

The Sisal Model of Functional Programming and its Implementation*

Jean-Luc Gaudiot
University of Southern California
Department of Electrical Engineering - Systems
Los Angeles, California
gaudiot@usc.edu

Wim Böhm and Walid Najjar
Colorado State University
Computer Science Department
Fort Collins, Colorado
{bohm, najjar}@cs.colostate.edu

Tom DeBoni and John Feo
Lawrence Livermore National Laboratory
Computer Research Group
Livermore, California
{deboni, feo}@llnl.gov

Patrick Miller
CaeSoft Development
Tracy, California
patmiller@sisal.com

Abstract

Programming massively-parallel machine is a daunting task for any human programmer and parallelization may even be impossible for any compiler. Instead, the functional programming paradigm may prove to be an ideal solution by providing an implicitly parallel interface to the programmer. We describe here the Sisal project (Stream and Iteration in a Single Assignment Language) and its goal to provide a general-purpose user interface for a wide range of parallel processing platforms.

1 Introduction

The history of computing has shown shifts from explicit to implicit programming. In the early days, computers were programmed in assembly language, mostly with the purpose of utilizing the available memory space as effectively as possible. This came at the cost of obscure, machine-dependent, hard to maintain programs, which were designed with high programming effort. Fortran was introduced to make programming more implicit, portable and less machine-dependent. With the advent of massively parallel computers and their promise of hundreds of gigaflops, we have seen a return to the explicit programming paradigm. Using for example C with explicit message passing library routines as "machine language," people attempt to utilize the available processing power to the largest extent, again at the cost of high programming effort, machine-dependent, and hard to maintain code. A compiler

for an implicitly parallel programming language alleviates the programmer from the task of partitioning program and data over the massively parallel machine.

It is our view that explicit parallel programming is a transition stage in the evolution of parallel computing and that implicit parallel programming languages will eventually become the norm as did high-level languages in the sequential paradigm. This will result in a tremendous improvement in *programming quality* in terms of programming effort, readability, portability, extendability and maintainability of parallel code. Another consequence will be the accessibility of parallel programming to a wider public that would make use of a wide spectrum of parallel computers: from a few processors on a chip to several thousand processor-machines.

On the other hand, functional programming [4] is an alternate programming paradigm which is entirely different from the conventional model: a functional program can be recursively defined as a composition of functions where each function can itself be another composition of functions or a primitive operator (such as arithmetic operators, etc.). This means that the programmer need not be concerned with explicit specification of parallel processes since independent functions are activated by the predecessor functions and the data dependencies of the program. This also means that control can be distributed. Further, no central memory system is inherent to the model since data is not "written in" by any instruction but is "passed" from one function to the next.

Sisal (Stream and Iteration in a Single Assignment Language) [20] is such a functional language which was originally designed by collaborating teams from the Lawrence Livermore National Laboratory, Colorado State University,

*This work is supported in part by DARPA Contract Number DABT63-95-0093

the University of Manchester and Digital Equipment Corporation. The goal of the project was to design a general-purpose implicitly parallel language for a wide range of parallel platforms.

The goal of this paper is to describe the last phases of this project as we are currently undertaking them. In section 2, a short tutorial will present the basic principles of Sisal. An early compiler implementation for shared memory systems is described in section 3. We turn our attention to Distributed Memory implementations in section 4, while section 5 introduces implementation of multithreading principles.

2 The Sisal language: a short tutorial

Sisal is a functional language that offers automatic exploitation and management of parallelism as a result of its functional semantics. In Sisal, and all functional languages, user-defined names are “identifiers” rather than variables, and they refer to values rather than memory locations. The values produced and used in a Sisal program are all dynamic entities, and their identifiers are defined, or bound to them, only for the duration of their existence in an execution. This is the dynamic of the data flow graph, in which graph nodes are operations, and values are carried on the arcs connecting the nodes. The extent of the existence of a value is the set of arcs on which it travels between the point of its definition and the point of its final consumption. The values that are defined by the graph arcs may or may not have names assigned to them within a program.

All Sisal expressions and higher-level syntactic elements evaluate to and return values based solely on the values bound to their formal arguments and constituent identifiers. This eliminates any possibility of side effects, and allows much richer analyses of program by the compiler than is typically the case for imperative languages.

To best illustrate these points, consider the following brief code fragment. It is written in Sisal 1.2, the language currently accepted by the Optimizing Sisal Compiler. The Sisal language is undergoing expansion and refinement, as discussed in other sections, but the syntax of version 1.2 will suffice for this example.

```
type OneDim = array [ real ];
type TwoDim = array [ OneDim ];

function generate( n : integer
                  returns TwoDim, TwoDim )
  for i in 1, n cross j in 1, n
    t1 := real(i) * real(j);
    t2 := real(i) / real(j)
  returns array of t1
         array of t2
  end for
end function % generate
```

The first two statements define type names for arrays. Note that no sizes are provided; all Sisal aggregate data items are dynamically created, resized, and de-allocated at runtime. Only the dimensionality and element types are relevant to the type specifications. The header for function “generate” shows that one integer argument, “n”, is expected, and two unnamed values will be returned. The returned values are arrays, but again, only typing and not sizing is specified. Names can be bound to these returned values at the site of invocation of function generate if the programmer wishes. An invocation of a function is semantically equivalent to the reproduction of the function code at that site, with appropriate argument substitution. This equivalence, called “referential transparency” is a fundamental property of functional languages, and is responsible for the strengths of the Sisal language.

All Sisal expressions, including whole functions and programs, evaluate to value sets. In the above case, the function evaluates to two arrays, which are the values of the expression contained in the function definition. The for-expression shown is a loop construct, which is an indicator of potential parallelism to the Sisal compiler. This loop has an index range defined as the cross product of two simpler ranges. This means that the body of the loop will be instantiated as many times as there are values in the index range, in this case $n*n$, and each body instantiation will be independent, since no data dependencies exist among them. The set of independent loop bodies can be executed in parallel or not, based on the compiler’s and the runtime system’s analyses of their costs, as well as on options specified by the programmer.

