

Extracting Queries by Static Analysis of Transparent Persistence*

Ben Wiedermann and William R. Cook

Department of Computer Sciences, The University of Texas at Austin

{ben,wcook}@cs.utexas.edu

Abstract

Transparent persistence promises to integrate programming languages and databases by allowing procedural programs to access persistent data with the same ease as non-persistent data. When the data is stored in a relational database, however, transparent persistence does not naturally leverage the performance benefits of relational query optimization. We present a program analysis that combines the benefits of both approaches by extracting database queries from programs with transparent access to persistent data. The analysis uses a sound abstract interpretation of the original program to approximate the data traversal paths in the program, and the conditions under which the paths are used. The resulting paths are then converted into a query, and the program is simplified by removing redundant tests. We study an imperative kernel language with read-only access to persistent data, and identify the conditions under which the transformations can be applied. This analysis approach promises to combine the software engineering benefits of transparent data persistence with the performance benefits of database query optimization.

1. Introduction

The effective integration of programming languages and databases is a long-standing and critical open problem. From a programming language viewpoint, databases manage *persistent* data, which has a lifetime longer than the execution of an individual program. Ideally a unified programming model should be applicable to both persistent and non-persistent data. This goal has been pursued for the last 30 years in numerous forms, including orthogonal persistence [2, 3, 4, 22, 27], object-relational mapping [14, 19, 25, 31], and object-oriented databases [11, 13, 24]. Despite differences in particular details, these approaches all share the goal of *transparent persistence*—a programming paradigm wherein the programmer need not distinguish between persistent and non-persistent values.

Transparent persistence can be added to most any language by extending the concepts of automatic memory management & garbage collection to the management of persistent data: by identifying a persistent root object, any object or value reachable from the root is also persistent [4]. For example, the Java program in Fig. 1 manipulates a collection of employee objects associated with a root object. If root identifies a persistent store of objects, then the employee objects may be loaded from that store. However, the program’s function remains independent of whether root is persistent or not.

This kind of transparent persistence does not easily leverage the power of database query optimization. Database optimizations work best when records are loaded in bulk and conditions for selecting records are executed in the database rather than the procedural program. The mismatch between one-at-a-time processing in procedural language and bulk data processing in query operations is called “impedance mismatch” [23]. To solve this problem,

```
for (Employee e : root.employees) {
  if (e.salary > 65000) {
    print (e.name + ": " + e.manager.name);
  }
}
```

Figure 1. A program using transparent persistence.

```
// define an explicit query
String query = "from Employee e
  left join fetch e.manager
  where e.salary > 65000";
// execute the query
List result = session.createQuery(query);
for (Employee e : result.list()) {
  // no test required: all elements already satisfy
  // the condition salary > 65000
  print (e.name + ": " + e.manager.name);
}
```

Figure 2. Explicit query execution using Hibernate.

many persistence models allow programmers to execute explicit queries. For example, Fig. 2 uses Hibernate, an object-relational mapping tool, and its query language HQL [19] to execute an explicit query. The query returns only employees with salary greater than \$65,000; the prefetch clause **left join fetch e.manager** indicates that each employee’s manager should also be loaded. The **if** statement in Fig. 1 is not needed in Fig. 2 because the query’s **where** clause ensures the query only returns employees for which the test is true.

Although the programs in Fig. 1 and Fig. 2 print the same results, they have different performance and software engineering benefits. In the transparent persistence version, all employees will be loaded even though only those with salary greater than \$65,000 are printed. Manager objects will be loaded individually, because the persistence layer cannot predict which ones will be needed. In the Hibernate version, the underlying relational query optimizer will likely use an index to locate all employees whose salary is greater than \$65,000. The optimized version runs in time proportional to the size of the query result, rather than the total number of employees and may be orders of magnitude faster [8].

Despite its performance benefits, there are some drawbacks to the Hibernate version. Query strings are not checked at compile time for syntax or type safety, and they reduce modularity and increase the complexity of programming. Proposals to address these problems [6, 10, 18] either reduce or do not address the transparency of persistence. There is also a subtle dependency between the query and the code: the prefetch clause is logically redundant with the use of the employee’s manager in the **print** method.

This paper describes a static analysis technique that allows a programming language with transparent persistence to leverage the

*This work was supported by the National Science Foundation under Grant No. 0448128.

$$\begin{aligned}
l &\in \text{Variable} \\
f &\in \text{Field} \\
e \in \text{Expression} &::= l \mid e.f \mid \text{op}_n(e_1, \dots, e_n) \\
\text{op}_0 \in \text{Constant} &::= \mathbf{true} \mid \mathbf{false} \mid \text{number} \mid \text{string} \\
\text{op}_1 &::= \neg \mid \mathbf{print} \\
\text{op}_2 &::= \wedge \mid \vee \mid > \mid < \mid = \mid \geq \mid \leq \mid \neq \\
c \in \text{Command} &::= \mathbf{skip} \mid l := e \mid c; c \\
&\quad \mid \mathbf{if } e \mathbf{ then } c \mathbf{ [else } c \mathbf{]} \\
&\quad \mid \mathbf{for } l \mathbf{ in } e \mathbf{ do } c
\end{aligned}$$

Figure 3. Syntax of a persistent data kernel language.

power of query optimization. Our approach automatically partitions programs by extracting data traversals and conditions into a query, and removing them from the program—essentially transforming the program in Fig. 1 into the program in Fig. 2. Our analysis consists of three parts. The first part (Section 3) identifies the traversals used in the program; these traversals specify the data that must be loaded by the query. The second part (Section 4) identifies the conditions under which data is used, for example by a **print** method, so that the conditions can be included in the query. In the final part (Section 5) the individual conditions on the use of fields are promoted to apply to entire records, and a query is created. This final step also modifies the program to use the results of the query, and eliminates the redundant **if** statement.

The primary contribution of this paper is a new approach to optimization of transparent persistence, by extracting queries from imperative programs. This result is based on a sound abstract interpretation of programs, together with techniques for converting the resulting abstract values into queries and simplifying the original program. We have developed a prototype implementation of the analysis and applied it to simple examples to demonstrate its viability. While this work re-opens an important line of research, there are many topics left to future work. In particular, we do not analyze the performance of the analysis or the transformed programs—although the performance gains from query optimization are well-known. We have not applied the analysis to large programs with procedures, or addressed the problem of identifying where in a large program the analysis should be applied. Complex query behaviors, like aggregation, exists queries, and database mutations (creations, updates, and deletions) are not considered. We expect that the current work will serve as a solid foundation for ongoing work on these problems, with the goal of combining the software engineering benefits of transparent persistence and the performance benefits of query optimization.

2. A Kernel Language with Persistent Data

We study a simple imperative language with records and access to persistent data. The persistent data is an instance of an Entity-Relationship Model [9], which provides natural mappings to both relational databases [5] and class models in UML/object-oriented programming [34]. A persistent value is a record, or labeled product, whose fields are either basic values or references to other records (these are called “attributes” and “relationships” in an ER model). A reference/relationship field may be either single-valued or multi-valued. Multi-valued relationships correspond to collection objects in object-oriented programming. The language expresses key concepts in practical orthogonally persistent object-oriented languages, but it also has several restrictions. Only the structural representation of data is considered, not behavioral meth-

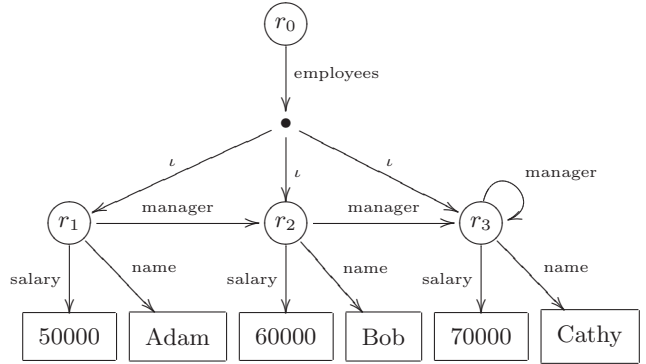


Figure 4. An object graph example.

ods, and the language contains no procedures. We do not model the three-valued logic of null values, but assume that a value is defined for every persistent element a program accesses. While the language supports imperative update of local variables, persistent data is read-only. We believe these restrictions to be reasonable, as the current work is designed to introduce a technique for extracting procedural queries. Section 7 discusses extensions to support interprocedural analysis, analysis of more complicated query idioms, and creation, update or, deletion of persistent data.

2.1 Syntax

The abstract syntax of the kernel language is defined in Fig. 3. The traversal expression $e.f$ projects a field f of a record e . The value of $e.f$ can be a simple value, or references to one or more records.

Persistent data is introduced through a special root identifier (variable) that refers to a record representing persistent data [4]. Any value that is reachable from the root is also persistent. As mentioned above, no constructs create or modify persistent records; all records are loaded from the persistent store.

Primitive functions op_n have a specified number of arguments n . Infix notation is used where appropriate.

The **for** command allows iteration over the elements of a collection. For simplicity, iteration is supported only for multi-valued database fields; however, the language could easily be extended to allow collections of basic program values.

A simple static type system for records is assumed for this language [29]; programs are assumed to be well typed.

2.2 Values

A program operates over the domain:

$$v \in \text{Value} = \text{Basic} + \text{RecordID} + \text{RecordID}^*$$

where *Basic* is the domain of basic values (integers, Booleans, and strings), and *RecordID* is the domain of *record identifiers* that reference persistent database values. When a program traverses a record identifier, a runtime function $\text{Load} :: \text{RecordID} \times \text{Field} \rightarrow \text{Value}$ retrieves the corresponding record’s field value(s). A special record identifier r_0 corresponds to the store’s root, and the store’s structure is the graph formed by the transitive closure of traversals from r_0 .

Figure 4 illustrates a persistent object graph, against which the program in Fig. 1 can be evaluated. The graph’s solid dot denotes a collection of record identifiers, where the target of each outgoing edge is a member of the collection. Each of these edges is implicitly labeled with an *iterator field name* l , which identifies distinct elements of the collection.

$$\begin{array}{c}
\langle l, \sigma \rangle \rightarrow \sigma[l] \quad (\text{S-VAR}) \\
\langle \text{skip}, \sigma \rangle \rightarrow \sigma \quad (\text{S-SKIP}) \\
\frac{\langle e, \sigma \rangle \rightarrow r}{\langle e.f, \sigma \rangle \rightarrow \text{Load}(r, f)} \quad (\text{S-TRAVERSE}) \\
\frac{\langle e_i, \sigma \rangle \rightarrow v_i \quad \text{for } i \in \{1, \dots, n\}}{\langle \text{op}_n(e_1, \dots, e_n), \sigma \rangle \rightarrow f_{\text{op}_n}(v_1, \dots, v_n)} \quad (\text{S-OP}) \\
\frac{\langle c_1, \sigma \rangle \rightarrow \sigma' \quad \langle c_2, \sigma' \rangle \rightarrow \sigma''}{\langle c_1; c_2, \sigma \rangle \rightarrow \sigma''} \quad (\text{S-SEQ}) \\
\frac{\langle e, \sigma \rangle \rightarrow v}{\langle l := e, \sigma \rangle \rightarrow [l \mapsto v]\sigma} \quad (\text{S-ASSIGN}) \\
\frac{\langle e, \sigma \rangle \rightarrow \text{true} \quad \langle c_1, \sigma \rangle \rightarrow \sigma'}{\langle \text{if } e \text{ then } c_1 \text{ else } c_2, \sigma \rangle \rightarrow \sigma'} \quad (\text{S-IFT}) \\
\frac{\langle e, \sigma \rangle \rightarrow \text{false} \quad \langle c_2, \sigma \rangle \rightarrow \sigma'}{\langle \text{if } e \text{ then } c_1 \text{ else } c_2, \sigma \rangle \rightarrow \sigma'} \quad (\text{S-IFF}) \\
\frac{\langle c, [l \mapsto r_i]\sigma_i \rangle \rightarrow \sigma_{i+1} \quad \text{for } i \in \{1, \dots, n\}}{\langle \text{for } l \text{ in } e \text{ do } c, \sigma_1 \rangle \rightarrow \sigma_{n+1}/l} \quad (\text{S-FOR})
\end{array}$$

Figure 5. Operational semantics of the kernel language.

