queries, while also offering interfaces for coupling with outside data sources, processes, and systems.

2 Running Example

We use a detailed fictitious example to motivate our plans for the Panda system. Consider *ClothCo*, a mail-order clothing company that is trying to decide which items to feature most prominently in its upcoming catalog. ClothCo runs an analytics workflow to predict which items will have the largest demand among its customers. The workflow is shown in Figure 1.

The initial input to the workflow is customer lists $\text{CustList}_1, \text{CustList}_2, \dots, \text{CustList}_n$, obtained from ClothCo’s own databases and from partnerships with other companies. The lists contain names, addresses, and possibly additional attributes (e.g., gender, income), all of which are used for list deduplication and for buying-behavior predictions. The workflow involves the following steps:

- **Dedup:** Records from the customer lists are *dedupli-*

*This work is supported by the National Science Foundation under grants IIS-0414762 and IIS-0904497.

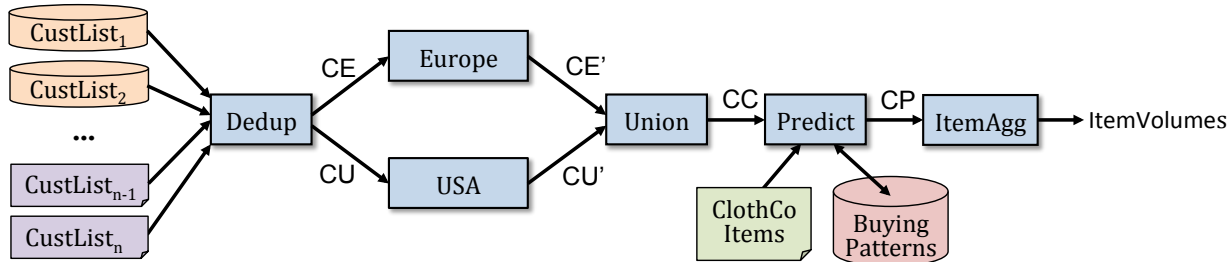


Figure 1: Analytics workflow example.

cated so that customers represented in multiple lists are counted only once. The deduplicated customer records are then partitioned into a European customer set (CE) and a USA customer set (CU). (If **Dedup** cannot determine that a customer is European, it by default assumes the customer is American.)

- **Europe / USA / Union:** The deduplicated records are routed either to **Europe** or **USA**, where the records' addresses are "canonicalized" into a standard form using existing software specialized for the region. The output lists CE' and CU' are unioned back to form canonicalized customer list CC.
- **Predict:** The next stage in the workflow predicts which ClothCo items customers are most likely to purchase. Specifically, for each customer in input list CC, and for each item ClothCo stocks, the **Predict** module consults a **Buying Patterns** database and attaches to the customer-item pair a likelihood that the customer will buy the item. **Predict**'s output CP is a set of customer-item-probability triples.
- **ItemAgg:** Finally, from the CP triples, **ItemAgg** aggregates the predicted demand for each of ClothCo's items. The output, **ItemVolumes**, contains a list of items and predicted sales volumes for each one.

Now suppose a ClothCo analyst runs the workflow and is surprised to find that the item in **ItemVolumes** predicted to have the highest demand is a cowboy hat, despite the fact that ClothCo's target customers rarely hail from the southern United States. Noticing this anomalous result, the analyst would like to find out why the predicted demand for cowboy hats is so high.

We envision that the Panda system will support the following interaction. Tracing the provenance of the cowboy-hat record in **ItemVolumes** back one step to CP yields a large data set, but with a simple provenance query the analyst discovers that the majority of the customers predicted to want the cowboy hat live in Paris, Texas. Something is clearly wrong: the population of Paris, Texas is only 25,000, and ClothCo doesn't cater to this demographic anyway.

Further tracing the provenance of the relevant CP records, the analyst discovers that most of them were

processed by **USA**, but they came originally from a French customer list. It turns out because the French list did not specify "France" explicitly in its addresses, **Dedup** mistakenly routed these customer records to **USA**, which added "Texas" during the canonicalization process.

To fix the error, a simple new module is inserted that appends "France" to the addresses from the problematic list. Using provenance, when the modified workflow is rerun with the new module, only the item predictions potentially affected by the modifications are recomputed.

While this example captures a large number of the most important capabilities we want to include in the Panda system, there are a few missing. For example, we want to also capture provenance from *human-generated* data edits and manipulations, along with the computer-generated ones. Also, we expect that capturing and managing *evolving versions* of data, including long-term and frequent evolutions, will be a critical feature.