The appearance of the names “t1” and “t2” within the body of the loop should not be considered a reuse of these names in the sense of the reassignment of a variable in an imperative program. Instead, the names are used to define the computation in the loop body, and in fact these names will likely have no real existence within the executing program. The important point here is that each instance of the loop body, containing specific values for i and j , will independently compute specific instances of the values defined as $\text{real}(i)*\text{real}(j)$ and $\text{real}(i)/\text{real}(j)$; then all these separate

values will be gathered together into a pair of arrays and returned. The positions of the values in the result arrays are determined by the loop's index ranges, as are the overall size and dimensionality of the returned arrays. In this case, two two-dimensional arrays are returned, with index ranges from 1 to n in each dimension. The use of loop-temporary names is optional, and the return-clause above could be rewritten as:

```
returns array of real(i)*real(j)
array of real(i)/real(j)
```

with no change in the ultimate results. The loop body, then, would appear to be empty, but in fact there would still be a potentially parallelizable set of independent computations to be performed by it.

Further syntactic elements of the Sisal language include let-in statements, which allow for name definition and use; if-statements, which allow conditional name definition; record and union types, which allow for flexible data aggregation; streams, which allow for producer-consumer computations; and sequential loops, which allow true iteration, with specified data dependencies existing between iterations. I/O in Sisal is performed by passing inputs as arguments to, and receiving outputs as the results of, the outermost function. The values used for inputs and returned as outputs obey a syntax called "Fibre", which allows the demarcation of dynamically sized aggregates.

The Optimizing Sisal Compiler translates source programs into executable memory images, including the runtime system components required to automatically manage memory, tasking, and I/O. The amount of parallelism to be exploited by a program can be controlled by user options, and once compiled, a program can be executed by any number of worker processes, by way of a single runtime parameter. Similarly, compiler optimization behavior and runtime performance can be observed and controlled by options applied at various points during compilation and execution.

3 An Early Implementation: The Optimizing Sisal Compiler

Early implementations of the Sisal language were basic proofs of concept. Various interpreters, like DI[33], TWINE[21], and SSI[22] but for greatest execution speeds, a compiled code was needed. Sisal was ported to novel architectures like the Manchester Dataflow Machine [5], but complete acceptance for the language required porting to newly emerging shared memory parallel machines then coming to market (HEP[1], Encore, Sequent[25], Cray).

3.1 Update in place and copy elimination

Several key obstacles emerged. The first of these came from Sisal's semantic concept of making "copies" to preserve single assignment and referential transparency. A fragment like

```
let
  A := array[1: 1,2,3];
  B := A[2 : 999];
in
  A,B
end let
```

returned [1: 1,2,3], [1: 1,999,3]. The original value of A had to be preserved, so a copy was made to enable the replacement.

Consider instead swapping two elements of an array. The semantics of Sisal call for making a copy of C to make the first replacement, and making another copy for the second replacement. A FORTRAN programmer would never do this, instead they would write:

```
ITEMP = IARRAY[3]
IARRAY[3] = IARRAY[4]
IARRAY[4] = ITEMP
```

This program has no (array) copies and is done in place. Similarly, many 'in-place' algorithms that are efficient in space and time have been designed in imperative languages that would run poorly if all data structures had to be copied. Clearly there was room for improvement[35]. Consider the following small modification in Sisal:

```
C := array[1: 1,2,3,4,5];
T0 := C[3]; T1 := C[4];
D := C[3: T1];
E := D[4: T0];
```

When the array replacement is done on line 3, the value in "C" is dead. When the replacement is done, the C value can be thrown away. Instead of making a copy for D and throwing C away, we can safely use C's container instead. Similarly, the second replacement on line 4 is the last use of D, so E can use D's container. The Optimizing Sisal Compiler (OSC) makes heavy use of "update-in-place" and copy elimination analysis to eliminate many unnecessary copies[7]. In its simplest sense, update-in-place migrates reader operations before writers. Here, the C[3] and C[4] readers were moved before the replacement operations.

3.2 Build in place

Other important optimizations had to be developed [28]. For instance, many functional programs work on pieces of a large structure and then "Glue" the computed fragments together. For instance:

```

L := F(0,A[1],A[2]);
R := F(A[N-1],A[N],0);
III := for i in 2,n-1 ...
LIII = array_addl(III,L);
LIIIR = array_addh(LIII,R);

```

Semantically, this says: 1) Build piece L, 2) Build piece R, 3) Build size n-2 array III, 4) Stick L.III into a size n-1 array, 5) Stick LIII.R into a size n array. This seems to require two allocation and a lot of meaningless data movement. OSC introduced the idea of “BUFFERS” and persistent memory to the BACKEND of the compiler – leaving the frontend unchanged. Using a buffer system, the same operation proceeds as follows:

```

% New container of size N
.....
L....
L...R
LIIIR

```

This trick can be played even if the left and right pieces are loops. The beauty of this “Build-in-place” system is that memory can be preallocated and parallel computations can simply stick values where they belong – even if the original computation parts were from distant parts of the computations.

3.3 Reference Counting Optimization

We have seen that we can take advantage of an object ending its life just as we would otherwise need to copy. Reference counts were introduced to help know when: 1) a value can be updated in place and 2) a value’s memory can be recycled. Reference counting can be a very expensive operation on sequential machines – on parallel machines it is much worse!!! Parallel reference counts must be updated in a critical section. This operation keeps banging on locks every few operations, swamping the machine. Luckily, programs tend to have simple patterns of use for aggregate values and OSC can cleverly eliminate[32] nearly all reference counting in a program through lifetime analysis and operation merging.