2.3 Semantics

Figure 5 defines a big-step operational semantics for the kernel language. A store σ maps variables to values. The evaluation relations for expressions $\langle e, \sigma \rangle \rightarrow v$ and commands $\langle c, \sigma \rangle \rightarrow \sigma'$ follow standard form. All programs begin computation with a store σ_0 that maps the variable *root* to the persistent value r_0 .

Rule S-VAR retrieves a variable’s value from the store. For operations, rule S-OP evaluates the operands, then applies the operator’s function to the result. The functions f_{true} , f_{\neg} , $f_{<}$, etc., have the standard, mathematical meanings. Note that these functions are defined so that they return only primitive values, not record identifiers. Rules S-SKIP, S-ASSIGN, S-SEQ, S-IFT and S-IFF are standard. The expression $[l \mapsto v]\sigma$ denotes σ updated so that $\sigma[l] = v$.

Rule S-TRAVERSE loads persistent data. If expression e evaluates to the record identifier r , then the expression $e.f$ evaluates to the result of calling $\text{Load}(r, f)$.

Rule S-FOR defines iteration over a collection of record identifiers. For each record identifier r_i in the collection, the **for** command’s body c is evaluated in a new store σ_i that maps the loop variable l to record identifier r_i . The result of the entire command is the final store produced σ_{n+1} . A loop variable is defined only in the loop’s body, so the variable is removed from the final store.

One subtle difference between the semantics of object-oriented programming languages and relational databases is that programming languages often impose an order on iterated collections, whereas databases do not. For our purposes, we assume a default order exists for every database collection and that programs iterate over collections in that order.

Output from the **print** function is modeled by a special variable **output**; the **print** function simply concatenates onto the end of this variable.

Evaluating the example program in Fig. 1 against the persistent data in Fig. 4 generates a final store with the following mappings:

$$\{\text{root} \mapsto r_0, \text{output} \mapsto \text{“Cathy : Cathy”}\}$$

2.4 Operational Semantics with Explicit Used-Set

Our analysis summarizes the set of persistent values a program uses. These values—which we refer to as the program’s *used-set*—can then be loaded in bulk before the program needs them. The operational semantics of the base language is extended in Fig. 6 to keep track of a computation’s used-set. The modified semantics has evaluation relations $\langle e, \sigma \rangle \rightarrow \langle v, \rho \rangle$ and $\langle c, \sigma \rangle \rightarrow \langle \sigma, \rho \rangle$ where ρ is the set of database values that were loaded during the entire computation. S-TRAVERSE is the only rule that loads database values, so the rule adds the newly loaded values ρ_f to the set.

All other rules are modified to collect the loaded values for any sub-computation, where $\bigcup \rho_i$ is shorthand for $\bigcup_{i \in \{1, \dots, n\}} \rho_i$. Evaluating our running example with the extended semantics generates the following set of database values:

$$\{r_0, r_1, 50000, r_2, 60000, r_3, 70000, \text{“Cathy”}\}$$

3. Analyzing Traversals

Traversal analysis is an abstract interpretation [12] of the operational semantics in which database values are replaced by paths. The *path* corresponding to a database value is the sequence of field names traversed to load that value. Note that many database values may have the same path; for example, in Fig. 4 the paths to record identifiers r_1 , r_2 , and r_3 are identical. Many paths may lead to the same database value; for example, the value $r_2.\text{name}$ can be reached by following either $\text{employees}.l.\text{name}$ or $\text{employees}.l.\text{manager.name}$.

Abstract interpretation uses abstraction and concretization functions to specify the relationship between abstract and concrete values. Given an abstract path, its concretization is the set of database values that can be reached by following the path. If the path includes a collection field, the concretized result includes all the traversals of items in the collection. Thus concretization corresponds to interpreting the path as a query against the database.

The analysis is conservative (sound), so that the concretization of a path may return a larger set of database values than actually appear in the concrete execution of the program. However, because a program typically operates over a small subset of a large database, the amount of data represented by the concretized paths should be small relative to the overall database size. Soundness justifies bulk loading of data before executing the program; precision gives better performance. The analysis in this section uses a loose approximation, but serves as a useful foundation for the more precise analysis in Section 4.

3.1 Abstract Value Domain

The abstract value domain is

$$\hat{v} \in \widehat{\text{Value}} = \wp(\text{Path}) + \top$$

where \wp is the powerset operator and Path is the finite set of paths a program can traverse. All non-path values are abstracted away as \top . The domain forms a finite, complete lattice ordered by the subset relation (\subseteq).

To ensure that the analysis terminates, we restrict Path to be a finite subset of the all possible sequences of fields Field^* . One possible finite subset is the set of paths in which each field name

$\langle l, \sigma \rangle \rightarrow \langle \sigma[l], \emptyset \rangle$	(U-VAR)
$\langle \mathbf{skip}, \sigma \rangle \rightarrow \langle \sigma, \emptyset \rangle$	(U-SKIP)
$\frac{\langle e, \sigma \rangle \rightarrow \langle r, \rho_e \rangle \quad \rho_f = \text{Load}(r, f)}{\langle e.f, \sigma \rangle \rightarrow \langle \rho_f, \rho_e \cup \rho_f \rangle}$	(U-TRAVERSE)
$\frac{\langle e_i, \sigma \rangle \rightarrow \langle v_i, \rho_i \rangle \quad \mathbf{for} \ i \in \{1, \dots, n\}}{\langle \text{op}_n(e_1, \dots, e_n), \sigma \rangle \rightarrow \langle f_{\text{op}_n}(v_1, \dots, v_n), \cup \rho_i \rangle}$	(U-OP)
$\frac{\langle c_1, \sigma \rangle \rightarrow \langle \sigma', \rho_1 \rangle \quad \langle c_2, \sigma' \rangle \rightarrow \langle \sigma'', \rho_2 \rangle}{\langle c_1; c_2, \sigma \rangle \rightarrow \langle \sigma'', \rho_1 \cup \rho_2 \rangle}$	(U-SEQ)
	$\frac{\langle e, \sigma \rangle \rightarrow \langle v, \rho_e \rangle}{\langle l := e, \sigma \rangle \rightarrow \langle [l \mapsto v] \sigma, \rho_e \rangle}$ (U-ASSIGN)
	$\frac{\langle e, \sigma \rangle \rightarrow \langle \text{true}, \rho_e \rangle \quad \langle c_1, \sigma \rangle \rightarrow \langle \sigma', \rho_c \rangle}{\langle \mathbf{if} \ e \ \mathbf{then} \ c_1 \ \mathbf{else} \ c_2, \sigma \rangle \rightarrow \langle \sigma', \rho_e \cup \rho_c \rangle}$ (U-IFT)
	$\frac{\langle e, \sigma \rangle \rightarrow \langle \text{false}, \rho_e \rangle \quad \langle c_2, \sigma \rangle \rightarrow \langle \sigma', \rho_c \rangle}{\langle \mathbf{if} \ e \ \mathbf{then} \ c_1 \ \mathbf{else} \ c_2, \sigma \rangle \rightarrow \langle \sigma', \rho_e \cup \rho_c \rangle}$ (U-IFB)
	$\frac{\langle e, \sigma \rangle \rightarrow \langle \{r_1, \dots, r_n\}, \rho_1 \rangle \quad \langle c, [l \mapsto r_i] \sigma_i \rangle \rightarrow \langle \sigma_{i+1}, \rho_{i+1} \rangle \quad \mathbf{for} \ i \in \{1, \dots, n\}}{\langle \mathbf{for} \ l \ \mathbf{in} \ e \ \mathbf{do} \ c, \sigma_1 \rangle \rightarrow \langle \sigma_{n+1}/l, \cup \rho_i \rangle}$ (U-FOR)

Figure 6. Operational semantics, extended to collect used-sets.

$\pi_f = \begin{cases} \top & f^t \in \pi_e \\ \{p.f^t \mid p \in \pi_e\} & \text{otherwise} \end{cases}$	(A-TRAVERSE)
$\frac{\langle e_i, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}_i, \pi_i \rangle \quad \mathbf{for} \ i \in \{1, \dots, n\}}{\langle \text{op}_n(e_1, \dots, e_n), \hat{\sigma} \rangle \hat{\rightarrow} \langle \top, \sqcup \pi_i \rangle}$	(A-OP)
$\frac{\langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}, \pi_e \rangle \quad \langle c_2, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{\sigma}_2, \pi_2 \rangle}{\langle \mathbf{if} \ e \ \mathbf{then} \ c_1 \ \mathbf{else} \ c_2, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{\sigma}_1 \sqcup \hat{\sigma}_2, \pi_e \sqcup \pi_1 \sqcup \pi_2 \rangle}$	(A-IF)
	$\langle l, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{\sigma}[l], \emptyset \rangle$ (A-VAR)
	$\frac{\langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}, \pi \rangle}{\langle l := e, \hat{\sigma} \rangle \hat{\rightarrow} \langle [l \mapsto \hat{v}] \sqcup \hat{\sigma}, \pi \rangle}$ (A-ASSIGN)
	$\frac{\langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \pi_e, \pi \rangle \quad \pi_l = \{p.l^l \mid p \in \pi_e\}}{(\hat{\sigma}', \pi') = \sqcup \{ \langle \mathbf{do} \langle c, l, \pi_l \rangle^n \hat{\sigma}, \emptyset \rangle \mid n \in \mathbb{N} \}}$ (A-FOR)
	$\frac{\langle c, [l \mapsto \pi_l] \sqcup \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{\sigma}', \pi' \rangle}{\mathbf{do} \langle c, l, \pi_l \rangle (\hat{\sigma}, \pi) = (\hat{\sigma}', \pi \sqcup \pi')}$ (A-DO)

Figure 7. Path-based abstract interpretation for approximating used-sets.

occurs at most once — any other paths would be abstracted as \top . With this domain, the expression `root.manager.manager` would be assigned abstract value \top , even though it is a finite path. On the other hand, in the following program `x` should be assigned value \top , because it produces a path of unbounded length:

```
for employee in root.employees do
  x := x.manager;
```

Note that this program is not very useful because there is usually not a meaningful relationship between the number of employees in a list and the depth of a manager traversal.

