

# A C++ Data Model Supporting Reachability Analysis and Dead Code Detection

Yih-Farn Chen, Emden R. Gansner, Eleftherios Koutsosios.

**Abstract**—A software repository provides a central information source for understanding and reengineering code in a software project. Complex reverse engineering tools can be built by analyzing information stored in the repository without reparsing the original source code. The most critical design aspect of a repository is its data model, which directly affects how effectively the repository supports various analysis tasks. This paper focuses on the design rationales behind a data model for a C++ software repository that supports reachability analysis and dead code detection at the declaration level. These two tasks are frequently needed in large software projects to help remove excess software baggage, select regression tests, and support software reuse studies. The language complexity introduced by class inheritance, friendship, and template instantiation in C++ requires a carefully designed model to catch all necessary dependencies for correct reachability analysis. We examine the major design decisions and their consequences in our model and illustrate how future software repositories can be evaluated for completeness at a selected abstraction level. Examples are given to illustrate how our model also supports variants of reachability analysis: impact analysis, class visibility analysis, and dead code detection. Finally, we discuss the implementation and experience of our analysis tools on a few C++ software projects.

**Keywords**—C++, Conceptual Modeling, Dead Code Detection, Program Database, Software Repository, Reachability Analysis, Reverse Engineering, Static Analysis.

## I. INTRODUCTION

There has been a growing trend [27][1][26] in building software repositories to help maintain structure information of existing legacy code. A software repository provides a central information source for understanding and reengineering code in a software project. While many variants of repository-based systems have been constructed, there is not a clear agreement on how the repository should be organized. One popular approach is to store variants of abstract syntax trees in the repository, such as those used in Reprise [25], ALF [22], Genoa [10], Cobol/SRE [23], PRODAG [24], Aria [11], and Rigi [20] in the IBM program understanding project [3]. Because of the nature of the representations, tree traversal routines are frequently used to generate various abstractions.

The other popular approach, which was adopted in the construction of the C/C++ Information Abstraction Systems [5][8][16] and in XREFDB [19], is to structure the repository as a relational database so as to reuse the large body of existing database technology. Complex reverse engineering tools can be built by applying queries to the generated database.

In principle, the two approaches are equally expressive

for a given level of granularity. (In practice, abstract syntax trees typically are built to contain information of much finer granularity.) In terms of which queries are most easily asked and most quickly answered, the two structures have very different behavior. Abstract syntax trees naturally support queries on local structure efficiently but queries that involve regular expressions on selected relationships may require the traversal of the entire tree.<sup>1</sup> Conversely, it is cumbersome to use relational databases to discern certain information on program structure such as finding all the dependency paths between two program entities. The problem can be partially remedied by exploiting the duality between directed graphs and entity-relationship databases with directed binary relationships. These databases can be converted to graphs, and vice versa, in order to support certain query operations efficiently.

In the database approach, the most critical design component of a repository is its data model, which directly affects how effectively the repository supports various analysis tasks. This paper focuses on the design of a data model for a C++ software repository that, among other goals, supports reachability analysis and dead code detection, two tasks that are frequently needed in large software projects. They serve as the basis for various reverse engineering tasks, such as detecting unnecessary include files [29], performing selective regression testing [9], and computing objective software reuse metrics [7]. This paper also examines how we use our model to implement reachability analysis and its variants (including dead code detection) in the context of the C++ programming language.

Researchers have found it frustrating to compare reverse engineering tools even on simple criteria such as how well they extract function call graphs [21]. The difficulty arises because different underlying models of these tools give different interpretations of what a function call means. On the other hand, the ability for a repository to support complete reachability analysis as defined in this paper is an objective criterion, for a selected abstraction level, that programmers and researchers can use to compare different repository implementations.

A complex, object-oriented language such as C++ makes constructing a data model adequate for reachability analysis a complicated and delicate process. In addition to the entities and relationships found in typical procedural languages, C++ introduces such additional relationships as inheritance, friendship, access adjustments and template

Y.-F. Chen, E. R. Gansner, and E. Koutsosios are with AT&T Labs - Research, 180 Park Ave., Florham Park, NJ 07932. USA. Email: {chen,erg,ek}@research.att.com

<sup>1</sup>Systems taking the abstract syntax tree approach will sometimes provide a symbol table as an auxiliary structure. If rich enough, this can act like a relational database. For another approach to combining the two techniques, see Horwitz and Teitelbaum [17].

instantiation, which affect the analysis in various ways. As an example, Figure 1 shows what we expect to obtain from a simple query like

*Find all C++ entities reachable from the class Pool*

`Pool` is a class in a C++ components library developed in AT&T that manages a set of same-size memory blocks. This picture reveals several key relationships used in C++ reachability analysis:

- *containment relationship*: `Pool`  $\rightarrow$  `Pool::alloc()`
- *friendship relationship*: `Pool`  $\rightarrow$  `Pool_element_header`
- *inheritance relationship*: `Pool`  $\rightarrow$  `Block_pool_ATTLC`
- *reference relationship*: `Pool::purge`  $\rightarrow$  `Pool::head`

These relationships and the template instantiation relationship will be examined in detail when we discuss our C++ model.

This paper is organized as follows. We start by presenting our C++ data model and explain the rationales behind many design decisions. We then discuss how the model supports various flavors of reachability analysis in C++ programs, followed by a description of our implementation and an experience report on some sample C++ code we collected from a few C++ software projects. Finally, we discuss our future plans for the C++ software repository and research opportunities.

## II. A C++ DATA MODEL

Our C++ data model is formulated using Chen's entity-relationship modeling [4] and it supports both the C and C++ programming languages. We consider a C or C++ program as a collection of source entities referring to each other, an entity representing a static, syntactic construct such as a macro, a type, a function or a variable. Since our focus is on creating a *complete* data model that supports reachability analysis and dead code detection, we need to provide a clear definition of *completeness*:

*Completeness*: A data model  $M$  of a programming language  $L$  is considered *complete* if, for any two entities  $a$  and  $b$  in the model, a dependency relationship  $a \rightarrow b$  also exists in  $M$  when one of the following two conditions holds:

- $C1$ : if the *compilation* of the entity  $a$  depends on the existence of a declaration of the entity  $b$ .
- $C2$ : if the *execution* of the entity  $a$  depends on the existence of the entity  $b$ .

For example, if  $a$  is a source file that includes a header file  $b$ , then  $a \rightarrow b$  should be captured according to  $C1$ . Similarly, if  $a$  is a variable initialized with a macro  $b$ , or a class that inherits from class  $b$ , or a template class instantiated from class template  $b$ , then  $a \rightarrow b$  should exist as well because  $a$  cannot be compiled without a declaration of  $b$ .

On the other hand, if a function  $a$  calls or refers to a function  $b$ , even if  $b$  is not declared (as is allowed in some C programs and shell scripts), then  $a \rightarrow b$  should exist in the model according to  $C2$ . For a discussion on conditions required (*well-defined memory* and *well-bounded pointer*) for static analysis tools to capture such relationships, directly or transitively, refer to [9].

A model that satisfies the completeness criterion allows us to define *reachable entity set* and *dead entity set* in the following way:

*Reachable Entity Set*: A reachable entity set  $R(a)$  is the set of entities reachable from an entity  $a$  through standard closure computations on the dependency relationships in the model.