3.4 Vectorization

On vector machines all the speed advantages come from routing array operations through temporary vector registers. OSC has fine control of loop placement so that reader/writer chains can be established. In imperative languages, this generally requires very careful writing of loops in order to clearly establish vector relationships between loops. The semantics of Sisal’s underlying dataflow representation make loops easy to move and so OSC can vectorize extremely well[9].

3.5 Loop Fusion, Double Buffering Pointer Swap, and Inversion

On scalar and scalar/parallel machines, loop overhead and memory fetch time tends to dominate computations. OSC can accommodate these machines by applying aggressive loop fusion. Fusion can rewrite loop code like

```

T0 := for i in 1,n returns
      array of A[i]*2
      end for;
T1 := for i in 1,n returns
      array of B[i]*3
      end for;
X := for i in 1,n returns
      array of T0[i] + T1[i]
      end for;

```

into

```

X := for i in 1,n returns
      array of A[i]*2 + B[i]*3
      end for;

```

eliminating the generation of two temporary arrays and setting values that can stream into internal registers from the cache. FORTRAN90 has similar semantics for its array operations: $X = A*2 + B*3$. Sisal’s OSC compiler can implement this more efficiently than a FORTRAN90 compiler because FORTRAN must know absolutely that neither A nor B is aliased to X. Sisal’s functional semantics insure that a left-hand-side of a definition is **never** an alias for a right-hand-side.

A typical scientific computation proceeds as follows:

```

for initial
  A := start_values()
while not done(A) repeat
  A := time_step(old A)
...

```

Here, a new version of A the same size is generated at each time step. A naive implementation of Sisal would allocate a new buffer for each time step and throw away the old even though it was the right size. OSC notices this and initially allocates a buffer outside the loop and pointer swaps the original and secondary buffers.

Consider a 1D smoothing function that averages values using a three point stencil – $X[i] = (A[i-1] + A[i] + A[i+1])/3.0$ At the endpoints a two point stencil is used instead. This is most easily expressed as:

```

X := for i in 1,n
  v := if i = 1 then (A[1]+A[2])/2.0
        elseif i = n then (A[n-1]+A[n])/2.0
        else (A[i-1] + A[i] + A[i+1])/3.0
        end if
  returns array of v end for

```

The if-tests appear to introduce a large overhead and to inhibit parallelism/vectorization/pipelining. The loop can be specialized doing the boundary computations separate from the inner computations. The inner computation is simple:

```
inner := for i in 2,n-1 returns
  array of (A[i-1]+A[i]+A[i+1])/3.0
```

Now we just need to glue on the lower bound computation and the upper bound computation. We need to be careful to handle zero trip loops here! This is done by producing an *array* of for the boundary values. In zero trip case, an empty array is generated. In all other cases, an array of size 1 is generated.

```
leftbound := for in 1,min(1,n) returns
  array of (A[1]+A[2])/2.0
end for
rightbound := for in max(2,n),n returns
  array of (A[i-1]+A[i])/2.0
end for
X := leftBound || inner || rightBound;
```

The max/min function calls make sure the zero trip cases are handled gracefully. The final catenation puts the results in the correct form. The catenations will actually be removed by build-in-place optimizations later in the optimization process.

4 Sisal90

The original Sisal definition has been extended and modernized. The new language includes language level support for complex values, array and vector operations, higher order functions, rectangular arrays, and an explicit interface to other languages like FORTRAN and C. See [12] for more details. The Foreign Language Interface (FLI) was included as it became clear that legacy codes were parallelized by changing only their kernels, leaving I/O and bookkeeping code intact. The FLI supports definition of functions **from** a language and **for** a language. The details of type interface, and parallel task management are all hidden.

5 A Prototype Distributed-Memory SISAL Compiler

In this section we present D-OSC, a prototype SISAL compiler for distributed-memory machines. D-OSC is an extension of OSC[10]. A new analysis phase for loop and array distribution has been added and the code generation phase has been modified to produce C plus MPI[13]

calls. The run-time system has been modified to support array distribution and communicating threads. Information needed to perform distributed memory optimizations is established by the analysis phase and provided to the code generator by decorating the appropriate IF2 nodes and edges.

The D-OSC model of execution is activation-based. A *master* process is responsible for dividing parallel loops into *slices* which will be executed by *slave* processes running in parallel. A slice is represented by an *activation record*, which contains a code pointer, the loop range, a unique loop identifier, input parameters to the slice, and destinations for values to be returned upon termination. Activation records are distributed over the machine and each processor maintains a local *activation record queue*. Upon completion of a slice, the slave process sends a completion message to the master and updates global results with locally-computed values. As a slice may contain a parallel loop, each slave can become a master and distribute its inner loop. Each processor must be able to receive a request for service from other processors, such as a read, write or allocate request. This is achieved by having a *listener* thread always active on every processor.

D-OSC is implemented in four phases, where each phase relies on the previous one.

- *Base*. This phase employs no analysis whatsoever, hence the code generated is very naive. Arrays and loops are distributed equally among processors. Message passing is used to access remote array elements. This compiler version serves as a reference for further implementations, providing useful information about the effectiveness of certain optimizations.
- *Rectangular Arrays*. The standard implementation of higher-dimensional arrays as arrays of arrays is replaced, where possible, by rectangular arrays with a single descriptor. Arrays and the loops creating or using arrays can be distributed by rows, block or columns. Not all loops are distributed.
- *Block Messages*. The reading and writing of remote array elements within certain loops is optimized by combining all the messages directed to the same processor into a single block message.
- *Multiple Alignment*. In previous phases arrays partitioning created disjoint sections of an array. In this phase *overlapping* array sections are created. This optimization reduces the number of messages passed, at the cost of using more space for the overlapping array sections.