To distinguish these cases, we create the domain *Path* by labeling each field in the program, and considering all paths in which each labeled field occurs at most once. With labels, the first expression `root.manager1.manager2` has a finite path, but `x` in the example above would still be assigned abstract value \top . More expressive abstractions for representing infinite paths would certainly be useful, for example in the analysis of recursive procedures. However, such abstractions are beyond the scope of the current work.

3.2 Abstract Semantics for Traversals

The operational semantics in Fig. 7 computes the paths a program may traverse. For brevity we omit the rules A-SKIP and A-SEQ, which merely collect values and paths for subcomputations.

The abstract semantics has evaluation relations for expressions $\langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}, \pi \rangle$ and commands $\langle c, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{\sigma}', \pi \rangle$, where \hat{v} is an abstract value, $\hat{\sigma}$ maps variables to abstract values, and π is the set of paths traversed by a computation.

Rule A-TRAVERSE defines how field traversal extends a path. In evaluating $e.f^t$, if e yields the set of paths π_e then the result of the traversal extends each path in π_e with the labeled field f^t . The traversal rule also includes a widening clause to ensure the analysis converges. If the field label f^t already appears in one of the paths that e may traverse, then the program traverses an invalid path in the sense described by Section 3.1. In this case, the analysis approximates the expression's traversals with \top .

Rule A-OP gives \top as the abstract value for any operation, because the analysis ignores basic values. In Section 4 we extend the analysis to include abstractions for basic values.

Rule A-IF combines the paths traversed in evaluating the condition and the two command branches of an **if** statement. The join $\hat{\sigma}_1 \sqcup \hat{\sigma}_2$ of two maps $\hat{\sigma}_1$ and $\hat{\sigma}_2$ is a map that includes all elements of both maps:

$$(\hat{\sigma}_1 \sqcup \hat{\sigma}_2)[l] = \hat{\sigma}_1[l] \sqcup \hat{\sigma}_2[l]$$

If $\hat{\sigma}_i$ is not defined for l , then $\hat{\sigma}_i[l] = \emptyset$.

Rule A-VAR retrieves a variable's abstract value from the store and does not generate any new paths.

Rule A-ASSIGN describes how a program binds a variable to an abstract value. To ensure soundness, the store maintains a may-be-bound-to relationship between variables and abstract values. Thus the binding operation is a join, rather than an overwrite.

Rule A-FOR evaluates the expression e to determine the paths π_e representing the possible collections to be iterated. To each of these paths, the analysis appends the iterator field name l^l that

stands for a particular element of the collection. Thus, if the path to the collection is $f_1.f_2$ and the collection's elements each has a field f_3 , then the path to one of the element's f_3 field is $f_1.f_2.l^l.f_3$. Each collection iteration (**for** loop) that appears in a program has a unique iterator field name l^l , where l is the corresponding loop variable. Iterator field names are useful for transforming the analysis results into a query, as discussed in Section 5.

The analysis approximates the loop body's concrete behavior by taking the transitive closure of abstractly executing the loop an arbitrary number of times. This value is the least upper bound of a function **do** that is specialized for a given command c , loop variable l , and set of iterator paths π_l . The function takes an initial store $\hat{\sigma}$ and set of paths π and evaluates c under an updated store that maps l to π_l to yield a new store $\hat{\sigma}'$ and a new set of paths π' . The result of the function is $\hat{\sigma}'$ and the combined path set $\pi \sqcup \pi'$.

The abstract evaluation of our running example generates a final store with the following mappings

$$\{\text{root} \mapsto \{\epsilon\}, \text{output} \mapsto \top\}$$

and generates the following set of paths:

$$\{\epsilon, \text{employees}, \text{employees}.l^e, \text{employees}.l^e.\text{salary}, \text{employees}.l^e.\text{name}, \text{employees}.l^e.\text{manager}, \text{employees}.l^e.\text{manager.name}\}$$

3.3 Soundness

The analysis is sound if it safely approximates the values a program loads. If the database stores a set of values V , and if executing a program causes the set of persistent values $\rho \subseteq V$ to be loaded, then the analysis should describe a set of values $\hat{\rho}$ such that $\rho \subseteq \hat{\rho} \subseteq V$. We formalize this relation between concrete and abstract values and prove that the operational semantics preserves the relation.

The set of concrete values described by an abstract path is the set of values reachable by following that path from the root. We can formalize this description by lifting the definition of *Load* to operate on paths:

$$\text{Load}(r, \epsilon) = \{r\} \quad (\text{P-LOAD}_1)$$

$$\frac{\text{Load}(r, f) = \{r_1, \dots, r_n\}}{\text{Load}(r, f.p) = \bigcup \text{Load}(r_i, p)} \quad (\text{P-LOAD}_2)$$

$$\text{Load}(r, l^l.p) = \text{Load}(r, p) \quad (\text{P-LOAD}_3)$$

Rule P-LOAD₁ states that traversing an empty path from record identifier r yields the set containing r . Rule P-LOAD₂ loads one level of the traversal hierarchy, then recursively loads the remainder of the hierarchy. Rule P-LOAD₃ removes an iterator field name from a path, essentially binding the name to record identifier r .

A set of paths π safely approximates a set of values ρ if the set of values reachable by following all paths in π is a superset of ρ :

$$\rho \mathcal{R} \pi \Leftrightarrow \rho \subseteq \bigcup_{p \in \pi} \text{Load}(r_0, p)$$

For our running example, $\bigcup_{p \in \pi} \text{Load}(r_0, p) =$

$$\{r_0, r_1, 50000, \text{"Adam"}, r_2, 60000, \text{"Bob"}, r_3, 70000, \text{"Cathy"}\}$$

which safely over-approximates the concrete results.

The abstract domain consists of sets of paths and \top , so we lift \mathcal{R} to relate a concrete value v to an abstract value \hat{v} as follows:

$$v \mathcal{R} \hat{v} \Leftrightarrow \begin{cases} \{v\} \mathcal{R} \hat{v} & v = r, \hat{v} = \pi \\ v \mathcal{R} \hat{v} & v = \{r_1, \dots, r_n\}, \hat{v} = \pi \\ \hat{v} = \top & \text{otherwise} \end{cases}$$

The first two cases relate record identifiers and paths, as above. The final case states that \top always safely approximates a basic program value. The relation is lifted to be defined on stores as follows:

$$\sigma \mathcal{R} \hat{\sigma} \Leftrightarrow \forall x \in \text{Dom}(\sigma) \cap \text{Dom}(\hat{\sigma}). \sigma[x] \mathcal{R} \hat{\sigma}[x]$$

By these definitions, the initial stores are *compatible*, in that the initial abstract store safely approximates the initial concrete store.

If we define an ordering on abstract stores as follows:

$$\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2 \Leftrightarrow \text{Dom}(\hat{\sigma}_1) \subseteq \text{Dom}(\hat{\sigma}_2) \wedge \forall x \in \text{Dom}(\hat{\sigma}_1) \cap \text{Dom}(\hat{\sigma}_2). \hat{\sigma}_1[x] \sqsubseteq \hat{\sigma}_2[x]$$

then the following lemma and its corollary state that if a concrete value (or store) relates to one abstract value—and if that abstract value is less than a second abstract value—then the concrete value also relates to the second abstract value.

Lemma 1. *For all $v, \hat{v}_1, \hat{v}_2: v \mathcal{R} \hat{v}_1 \wedge \hat{v}_1 \sqsubseteq \hat{v}_2 \Rightarrow v \mathcal{R} \hat{v}_2$.*

Corollary 1. *For all $\sigma, \hat{\sigma}_1, \hat{\sigma}_2: v \mathcal{R} \hat{\sigma}_1 \wedge \hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2 \Rightarrow \sigma \mathcal{R} \hat{\sigma}_2$.*

The lemma and corollary are trivially proved by examining the possible values for v (σ) and \hat{v}_i ($\hat{\sigma}$) and applying the appropriate definition of \mathcal{R} .

The proof that the abstract semantics preserves \mathcal{R} requires that we first prove the rules to be monotone. The following two properties of lattices will be useful for proving monotonicity:

$$\hat{v}_1 \sqsubseteq \hat{v}_2 \wedge \hat{v}_3 \sqsubseteq \hat{v}_4 \Rightarrow \hat{v}_1 \sqcup \hat{v}_3 \sqsubseteq \hat{v}_2 \sqcup \hat{v}_4 \quad (1)$$

$$(\forall \hat{v}_1 \in V_1, \forall \hat{v}_2 \in V_2. \hat{v}_1 \sqsubseteq \hat{v}_2) \Rightarrow \bigsqcup V_1 \sqsubseteq \bigsqcup V_2 \quad (2)$$

The first property states that the upper bound of a pair of “lower” elements (\hat{v}_1, \hat{v}_3) is less than the upper bound of a pair of “higher” elements (\hat{v}_2, \hat{v}_4). The second property describes the same relationship for sets of lattice elements, rather than pairs.

Theorem 1 (Monotonicity of Expression Evaluation). *For all $\hat{\sigma}_1, \hat{\sigma}_2, e$:*

$$\frac{\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2 \quad \langle e, \hat{\sigma}_1 \rangle \hat{\rightarrow} \langle \hat{v}_1, \pi_1 \rangle \quad \langle e, \hat{\sigma}_2 \rangle \hat{\rightarrow} \langle \hat{v}_2, \pi_2 \rangle}{(\hat{v}_1, \pi_1) \sqsubseteq (\hat{v}_2, \pi_2)}$$

where $(\hat{v}_1, \pi_1) \sqsubseteq (\hat{v}_2, \pi_2)$ means $(\hat{v}_1 \sqsubseteq \hat{v}_2) \wedge (\pi_1 \sqsubseteq \pi_2)$.

Proof. By induction on the structure of e

Case $e \equiv \llbracket l \rrbracket$ Rule A-VAR gives $(\hat{v}_1, \pi_1) = (\hat{\sigma}_1[l], \emptyset)$ and $(\hat{v}_2, \pi_2) = (\hat{\sigma}_2[l], \emptyset)$. The premise $\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2$ yields the desired result.

Case $e \equiv \llbracket \text{op}_n(e_1, \dots, e_n) \rrbracket$ Rule A-OP gives $(\hat{v}_1, \pi_1) = (\top, \bar{\pi}_1)$ and $(\hat{v}_2, \pi_2) = (\top, \bar{\pi}_2)$, where $\bar{\pi}_i$ is the results of analyzing the operand expressions. The induction hypothesis and property (2) yield $\bigsqcup \bar{\pi}_1 \sqsubseteq \bigsqcup \bar{\pi}_2$.

Case $e \equiv \llbracket e.f^t \rrbracket$ Rule A-TRAVERSE first evaluates e which gives results $\hat{v}_{e_1} = (\pi_{e_1}, \pi'_1)$ and $\hat{v}_{e_2} = (\pi_{e_2}, \pi'_2)$. There are four cases to consider for the final results of the analysis \hat{v}_1, \hat{v}_2 :

Case 1 $f^t \in \pi_{e_1}, f^t \in \pi_{e_2}$: Then $\hat{v}_1 = \hat{v}_2 = \top$.

