



Gio Wiederhold

Peter Wegner

Stefano Ceri

Toward Megaprogramming

Megaprogramming is a technology for programming with large modules called *megamodules* that capture the functionality of services provided by large organizations like banks, airline reservation systems, and city transportation systems. Megamodules are internally homogeneous, independently maintained software systems managed by a community with its own terminology, goals, knowledge, and programming traditions. Each megamodule describes its externally accessible data structures and operations and has an internally consistent behavior. The concepts, terminology, and interpretation paradigm of a megamodule is called its *ontology*.

The problem of scaling up software engineering to handle very large software systems is recognized as a bottleneck in software engineering [7]. The term *megaprogramming* was introduced by DARPA (see [4]) to motivate research on this problem. This article proposes a framework for megaprogramming in terms of software components called *megamodules* that capture the functionality of services provided by large organizational units such as banks, airline reservation systems, and city transportation systems (see Figure 1).

Megamodules realize greater abstraction power than traditional modules by stronger encapsulation mechanisms: they can encapsulate not only procedures and data, but also types, concurrency, knowledge, and ontology. Programming within a megamodule can be handled by traditional technology, while computations spanning several megamodules are specified by *megaprograms* in a *megaprogramming language*. Megaprograms provide the "glue" for megamodule composition and typically involve communication over networks.

Encapsulation Power of Software Components

Software components provide their clients with services specified by their interface and encapsulate (hide) local structures that implement their services [21], along the following progression:

- Functions and procedures: encapsulate statements and expressions
- Objects and classes: encapsulate data and procedures
- Megamodules: encapsulate behavior, knowledge, concurrency, ontology

Functions and procedures encapsulate statements and expressions. Objects and classes realize greater expressive power by encapsulating data (instance variables) as well as procedures (methods) [26]. Megamodules support effective hiding of domain-specific information by encapsulating behavior, knowledge, and concurrency. A better encapsulation concept allows megamodules to model very large organizations as well as software communities having a local language, culture, and traditions.

Megaprogramming is a form of *programming in the large* that encom-

passes not only largeness of size but also persistence (largeness in time), diversity (large variability), and infrastructure (large capital investment):

- Size: many lines of code
- Persistence: data survive the execution of programs
- Diversity: of concepts, people, knowledge, traditions, software communities
- Infrastructure: education, interfaces, tools, networks

Traditional programming in the large, as defined in [9], is concerned primarily with largeness of size and does not substantially address the requirements of persistence, diversity, and infrastructure. To achieve the needed scaleup in software design and management, much current attention is placed on modeling the process and on enhancement of computer-aided software engineering (CASE) tools [3, 15].

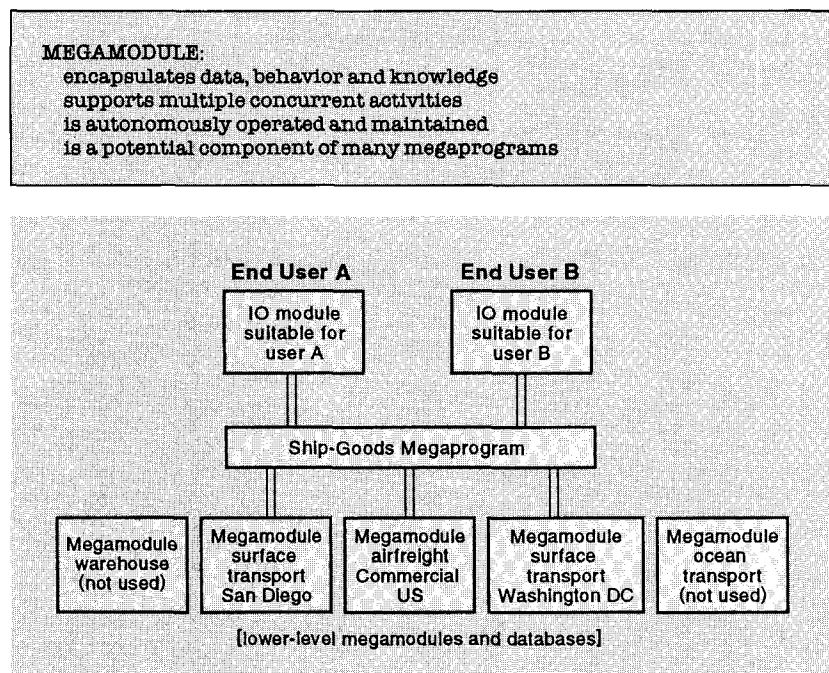
Persistence requires great attention to the management of change. Diversity of concepts and people requires the coordination of diverse languages, multiple programming paradigms, alternate information representation, and heterogeneous components. An adequate infrastructure requires powerful tools and programming environments, as well as support of networks, system evolution, and education [4].

The megaprogramming paradigm, modeling real-world services, is naturally parallel and distributed; it should support dynamic changes of the interface behavior and structure of constituent megamodules; for example, a bank may provide new services to its customers, while a corporation may create new manufacturing and service departments.

Megaprogramming is concerned with computations that span several software components. Its roots may be found in the work of [9], which introduced *module interconnection languages* (MILs) in order to facilitate the construction and management of programs consisting of collections of software components, thereby recognizing that structuring collections of modules is an essentially different intellectual activity from constructing individual modules. Recent work [1]

Figure 1. Properties of a megamodule

Figure 2. Example of use of megamodules by a megaprogram



on *module interface formalisms* (MIFs) provides the basis for the specification and implementation of MILs.

Megamodules require stronger module interconnection mechanisms than traditional modules because of their stronger encapsulation facilities and their support of heterogeneous interfaces and dynamic evolution. Therefore, megaprogramming languages (MPLs) extend MILs' expressive power by handling information transfer between heterogeneous modules, between dynamic queries and updates by users, between distributed network communication protocols, and by dynamically changing the specifications of interfaces.