5.1 Base Compiler

In OSC, the representation of arrays consists of an array descriptor, which contains information such as bounds, reference count, size, and other information, and a pointer to the physical array. OSC assumes a shared-memory model, and the pointers to the array descriptor provides a unique array identifier. An evident problem on a distributed-memory machine is that the descriptor pointer cannot be used as a unique identifier, since the address of the array descriptor is different for each processor. Hence a unique array identifier is created explicitly as the index in an array table that exists on each processor. The design of the array table permits a great deal of compatibility with existing array operations since the OSC array descriptor is preserved.

Arrays are partitioned according to the distribution of the creating loop, such that each array dimension is distributed equally among processors. To create the unique identifier for distributed arrays, the array creating process that creates the loop slices, allocates the array identifier and sends it as part of the activation message. Each slice then proceeds in parallel and updates its local entry in the array table.

Array access in the base compiler is straight-forward. The processor that owns the array element is determined. If the owner is the local processor, the array element is read directly from local memory, otherwise a request message is sent to the listener thread of the processor that owns the array element.

5.2 Rectangular Arrays

Rectangular arrays have only one descriptor per array, regardless of its dimensionality. Only one possibly remote memory access to fetch the array element is needed, where an arrays of arrays implementation requires n memory accesses to fetch an element. With one array descriptor per array traditional distributions, such as row, block and column, are easier to implement. A disadvantage of rectangular arrays is that sub-arrays cannot be shared. However, sharing also has disadvantages since update-in-place cannot be performed. Another disadvantage of rectangular arrays is that ragged arrays cannot be represented.

Arrays are created using IF2 `AGather` nodes. Consider the case of a SISAL *triple cross product* for loop that returns a three-dimensional array. In the original IF2, `AGather` nodes in the result graphs of all three nested loops create arrays. In the rectangular array case, the actions that `AGather` nodes perform are different. The outermost `AGather` node must perform the allocation of the physical space for the whole 3D-array, and the allocation of the single array descriptor. The innermost `AGather` node

fills in the elements of the array. The `AGather` node in the middle loop does not perform any action.

In the original IF2 an arrays of arrays access consists of multiple `AElement` nodes scattered over the dependence graph, each with one index input. For a rectangular array this must be transformed into one `AElement` node with all indices as input. The analysis phase identifies the *tree* of `AElement` nodes that is spanned by the output edge of a root `AElement` node and marks these nodes with information such as the level of the node in the tree and back-edges to ancestor nodes.

5.3 Block Messages

The implementation of array access operations described above is not always efficient for array references in loop bodies, as performing remote exchanges for individual elements is less efficient than performing at most one block exchange per producer-consumer processor pair. Our algorithm for obtaining block messages is a modification of the algorithm presented in [14].

5.4 Multiple Alignment

The last phase of the compiler implements the overlapping allocation of array sections presented in [15] for one-dimensional arrays. Overlapping allocation is applied to loops with *restricted affine* references as in the following loop model, where the `cjs` are constants.

```
for i in lo, hi
returns array of
    f(B1[i+c1], ..., Bm[i+cm])
end for
```

In the case of *single alignment*, i.e. $m = 1$, the first element of the consumer array is aligned with element $1 + c_1$ of the producer array. For the general case, the analysis phase identifies restricted affine loops, that create one-dimensional arrays while accessing elements from other one-dimensional arrays. Multiple alignment is achieved by identifying all the unaligned references required, and the maximum and minimum offsets of these with respect to the consumer index. The contiguous set of indices thus obtained is a superset of the producer array elements needed. Loops are marked *RightOverlap* and *LeftOverlap* to be used in the code generation phase to determine the upper and lower bounds for each slice.

5.5 Results

The benchmark programs used here to assess the effectiveness of the various optimization phases are Livermore

loops 1, 2, 3, 6, 7, 9, 12, 21, and 24, run on a network of four workstations. Since the initial objective is to reduce communication, we measure the total number of messages exchanged - the first number in table 1, and the total volume of communication - the second number in table 1.

Rectangular arrays decrease the number of messages exchanged for some of the programs that use 2-D arrays. However, sometimes the number of messages increases, as in loop 21. The reason for this is that the partitioning of the loops and arrays performed by the base compiler matches the accesses of the array elements better than the rectangular arrays implementation.

Most of the programs that access arrays benefit greatly from the implementation of block messages. The greatest improvements occur for loops 1 and 21. Loops 2 and 24 are sequential and the current implementation only generates block messages for references accessed in parallel loops. Loop 6 contains subscript expressions that use non loop variables.

Multiple alignment reduces the number of messages for the programs with producer consumer relations of one-dimensional arrays, such as loops 1, 7 and 12.

The volume of communication does not always decrease and varies with program characteristics. In loop 24, where the number of messages exchanged remains the same for all the compiler phases, the volume of communication increases. This is because the implementation of rectangular arrays increases the size of messages required to access array elements in order to accommodate the multiple indices of rectangular arrays.

5.6 Further Work

D-OSC is a prototype implementation that helps us to quantify compiler optimizations for distributed-memory machines. The following are some of the tasks that must be performed to improve DOSC. A *more efficient run-time system* is needed. There are situations where *run-time reference counting* is necessary. If one processor owns a reference count, each remote processor that updates the reference counter must contact this processor. When deallocating an array, the responsible processor must notify all processors that have partial copies of the array to deallocate the space. The implementation of *function call parallelism* is very easy under the activation-based model. However, inter-functional analysis is required to determine when and where to spawn functions. Currently loops are always distributed over all processors. If an analysis phase can *estimate the computation cost of a loop body*, then it is possible to generate code that decides the number of processors to be used. *Parallel I/O* must be implemented.

6 Architecture Support for Multithreaded Execution

Multithreaded execution has been proposed as a model for parallel program execution. As a model, or rather a family of models, multithreading views a program as a collection of concurrently executing sequential threads that are asynchronously scheduled based on the availability of data. This definition is intentionally wide in that it attempts to capture the common features among various multithreaded execution models proposed to date. It is important to note that in this definition the multithreaded execution model does not specify any form of memory hierarchy (it is common though to expect a single logical address space mapped over several nodes), any specific language feature, whether threads are user specified or compiler generated, the mechanism for communication and/or synchronization among threads, or the order of thread execution. There is no standard definition of a thread. In this document we will define a thread as the set of sequential instructions executed between two synchronization points. This does not preclude exploiting any instruction level parallelism within a thread.