*Source Entity Set*: A source entity set  $S$  is simply the set of all entities in a program according to the model.

*Dead Entity Set*: A dead entity set  $D(r)$  is simply the difference between  $S$  and  $R(r)$ , where  $r$  is the entity that serves as the starting point of the program execution.  $D(r)$  is the set of program entities that are not needed for the compilation or execution of the program. The notion of a dead entity set can be generalized by making  $r$  a set when there is more than one entry point to a system.

Note that a function referenced through a pointer (such as in C and C++) does not present a problem to the computation of a reachable or dead entity set because the pointer assignment has to occur somewhere on the dependency paths from `main` to the point where a function is called.

The first design choice we have to make in designing a *complete* model is the entity granularity. There are several possibilities, moving from coarser to finer granularity:

- *file*: This is the granularity used by most source code control systems such as RCS [28] and configuration management tools such as *nmake* [13].
- *top-level declaration*: This creates entities for all constructs not defined within function bodies, *plus* the nested components of any entity representing a type. This level could be expanded to include entities for all declarations.
- *atomic*: This models all program information, down to the level of statements and expressions. It captures the complete static syntactic and semantic information of the program. This is used by syntax-directed editors, and is the basis of many commercial software browsers and debuggers.

Each choice has its own consequences. For example, the file granularity does not allow dead entity declarations in a source file to be detected, while the declaration granularity implies that questions concerning detailed flow control cannot be answered. In addition, as the granularity becomes finer, the size of the repository and the time of queries can be prohibitive for real software projects. Whatever entity granularity is selected, all relationships among entities at that level of abstraction must be captured in the repository in order for reachability analysis to be complete and accurate.

Our model uses the granularity of top-level declarations. This includes entities for types, functions, variables, macros and files. We feel that this level provides adequate information for the vast majority of analyses pertaining to issues of software engineering, while avoiding the excessive overhead of finer granularities. It captures the principal structural artifacts of a program, especially those used across modules, functions, and classes.

Our C++ model significantly expands and cleans up an earlier model proposed in [16]. In particular, we have added

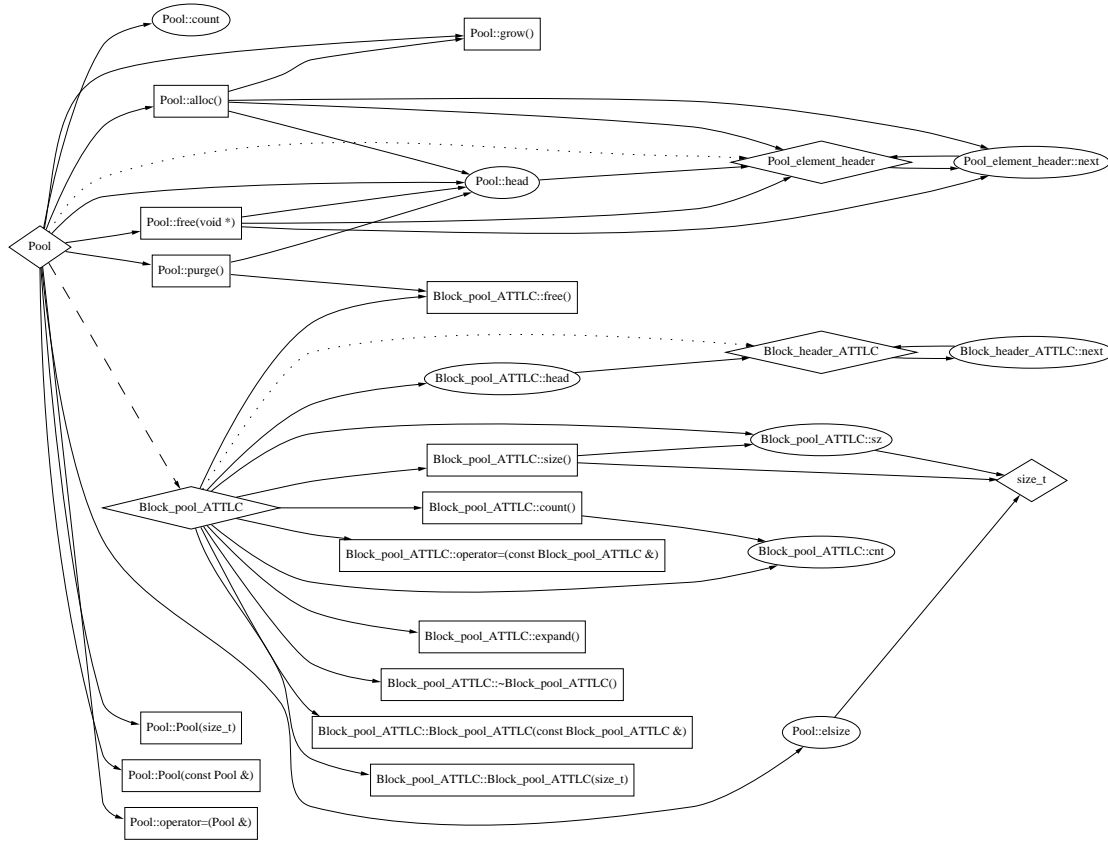


Fig. 1. The reachability graph of the class `Pool`. Boxes stand for functions, diamonds for types, and ovals for variables. Class inheritance relationships are shown in dashed lines, friendships in dotted lines, and all other relationships in solid lines.

support for template-related entities and relationships, and enforce consistent reference relationships in nested class declarations. Both are required for complete reachability analysis.

In the following, we first discuss the basic attributes shared by all C++ entities and then discuss additional attributes that are required for each entity kind. Many of these attributes are used to modify the behaviors of variants of reachability analysis.

#### A. Common Entity Attributes

The attributes that are shared by all C++ (and C) entities in our model include *unique id*, *entity kind*, *entity name*, *source file*, *location*, *definition/declaration flag*, and *checksum*, a numeric value associated with the entity's contents. All C++ entities, including functions, variables, types, macros, and files, have these common attributes.

Note that even such a simple model, without additional attributes discussed later, already implies that several queries are possible:

- count entities: we can count the number of entities of each *kind*.
- search entities: we can find out if an entity with a certain name pattern exists in the database.
- retrieve entity source: we can retrieve and view the source code of each entity by using the *source file* and *lo-*

*cation* attributes. This is useful if software entities are to be rearranged or packaged for reuse.

- detect entity changes: With two versions of a database, we can find out the lists of entities that are deleted, added, or changed from the old version to the new by examining their corresponding checksums.

For example, a query like “count the number of deleted function definitions from version 1 to version 2” can easily be handled by a difference database, which has an additional entity attribute that specifies whether an entity is deleted, added, or changed.

#### B. Principal Entities

The entities for types, functions and variables form the basis for most significant program analyses. In addition, C++ allows the constructs to be declared as members inside classes and structs. We call these entities *principal entities*, which require three additional attributes:

- *scope*: A member can be either *private*, *protected*, or *public*. The scope of a non-member entity is either *extern* (global scope) or *static* (file scope). The scope of an entity affects the *visibility analysis* discussed later.
- *parent*: This attribute records the containing parent class or struct of a member, if any.
- *subkind*: In C++, a type entity can be a class template, template class, typedef, struct, class, union or an enum. A

function entity can be a function template, template function, or just an ordinary function. Instead of creating an entity kind for each of these variations, we use this attribute to distinguish among them.