Case 2 $f^t \in \pi_{e_1}, f^t \notin \pi_{e_2}$: Not possible, because the induction hypothesis asserts $\pi_{e_1} \sqsubseteq \pi_{e_2}$, which implies $f^t \in \pi_{e_2}$.

Case 3 $f^t \notin \pi_{e_1}, f^t \in \pi_{e_2}$: In this case, $\hat{v}_1 = \{p.f^t \mid p \in \pi_{e_1}\} \sqsubseteq \hat{v}_2 = \top$.

Case 4 $f^t \notin \pi_{e_1}, f^t \notin \pi_{e_2}$: In this case, the induction hypothesis assures $\hat{v}_1 = \{p.f^t \mid p \in \pi_{e_1}\} \sqsubseteq \hat{v}_2 = \{p.f^t \mid p \in \pi_{e_2}\}$.

Finally, (2) gives $\pi_1 = \hat{v}_1 \sqcup \pi'_1 \sqsubseteq \hat{v}_2 \sqcup \pi'_2$. \square

Theorem 2 (Monotonicity of Command Evaluation). *For all $\hat{\sigma}_1, \hat{\sigma}_2, c$*

$$\frac{\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2 \quad \langle c, \hat{\sigma}_1 \rangle \dot{\rightarrow} \langle \hat{\sigma}'_1, \pi_1 \rangle \quad \langle c, \hat{\sigma}_2 \rangle \dot{\rightarrow} \langle \hat{\sigma}'_2, \pi_2 \rangle}{(\hat{\sigma}'_1, \pi_1) \sqsubseteq (\hat{\sigma}'_2, \pi_2)}$$

Proof. By induction on the structure of c . We omit the proofs for **skip**, **if** commands, and sequences of commands, because these proofs merely appeal to the induction hypothesis and properties (1) and (2).

Case $c \equiv [l := e]$ Rule A-ASSIGN first evaluates e and gives results (\hat{v}_1, π_1) and (\hat{v}_2, π_2) . The rule then gives $\hat{\sigma}'_1[l] = \hat{v}_1 \sqcup \hat{\sigma}_1[l]$ and $\hat{\sigma}'_2[l] = \hat{v}_2 \sqcup \hat{\sigma}_2[l]$. Theorem 1, the induction hypothesis, and (1) give the full result $(\hat{\sigma}'_1, \pi_1) \sqsubseteq (\hat{\sigma}'_2, \pi_2)$.

Case $c \equiv [\text{for } l \text{ in } e \text{ do } c]$ Rule A-FOR first evaluates the collection expression, giving results (π_{e_1}, π_1) and (π_{e_2}, π_2) . The rule then creates iterator paths sets π_{l_1} and π_{l_2} . The premise states that $\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2$ and the induction hypothesis asserts $\pi_{l_1} \sqsubseteq \pi_{l_2}$. Next it must be shown that **do** is monotonic:

$$\frac{\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2 \quad \pi_{l_1} \sqsubseteq \pi_{l_2} \quad \pi_1 \sqsubseteq \pi_2 \quad \begin{array}{l} \mathbf{do}\langle c, l, \pi_{l_1} \rangle(\hat{\sigma}_1, \pi_1) = (\hat{\sigma}'_1, \pi'_1) \\ \mathbf{do}\langle c, l, \pi_{l_2} \rangle(\hat{\sigma}_2, \pi_2) = (\hat{\sigma}'_2, \pi'_2) \end{array}}{(\hat{\sigma}'_1, \pi'_1) \sqsubseteq (\hat{\sigma}'_2, \pi'_2)}$$

This result can be achieved because **do** makes the transition $\langle c, [l \mapsto \pi_{l_i}] \sqcup \hat{\sigma}_i \rangle \dot{\rightarrow} \langle \hat{\sigma}'_i, \pi_{c_i} \rangle$. The first two premises of the above rule and (1) give $[l \mapsto \pi_{l_1}] \sqcup \hat{\sigma}_1 \sqsubseteq [l \mapsto \pi_{l_2}] \sqcup \hat{\sigma}_2$. The main induction hypothesis then yields $\hat{\sigma}'_1 \sqsubseteq \hat{\sigma}'_2$ and $\pi_{c_1} \sqsubseteq \pi_{c_2}$. The premise $\pi_1 \sqsubseteq \pi_2$ and (1) give $\pi'_1 = \pi_1 \sqcup \pi_{c_1} \sqsubseteq \pi'_2 = \pi_2 \sqcup \pi_{c_2}$. Thus **do** is monotonic.

Returning to the **for** case: Because **do** is monotonic and because $\hat{\sigma}_1 \sqsubseteq \hat{\sigma}_2$ and $\pi_{l_1} \sqsubseteq \pi_{l_2}$, (2) yields the conclusion:

$$\begin{aligned} (\hat{\sigma}'_1, \pi'_1) &= \bigsqcup \{ (\mathbf{do}\langle c, l, \pi_{l_1} \rangle)^n(\hat{\sigma}_1, \emptyset) \mid n \in \mathbb{N} \} \\ &\sqsubseteq \\ (\hat{\sigma}'_2, \pi'_2) &= \bigsqcup \{ (\mathbf{do}\langle c, l, \pi_{l_2} \rangle)^n(\hat{\sigma}_2, \emptyset) \mid n \in \mathbb{N} \} \end{aligned}$$

The full result $(\hat{\sigma}'_1/l, \pi_1 \sqcup \pi_{l_1} \sqcup \pi'_1) \sqsubseteq (\hat{\sigma}'_2/l, \pi_2 \sqcup \pi_{l_2} \sqcup \pi'_2)$ is achieved by applying (2). \square

We now proceed to show that the computation semantics preserve \mathcal{R} . Computation soundness requires the following lemma, which states that compatibility is maintained when combining the results of compatible subcomputations.

Lemma 2 (Subcomputation compatibility). *If $(\rho_1 \mathcal{R} \pi_1)$ and $(\rho_2 \mathcal{R} \pi_2)$, then $(\rho_1 \cup \rho_2) \mathcal{R} (\pi_1 \sqcup \pi_2)$.*

Proof. If $\pi_1 \sqcup \pi_2 = \top$, then the relation trivially holds. Otherwise, by the definition of \mathcal{R} , $\rho_1 \subseteq \bigcup_{p \in \pi_1} \text{Load}(p)$ and $\rho_2 \subseteq \bigcup_{p \in \pi_2} \text{Load}(p)$. Since $\rho_1 \cup \rho_2 \subseteq \bigcup_{p \in \pi_1 \cup \pi_2} \text{Load}(p)$, the relation holds. \square

Theorem 3 (Soundness of expression evaluation). *For all $\sigma, \hat{\sigma}, e$,*

$$\frac{\langle e, \sigma \rangle \rightarrow \langle v, \rho \rangle \quad \langle e, \hat{\sigma} \rangle \rightarrow \langle \hat{v}, \pi \rangle \quad \sigma \mathcal{R} \hat{\sigma}}{(v, \rho) \mathcal{R} (\hat{v}, \pi)}$$

Proof. By induction on the structure of e .

Base case $e \equiv [l]$ The concrete and abstract semantics rules for this case give $(v, \rho) = (\sigma[l], \emptyset)$ and $(\hat{v}, \pi) = (\hat{\sigma}[l], \emptyset)$. The premise $\sigma \mathcal{R} \hat{\sigma}$ yields the desired result.

The induction hypothesis asserts that evaluating subexpressions produces sound results. It remains to show that evaluating operators and traversals produces sound results.

Case $e \equiv [\text{op}_n(e_1, \dots, e_n)]$ In this case, $\hat{v} = \top$, so $v \mathcal{R} \hat{v}$ by definition. The rules state $\rho = \bigcup \rho_i$ and $\pi = \bigsqcup \pi_i$. The induction hypothesis and Lemma 2 achieve $\rho \mathcal{R} \pi$.

Case $e \equiv [e.f^t]$ If rule A-TRAVERSE analyzes the record expression e and gives \top for π_f , then the case is proved. Otherwise, $\langle e, \sigma \rangle \rightarrow \langle r, \rho \rangle$, $\langle e, \hat{\sigma} \rangle \dot{\rightarrow} \langle \pi_e, \pi \rangle$, rule U-TRAVERSE gives $\rho_f = \text{Load}(r, f)$ and rule A-TRAVERSE gives $\pi_f = \{p.f^t \mid p \in \pi_e\}$. The induction hypothesis gives $r \mathcal{R} \pi_e$. A simple inductive argument shows that P-LOAD guarantees the following property, which states that if record identifier r can be reached by following path p from the root, then the record identifiers that correspond to r 's field f can be reached from the root by following p extended with f .

$$\frac{r \in \text{Load}(r_0, p)}{\text{Load}(r, f) \subseteq \text{Load}(r_0, p.f)} \quad (3)$$

Thus, if $r \in \rho$ then $\rho_f \mathcal{R} \pi_f$. This result, the induction hypothesis, and Lemma 2 give $(\rho_f \cup \rho) \mathcal{R} (\pi_f \sqcup \pi)$. \square

Theorem 4 (Soundness of command evaluation). *For all $\sigma, \hat{\sigma}, c$,*

$$\frac{\langle c, \sigma \rangle \rightarrow \langle \sigma', \rho \rangle \quad \langle c, \hat{\sigma} \rangle \rightarrow \langle \hat{\sigma}', \pi \rangle \quad \sigma \mathcal{R} \hat{\sigma}}{(\sigma', \rho) \mathcal{R} (\hat{\sigma}', \pi)}$$

Proof. By induction on the structure of c .

Base case $c \equiv [\text{skip}]$ The premises suffice, because the command neither alters the store nor generates any database loads.

The induction hypothesis asserts that analyzing subcommands produces sound results. We omit the proof case for sequences, because that case merely invokes the induction hypothesis.

Case $c \equiv [l := e]$ The semantic rules give $\langle e, \sigma \rangle \rightarrow \langle v, \rho \rangle$ and $\langle e, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{v}, \pi \rangle$. The conclusion $\rho \mathcal{R} \pi$ is given by Theorem 3, which also gives $v \mathcal{R} \hat{v}$. The premise $\sigma \mathcal{R} \hat{\sigma}$ guarantees that σ' and $\hat{\sigma}'$ relate for all variables that are not l . It remains to show that $\sigma'[l] \mathcal{R} \hat{\sigma}'[l]$. Proceed by cases:

Case 1: $\sigma[l]$, $\hat{\sigma}[l]$ undefined. In this case, $\sigma'[l] = [l \mapsto v]\sigma$ and $\hat{\sigma}'[l] = [l \mapsto \hat{v}]\hat{\sigma}$. The result $v \mathcal{R} \hat{v}$ suffices.

Case 2: $\sigma[l] = v_0$, $\hat{\sigma}[l] = \hat{v}_0$. In this case, $\sigma' = [l \mapsto v]\sigma$ and $\hat{\sigma}' = [l \mapsto \hat{v} \sqcup \hat{v}_0]\hat{\sigma}$. The premise and Lemma 1 guarantee $v \mathcal{R} (\hat{v} \sqcup \hat{v}_0)$.

