A Model for Compound Type Changes Encountered in Schema Evolution*

Barbara Staudt Lerner Computer Science Department University of Massachusetts, Amherst

June 12, 1996

Abstract

Schema evolution is a problem that is faced by long-lived data. When a schema changes, existing persistent data can become inaccessible unless the database system provides mechanisms to access data created with previous versions of the schema. Existing systems that support schema evolution focus on changes local to individual types within the schema, thereby limiting the changes that the database maintainer can perform. We have developed a model of type changes incorporating changes local to individual types as well as compound changes involving multiple types. The model describes both type changes and their impact on data by defining derivation rules to initialize new data based on the existing data. The derivation rules can describe local and non-local changes to types to capture the intent of a large class of type change operations. We have built a system called Tess (Type Evolution Software System) that uses this model to recognize type changes by comparing schemas and then produces a transformer that can update data in a database to correspond to a newer version of the schema.

1 Motivation

Databases frequently have long lives. During a database's lifetime, the database schema is likely to undergo significant change as new demands are placed on the data. The database schema serves two purposes. First, it defines an interface for programs and users to query the data contained within the database. Second, it determines how the database management system physically stores the data on the disk. When the schema is changed so that the data can be used for a new purpose, this also impacts the way data is physically stored. The goal of schema evolution research is to allow schema definitions to change while maintaining access to data that has already been stored to disk.

There are two major issues involved in schema evolution. The first issue is understanding how a schema has changed. The second issue involves deciding when and how to modify the database to address such concerns as efficiency, availability, and impact on existing code. Most research efforts have been aimed at this second issue and assume a small set of schema changes that are easy to support, such as adding and removing record fields, while requiring the maintainer to provide translation routines for more complicated changes. As a result, progress has been made in developing the backend mechanisms to convert, screen, or version the existing data, but little progress has been made on supporting a rich collection of changes. The

^{*}This work was supported in part by the Air Force Materiel Command, Rome Laboratory, and the Defense Advanced Research Projects Agency under Contract F30602-94-C-0137 and in part by the National Science Foundation, Grant No. CCR-9504170.

purpose of this work is to enrich the collection of changes supported, independent of the backend mechanism used to manage the data.

Existing database systems that provide schema evolution support changes isolated to individual types within a schema, such as adding a field to a record. More radical changes of representation, such as combining two records are either difficult or impossible with existing database systems. Changes isolated to individual types are not always sufficient, however. A new record may be created to combine the information of several related records, or a large record may be decomposed into several simpler ones. Such changes will clearly impact the representation of persistent data.

The flexibility in data type definition offered by object-oriented databases and persistent programming languages admits the possibility of more complicated changes than those typically encountered in relational database systems, making schema evolution a more difficult problem.

With persistent programming languages the evolution problem is more pervasive than with databases. When using a database, those types that have persistent data are defined in the schema, while transient types are defined in traditional programming languages that interoperate with the database. The transient types can be changed without impacting the persistent data. With persistent programming languages, there is typically no distinction in the programmer's eyes between transient and persistent types. In particular, some persistent programming languages, such as PGraphite, Pleiades, Napier-88, and PS-Algol [WWFT88, TC93, DCBM89, ABC⁺83], treat persistence orthogonally to types. With these languages, an instance of any type can be made persistent dynamically. This approach is very powerful and flexible, since it allows programs to manipulate data uniformly without being concerned about whether it is persistent or transient data. It aggravates the evolution problem, however, because every type potentially has persistent data associated with it. Modifying any type definition can make some persistent data inaccessible.

Unfortunately, there is little published data [GKL94, Sjø93] about how persistent or transient types change during maintenance. Researchers studying maintenance of object-oriented hierarchies, which are not necessarily persistent, cite modifying types in the hierarchy and reorganizing the hierarchy as frequently desirable activities [JO93, OJ93, LBSL91, OJ90, Cas90, MS92a]. One can expect, however, that maintainers are reluctant to make radical changes to an object-oriented hierarchy or any other persistent type or schema definitions if those changes make it difficult or impossible to access existing data. As a result, the maintainers may sacrifice other desirable properties for their schemas and type definitions such as appropriateness of abstractions, modularity, efficiency, etc.

Our goal is to facilitate schema evolution involving complex type changes to allow more natural evolution of persistent types. We have developed a model of type changes incorporating changes local to individual types as well as compound changes involving multiple types. The model describes both type changes and their impact on data by defining derivation rules to initialize new data based on the existing data. The derivation rules can describe local and non-local changes to types to capture the intent of a large class of type change operations.

2 Overview

One way to design a schema evolution system is to define schema modification commands to implement each type change that we want to support. The advantage of this approach is that the maintainer can explicitly inform the schema evolution system of the changes. For each type change, the system defines the effect that the change will have on the data. By choosing the appropriate commands, the maintainer simultaneously modifies the schema and develops a transformer to update the existing data. For example, there may be a command to add fields to a record whose effect on the data is to create a new field set to a default value.

Another command may be used to delete a field from the type. Its effect is to delete the data corresponding to that field. In this type of system, there will be many specialized commands to accomplish all the supported type changes such as changing the size of an array, adding a value to an enumerated type, etc. When dealing with type changes isolated to individual types, it is possible for us to enumerate all the changes that may occur, provide a command for each of these, and define precisely what the effect on existing data will be.

Now, let's consider a more complex type change. Suppose the maintainer wants to move a field from one record to another. If the maintainer applies the *delete-field* command on the original type followed by the *add-field* command on the new type, this will be treated as two separate commands. The semantics of the *delete-field* command will result in deletion of the associated data. The semantics of the *add-field* command will result in addition of a new field set to some default value. To solve this problem we could introduce a new command, *move-field*. Now, there are two ways in which the maintainer can modify the types. If one examines the type definitions after the changes using the two approaches, the definitions are identical. The effect on the persistent data is quite different, however. When using the *delete-field* and *add-field* commands, data is lost. When using the *move-field* command, data is preserved. The maintainer needs to understand this difference and needs to be careful in choosing which commands to use when modifying the types. The problem is that the command approach focuses on the editing *process* rather than the editing *result*. Furthermore, the number of commands would proliferate and the complexity of using the schema evolution system would increase as more complex type changes are supported.

Another alternative is to allow the maintainer to modify the types as necessary and then compare the two versions of the schemas to identify the changes, thereby focusing on the editing result rather than the editing process. The advantage is that the maintainer can edit the schemas using a normal editor, focusing on producing the correct new type definitions without worrying about the exact process used to create those new definitions. The disadvantage is that the system must now infer how the types have changed instead of being explicitly told. In this research, we have developed algorithms to perform these inferences by comparing successive versions of a schema to identify the changes. The schema comparison algorithms use naming similarities, structural similarities, and interrelationships among the types from successive versions to infer the type changes. Experimentation with these algorithms has demonstrated that they can identify a wide variety of type changes successfully.

Of course, it is also possible that the algorithms will make incorrect inferences. As a result, it is important for the maintainer to be involved in the type comparison process. We have incorporated maintainer control into the type comparison algorithms in several ways. First, the maintainer can control which types are compared, if desired. Second, the algorithms can generate multiple inferences of observed type changes from which the maintainer can choose. Third, the algorithms associate a qualitative assessment with each inference indicating the complexity of the change and its impact on the data. The maintainer can use these assessments to set thresholds on the comparison algorithms or to focus attention on the more complex inferences or those with greater impact on the data. Finally, the maintainer can ignore the inferences generated by the algorithm and explicitly tell the system what the impact on the data should be. Using these techniques the maintainer can guarantee that the type changes have the appropriate impact on the existing data. Note that this ability to review the anticipated impact on existing data is useful to remove ambiguity even if schema modification commands are used, particularly in the case where there are multiple commands that lead to the same result as in the *move-field* example earlier.

As another example of a type change that requires understanding the impact on multiple types simultaneously, consider adding a new type to a schema. In most database systems this change is understood in isolation from other changes. When a new type is added, it has no impact on existing data. It is quite likely, however, that the addition of a new type actually represents the reorganization of other types in the

schema. For example, an individual type may be split into two types. The desired effect on the data is to move the data associated with the fields of the new type from objects of the old type to new objects. We could define a *split-type* command to accomplish this, further complicating the maintainer's job. Instead, we develop algorithms to recognize that a type has been added and then look for modified types that may act as sources of information for objects of this new type. Similarly, in most database systems, type deletion results in deletion of objects of the type. Instead, our algorithms look for other modified types that may serve as destinations for the data that would otherwise be deleted.

