# LaSSIE:

## A
## Knowledge-Based
## Software
## Information
## System

SOFTWARE

The growing cost of software development, particularly in larger systems, is well documented. Efforts at containing this growth have met with limited success. In his classic paper, "No Silver Bullet" [11], Brooks attempts to explain this; he argues that existing developments in software engineering, such as high-level languages, time-sharing and integrated programming environments have overcome certain *accidental* difficulties in the software development process. However, certain other problems that are the *essence* of the software process remain. He identifies four essential difficulties in building large soft-

ware systems as follows:

- *Complexity:* The complexity of the development process, the application domain, the internals of the system, etc.
- *Conformity:* Software cannot have a regular formal structure (like Mathematics or Physics) because it must adapt to interact with people and institutions according to their diverse needs.
- *Changeability:* In any system, the software component is usually the easiest one to modify.
- *Invisibility:* The structure of software (unlike that of buildings or automobiles) is hidden. The only external evidence we have of software is its behavior when executing.

Brooks argues that there is no single, immediate solution to these difficulties, and thus, no silver bullet to solve the software engineering crisis. In the case of conformity and changeability, we must agree; it appears that software exists to fill market needs, and there is no denying that these needs change. Thus, the extreme malleability of software presents a compelling reason to change the software part of the system to meet these evolving needs. These two difficulties seem inherent in the interactions of software-driven systems with the marketplace.

In this article we focus on the other two essential difficulties— complexity and invisibility. We explore the relationship between

them, and their causes and effects. We also discuss existing approaches to these problems, and the shortfalls thereof. We then present a system called Large Software System Information Environment (LaS-SIE) that incorporates a large knowledge base, a semantic retrieval algorithm based on formal inference, and a powerful user interface incorporating a graphical browser and a natural language parser. LaSSIE is intended to help programmers find useful information about a large software system— in our case the AT&T Definity™ 75/85[1] [3]. Finally, we evaluate the contributions made by LaSSIE to the technology of information systems that deal with invisibility.

## Complexity and Invisibility
Complexity is a common feature of software systems, especially large ones. It has various causes: a large, complex (perhaps even inconsistent) set of requirements, long system lifetimes (due to prohibitive costs of redevelopment), difficulties in interpersonal and interorganizational relationships (see [14] for a comprehensive survey), etc. For a descriptive analysis of the nature of complexity in large software systems, see Perry [35]. For our purposes, it is instructive to consider how invisibility and complexity interact in large systems.

[1]Definity 75/85 is a scalable PBX product with a flexible and powerful feature set. The software controlling this switch is of the order of a million lines of NCSL of C.

In many cases, the *initial* design and architecture of a system is derived after careful analysis of the requirements. A small group of professionals are involved in the beginning and they usually communicate well with each other. This group is generally made up of experienced designers and architects who have had considerable exposure to the application domain; they have a deep understanding of the relationship between the needs of the domain and the design rationale. The initial architecture often structures a system in a virtual machine (VM) style [31]: the system is implemented as a layered set of VMs, each VM implementing reusable primitives for the VMs above it.

These primitives have a simple, well-defined relationship to the needs of the application domain. The layered structure is conceived to promote the simplicity and understandability of the system, and consequently, programmer productivity. In short, the design embodies certain architectural principles that, when carefully followed by developers, keep the system quite simple and elegant, and relative to its size, not all that complex. The design of the AT&T Definity 75/85 is a good example.

Unfortunately, the structure and principles underlying large systems are invisible; there is no clear presentation of the architecture for programmers to examine, and thus to honor, as they extend the system.

Premkumar Devanbu, Ronald J. Brachman, Peter G. Selfridge, and Bruce W. Ballard

In [11], Brooks points out that in civil engineering, geometric abstractions such as scale drawings and stick-figure models capture the essence of the design for everyone involved in a large construction project. He goes on to point out why this is much harder for software systems:

> "As soon as we attempt to diagram software structure, we find it to constitute not one, but several graphs superimposed one upon another.[2] These graphs represent the flow of control, the flow of data, patterns of dependency, time sequence, name-space relationships. . . . In spite of the progress in restricting and simplifying the structures of software, they remain essentially unvisualizable . . ." [11]

Invisibility is thus an inherent property of large software projects; programmers find it difficult to get needed information. Because of this, as the project evolves in response to the marketplace, and grows in size (both people and code) and complexity, the initial architecture, as well as the advantages afforded by it, slowly dissipate. Invisibility has other adverse effects.

### Effects of Invisibility: The Discovery Task

Invisibility complicates the task of acquiring the knowledge—a complex melange of application domain knowledge, architectural and design assumptions, and project-specific development practice—that a programmer needs to do his/her job. This phenomenon has recently been the subject of study at AT&T Bell Laboratories [28] and elsewhere [44]. Modica [28] has called this the *discovery* task. Discovery has impacts on both productivity and quality of the software produced.

Modica has found from surveys

that the discovery task can take anywhere from 30 to 60% of a developer's time, depending on the nature of the project, and the expertise of the person involved. Thus, discovery adversely impacts a developer's productivity. The quality impact of discovery arises from a developer's action upon incorrect or insufficient information. Soloway et al., in [44], document cases (actually in a very small task) where programmers made mistakes when they used documentation that presented information in one way, but did not make these mistakes when it presented information in another form. In two different studies, Soloway, Adelson, and Ehrlich showed how a lack of knowledge affects the quality of the results in design tasks [1] and in small programming-related tasks [43]. Without proper knowledge of the architecture of the system, developers can make mistakes and produce incorrectly functioning code.

In large systems, the impact of poor discovery takes a different form, and is much more insidious. Developers can implement their subsystems in a manner that *violates* the architectural principles of the large system they are modifying. The code may work correctly, and consequently pass system test; however this kind of violation, carried out at various points over a period of time, results in a loss of architecture that gradually erodes the original simplicity of the system. Thus, this lack of knowledge among developers leads to a vicious cycle where the system becomes progressively more complex, and thus harder to know. Invisibility and complexity can be seen to feed off each other.

### Effects of Invisibility: Reduced Reuse

As we described earlier, the initial layered design of large systems is conceived to provide a large number of reusable primitives. Unfortunately, a developer whose task it is to implement, modify or add a special operation to the system

often cannot determine if it has already been done, and if not, whether there is a particular way of doing it that would conform to the initial, intended architecture of the system. Because of this difficulty, instead of reusing existing primitives, programmers reimplement them, which results in lowered productivity.

Reduced reuse can also have a negative impact on quality. Conversations with developers revealed several cases in which programmers, unaware of a virtual machine primitive for an operation, repeatedly reimplemented the same operation—in one case, ten times! When a bug was found in the operation *every single implementation* had to be successively found and fixed. This is a case in which invisibility led to reduced code reuse, and then to increased complexity. Proper knowledge of the system could have prevented this. Thus, we would like to build into our system a library of reusable parts, along with a helpful access mechanism. These points are discussed in more detail in [16].

In this section, we have described the problem of invisibility and its deleterious effects. There have been some attempts to deal with this problem; but before we explore them we must describe an important aspect of invisibility that many systems fail to address.

### Multiple Points of View

The discovery process can be thought of as a series of questions that developers must ask and answer, to gain a proper understanding of the part of the system they are working on. Here are some typical questions, gathered from conversations with developers in the AT&T Definity 75/85 project:

**Q1.** How do I allocate an international toll trunk?

**Q2.** What messages are sent by a process in the network layer when an attendant pushes a button to activate the "Hold" feature?

**Q3.** What C functions enable a Call Forwarding feature activation?

**Q4.** What functions in the Line

---

[2] Brooks suggests that software must be understood from *multiple points of view*.

Manager Process access global variables defined in "/usr/pgs/gp/tgpall/profum.h"?

These queries require different kinds of answers, which depend on knowledge associated with at least four different views of the system:

- A *domain model* view (Q1)—What is the code doing relative to the conceptual objects (trunks) and actions (allocating, releasing, connecting, signaling) in the switching software domain?

- An *architectural* view (Q2)—Q2 makes reference to the network layer, which is a part of the architecture of Definity 75/85.

- A *feature* view (Q3)—How are basic system functions associated with customer features such as "Call Forwarding"?—Q3 combines the feature view with the code view.

- A *code* view (Q4)—How do the code-level components (source files, header files, functions, declarations, etc.) relate to each other?