Since a type can be a class template or an instantiated template class, an additional attribute *type param* is necessary to store the formal parameters in the former case and actual parameters in the latter case. For example, the following shows a class template **Map** and an instance of it retrieved from one of our C++ databases for standard C++ headers:

```
<class S, class T>Map
Map<Set_or_Bag_hashval, unsigned int>
```

Similarly, since a function can be a function template or an instantiated template function, we need an additional attribute *function param* to store the template or actual arguments. The following example shows both the function template declaration and a template function of **remove** instantiated with a **String** type:

```
<class T>remove(T * array, int sz, const T & val)
remove(String *, int, const String &)
```

The definition of C++ limits what a model can provide concerning templates. The macro nature of templates in the language, with non-lexical scoping of free identifiers, means that some relationships simply cannot be resolved at the point of definition. Given this, our model supports various attributes, such as formal template parameters, and containment relationships in templates, such as those between a parent class template and its member function templates, but largely foregoes more complete analysis of template definitions.<sup>2</sup> On the other hand, when a template is used, i.e., instantiated, all the relationships involving the template class or function are captured.

Note that the kind and subkind fields are useful in determining whether some operations are applicable. For example, a tool that detects unnecessary include files needs only be applied to file entities, while *visibility analysis*, discussed later, can only be applied to classes and their members.

### C. C++ Relationships

There are several possible relationships in C++: inheritance, friendship, containment, instantiation, and reference relationships. We examine each relationship in detail and explain how it affects reachability analysis.

#### C.1 Inheritance Relationship

Figure 2 shows an inheritance structure in the C++ **iostream** library. Note that if an entity refers to **iostream\_withassign**, then it also depends on **iostream**, **ostream**, **istream**, and **ios** for compilation. An inheritance relationship can be *private*, *protected*, or *public*. Also, an inheritance relationship can be *virtual*. Two additional

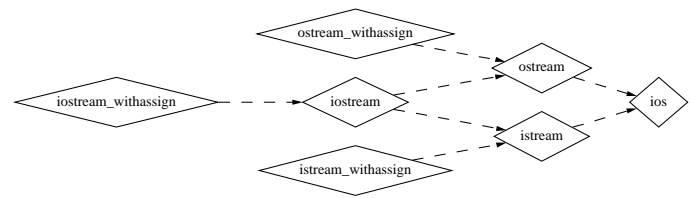


Fig. 2. Class inheritance structure of the **iostream** library

attributes are created to handle these variations: *protection kind* and *virtual flag*. The *protection kind* attribute affects the visibility analysis described later.

#### C.2 Friendship Relationship

There is a friendship relationship from *class A* to *class B* if *class B* declares *class A* as a friend. The relationship direction is set this way because members in *class A* may access members in *class B* and therefore depend on *class B* as far as the direction of reachability analysis is concerned. For example, **ios** depends on **Iostream\_init** in the following piece of code because members of **ios** are allowed to access members of **Iostream\_init**.

```
static class Iostream_init {
    static int      stdstatus ;
    static int      initcount ;
    friend class    ios ;
public:
    Iostream_init() ;
    ~Iostream_init() ;
} iostream_init ;
```

#### C.3 Containment Relationship

There is a containment relationship between every parent class or struct and a member. For example, Figure 3 shows that there are 13 member functions (boxes) and two member variables (ovals) contained in the template class **Array<int>**.

Containment relationships may or may not be walked through depending on the purpose of the reachability analysis. We shall elaborate on this in the next section.

#### C.4 Instantiation Relationship

An instantiation relationship exists if entity *A* is an instance of template *B*. *A* depends on *B* for compilation and linking. For example, in Figure 4, the template class **Node<int>** is an instance of the class template **<class T>Node** and the template function **sort(String \*,int)** is an instance of **<class T>sort(T\* array, int sz)**.

#### C.5 Reference Relationship

Formally, a reference relationship exists between entity *A* and entity *B* if (a) it is not one of the above relationships, and (b) entity *A* refers to entity *B* in its declaration or definition. Again, entity *A* cannot be compiled and linked

<sup>2</sup>We could extend the model to provide entities for the components (e.g., member functions and variables) of a class template, but even this is problematic, as the types associated with these entities may well involve unresolvable type identifiers.

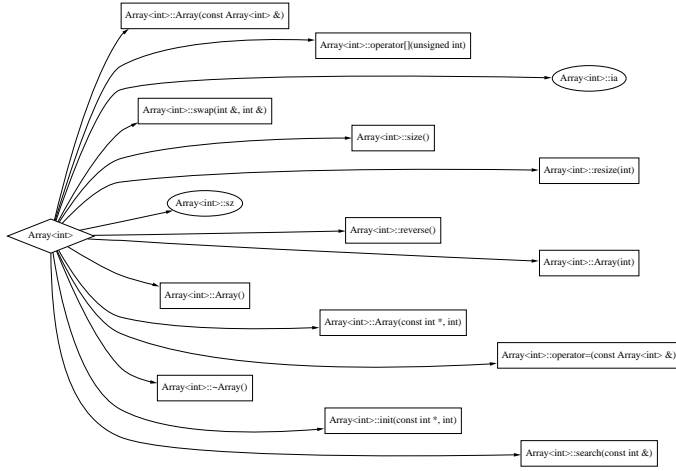


Fig. 3. Containment relationships between `Array<int>` and its members

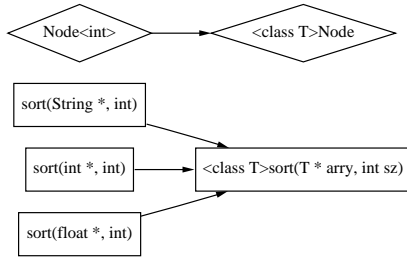


Fig. 4. Template instantiation relationships

without the declaration or definition of entity B. The following example shows several reference relationships:

- `ios::bp` → `streambuf`
- `ios::setstate` → `ios::state`
- `ios::setstate` → `ios::type-230-9::skipping`<sup>3</sup>

```
class ios {
    ...
    enum { skipping=01000, tied=02000 } ;
    streambuf* bp;
    void setstate(int b)
    {
        state |= (b&0377) ;
        ispecial |= b&~skipping ;
    }
    int state;
    ...
};
```

#### D. Generality

Although our model is constructed for C++, we feel much of the model can be carried over to many object-oriented languages. There is broad agreement on what language aspects are basic to object-oriented programming.

<sup>3</sup>`type-230-9` is a surrogate name created for the anonymous enum type.

Languages providing class-based objects typically follow this consensus, which is also captured in our model. For some languages, we need to contract or extend the model, as in Java to capture the distinction between subtyping and inheritance. For other languages, such as Ada 95, the non-object part must be significantly extended. In general, though, we believe our model captures much of what is intrinsic to object-oriented languages. Supporting this belief, our model has been used as the basis for a Java repository.