Case 3: $\sigma[l]$ undefined, $\hat{\sigma}[l] = \hat{v}_0$. Same as Case 2.

Case 4: $\sigma[l] = v_0$, $\hat{\sigma}[l]$ undefined. This case cannot occur, because the semantic rules guarantee $\text{Dom}(\sigma) \subseteq \text{Dom}(\hat{\sigma})$.

Case $c \equiv [\text{if } e \text{ then } c_1 \text{ else } c_2]$ If $\langle e, \sigma \rangle \rightarrow \langle \text{true}, \rho \rangle$, then the concrete semantics gives $\langle c_1, \sigma \rangle \rightarrow \langle \sigma', \rho \rangle$ and the abstract semantics gives $\langle c_1, \hat{\sigma} \rangle \rightarrow \langle \hat{\sigma}_1, \pi_1 \rangle$ and $\langle c_2, \hat{\sigma} \rangle \rightarrow \langle \hat{\sigma}_2, \pi_2 \rangle$. By the structural induction hypothesis, $\rho \mathcal{R} \pi_1$ and $\sigma' \mathcal{R} \hat{\sigma}_1$. Lemma 2 yields $\rho \mathcal{R} \pi_1 \sqcup \pi_2$ and $\sigma' \mathcal{R} \hat{\sigma}_1 \sqcup \hat{\sigma}_2$. The argument for $\langle e, \sigma \rangle \rightarrow \langle \text{false}, \rho \rangle$ proceeds similarly.

Case $c \equiv [\text{for } l \text{ in } e \text{ do } c]$ Theorem 2 states that **do** is monotone. The Knaster-Tarski theorem thus asserts that a fixpoint exists, i.e. there is a number m such that $\mathbf{do}_{c,l,\pi_e}^m(\hat{\sigma}, \emptyset) = \bigsqcup \{ \mathbf{do}\langle c, l, \pi_i \rangle^i(\hat{\sigma}, \emptyset) \mid i \in \mathbb{N} \}$. Because this value is the greatest fixpoint, we can conclude $\mathbf{do}\langle c, l, \pi_i \rangle^n(\hat{\sigma}, \emptyset) \sqsubseteq \mathbf{do}\langle c, l, \pi_i \rangle^m(\hat{\sigma}, \emptyset)$, where n is the number of iterations that take place when the program actually executes. \square

4. Analyzing Traversal Conditions

The precision of the analysis can be significantly increased by analyzing the conditions under which a program traverses its paths. For example, the analysis in the previous section conservatively estimates that the program in Fig. 1 needs the name field for every employee even though the program traverses the name field only if the employee's salary is greater than \$65,000. We extend our analysis to identify and include such conditions, so that they may be expressed in a database query.

4.1 Query Conditions

A *query condition* is a conditional program expression that can be expressed as a part of a query and evaluated by a database. Databases typically only allow conditions that operate on individual records of a set; for example, the **select** operator in relational algebra evaluates a condition separately for each tuple in a relation. An expression in the kernel language is a query condition if it satisfies the following requirements, related to the requirements for parallelizing code in a parallelizing compiler [1, 30]. The requirements are illustrated in Fig. 8, where the notation $C[l]$ means that condition C can depend upon variable l .

1. The database must be able to evaluate all operations (e.g., $>$) that appear in the expression. The database evaluation should produce the same result as program evaluation.
2. The expression can contain no loop-carried dependences [1]. A loop-carried dependence occurs when the evaluation of an expression in a loop depends upon variables assigned in previous iterations of a loop. The condition in (a) is a query condition because x is redefined in each iteration of the loop. Example (b) is not because x has a loop-carried dependence.
3. The paths the expression traverses can refer to no more than one element from a given collection. If an expression meets this requirement, its paths are *distinct*. Both conditions of example (c) are query conditions, because each expression depends on only one element of a given collection. Condition C_2 of example (d) is not a query condition, because it depends on a variable bound in a different iteration.
4. The expression may traverse paths that refer to elements of more than one collection, but only if the expression refers to elements of nested collections that correspond with the nested bindings of loop variables that a query can express. This restriction is satisfied when each nested loop iterates over a path that extends the path of its outer loop(s). Example (e) is a query condition because the expression depends only on paths that satisfy this condition. Example (f) is not a query condition because the expression depends on two unrelated collections.

These restrictions do not prevent common programming idioms used in data-intensive applications. A more powerful query translation could support more complex conditions.

4.2 Data Dependences

A program's data dependences [1] provide information about which persistent values the program must retrieve. If a persistent value affects the contents of the final store, the program must retrieve that value. Assignment statements introduce data dependences, because any assigned value may affect the contents of the final store. Loop variables, however, do not directly induce a data dependence on the final store, because these variables are removed from the store after the loop terminates. We extend our analysis to collect information about which paths induce data dependences. The query creation algorithm in Section 5 uses this information to ensure retrieval of all values represented by data-dependent paths.

<pre>for l₁ in p x := E[l₁] if C[l₁,x] then S</pre> <p>(a) C is a query condition</p>	<pre>for l₁ in p x := E[l₁] + x if C[l₁,x] then S</pre> <p>(b) C is not a query condition</p>
<pre>for l₁ in p if C₁[l₁] then S₁ for l₂ in p if C₂[l₂] then S₂</pre> <p>(c) $C_{1,2}$ are query conditions</p>	<pre>for l₁ in p if C₁[l₁] then x := e[l₁] for l₂ in p if C₂[l₂,x] then S</pre> <p>(d) C_2 is not a query condition</p>
<pre>for l₁ in p for l₂ in l₁.f if C[l₁,l₂] then S</pre> <p>(e) C is a query condition</p>	<pre>for l₁ in p₁ for l₂ in p₂ if C[l₁,l₂] then S</pre> <p>(f) C is not a query condition</p>

Figure 8. Examples of conditions and iterations.

4.3 Domains for Paths with Conditions

A path $p[k]$ represents a query of the database for values located at path p for which the condition k is true. The condition is expressed as an operation on abstract values, including other paths. The domain of abstract values is extended to include conditions:

$$\begin{aligned}
 k \in \text{Condition} & ::= \text{op}_n^t(\hat{v}_1, \dots, \hat{v}_n) \\
 cp \in \text{CPath} & ::= p[k] \mid p[k]^* \\
 \hat{v} \in \widehat{\text{Value}} & ::= \wp(\text{CPath}) + \wp(\text{Condition}) + \top
 \end{aligned}$$

A path marked $p[k]^*$ is involved in a data dependence. A non-conditional path p is lifted to a conditional path $p[\text{true}]$ signifying that the path is always traversed. The label t in a condition is a syntactic label constructed in the same fashion described for traversals in Section 3.1. The syntactic labels of a given program restrict the domain of CPath to be finite.

4.4 Abstract Semantics for Conditional Traversals

Figure 9 extends the abstract semantics to include the condition k under which a program traverses a path. The evaluation relation $k, I \vdash \langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}, \pi \rangle$ means that e evaluates to \hat{v} under condition k in abstract store $\hat{\sigma}$, with a list of enclosing iteration variable paths I . Initially k is set to *true*, and I is empty. We omit rules K-VAR, K-SKIP, and K-SEQ because these rules only collect the results of subcomputation and do not alter the context.

Rule K-IF₁ identifies query conditions. The analysis first determines that e contains no loop-carried dependences. This determination is the result of a standard analysis; for brevity, we omit its details. The analysis also checks that all paths in the condition's abstract value are *distinct*, in the sense that they do not traverse the same set of fields with different iteration variables. Finally, it checks that all the paths used in the expression are based on lexically enclosing iteration paths.

The lexically enclosing iterations needed by K-IF₁ are created by K-FOR. The rule appends a new set of paths π_ι to the list of iteration paths I only if all paths in π_ι extend some path in I_n , the list's most recently added member. In this way, K-FOR maintains the constraint that I is a list of lexically enclosing iteration paths. Other than imposing this constraint on I , K-FOR is the same as A-FOR (Fig. 7).

Query conditions are used in the true and false branches of the **if** command. Given the abstract value \hat{v} of e , the true-branch body

$$\begin{array}{c}
\text{e contains no loop-carried dependences} \\
\frac{k, I \vdash \langle e, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{v}, \pi_e \rangle \quad \text{Distinct}(\text{Paths}(\hat{v})) \quad \text{Trim}(\text{Paths}(\hat{v})) \subseteq I}{(k \wedge \hat{v}), I \vdash \langle c_1, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}_1, \pi_1 \rangle} \\
\frac{(k \wedge \neg \hat{v}), I \vdash \langle c_2, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}_2, \pi_2 \rangle}{\pi' = \pi_e \sqcup \pi_1 \sqcup \pi_2} \\
\frac{}{k, I \vdash \langle \text{if } e \text{ then } c_1 \text{ else } c_2, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}_1 \sqcup \hat{\sigma}_2, \pi' \rangle} \quad (\text{K-IF}_1)
\end{array}$$

$$\frac{\text{K-IF}_1 \text{ does not apply}}{k, I \vdash \langle c_1, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}_1, \pi_1 \rangle} \\
\frac{k, I \vdash \langle c_2, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}_2, \pi_2 \rangle}{\pi' = \pi_e \sqcup \pi_1 \sqcup \pi_2} \\
\frac{}{k, I \vdash \langle \text{if } e \text{ then } c_1 \text{ else } c_2, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}_1 \sqcup \hat{\sigma}_2, \pi' \rangle} \quad (\text{K-IF}_2)$$

$$\frac{k, I \vdash \langle e, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{v}, \pi \rangle}{k, I \vdash \langle l := e, \hat{\sigma} \rangle \dot{\rightarrow} \langle [l \mapsto \hat{v}] \sqcup \hat{\sigma}, \pi^* \sqcup I_n[k]^* \rangle} \quad (\text{K-ASSIGN})$$

$$\frac{\forall p_1, p_2 \in \pi : \text{Erase}(p_1) = \text{Erase}(p_2) \Rightarrow p_1 = p_2}{\text{Distinct}(\pi)}$$

$$\begin{array}{l}
\text{Erase}(\bar{f}_1.l^{l_1} \dots .l^{l_n}.\bar{f}_{n+1}) = \bar{f}_1 \dots \bar{f}_{n+1} \\
\text{Trim}(\pi) = \{p.l^l \mid p.l^l.f \in \pi\} \\
\text{Paths}(\pi) = \pi \\
\text{Paths}(\text{op}_n(e_1, \dots, e_n)) = \text{Paths}(e_1) \sqcup \dots \sqcup \text{Paths}(e_n)
\end{array}$$

$$\frac{k, I \vdash \langle e, \hat{\sigma} \rangle \dot{\rightarrow} \langle \pi_e, \pi \rangle}{\pi_f = \begin{cases} \top & f^t \in \pi_e \\ \{p.f^t[k] \mid p \in \pi_e\} & f^t \notin \pi_e \end{cases}} \\
\frac{}{k, I \vdash \langle e.f^t, \hat{\sigma} \rangle \dot{\rightarrow} \langle \pi_f, \pi \sqcup \pi_f \rangle} \quad (\text{K-TRAVERSE})$$