To maintain the system competently, a developer must be able to explain the structure and behavior of the system; in order to construct such an explanation, the developer has to draw upon at least the four different kinds of knowledge outlined earlier—how does this piece of the system relate to other components, fit in with the architecture, reuse the built-in primitives, and satisfy the customer's need? The pursuit of this sort of knowledge is hindered by the complexity and invisibility of the system's architecture. Field studies by Curtis, Krasner and Iscoe [14] have shown that programmers by and large fail to develop this kind of understanding; they talk about the problem of thin spread of application domain knowledge. They find that a deep understanding of the application domain and its relationship to system architecture is not widespread in the programming workforce. The majority of programmers do not understand the system from

"multiple points of view."[3]

Invisibility can then be thought of as the problem of getting knowledge from the few experts that have it, to the large number of developers who do not.[4] As we shall now see, some attempts have been made to address this information management issue.

## Software Information Systems

There are several systems that provide a variety of information about an existing system to a programmer. They seem to fall into three main categories: *relational code analyzers*, *project management databases*, and *reuse librarians*.

### Relational Code Analyzers

Relational Code Analyzers (RCAs) such as MasterScope [47], CScope [46], and CIA [12] derive certain code-level relationships (such as function-calls-function, function-uses-variable, etc.) directly from source code. This information is generally stored in a relational database. One can query this database using a typical relational language. Questions such as "What functions that call the function flash-button refer to a global variable that is also referred to by the function display-number?" can be asked.

RCAs can answer a number of useful questions for programmers, and are widely used. The information that is retrieved is current, and reflects the actual state of the code base. However, they fail to address certain important aspects of invisibility. First, they simply do string-matching during retrieval; they do not capture the semantics of the strings used in the query. For example, in the previously cited query

---

[3]This should not be confused with the "multiple views" problem in databases where different views are provided for different users. Our goal here is to illustrate the fact that different kinds, or *ontologies* of knowledge, are needed to understand how a large system works.

---

[4]Often-asked questions in the Definity 75/85 project, such as the ones mentioned, can usually be answered by only a few human experts whose time is both valuable and limited.

---

the meaning of flash-button or display-number (or indeed, their relationship to hold-button and display-name respectively) would be opaque to an RCA. Users cannot query RCAs using *descriptions* of components; if one gets a name wrong, there is no recourse.

Secondly, the RCA's query-processing ability is strictly limited to the code view discussed earlier. RCAs cannot answer queries that incorporate other views. Though they capture the syntactic structure of the code, they do not describe the relationship between the structure of the code and the architecture of the system or the domain of application. As Brooks points out, the difficulty of integrating these different views is one reason the structure of large systems remains invisible.

### Project Management Databases

Project management databases (PMDBs) are primarily designed to serve as repositories of all the artifacts generated and used during the life of a software development project, such as documents, programs, test scripts, problem reports, personnel, tools, milestones, and accounting charts. They are intended to provide database support for all lifecycle activities, beginning with requirements specification right through to system test.

PMDB [34] is a typical example of such a system. It aims to collect almost all of the information about a project into one central information system, which can be queried and updated. The information system is based on the entity-relationship model to characterize the artifacts and activities of a software project—with their attributes and interrelationships. Another example is the SODOS system [21] which is intended to "support the definition and manipulation of documents" used in developing software. For this purpose, source code files are considered documents. Features include the management of documentation templates and cross-indexing documents based on

their life cycle interrelationships.[5] There is a query-by-example retrieval mechanism which combines the use of relations such as "implemented-by" with keyword matching.

ALMA [48] is a generic or meta-level kernel that can be used to implement customized database support for a variety of software-engineering environments. It comprises an entity-relationship metamodel that can be instantiated to produce a schema for a specific project database. It has facilities to generate various tools for manipulating lifecycle objects, such as syntax-directed editors, updaters, reporters etc. It includes version management facilities. ALMA can thus be viewed as an applications generator-type tool that can be used to generate a range of different kinds of PMDBs, suited to the needs of particular projects.

None of the existing project management systems, however, attempt to capture knowledge about the application domain of the system, or how the needs of the application domain are manifested in the architecture. The type of information captured in PMDBs is more or less independent of the exact nature of the system under development. They have been used to provide facilities to track the various activities, resources, and artifacts of the software development process. They usually do not provide a description of the structure of the system itself that programmers can use to combat invisibility. Also, the query mechanisms used in these systems are usually quite simple; they do not include any inferential ability. We will address the importance of inference after describing how LaSSIE works.

### Reuse Librarians

Reuse librarians retrieve software components for possible reuse. These systems usually comprise a library of components (object mod-

ules, source modules, specification documents, etc.) and a query/retrieval mechanism. A variety of libraries and retrieval methods are available. The retrieved components can either be directly reused, or adapted as necessary. The retrieval mechanisms in these systems are usually of three types: *keyword, faceted-index*, or *semantic-net*. CATA-LOG [20] is an example of the keyword-based approach typical of standard information retrieval systems. Each library item has an associated set of keywords; the retrieval mechanism takes a set of keywords that the user specifies and matches it with the stored keywords in various ways, and retrieves a matching set. Depending on the choice of keywords used in storage and retrieval, keyword systems may provide adequate performance. However, the semantics of the keywords used in retrieval is unavailable to either the storage or retrieval algorithms of such a system. Because of this, they can neither organize the components in a way that facilitates query reformulation or browsing, nor can they in any way infer the "meaning" of the special set of keywords used in the query. In Definity 75/85, for example, we found that the words "connect," "cut-through," and "terminate" are all used to mean the same thing; more interestingly, combinations of keywords can mean something substantially different in different contexts. For example, the meaning of "termination" is quite different in "process termination" (stop a running process) than it is in "call termination" (connect a call to the designated extension). Given the size of software systems, and the variety of naming styles, it is important for a library of reusable components of a large system to provide semantic retrieval, i.e., retrieval based on meaning.

The classification and retrieval library scheme suggested by Prieto-Diaz and Freeman [37] involves the construction of a domain model and its subsequent use in query formulation/reformulation. In their

paper, Prieto-Diaz and Freeman developed a taxonomic domain model for the set of data operations embodied in a library of software components, categorized along different facets. For example, the facet "Function" can have values such as add, append, create, and evaluate, and the facet "Object" has values such as array, expression, and file. The library is queried by specifying facets and values. It is more amenable to a query-modify retrieval cycle than a pure keyword description. Once designed, however, the classification scheme is static and fixed. Additionally, the representation scheme used here is rather weak; one cannot express constraints (e.g., "password files can *only* be changed by a process with root privileges") in this faceted language.

The abstract data-type (ADT) library proposed by Embley and Woodfield [18] uses aspects of faceted indexing as well as keywords—they propose a software library consisting of a collection of general purpose ADTs, each with a special descriptor. The descriptor includes facets such as "Domain" and "Operations"; it lists several aliases, and descriptive keywords. The user can define explicit relationships between ADTs in the associated descriptors; certain relationships such as close-to, depends-on, or generalizes can also be derived automatically from the values of their facets, and the given keywords. Unfortunately, these derivations are not inferences, in the sense that they are not formally defined, or based on semantics.

The AIRS [30] system, and other semantic-net-based systems such as that of Woods and Somerville [50], and the RLF [41] work at UNISYS, all provide some version of a structural representation of knowledge. In the case of AIRS and RLF, the representational framework is based on KL-ONE [10].

All of the systems described here are intended to be general-purpose software librarians. They are not specifically intended to promote

---

[5]Thus, a feature description in a requirements document might be tied by an *implemented-by* relation to a source code module.

reusability within the framework of a large system; as such, they do not address the issue of invisibility. In addition, none of these systems support a *classification*-style inference. Classification is helpful in dealing with invisibility[6], as we shall argue later.

An important, often-ignored, aspect of reuse libraries is that the index of reusable components of a large software system must reflect the programmer's view of the domain of application of the library. As Curtis argues:

"The effective use of a reusable library will require an indexing scheme similar to the knowledge structures possessed by most programmers working in an application area." [13]

Thus, a library of reusable components within a large system needs an *intelligent indexing scheme* that includes *specific* knowledge about the architecture within which these components are embedded. In addition, Bellin [6] suggests that the description of each component should include constraints on how it should be reused; without this information, the reuser is liable to violate the architecture and contribute to increased complexity of the system. In all of the systems mentioned above, with the exception of [41], the languages can be used to specify mere associations; they are not intended to specify constraints. For example, in these schemes, one cannot specify that update operations must be done *only* by the database manager, or that password files can *only* be changed by a process with root privileges. These constraints are needed to specify architectural principles, as discussed earlier.