### III. REACHABILITY ANALYSIS

Reachability analysis can have different slants depending on the purpose of the analysis. We shall discuss several variants and describe how our model and the implementation, *Acacia*, supports each one of them.

Note that the reachability analysis supported by our model is conservative, in that the set of entities returned may be a superset of the minimal closure based on the actual source. This follows from the fact that we are only capturing static information concerning top-level declarations. For example, an analysis of a program at the expression level may indicate that a call to a function only occurs in a branch that is never executed, and hence the function is never called, whereas *Acacia* would report the function as needed. Based on our experience, this appears to be reasonable for typical software engineering tasks.

#### A. Forward Reachability Analysis

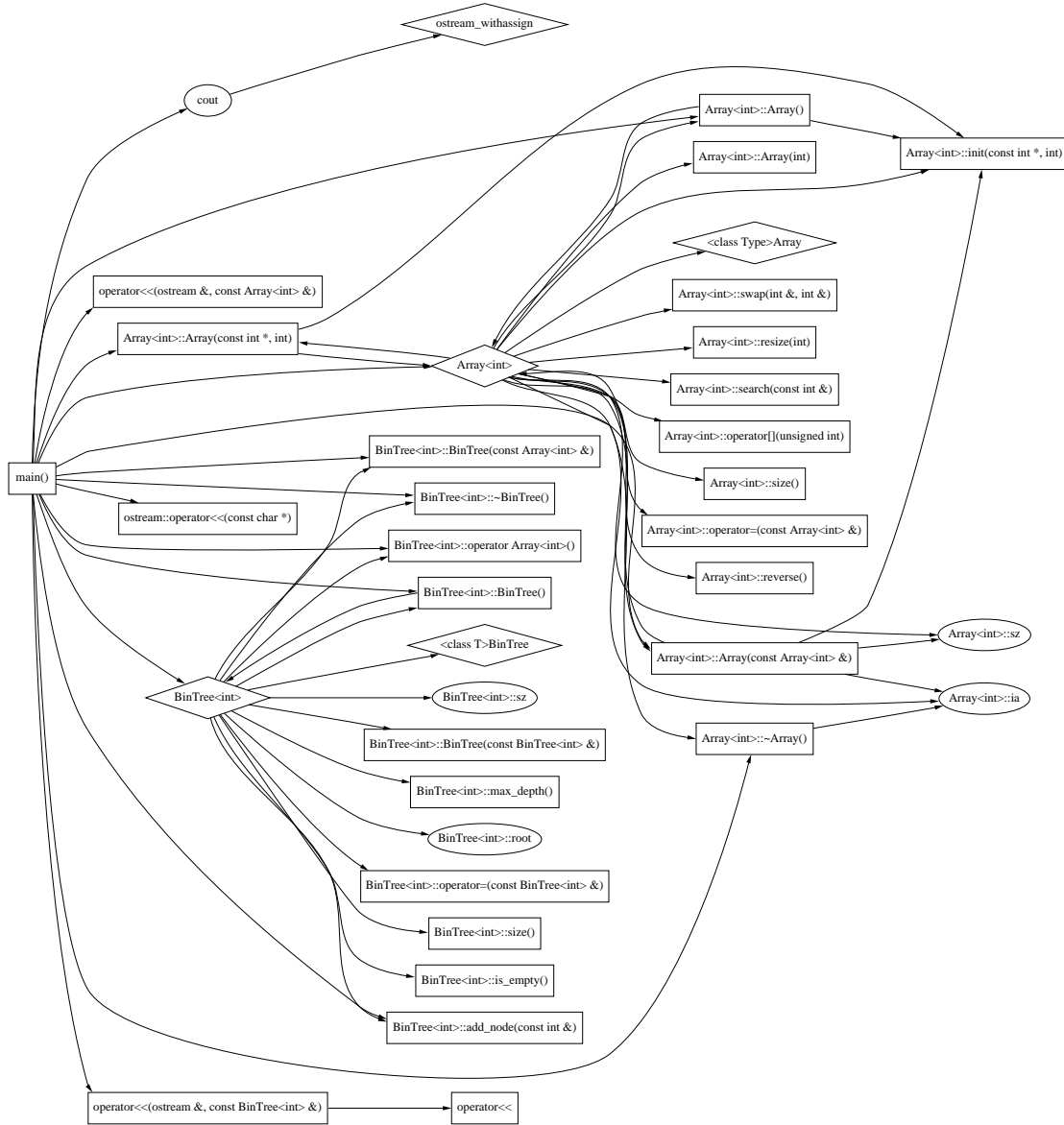
Forward reachability analysis, which computes a *Reachable Entity Set* as defined earlier, is the basis for detecting dead code, packaging reusable software entities, and computing software reuse metrics.

For many tasks, computing the simple transitive closure is sufficient. In some cases, such as software reuse, the task also requires computing certain indirect relations by doing selective reverse reachability computations. For example, class member declarations cannot exist on their own for compilation and therefore we must also capture the *containing* parent declarations. In general, the model explicitly or implicitly contains complete reachability information, so that indirect relations can be generated using appropriate queries over the database.

Figure 5 shows the first three layers of C++ entities transitively reachable from the `main` function in a sample program. The complete closure set includes 175 functions, 17 types, and 72 variables and is too large to fit on a single page.

Note that there is an *instantiation relationship* between `Array<int>` and `<class Type>Array`. Also, while `ostream::operator <<(const char*)` is considered referenced by `main`, the class `ostream` itself is not. We have two choices, depending on the purpose of reachability analysis:

- *software reuse*: In this case, the parent of `ostream::operator <<(const char*)`, `ostream`, is considered referenced. We then include the forward reachability analysis on `ostream` in the closure. While this process may include some class members that are not used in a

Fig. 5. The first three layers of C++ program entities reachable from `main`

particular application, we assume that a class should be treated as a unit.<sup>4</sup>

- *dead code detection*: In this case, even (forward) containment relationships can be excluded if a class is not treated as a unit. This is discussed next.

### B. Dead Code Detection

Many large software projects suffer from a syndrome called *excess baggage* that has one or more of the following symptoms:

- *unnecessary include files*: Many declarations in the header files are never used, but are compiled repeatedly for the source file that includes them.

<sup>4</sup>If space is an issue, the analysis can be modified to only include class member functions that are used in the given program.

- *dead program entities*: Due to program evolution, many program entities usually become obsolete, but programmers either cannot locate them or are afraid to delete anything because they cannot predict the consequences.

To remove excess baggage, we start from the entry points of a program and find the closure set of entities reachable. Containment relationships do not have to be expanded if we want to detect dead member entities for a particular application. This is sometimes critical for applications with strict memory requirements. As described previously in the definition of *dead entity set*, by comparing the closure set against the complete set of program entities in the database, we get a list of unused program entities. Usually, the user is only interested in dead entities in their own code and ignores dead ones in system header files. Our dead code detection tool creates a database of dead

program entities; queries can be used to filter out or focus on particular subsets.