The traditional CALL statement becomes overloaded when faced with this array of tasks. Work on parallel languages has already replaced a single statement with multiple segments [23]. We generalize on that concept. We present a particular MPL that is relatively low level as a programming language but high level when viewed as a networking language.

Ontologies

The term "ontology" derives from logic and artificial intelligence, where it denotes the primitive concepts, terminology, and interpretation paradigm of a domain of discourse [13]. Ontologies provide a conceptual framework for talking about an application domain and an implementation framework for problem solving. They provide a partitioning of context so that systems that would otherwise be too complex to understand become manageable [17, 29].

Our paradigm for megaprogramming is based on partitioning of tasks into large components, called megamodules, that provide greater spatial, temporal, and conceptual independence than traditional modules. Each megamodule should have an internally consistent ontology to reflect a uniform conceptual framework and problem-solving paradigm. Whereas ontologies in logic are specified by systems of axioms, the ontology of a megamodule should be determined by declarations specifying the properties of named entities accessible within the megamodule. Each declaration determines an *ontological commitment* to use the defined

term in the manner specified by the declaration. The collection of all ontological commitments determines the ontology. Declarations may be classified by the nature of their ontological commitment into declarations for variables, operations, behavior, and knowledge:

1. declare(variable, x , \exists)
2. declare(operation, successor, fun(x) . ($x+1$))
3. declare(behavior, vehicleclass, vehiclebehavior)
4. declare(knowledge, time, ontology-of-time)

Variables, operations, and behavior can be introduced by traditional declarations. Ontological commitments for knowledge are beyond the current state of the art and beyond the scope of this article, but are a central goal of artificial intelligence [13, 25].

Fortunately, we can separate the question of ontological commitment within megamodules from that of communication among megamodules. Since we are largely concerned with megamodule composition rather than with their internal structure, the mechanisms of ontological commitment, though important to the practical realization of megaprogramming, are not a primary concern of this article: the specification of ontologies local to a component remains a local issue.

An Example: Transportation of Goods between Two Cities

As an illustration, consider the transportation of goods between two cities. The general problem is that of shipping goods between any pair of cities by various forms of transportation, such as air, rail, or truck. Local transportation within each city is managed by megamodules, and intercity transportation by air, rail, and truck is managed by other megamodules. The complexity of each megamodule is substantial, and the total number of megamodules can be moderately large. We examine a specific application in greater detail.

Consider a logistics application for the shipment of goods between two U.S. Navy installations, NOSC in San Diego to NRL in Washington DC. This **Ship-Goods** task can be realized by a megaprogram using three megamodules for:

1. Surface transport in San Diego from NOSC
2. Commercial air freight
3. Surface transport in Washington DC from one of its airports to NRL

The megamodules (1) and (3) will be similar in structure and processing, but will use different databases. Local experts in these cities will maintain their own databases. Megamodule (2) will differ in ontology as well, as it represents and processes data according to its particular needs and according to the unique characteristics of the environment.

Megamodule (1), which deals with San Diego surface transport, will use tables that enumerate the trucks available to NOSC, the times when they are available, their capacities, and their loading facilities. Also, it will contain a San Diego roadmap and transit times for road segments at different times of day. Programs within this megamodule will deal with reserving trucks, moving them to NOSC, their loading, and travel to the airport. They will also be able to produce cost and time estimates in these tasks. However, when the departure time from NOSC is not known, the estimates may have excessive uncertainty or require very large tables of alternatives for various conditions.

Megamodule (2), for planning delivery of goods via commercial air freight, is itself a substantial application. It has to be able to choose among multiple airlines, obtain schedules and tariffs, and check for available space. This megamodule may itself invoke computational codes shared by other megamodules through conventional interfaces to program libraries.

Megamodule (3), for the Washington DC surface transport, will need to know routes from each of the three Washington DC airports. Rush hours differ among cities, and in Washington DC some critical roads are closed in one direction during rush hours. Geography makes a difference as well. Heuristics that are effective in San Diego, such as 'take a local road if the main road is blocked', fail when dealing with Washington DC bridges.

The megaprogram will coordinate

the activities of megamodules (1), (2), and (3) and present to its end user the *best* plan for the shipment. The execution sequence of megamodules may be affected by a variety of application requirements and constraints. If the *time-of-arrival* at NRL is the critical factor, then execution should probably commence with the Washington D.C. megamodule and proceed backwards. If the *time-ready* of the shipment at NOSC is supplied to the San Diego megamodule, then forward computation determines the arrival time. Without such requirements, the megaprogram might focus first on the most economical flight.

Megaprograms and Megaprogramming Languages

Megaprograms compose and schedule the computations performed by megamodules. Megaprogramming languages must allow flexible composition of megamodules and support both synchronous and asynchronous coordination schemes, decentralized data transfer, parallelism, and conditional executions.

In this section, we propose a particular MPL [30]; many concepts presented here, though specifically related to one proposed language, are quite general. The language is used for introducing functionalities to be implemented by megaprograms and megamodules. The reader should be aware that the language and its functionalities correspond to our vision; research and development are required to turn this vision into reality. The design principles of MPL can be summarized as follows:

- Traditional CALL statements for

invoking remote procedures and passing parameters are overloaded when used to compose large systems. MPLs segment their functions by separating parameter management in input and output from the invocation.

- Flexibility and asynchrony are enhanced by allowing autonomous execution of megamodules, while the invoking megaprograms retain great flexibility, since MPL statements can monitor, inspect, and constrain the execution of megamodules and can accept partial results.

- We are inspired by our database background; databases provide services autonomously; input and output parameters of megamodules are presented through database-like schemas, and the ideas behind megamodule optimization have several common features with query optimization.