The knowledge-based approach of the LaSSIE system is intended to meet these various needs. In a knowledge base, we capture de-

scriptions of reusable objects, along with the architectural constraints embodied in the design of a sizable body of software controlling the AT&T Definity 75/85. Then we use inference to answer the programmer's queries about reusable components. In the following sections, we discuss how the LaSSIE system approaches the invisibility problem.

## An Attack on Invisibility—LaSSIE

The LaSSIE system is an attempt to build a software information system that integrates architectural, conceptual, and code views of a large software system into a knowledge base (KB) for use by developers. The KB is built using a classification-based knowledge representation language, KANDOR [32], to provide semantic retrieval. Besides serving as a repository of information about the system, the KB serves as an intelligent index for reusable components. LaSSIE also provides a user interface with a graphical browser and a natural language query processing system.

The first task we have to undertake is to find a body of knowledge that embodies the architecture and function of a target software system. Since this is a task of formalizing knowledge about a preexisting system, we call this process "reverse knowledge engineering." After the knowledge is identified, we would then try to embed this knowledge in a suitable formal knowledge representation language.

### What Goes in the KB

The first task in building any KB, of course, is to determine the form of knowledge to be represented—one has to identify the *ontology* of the domain. In other words, what will the KB that is produced by reverse knowledge engineering talk about? In the LaSSIE system, we have chosen an ontology that focuses on the actions and objects of the Definity 75/85 architecture. Our choice was influenced by three different factors: existing work in the field of Domain Analysis, em-

pirical studies of programming knowledge, and vocabulary used by programmers in Definity 75/85.[7]

Reverse knowledge engineering is closely related to the notion of Domain Analysis [2, 36, 29]. There are some differences: the domain analyst looks at a variety of systems that service the same application domain and produces an external description of a set of reusable components, which can be used to build applications for that domain. The reverse knowledge engineer, on the other hand, looks at the architecture of a specific large system and produces an internal description of the components of the system, and the architectural framework that unites them.[8] Despite the differences, the results of both kinds of analyses include a *description of a set of components*. Thus the terminology used by domain analysts to describe the results of their efforts is pertinent to a reverse knowledge engineering effort. Both Neighbours [29] and Prieto-Diaz [36] suggest a model based on *objects* and *actions*.

Our KB is concerned with storing knowledge about a large software system, and using this knowledge as an index into a library of reusable components. As Curtis [13] suggests, the way programmers think about the system is an important determinant of the structure of the KB. Shneiderman [40] has conducted empirical studies of programmers, and developed a syntactic/semantic model of programming knowledge. The syntactic part consists mainly of the structural details of the programming language itself. For the semantic part, he proposes a cognitive model based on the action-object paradigm. This model has also been suggested by educators responsible for training new developers on large software projects within AT&T Bell Laboratories. They

---

[6]We defer this discussion until after we have described our classification-based knowledge representation technique and its application LaSSIE.

---

[7]The discussion on ontology is abbreviated here; a longer discussion can be found in [9].

---

[8]The precise role played by the domain model in the knowledge base is described later.

have suggested that a common, unified view of the architecture should be developed and taught to all the developers in the project. Thus, a KB describing the architecture in these terms could serve as a common baseline of information for all developers; such a view has also been advocated by Fischer and Schneider [19].

Descriptions that programmers generate and use also reflect this model. As a normal coding practice, many source files in Definity 75/85 include a short one-line description of the function implemented in that source file. Here are some sample comments:

- This function terminates a call to a user.
- Conference the held user to the current call.
- Update the attendant's status light for the night service feature.
- Cut-through a station user to a call because of answer, originate, or unhold.

Most of these describe a simple action, sometimes caused by other actions, with operands, actors, and contexts. In these, the contexts simply amount to the customer feature that this action supports or plays a part in. In the following section, we describe the representational framework used to build a KB based on the view we just discussed, and then we describe the LaSSIE KB itself.

## KANDOR: the Knowledge Representation Language

Any knowledge-based system needs a language in which the KB is to be expressed. Since the system's knowledge is to be made explicit and available for inspection by the programmer, as well as used for computing inferences, the *knowledge representation language* must be precise and formal—any language with the same type of precision afforded by classical predicate logic is a potential candidate. It is important that the language be expressive enough to allow us to say exactly what needs to be said about a large

software system. It is equally important that the computation of inference with the language be tractable; thus the language will need to be restricted in some way to avoid theorem-improving that takes exponential or infinite time.

There are a number of currently popular styles of knowledge representation language. Standard logical languages are good in situations in which knowledge is incomplete, and in which one needs to express mainly factual assertions. Production systems are useful when the knowledge is heuristic and can best be expressed in a condition/action style (e.g., in medical diagnosis). Frame systems are best when we want to describe sets of objects with complex relational structure, and when the domain exhibits strongly hierarchical, taxonomic categories of objects.

Given our emphasis on retrieving descriptions of parts of a large software system, and the centrality of the object/action paradigm, we have chosen a frame-based representation for the LaSSIE knowledge base. In this style of representation, frames represent classes of objects; they do so in an object-centered way, by clustering all information about a given object in one place. In KANDOR [32], our language, a frame is considered a complex description, which expresses constraints on members of the class that it denotes. A frame's definition is very much like a standard dictionary definition, starting with a set of super-frames (more general classes) and then adding restrictions to those parent frames, thereby creating a more specialized class. It should be noted that frames formed in this way induce a *taxonomy* of classes, with the most general class **THING** at the top, and the most specific classes at the bottom.

The restrictions in frame definitions are usually specified in terms of slots, which are two-place relations that describe the attributes of class members. Restrictions can be formed by limiting the type of slot-filler to be expected (e.g., "a direc-

tory *all* of whose files are header-files") or by specifying the maximum and minimum number of fillers expected (e.g., "a function with *exactly two* arguments"). In order to be considered a member of the class defined by a frame, an individual object must provably be describable by all of the super-frames and must provably satisfy each of the restrictions.

In KANDOR, *individuals* are the objects that represent individual items in the domain. Individuals have their own identities, regardless of their descriptions, but, as implied, an individual will be described by a set of frames. In addition to the general description of an object by its set of parent frames, the representation of an individual is completed by expressing the relationships between it and other individuals (this is done using instantiated versions of the slots). Thus, a particular file might be considered to be describable by the frames C-FILE, SOURCE-FILE, and CALL-PROCESSING-FILE, and might have the **name** slot filled in by **Capro.c**, and the **directory** slot filled in by /Usr/Callpro. These two items themselves would be KANDOR individuals with their own descriptions.

Frame systems like KANDOR perform two key inferences that allow them to be used in retrieval applications like LaSSIE. First, and most commonly, there is *inheritance of properties*, which allows a subframe to inherit all properties of its superframes (this is a simple consequence of the sub/superframe relationship representing a universally quantified conditional). Typically inheritance is considered to contribute efficiency of representation, since one need only represent in a single place a general property that applies to many classes. The second key inference—performed by only a small number of systems like KANDOR, which interpret frames as descriptions (rather than as sentences)—is *classification*, in which, for any individual, all descriptions that can provably apply to that indi-

vidual are found. If KANDOR descriptions can be interpreted as representing necessary and sufficient conditions for class membership, then inheritance can be considered to use the necessary conditions in a "forward" direction, and classification can be considered to use the sufficient conditions in a "backward" direction. Classification can also be applied to frames, wherein all appropriate super- and subframes for a frame are determined automatically.