As an example, we applied our analysis tool to a C++ program written by Andrew Koenig that illustrates the concept of dynamic binding [18]. One of the key classes is **Tree**:

```
class Tree {
public:
    Tree(int);
    Tree(char*,Tree);
    Tree(char*,Tree,Tree);
    Tree(const Tree& t){ p = t.p; ++p->use; }
    ~Tree() { if (--p->use == 0) delete p; }
    void operator=(const Tree& t);
private:
    friend class Node;
    friend ostream& operator<<
        (ostream&, const Tree&);
    Node* p;
};
```

We would like to determine if the sample test program (shown below) exercises all member entities in the **Tree** class.

```
main()
{
    Tree t = Tree ("*", Tree("-", 5),
                  Tree("+", 3, 4));
    cout << t << "\n";
    t = Tree ("*", t, t);
    cout << t << "\n";
}
```

While it may not be immediately obvious for some users, this small test program does exercise all member functions of **Tree**, including the destructor, which is called implicitly on the local variable **t** just before the function exits, as well as the constructor **Tree(int)**, which is used as an implicit conversion operator from values of type **int** to objects of type **Tree**. Figure 6 shows the entities reachable from **main** of the test driver in the first three layers, excluding containment relationships.

On the other hand, if we replace the test driver with the following piece of code:

```
main()
{
    Tree t = Tree (5);
    cout << t << "\n";
}
```

then the dead code analysis tool reports that the following three member functions of **Tree** are not exercised by the new test driver, as visualized by Figure 7.

```
Tree::Tree(const Tree &)
Tree::Tree(char *, Tree)
Tree::Tree(char *, Tree, Tree)
```

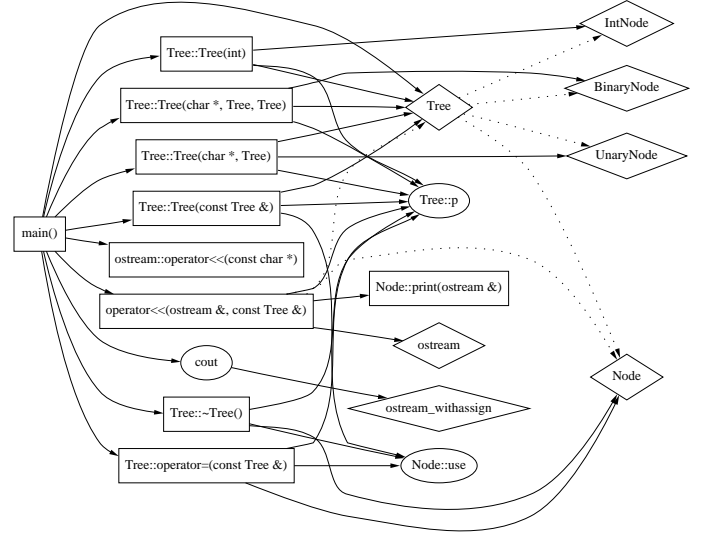


Fig. 6. The test driver exercised all members of the **Tree** class

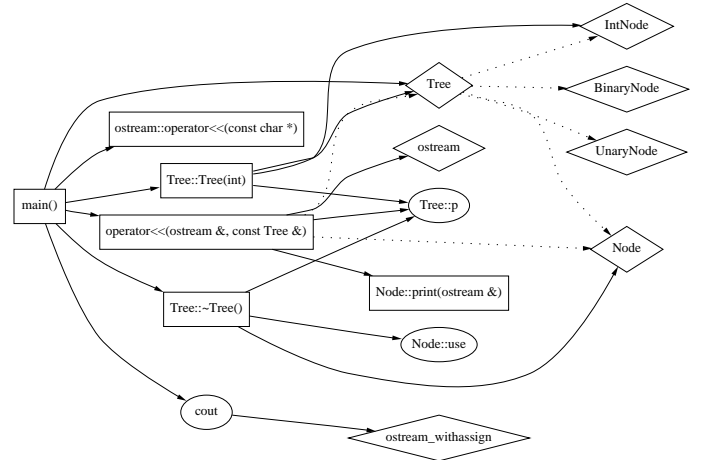


Fig. 7. The other test driver misses a few members of the **Tree** class

### C. Reverse Reachability Analysis

Before software changes are made, it is frequently desirable to find all program entities that potentially can be affected. Reverse reachability analysis allows programmers to find all program entities that depend on an entity directly or transitively through the dependency relationships recorded in the database. For example, Figure 8 shows that if **BinTree::sz** is changed, then all the other entities in the graph potentially can be affected. Such an impact analysis is the basis for selecting regression tests [9] after a change is made in the source code. For example, if the data type of **BinTree::sz** is altered, then all test cases selected for retesting must cover at least the program entities shown in Figure 8.

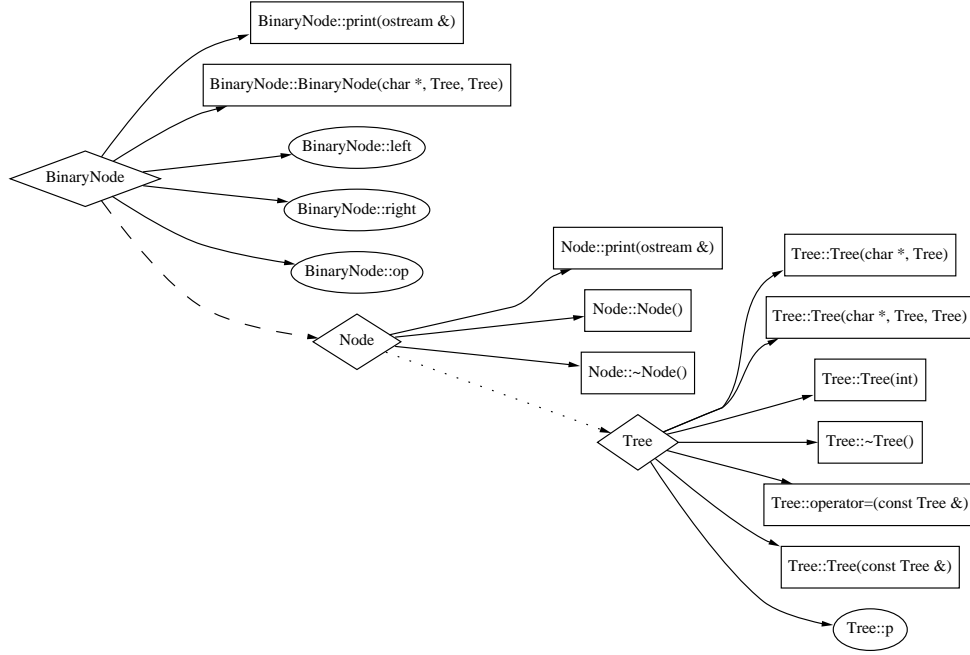


Fig. 9. Visibility analysis of **BinaryNode** and **Node**. In this graph, a dashed edge represents an inheritance relationship, while a dotted edge represents a friendship relationship. The public and protected members of **Node** are visible to members of **BinaryNode** because of the public inheritance relationship. The subgraph in the lower right is visible only to members of **Node**, not **BinaryNode**.