$$\frac{k, I \vdash \langle e_i, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{v}_i, \pi_i \rangle \quad \text{for } i \in \{1, \dots, n\}}{\hat{v} = \begin{cases} \top & \text{op}_n^t \in \hat{v}_i \\ \text{op}_n^t(\hat{v}_1, \dots, \hat{v}_n) & \text{op}_n^t \notin \hat{v}_i \end{cases}} \\
\frac{}{k, I \vdash \langle \text{op}_n^t(e_1, \dots, e_n), \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{v}, \sqcup \pi_i \rangle} \quad (\text{K-OP})$$

$$\frac{k, I \vdash \langle e, \hat{\sigma} \rangle \dot{\rightarrow} \langle \pi_e, \pi \rangle}{\pi_l = \{p.l^l \mid p \in \pi_e\}} \\
\frac{I' = \text{Extend}(I, \pi_l)}{(\hat{\sigma}', \pi') = \sqcup \{(\text{do}(k, I', c, l, \pi_l))^n(\hat{\sigma}, \emptyset) \mid n \in \mathbb{N}\}} \\
\frac{}{k, I \vdash \langle \text{for } l \text{ in } e \text{ do } c, \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}'/l, \pi \sqcup \pi_l \sqcup \pi' \rangle} \quad (\text{K-FOR})$$

$$\frac{k, I \vdash \langle c, [l \mapsto \pi_l] \sqcup \hat{\sigma} \rangle \dot{\rightarrow} \langle \hat{\sigma}', \pi' \rangle}{\text{do}(k, I, c, l, \pi_l)(\hat{\sigma}, \pi) = (\hat{\sigma}', \pi \sqcup \pi')} \quad (\text{K-DO})$$

$$\text{Extend}(I, \pi_l) = \begin{cases} \pi_l & n = 0 \\ I, \pi_l & \text{Prefixes}(I_n, \pi_l) \\ I & \text{otherwise} \end{cases}$$

$$\frac{\forall p_2 \in \pi_2 : (\exists p_1 \in \pi_1, \exists p' \in \text{Field}^* : p_1.p' = p_2)}{\text{Prefixes}(\pi_1, \pi_2)}$$

Figure 9. Abstract interpretation with paths and conditions.

is evaluated under the condition $k \wedge \hat{v}$ and the false-branch body under the condition $k \wedge \neg \hat{v}$. When the program makes a traversal, rule K-TRAVERSE attaches the condition k to the path generated by the traversal.

If the condition does not satisfy the requirements of a query condition, rule K-IF₂ does not augment the paths with conditions.

Rule K-ASSIGN performs assignments but also marks all paths in the bound expression as having data dependences. π^* means the marking of all paths in π with $*$, and $I_n[k]^*$ means marking all paths in I_n with condition k and $*$. The iteration variable paths themselves are marked as a data dependence because the execution of any assignment can depend upon the *existence* of an element in an iteration, even if no fields of the iteration variable are used. For example, in the program **for** x **in** p **do** $y := y + 1$, the variable x is never used, yet there is still a data dependence upon it because its elements must be enumerated.

Rule K-OP defines the semantics of operations on abstract values. The operands are evaluated and the operator is retained in the result. The rule also includes a widening clause to ensure convergence to a fixed-point. The rule is similar to the one for traversals: If the syntactic use of the operator op_n already occurs in one of the \hat{v}_i , then the expression evaluates to \top .

The abstract evaluation of our running example generates a final store with the following mappings:

root $\mapsto \{\epsilon\}$, **output** \mapsto
 $\{\text{print}(\text{employees}.l^e.\text{name}\{\{\text{employees}.l^e.\text{salary}\} > 65000\} +$
 $\text{employees}.l^e.\text{manager.name}\{\{\text{employees}.l^e.\text{salary}\} > 65000\})\}$

and the following set of paths:

$\{\epsilon, \text{employees}, \text{employees}.l^e, \text{employees}.l^e.\text{salary},$
 $\text{employees}.l^e.\{\{\text{employees}.l^e.\text{salary}\} > 65000\}^*,$
 $\text{employees}.l^e.\text{name}\{\{\text{employees}.l^e.\text{salary}\} > 65000\}^*,$
 $\text{employees}.l^e.\text{manager}\{\{\text{employees}.l^e.\text{salary}\} > 65000\}^*,$
 $\text{employees}.l.\text{manager.name}\{\{\text{employees}.l^e.\text{salary}\} > 65000\}^*\}$

At this stage of the analysis, the conditions apply to the final attributes loaded by a path. In Section 5 we further analyze the conditions and paths to avoid loading entire records.

4.5 Soundness

Proceeding as before, we define the load operation for conditional paths and define the relations between concrete and abstract domains. We then prove that evaluation preserves these relations. This proof relies on the previous soundness proof (Section 3.3).

The load operation for conditional paths is defined in Fig. 10. Function CLoad_i loads all records reachable by a path p , provided the path's condition k may be *true*. A mapping ϕ binds an iterator field name to a specific record identifier, to be referenced in the evaluation of the path's condition. CLoad_1 creates iterator field name bindings, and CLoad_2 uses the bindings.

Function *eval* defines condition evaluation. Operator evaluation calls *eval* on the operands and applies f'_{op_n} to the results, where

$$f'_{\text{op}_n}(\hat{v}_1, \dots, \hat{v}_n) = \begin{cases} \top & \top \in \{\hat{v}_1, \dots, \hat{v}_n\} \\ f_{\text{op}_n}(\hat{v}_1, \dots, \hat{v}_n) & \text{otherwise} \end{cases}$$

Note that because the underlying operators are monotonic, all functions f' are also monotonic.

$$\begin{aligned}
CLoad_i(r, \epsilon[k], \phi) &= \begin{cases} \{r\} & true \sqsubseteq eval(k, \phi) \\ \emptyset & otherwise \end{cases} \\
CLoad_i(r, f.p[k], \phi) &= \bigcup_{r' \in Load(r, f)} CLoad_i(r', p[k], \phi) \\
CLoad_1(r, l.p[k], \phi) &= CLoad_1(r, p[k], [l^l \mapsto r]\phi) \\
CLoad_2(r, l.p[k], \phi) &= CLoad_2(\phi(l^l), p[k], \phi) \\
\\
eval(p[k], \phi) &= \bigsqcup_{\hat{v} \in S} CLoad_2(r_0, p[k], \phi) \\
eval(\text{op}_n(\hat{v}_1 \dots \hat{v}_n), \phi) &= f'_{\text{op}_n}(eval(\hat{v}_1, \phi), \dots, eval(\hat{v}_n, \phi)) \\
eval(S, \phi) &= \bigsqcup_{\hat{v} \in S} eval(\hat{v}, \phi)
\end{aligned}$$

Figure 10. Conditional record loading.

Path evaluation calls $CLoad_2$ on the path, passing bindings ϕ for any iterator field names that appear in the path. Note that, because the evaluated path can appear only in an operation expression, the result of path evaluation must be a set that contains a single basic value. The least upper bound operation (\sqcup) retrieves this value from the set. Evaluating a set of abstract values yields the least upper bound of evaluating each value in the set.

A set π of conditional paths safely approximates the persistent values a program loads if it describes a superset of those values. We modify the definitions of \mathcal{R} to relate concrete values and conditional paths:

$$\begin{aligned}
\rho \mathcal{R} \pi &\Leftrightarrow \rho \subseteq \bigcup_{p[k] \in \pi} CLoad_1(r_0, p[k], \emptyset) \\
(v, \sigma) \mathcal{R} \hat{v} &\Leftrightarrow \begin{cases} \{v\} \mathcal{R} \hat{v} & v = r, \hat{v} = \pi \\ v \mathcal{R} \hat{v} & v = \{r_1, \dots, r_n\}, \hat{v} = \pi \\ v \sqsubseteq eval(\hat{v}, \phi_\sigma) & v \in Basic, \hat{v} = k \\ \hat{v} = \top & otherwise \end{cases} \\
\sigma \mathcal{R} \hat{\sigma} &\Leftrightarrow \forall x \in Dom(\sigma) \cap Dom(\hat{\sigma}). (\sigma[x], \sigma) \mathcal{R} \hat{\sigma}[x]
\end{aligned}$$

The relation between concrete and abstract values is defined only in the context of a store, because the store provides a binding for any iterator field names that may appear in paths and conditions. When \hat{v} is an abstract operation, evaluating \hat{v} must approximate v . If L is the set of loop variables that appear in the entire program, $\phi_\sigma = \bigcup_{l \in L} [l^l \mapsto \sigma[l]]$, where $\sigma[l] = \top$ if $\sigma[l]$ is undefined. As before, \mathcal{R} is lifted to be defined on stores.

To prove soundness, we first show that expression evaluation gives the same results as the analysis in Section 3.3, *assuming that evaluating every path condition may give the value true*. We then show that the analysis only constructs conditions that satisfy this assumption. Proof of soundness for command evaluation follows trivially.

As before, we state a lemma for soundness of subcomputation. The lemma's proof proceeds similarly to that of Lemma 2.

Lemma 3 (Subcomputation compatibility). *If $(\rho_1 \mathcal{R} \pi_1)$ and $(\rho_2 \mathcal{R} \pi_2)$, then $(\rho_1 \cup \rho_2) \mathcal{R} (\pi_1 \sqcup \pi_2)$.*

Theorem 5 (Soundness of expression evaluation). *For all $e, \sigma, \hat{\sigma}, k$:*

$$\frac{\langle e, \sigma \rangle \rightarrow \langle v, \rho \rangle \quad k, I \vdash \langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}, \pi \rangle \quad \sigma \mathcal{R} \hat{\sigma} \quad true \sqsubseteq eval(k, \phi_\sigma)}{(\langle v, \sigma \rangle, \rho) \mathcal{R} \langle \hat{v}, \pi \rangle}$$

Proof. By induction on the structure of e .

Base case $e \equiv [l]$ In this case, $(v, \rho) = (\sigma[l], \emptyset)$ and $(\hat{v}, \pi) = (\hat{\sigma}[l], \emptyset)$. The premise $\sigma \mathcal{R} \hat{\sigma}$ gives the desired result.

The induction hypothesis asserts that evaluating subexpressions under condition k produces sound results. It remains to show that evaluating operators and traversals under condition k produces sound results.