When supporting changes local to an individual type, the appropriate object changes can be performed by modifying each object in isolation. For example, if a field is deleted from a type, each object of that type can be modified independently of all other objects. The same is not true when supporting changes that affect multiple types. Implementing a single change may require modifying more than one object. Consider moving a field from one type to another again. To implement this change correctly, we must move data from one object to another. This implies that we must identify pairs of objects to operate on. Our algorithms identify collections of objects in two ways. Objects may be related *structurally*. That is by dereferencing fields of one object transitively we may reach other objects that we need to modify. Alternatively, objects may be related by having a common *value*. This is similar to a relational join operation. Identifying these collections of objects is key to being able to implement complex type changes.

Most schema evolution research has addressed the problem of how to update existing data efficiently assuming the type changes are well understood. The emphasis of the research described in this paper is to understand how schemas change during evolution and to develop algorithms that can recognize those changes. Our goal is to represent the schema changes that occur in such a way that their effect on existing data can be accomplished using a variety of data translation mechanisms. For example, in small databases that may belong to an individual user, we can make the database unavailable temporarily and transform all data in the database at once. For large, shared databases, we can employ more sophisticated algorithms such as those developed by Ferrandina [FMZ94] to transform individual objects as they are accessed in order to maintain high availability. In situations where the data is shared by many programs, schema changes may also impact a great deal of code. In those cases, we can use the inferences we produce to define views on the data. Thus, the emphasis of this research is to develop algorithms that can recognize complex type changes made by a maintainer. Instead of constraining the maintainer to perform only supported type changes using a small set of primitive type change commands, we give the maintainer great flexibility in how to change the types. We are addressing the front-end problem of understanding schema changes in a flexible manner to allow integration with a variety of data translation mechanisms.

The remainder of the paper is organized as follows. In Section 3 we describe related work in schema evolution. In Section 4 we present the type model used in this research. In Section 5 we describe a model of how data changes in response to schema changes. In Section 6, we present a model for simple type changes. In Section 7, we present our model of compound type changes. In Section 8, we present an example schema evolution that uses compound type changes. In Section 9, we describe some algorithms developed to compare schemas. In Section 10, we describe some experimental results with the type comparison algorithms. In Section 11 we describe future directions of our research and conclude in Section 12.

3 Related Work

The problem of schema evolution was first addressed with respect to traditional database systems. While many database systems support a few simple changes automatically, such as adding or deleting record fields,

only a few systems [SHL75, Nav80, ST82] support more general transformations. In these cases, the maintainer is responsible for explicitly describing how to convert the data from its old format to its new format using a special-purpose data translation language. This approach is a powerful one, but creation of the transformer is a manual process.

More recent database systems generate transformation functions based upon the changes made to the type definitions. Orion [BKKK87, KK88] and GemStone [PS87] are object-oriented database systems that provide some evolution support. In these systems, evolution is defined in terms of primitive operations that change individual type definitions, such as adding instance variables to a class, removing instance variables from a class, and renaming instance variables. Some type changes are completely automated, but at the expense of limiting the ways in which a maintainer can change type definitions. For example, in Orion the type of an instance variable can only be replaced by a supertype in the type hierarchy. More complex type changes, such as combining two records, are not supported directly. Instead this change is accomplished as several independent changes as follows. The maintainer deletes each instance variable individually from one of the types. The maintainer adds an equivalent instance variable to the second type for each instance variable deleted from the first type. The maintainer modifies all references to the first type to refer to the second type. Finally, the maintainer deletes the first type. Since each change is treated individually rather than as a collection of related changes, deleting the instance variables results in deleting the data contained in those instance variables. To preserve the data, the maintainer must develop code to move the data explicitly. In GemStone the maintainer directly extends the transformer, while in Orion the maintainer must develop and execute programs to move the data prior to deleting the instance variables containing the data.

 O_2 [BFK95] is another object-oriented database system that supports evolution through the use of operations. In addition to primitives similar to those of Orion and GemStone, Q provides high-level operations to manipulate the class hierarchy. These high-level operations are defined as a composition of primitive operations. As a result, they provide better support for the maintainer in expressing type changes and preserving data. For example, the Abstraction-Generalization operation can be used to create a new superclass that generalizes a set of existing classes. While these high-level operations support more complex changes than previous systems, defining type changes via a pre-defined set of operations necessarily restricts the kinds of type changes that are supported. In particular, while they have numerous operations to allow the definition of new classes and migration of existing objects to these new classes, none of their operations allow simultaneous modification of multiple types such as moving a field from one existing type to another.

Another approach to schema evolution relies on the simultaneous maintenance of multiple versions of types and data[SZ86, Cla94, Bra92, MS92b, TS92]. With these approaches, multiple versions of the same type exist within a single database. The advantage is that old and new code can operate on old and new data without requiring either to be changed. The disadvantage is that the maintainer must provide routines to make data appear to be of the version of the type that the code is expecting. This approach admits more general changes, but it still limits changes to be isolated to individual types. It also results in significant overhead (in both space and time) for maintaining and accessing multiple type and data versions. Odberg [Odb94] extends the versioning approach to the entire schema, which is versioned when a type is modified. This allows the description of changes that simultaneously affect multiple types, but still requires the maintainer to define the translation routines between versions.

TransformGen [GKL94] is a system to support evolution of abstract syntax grammars used by Gandalf programming environments [HN86, HGN91]. The abstract syntax grammars are analogous to type

¹Note that if the type being deleted is used as the type of an instance variable, we cannot in general replace its type with the second type since the second type is not necessarily a supertype of the first type. In that case, we would need to delete the instance variable and create a new instance variable of the desired type.

definitions; they define the format of the abstract syntax trees stored in databases maintained by Gandalf environments. The abstract syntax changes for which TransformGen automatically generates transformation routines are analogous to the type changes supported by Orion and GemStone. TransformGen goes beyond these two systems, however, by allowing the maintainer to modify the generated transformation using a declarative data manipulation language. In this way, the maintainer can perform complex type changes using the primitive operations provided and then easily fix the generated transformations to have the intended effect. The significance of this extensibility is that the resulting transformers can handle arbitrary type changes, including those involving multiple types, but without requiring the maintainer to write transformation routines. While the maintainer can extend the transformer, there is little guidance in identifying the limits of the generation process and the situations that require extension. OTGen [LH90] is a system designed using the concepts developed in TransformGen to support flexible transformation of object-oriented databases. As such it has many of the features and limitations of TransformGen, but is aimed at a more general type system.

4 Type Model

Before discussing the details of the type change model, we must first present the type model that we use. The type model is a language-independent type model that captures features common to many programming languages. We are concerned with the structural aspects of the type model as those are most relevant to understanding the impact of schema changes on persistent data. As a result, we do not treat the types as abstract types, although they may, in fact, be implemented abstractly. In our examples, we therefore present the type representations used, but not the interfaces or operations belonging to those types.

A schema consists of a collection of type definitions. Schema changes are performed by editing types within the schema. Editing a schema is treated as an atomic operation, independent of how many types are modified in the process.

The type model includes the predefined types of character, integer, string, and boolean. Programmers can define new types using the following constructors: record, bounded and unbounded array, set, multi-set, union, enumeration, subrange, pointer, and alias.

Data is organized into objects. Each instantiation of a type results in the creation of a new object. Each object has an object identifier to allow objects to reference each other. Each object is tagged with its type. An object can be made persistent at any time. When an object is made persistent, all other objects reachable from that object are also made persistent. Any object can serve as the root of such a persistent structure.

Because the type model does not a priori restrict persistence in any way, the schema evolution support must be very general as it must support changes to any type within a schema.

5 Object Changes

What makes schema evolution an interesting and difficult problem is not that types change, but rather the impact that those type changes have on persistent objects. Therefore we begin by presenting a model of how objects can change as a result of schema evolution. Following that, we describe how types can change and relate type changes to object changes.

²This persistence model could be changed without impacting the research presented here significantly. For example, the model could allow the maintainer to restrict persistence to a subset of the types, allow only a subset of the types to be roots of persistent structures, or not automatically make all objects reachable from a persistent instance persistent.

There are fundamentally three object operations associated with evolution: initialization, derivation, and deletion. New objects can be *initialized* to a default value. New objects can be *derived* from existing objects. Existing objects can be *deleted*. As derivation is the only technique that involves both existing and new objects, it is of greatest interest.