### The KB—Actions, Objects, and Doers

Figure 1 shows the root of the LaSSIE taxonomy (THING) and the two levels immediately below it. The four principal object types of concern in our domain are ACTION, OBJECT, DOER, and STATE. The edges of the taxonomy have their usual "IS-A" interpretation.[9] DOER represents those THINGs in the system that are capable of performing actions. Nodes below DOER and OBJECT represent the architectural components of the system, i.e., its hardware and software components. Nodes below ACTION represent the system's functional component, i.e., the operations that are performed on or by the system. The relationship between the two components is captured by various slot-filler relationships between ACTIONs, OBJECTs and DOERs. Each action description combines the top-level concepts in various ways using KANDOR descriptions. Here is a typical one:

```
1(define concept
2   USER-CONNECT-ACTION
3   NETWORK-ACTION
4   CALL-CONTROL-ACTION
5   defined
6   (exists has-actor
7       (generic PROCESS))
8   (exists has-agent
9       (value Bus-Controller))
10  (all    has-operand    (generic
```

[9]In particular, a TRUNK IS-A COMMUNI-CATIONS-DEVICE, a RESOURCE-OBJECT, and a DOER.

```
   USER))
11  (exists has-environment
12      (generic CALL-STATE))
13  (exists has-result
14      (value Talking-State)))
```

In this example, words in italics represent reserved words in the formal representation language. These perform logical functions like stating that a certain individual plays a certain role with respect to a general class of items. For example, (exists has-agent (value Bus-Controller)) says that the bus controller process is always an agent for every USER-CONNECT-ACTION. Among the other items, those in all capitals are frames; they represent general concepts like users and call states. Items with initial capitals are individuals; they stand directly for individual objects in the domain, such as the bus controller process. Items in all lower case are roles; they represent relationships between individuals. Thus, the interpretation of this frame is as follows: USER-CONNECT-ACTION [Line 2] is by definition [5] a network action [3] and a call-control-action [4] that is performed by a process [6, 7], using the bus-controller [8, 9] on a user [10], and which takes the user from
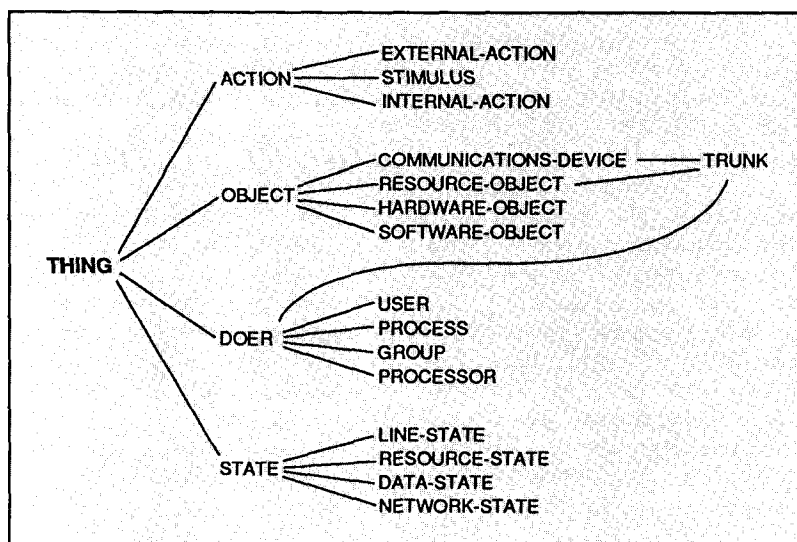
some call state [11, 12] to the talking state [13, 14]. LaSSIE's KB contains 102 action concept descriptions of this type, which are classified into a conceptual hierarchy. Further down in the hierarchy, the action concepts become very specific.

The most specific action types, which each correspond to a particular function/source file,[10] are coded as individuals. For example,

```
1(define individual
2   Add-User-Action
3   (ACTION)
4   (has-actor
5       Call-Control-Process)
6   (has-agent Bus-Controller)
7   (has-operand Generic-User)
8   (has-recipient Generic-Call)
9   (has-environment
10      Generic-Call-State)
11  (has-result Talking-State)
12  (implemented-by
13  /Usr/Pgs/Gp/Tgpall/Profum.c)
14  (calls-function
15      Signal-Add-Error-Action)
16  (accesses-variable
17  *Call-Record*))
```

In other words, Add-User-Action [2] is an action [3] that is performed

[10]Definity 75/85 has only one function per source file.



**Top Three Levels of the LaSSIE KB**

by the call control process [4, 5] using the bus-controller [6]; its operand is any user [7], and its recipient is a call [8]; it takes its operand from any call state [9, 10] to the talking state [11]; it is implemented by the source file /Usr/Pgs/Gp/Tgpall/Profum.c [12, 13], calls the function Signal-Add-Error-Action [14, 15], and uses the global variable *Call-Record* [16, 17]. It should also be noted here that the KANDOR classification algorithm will ensure that this individual gets classified under the frame USER-CONNECT-ACTION. It is this kind of classification that organizes the large number of frames and individuals in LaSSIE into a usable form.

Concepts in the KB, such as these examples and USER, RESOURCE-OBJECT, GROUP, etc., are concepts that are specific to Definity 75/85 and have a well-defined meaning within the architecture. They are referred to repeatedly in the design and specification documents. We show further details of the portion of the laSSIE KB under the OB-JECT concept in Figure 2. Notice that some concepts appear in **bold-face.** These indicate that a concept has more than one parent in the taxonomy (i.e., the taxonomy is a tangled hierarchy); For instance, PROCESS is both a RESOURCE-OBJECT and a SOFTWARE-OBJECT.

The operations that can be performed on these objects are also well defined: they constitute some of the reusable primitives of the architecture. The descriptions of these operations, the way in which they are implemented in the system, and the conditions under which they are performed constitute some of the architectural principles that promote simplicity and understandability of the system. The domain model for Definity 75/85 plays a crucial role in structuring this knowledge.

### The Role of the Domain Model

Definity 75/85 and other large software systems are often designed in

a layered virtual machine style, each layer providing reusable primitives for the layers above. In Definity 75/85, the key ingredients of these layers are communicating independent processes; the reusable primitives are the functions that can be performed when messages are sent to these processes. Examples include "Connect an ISDN trunk to this call," "Translate the dialed digits," and "Display a number on the Attendant's console." These key primitive building blocks are combined and sequenced in various ways to implement the customer features; the descriptions of these operations, with their operands and precise effects constitute the "domain model" (in the spirit of Prieto-Diaz [36]) of the PBX switching features supported by the Definity 75/85.

Though there was no explicit discussion of the domain model anywhere in the documentation, a detailed study of the architecture indicates that the designers carefully considered the feature set and isolated a set of primitives that could be used to implement them; the architecture is designed to implement these primitives in a reusable manner. The architecture also includes some constraints (strict layering, deputation of certain tasks to certain processes, etc.) that govern the manner in which the primitives in the domain model are embedded in the architecture. Definity 75/85 and related products are typically described at the level of the process/message architecture, or at the level of the customer services performed by the switch.[11] The domain model is only implicit in the system, and while it may constitute an important part of the programmer's understanding of the system, it is left for each programmer to divine in his/her own way.

It is our view that in an integrated software information system, this key level of description

---

[11]For example, "plain old telephone service" (POTS), call forwarding, call waiting, conference calls—these services are usually referred to as "features."

should be made explicit, as should its connections to the features built from its primitives and its connections to the architecture that implements it. Thus, in LaSSIE the domain model serves a key integrative function. It can also help the programmer access the rest of the knowledge base, since it maps most closely onto the common conception of the overall function being performed by the switch. Queries about various functions of the switch—especially by programmers newer to the project—will naturally be formed in terms of actions like connecting two trunks rather than in terms of certain messages received by the "connection manager" (these latter types of queries will also need to be answered, of course). In LaSSIE a concept like TRUNK-CONNECT-ACTION denotes a domain model primitive, i.e., a common, important operation that can be reused in different features. In the Definity 75/85 architecture, it is executed by the Timeslot-Controller, which is a process in the Network-Layer. The definition of this concept will include all this information, as well as a pointer to the code object (subroutine or function) implementing this operation, and a list of customer features that use this operation.

The domain model and its relationship to the architecture of the system is built manually into the LaSSIE KB. In addition to this, we also capture the code view of the system, which is, fortunately, mostly automatically derived, using CScope [46].