The distinguishing feature of the proposed MPL is the substitution of the traditional CALL statement, which exerts control and provides communication between software components by three distinct statements: SUPPLY, INVOKE, and EXTRACT.

a. SUPPLY provides global arguments from a megaprogram to the megamodule.

b. INVOKE causes the megamodule to actually process these arguments and prepare results.

c. EXTRACT permits the megaprogram to extract results from the megamodule.

Further statements are used for inspecting and controlling megamodule execution. The most impor-

tant follow.

d. EXAMINE: to allow the megaprogram to check on the progress being made in obtaining results.

e. ESTIMATE: to cause the megamodule to provide a performance estimate, so the megaprogram can judge whether invocation is likely to be effective.

To illustrate our MPL, we give a simple megaprogram that fills an order by invoking two megamodules Producer and Consumer, supplying and extracting data for each megamodule and transducing the output of the producer into the format required by the consumer. Keywords of MPL are entirely in capital letters; megamodule and megaprogram names have an initial capital letter, and data names are entirely in small letters.

Our MPL is not complex from a programming language point of view, but is at a higher level than a module interconnection language. In the following, we examine the features of our MPL in greater detail; first we present the data structures for the supply and extract operation, next the operations for megamodule interaction, then the operations for inspecting at compile time the interfaces and content of megaprograms, and finally those for examining the run-time status of computations. Compile time and execution time can easily overlap.

Data Structure Interfaces for Supply and Extract

The *data structures* supported by a megaprogramming declaration enable an application to provide input and output parameters, of arbitrary complexity, between a megaprogram and megamodules; data structures are defined in the context of megamodules and are accessible by megaprograms. The MPL type system should be powerful enough to support complex data structures; therefore, it should include generalized type constructors, such as records, sets, multisets, lists, sequences, updatable arrays, and sparse arrays.

Data structures of megamodules are designed independently; they have arbitrary internal type systems and are defined asynchronously

Figure 3. A simple megaprogram

```

MEGAPROGRAM Fill-order(Producer, Consumer)
Input product, list-of-parts
SUPPLY TO Producer, list-of-parts
INVOKE Producer.product
EXAMINE (Status) UNTIL Status = DONE
EXTRACT FROM Producer, list-of-produced-items
PERFORM Transduce(list-of-produced-items, list-of-consumer-items)
SUPPLY TO Consumer, list-of-consumer-items
INVOKE Consumer.purchase
EXAMINE (Status) UNTIL Status = DONE
EXTRACT FROM Consumer, list-of-purchased-items cost
Output cost, list-of-purchased-items

```

from megaprogram definitions. Megamodules also provide self-describing EXPORT data structures that should be compatible with the MPL type system. However, a megamodule need not be able to handle all of the MPL type system; its capability must only be adequate to support all of its EXPORT data structures.

The technology we propose to borrow and adapt here derives from database schemas: the description of database contents and formats that permits many transaction programs to share a database and obtain data selectively and associatively. Accordingly, data structures of megamodules are declared in terms of a schema definition, but significant extension of database schema concepts will be required, since:

1. The data structures permissible in general megaprograms are more complex than the simple relations seen in most DBMSs [5].
2. Enough information must be provided to permit transduction of information between data structures that differ. Transduction occurs when data extracted from a megamodule are submitted to a successor megamodule: optimizations, described later, combine consecutive EXTRACT and SUPPLY operations to generate a direct flow of information among megamodules.

Operations for Megamodule Interaction

These operations split the traditional CALL module; statement into three parts to allow adequate management of large megamodules. Also, this split fosters asynchronous operation and parallel execution of megamodules. **SUPPLY.** The SUPPLY statement provides information to the megamodule as a subset of the megamodule's data structure. When this statement is compiled by the megamodule, the type of the MPL data structure is compared with the type of the data structure definition previously exported by the megamodule; the megamodule is responsible for managing type coercion transformations needed for ensuring compatibility of internal data structures with the MPL type system. Problems of this translation inside a

megamodule are not discussed in this article.

The supply of data values by megaprograms conforming to this general schema may differ. In particular:

- Many values may not be supplied by some megaprograms; default values are assumed by the megamodules for SUPPLY variables that have not been supplied.
- Some elements of the SUPPLY data structures may be defined to be PARAMETERS. Such parameter values are specified with the actual invocation of the megamodule. During execution of the SUPPLY statement the megamodule retains control. This assures that the data transmission to the megamodule is synchronized properly.

INVOKE. The statement "INVOKE megamodule WITH parameters" causes the initiation of asynchronous execution of the megamodule with any parameter values listed in the WITH portion. The parameter values are bound to those data structure items that were defined as parameters in the SUPPLY statement. After the INVOKE statement is processed, control is returned to the megaprogram; thus, the megaprogram can schedule other work in parallel.

EXTRACT. An EXTRACT statement provides megamodule results to the megaprogram. All variables that are expected to be extracted as a result of the megamodule execution are made available to the MPL environment as data structures accessed through the MPL statement. The result is fully valid when the megamodule state is DONE (corresponding to completion of the megamodule activity, as shown in Figure 3), but we place no constraints on earlier extraction of incomplete information.

During execution of the EXTRACT statement, the megamodule retains control again. This assures in this phase, as during the execution of SUPPLY, that the data transmission is synchronized by the megamodule.

Note that INVOKE is typically executed asynchronously, with the parameters of the invoke being part of the activation message; however, megamodule operation can remain

independent. For instance, a megamodule may always be active and keep the EXTRACT data structure up to date. Then, the INVOKE message only provides a synchronization handle and the EXTRACT operation may be immediately executed. A megamodule might be active periodically, based on external triggers. An EXTRACT operation then provides potentially not up-to-date snapshots; the EXTRACT operation may be executed at any time, since it does not affect the behavior of the megamodule.