Derivation rules define how to derive new objects from existing objects. A derivation rule specifies a source type, a destination type, and a derivation function. The *source type* is a type from the schema before modification. It identifies the type of an existing object to transform. The *destination type* is a type from the schema after modification. It identifies the type of the new object to create. The *derivation function* is a function to apply to a source object to create a destination object. The simplest derivation function simply copies an existing object unmodified. A more complicated function might traverse the persistent structure starting at the source object to perform a more complex derivation such as summing a collection of values to produce a total, or selecting the median from a collection of values.

When evolving an object, we apply the derivation rule associated with the type of the existing object to create a new object. A derivation rule for a structured type, such as a record, is typically defined using other derivation rules. For example, to derive a new record object, it is necessary to assign a value to each new record field. The fields may be initialized to a default value or themselves derived from existing objects.

6 Simple Type Changes

We categorize simple type changes as being either local type changes or reference type changes (Figure 1 defines a complete list of all simple type changes in our type model.) A *local type change* affects the structure of an individual type, such as adding a record field or changing the bounds of a subrange. A local type change affects data local to individual objects. The effects of local type changes can be expressed with derivation rules that derive each new object from a single old object. For example, a derivation rule for a record type can capture all local changes to records by initializing new fields, deleting fields no longer belonging to the record type, and providing a one-to-one mapping between the fields present in both the old and new versions of the record.

A reference type change replaces a type used within a type constructor with another type, such as changing the type of a record field or an array element. To fully understand how the constructed type is changed, it is necessary to understand the relationship between the old and new reference types. The effects of each reference type change are described with a derivation rule from the old reference type to the new reference type. This separation of concerns makes the derivation rules easier to understand since each derivation rule describes changes local to an individual object. For example, when deriving a new record field from an old one, we would refer to a derivation rule defined between the type of the old field and the type of the new field.

Existing database systems that support schema evolution interpret all type changes as simple type changes similar to those outlined above. Changes that make objects of a type smaller, such as deleting a record field, result in deletion of data. Those that make objects of a type larger, such as adding a record field, result in fields initialized to a default. Reference type changes result in application of the derivation rule for the reference type. The only derivation rules produced by these systems are rules that derive new objects by copying values local to the corresponding old object. Definition of non-local derivation rules is left to the maintainer.

• Local type changes:

- Creating or deleting a type
- Changing the name of a type
- Changing the type constructor of a type
 - * Changing an array type to a set or multi-set type, or vice versa.
 - * Changing a set type to a multi-set type, or vice versa.
 - * Replacing one scalar type with another.
- Changing a type constructor argument
 - * Adding an enumeration value, deleting an enumeration value, reordering enumeration values, or renaming an enumeration value.
 - * Changing the lower or upper bounds of a subrange type.
 - * Adding a record field, deleting a record field, reordering record fields, or changing the name of a record field.
 - * Adding an array dimension, deleting an array dimension, changing the bounds of an array dimension, or reordering array dimensions.

• Reference type changes:

- Changing the type referenced by a pointer or alias type.
- Changing the type a subrange is defined over.
- Changing the type of a record field.
- Changing the index type of an array dimension.
- Changing the type of array, set, or multi-set elements.

Figure 1: Simple Type Changes

7 Compound Type Changes

For database systems to support non-local derivation rules, they must have a richer model of type changes. These *compound type changes* modify more than one type and as a result affect more than one object. Compound type changes compose three basic kinds of type operations, type deletion, type creation, and type modification, to produce more complex type changes. As with simple type changes, the fundamental *object* change that is desired is derivation of the value for a new field from the value of one or more old fields. In the case of compound type changes, however, the old and new fields belong to different types, not different versions of the same type. Each compound type change could be modeled as a collection of simple type changes, where old fields are deleted from their types and the new field is added to a different type. In doing so, however, the ability to describe non-local derivation is lost. In our model we include compound type changes whose effects on objects are defined with non-local derivation rules. In Figure 2, we list the compound type changes in our model. In this section we define the compound type changes provided by our model.

- **Inline** Replacing a type reference with its type definition.
- **Encapsulate** Creating a new type by encapsulating parts of one or more old types.
- Merge Replacing two or more type definitions with a new type that merges the old type definitions.
- Move Moving part of a type definition from one type to another existing type.
- **Duplicate** Duplicating part of a type definition in another type definition.
- **Reverse link** Reversing the connection between two types.
- Link addition Adding a link between two existing types.

Figure 2: Compound Type Changes

7.1 Type Deletion

When deleting a type, a database maintainer is either reorganizing the type system or removing functionality. In the former case, the fields of the deleted type most likely become associated with another type, either a new type or an existing type. In these cases, the data associated with the deleted type should be moved to existing instances of the modified/created type.

One kind of compound change involving type deletion is inlining. *Inlining* involves replacing a use of a type with the type definition. Figure 3 provides an example of inlining. Here the address field is replaced with a collection of fields previously contained in the Address type. The new field values are derived from fields of the old Address object. If this compound type change were viewed as a collection of simple type changes, the new fields would be uninitialized and the old Address object would be deleted.

Another compound type change involving type deletion is merging. *Merging* deletes two or more object types and creates a new type that represents the integration of the deleted types. Figure 4 provides an

Old version: New version:

type Person is
name: string;
address: Address;
end Person;

type Person is
name: string;
street: string;
city: string;
state: string;
type Address is

type Person is
name: string;
street: string;

street: string; end Person; city: string; state: string; zipcode: integer;

end Address;

example of a merge type change. Here two or more objects must be located and combined to define a new object. In the example, PersonalInfo and EmployeeInfo objects that have the same value in their name field will be combined. This merge change finds its pairs of objects by joining on the name field. If the name field does not serve as a key for the two types, the results are ambiguous.

Figure 3: Inlining

As these two examples indicate, for complex type changes to be integrated into a schema evolution system, it must be possible to identify collections of objects to modify instead of individual objects as with simple type changes. The inlining example showed a relationship between objects based on a structural connection, while the merging example showed a relationship based on equivalent values. A database maintainer could define other relationships as well.

7.2 Type Creation

There are two type creation operations that are analogous to the type deletion operations. The merging compound type change discussed above also involves type creation. The second type creation operation is encapsulation. *Encapsulation* produces the opposite effect of inlining. Here one or more fields are replaced with a single field. The type of the new field includes the old field type(s) as a reference type(s). An example of encapsulation can be seen by swapping the old and new versions in Figure 3. As with inlining, the relationship between objects is structural.

7.3 Type Modification

Compound type changes may involve the modification of types without requiring types to be created or deleted. There are four kinds of type changes fitting this description: *moving*, *duplication*, *link reversal*, and *link addition*.

Both moving and duplication involve deriving a new field from an old field. The difference is that moving deletes the original field while duplication maintains the original field. As with simple type changes, the derivation associated with moving and duplication may derive a new value, not just copy the old value.

Old version:

New version:

type PersonalInfo is
 name: string;
 address: Address;
 phone: Phone;
 marital_status: MaritalStatus;
 num_children: integer;
end PersonalInfo;

type EmployeeInfo is
 name: string;
 id: integer;
 salary: integer;
end EmployeeInfo;

type Person is name: string; address: Address; phone: Phone;

marital_status: MaritalStatus; num_children: integer;

id: integer;
salary: integer;

end Person;

Figure 4: Merge

Figure 5 shows the *address* and *phone* fields being moved from the *Personal Info* type to the *Person* type. In this case, the corresponding objects are identified using their structural relationship, specifically, the *Personal_Info* and *Person* objects that are connected using the *Person.personal* field are modified together. Figure 6 shows duplication between objects with a value relationship. Here the *id* field is duplicated from the *EmployeeInfo* object to the *PersonalInfo* object with the same value in their *name* fields.

Link reversal involves reversing the direction of a pointer. For example, consider Figure 7. Originally, the *PersonalInfo* type has a pointer to the *EmployeeInfo* type. In the modified version, *EmployeeInfo* has a pointer to the *PersonalInfo* type. Here we are reversing the structural relationship between two types.

Link addition involves adding a link between two existing types. The difference between this change and the simple type change of adding a record field is that in the former case we expect the value of the new link field to be an existing object, while in the latter case we expect to create a new value for the new field. Once again, we can use either structural or value relationships to identify pairs of objects to add a link between. For example, Figure 8 shows the addition of an inverse link between two structurally connected types.