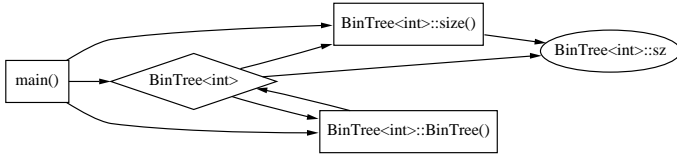


Fig. 8. Impact Analysis: C++ entities that depend on **BinTree::sz** directly or indirectly

#### D. Visibility Analysis

It is frequently necessary to determine what member variables and functions in a class inheritance hierarchy are visible to a derived class. For example, to find all member functions, variables, and types visible to class **BinaryNode** in Koenig's example [18], we can perform a reachability analysis on the containment relationships in the inheritance tree starting from **BinaryNode** and limit the search to members of the proper scope. All members in **BinaryNode** are obviously visible to itself; all *public* and *protected* members from **Node** are also visible because **BinaryNode** has a *public* inheritance relationship with **Node**, but the *private* member variable of **Node** (**Node::use**) is not included. On the other hand, **Node** is a friend of **Tree** and therefore all members of **Tree** are visible to **Node**, but none of them are visible to **BinaryNode** because friendship cannot be inherited. Figure 9 shows the result obtained.

## IV. IMPLEMENTATION

We have implemented a system called Acacia that implements the data model described above. This system consists of a collection of tools for analyzing C++ source, plus an instantiation of the CIAO software visualization system [6] based on our C++ model. Acacia uses **cql** [14] for query and closure computations, and **dot** [15] for automatic graph layouts. In this section, we briefly describe the implementation of the major components in Acacia.

### A. Repository Creation

We built the C++ database extraction tool using the Edison Design Group's (EDG) [12] compiler front end. The front end preprocesses, parses, and typechecks the source, producing a fairly detailed representation, essentially corresponding to a high-level abstract syntax tree, in an intermediate language. In addition to semantic information, this representation also contains declaration, file and source data. This latter detail is crucial for constructing a source level view of the code.

Given the intermediate language representation, the extraction tool traverses the data structure to generate entities for all source items that are not nested within a function scope. After creating the entities, the extraction tool then performs a second pass over the intermediate language representation to generate the required relationships. For C++, this analysis must record implicit uses, such as calls to copy constructors, destructors, and assignment or cast operators, that are not syntactically evident but are generated due to C++ semantics.



The extraction tool produces a repository based on the information within a single source file, analogous to a compiler's producing a single object file corresponding to a single source file. A second tool in the Acacia suite plays a role analogous to the linker, which combines multiple object files into a single executable image. In Acacia, this tool combines the individual repositories into a single repository representing an entire program or subsystem. This integration largely involves replacing references to declarations with references to the corresponding definition, if found.

Since the analysis in Acacia is performed at the source code level, the use of an imported library, without access to its source code, limits the completeness and accuracy of the analysis. Library interfaces, typically specified in source include files, can provide much of the relevant entity and relation information. By the nature of a library, most direct relationships involve external code using something in a library, and not the other way around. When the library does access an entity from a higher level, such as a call-back function, this usually involves the higher level providing a pointer to this entity. Since we cannot analyze how this entity is used within the library, we can only make the conservative assumption that it is used. In the cases when external libraries are involved for which Acacia repositories are not available, definition entities will not exist. In this case, all references to differing declarations of a single entity are replaced by a reference to a single, representative declaration entity.

### B. Instantiation of CIAO for C++

The query and visualization subsystem of Acacia is built by constructing a C++ instance of the CIAO system [6] using an *instance compiler* that takes a specification file for a new language or document type and generates the complete query and visualization environment automatically.

The specification file has five sections:

- *schema*: It maps our data model to the physical *cql* [14] database schema by enumerating the entity and relationship fields. Typically, a field is either an integer or string, but a data type of *entity pointer* allows an entity record to refer to another entity. For example, the *parent type* field of a member entity stores the entity id of its parent class.
- *database view*: This section defines how different entity and relationship records are to be presented as a query result in CIAO's database mode. Each entity kind can have a customized format. For example, the printed name of a member entity is *parent\_name::member\_name*. If an entity is a template instance, the template arguments are attached after the name. Otherwise, the name is just the plain name of the entity. This is important in C++ since it is very common to have members of the same or different classes share a common name due to operator overloading, redefining of inherited methods, etc.
- *source view*: The third section defines what fields are needed to locate the source file and position within that file where an entity appears. A standard CIAO source view tool can use these pointers to output the actual text.
- *graph view*: The fourth section defines how to represent

each entity and relationship when CIAO displays the results of a query as a graph. Entities are represented as nodes of various shapes, colors, and fonts. Relationships are represented as edges of various edge styles and colors.

- *GUI front end*: The fifth section defines the appearance and functionality of the graphical front end. It also defines which queries are appropriate for each kind of entity. For example, in the C instance of CIAO, the query *incl* is marked as only being appropriate for file entities.

Our specification file for C++ consists of only 284 lines. The complete suite of query, visualization, and generic reachability analysis tools in Acacia was generated from this specification file. Only dead code detection and visibility analysis tools require special *cql* [14] query code to handle customized closure computations.

## V. PERFORMANCE AND EXPERIENCE

This section examines speed and storage requirements of our tools and reports our experience in applying them to a few C and C++ software projects.

### A. Storage Requirements and Speed

To evaluate the storage requirements of our databases, we compared the source code and database sizes of three programs. The results are shown in Table I. In general, the size of database and index files together is usually between 1.5 to 2.5 times the size of the source code. On the other hand, we have seen software vendors that create databases in the order of 150 times bigger than the original source.

To see the effect of entity granularity on database size, we compared our database, which is based on top-level declaration granularity, with the intermediate language representation (IR) of EDG, which keeps detailed program information – equivalent to a typical abstract syntax tree. Since the EDG format is only available for each compilation unit, which contains information about a C file with all nested include files, we compare that with the equivalent Acacia repository for each source module. The results are shown in Table II. Note that the EDG IR representation is stored in the binary format, while the Acacia database is stored in ASCII text, and the EDG format is still much larger than the Acacia representation. If we simply compress the Acacia database and use that as our binary format, then the EDG format could use 10 times more space than that of the compressed Acacia format. Of course, in this case, we need to pay the price of uncompressing the database whenever a query is made. On our SGI Challenge L server with four 200 MHZ MIPS R4400 processors, the uncompressing time was less than 0.12 seconds for each of the three modules listed in Table II; so we consider the approach of uncompressing the database before every query a viable alternative for small to medium databases.

To evaluate the speed of our C++ database generator, we compared it to our local C++ compiler, which is also based on EDG's front end. The sample C++ source program consists of 59 C++ source and header files and 9,558 lines of C++ code, with a total size of 239 KB. It took 25.68 CPU seconds (sys time + user time) for the EDG compiler

<i>Program</i>	line count	source size (src)	database and index size (db)	ratio of db/src
<i>A</i>	9,558	239 KB	624 KB	216%
<i>B</i>	15,569	345 KB	524 KB	152%
<i>C</i>	18,182	533 KB	922 KB	173%

TABLE I  
COMPARISON OF SOURCE AND DATABASE SIZES

<i>Source Module</i>	EDG IR size (ir)	Acacia db size (adb)	Acacia compressed db size (cdb)	ratio of ir/adb	ratio of ir/cdb
<i>S1</i>	38.4 KB	22.7 KB	5.04 KB	169%	762%
<i>S2</i>	773 KB	308 KB	69.4 KB	253%	1,114%
<i>S3</i>	958 KB	315 KB	71.5 KB	304%	1,340%

TABLE II  
SIZE COMPARISON OF EDG INTERMEDIATE REPRESENTATION AND ACACIA DATABASE

to compile this program on an 150 MHZ SGI Indy. Our C++ database generator spent 36.71 seconds on the same piece of code. So it took roughly 43% more time.