### Integrating Code-level Information

We enhanced the conceptual action-object-actor level of representation with certain simple but important code-level relationships such as function-*calls*-function, sourcefile-*includes*-headerfile, etc., which are syntactically derived by scanning the code. Thus questions such as "what functions call 'apost' and include 'errproc.h'?" can be processed.
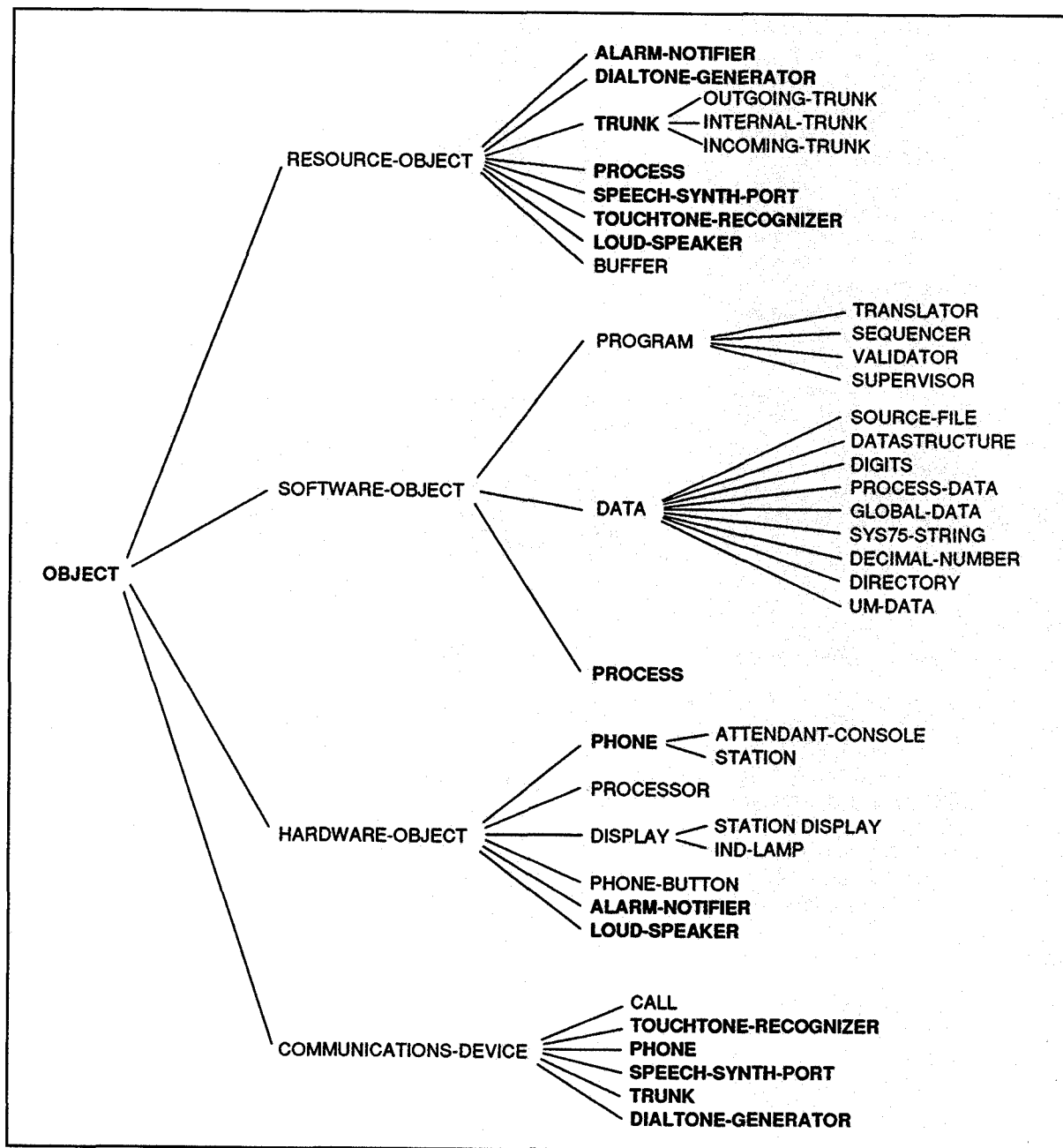
To accomplish this, we developed a general representation of code objects and their interrelationships and added them, codified in KANDOR, into the knowledge base. The instances of these concepts were derived as follows: first the source files were scanned with CScope and several relations gener-

ated. Then these relations were organized around specific instances of functions and source files. These collections were translated directly into KANDOR function calls to add the appropriate information to existing information from the manually generated KB concepts. Thus the syntactically derived code-level

information is smoothly incorporated with the conceptual and architectural views, so one can pose queries such as "What global variables are accessed by a function that flashes a display lamp at an attendant's console?"

A frame-based description of a function can combine the architec-



The LaSSIE Taxonomy under **OBJECT**

tural/domain aspects as well as its code-level relationships into one unit that can be manipulated by the classification algorithm. This facility for combining different kinds of information about an object provides a way to integrate the code, conceptual and architectural views of the system.

### Querying and Browsing

Most of the LaSSIE KB describes actions; the paradigmatic usage is a query from a user that describes an operation that she wants to see examples of, to understand in more detail, or to reuse; the LaSSIE system will retrieve instances. LaSSIE also has an interactive graphical user interface that includes a modified version of the ARGON [33] system, and the ISI-Grapher [38]. In this section, we present examples that illustrate how the retrieval system works. The core of the retrieval system is the classification algorithm; it simplifies the task of querying a large KB.

A LaSSIE query[12] is simply a description of an action, such as the following:

```
CONNECT-ACTION
    has-actor    DOER
    has-operand OBJECT
    has-cause    ACTION
                 has-actor USER
```

This query will retrieve all individual CONNECT-ACTIONs performed by a DOER on an OBJECT, because of an ACTION by a USER. The query processing is carried out in two stages. First the query is placed in the LaSSIE taxonomy by the classification algorithm, using the description in the query and descriptions of the frames in the taxonomy; then, the matching instances are the instances of those frames that are subsumed by the classified query. The user who generated the query does not have to know that PROCESSes, USERs and TRUNKs are DOERs (see Figure 1);

neither does she need to know how many processes there are, or what their names are. The classification algorithm takes care of all of that. Thus, it is possible for a user to retrieve answers to a query without knowing the structure of the taxonomy of a KB, or the individuals therein, using the semantic retrieval provided by the classification algorithm. For example, the following individual would be retrieved by this query:

```
1 (define individual
2     Attd-Merge-Call-Action
3     (CALL-MERGE-ACTION
4     ATTD-CAUSED-ACTION)
5     (has-actor
6         Attd-Monitor-Process)
7     (has-agent
8     System-Fabric-Mgr-Process)
9     (has-operand Generic-Call)
10    (has-recipient Generic-Call)
11    (has-environment
12        Generic-Call-State)
13    (has-result Talking-State)
14    (has-cause Attd-Button-Push)
15    (implemented-by
16    /Src/Ugp/Attd/Atd-Au-Im.c)
17    (calls-function
18        Signal-Error-Action)
19    (accesses-variable
20        *Call-Record*))
```

The advantages of semantic retrieval are illustrated in this example.[13] Though the action is called a MERGE-ACTION, the classifier recognizes it as a CONNECT-ACTION based on the descriptions of each. Additionally, from the description of Attd-Button-Push (not shown here) it realizes that this is an ACTION by an ATTENDANT, which is defined to be a specialization of USER.

The ARGON interface displays all the retrieved individuals, and the user can select any one of them

for detailed display, with all its slots and fillers; each of these fillers and slots can, in turn, be selected, and further information about them is available. Thus, the user can find out that this action is performed by the Attd-Monitor-Process (selecting that will retrieve its description, which tells the user that this process is in the Service-Layer of the architecture), and that it uses the System-Fabric-Mgr-Process as an agent. Furthermore, if one selects the concept ATTD-CAUSED-ACTION (which denotes any action by the system caused by an action by an attendant), one would find that it is always performed by the Attd-Monitor-Process, i.e., this process, and only this process, responds to actions by the attendant.[14]

Thus classification makes it possible that a merge action caused by a button push by an attendant gets retrieved by a query that asks for a connect action caused by an action by a user. Furthermore, the retrieved individual carries with it a wealth of information about the architectural, conceptual, and code aspects of the operation. Thus, retrieval based on the semantics, along with integration of multiple views helps in the discovery of useful information about Definity 75/85.

In the event that the query retrieves too many, or too few relevant instances, the user can reformulate it, taking advantage of the information in the taxonomy. For example, assume our query retrieved too many instances—there may be too many DOERs doing too many things to too many OBJECTs. The user can specialize the query by replacing either OBJECT or DOER with its children in the taxonomy, either by using a menu choice, or by looking it up in the graph. This is the advantage gained by having an intelligent index. Thus, if there are no (or an insuffi-

---

[12]The actual query language in ARGON [33], the system that LaSSIE is based on, includes negation, disjunction, etc. . . Here we illustrate just a simple query.