7.4 Limitations of the Compound Type Change Model

While this model of compound type changes allows a schema evolution system to develop non-local derivation rules, the maintainer still needs to be involved directly in the definition of derivation rules for two reasons. First, the default for both local and non-local derivation rules is to copy old values. If the maintainer wants to use a different function, such as summing a collection of values, or finding a median, the maintainer must provide this function explicitly. Second, the merge, move, duplicate, and link addition type changes require finding collections of old objects to operate on. It may be necessary for the maintainer to indicate how to find matching objects to operate on, particularly if the relationships are not structural or by equivalent values.

Old version:

New version:

type Person is	type Person is
name: string;	name: string;
<pre>personal: Personal_Info;</pre>	address: Address;
end Person;	home_phone: Phone;
	personal: Personal_Info;
type Personal_Info is	end Person;
address: Address;	

phone: Phone;
marital_status: MaritalStatus;
num_children: integer;
end Personal_Info;

type Personal_Info is marital_status: MaritalStatus;

num_children: integer;

end Personal_Info;

Figure 5: Moving Using a Structural Relationship

Old version:

New version:

```
type PersonalInfo is
                                         type PersonalInfo is
    name: string;
                                             name: string;
    address: Address;
                                             id: integer;
    phone: Phone;
                                             address: Address;
    marital_status: MaritalStatus;
                                             phone: Phone;
    num_children: integer;
                                             marital_status: MaritalStatus;
end PersonalInfo;
                                              num_children: integer;
                                         end PersonalInfo;
type EmployeeInfo is
    name: string;
                                         type EmployeeInfo is
    id: integer;
                                              name: string;
    salary: integer;
                                             id: integer;
end EmployeeInfo;
                                              salary: integer;
                                         end EmployeeInfo;
```

Figure 6: Duplication Based on Value Relationship

Old version:

New version:

type PersonalInfo is	type PersonalInfo is
name: string;	address: Address;
address: Address;	phone: Phone;
phone: Phone;	marital_status: MaritalStatus;
marital_status: MaritalStatus;	num_children: integer;
num_children: integer;	end PersonalInfo;
emp_info: EmployeeInfo;	
end PersonalInfo;	type EmployeeInfo is
	name: string;
type EmployeeInfo is	id: integer;
id: integer;	salary: integer;
salary: integer;	<pre>private_info: PersonalInfo;</pre>

Figure 7: Link Reversal

Old version:

end EmployeeInfo;

New version:

end EmployeeInfo;

```
type PersonalInfo is
                                        type PersonalInfo is
    name: string;
                                             name: string;
    address: Address;
                                             address: Address;
                                             phone: Phone;
    phone: Phone;
    marital_status: MaritalStatus;
                                            marital_status: MaritalStatus;
    num_children: integer;
                                             num_children: integer;
    emp_info: EmployeeInfo;
                                             emp_info: EmployeeInfo;
end PersonalInfo;
                                        end PersonalInfo;
type EmployeeInfo is
                                        type EmployeeInfo is
    id: integer;
                                             id: integer;
    salary: integer;
                                             salary: integer;
end EmployeeInfo;
                                             private_info: PersonalInfo;
                                        end EmployeeInfo;
```

Figure 8: Link Addition

8 Example

Compound type changes can be combined to produce interesting schema changes whose effects on existing data can be understood following the model given. In this section, we provide an example of a real schema evolution and describe how it fits into the compound type change model.

Figure 9 shows consecutive versions of a collection of interrelated types extracted from TAOS, a software testing tool [Ric93]. In this example, we see three modified types and four new types.

Old version:

New version:

```
type SaveTestCases is (nada, todo);
                                                  type TestCaseState is (Pass, Fail, Untested);
type RandomTestInfo is
                                                  type SaveTestCases is array ( TestCaseState )
    MinLength: natural := 0;
                                                       of boolean;
    MaxLength: natural := 0;
    NumberRequired: positive := 1;
                                                  type RandomTestInfo is
    Persistence: SaveTestCases := todo;
                                                       MinLength: natural := 0;
    NumberNonPersistentPassed: natural := 0:
                                                       MaxLength: natural := 0;
    NumberNonPersistentFailed: natural := 0;
                                                       NumberRequired: positive := 1;
end:
                                                  end:
type TestClass is
                                                  type Saved is ( persistent, nonpersistent );
    ExtraInfo: RandomTestInfo;
end:
                                                  type TestCaseCounts is
                                                       array ( Saved, TestCaseState ) of natural;
                                                  type TestCasesInfo is
                                                       PersistencePreferences: SaveTestCases :=
                                                           Default_Persistence;
                                                       NumTestCases: TestCaseCounts :=
                                                           Default_Counts;
                                                  end:
                                                  type TestClass is
                                                       TestSetInfo : TestCasesInfo := Create;
                                                       ExtraInfo: RandomTestInfo;
                                                  end;
```

Figure 9: Schema Evolution in TAOS

If we consider each type in isolation, we see the following simple type changes: three fields have been deleted from RandomTestInfo, one field has been added to TestClass, SaveTestCases has changed from an enumerated type to an array of booleans, and four new types have been created. Treating these as simple type changes would result in the deletion of the values associated with the deleted fields of RandomTestInfo, the initialization of the new field in TestClass to its default value, and the deletion of values of the SaveTestCases type.

Now, let's reconsider the example as a sequence of compound type changes.

• Moving: The Persistence, NumberNonPersistentFailed, and NumberNonPersistentPassed fields are moved to the TestClass type using a structural relationship..

```
type RandomTestInfo is
    MinLength: natural := 0;
    MaxLength: natural := 0;
    NumberRequired: positive := 1;
end;

type TestClass is
    Persistence: SaveTestCases := todo;
    NumberNonPersistentPassed: natural := 0;
    NumberNonPersistentFailed: natural := 0;
    ExtraInfo: RandomTestInfo;
end;
```

• Encapsulation: The NumberNonPersistentFailed and NumberNonPersistentPassed fields are encapsulated into a new field named NumTestCases whose type is the new TestCaseCounts type. Specifically, the value of the NumberNonPersistentFailed field is moved to the TestCaseCounts element indexed by (nonpersistent, Fail). The value of the NumberNonPersistentPassed field is moved to the TestCaseCounts element indexed by (nonpersistent, Pass).

```
type TestCaseState is (Pass, Fail, Untested);
type Saved is ( persistent, nonpersistent );
type TestCaseCounts is array ( Saved, TestCaseState ) of natural;
type TestClass is
    Persistence: SaveTestCases := todo;
    NumTestCases: TestCaseCounts := Default_Counts;
    ExtraInfo: RandomTestInfo;
end:
```

• Encapsulation: The Persistence field is encapsulated into an attribute named PersistencePreferences whose type is the new SaveTestCases. The value of the field is duplicated in each element of the SaveTestCases array, translating nada to false and todo to true.

```
type TestCaseState is (Pass, Fail, Untested);

type SaveTestCases is array ( TestCaseState ) of boolean;

type TestClass is
    PersistencePreferences: SaveTestCases := Default_Persistence;
    NumTestCases: TestCaseCounts := Default_Counts;
    ExtraInfo: RandomTestInfo;
end:
```

• Encapsulation: PersistencePreferences and NumTestCases are encapsulated into a new field named TestSetInfo whose type is the new TestCasesInfo type.

```
type TestCasesInfo is
    PersistencePreferences: SaveTestCases := Default_Persistence;
    NumTestCases: TestCaseCounts := Default_Counts;
end;

type TestClass is
    TestSetInfo : TestCasesInfo := Create;
    ExtraInfo: RandomTestInfo;
end;
```

This example demonstrates the type change model. It also demonstrates that describing these type changes via editing commands would be cumbersome. We have developed type comparison algorithms to support such changes without requiring the user to specify them with explicitly.

9 Type Comparison

For our type comparison approach to be feasible, we assume that between successive versions of a system, most type definitions remain mostly the same. We rely heavily upon the similarities that exist to quickly prune the space of types that must be compared. Since the types in databases tend to experience evolutionary change, rather than revolutionary change, we do not expect this to be a significant problem for most situations. In those situations in which revolutionary changes occur, the maintainer can and should provide more guidance rather than relying on the fully automated control algorithm. In Section 10.1, we describe how the maintainer can provide guidance in our implementation of the type comparison algorithms.

In this section, we describe derivation rules in more detail. Next, we describe the algorithm that controls which types are compared. Following that, we describe the algorithm to recognize simple changes that

may occur in a record definition. Then, we present the algorithm to identify movement of fields between structurally-connected records, including encapsulation and inlining. Finally, we explain how derivation rules could be used with a variety of data translation mechanisms.