The storage and speed overhead of Acacia is modest and is acceptable to the software projects that we have been working with inside AT&T.

### B. Experience with Dead Code Detection

To find out the nature of dead code in real projects, we applied the dead code detection algorithm to four projects, two in C and two in C++. Results are shown in Table III. For the two C programs, no containment relationships are traversed and structure members are not considered variables. For the two C++ programs, containment relationships are traversed in determining dead program entities. All entities defined or declared in system headers were not included in the entity counts in this table. A tool called *ciao\_unreach* was applied to detect the dead entities and save the results in a database that allows us to further analyze or filter the dead entity sets.

The lesson we learned immediately was that not all dead entities reported by our tools should be deleted. In project *P1*, many dead types and functions are either from a source module shared with another program or are generated through programs like parser-generators. While they can be deleted (or skipped) to create a smaller executable for *P1*, the programmer may decide to keep them around for ease of maintenance and sharing. On the other hand, most of the dead variables in *P1* are indeed removable because they became unnecessary after several iterations of program modifications. There was only one exception: a variable called *ident*, not referenced by any other program entity, but defined as a string with the program name and organization id so as to embed the owner information in the executable.

In project *P2*, there are only a few monitoring routines that probably should be kept for future use (invoked through different conditional compilation flags). To see

the impact on the code size, we decided to delete all the 37 dead functions. The code was recompiled and linked without any missing references. The source code size was reduced from 17,422 lines to 14,115 lines – 3,307 lines or 19% of the code was removed. Moreover, the executable was reduced from 414 KB to 307 KB – a 25.8% reduction.

Project *P3* is a nice piece of code with very little dead code. All the dead types and variables (including member variables) detected were for error detection or future use. Therefore, none was removed.

Project *P4* is very special in that the C++ software subsystem was merged from two previous and similar projects to create a new subsystem and is expected to have a significant amount of unnecessary code. The complete program database consists of 22,115 entity records and 10,309 relationships. The database and index size are 2.42MB and 1.75MB, respectively, for 193 source files and 71,669 lines of C++ code. More detailed statistics of project *P4* are shown in Table IV. Not counting system headers, the source code itself consists of only 80 source files and 34,360 lines of code, representing 13,636 entity records and 5,534 relationships, as shown in the first row of Table IV. Among the entities, there are 481 types, 209 functions, and 2,668 variables.

Since the new subsystem requires only those entities and relationships reachable from the **main** function, we performed a closure computation on **main** and discovered that only a very small percentage of the entities and relationships are actually needed to carry out the task, as shown in the second row of Table IV. It turns out that only 56 types, 36 functions, and 328 variables are needed – a significantly smaller set of entities and relationships than those of the original merged program. The closure computation took only 5.80 CPU seconds to run on the same SGI Challenge L server.

Without the assistance of automatic analysis tools, it would have been very difficult for programmers to detect all the unnecessary components when two programs are merged to create a new program.

<i>Project</i>	Language	type count	function count	variable count	dead type count	dead function count	dead variable count
<i>P1</i>	C	119	105	87	29	14	5
<i>P2</i>	C	70	229	445	23	37	161
<i>P3</i>	C++	17	74	142	5	0	21
<i>P4</i>	C++	481	209	2,668	425	173	2,340

TABLE III  
COMPARISON OF DEAD ENTITY COUNTS IN FOUR PROJECTS

<i>Version of P4</i>	file count	line count	entity count	relationship count	type count	function count	variable count
<i>Merged</i>	80	34,360	13,636	5,534	481	209	2,668
<i>Closure</i>	21	10,952	1,356	1,970	56	36	328

TABLE IV  
COMPARISON OF THE MERGED PROGRAM AND A SUBSET THAT CONSISTS OF ENTITIES AND RELATIONSHIPS REACHABLE FROM THE **main** FUNCTION

## VI. SUMMARY AND FUTURE WORK

The growing body of C++ code and its language complexity have been presenting challenging maintenance tasks and generating research opportunities in the software engineering community. Reachability analysis is the fundamental building block that supports many complex analysis tasks such as dead code detection, software reuse, and selective regression testing. This paper presents the data model of a C++ software repository and discusses how design decisions made in our C++ model affect variants of reachability analysis. It is crucial that the model be complete at the selected level of abstraction so that the analysis can be performed accurately. Efforts in our implementation of Acacia were greatly reduced due to two factors: the use of EDG's mature compiler front end and CIAO's new instance compiler, which generates the complete query and visualization environment from a small specification file. Due to the entity granularity and tradeoffs we selected, the performance and storage overhead of our implementation is quite acceptable to real software projects. We feel our work and experience on the C++ model can benefit future repository builders for other object-oriented languages such as Java, Eiffel or Ada 95.

As Acacia develops, we have several tasks planned:

- *Incl++*: Besides dead code detection, we would like to create a version of *incl* [29] for C++ to detect unnecessary C++ header files. Not only do we expect the tool to help save a significant amount of compilation time in certain large C++ projects, we also would like it to serve as an additional check on the *completeness* of our C++ model.
- *Automatic program transformation*: While we were able to find a large set of dead entities in *P2* and *P4* (see previous section), we still need tools to automate the code transformation process. Either dead code should be removed or the closure set should be retrieved and packaged in the right sequence for correct compilation.
- *Regeneration of source code*: Currently, the program database does not have sufficient information to regenerate

the source code. As storage costs go down, we expect to begin exploring the possibility of building such a database that would allow us to perform program transformations more easily. This is not just an issue of entity granularity: we need to keep track of all the relationships in a way that is consistent with our current top-level declaration model.

- *Web-based reverse engineering service*: Several CIAO components are now being rewritten to take advantage of new web technologies. Instead of installing CIAO and Acacia on each user's machine, we can create a web service for developers who would like to share the understanding of a particular program but are located in different parts of the world. The users can generate graph views using a Java-based graph applet [2], run database queries through Java-based database connection (JDBC), and view source code using HTML or the emerging XML with embedded hyperlinks for entity references. This would allow us, for example, to analyze the Netscape source code (available through <http://www.mozilla.org>) and to allow programmers all over the world to analyze and visualize the source code without going through the troubles of obtaining both the Netscape source and the CIAO tools. Moreover, with a standard program database format, such services will allow researchers to freely analyze and experiment with public source code and exchange research findings easily.

## VII. AVAILABILITY

Acacia is available for experiments to educational institutions Please visit

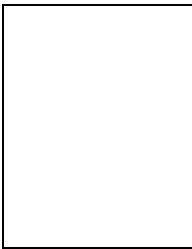
<http://www.research.att.com/sw/tools/Acacia>

for information on how to obtain the package.