Because of its functional properties, the Sisal language is particularly well suited as a source for multithreaded code. In this Section we present some results related to the evaluation of multithreaded execution. The performance of multithreaded execution is determined by the complex interaction of a number of inter-related architectural and compilation issues such as code generation, thread firing rules, synchronization schemes and thread scheduling. The relation between these issues and the tradeoffs between various alternatives for each of these issues is complex and requires extensive experimental evaluation. For example, the thread firing rule (which determines when threads are enabled) can be based on either a *blocking* or a *non-blocking* strategy. The blocking strategy is adopted in Iannucci's Hybrid Architecture [18], the Tera MTA [2] and the EARTH machine [17]. The non-blocking strategy is adopted in Monsoon [26, 27], *T [23] and the EM-4 [31] among others. The Threaded Abstract Machine (TAM) [11] is a software implemented multithreaded execution that has been ported to a number of platforms (such as the TMC CM-5 and the Cray T3D), it implements the non-blocking model.

In this section we present an experimental and quantitative evaluation of these two execution models. The evaluation includes their respective code generation strategies, its implications on data distribution and access and the performance of their respective storage hierarchies.

Table 1: Number of Messages, Communication Volume (4 PEs).

<i>Program</i>	<i>Type</i>	<i>Base</i>	<i>RecArrays</i>	<i>BlkMssgs</i>	<i>Mulal</i>
ll1	1D	6605, 132132	6603, 211368	603, 31368	303, 12168
ll2	1D	6443, 126656	6443, 213128	6443, 213128	6443, 213128
ll3	1D	3, 96	3, 168	3, 168	3, 168
ll6	1D, 2D	10533, 213036	13223, 430408	13223, 430408	13223, 430408
ll7	1D	4807, 86568	7503, 225168	953, 18968	303, 8568
ll9	2D	5883, 117136	2403, 76968	603, 28968	603, 28968
ll12	1D	9005, 180132	3003, 96168	1503, 24168	3, 168
ll21	2D	471, 8520	14403, 460968	123, 58728	123, 58728
ll24	1D	29703, 594096	29703, 950568	29703, 950568	29703, 950568

6.1 Blocking and Non-Blocking Models

The two multithreaded execution models considered here are based on data-driven dynamic execution with statically generated threads. A detailed description of these models can be found in the following:

Blocking Thread Model. In this model a thread may be suspended and then resumed later. It requires the architecture to support context switching: the saving of thread state, the selection of a new thread and moving threads between the different states. A thread is suspended after initiating a remote memory access.

In this model the synchronization and storage mechanisms rely on the Frame model: A frame represents a storage segment associated with each *invocation* of a code-block¹. All the threads within the code-block instance refer to its associated frame to store and load data values. Frames are of variable size and contiguously allocated in the virtual address space. The size of a frame is determined based on the maximum number of data values associated with the code-block. The Frame model is used in several multithreaded machines (e.g. TAM [11], StarT-NG [3] and the EM-4 and EM-X [19]). When an instance of a particular code-block is invoked, a frame is first allocated in a given processor's frame store and all the tokens generated within that code-block instance will be stored in that frame. The virtual address carried by a token is of the form:

$$\langle \text{frame pointer}, \text{frame offset} \rangle$$

A synchronization slot in the frame is associated with each thread. The synchronization slot is initialized with the count of the number of the inputs to the thread and is decremented with the arrival of each input. The thread is ready when the count reaches zero: a continuation corresponding to that thread is placed on the Ready Queue by the Synchronization Unit. A data value that is shared (i.e. read) by several threads in the same frame occupies only one location.

¹ A code-block is a semantically distinguishable unit of code such as a loop or function body.

The content of the frame is accessed by the Execution Unit via a read-only cache. Data values generated by the executing threads are sent to the Synchronization Unit which writes them in the frame. The frame is deallocated when all the threads in the code-block have terminated. The cache organization is that of a conventional write-through cache. Because of their variable size, frames are not aligned with cache blocks.

Non-blocking model. In the non-blocking model, once a thread starts its execution it runs until termination. Hence the thread is activated only when all the inputs to the thread are available. All memory accesses are performed as split-phase accesses: the request is issued by a thread but the result is returned to another thread. In this mode the thread never has to block waiting for a remote memory access.

The synchronization and storage mechanisms for the non-blocking threads is the Framelet model. A framelet is a fixed sized unit of storage that is associated with each thread instance. Each framelet has one synchronization slot for that thread instance. In the Framelet model a data value that is shared among several threads within a same code-block would be replicated in the framelet of each thread instance. The framelet is deallocated when the thread instance completes its execution. Because their size is fixed, framelets are aligned with cache blocks. The virtual address of a data value in the Framelet model is of the form:

$$\langle \text{context \#, thread \#, framelet offset} \rangle$$

Example. A code-block consisting of three threads is shown in Figure 1. The corresponding Frame memory model is shown in the Figure 2. The input z which is used by both threads A and B is stored at only one place in the frame memory. Each of the values in the frame memory is accessed by the frame base address and the offset into the frame. The first three slots are the counters for the three threads. Thus when the value y is stored only Counter A is decremented. But when z is stored both Counter A and

Counter B are decremented but only one copy of z is stored in the frame.

The Framelet memory model corresponding to the same code block is shown in Figure 3. There are three separate framelets. Each framelet contains the counter for the corresponding thread. Each framelet contains a memory location for all the inputs to the corresponding thread. Hence framelet A corresponds to one particular activation of thread A . The z is stored in the framelets of both threads A and B and both counters are decremented. This accomplished as two separate store operations.