[13]The retrieved individual is an instance of a Call-Merge-Action [3] done by the Attd Monitor process [5, 6], using the System-Fabric-Mgr [7, 8] process, which connects one CALL [9] to another [10], leaving both calls in a Talking-State [13]. This action is caused by the attendant pushing a button [14]; it is implemented by the source file listed in line 10; it calls the function Signal-Error-Action [18, 19] and uses the global variable *Call-Record* [20, 21].

[14]This would be specified by the inclusion of the restriction (all has-actor (value Attd-Monitor-Process)) in the concept definition of ATTD-CAUSED-ACTION.

cient number of) answers, the user can generalize parts of the query.

In addition to its graphical interface, LaSSIE also has a natural language query interface based on TELI [4], which we now describe.

### Adding a Natural Language Interface

To provide a natural language interface for LaSSIE, we customized the TELI system, which maintains data structures for each of several types of knowledge [4]. This information includes (1) a taxonomy of the domain, which enables the parser to perform several types of disambiguation; (2) a lexicon, which lists each word known to the system, along with information about it; and (3) a list of compatibility tuples, which indicate plausible associations among objects and thus reflect the semantics of the domain at hand. For example, an agent can perform an action on a resource, but actions cannot be performed on agents, resources cannot perform actions, etc.

In LaSSIE, KANDOR individuals generally correspond to proper nouns (i.e., names), while a frame may correspond to either a verb or a common noun. Generally, frames under ACTION correspond to verbs describing actions, while nodes under OBJECT or DOER correspond to nouns. For example, the frame ALLOCATE-ACTION maps to "allocate," "reserve," and "grab," and PROCESS maps to the noun "process." Individuals are usually associated with one or more proper nouns in an obvious way. For example, the individual process BUS-CONTROLLER is named "bus controller."

As explained, action frames include slot restrictions corresponding to case roles including the actor, the operand, the recipient, the cause of the action, etc. To each of these, there naturally correspond one or more English prepositions. Thus, each slot associated with an action frame gives rise to compatibility tuples as described above. For example, consider the frame defi-

nition[15] shown below, with its associated verb *connect:*

```
1(verbframe
2CALL-CONNECT-TRUNK-
   ACTION
3   (connect)
4   (ACTION)
5   (exists has-operand
6         (generic TRUNK))
7   exists has-recipient
8         (generic CALL))
9   (exists
10   has-actor
11   (value
12      CALL-CONTROL-PROCESS)))
```

For this frame and its slots, the following compatibility tuples are generated:

⟨CALL-CONTROL-PROCESS
   connect TRUNK⟩
⟨CALL-CONTROL-PROCESS
   connect to CALL⟩

The annotation of the KB was done manually, after which the conversion to the TELI data structures was automatic. The resulting compatibility tuples for LaSSIE include 167 verb case frames, corresponding to a total of 40 verbs. The lexicon contains 882 entries, including 193 common nouns and 260 proper nouns.

To process a query such as "What actions by the bus controller are caused by an action of an attendant?," TELI parses the input, making intimate use of the compatibility tuples and the taxonomy to ensure globally consistent case bindings. The final parse tree is then converted into a semantic structure resembling a first-order logical form, which is sent to a LaSSIE-specific filter to strip out quantifiers associated with words such as "a" and "the." The resulting structure is then passed back to LaSSIE for translation into a query that is executed (thus performing a retrieval); this query can also be directly edited if reformulation is necessary.

For example, TELI's output for

---

[15]Actually, the form shown generates a table entry for the action, associating it with a verb name, and generates a standard function call to define a KANDOR frame.

the above query is:

```
(set A1 (ACTION A1)
   ((ACTION BY AGENT)
    A1 Bus-Controller)
   ((ACTION CAUSE ACTION)
    A2 A1)
   ((ACTION BY AGENT) A2 P1)
   (ATTENDANT P1))
```

This is then translated into the following editable ARGON query:

```
ACTION
   has-actor  Bus-Controller
   has-cause ACTION
            has-actor
               ATTENDANT
```

It should be noted that the user of LaSSIE need not know the details of the underlying KB in order to pose questions in English but, by seeing the associated ARGON query, may well learn something about the KB when the input is processed. For example, the query "What actions by a process reserve a touch tone recognizer because of a pickup by a user?," would be translated to

```
ALLOCATE-ACTION
   has-actor  PROCESS
   has-operand
      TOUCH-TONE-RECOGNIZER
   has-cause
      OFF-HOOK-ACTION
         has-actor  USER
```

In this case, the user would learn that the action verb "reserve" corresponds to ALLOCATE-ACTION, "pickup" to OFF-HOOK-ACTION, and also that actors of and causes of ACTIONs respectively are specified by using the HAS-ACTOR and HAS-CAUSE slots.

### Why Frame-Based KBs with Classification are Important

LaSSIE accrues several key advantages by using a frame system like KANDOR, which incorporates classification as its foundation. These include the following:

● *Aggregation of information about individuals:* The object-centered approach of a frame language allows all information about an in-

dividual to be concentrated in one place. This allows us to integrate information obtained from different sources and different points of view into one description.

• *Semantic retrieval:* Classification helps during query processing by reducing the amount of information the user has to know about the knowledge base. In purely syntactic information systems, where the terms used have no descriptions (e.g., relational databases or keyword systems), the user has to know the exact terms that are used in the database. Classification reduces the number of terms that the user needs to know. For example, when searching for a general category of items (e.g., FILE), all subcategories (as expressed in the generalization hierarchy induced by classification of frames; e.g., HEADER-FILE) can also be retrieved. This point has been made by Beck, Gala and Navathe (and others):

". . . the user can describe a query in terms which may be different from the exact terms under which the desired information is stored, as long as the meaning is similar." [5]

Essentially, classification is a logically defined operation that allows semantic retrieval in a simple yet principled manner.

• *Use of classification and inheritance to support updates:* When new information is added to the KB, the system automatically classifies the new items, thereby giving the developer an integrity check. Since all categories implied by the new description are found, the developer can see if there were any omissions or accidental inferences to be made. Inheritance also makes the addition of new information easier, since many of the properties of a new item can be derived from its membership in existing classes.

• *Use of the KB as an index:* In a way, the frame KB acts as a general schema over the particular code data stored in the system. The general knowledge encoded in the frames can be used to guide the user in his/her search for particular items of interest. LaSSIE has several means for using the frame KB in this way (including a graphical display of the frame hierarchy).

We can now compare our work to the other semantic-net reuse library systems listed earlier. RLF does not have a built-in classifier; the KNET language is used to describe the components, and the retrieval/browsing algorithm is left to the application developer. The AIRS [30] system uses a heuristic retrieval algorithm based on a numerical conceptual-distance measure, which the user has to specify. Woods and Somerville use conceptual dependency (CD) diagrams. Their query mechanism is based on a set of *verbs*. They use word associations specified by the user to identify the conceptual dependency graph for a given verb. The user is then prompted for further information based on the selected conceptual dependency graph to further narrow the search for components. For example, the verb *send* might be associated with the CD graph named *communicate*. Then the structure of the communicate concept would be displayed on the terminal, and the user would be prompted to fill in keywords for objects associated with this structure; these would then be matched with stored descriptions to find close matches. The algorithm is implemented in Prolog; the retrieved components are listed in order of closest match, depending on how many keywords matched.

None of these systems has a well-defined notion of classification; this makes them ill-suited for our application, as we now explain. Building a KB of the size required for LaSSIE (about 200 frame descriptions) involves constructing a large number of KANDOR descriptions (several example concept descriptions can be found above). Once a description is input to KANDOR, the classification algorithm kicks in and places the concept in the proper relation to other descriptions already in the KB. Without classification, the knowledge engineer would have to figure out the proper placement by hand—a daunting task with such a large KB. Thus, classification plays a crucial role in assembling the taxonomic KB that provides intelligent indexing, as called for by Curtis [13]. Constructing a very large KB in the KL-ONE-style language used in AIRS and RLF would be easier with a classifier. Unfortunately, because of the power of KL-ONE [39], a complete classification algorithm is impossible. The CD diagrams used in Woods and Somerville do not have a taxonomic organization; the matching is done by keyword associations. Thus in this case, in addition to constructing the CD diagrams, it is necessary for the knowledge engineer to specify keyword associations. Likewise, with AIRS, in addition to describing the components in a KL-ONE-like language, the knowledge engineer has to specify a numerical measure of conceptual distance. The simplicity of the KANDOR language, and the built-in classifier, simplifies the knowledge engineer's task, while continuing to provide semantic retrieval as well as a fairly rich knowledge structure for browsing, navigation, and query reformulation.