9.1 Derivation Rules

A *derivation rule* describes how to translate data created using one type definition to a different type definition. For simple values, such as integers and enumerated types, the derivation rule defines a function to apply to the old value to compute the new value. In the simplest derivation rules, the function is simply an identity function. For example, suppose we have a *Counter* type in our old and new schema. Assume that this *Counter* type is unmodified. The corresponding derivation rule is the following:

```
Counter \Rightarrow Counter: new := old:
```

Note the use of the keywords *old* and *new*. *old* refers to the existing data that we are translating from. *new* refers to the new data that we are creating. In this case, the new data has the same value as the old data.

For structured types, such as records, the function specifies how to compute the value for each new substructure of the new type. Usually, the value of a new substructure is defined in terms of the value of existing data. As a result, the derivation of most substructures is performed by applying the derivation rule defined between the types of the corresponding substructures. A new substructure may be defined using a constant value or a user-supplied function. For example, Figure 10 is the derivation rule that corresponds to Figure 5:

```
Person \Rightarrow Person:
```

new.name: derive from old.name;

new.address : derive from old.home.address;
new.home_phone : derive from old.home.phone;

new.home : derive from old.home;

Figure 10: Derivation Rule for a Structured Type

In this case, the new fields are all derived from existing data. For example, the value for the new *home* field is computed by applying the derivation rule from the old *Personal Info* type to the new *Personal Info* type.

In some cases, we may want to use slightly different derivation rules between a pair of types depending upon the state of the database. For example, suppose we replace one type with a collection of types. The intent may be to partition the existing values so that each value belongs to one of the new types. To support this, we add conditionals to our derivation rules. Consider the type change and corresponding derivation rule shown in Figure 11. Here we will create a different type of new object depending on the value of an existing field.

A *similarity metric* is associated with each derivation rule. A similarity metric is a qualitative description of the impact that applying the derivation rule would have on existing persistent data. For example, derivation rules between record types have one of the following similarity metrics:

type Plane is type Jet is engine: EngineType; num_passengers: positive; num_passengers: positive; max_speed: positive; max_speed: positive; end: end; type PropellorPlane is type PlaneFleet is set (Plane); num_passengers: positive; max_speed: positive; end: type Glider is num_passengers: positive; max_speed: positive; end; **type** *Plane* **is union of** (*Jet*, *PropellorPlane*, *Glider*); **type** *PlaneFleet* **is set** (*Plane*); $PlaneFleet \Rightarrow PlaneFleet$ for each old_plane in old $if {\it old_plane.engine} = JetEngine$ **let** *new_plane* = *Jet* **derived from** *old_plane* **else if** *old_plane.engine* = *PropellorEngine* **let** new_plane = PropellorPlane **derived from** old_plane **else if** *old_plane.engine* = *None* **let** new_plane = Glider **derived from** old_plane end if; insert new_plane in new end for;

New version:

Old version:

Figure 11: Conditionals in Derivation Rules

Similarity Metric Meaning

Identical No changes to the type

FieldOrderChange Fields appear in a different order.

FieldTypeNameChange Name of the type of a field has changed.

FieldNameChange
NewField
New type has an extra field.
DeletedField
Old type has an extra field.

Each derivation rule has a single metric that describes the worst effect of applying the rule. Thus a rule with a *NewField* metric may also have fields whose names have changed, but it will not have any deleted fields. Similarity metrics are used within the comparison algorithms to prune the space of comparisons considered. (In Section 10.1, we will also describe how similarity metrics are used to focus the maintainer's attention on the derivation rules with greatest impact on the data.)

9.2 The Type Comparison Control Algorithm

The input to the type comparison algorithms is the set of type definitions of consecutive schemas. The algorithms selectively compare the types to identify how the types have changed and output derivation rules describing how to transform instances of the old version into instances of the new version. The type comparison control algorithm is responsible for determining which types to compare, based primarily on the results of comparisons done thus far and on naming similarities between old and new types, as well as which comparison algorithms to use based on the type constructors used by the types being compared. The algorithms ignore changes to white space and comments and the order in which the type definitions appear in the schema.

The fully-automated type comparison control algorithm is shown in Figure 12. It proceeds through three stages. First, in the *name comparison* stage, old and new types that have the same names in both versions are compared. For structured types, such as records and arrays, this may result in further type comparisons. For example, a derivation rule that derives a new array from an old array requires a derivation rule from the old array element type to the new array element type. Comparing these element types is called *component comparison*. In the second stage, called *use site comparison*, types that use types that have been successfully compared are compared. In the final stage, called *exhaustive comparison*, each old type that does not already have a derivation rule is compared to each new type, first considering only those new types that use the same type constructor and, if that fails to produce an acceptable derivation rule, considers all remaining new types. The exhaustive comparison algorithm also performs component comparisons and use site comparisons as derivation rules are generated. Thus if a derivation rule is found by exhaustive comparison, the algorithms immediately compare pairs of types used by the matched type pair as well as pairs of types using the matched type pair. This further reduces the search for matching types.

To better understand the stages of the algorithms, consider the type definitions in Figure 13. *old record* is an old type and *new_record* is the corresponding new type. Since they have different type names they are not compared during the name comparison phase. The two versions of *field2 type* (not shown) are compared in this phase. Assuming a derivation rule is found between these types, the use site comparison stage searches for pairs of types that use *field2 type*. It finds *old_record* and *new_record* and *new_record* and *new_field1 type* are compared during component comparison since they have the same field name.

```
procedure CompareTypes (old_types, new_types) is
begin

    Compare types with the same name.

    for each type o in old_types
        let n = \text{type in } new\_types with the same name as o
         if Compare(o, n) finds a derivation rule then
             add (o, n) to TypePairList
         end if;
    end for ;
    - Check the use sites for each pair of types that have a derivation rule.
    for each type pair tp in TypePairList
         let o = \text{old type in } tp
        let n = \text{new type in } tp
        - Find where the old and new type are used.
        let old_uses = set of types in old_types that use o
        let new_uses = set of types in new_types that use n
         - Compare each pair of use sites.
         for each type o_u in old_uses
             for each type n_u in new_uses
                  if Compare (o_u, n_u) finds a derivation rule then
                      add (o_u, n_u) to TypePairList
                  end if;
             end for;
         end for :
    end for;
    - Exhaustive search
    for each type o in old_types
         - Make sure we have at least one derivation rule for each old type
        if there is no derivation rule from o to any type in new\_types then
             - Compare to new types with the same type constructor.
             for each type n in new_types with the same type constructor
                  if Compare(o, n) finds a derivation rule then
                       compare the use sites of o and n
                  end if:
             end for;
             - Compare to new types with different type constructors.
             for each type n in new_types with a different type constructor
                  if Compare(o, n) finds a derivation rule then
                      compare the use sites of o and n
                  end if:
             end for :
         end if:
    end for;
end;
```

Figure 12: Type Comparison Control Algorithm

type old_record is record
 field1: old_field1_type;
 field2: field2_type;
end record;

type new_record is record
 field1: new_field1_type;
 field2: field2_type;
end record;

Figure 13: Component and Use Site Comparisons

9.3 Recognizing Simple Type Changes: A Sample Algorithm

When looking for simple changes between two types, the algorithm varies depending upon the type constructors that the types use. For example, a different algorithm is used to compare two enumerated types than to compare two record types. We have also developed algorithms to compare two types that use different type constructors, such as sets and arrays.

In Figure 14 we show the algorithm that compares two records to give more insight into how the type comparisons proceed. The input to this algorithm is the type definitions of two record types. The output is a derivation rule between those record types, such that each record field of the new type is either derived from an old record field or is initialized to a default value, and each record field of the old type that is not used in a derivation is explicitly identified as being deleted.

The record comparison algorithm is quite similar to the algorithm used to compare the sets of type definitions. First, it compares record fields with the same name. Next, it compares old unmatched fields with new unmatched fields with the same type name. Finally, it compares each old unmatched field to each new unmatched field. When it compares fields, it compares the field names and recursively compares the field types. If type definitions are recursive, as with linked lists for example, recursive comparison of field types leads to an infinite loop. To avoid this, we use an algorithm similar to the one used by Amadio and Cardelli to check subtyping of recursive types [AC93], which limits the recursion performed when comparing recursive types. We cache the results of type comparisons in a matrix so that we can look up the results of previous comparisons instead of repeating them.

Using algorithms such as the one described here we can recognize changes equivalent to those supported by databases that provide automatic support for schema evolution, including Orion and GemStone.