6.2 Code Generation

The source language used for the generation of multithreaded code is Sisal. The compilation process converts the programs into two intermediate forms: MIDC-2 (non-blocking) and MIDC-3 (blocking) which are both derived from the *Machine Independent Dataflow Code (MIDC)* [30]. MIDC is a graph structured intermediate format: The nodes of the graph correspond to the von Neumann sequence of instructions and the edges represent the transfer of data between the nodes. MIDC has been used to generate the executable code for other multithreaded machines (e.g Monsoon and EM-4). Both MIDC-2 and MIDC-3 are highly optimized codes with optimization done both at the inter- and intra-thread level.

The code generation compiler is guided by: (1) Minimize synchronization overhead. (2) Maximize intra-thread locality. (3) Assure deadlock-free threads. (4) Preserve functional and loop parallelism in programs. The first phase of the code generation is the same for both models, it involves compiling the Sisal programs to IF2 using *OSC* [8].

The second phase differs for the two models in the handling of structure store accesses and the data storage models (frames or framelets). The long latency operations consist of remote memory reads, memory allocations, function calls and remote synchronizations. The remote memory references can be handled either as *Split-phase* access or *Single-phase* access. In the split-phase access the request is sent by one thread and the result is forwarded to another thread. In the single-phase access the result is returned to the same requesting thread. In the non-blocking model all remote accesses are split-phase. The blocking model uses both types of accesses: the code is analyzed at compile time to identify remote and local accesses. Remote accesses are implemented by split-phase operations while local accesses are regular memory access.

- In the non-blocking model (MIDC-2 form) all structure store accesses are turned into split-phase accesses. A split-phase access terminates a thread: the request

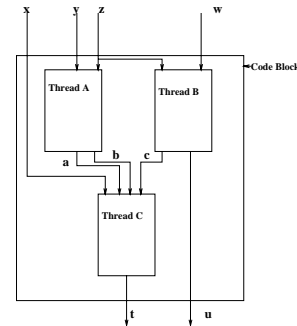


Figure 1: Code block with Three threads.

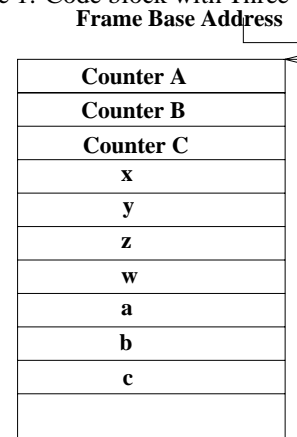


Figure 2: Frame Memory Representation.

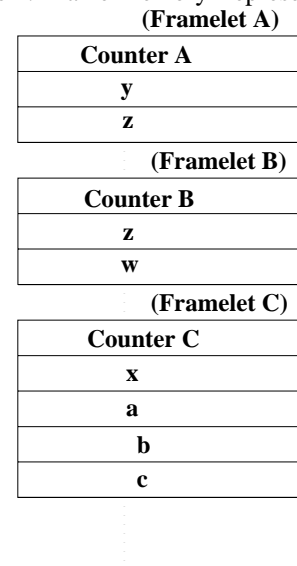


Figure 3: Framelet Memory Representation.

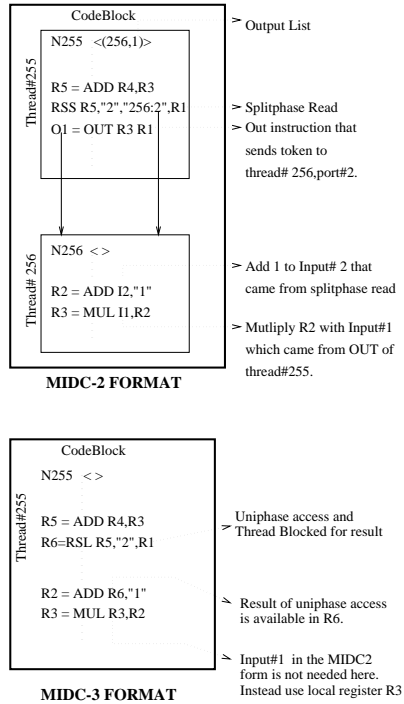


Figure 4: MIDC-2 and MIDC-3 code examples.

is sent by a thread but the result is returned to another thread. In this model a thread has never to block on a remote memory access. This model does not make any assumption regarding data structure distribution.

- In the blocking model (MIDC-3) the IF2 graph is statically analyzed to differentiate between local and remote structure store accesses: a local access does not terminate a thread while a remote one does. If the result of a structure store access is used within the same code-block where the access request is generated, the access is considered local. In this case, the thread will block until the request is satisfied. This model relies on a static data distribution to enhance the locality of access. Note that a data structure is often generated in one code block and used in several others in which case only one of the consumer code-blocks would have a local access.

Example The example in Figure 4 demonstrates the difference between MIDC-2 and MIDC-3. In MIDC-2, *Thread 255* performs a structure memory read operation. The read is performed as a split-phase access where the result is sent to *Thread 256*. *Thread 255* does not block, it continues execution until termination. When the results of the split-phase read is available it is forwarded to *Thread 256* which starts execution when all its input data

is available. There are no restriction on the processor on which *Thread 255* and *Thread 256* are executed.

In the MIDC-3 code, *Thread 255* and *256* belong to the same code-block. The read structure memory operation is a local single phase operation. Hence, the two threads become a single thread. The thread blocks when the read operation is encountered and waits for the read request to be satisfied.

Discussion of the Models. The main differences between the blocking and non-blocking models lie in their synchronization and thread switching strategies. The blocking model requires a complex architectural support to efficiently switch between ready threads. The frame space is deallocated only when all the thread instances associated with its code block have terminated execution which is determined by extensive static program analysis. The model also relies on static analysis to distribute the shared data structures and therefore reduce the overhead of split-phase accesses by making some data structure accesses local. The non-blocking model relies on a simple scheduling mechanism: data-driven data availability. Once a thread completes execution, its framelet is deallocated and the space is reclaimed.