Operations for Inspecting the Interface and Content of Megaprograms

The following operations, provided by megamodules, are *public* for use in a megaprogramming environment; they can be requested by megaprogrammers, and their result helps in the specification and design of megaprograms. Also, they can be used by MPL compilers to obtain information and may be used at execution time in interpretive situations. These operations permit megamodules to specify the data structures and defaults that they import and export; some megamodule-specific information can also be made available for inspection. Note that the information flows up, while traditionally subroutines depend for their environment on the invoking programs and thus are difficult to reuse.

IMPORT SUPPLY. The IMPORT SUPPLY statement causes the megamodule to export the supply data structures. Standard names are used to describe data structures, formats, extents, and sizes, according to MPL types. Elements of data structures may have default values that can be retrieved by an INSPECT command.

IMPORT EXTRACT. The IMPORT EXTRACT statement, similar to the IMPORT SUPPLY, causes the megamodule to export the EXTRACT data structure.

INSPECT. Inspection of a megamodule provides information about its ontology and specific parameters; the result of an INSPECT operation does not depend on the current state

of the module being inspected. Inspection could yield documentation, version information, or formal descriptions of methods and scope of applicability.

Operations for Examining the Status of Megamodules

A number of operations permit interaction of megaprograms and megamodules. Their function is to aid in execution to attain a high level of efficiency. They may provide, for instance, information that leads to a rescheduling of operational sequences of megamodule invocation. Since we assume a parallel computing environment throughout, such flexibility is essential to exploit parallel capabilities in a dynamic manner.

EXAMINE. In addition to the data structures used to supply and extract data, a megamodule maintains a

number of state variables. The **EXAMINE** statement accesses the most critical variable, namely **STATE**. Reasonable values for **STATE** variables in megamodules are: **DONE**, **IN_PROGRESS**, **IDLE**, **IN_ERROR**, or **WAITING**. For instance, **STATE = DONE** indicates that the megaprogram can safely extract results from the megamodule. Other variables that may be provided by a megamodule include, for instance, measures of nearness to a solution for an iterative program, the number of choices being considered for a search-type megamodule, and so on. These variables may be accessed at any time for a megamodule that has been invoked. The operation of **EXAMINE** is synchronous with respect to the megaprogram.

ESTIMATE. For megaprogramming optimization, estimates of exe-

cutation cost of megamodules may be needed [27, 28]; because of the encapsulation provided by megamodules, such estimation can be provided only from inside the megamodule. The result of execution of an **ESTIMATE** statement is a simple estimate of execution cost (time, dollars, result volume), with the degree of certainty of that estimate (0..1).

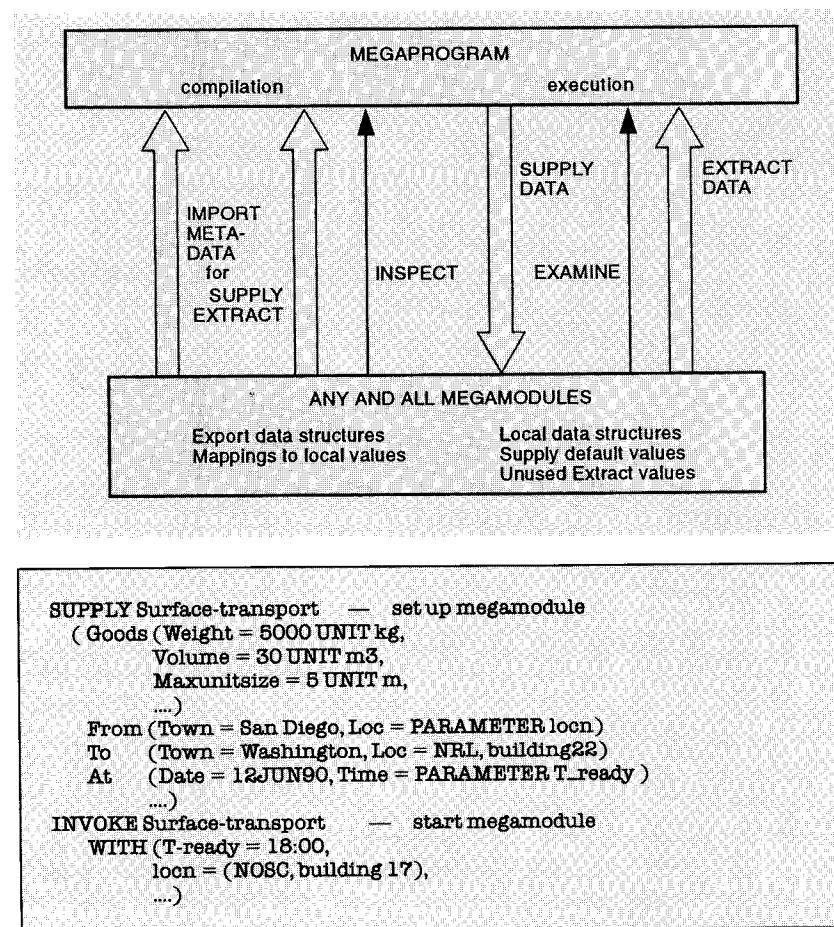
CONSTRAIN. The **CONSTRAIN** statement permits the resetting, by a megaprogram, of internal megamodule parameters obtained by inspection.

LIMIT. The **LIMIT** function constrains megamodule execution in terms of time or cost (as a surrogate for computational time spent). Not all megamodules will be able to accommodate such a request, but if they can, certain megaprogram tasks would be greatly simplified.

The interaction between megaprograms and megamodules is illustrated in Figure 4.

Figure 4. Information flow in megaprogramming

Figure 5. Example of a megaprogram invoking one megamodule



Composition and Interaction of Megamodules

Function composition in traditional languages can be specified by the notation $f(g(x))$, where f operates on the value produced by $g(x)$. Procedure composition is realized by executing a sequence of procedure statements $P; Q; R$, where each transforms a global state.

The composition of megamodules requires more than simple catenation. Information has to be transduced: collected, transformed, and forwarded among megamodules and input-output modules. To elucidate these MPL concepts, we can compare how direct communication among programs is specified in a Unix environment. MPL concepts can be viewed as a generalization of such Unix concepts.