3 Processing Nodes and Provenance Capture

The workflow in our example is representative of the variety of processing node types that Panda should be able to support. Specifically, we want to handle processing nodes that vary from fully transparent (e.g., relational queries) to fully opaque (e.g., private code, calls to external services). We plan to define an interface by which any type of processing node (including a human) can write out provenance information in a well-specified uniform fashion, so that it can be used by the Panda system.

In our example, the **Union** and **ItemAgg** nodes perform standard relational operations. For relational nodes, there is a great deal of past work that can be applied for capturing and tracing provenance automatically and efficiently (see [4] for a survey). Consider **ItemAgg**, for example, which is a standard relational *group-by aggregation* operator. As defined by past work, the provenance of an item *I* output by **ItemAgg** consists of all items in CP with the same *group-by* attribute as *I*.

Now consider **Dedup**, an opaque processing node. We cannot rely on the automatic methods for relational queries mentioned above to capture and trace provenance

for **Dedup**. Either **Dedup** must be instrumented to write out provenance information as it executes (presumably in the form of mappings between deduplicated records and their original customer records), or it must provide some sort of procedure for provenance-tracing, or in the worst case we cannot determine fine-grained provenance at all. Since the Panda system is designed specifically for provenance and data, we prefer the first approach.

One of the most important first steps in Panda will be to formalize a model that its provenance information must conform to. On one end, our model could be a simple bipartite graph structure connecting input and output data elements. Or, we could adopt the much higher-level and more “semantic” *Open Provenance Model* [2]. A likely scenario is some combination of these two approaches, in keeping with our goal of combining *data-based* and *process-based* provenance as discussed in Section 1. Combining these two approaches in some fashion may also serve our goal (also discussed in Section 1) of allowing provenance to be captured at a variety of granularities, ranging from individual data elements, to schema-level data sets, to processing-node signatures. Once the provenance model is defined, we then must define a uniform interface by which all types of processing nodes can create and manipulate provenance, both manually and automatically.

Once provenance has been captured, what are we going to do with it? Our overall goal is to support the features mentioned in Section 1: explanation, verification, and recomputation. To do so, we will introduce a set of basic operations (Section 4) and an ad-hoc query language (Section 5).

4 Provenance Operations

We plan to offer at least two methods of using the provenance captured in the Panda system:

- A set of built-in *operations* that can be used on their own or as building blocks for higher-level functionality.
- A full-featured *query language* that can be used to pose ad-hoc queries and analyses over provenance information, and over provenance information combined with relevant data.

For basic built-in operations, we envision at least the following two:

- **Backward tracing.** Given a derived data element D , where did D come from? That is, what data elements and/or processing contributed to D ? In our running example, we used backward tracing to go from the output cowboy-hat record C to the record-set \mathcal{R} in input list **CP** from which C was derived. Then we used further backward tracing to determine that most records producing \mathcal{R} came through the **USA** process-

ing node and originated from a specific French customer list.

- **Forward tracing.** Given an input or derived data element D , where did D subsequently go? That is, what processing nodes did D later pass through and what data elements were produced by it? In our running example, we can use forward tracing to determine all of the item predictions that were affected by French customers erroneously passing through the **USA** processing nodes.

From these two basic operations, we can layer on additional functionality, for example:

- **Forward propagation.** If an input or derived data element D changes, propagate the change to everything it affects. Clearly this function will rely on forward tracing. In our running example, once we correct the problem with our French customer list, we can use forward propagation to recalculate only those item predictions affected by the correction.
- **Refresh.** Given a derived data element D , check if D is still valid. If it is not, refresh it to its new valid value. Clearly this function will rely on both backward tracing and forward propagation. In our running example, suppose we correct the error with the French addresses, but ironically a celebrity then begins wearing cowboy hats. Suddenly cowboy hat sales soar, with corresponding modifications in our **Buying Patterns** database. We decide to once again investigate the cowboy hat prediction in **ItemVolumes**, this time asking for it to be refreshed to ensure we get the latest predicted demand.

5 Provenance Queries

The operations in the previous section are not specific to a particular application, but they still encode very specific provenance-based functionality. We believe that in addition to a set of common operations, it will be very useful for a general-purpose provenance system to support a declarative ad-hoc query language, similar to what is provided by database management systems.