9.4 Recognizing Compound Type Changes: A Sample Algorithm

We have also developed algorithms to recognize compound changes. This allows us to support type changes not supported by other databases. Figure 15 shows the algorithm for recognizing movement and encapsulation of fields from one record type to another where the types are structurally related. This algorithm is passed an old record type, a new record type, and the derivation rule constructed by the algorithm to detect simple changes in record types. In this initial derivation rule, old fields that may be sources of movement are marked as deleted while new fields that may be the destinations of movement are marked as uninitialized. After applying this algorithm, the derivation rule identifies the structural connections that must be traversed to move the data from old fields to the corresponding new fields.

To accomplish this task, the algorithm identifies all the unused fields of the old type in the derivation rule. It also transitively finds all unused subfields of any field of the old type. These are fields that might be moved.

```
function CompareRecords (old_record, new_record) return derivation_rule is
begin
    let r = new derivation rule from old_record to new_record;

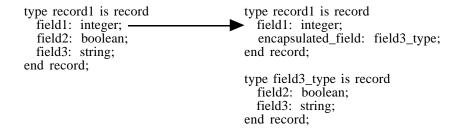
    Compare fields with the same name

    for each field o_f in old_record
         for each field n_f in new_record
              if o_f and n_f have the same names
                   if Compare (type of o_{\underline{f}}, type of n_{\underline{f}}) finds a derivation rule
                        map o_{-}f to n_{-}f in r;
                        mark o_f and n_f as used;
                   end if:
              end if:
         end for:
    end for;
     - Compare fields with the same type name
    for each field o_f in old_record
         if o_f is not used
              for each field n_f in new_record
                   if n_{\underline{f}} is not used and
                        o_f and n_f have different names and
                        o_f and n_f have the same type names then
                        if Compare (type of o_f, type of n_f) finds a derivation rule
                            map o_f to n_f in r;
                            mark o_f and n_f as used;
                        end if:
                   end if:
              end for;
         end if:
    end for;
    - Compare unmatched old fields to unmatched new fields
    for each field o_f in old_record
         if o_f is not used
              for each field n_f in new_record
                   if n_f is not used and
                        o_f and n_f have different names and
                        o_f and n_f have different type names then
                        if Compare (type of o_{\underline{f}}, type of n_{\underline{f}}) finds a derivation rule
                            map o_f to n_f in r;
                            mark o_f and n_f as used;
                        end if:
                   end if;
              end for;
         end if:
    end for;
end;
```

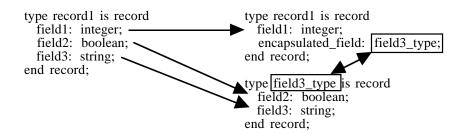
Figure 14: Record Comparison Algorithm

```
procedure RecordFieldMove (old_record, new_record, deriv_rule) is
begin
    create an empty record type o
    for each field o_f in old_record
         if o_f is unused in deriv_rule then
             add o_f to o;
         end if;
         add unused fields and subfields of o_f in deriv_rule to o
    end for
    create an empty record type n
    for each field n_f in new_record
         if n_f is unused in deriv_rule then
                add n_f to n;
         end if;
         add unused fields and subfields of n_f in deriv_rule to n
    end for
    move\_rule = CompareRecords (o, n);
    for each field mapping fm in move_rule
         copy fm to deriv_rule
    end for;
end;
```

Figure 15: Record Field Move Algorithm



Before Checking for Encapsulation



After Checking for Encapsulation

Figure 16: Encapsulation Example

It then constructs a dummy record definition whose fields are these unused fields and subfields. The field names used in these dummy record definitions encode the path to the real subfield so that this information can be used to identify the source of a moved field. In a similar manner, a second dummy record definition is created to hold all the unused fields and subfields of the new type, again encoding the path to the real subfield. Next, the algorithm applies the record comparison algorithm for recognizing simple type changes given in Figure 14. Field mappings identified by comparing these two dummy types necessarily involve a subfield from the old, new, or perhaps both types, since all mappings between fields of the old and new type have been identified prior to calling the compound type comparison algorithm. These mappings correspond to compound type changes. Each field mapping identified by this algorithm is then merged into the original derivation rule. In this way, the derivation rule now encodes both local and non-local changes.

Figure 16 graphically shows the derivation rules that the encapsulation algorithm finds for a particular set of type definitions. The top of the figure shows the result of comparing the old and new versions of record1 looking for only simple type changes. The field2 and field3 fields of the old version are unused; the nested_field of the new version is unused. After looking for compound type changes, field2 is mapped to

encapsulated_field.field2 and field3 is mapped to encapsulated_field.field3.

Using algorithms such as this one we can recognize compound type changes where the relationships between the source and destination types is structural.

9.5 Using Derivation Rules to Perform Data Translation

The derivation rules do not specify when objects are updated or whether those updates are persistent. This separation is done deliberately to allow derivation rules to be used by a variety of data translation mechanisms. In this section, we briefly describe how the derivation rules would be used by several different data translation mechanisms.

9.5.1 Conversion

In a conversion backend, data are read from their old format, converted to the new format, and written back to the database. GemStone and O_2 are examples of object-oriented database systems that perform conversion. Conversions can be performed by taking the database offline, starting at each root object, visiting each object in turn, converting it, and writing the result to the database.

If the database is large or high availability is required, it may not be feasible to take the database off-line. In these cases, lazy conversion can be done as in O₂ [FMZ94]. With lazy conversion individual objects are converted as they are accessed. To work with conversion, we would apply a derivation rule locally, but would only convert components of structured objects as they were accessed. Ferrandina presents a solution to the problem of ordering conversions in lazy conversion to ensure that data is not deleted before it is transformed. Our derivation rules can be used with his algorithms.

9.5.2 Screening

With screening, information is never deleted from objects. Instead the accessing functions hide the appropriate information based upon the version of the code that is accessing the object. Orion [BKKK87, KK88] uses screening to evolve its objects. To do so, it uses a clever object representation that allows fast access to objects even after the object's type has been changed. It also restricts the kinds of changes to a subset of our local change model. The changes that they allow minimize the impact on the persistent objects, but reduce the flexibility available to the maintainer.

Derivation rules can extend the screening approach to more complex type changes. This is accomplished by applying the derivation rules as objects are accessed. For example, suppose we want to move data from one object to another conceptually. When an object is accessed from which data has moved, the data should be hidden by the accessing function, just as deletion is managed currently. When an object is accessed to which data has moved, the accessing function must apply the portion of the derivation rule that defines where the data comes from to access the data. The advantage of this approach is that it has minimal impact on the data, just as current screening techniques. Furthermore, there is no need to take the database off-line. The disadvantage is that accessing objects associated with these complex type change operations will pay a penalty on each access. Additionally, we must be careful to not delete objects that contain data that serves as the source of a data movement operation.

9.5.3 Versioning

In a versioning backend, such as Encore [SZ86] provides, multiple versions of an object may exist at one time. The runtime system compares the version of code accessing an object with the versions available so that the correct version can be returned. If the correct version does not exist, it is created dynamically by applying the appropriate derivation rules. Current systems that use versioning only support local type

changes. In order to derive a new version of an object, one only needs to access an existing version of that object. Our derivation rules could be used in a similar manner to support non-local derivations.

To support this, derivation rules must be able to translate from newer versions of a type system to older versions. This could be done by applying the same comparison algorithms but changing which is being used as the source and which as the destination. A more straightforward technique would be to define reverse transformations directly by analyzing the existing derivation rules.

Thus, we see that the basic notion of derivation rules is quite flexible and could be used with a variety of data translation mechanisms to address different database concerns such as minimizing access times, maintaining high availability, and reducing the impact on existing code.

10 TESS: An Experimental Type Comparison System

We have implemented type comparison algorithms in a tool called Tess. Tess can automatically generate derivation rules for all simple type changes listed in Figure 1. It can also generate derivation rules for inlining and encapsulation, as well as merging of structurally-connected types and moving fields between structurally-connected types. The maintainer has extensive control over what is automated: Tess can operate in modes ranging from completely manual, as in early database systems, to fully automated. The combination of powerful derivation rule generation algorithms and flexible application of those algorithms as determined by a maintainer leads to a synergism not found in existing systems. In this section, we describe the maintainer's role in running Tess, how Tess assures that all necessary derivation rules have been provided, and then discuss some experimental results.

10.1 Maintainer Control over Type Comparison

We provide an interactive user interface to Tess that gives the maintainer control over how type comparisons proceed. There are three dimensions that the maintainer has control over. First, the maintainer can control which stages of the comparison control algorithm are used (name comparison, use site comparison, and exhaustive search).