The Unix shell plays the role of a megaprogramming language for Unix processes. Unix pipes provide direct communication among processes that transform a stream of input data into output data. A Unix command of the form $P1|P2$ allows the results of the process $P1$ to be incrementally piped to $P2$, so that $P2$ can process the output stream of $P1$

in parallel with further computation by *P1*.

The generalizations provided by MPL require that we identify each MPL operation with an explicit megamodule name, since the compositions of **IMPORT**, **SUPPLY**, **INVOKE**, **EXAMINE**, **EXTRACT**, and so forth are arbitrary and since any of the megamodules may be addressed. In MPL, megaprogram execution will be asynchronous, similar to what Unix supports with the form *P&*. However, the megamodules remain accessible through the **EXAMINE** statement so that their progress can be monitored. The combination of asynchrony and control adds much power to the megaprogramming language paradigm. Research will be needed to phrase a syntactically attractive language that permits specification of the parallelism and asynchrony implicit in the megaprogramming paradigm.

A simple pattern of interaction between the megaprogram and the megamodules can be realized by **SUPPLY**, **INVOKE**, and **EXTRACT**. Note that even this simple pattern of interaction improves the conventional development of applications. The **SUPPLY** and **EXTRACT** data structures provide a higher level of abstraction with respect to direct linking of software modules, enabling transduction of data structures among megamodules according to mechanisms that are defined in the megaprogram. Selective specification of parameters reduces setup and communication costs.

More complex interaction, yet still under the direct control of the megaprogrammer, is possible through the use of control structures **IF** . . . **THEN** . . . , whose effect is based on returned parameters from the **INSPECT**, **IMPORT**, **EXAMINE**, and **ESTIMATE** statements. By being able to assess the behavior of a megamodule in execution, it is possible to achieve important flexibility at execution time. In particular, execution of an unproductive module may be aborted, and a different execution pattern may be tried then. When computational resources are plentiful, multiple similar megamodules may be started in parallel, and only the one giving the best response will

be retained. This approach requires that the **INVOKE** primitive be non-blocking and also requires explicit time control within the MPL. The MPL will need an internal time-driven interrupt, so that it can decide when situations equivalent to **TIMEOUT** in conventional control of independent processes occur.

The access to a megamodule might be subject to limits, and the *megaprogram composer* might be able to check that limits are feasible through the **ESTIMATE** command. This merges the traditional compilation and execution phases into a complex interaction.

Compiling Megaprograms

Compilation of megaprograms is important to reduce the cost of extract, transduce, and supply sequences among modules. This cost is especially high if the megamodules are widely distributed. For compilation and initial optimization all interactions require communication with the site where the megaprogram resides. After compilation, information can be directly conveyed among the megamodules, and the load on the network as well as on the originating site of the megaprogram can be considerably reduced.

Many of the metainformation statements, previously described as aiding the execution of an interpreted megaprogram, will be used by the megaprogram compiler to provide the input for compilation decisions. Compilation enables *megaprogram optimization*, where the entire megaprogram is submitted to an optimization module. Such optimization would be responsible for:

- Combining consecutive **SUPPLY** and **EXTRACT** operations by generating a direct flow of information among megamodules
- Deciding the order of invocation of megamodules, in particular invoking them in the "optimal" (e.g., fastest, least expensive) order
- Introducing parallelism and asynchrony in the schedule of invocations of modules

These optimizations might seem novel in the field of software engineering, but have been successful in database systems. The information

provided by **INSPECT**, **EXAMINE**, and **ESTIMATE** is comparable to the information available to a database query optimizer.

The composition for our initial example, in the case that the **ready-time** at NOSC is given, might consist of:

1. Invoking the San Diego module with information about the **ready-time** and the load. Extracted from that module would be cost and time-required information and any needed schedule and routing data. Now the earliest arrival time of the load at Lindberg field in San Diego can be computed.
2. The airfreight megamodule is now supplied with these specifications, invoked, and similar results are extracted. The arrival time at Washington D.C.'s Dulles airport, for example, for the best flight is now computed.
3. The Washington D.C. surface module is now invoked with that time, and the final arrival time at NRL is now produced.
4. Desired output, extracted from these modules, is supplied to an output module by the megaprogram.

Many versions of such megaprograms can be envisaged. A megaprogram may iterate through the alternate airports, given as choices by the airfreight module, to optimize the Washington D.C. end of the transport task. Also, as indicated earlier, if the due time at NRL is the critical factor, then the megamodules will be invoked in the reverse order.

If minimal cost is the criterion, step 2 will likely gain priority and be scheduled first with minimal constraints.

Megaprogramming System Architecture

The architecture of a megaprogramming system consists of a collection of geographically distributed megamodules linked by high-capacity networks; the megaprogramming environment should include a repository and dictionary and should support megamodule execution and maintenance.

Megamodule Repository and Dictionary

Information about the collection of

megamodules in the megaprogramming environment is specified in a data dictionary that includes the following information about each megamodule:

1. The names of megamodules that are currently available.
2. A description of the interface of the megamodule (simply stated, some formal description of the function performed by the module).
3. The information needed to be supplied and available for extraction, expressed as data structures in the MPL type system.
4. Internal information about megamodules; these include the programming language used by the megamodule, its internal modularization, the schemas of databases accessed by the megamodule, the level of consistency provided by the megamodule applications in accessing the databases, and real-time constraints on megamodule execution.
5. An indication of the availability and the cost of using the megamodule.

Database techniques may provide useful concepts in organizing this information and making it available to the megaprogrammer through conventional and *ad hoc* query interfaces. In particular, the `IMPORT SUPPLY`, `IMPORT EXTRACT`, and `INSPECT` statements in MPL may be considered as formal interfaces to this dictionary.