## Limitations of LaSSIE

LaSSIE is primarily intended to process queries about actions. Its KB consists of a large collection of action descriptions in KANDOR, and the user can query this collection by specifying an action description that she is interested in; the classification algorithm retrieves the instances thereof.

The strengths of LaSSIE are discussed in the previous section. LaSSIE's essential limitations are those caused by the limitations of KAN-

DOR. These limitations can be divided into two classes:

The first class of limitations is rooted in the fact that KANDOR is a domain-independent language, not specifically designed to represent knowledge about real-time software in terms of objects and actions. So there are various aspects of each that cannot be expressed adequately within its representational framework. With respect to ACTIONs, one must express all aspects using the same undifferentiated slot mechanism. This means that special inferences based on particular slots such as has-cause,[16] has-result and has-environment,[17] or has-service[18] cannot be performed. With respect to OBJECTs, it would be useful to represent part-of hierarchies. For example, a light is a part-of a button which in turn is a part-of a telephone. KANDOR does not support reasoning based on meronymic (part-of) hierarchies. Meronymic reasoning is also important for dealing with plan-like groups of ACTIONs.

The second class of limitations is due to certain expressive limitations made in KANDOR to make the classification algorithm faster and easier to implement. For example, one cannot specify a relationship between the fillers of different slots in a concept. Thus, in the concept describing the "call waiting" feature, one might want to specify that the person being called (and currently off-hook) is the same person who receives the call-waiting signal. KANDOR is too weak to specify this precisely.

## Further Work
The action-based representation used in LaSSIE has one crucial failure; the *contexts* within which operations are performed are not well

---

[16]e.g., transitive closure over causality

[17]e.g., relating the result postcondition of one action to the environment precondition of another

[18]in regard to establishing the context or rationale for an action

---

established. Thus, the following kinds of questions cannot be answered:

- Why is this action being performed?
- How can this operation be performed?
- Can this operation be done in that context?
- Is this operation involved in more than one customer feature?

Anyone familiar with the planning literature in AI will recognize that all of these questions relate to plan-like information. This is a considerable body of literature linking plan knowledge to the cognitive processes involved in software design/development. Soloway, et al. [42–44] have established the importance of plans in programming knowledge; Letovsky [26, 25] has shown that plans are important for program comprehension and inspection. To answer "Why" questions,[19] programmers give answers that involve describing a plan that contains the operation; "How" questions[20] are answered by describing a plan that satisfies the goals of the given operation.

How can a LaSSIE-like KB be enhanced by including plan knowledge? With more information about the context of an action within a plan that implements either a higher-level goal, or a customer feature, more questions could be answered. Thus, one can imagine a KB that contains a large number of plans, describing various operations, features, etc., of Definity 75/85. These plans would each be composed of a series of actions. Inference algorithms would run over these plan data structures that would answer questions of the sort posed at the beginning of this section; others may assist a user to insert a new plan into the KB and properly classify it in relation to previously stored plans. Knowledge

---

[19]e.g., Why is a *busy-wait-alert* operation being performed here?

[20]e.g., How is a *busy-wait-alert* operation implemented?

---

about the software *process*, that is, information about the tools and methods used in development which is not currently available in LaSSIE, could also be stored as plan knowledge (see Huff [22] for a plan-based approach to software process representation).

We have built a plan classification system [15] to experiment with various relationships that can be defined between plans, and inference algorithms that can derive these relationships.

## Knowledge Acquisition
We have discussed how the invisibility problem can be addressed by properly storing and retrieving knowledge of various kinds about the architecture of a system. In the current LaSSIE system, the concepts describing the actions and objects in the domain were generated by reading and understanding the architecture documents, and the comments in the source files. This is an entirely manual process. Though the end result of this effort is a KB that can be queried by developers in a manner that has an impact on the visibility of a system, building this KB is an intensive task, the automation of which would be a difficult but worthy goal.

Various researchers [7, 24, 49] have addressed the task of scanning code automatically and deriving various kinds of information about the function of the code. Letovsky [24] and Wills [49] are concerned with recognizing pieces of code that correspond to simple algorithmic fragments, such as loops, accumulations, etc. Biggerstaff [7] attacks a different problem; he seeks to mimic the process by which a programmer experienced in writing, for instance, device drivers, might scan the code for a new device driver (previously unknown to him/her) and use various clues such as typical variable names, characteristic data structures, etc., to explain the function of the new device driver.

Even if this problem remains

unsolved, we believe that the amelioration of the invisibility problem, the consequent reduction in architectural erosion, and the possible increase in code reuse should offset the cost of building the KB by hand.

## Summary

We have argued that the problem of invisibility and the related problem of complexity are caused by information barriers. We briefly discussed the limitations of existing attempts to solve these problems. We then presented a system called LaSSIE, which uses explicit knowledge representation and reasoning to address these limitations. The LaSSIE system is unique in its ability to provide semantic retrieval from an information base that integrates multiple views. Additionally, LaSSIE's KB provides an intelligent index (over a library of reusable components) that is attuned to the programmer's view of a large system. This approach shows promise, and also shows up some inherent limitations in the classification approach. The work is continuing; we are now focusing on both the representation and the knowledge acquisition issues. ◪

## Acknowledgments

### References

1. Adelson, B. and Soloway, E. The role of domain experience in software design. *IEEE Trans. Softw. Eng. 11*, 11 (1985).
2. Arango, G. Domain Analysis: From art form to engineering discipline. In *Proceedings of the Fifth International Workship on Software Specification and Design* (Pittsburgh, Pa., May, 1989).
3. *AT&T Technical Journal 64*, 1 (Jan. 1985), Part 2. Special Issue on the System 75 Digital Communications System.
4. Ballard, B.W. A lexical, syntactic, and semantic framework for a user-customized natural language question-answering system. In *Lexical-Semantic Relational Models*, Martha Evens, Ed., Cambridge University Press, 1988.
5. Beck, H.W., Gala, S.K., Navathe, S.B. Classification as a query processing technique in the CANDIDE semantic data model. In the *Fifth International Conference on Data Engineering* (Los Angeles, Calif., 1989).
6. Bellin, S. Personal Communication.
7. Biggerstaff, T.J. Design recovery for maintenance and reuse. MCC Tech. Rep. STP-378-88, Austin, Tex., 1988.
8. Borgida, A., Brachman, R.J., Resnick, L.A., and McGuinness, D.L. CLASSIC: A structural data model for objects. In *Proceedings of ACM SIGMOD-89* (Portland, Oreg., May–June 1989).
9. Brachman, R. and Devanbu, P. Domain modeling in a software information system. In *Proceedings of OOPSLA Workshop on Domain Modeling in Software Engineering* (New Orleans, La., 1989).
10. Brachman, R.J. and Schmolze, J.G. An overview of the KL-ONE knowledge representation system. *Cog. Sci. 9*, 1985, 171–216.
11. Brooks, F.P. No silver bullet: Essence and accidents of software engineering. *IEEE Comput. Mag.* (Apr. 1987).
12. Chen, Y.-F., Nishimoto, M., and Ramamoorthy, C.V. The C information abstraction system. *IEEE Trans. Softw. Eng.* (March 1990).
13. Curtis, W. Cognitive Issues in Reusing Software Artifacts. Vol II, *Software Reusability*. T.J. Biggerstaff and A.J. Perlis, Eds. ACM Press, N.Y., 1989.
14. Curtis, W., Krasner, H., and Iscoe, N. A field study of the software design process for large systems. *Commun. ACM 31*, 11 (Nov. 1988).
15. Devanbu, P. and Litman, D. Plan-based terminological reasoning. In *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning* (Boston, Mass., 1991).
16. Devanbu, P., Brachman, R., and Selfridge, P. Inference in support of retrieval for reuse in large software systems. In *Proceedings of the IEEE/SPS Workshop on Software Reuse* (Indialantic, Fla., 1989).
17. Doyle, J. and Patil, R.S. Language restrictions, taxonomic classification, and the utility of representation services. MIT/LCS/Tech. Mem. 387, 1989.
18. Embley, D.W. and Woodfield, S.N. A knowledge structure for re-using abstract data types. In *Proceedings of the Ninth International Conference on Software Engineering* (Monterey, Calif., 1987).
19. Fischer, G. and Schneider, M. Knowledge-based communication processes in software engineering. In *Proceedings of the Seventh International Conference on Software Engineering* (Orlando, Fla., 1984).
20. Frakes, W.B. and Nejmeh, B.A. An information system for software reuse. In *Proceedings of the Tenth Minnowbrook Workshop on Software Reuse* (1987), 142–151.
21. Horowitz, E. and Williamson, R. SODOS—A software documentation support environment: Its use. In *Proceedings of the Eighth International Conference on Software Engineering* (London, 1985).
22. Huff, K.E. Software process instantiation and the planning paradigm. In *Proceedings of the Fifth International Software Process Workshop* (Kennebunkport, Maine, Oct. 1989).
23. Johnson, W. and Soloway, E. PROUST: Knowledge-based program understanding. *IEEE Trans. Softw. Eng.* (Mar. 1981).
24. Letovsky, S. Plan analysis of programs. Ph.D. thesis, Yale University, New Haven, Conn., 1988.
25. Letovsky, S. Cognitive processes in program comprehension. In *Proceedings of the Second Workshop on Empirical Studies of Programmers*, E. Soloway and S. Iyengar, Eds., Ablex Publishers, Norwood, N.J., 1986.
26. Letovsky, S., Pinto, J., Lampert, R., and Soloway, E. A cognitive analysis of a code inspection. In *Proceedings of the Second Workshop on Empirical Studies of Programmers*, E. Soloway and S. Iyengar, Eds., Ablex Publishers, Norwood, N.J., 1986.
27. Levesque, H.J. and Brachman, R.J. Expressiveness and tractability in knowledge representation and reasoning. *Comput. Intell. 3* (1987), 78–93.