The second dimension that the maintainer can control is which types get compared. Here there are three options. All types can be considered at once (the fully-automated control algorithm shown in Figure 12), a specific old and new type can be compared, or an individual old type can be compared with all new types. When schema changes are minor, it is reasonable for the maintainer to use all stages of type comparison and allow all types to be compared. When the schema has changed radically, it would be better to use only the name comparison algorithm on all types to identify the unchanged types and obvious changes and then complete the transformer by specifying pairs of types for Tess to develop derivation rules for.

The third dimension involves determining which derivation rules are automatically accepted as correct and which must be presented to the maintainer for manual acceptance. This is done by defining a threshold value for the similarity metric. The most conservative approach accepts only derivation rules between unchanged types, that is, simple identity rules. A more liberal policy accepts changes in which all old values still belong to the new type, such as increasing the size of a subrange. The most liberal, yet still sensible, policy automatically accepts those changes that affect the representation of the type within the database, but not its use within a program, such as reordering the fields of a record. If a schema contains many unchanged types or trivial changes, the use of similarity metrics allows the user to quickly focus on the interesting changes.

10.2 Assuring Completeness of Derivation Rules

Since generation of derivation rules is a separate activity from updating the persistent data, it is important that all the necessary derivation rules are produced so that they are available when we later attempt to access old data. In this section, we describe how we ensure this.

Recall that our model of persistence assumes that any type can be the root of a persistent structure. As a result, we require a *root derivation rule* for each old type. The root derivation rule is used to translate an instance of a type that appears as an old root into an object of the new schema. Frequently, the old and new types of a root derivation rule have the same type name, but this is not necessarily so.

Recall also that when we apply a derivation rule that creates a structured type, it generally assigns values to the components of the structure by applying another derivation rule (as in Figure 10). Unlike root derivation rules, Tess can determine which old and new types the derivation rule must operate on by examining the types of the components that are paired. These derivation rules are referred to as *reference derivation rules*. For each accepted derivation rule (root or reference), we examine the derivation rules used within the accepted rule to determine what pairs of types require reference derivation rules.

Using this information, Tess keeps track of which types still require derivation rules. A user may decide that not all types actually are used as the roots of structures and thus some types might not require root derivation rules. In contrast, the analysis of which reference derivation rules are required is precise. If a reference derivation rule is missing, a runtime error would occur if the derivation rule that used the missing reference derivation rule was applied. Tess displays this status information to the user, indicating which old types do not yet have root derivation rules and which type pairs referenced by accepted derivation rules do not have derivation rules. Requiring completeness ensures that we will be able to transform any old data that we might encounter.

10.3 Experimentation

Figure 17 shows the results of applying Tess to the example shown in Figure 9. The compound change from the old *TestClass* type to the new *TestClass* type is correctly identified. If the compound change algorithm had not been applied, the *TestClass* to *TestClass* derivation rule would have initialized the *TestSetInfo* field to a default value. With the compound change algorithm, we see that the *Persistence*, *NumberNonPersistentPassed*, and *NumberNonPersistentFailed* fields are moved from the *RandomTestInfo* type to the *TestCasesInfo* type. Objects that the data should move between are connected structurally through the old and new *TestClass* type. The data moves from the old *TestClass.ExtraInfo* to the new *TestClass.TestSetInfo* field. Of particular interest is the movement of data from *TestClass.ExtraInfo.Persistence* to *TestClass.TestSetInfo.PersistencePreferences*. The movement is accomplished by applying the reference derivation rule between the old and new *SaveTestCases* types. The definition of the *SaveTestCases* type has changed considerably, however. In the old version, it was an enumerated type of two values. In the new version, it is an array of booleans. The derivation rule generated by Tess specifies that the old value should be placed in the first element of the new array, applying the derivation rule between *SaveTestCases* and *boolean* to compute the new value.

This experiment demonstrates Tess's flexibility in recognizing simple and compound changes as well as the need for continued human involvement in the development of powerful derivation rules. While Tess generates the correct derivation rule for the compound type change. It does not generate the correct rule for the *SaveTestCases* type. The generated rule takes the old value and places it in the first element of the new array. The correct derivation rule would duplicate the old *SaveTestCases* value into all elements of the new *SaveTestCases* array. It is not possible to distinguish between the case of insertion into a single array

```
TestClass ⇒ TestClass: Compound Change

Handle each field as follows:

ExtraInfo ⇒ ExtraInfo

See the mapping from RandomTestInfo to RandomTestInfo

ExtraInfo.Persistence ⇒ TestSetInfo.PersistencePreferences

See the mapping from SaveTestCases to SaveTestCases

ExtraInfo.NumberNonPersistentPassed ⇒ TestSetInfo.NumTestCases(persistent, Pass)

ExtraInfo.NumberNonPersistentFailed ⇒ TestSetInfo.NumTestCases(persistent, Fail)

RandomTestInfo ⇒ RandomTestInfo: Deletes Old Component
```

Handle each field as follows:

MinLength ⇒ MinLength

MaxLength ⇒ MaxLength

NumberRequired ⇒ NumberRequired

SaveTestCases ⇒ SaveTestCases: Requires New Component Element indexed by Pass: See the derivation rule from SaveTestCases to boolean Elements indexed by Fail and Untested will be uninitialized.

```
SaveTestCases ⇒ boolean: Value Change
Handle each value as follows:
nada ⇒ false
todo ⇒ true
```

Figure 17: The Derivation Rules Generated by Tess for the TAOS Example

element and duplication in more than one array element by looking at the type definitions alone. It requires a more semantic understanding of the change and thus we expect the maintainer to provide this information.

The similarity metrics on the derivation rules shown result in Tess requiring approval of the derivation rules before they would be applied (for any reasonable threshold for automatic acceptance). Out of a total of 37 derivation rules generated by the complete example, only two other derivation rules required review and both of those derivation rules were correct. Thus, even though the totally-automated algorithm did not produce a completely correct set of derivation rules, it did focus the maintainer's attention on the few complicated situations that existed. Even in the case where the derivation rule was wrong, the changes required to correct the derivation rule were quite minor relative to the overall complexity of the derivation rules.

This example demonstrates capabilities for which existing evolution systems provide no automated support. The change of *SaveTestCases* from an enumerated type to an array of booleans cannot be done in existing automated systems. The movement of fields from *RandomTestInfo* to the *TestSetInfo* field of the *TestClass* type would result in deletion of the associated data with existing automated evolution systems. Systems that require the maintainer to provide the transformation routines would allow proper handling of these transformations, but development of those routines would be entirely manual.

11 Future Work

The type comparison algorithms currently implemented in Tess compare types based upon the structure of those types. Zaremski and Wing have demonstrated the use of type comparison to locate components in a library for reuse [ZW95a, ZW95b]. The type comparison algorithms that they use rely on type signatures and formal specifications. Since signatures and formal specifications generally change less frequently than representations, incorporating these algorithms into Tess may improve Tess's ability to find matching types in old and new versions of a system. The algorithms to compare types at the representational level are still required to produce the transformers between the types. Using signatures and formal specifications in comparisons may also make it apparent that the database must evolve to respond to changing semantics of the types, even when the representations are unmodified. For example, if a list type is changed from an unsorted list to a sorted list, the representation would not be changed, but the existing values would still not be appropriate to use with the new definition.

Type change is also an issue for dynamic module replacement systems whose goal is to replace program components without stopping execution of a program. In this case there is existing data that may need to be transformed even though it is not necessarily persistent data. Existing systems (such as [Fab76, FS91]) recognize the need for such transformation functions, but leave the development of those functions to the maintainer. Tess's comparison algorithms could be used to generate these transformation functions.

Another situation in which type comparison may be applicable is schema integration. Here the goal is to develop derivation rules between the types defined in interoperating databases in order that they can share data. In this scenario, the role of the maintainer will become more important as the assumption of naming similarities will most likely be violated. Also, the maintainer would be able to provide valuable guidance in distinguishing between types whose data should be shared and types whose data should remain encapsulated within one database. With the maintainer's guidance, derivation rules could be developed between schemas to allow the necessary data sharing to occur.

12 Conclusions

During software maintenance, a maintainer is typically expected to increase the functionality of software and improve its performance while maintaining backward compatibility. Backward compatibility is required so that existing users will not need to be retrained to use the new version of the system, and so that existing persistent data can continue to be used. With traditional approaches to managing persistent data, it is typically impractical to make major changes to types for which there is persistent data. This restriction in changing type definitions complicates the design and implementation of the desired functionality and performance modifications.