Support for Megaprogram Execution

The data dictionary (library) is part of a program execution environment that may be replicated at many sites, including sites associated with individual megamodules, to permit distributed specification and execution of megamodules. The responsibilities of the megaprogramming environment in supporting execution include:

1. Type checking the consistency of `IMPORT` pairs for successive `EXTRACT` and `SUPPLY` operations. If the type systems are different, provide required routines for type coercion.
2. Generate executable code for megaprograms after compilation,

where this code might consist of several remote procedure calls and support for data transduction.

3. Generate compiled versions of the megamodules used in the context of a specific megaprogram, with optimizations due to the values provided in the `IMPORT` data structure.
4. Provide late type checking for run-time parameters.
5. Handle concurrent execution and synchronization of multiple megaprograms, initiated at different sites of the megaprogramming system.
6. Deal with the notion of transaction atomicity in the context of megaprograms.

Access to Databases

We assume that databases are part of megamodules. The semantics of the database being used is completely captured within the "ontology" of the megamodules and does not need to become "public." Instead, data structures for `SUPPLY` and `EXTRACT` do become public; they will be influenced by the schemas of databases that are stored within a megamodule.

We do not focus on fully transparent distributed databases where access to any data item can be made from any site regardless of the logical fragmentation and physical distribution; this issue is not the concern of megaprogramming. If used, their locking schemes will constrain megamodule execution parallelism. Our view is consistent with autonomous and independent access; we can lower consistency requirements to match those seen in federated databases [19] because we do not expect to provide comparable semantics across databases that are managed within different megamodules.

Transaction Management and Concurrency Control

We safely assume that durability (e.g., permanence of effects of committed transactions), local serializability (i.e., equivalence of each local schedule to some serial local schedule), and isolation will be provided by each local database management system. Therefore, the only remaining issues concern global atomicity and serializability of megaprograms, seen as large transactions operating as a unit over dispersed, heterogeneous

databases. This problem is quite difficult, but some practical solutions have been studied; they are briefly mentioned later.

Maintenance

A major motivation for our approach to megaprogramming is to improve the management of maintenance tasks. The development and maintenance of each megamodule are entirely under the authority of a local software development team. This provides freedom for independent development and maintenance of each megamodule, but also it confers a responsibility for providing reliable service to clients. Reliability can be ensured by globally enforced *acceptance testing* for both new modules and new versions (releases) of already existing modules. Some megamodules may have requirements for continuous execution of processes that monitor ongoing real-world activities and cannot be taken out of service.

Maintenance is greatly facilitated by the separation between megamodules and megaprograms. Maintenance is partitioned at both the megaprogram and megamodule levels.

- Maintenance within each megamodule is done by an expert of the megamodule's "ontology." We distinguish two types of changes. *Local changes* do not affect the `SUPPLY/EXTRACT` data structures and therefore do not require megaprogram recompilation. *Global changes* affect the `SUPPLY/EXTRACT` data structures and therefore cause megaprogram recompilation.
- Maintenance of a megaprogram is done by megaprogrammers; this may be due to changes of the invocation strategies of megamodules or to global changes to megamodules.
- No maintenance operation should span two megamodules.
- While the maintenance of a specific megamodule is taking place, it is advisable to support the older version together with the new version, so that use of existing megaprograms can continue.

Comparison with Related Work
Megaprogramming is primarily a

software engineering activity, but borrows heavily from the technologies of databases and of artificial intelligence.

The goals of megaprogramming, developing a component-based software technology for programming in the large, are also central goals of software engineering. Megaprogramming builds on earlier work on data abstraction [18, 21], object-oriented programming [12], and distributed computing [6].

Scaling up from objects to megamodules goes beyond object-oriented notions of encapsulation and message passing; techniques for managing megamodules and communicating among them are qualitatively different from those for objects. Because megamodules are so much more substantial than traditional modules, techniques for object management through types, classes, and inheritance are less applicable; megamodules require special database retrieval, browsing, and instrumentation methods.

Reusability is another related area in which traditional mechanisms [2] need to be supplemented to handle the larger granularity of megamodules. We have defined megamodules as the unit of use and reuse by megaprograms in a manner that use and reuse become indistinguishable. Removing the modules from the invoking environment is critical.

Kron and DeRemer [9] and subsequent work [1] on module interconnection languages and programming in the large are very influential, though megamodule interconnection mechanisms differ from traditional point-to-point communication mechanisms of object-based languages.

The need for a distinct language to manage module distribution alternatives of an application was recognized by Putilo in his thesis and implemented in his Polyolith software bus [22]. Our megaprogramming paradigm extends this view by proposing a compilable language for information flow control, conversion, and optimization (*transduction*), avoiding programmer involvement as system configuration and application needs change.

Databases and Transaction Processing

Work in databases has taken a number of recent directions [24] that are relevant to megaprogramming, including enhanced functionalities for data distribution and richer data models through the use of rich type systems for data definition. In this article we suggest that the ideal data model for megaprogramming should have a rich type system, but we deliberately avoid taking a position on whether it must be object-oriented. Indeed, the usefulness of concepts such as object identity and generalization hierarchies with distributed data, autonomously managed by each megamodule, is being debated.

Databases must already deal with heterogeneity, sometimes caused by finding that distinct information resources are located at distinct sites and under distinct management. The query program has no control then over the internals of the resource representation and very little control over the interface. These systems go by the name of *federated systems*. Languages to deal with combining information from federated systems have been developed [8, 19].

Transaction processing is another field of concern for megaprogramming. Classical transactions are minuscule tasks that only interact with one another through the database; typical workloads reach thousands of these small transactions per sec. However, it has become clear that not all database applications fit into this basic paradigm, so that new models of transactions are emerging, including *long transactions* [16] and *nested transactions* [20]. Long-lived transactions typically release their locks and run at a lower level of consistency with respect to full serializability. *Sagas* [10] are particular nested transactions that achieve global atomicity but permit the autonomous commitment of component subtransactions, whose effect can be undone by compensation (typically, by firing another transaction that will reconstruct the initial state). These mechanisms may find an application in megaprogramming, where the autonomy of each megamodule is a strict requirement.