28. Modica, L. Personal Communication, 1989.

29. Neighbors, J. Software construction using components. Ph.D. thesis, University of California, Irvine, Calif., 1981.

30. Ostertag, E. and Hendler, J.A. AIRS: An AI-based Ada reuse system. Tech. Rep. CS-TR 2197, Computer Science Center, University of Maryland, 1987.

31. Parnas, D.L. Designing software for ease of extension and construction. *IEEE Trans. Softw. Eng. SE-5*, 2 (Mar. 1979).

32. Patel-Schneider, P.F. Small can be beautiful in knowledge representation. In *Proceedings of the IEEE Workshop on Principles of Knowledge-Based Systems* (Denver, Dec. 1984). Extended version appears as AI Tech. Rep. 37, Schlumberger Palo Alto Research, Palo Alto, Calif., Oct. 1984, pp. 11–16.

33. Patel-Schneider, P.F., Brachman, R.J., and Levesque, H.J. Argon: Knowledge representation meets information retrieval. In *Proceedings of the First Conference on Artificial Intelligence Applications* (1984), pp. 280–286.

34. Penedo, M.H. Prototyping a project master database for software engineering environments. In *Proceedings of the ACM SIGSOFT/SIGPLAN Software Engineering Symposium on Practical Software Development Environments* (Palo Alto, Calif., Dec. 1986).

35. Perry, D.E. Industrial strength software development environments. In *Proceedings of the IFIP Congress '89, The 11th World Computer Congress* (Aug./Sept. 1989, San Franciso).

36. Prieto-Diaz, R. Domain analysis for reusability. In *Proceedings IEEE COMPSAC-87*, Tokyo, Japan, (Oct. 1987).

37. Prieto-Diaz, R. and Freeman, P. Classifying software for reusability. *IEEE Softw. 4* (Jan. 1987), 6–16.

38. Robins, G. The ISI grapher manual, University of Southern California Information Sciences Institute, Tech. Manual ISI/TM-88-197, Feb. 1988.

39. Schmidt-Schauss, M. Subsumption in KL-ONE is undecidable. In *Proceedings of the First International Conference on Principles of Knowledge Representation and Reasoning* (Toronto, Canada, 1989).

40. Schneiderman, B. Empirical Studies of Progammers. In *Proceedings of the Second Workshop on Empirical Studies of Programmers*, E. Soloway and S. Iyengar, Eds., Ablex Publishers, Norwood, N.J., 1986.

41. Solderitsch, J., Wallnau, K., Thalhamer, J. Constructing domain-specific Ada reuse libraries. In *Proceedings of the Seventh Annual National Conference on ADA Technology* (Atlantic City, Mar. 1989).

42. Soloway, E. Learning to program = learning to construct mechanisms and explanations. *Commun. ACM 29*, 9 (1986).

43. Soloway, E. and Ehrlich, K. Empirical studies of programming knowledge. *IEEE Trans. Softw. Eng. SE-10*, (Sept. 1984).

44. Soloway, E., Pinto, J., Letovsky, S., Littman, D., and Lampert, R. Designing documentation to compensate for de-localized plans. *Commun. ACM 31*, 11 (Nov. 1988).

45. Smoliar, S.W. and Swartout, W. A report from the frontiers of knowledge representation. Draft Manuscript, USC/ISI, 1988.

46. Steffen, J.L. Interactive examination of a C program with $C_{scope}$. In *Proceedings of the USENIX Association* (Winter Conf., Jan. 1985).

47. Teitelman, W. The INTERLISP reference manual, Bolt, Beranek and Newman, 1974. Section 20 describes MasterScope, which was written by L.M. Masinter.

48. Van Lamsveerde, A., Delcourt, B., Delor, E., Schayes, M-C., and Champagne, R. Generic lifecycle support in the ALMA environment. *IEEE Trans. Softw. Eng. 14*, 6 (June 1988).

49. Wills, L.M. Automated program recognition, MIT Tech. Rep. 904, Cambridge, Mass., 1987.

50. Woods, M. and Somerville, I. An information system for software components. In *Proceedings of ACM SIGIR Forum, 22*, 3 (Spring/Summer, 1988).

CR Categories and Subject Descriptors: D.2.2 [**Software Engineering**]: Tools and Techniques—*Software libraries*; D.2.7 [**Software Engineering**]: Distribution and Maintenance; H.3.1 [**Information Storage and Retrieval**]: Content Analysis and Indexing; H.3.4 [**Information Storage and Retrieval**]: Systems and Software—*question-answering (fact retrieval) systems.*

**About the Authors:**

**PREM DEVANBU** is a member of Technical Staff at AT&T Bell Laboratories. His research interests include the use of knowledge representation and reasoning to help programmers understand large software systems and techniques for automatic recovery of useful knowledge from existing software. **Authors Present Address:** Premkumar T. Devanbu AT&T Bell Laboratories, 3C-441, Murray Hill, NJ 07974; prem-@research.att.com. Also at the Department of Computer Science, Rutgers Univ., Piscataway, N.J. 08854.

**RON BRACHMAN** is the Head of the Artificial Intelligence Principles Research Department at AT&T Bell Laboratories in Murray Hill, NJ. He has worked for many years on knowledge representation theory and applications, highlighted by early work on KL-One and recent work on CLASSIC. His work with others on the ARCON System provided the foundation for LaSSIE. **Authors Present Address:** Ronald J. Brachman, AT&T Bell Laboratories, Room MH3C-439, Murray Hill, NJ 07974; rjb@research.att.com

**PETER G. SELFRIDGE** is a member of Technical Staff at AT&T Bell Laboratories. His research interests center around applying knowledge representation technology to problems in software engineering, especially to support and understand the process of understanding large systems. **Authors Present Address:** Peter G. Selfridge, AT&T Bell Laboratories, Room 3C-441, Murray Hill, NJ 07974; pgs@research.att.com

**BRUCE BALLARD** is a member of Technical Staff at AT&T Bell Laboratories. His general interest is making computers easily accessible to all people. His research interests include natural language processing, artificial intelligence, and human interfaces. **Authors Present Address:** Bruce W. Ballard, c/o Helen Surridge, AT&T Bell Laboratories, MH3C-431, Murray Hill, NJ 07974; bwb@research.att.com

Definity is a trademark of AT&T