The main difference between the Frame model and the Framelet models of synchronization is the token duplication. The Framelet model does require that variables which are *shared* by several threads *within* a code block be replicated to all these threads while in the Frame model these variables are allocated only once in the frame. The advantage of the Framelet model is that it is possible to design special storage schemes [29] that can take advantage of the locality of the inter-thread and intra-thread locality and achieve a cache miss rate close to 1%.

6.3 Effect of Network and Memory

This section studies the effect of cache misses and network latencies on the overall performance.

Effect of Network Latency. The plots in Figure 5 show the execution time versus average network latency with five and ten processors. These plots take into consideration the effect of cache misses on the execution time.

The performance of the blocking model degrades more than the non-blocking model with a non-ideal memory system. This is due to the higher miss rates of the frame memory used for the blocking execution model. Since the number of tokens generated both in the blocking and non-blocking models are almost the same, as the miss rate is higher for the blocking model, the memory access time is also higher and hence the performance degrades.

With a non-ideal memory at a network latency of 100 cycles the performance of the non-blocking model is comparable to the blocking model with a 50% success rate.

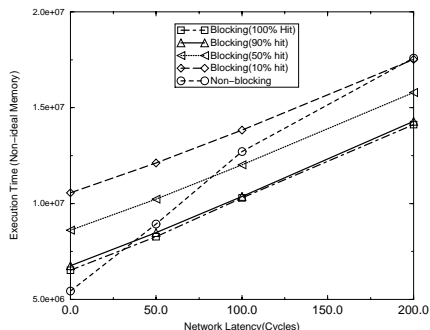


Figure 5: Execution time versus average network latency(10 processors and 8 cycle miss penalty)

Effect of Miss Penalty. Figure 6 shows the effect of varying miss penalty on the execution time. Three miss penalties (8, 25, 50 cycles) are used in this figure. A network latency of 50 cycles is assumed here. The non-blocking model masks higher miss penalties effectively compared with the blocking model. The slopes of all the four plots corresponding to the blocking model are the same. The slope is higher than the non-blocking model. Since the miss rates are higher for the blocking model, the miss penalties affect the execution time more severely than the non-blocking model.

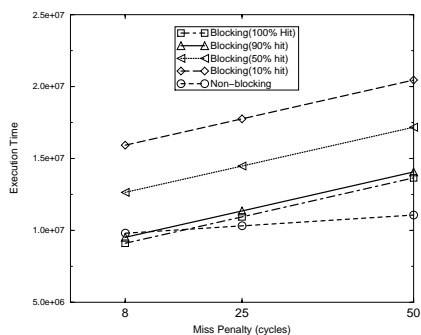


Figure 6: Execution time versus miss penalty(5 processors, 50 cycle network latency)

7 Conclusions and Future Research

The functional model of computation is one attempt at providing an implicitly parallel programming paradigm². Because of the lack of state and its functionality, it allows the compiler to extract all available parallelism, fine and coarse grain, regular and irregular, and generate a partial evaluation order of the program. In its pure form (*e.g.*, pure Lisp, Sisal, Haskell), this model is unable to express algorithms that rely explicitly on state. However, extensions to these languages have been proposed to allow a limited amount of stateful computations when needed. Instead, we are investigating the feasibility of the declarative programming style, both in terms of its expressibility and its runtime performance, over a wide range of numerical and non-numerical problems and algorithms, and executing on both conventional and novel parallel architectures. We are also evaluating the ability of these languages to aid compiler analysis to disambiguate and parallelize data structure accesses.

On the implementation side, we have demonstrated how multithreaded implementations combine the strengths of both the von Neumann (in its exploitation of program and data locality) and of the data-driven model (in its ability to hide latency and support efficient synchronization). New architectures such as TERA [2] and *T [24] are being built with hardware support for multithreading. In addition, software multithreading models such as TAM [11] and MIDC [6], are being investigated.

We are currently further investigating the performance of both software-supported and hardware-supported multithreaded models on a wide range of parallel machines. We have designed and evaluated low-level machine independent optimization and code generation for multithreaded execution. The target hardware platforms will be stock machines, such as single superscalar processors, shared memory, and multithreaded machines. We will also target more experimental dataflow machines, (*e.g.*, Monsoon [16, 34]).

References

- [1] S. J. Allan and R. R. Oldehoeft. Parallelism in SISAL: Exploiting the HEP architecture. In *19th Hawaii International Conference on System Sciences*, pages 538–548, January 1986.
- [2] R. Alverson, D. Callahan, D. Cummings, B. Koblenz, A. P. ortfield, and B. Smith. The Tera computer system. In *Proceedings 1990 Int. Conf. on Supercomputing*, pages 1–6. ACM Press, June 1990.

²Other attempts include the vector, data parallel and object-oriented paradigms.