Artificial Intelligence

Independence of modules has been a long-standing concern in artificial intelligence. Specifically, the paradigm of independent actors and agents has captured attention from [14] to [11], proposing various models for building complex computations and knowledge bases starting from knowledge that is independently provided in different contexts. Since programming was not an objective of AI researchers, this work was focused on bringing intelligence within modules and on providing loose and flexible control of module interconnection. In contrast, we propose mechanisms for megamodule control that, though maximally flexible, are fully specified through the MPL instructions; we do not expect that megamodules propose and negotiate their own computations, as is often expected in AI models.

The notion of megamodule ontology and the importance of ontological commitment (i.e., that entities within megamodules conform to their declarations) were developed in the context of knowledge representation models [13]; megaprogramming borrows this notion in order to define the expected qualities of megamodules, which however remain out of megaprogramming control!

Conclusion

We have defined megaprogramming as a departure from standard top-down programming practice, which has not been able to deal well with the problems encountered when building and maintaining large software systems. Our objective has been to provide a conceptual framework for the study of megaprogramming and to identify technical problems to be addressed in developing an effective technology of megaprogramming.

Concrete application of a comprehensive megaprogramming technology along the directions outlined in this article is clearly still far away. A number of problems need to be solved, not only related to megamodule integration, but also in the building of individual components, especially if we aim to enhance their portability and reusability. Improving "programming in the medium" is

a premise for the effective development of megaprogramming. Some of its principles can be incorporated in handcrafted systems.

Another critical factor to the success of megaprogramming is its integration with current standardization efforts. Emerging standards for supporting remote database access and distributed commitment protocols may provide useful concepts for the technological infrastructure needed for supporting megaprogramming requirements on a computer network (based on transfer and transduction of data structures among megamodules).

Though we see an existing trend toward the design of powerful MPLs and module interface formalisms, we must significantly raise the level of expected functionalities compared to the current state of the art; feasibility and applicability of these ideas must be demonstrated through further research and experiments. This article identifies a number of open research issues:

- Design of a megaprogramming language and type system.
- Design and implement the megaprogramming support environment.
- Design optimizers for megaprograms, with an increasing degree of sophistication, including the minimization of coercion and transduction costs.
- Establish trade-offs between pre-execution compilation and dynamic scheduling of megaprograms.
- Understand the role of emerging standards in the context of databases and transaction-processing systems for the implementation of megaprogramming interfaces.
- Understand deeply the nature of "ontologies" as a basis for programming megamodules.
- Understand and disclose the methodological issues concerned with building individual megamodules so as to assure their reusability in the development of new applications.
- Understand and disclose the methodological issues concerned with building complex megaprogramming applications. The switch to a service from a subroutine paradigm refocuses many of these issues.

Acknowledgments

We have received useful comments on this work from Tom Dean, Tom Doepfner, Steve Reiss, and Peter Rathmann. A long-standing interest in new programming paradigms has been encouraged by Sheldon Finkelstein, Witold Litwin, Carl Hewitt, Yoav Shoam, and others. Work of Stefano Crespi-Reghizzi, Letizia Tanca, and Roberto Zicari contributes insight into the applicability of these concepts. The initial motivation for this formulation was provided at the 1990 DARPA-ISTO Software Engineering PI meeting, where the need for new directions in programming in the large was explicated and where the term *megaprogramming* was initially introduced. ■

References

1. Belz, F.C., Luckham, D.C. and Purtilo, J.M. Application of ProtoTech technology to the DSSA program. In *Proceedings of the DARPA Software Technology Conference 1992* (Los Angeles, Calif., Apr. 28–30). Meridien Corp., Arlington, Va., 1992.
2. Biggerstaff, T.J. and Perlis, A.J. *Software Reusability*. Addison-Wesley, Reading, Mass., 1989.
3. Blum, B.I. Modeling-in-the-large. Tech. Rep. T.R.RMI-90-023, Milton S. Eisenhower Research Center, Johns Hopkins Univ., 1990.
4. Boehm, B. and Scherlis, B. Megaprogramming. In *Proceedings of the DARPA Software Technology Conference 1992* (Los Angeles, Calif., Apr. 28–30). Meridien Corp., Arlington, Va., 1992.
5. Ceri, S., Crespi, S., Zicari, R., Lamperti, G. and Lavazza, L. Algres. An advanced database system for complex applications. *IEEE Softw.* (July 1990).
6. Colouris, G.F. and Dollimore, J. *Distributed Systems: Concepts and Design*. Addison-Wesley, Reading, Mass., 1988.
7. Computer Science and Technology Board Scaling up: A research agenda for software engineering. *Commun. ACM* 33, 3 (Mar. 1990).
8. Dayal, U. Query processing in a multidatabase system. In *Query Processing in Database Systems*. Springer, New York, 1985, pp. 81–108.
9. DeRemer, F. and Kron, H.H. Programming in the small versus programming in the large. *IEEE Trans. Softw. Eng.* (June 1976). Also in *Tutorial on Software Design Techniques*, P.