Case $e \equiv [\text{op}_n^t(e_1 \dots e_n)]$ If the abstract semantics gives $\hat{v} = \top$ for e , then this case is trivially proved. Otherwise, it must be shown that if, for each e_i , the concrete semantics gives (v_i, ρ_i) and the abstract semantics gives (\hat{v}_i, π_i) , then:

$$f_{\text{op}_n}(v_1, \dots, v_n) = f'_{\text{op}_n}(v_1, \dots, v_n)$$

$$eval(\text{op}_n(\hat{v}_1, \dots, \hat{v}_n), \phi_\sigma) = f'_{\text{op}_n}(eval(\hat{v}_1, \phi), \dots, eval(\hat{v}_n, \phi))$$

Because f' is monotonic, it suffices to show that if $(v_i, \sigma) \mathcal{R} \hat{v}_i$, then $v_i \sqsubseteq eval(\hat{v}_i, \phi_\sigma)$. If v_i is a basic value, then the definition of \mathcal{R} suffices. It remains to be shown that if v_i is a record identifier,

$$v_i \sqsubseteq \bigsqcup_{p[k] \in \hat{v}_i} \left\{ \bigsqcup CLoad_2(r_0, p[k], \phi_\sigma) \right\}$$

Because $v_i \mathcal{R} \hat{v}_i$, there exists some paths $\pi' \subseteq \hat{v}_i$ such that $v_i \in CLoad_1(r_0, p', \emptyset)$, where $p' \in \pi'$. Calling $CLoad_1$ on these paths generates a set of iterator field bindings Φ that includes ϕ_σ ; therefore $r \in \bigcup_{p[k] \in \hat{v}_i} CLoad_2(r_0, p[k], \phi_\sigma)$. Hence, the desired result that evaluating all paths in \hat{v}_i with bindings ϕ_σ approximates r . Lemma 3 gives $\bigcup \rho_i \mathcal{R} \bigsqcup \pi_i$.

Case $e \equiv [e.f^t]$ Rules U-TRAVERSE and K-TRAVERSE and the induction hypothesis give $(r, \rho_e) \mathcal{R} (\pi_e, \pi)$ for the subexpression e . For the entire expression, the rules give $\rho_f = Load(r, f)$, $\pi_f = \{p.f^t[k] \mid p \in \pi_e\}$. If $true \sqsubseteq eval(k, \phi_\sigma)$, then $\pi_f = \{p.f^t \mid p \in \pi_e\}$. Section 3.3 proved soundness for this case. Lemma 3 gives $(\rho_e \cup \rho_f) \mathcal{R} (\pi \sqcup \pi_f)$. \square

Theorem 6 (Condition evaluation approximates true). *For all $\sigma, \hat{\sigma}$, conditions k produced by the analysis:*

$$\frac{\sigma \mathcal{R} \hat{\sigma}}{true \sqsubseteq eval(k, \phi_\sigma)}$$

Proof. By induction on the structure of k .

Base case $k = true$ Trivial, because $eval(true, -) = true$.

The induction hypothesis asserts that evaluating subconditions approximates $true$. It remains to prove the theorem for any condition k' the analysis creates.

Case $k' \equiv [k \wedge \hat{v}]$, \hat{v} is a query condition In this case, the analysis attaches k' to all paths generated by the true-branch of an **if**. So, it must be shown:

$$\frac{\langle e, \sigma \rangle \rightarrow \langle true, \rho \rangle \quad k, I \vdash \langle e, \hat{\sigma} \rangle \hat{\rightarrow} \langle \hat{v}, \pi \rangle \quad \sigma \mathcal{R} \hat{\sigma}}{true \sqsubseteq eval(k \wedge \hat{v}, \phi_\sigma)}$$

The induction hypothesis states $true \sqsubseteq eval(k, \phi_\sigma)$, so it remains to show $true \sqsubseteq eval(\hat{v}, \phi_\sigma)$. The induction hypothesis also enables the invocation of Theorem 5, which gives $(true, \sigma) \mathcal{R} \hat{v}$ which is defined to mean $true \sqsubseteq eval(\hat{v}, \phi_\sigma)$.

Case $k' \equiv \llbracket k \wedge \neg \hat{v} \rrbracket$, \hat{v} is a query condition In this case, the analysis attaches k' to all paths generated by the false-branch of an **if**. So, it must be shown:

$$\frac{\langle e, \sigma \rangle \rightarrow \langle \text{false}, \rho \rangle \quad k, I \vdash \langle e, \hat{\sigma} \rangle \rightarrow \langle \hat{v}, \pi \rangle \quad \sigma \mathcal{R} \hat{\sigma}}{\text{true} \sqsubseteq \text{eval}(k \wedge \neg \hat{v}, \phi_\sigma)}$$

Proceeding as above, Theorem 5 gives $\text{false} \sqsubseteq \text{eval}(\hat{v}, \phi_\sigma)$. Since f'_σ is monotonic, $\text{true} \sqsubseteq \neg \text{eval}(\hat{v}, \phi_\sigma)$. A simple analysis on the domain of f'_σ gives $\neg \text{eval}(\hat{v}, \phi) = \text{eval}(\neg \hat{v}, \phi)$, and the desired conclusion is reached. \square

Note that the transfer functions of this extended semantics are also monotonic. Thus the proof of soundness for commands is similar to that of Theorem 4, with appropriate applications of Theorems 5 and 6.

5. Query Creation and Program Simplification

The results of static analysis can now be employed to partition the original program into a query and its client. The query retrieves a subset of the database on which the client program executes. In some cases, the client program may be simplified by removing conditional tests that become redundant when executed on the data subset.

5.1 Query Creation

The concretization of a conditional path corresponds to a query against the database. In its current form, the conditions are associated with the use of individual attributes, yet conditionally loading an attribute is not nearly as useful as conditionally loading an entire record. To avoid loading records, the conditions on individual attributes are *promoted* to apply to iteration fields.

We provide an informal argument, rather than a formal proof, for the validity of promotion. Condition promotion depends on the connection between individual attributes and object loading: An object does not need to be loaded if none of its attributes are needed and if the object does not affect the final outcome of the program. Thus the condition for loading an object is the union (disjunction) of the conditions of all uses of its attributes. In this way, the conditions on attributes are promoted to be conditions on elements of a collection. If any of the paths has the condition *true*, then all elements of the collection will be loaded. In the example, only employees with salary greater than \$65,000 should be loaded.

One important point is the difference between marked and unmarked paths. A marked path $p[k]^*$ is a data dependence of the final state of the program store. An unmarked path only affects the program's control flow—for example, in a condition. Only conditions on marked paths are promoted. Conditions on unmarked paths are ignored, because these paths do not affect that final store.

Queries are created in a variant of the Object Query Language (OQL) [7]. The syntax is:

$$q ::= \mathbf{struct} (f_1 = q_1 \dots, f_n = q_n) \\ | \mathbf{select} \ q \ \mathbf{from} \ q \ \mathbf{as} \ x \ \mathbf{where} \ e \\ | x.\bar{f}$$

where f names a record field, \bar{f} is a sequence of field names, and x is a variable name. We restrict our use of OQL to queries that return a structural subset of the original database. This mirrors the capabilities of commercial products like Hibernate and EJB. Consideration of other query translations is an area for future research.

Figure 11 builds a query from the set of paths that result from the analysis in the previous section. Function Q takes a path p that represents a common prefix for a set of paths π ; it returns a query for the elements reachable by following each suffix from the given

$$Q(p, \{\epsilon\}) = p \quad (\text{Q-PATH})$$

$$\frac{\pi \langle \epsilon \rangle = f_1.\pi_1 \cup \dots \cup f_n.\pi_n \quad f_i \text{ distinct} \\ q_i = Q(p.f_i, \pi_i)}{Q(p, \pi) = \mathbf{struct} (f_1 = q_1, \dots, f_n = q_n)} \quad (\text{Q-FIELDS})$$

$$\frac{\pi = \iota.\pi' \cup \{ \iota.\bar{f}_1[k_1]^*, \dots, \iota.\bar{f}_n[k_n]^* \} \quad \iota.\bar{f}[k]^* \notin \pi' \\ q = Q(\iota, \pi' \cup \{ \bar{f}_1, \dots, \bar{f}_n \}) \\ c = T(k_1) \vee \dots \vee T(k_n)}{Q(p, \pi) = \mathbf{select} \ q \ \mathbf{from} \ p \ \mathbf{as} \ l \ \mathbf{where} \ c} \quad (\text{Q-ITER})$$

$$T(\pi) = \{ \iota.\bar{f} \mid p.\iota.\bar{f} \in \pi \} \\ T(\text{op}_n(e_1, \dots, e_n)) = \bigvee \{ \text{op}_n(e'_1, \dots, e'_n) \mid e'_i \in T(e_i) \}$$

Figure 11. Transforming conditional path sets to queries

prefix. In rule Q-PATH, the prefix is the query when the suffixes are empty.

Rule Q-FIELD handles the case where the suffixes all start with a distinct field name f_i . The query result is a **struct** where each field name is bound to a sub-query for that field. Each field name's query q_i is constructed by appending f_i to current prefix.

In rule Q-ITER, all suffixes start with an iterator field name ι , and the query result selects elements of the collection to which ι refers. If the suffixes begin with different iterator field names, each field name represents a different iteration of the collection. In this treatment we only consider queries that mirror the structure of the database, so only one collection can be returned for a given multi-valued field. Therefore the rule combines queries for multiple iterations. Function Q creates a query for the collection by partitioning suffixes into a set for which ι is the last iterator field name and a set of suffixes π' in which other iterator field names appear. The **select** clause is obtained by forming a query from the prefix ι and the set of suffixes that follow ι . The conditions are removed from the suffixes where no further iterator field names occur and instead are disjoined to form the **select** query's **where** clause.

Function T transforms any paths that may appear in a condition. If a path contains an iterator field name ι , T removes the prefix that appears before ι . Because this field name must be the last to appear in a path, it will be properly scoped by the **as** clause of the **select** query.

Function T also expands sets of paths that may appear in operations. Therefore T disjoins the cross-product achieved by applying the operation to each possible combination of operands.

5.2 Client and Query Simplification

In the next step of the analysis, the data constraints ensured by the query are used to simplify the program, and consequently the data elements in the result of the query. If the program tests a property of the data which is guaranteed by the query, the program test can be removed. Any data that is only used in such tests can then be removed from the query results.

The following two rules are used to simplify the client program:

$$\frac{\Gamma \vdash \hat{v}(e)}{\Gamma \vdash \langle e \rangle \rightarrow \text{true}} \quad \frac{\Gamma \vdash \neg \hat{v}(e)}{\Gamma \vdash \langle e \rangle \rightarrow \text{false}}$$

Each rule includes a context Γ , which is a set of constraints on the persistent data a query returns. The constraints are obtained by taking the conjunction of all the query's **where** clause conditions. The term $\hat{v}(e)$ means the abstract value for e produced by the

rules of Fig. 9. If a context *entails* an expression’s abstract value—written $\Gamma \vdash \hat{v}(e)$ —then the expression can be re-written as *true*. Similarly if the context entails the negation of an expression’s abstract value, the expression can be re-written as *false*. Entailment can be determined by a SAT solver. The rules are applied repeatedly until no more reductions are possible. Once the client has been simplified, a further analysis can remove trivial tests and dead code [36].

After simplifying the client, the query can be simplified by applying the analysis of Fig. 9 to the client and composing the results with the original query. The composite query does not retrieve values that appear only in **where** clauses. The partition for the example program in Fig. 1 is:

```
// define an explicit query
String query =
  "select struct (
    name = e.name,
    manager = struct (name = e.manager.name))
  from Employee as e
  where e.salary > 65000";
// execute the query
List result = session.createQuery(query);
for (Employee e : result.list()) {
  // no test required: all elements already satisfy
  // the condition
  print(e.name + ": " + e.manager.name);
}
```

The function `executeQuery` queries the database and returns a new structure that contains only the data retrieved by the query. The query does not retrieve the employee’s salary, and the program does not test for that value. Instead, the query retrieves only employees whose salary is greater than \$65,000.