## REFERENCES

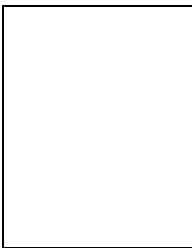
- [1] R. S. Arnold. Software Reengineering: A Quick History. *Commun. ACM*, 37(5):13–14, May 1994.
- [2] N. S. Barghouti, J. Mocenigo, and W. Lee. Grappa: A Graph Package in Java. In *Fifth International Symposium on Graph Drawing*, pages 336–343. Springer-Verlag, Sept. 1997.
- [3] E. Buss, R. D. Mori, W. Gentleman, J. Henshaw, J. Johnson, K. Kontogianis, E. Merlo, H. Müller, J. Mylopoulos, S. Paul,

- A. Prakash, M. Stanley, S. Tilley, J. Troster, and K. Wong. Investigating Reverse Engineering Technologies for the CAS Program Understanding Project. *IBM Systems Journal*, 33(3):477-500, 1994.
- [4] P. P. Chen. The Entity-Relationship Model – Toward a Unified View of Data. *ACM Transactions on Database Systems*, 1(1):9-36, Mar. 1976.
  - [5] Y.-F. Chen. Reverse engineering. In B. Krishnamurthy, editor, *Practical Reusable UNIX Software*, chapter 6, pages 177-208. John Wiley & Sons, New York, 1995.
  - [6] Y.-F. Chen, G. S. Fowler, E. Koutsofios, and R. S. Wallach. Ciao: A Graphical Navigator for Software and Document Repositories. In *International Conference on Software Maintenance*, pages 66-75, 1995.
  - [7] Y.-F. Chen, B. Krishnamurthy, and K.-P. Vo. An Objective Reuse Metric: Model and Methodology. In *Fifth European Software Engineering Conference*, 1995.
  - [8] Y.-F. Chen, M. Nishimoto, and C. V. Ramamoorthy. The C Information Abstraction System. *IEEE Transactions on Software Engineering*, 16(3):325-334, Mar. 1990.
  - [9] Y.-F. Chen, D. Rosenblum, and K.-P. Vo. TestTube: A System for Selective Regression Testing. In *The 16th International Conference on Software Engineering*, pages 211-220, 1994.
  - [10] P. Devanbu. Genoa—a language and front-end independent source code analyzer generator. In *Proceedings of the Fourteenth International Conference on Software Engineering*, pages 307-317, 1992.
  - [11] P. Devanbu, D. Rosenblum, and A. Wolf. Generating Testing and Analysis Tools with Aria. *ACM Trans. Software Engineering and Methodology*, 5(1):42-62, 1996.
  - [12] Edison Design Group. <http://www.edg.com>.
  - [13] G. Fowler. A Case for make. *Software – Practice and Experience*, 20:35-46, June 1990.
  - [14] G. Fowler. cql – A Flat File Database Query Language. In *USENIX Winter 1994 Conference*, pages 11-21, Jan. 1994.
  - [15] E. R. Gansner, E. Koutsofios, S. C. North, and K.-P. Vo. A Technique for Drawing Directed Graphs. *IEEE Transactions on Software Engineering*, pages 214-230, Mar. 1993.
  - [16] J. Grass and Y. F. Chen. The C++ Information Abtractor. In *The Second USENIX C++ Conference*, Apr. 1990.
  - [17] S. Horwitz and T. Teitelbaum. Generating editing environments based on relations and attributes. *ACM Trans. Programming Languages and Systems*, 8(4):577-608, 1986.
  - [18] A. Koenig. An Example of Dynamic Binding in C++. *Journal of Object-Oriented Programming*, 1(3), Aug. 1988.
  - [19] M. Lejter, S. Meyers, and S. P. Reiss. Support for Maintaining Object-Oriented Programs. *IEEE Transactions on Software Engineering*, 18(12):1045-1052, Dec. 1992.
  - [20] H. Müller, M. A. Orgun, S. Tilley, and J. S. Uhl. A Reverse Engineering Approach to Subsystem Structure Identification. *Journal of Software Maintenance*, 5(4):181-204, 1993.
  - [21] G. Murphy, D. Notkin, and E.-C. Lan. An Empirical Study of Static Call Graph Extractors. In *The 18th International Conference on Software Engineering*, pages 90 - 99, 1996.
  - [22] R. Murray. A Statically Typed Abstract Representation for C++ Programs. In *Proceedings of the USENIX C++ Conference*, pages 83-97, Aug. 1992.
  - [23] J. Q. Ning, A. Engberts, and W. Kozaczynski. Automated Support for Legacy Code Understanding. *Commun. ACM*, 37(5):50-57, May 1994.
  - [24] D. Richardson, T. O'Malley, C. Moore, and S. Aha. Developing and Integrating PROTAG in the Arcadia Environment. In *Fifth ACM SIGSOFT Symp. Software Development Environments*, pages 109-119, Dec. 1992.
  - [25] D. Rosenblum and A. Wolf. Representing Semantically Analyzed C++ Code with Reprise. In *USENIX C++ Conference Proceedings*, pages 119-134, Apr. 1991.
  - [26] D. Sharon and R. Bell. Tools that Bind: Creating Integrated Environments. *IEEE Software*, 12(2):76-85, Mar. 1995.
  - [27] I. Thomas. PCTE Interfaces: Supporting Tools in Software-Engineering Environments. *IEEE Software*, 6(6):15-23, Nov. 1989.
  - [28] W. F. Tichy. RCS—a system for version control. *Software – Practice and Experience*, 15(7):637-654, July 1985.
  - [29] K.-P. Vo and Y.-F. Chen. Incl: A Tool to Analyze Include Files. In *Summer 1992 USENIX Conference*, pages 199-208, June 1992.



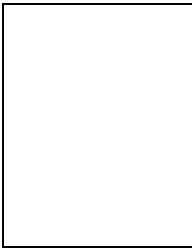
Yih-Farn Robin Chen received the B.S. degree in Electrical Engineering from National Taiwan University, Taiwan, the M.S. degree in Computer Science from University of Wisconsin, Madison, and his Ph.D. degree in Computer Science from University of California, Berkeley. He joined Bell Laboratories in 1987 and was a Member of Technical Staff in the Software Engineering Research Department. He is currently a Technology Consultant in the Network Services Research Center at AT&T Labs.

His research interests include data modeling, reverse engineering, software reuse, database visualization, and tracking changes in software and web repositories.



Emden R. Gansner received a Ph.D. in mathematics from MIT in 1978. After teaching at the University of Illinois, he joined Bell Laboratories in 1980 and was a Distinguished Member of Technical Staff in the Software Engineering Research Department. At present, he is a Technology Consultant in the Information Visualization Research Department at AT&T Labs and an adjunct professor with the Department of Computer Science at Stevens Institute of Technology. His research interests include

graphical user interfaces, programming tools and environments, programming languages, and graph drawing.



Eleftherios Koutsofios received his Diploma in Electrical Engineering from the National Technical University of Athens, Greece, and his M.A. and Ph.D. degrees in Computer Science from Princeton University. He joined AT&T as a member of the technical staff in 1990. Dr. Koutsofios's research interests are in the area of interactive techniques and visualization. He has worked on graph layouts, programmable graphics editors, tools for visualization of large datasets, and program animation.