We can only begin to develop a query language once our provenance model is fully defined, but we do have a few guiding principles. One is that we want our query language to operate over provenance and data seamlessly, so that they can be combined in useful ways. Second, we want our query language to be compact and intuitive for basic queries, with specific features for additional expressiveness, similar in spirit to SQL. Finally, the language must be amenable to finding efficient query execution plans, again in a database-system style.

Using our running example, here are a few queries in English to demonstrate the type of functionality we would like to support in a query language:

- From our specific example: Aggregate the data elements in CP contributing to the cowboy hat output element in ItemVolumes.
- Determine which customer list contributes the most to the top 100 predicted items.
- Considering only customers from a specific list, which items have significantly higher demand among those customers than among the general population?
- Which customers have more duplication in our original customer lists—those that are eventually processed by USA, or those that are processed by Europe?

6 System Issues

Building a full-featured generic system for managing provenance and data together will require a significant effort just for the basic functions: provenance capture and storage, built-in operations like those in Section 4, and general-purpose query processing. In addition, we see a number of specific opportunities and challenges.

Query-Driven Capture

The most general approach is to capture provenance in support of any possible operation or query that may subsequently be performed using it. However, if there is a known restricted set of operations and queries to be performed, we may be able to streamline what is captured.

Eager vs. Lazy

Even when supporting arbitrary operations and queries, there may be a choice between *eager* and *lazy* provenance capture. It's a typical space-time tradeoff: eager capture occupies space but speeds up operations and queries. (Note that capturing all possible provenance at a fine-grained level could entail enormous amounts of space.) Also, similar to materialized views, eager capture may incur an update cost if we wish to keep provenance up-to-date. The choice between eager and lazy might be made on a per-processing-node basis. For example, eager capture makes sense for **Dedup**, since the overhead of recording the source customer records contributing to each deduplicated record may be fairly low, while recomputing that information may be very expensive (perhaps requiring deduplicating the entire data set again). On the other hand, for **ItemAgg**, if we have access to intermediate data set CP, then computing the provenance of an output item in **ItemAgg** is a simple selection query over CP.

Intermediate Results

Along similar lines, in a workflow such as our ClothCO example, we may or may not wish to store all of the in-

termediate data sets. Storing these data sets may be very helpful for provenance operations (such as finding the provenance of an **ItemAgg** output element by querying CP, as above) but the storage overhead may be high, and intermediate results may not always be necessary. For example, if we only wish to consider provenance from final output elements back to initial input data, and we take a fully eager approach, then no intermediate data sets are necessary. Overall, we envision a suite of interesting optimization problems involving decisions about intermediate data sets and eager versus lazy computation.

Fine-Grained vs. Coarse-Grained

Another trade-off, perhaps intertwined with the previous two, is between fine-grained and coarse-grained provenance. Even if it is prohibitive to compute provenance eagerly at the data-element level, it may be helpful to record some provenance at the data-set level: some operations and queries may only need coarse-grained provenance, or perhaps coarse-grained provenance can be used to support later computation of fine-grained provenance. Suppose, for example, for each input list we record only the percentage of its customers that are eventually sent to processing node **Europe**, versus those sent to processing node **USA**. Queries may ask for this information directly, or in some cases (e.g., lists that are 100% **Europe** or 100% **USA**), this schema-level information may be used to optimize provenance queries involving individual records.

Approximation

Another avenue for dealing with prohibitively large fine-grained provenance may be to support *approximate* provenance capture, operations, and/or query results.

Acknowledgments

We are grateful to Parag Agrawal, Diana MacLean, Abhijeet Mohapatra, Raghotham Murthy, Aditya Parameswaran, Hyunjung Park, and Alkis Polyzotis, for many useful Panda meeting discussions that have been instrumental in forming the direction for our system. We thank Philip Guo for his helpful comments on the paper.

References

- [1] <http://i.stanford.edu/panda>.
- [2] The Open Provenance Model (v1.01). July 2008. <http://eprints.ecs.soton.ac.uk/16148/>.
- [3] BOSE, R., AND FREW, J. Lineage retrieval for scientific data processing: a survey. *ACM Comput. Surv.* 37, 1 (2005), 1–28.
- [4] CHENEY, J., CHITICARIU, L., AND TAN, W.-C. Provenance in databases: Why, how, and where. *Foundations and Trends in Databases* 1, 4 (2009), 379–474.
- [5] SIMMHAN, Y. L., PLALE, B., AND GANNON, D. A survey of data provenance in e-science. *SIGMOD Rec.* 34, 3 (2005), 31–36.