6. Related Work

Our path-based approach is similar to research on approximating the shape of pointer data structures [15, 17, 37]. However, we limit ourselves to intraprocedural analysis and focus on the traversal of read-only data structures, not mutation.

Vitenberg et al. describe a path-based abstract interpretation for predicting the persistent values a program may need [35]. Their approach supports runtime improvement of transaction lock scheduling. Kvikval and Singh use shape analysis to dynamically hoard (prefetch) remote data for mobile clients [21]. Their work reduces the effect of disconnections in mobile computing environments. Ours is a fully static approach that supports program transformation to bulk-load persistent data. Our analysis is also unique in that it identifies traversal conditions.

Neubauer and Theimann partition a sequential program run at one location into semantically equivalent, independent, distributed processes [28]. Their approach provides software engineering benefits similar to ours, except for multi-tier applications.

A common technique for integrating programming languages and databases is to make queries first-class values of a programming language. `C#` has been extended to incorporate relational constructs and structured data in middle-tier applications [6]. Willis et al. propose extensions to Java to support database-style optimizations for operations on collections of objects [38]. Safe queries describe queries with classes whose instances are translated into a form that can be executed on a remote database [10]. Unlike our proposal, each of these solutions reduce persistent transparency because they require explicit queries to be written in an extended programming language syntax.

The DBPL language [33] and its successor Tycoon [26] explored optimization of search and bulk operations within the framework

of orthogonal persistence. Tycoon proposed integrating compiler optimization and database query optimization [16]. Queries that cross modular boundaries were optimized at runtime by dynamic compilation [32]. The languages included explicit syntax for writing queries or bulk operations on either persistent or non-persistent data. We do not know of any published formal account of the optimizations used in Tycoon, or any evaluation of its performance or usability.

AUTOFETCH uses a profile-guided dynamic analysis that automatically inserts prefetch directives in queries executed in an object persistence architecture [20]. Because our work generates queries which could be used in object persistence architectures, the two techniques could be combined to achieve further performance benefits.

7. Future Work

While the current analysis provides a unique technique for extracting implicit queries from imperative programs, it contains several restrictions, which we hope to remove or diminish with future work. The imperative language we studied contains no procedures. We are currently extending the analysis to analyze whole programs with behavioral methods and recursive procedures.

Employing standard static analyses (e.g., range analysis) can improve the quality of the extracted queries. These analyses should also allow us to identify and extract aggregation and “exists” sub-queries.

To transform complete programs, more work is needed to identify where the analysis should be applied. Currently a new query is created each time the special variable `root` is used. In some cases it may be more efficient to break a query into parts, so that a result of one query becomes the root of a nested query. Multiple queries could also be used to transform programs in which an outer loop introduces a loop-carried dependence. The expressive power of the target query language also affects these decisions. Other strategies for promoting conditions may also be considered.

Key differences between programming languages and database semantics must be overcome to successfully integrate the two domains. In this paper, we identified two artifacts of the database domain—the three-valued logic of *null* values and the implicit ordering of database sets—that must have appropriate analogues in the programming languages domain.

The current work analyzes only data queries. Future work will extend this analysis to include updates to persistent data. If updates are performed immediately, the resulting aliasing may make it impossible to define a useful transformation for updates. Alternatively, it may be possible to delay the updates until a transaction boundary, at which point all database references must be released.

Finally, the technique must be applied to realistic programs to measure the performance of the analysis and effectiveness of the transformation.

8. Conclusion

We have formalized a new approach for optimizing transparent persistence. This approach extracts a query from an imperative program, then simplifies the program to operate over the bulk-load query results. This technique promises to combine the software engineering benefits of transparent persistence with the performance benefits of query optimization. We expect the current work to serve as a useful foundation for ongoing research into the long-standing effort to integrate programming languages and databases.

References

- [1] J. R. Allen and K. Kennedy. Automatic loop interchange. In *Proc. of the Symp. on Compiler Construction (CC)*, pages 233–246, 1984.

- [2] M. P. Atkinson. Programming languages and databases. In *Proc. of the Intl. Conf. on Very Large Data Bases (VLDB)*, pages 408–419. IEEE Computer Society, 1978.
- [3] M. P. Atkinson, L. Daynès, M. J. Jordan, T. Printezis, and S. Spence. An orthogonally persistent Java. *SIGMOD Rec.*, 25(4):68–75, 1996.
- [4] M. P. Atkinson and R. Morrison. Orthogonally persistent object systems. *VLDB Journal*, 4(3):319–401, 1995.
- [5] C. Batini, S. Ceri, and S. B. Navathe. *Conceptual Database Design - An Entity-Relationship Approach*. Benjamin Cummings, 1992.
- [6] G. M. Bierman, E. Meijer, and W. Schulte. The essence of data access in ω . In *Proc. of the European Conference on Object-Oriented Programming (ECOOP)*, pages 287–311, 2005.
- [7] R. G. G. Cattell, D. K. Barry, M. Berler, J. Eastman, D. Jordan, C. Russell, O. Schadow, T. Stanienda, and F. Velez, editors. *The Object Data Standard ODMG 3.0*. Morgan Kaufmann, January 2000.
- [8] S. Chaudhuri. An overview of query optimization in relational systems. In *Proc. of Symp. on Principles of Database System (PODS)*, pages 34–43, 1998.
- [9] P. P. Chen. The entity-relationship model - toward a unified view of data. *ACM Transactions on Database Systems (TODS)*, 1(1):9–36, 1976.
- [10] W. R. Cook and S. Rai. Safe query objects: Statically typed objects as remotely executable queries. In *Proc. of the Intl. Conf. on Software Engineering (ICSE)*, pages 97–106, 2005.
- [11] G. Copeland and D. Maier. Making smalltalk a database system. In *Proceedings of the 1984 ACM SIGMOD international conference on Management of data*, pages 316–325. ACM Press, 1984.
- [12] P. Cousot and R. Cousot. Systematic design of program transformation frameworks by abstract interpretation. In *Proc. of the ACM Symp. on Principles of Programming Languages (POPL)*, pages 178–190, 2002.
- [13] O. Deux. The O2 system. *Commun. ACM*, 34(10):34–48, 1991.
- [14] J.-A. Dub, R. Sapir, and P. Purich. Oracle Application Server TopLink application developers guide, 10g (9.0.4). Oracle Corporation, 2003.
- [15] M. Emami, R. Ghiya, and L. J. Hendren. Context-sensitive interprocedural points-to analysis in the presence of function pointers. In *PLDI '94: Proceedings of the ACM SIGPLAN 1994 conference on Programming language design and implementation*, pages 242–256, New York, NY, USA, 1994. ACM Press.
- [16] A. Gaweckı and F. Matthes. Integrating query and program optimization using persistent CPS representations. In M. P. Atkinson and R. Welland, editors, *Fully Integrated Data Environments*, ESPRIT Basic Research Series, pages 496–501. Springer Verlag, 2000.
- [17] R. Ghiya and L. J. Hendren. Is it a tree, a DAG, or a cyclic graph? a shape analysis for heap-directed pointers in C. In *Proc. of the ACM Symp. on Principles of Programming Languages (POPL)*, pages 1–15, 1996.
- [18] C. Gould, Z. Su, and P. Devanbu. Static checking of dynamically generated queries in database applications. In *Proc. of the Intl. Conf. on Software Engineering (ICSE)*, pages 645–654, 2004.
- [19] Hibernate reference documentation. http://www.hibernate.org/hib_docs/v3/reference/en/html, May 2005.
- [20] A. Ibrahim and W. Cook. Automatic prefetching by traversal profiling in object persistence architectures. In *Proc. of the European Conference on Object-Oriented Programming (ECOOP)*, 2006.
- [21] K. Kvilekval and A. Singh. SPREE: Object prefetching for mobile computers. In *Distributed Objects and Applications (DOA)*, Oct 2004.
- [22] B. Liskov, A. Adya, M. Castro, S. Ghemawat, R. Gruber, U. Maheshwari, A. C. Myers, M. Day, and L. Shriram. Safe and efficient sharing of persistent objects in Thor. In *Proceedings of the Intl. Conf. on Management of Data (SIGMOD)*, pages 318–329, 1996.
- [23] D. Maier. Representing database programs as objects. In F. Bancilhon and P. Buneman, editors, *Advances in Database Programming Languages*, pages 377–386. New York, NY, 1990.
- [24] D. Maier, J. Stein, A. Otis, and A. Purdy. Developments of an object-oriented DBMS. In *Proc. of ACM Conf. on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 472–482, 1986.
- [25] V. Matena and M. Hapner. Enterprise Java Beans Specification 1.0. Sun Microsystems, 1998.
- [26] F. Matthes, G. Schroder, and J. Schmidt. Tycoon: A scalable and interoperable persistent system environment. In M. Atkinson, editor, *Fully Integrated Data Environments*. Springer-Verlag, 1995.
- [27] R. Morrison, R. C. H. Connor, G. N. C. Kirby, D. S. Munro, M. P. Atkinson, Q. I. Cutts, A. L. Brown, and A. Dearle. The Napier88 persistent programming language and environment. In M. P. Atkinson and R. Welland, editors, *Fully Integrated Data Environments*, pages 98–154. Springer, 1999.
- [28] M. Neubauer and P. Thiemann. From sequential programs to multi-tier applications by program transformation. In *Proc. of the ACM Symp. on Principles of Programming Languages (POPL)*, pages 221–232, 2005.
- [29] B. C. Pierce. *Types and Programming Languages*. MIT Press, 2002.
- [30] T. Rus and E. Van Wyk. A formal approach to parallelizing compilers. In *Proc. of the SIAM Conf. on Parallel Processing for Scientific Computation*, March 14 1997.
- [31] C. Russell. Java Data Objects (JDO) Specification JSR-12. Sun Microsystems, 2003.
- [32] J. Schmidt, F. Matthes, and P. Valduriez. Building persistent application systems in fully integrated data environments: Modularization, abstraction and interoperability. In *Proceedings of Euro-Arch'93 Congress*. Springer Verlag, Oct. 1993.
- [33] J. W. Schmidt and F. Matthes. The DBPL project: advances in modular database programming. *Inf. Syst.*, 19(2):121–140, 1994.
- [34] R. Software. Whitepaper on the UML and Data Modeling, 2000.
- [35] R. Vitenberg, K. Kvilekval, and A. K. Singh. Increasing concurrency in databases using program analysis. In *Proc. of the European Conference on Object-Oriented Programming (ECOOP)*, pages 341–363, 2004.
- [36] M. N. Wegman and F. K. Zadeck. Constant propagation with conditional branches. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 13(2):181–210, 1991.
- [37] R. Wilhelm, S. Sagiv, and T. W. Reps. Shape analysis. In *Computational Complexity*, pages 1–17, 2000.
- [38] D. Willis, D. J. Pearce, and J. Noble. Efficient object querying in Java. In *Proc. of the European Conference on Object-Oriented Programming (ECOOP)*, Nantes, France, 2006.