- Freeman and A.I. Wasserman, Eds. IEEE Computer Society Press, 1983.
10. Garcia-Molina, H. and Salem, K. Sagas. In *Proceedings of the ACM-SIGMOD International Conference On the Management of Data*. San Francisco, Calif.
11. Genesereth, M.R. Proposal for research on informable agents. Logic Group Rep. 89–9, Stanford Univ., 1989.
12. Goldberg, A. and Robson, D. *Smalltalk, The Language and its Implementation*. Addison-Wesley, Reading, Mass., 1983.
13. Gruber, T.R. The role of common ontology in achieving sharable, reusable knowledge bases. In *Principles of Knowledge Representation and Reasoning*, J.A. Allen, R. Fikes, and E. Sandewall, Eds. Morgan-Kaufmann, San Mateo, Calif., 1991.
14. Hewitt, C., Bishop, P. and Steiger, R. A universal modular Actor Formalism for artificial intelligence. *IJCAI* 3 (Aug. 1973), 235–245.
15. Humphrey, W.S., Characterizing the software process. *IEEE Softw.* 5, 2 (Mar. 1988), 73–79.
16. Katz, R.H. Managing the chip design database. *IEEE Comput.* (Dec. 1983), 16–36.
17. D. Lenat, R.V. et al. Cyc, towards programs with common sense. *Commun. ACM* 33, 8 (Aug. 1990).
18. Liskov, B. and Guttag, J. *Abstraction and Specification in Programming Languages*. MIT Press, Cambridge, Mass., 1986.
19. Litwin, W. and Abdellatif, A. Multidatabase interoperability. *IEEE Comput.* 19, 12 (Dec. 1986), 10–18.
20. Moss, J.E.B. *Nested Transactions, An Approach to Reliable Distributed Computing*. MIT Press, Cambridge, Mass., 1985.
21. Parnas, D.L. On the criteria to be used in decomposing systems into modules. *Commun. ACM* 1972. Also in *Tutorial on Software Design Techniques*, P. Freeman and A.I. Wasserman, Eds. IEEE Computer Society Press 1983.
22. Purtilo, J.M. The Polyolith software bus. Tech. Rep. 2469, Univ. of Maryland, 1990. To be published in *ACM Trans. Prog. Lang. Syst.*
23. Rinard, M.C. and Lam, M.S. Semantic foundations of Jade. In *Proceedings of the 19th Annual ACM Symposium on Principles of Programming Languages* (Jan. 1992).
24. Silberschatz, A., Stonebraker, M. and Ullman, J. Database systems: Achievements and opportunities. *Commun. ACM* 34, 10 (Oct. 1991), 110–120.

25. Sowa, J.F. Issues in knowledge representation. In *Semantic Networks: Explorations in the Representation of Knowledge*, J.F. Sowa, Ed. Morgan-Kaufmann San Mateo, Calif., 1991.
26. Wegner, P. Concepts and paradigms of object-oriented programming. *Object-Oriented Messenger 1*, 1 (Aug. 1990).
27. Wiederhold, G. Mediators in the architecture of future information systems. *IEEE Comput.*, 25, 3 (Mar. 1992), 38-49.
28. Wiederhold, G. Model-free optimization. In *Proceedings of the DARPA Software Technology Conference 1992* (Los Angeles, Calif., Apr. 28-30). Meridian Corp., Arlington Va., 1992, pp. 82-96.
29. Wiederhold, G., Rathmann, P., Barsalou, T., Lee, B.S. and Quass, D. Partitioning and composing knowledge. *Inf. Syst.* 15, 1 (1990), 61-72.
30. Wiederhold, G., Wegner, P. and Ceri, S. Towards megaprogramming. Stanford Univ. Rep. STAN-CS-90-1341, Brown Univ. Rep. 90-20, and Politecnico di Milano, Dipartimento di Elettronica, Rep. 90-055, 1990.

CR Categories and Subject Descriptors: D.2.7 [Software Engineering]: Distribution and Maintenance; D.2.9 [Software Engineering]: Management; D.3.3 [Programming Languages]: Language Constructs and Features—*Modules, packages*; H.2.5 [Database Management]: Heterogeneous Databases; I.2.4 [Artificial Intelligence]: Knowledge Representation Formalisms and Methods

General Terms: Languages, Management, Standardization

Additional Key Words and Phrases: Mega programming languages and systems

About the Authors:

GIO WIEDERHOLD a professor of computer science and medicine at Stanford University, is on leave performing program management at DARPA. His research interests are the application and development of knowledge-based techniques for database management.

Author's Present Address: Computer Science Dept., Stanford University, Stanford, CA 94305; wiederhold@cs.stanford.edu; fax: (415) 725-7411

PETER WEGNER is a professor of computer science at Brown University. His research interests include programming languages and software engineering. He has authored or edited a dozen books including the first book on Ada in 1980, and books on research directions in software engineering and object-oriented

programming.

Author's Present Address: Computer Science Dept., Brown University, Providence, RI 02912; pw@cs.brown.edu or wegner@cs.brown.edu

STEFANO CERI is a professor in the Dipartimento di Elettronica e Informazione, Politecnico di Milano. His major fields of research interests include distributed databases, methodologies and tools for designing database applications, and deductive and active databases. **Author's Present Address:** Dipartimento di Elettronica e Informazione, Politecnico di Milano, Milano, Italy.

Unix is a trademark of Unix System Laboratories, Inc.

ACM 0002-0782/92/1100-000 \$1.50

Background research for this work has been partially supported by the DARPA contract N39-84-C-211. Earlier, Gio Wiederhold was supported by IBM Germany. Recent support includes a grant on coordinating updates to replicated copies on federated databases (FAUVE Contract IRI-9007753, NSF), a grant on knowledge-based mediators from the IBM Knowledge Systems Lab, Menlo Park, Calif., and a DARPA contract to ISI for development of standards for knowledge querying and manipulation. Stefano Ceri is partially supported by the ESPRIT Projects "Stretch" and "Idea" and by the CNR Project "Logidata+."

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

ALL NEW FROM IRWIN

DeLILLO

A First Course in
Computer Science with Ada

Includes IntegrAda™

a fully validated Ada compiler from AETECH, INC!

ALSO NEW

MACCABE

Computing Systems
Architecture, Organization, Programming

SHAPIRO

How to Program Well
A Collection of Case Studies

SCHACH

Software Engineering, 2/e

IRWIN

1818 Ridge Road Homewood, IL 60430

Circle #77 on Reader Service Card