Our research into persistent type evolution addresses the problem of modifying types for which persistent data exists. Specifically, we have defined a model of type changes that describes the complex type changes we have observed in maintenance histories of real systems. We have developed algorithms to recognize these type changes and to generate derivation rules that can translate data from an old representation to the new representation. By doing so, we offer the maintainer much greater flexibility in the modification of persistent types than traditional database systems do.

Acknowledgments

I would like to thank Lori Clarke for her support of this work. Additionally, I would like to thank Lori, Peri Tarr, Lee Osterweil, and Rick Lerner for their helpful comments on earlier versions of this paper.

References

- [ABC⁺83] M.P. Atkinson, P.J. Bailey, K.J. Chisholm, W.P. Cockshott, and R. Morrison. An approach to persistent programming. *The Computer Journal*, 26(4), 1983. Also published in *Readings in Object-Oriented Database Systems*, Stanley B. Zdonik and David Maier, eds., Morgan Kaufman, San Mateo, California, 1990.
- [AC93] Roberto M. Amadio and Luca Cardelli. Subtyping recursive types. *ACM Transactions on Programming Languages and Systems*, 15(4):575–631, September 1993.
- [BFK95] Philippe Breche, Fabrizio Ferrandina, and Martin Kuklok. Simulation of schema change using views. In *Proceedings of the 6th International Conference on Database and Expert Systems Applications*, London, September 1995.
- [BKKK87] Jay Banerjee, Won Kim, Hyoung-Joo Kim, and Henry F. Korth. Semantics and implementation of schema evolution in object-oriented databases. In *Proceedings of the ACM SIGMOD 1987 Annual Conference*, pages 311–322, San Francisco, May 1987.
- [Bra92] S.E. Bratsberg. Unified class evolution by object-oriented views. In *Proceedings of the 11th International Conference on the Entity-Relationship Approach*, pages 423–439, Karlsruhe, Germany, October 1992.
- [Cas90] Eduardo Casais. Managing class evolution in object-oriented systems. In Dennis Tsichritzis, editor, *Object Management*, pages 133–195. Université de Genève, Switzerland, 1990.

- [Cla94] Stewart M. Clamen. Schema evolution and integration. *Distributed and Parallel Databases: An International Journal*, 2:101–126, 1994.
- [DCBM89] Alan Dearle, Richard Connor, Fred Brown, and Ron Morrison. Napier88—a database programming language? In *Proceedings of the Second International Workshop on Database Programming Languages*, pages 179–195. Morgan Kaufmann, 1989.
- [Fab76] R. Fabry. How to design a system in which modules can be changed on the fly. In *Proceedings* of the International Conference on Software Engineering, pages 470–476, Los Alamitos, CA, 1976.
- [FMZ94] Fabrizio Ferrandina, Thorsten Meyer, and Roberto Zicari. Implementing lazy database updates for an object database system. In *Proceedings of the 20th International Conference on Very Large Databases*, pages 261–272, Santiago, Chile, September 1994.
- [FS91] O. Frieder and M. Segal. On dynamically updating a computer program: From concept to prototype. *Journal of Systems and Software*, pages 111–128, February 1991.
- [GKL94] David Garlan, Charles W. Krueger, and Barbara Staudt Lerner. TransformGen: Automating the maintenance of structure-oriented environments. *ACM Transactions on Programming Languages and Systems*, 16(3):727–774, May 1994.
- [HGN91] Nico Habermann, David Garlan, and David Notkin. Generation of integrated task-specific software environments. In Richard F. Rashid, editor, *CMU Computer Science: A 25th Anniversary Commemorative*, Anthology Series, chapter 4, pages 69–97. ACM Press, Reading, Massachusetts, 1991.
- [HN86] A. Nico Habermann and David Notkin. Gandalf: Software development environments. *IEEE Transactions on Software Engineering*, SE-12(12):1117–1127, December 1986.
- [JO93] Ralph E. Johnson and William F. Opdyke. Refactoring and aggregation. In *Proceedings of ISO-TAS '93: International Symposium on Object Technologies for Advanced Software*, November 1993.
- [KK88] Hyoung-Joo Kim and Henry F. Korth. Schema versions and DAG rearrangement views in object-oriented databases. Technical Report TR-88-05, University of Texas at Austin, February 1988.
- [LBSL91] Karl J. Lieberherr, Paul Bergstein, and Ignacio Silva-Lepe. Abstraction of object-oriented data models. In H. Kangassalo, editor, *Entity-Relationship Approach: The Core of Conceptual Modelling*, pages 89–102. Elsevier Science Publishers B.V., 1991. Similar to Johnson and Opdyke work on refactoring.
- [LH90] Barbara Staudt Lerner and A. Nico Habermann. Beyond schema evolution to database reorganization. In *Proceedings of the Joint ACM OOPSLA/ECOOP '90 Conference on Object-Oriented Programming: Systems, Languages, and Applications*, pages 67–76, Ottawa, Canada, October 1990.
- [MS92a] S. Mellor and S. Shlaer. *Object Lifecycles: Modeling the World in States*. Yourdon Press Computing Series. PTR Prentice-Hall, 1992.

- [MS92b] S. Monk and I. Sommerville. A model for versioning classes in object-oriented databases. In *Proceedings of the Tenth British National Conference on Databases*, Aberdeen, Scotland, 1992.
- [Nav80] Shamkant B. Navathe. Schema analysis for database restructuring. *ACM Transactions on Database Systems*, 5(2):157–184, June 1980.
- [Odb94] Erik Odberg. MultiPerspectives: The classification dimension of schema modification management for object-oriented databases. In *Proceedings of TOOLS-USA '94*, Santa Barbara, California, August 1994.
- [OJ90] William F. Opdyke and Ralph E. Johnson. Refactoring: An aid in designing application frameworks and evolving object-oriented systems. In *Proceedings of 1990 Symposium of Object-Oriented Programming Emphasizing Practical Applications*, 1990.
- [OJ93] William F. Opdyke and Ralph E. Johnson. Creating abstract superclasses by refactoring. In *Proceedings of CSC '93: The ACM 1993 Computer Science Conference*, February 1993.
- [PS87] D. Jason Penney and Jacob Stein. Class modification in the GemStone object-oriented DBMS. In *Proceedings of the ACM Conference on Object-Oriented Programming Systems, Languages, and Applications*, pages 111–117, Orlando, Florida, October 1987.
- [Ric93] Debra J. Richardson. TAOS: Testing with analysis and oracle support. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA) '94*. ACM Press, June 28-30 1993.
- [SHL75] Nan C. Shu, Barron C. Housel, and Vincent Y. Lum. CONVERT: A high level translation definition language for data conversion. *Communications of the ACM*, 18(10):557–567, October 1975.
- [Sjø93] D. I. K. Sjøberg. *Thesaurus-Based Methodologies and Tools for Maintaining Persistent Application Systems*. PhD thesis, University of Glasgow, July 1993.
- [ST82] Ben Shneiderman and Glenn Thomas. An architecture for automatic relational database system conversion. *ACM Transactions on Database Systems*, 7(2):235–257, June 1982.
- [SZ86] Andrea H. Skarra and Stanley B. Zdonik. The management of changing types in an object-oriented database. In *Proceedings of the ACM Conference on Object-Oriented Programming Systems, Languages, and Applications*, pages 483–495, September 1986.
- [TC93] Peri Tarr and Lori A. Clarke. Pleiades: An object management system for software engineering environments. In *Proceedings of ACM SIGSOFT '93: Symposium on the Foundations of Software Engineering*, Los Angeles, CA, December 1993.
- [TS92] Markus Tresch and Marc H. Scholl. Meta object management and its application to database evolution. In *Proceedings of the 11th International Conference on the Entity-Relationship Approach*, pages 299–321, Karlsruhe, Germany, October 1992.
- [WWFT88] Jack C. Wileden, Alexander L. Wolf, Charles D. Fisher, and Peri L. Tarr. PGRAPHITE: An experiment in persistent typed object management. In *Proceedings 3rd Software Development Environments Conference*, pages 130–142, December 1988.

- [ZW95a] Amy Moorman Zaremski and Jeannette M. Wing. Signature matching: A tool for using software libraries. *ACM Transactions on Software Engineering and Methodology*, 4(2):146–170, April 1995.
- [ZW95b] Amy Moorman Zaremski and Jeannette M. Wing. Specification matching of software compo. In *Proceedings of SIGSOFT '95: 3rd ACM SIGSOFT Symposium on the Foundations of Software Engineering*, 1995.