- [3] B. S. Ang, Arvind, and D. Chiou. StarT the Next Generation: Integrating Global Caches and Dataflow Architecture. Technical Report 354, LCS, Massachusetts Institute of Technology, August 1994.
- [4] J. Backus. Can programming be liberated from the von Neumann style? *Communications of the ACM*, 21(8):613–641, 1978.
- [5] A. Böhm and J. Sargeant. Efficient dataflow code generation for sisal. Technical report, University of Manchester, 1985.
- [6] A. P. W. Böhm, W. A. Najjar, B. Shankar, and L. Roh. An evaluation of coarse-grain dataflow code generation strategies. In *Working Conference on Massively Parallel Programming Models*, Berlin, Germany, Sept. 1993.
- [7] D. Cann. *Compilation Techniques for High Performance Applicative Computation*. PhD thesis, Colorado State University, 1989.
- [8] D. Cann. Compilation techniques for high performance applicative computation. Technical Report CS-89-108, Colorado State University, 1989.
- [9] D. Cann. Retire FORTRAN? a debate rekindled. *CACM*, 35(8):pp. 81–89, Aug 1992.
- [10] D. Cann. Retire Fortran? a debate rekindled. *Communications of the ACM*, 35(8):81–89, 1992.
- [11] D. E. Culler et al. Fine grain parallelism with minimal hardware support: A compiler-controlled Threaded Abstract Machine. In *Proc. 4th Int'l Conf. on Architectural Support for Programming Languages and Operating Systems*, April 1991.
- [12] J. T. Feo, P. J. Miller, and S. K. Skedzielewski. Sisal90. In *Proceedings of High Performance Functional Computing*, April 1995.
- [13] M. Forum. *MPI: A Message-Passing Interface Standard*, 1994.
- [14] G. Fox, S. Hiranandani, K. Kennedy, U. Kremer, C. Tseng, and M. Wu. Fortran D language specification. Technical Report CRPC-TR90079, Center for Research on Parallel Computation, Rice University, P.O. Box 1892, Houston, TX 77251-1892, 1990.
- [15] D. Garza-Salazar and W. Böhm. Reducing communication by honoring multiple alignments. In *Proceedings of the 9th ACM International Conference on Supercomputing (ICS'95)*, pages 87–96, Barcelona, 1995.
- [16] J. Hicks, D. Chiou, B. Ang, and Arvind. Performance studies of Id on the Monsoon dataflow system. *Journal of Parallel and Distributed Computing*, 3(18):273–300, July 1993.
- [17] H. Hum, O. Macquelin, K. Theobald, X. Tian, G. Gao, P. Cupryk, N. Elmassri, L. Hendren, A. Jimenez, S. Krishnan, A. Marquez, S. Merali, S. Nemawarkar, P. Panangaden, X. Xue, and Y. Zhu. A design study of the EARTH multiprocessor. In *Parallel Architectures and Compilation Techniques*, 1995.
- [18] R. Iannucci. A Dataflow/von Neumann Hybrid Architecture. Technical Report 418, Ph. D Dissertation Technical Report TR-418, Laboratory for Computer Science, MIT, Cambridge, MA, June 1988.
- [19] Y. Kodama, H. Sakane, M. Sato, H. Yamana, S. Sakai, and Y. Yamaguchi. The EM-X parallel computer: Architecture and basic performance. In *Proceedings of the 22th Annual International Symposium on Computer Architecture*, pages 14–23, June 1995.
- [20] J. McGraw, S. Skedzielewski, S. Allan, D. Grit, R. Oldehoeft, J. Glauert, I. Dobes, and P. Hohensee. SISAL-Streams and Iterations in a Single Assignment Language, Language Reference Manual, version 1. 2. Technical Report TR M-146, University of California - Lawrence Livermore Laboratory, March 1985.
- [21] P. Miller. TWINE: A portable, extensible sisal execution kernel. In J. Feo, editor, *Proceedings of Sisal '93*. Lawrence Livermore National Laboratory, October 1993.
- [22] P. Miller. Simple sisal interpreter, 1995. <ftp://ftp.sisal.com/pub/LLNL/SSI>.
- [23] R. S. Nikhil, G. M. Papadopoulos, and Arvind. *T: A multithreaded massively parallel architecture. In *Proceedings of the 19th Annual International Symposium on Computer Architecture*, pages 156–167, May 1992.
- [24] R. S. Nikhil, G. M. Papadopoulos, and Arvind. *T: A multithreaded massively parallel architecture. In *Proceedings of the 19th Annual International Symposium on Computer Architecture*, pages 156–167, Gold Coast, Australia, May 19–21, 1992. ACM SIGARCH and IEEE Computer Society TCCA. *Computer Architecture News*, 20(2), May 1992.
- [25] R. Oldehoeft and D. Cann. Applicative parallelism on a shared-memory multiprocessor. *IEEE Software*, January 1988.
- [26] G. Papadopoulos. Implementation of a general-purpose dataflow multiprocessor. Technical report TR-432, MIT Laboratory for Computer Science, August 1988.
- [27] G. M. Papadopoulos and D. E. Culler. Monsoon: an explicit token-store architecture. In *Proceedings of the 17th Annual International Symposium on Computer Architecture*, pages 82–91, June 1990.
- [28] J. Rannelletti. *Graph Transformation algorithms for array memory optimization in applicative languages*. PhD thesis, U. California, Davis, 1987.
- [29] L. Roh and W. Najjar. Design of storage hierarchy in multithreaded architectures. In *IEEE Micro*, pages 271–278, November 1995.
- [30] L. Roh, W. A. Najjar, B. Shankar, and A. P. W. Böhm. An evaluation of optimized threaded code generation. In *Parallel Architectures and Compilation Techniques*, Montreal, Canada, 1994.
- [31] S. Sakai, K. Hiraki, Y. Yamaguchi, and T. Yuba. Optimal Architecture Design of a Data-flow Computer. In *Japanese Symposium on Parallel Processing*, 1989. in Japanese.
- [32] S. Skedzielewski and R. Simpson. A simple method to remove reference counting in applicative programs. In *Proceedings of CONPAR 88*, Sept 1988.
- [33] S. K. Skedzielewski, R. K. Yates, and R. R. Oldehoeft. DI: An interactive debugging interpreter for applicative languages. In *Proceedings of the ACM SIGPLAN 87 Symposium on Interpreters and Interpretive Techniques*, pages 102–109, June 1987.
- [34] K. Traub. Monsoon: Dataflow Architectures Demystified. In *Proc. Imacs 91 13th Congress on Computation and Applied Mathematics*, 1991.
- [35] M. Welcome, S. Skedzielewski, R. Yates, and J. Rannelletti. IF2: An applicative language intermediate form with explicit memory management. Technical Report TR M-195, University of California - Lawrence Livermore Laboratory, December 1986.