

4

REPORT DOCUMENTATION PAGE

READ INSTRUCTIONS BEFORE COMPLETING FORM

1. REPORT NUMBER AI Memo 1083		2. GOVT ACCESSION NO		3. RECIPIENT'S CATALOG NUMBER	
4. TITLE (and Subtitle) Optimization of Series Expressions: Part II: Overview of the Theory and Implementation				5. TYPE OF REPORT & PERIOD COVERED Memorandum	
7. AUTHOR(s) Richard C. Waters				6. PERFORMING ORG. REPORT NUMBER	
9. PERFORMING ORGANIZATION NAME AND ADDRESS Artificial Intelligence Laboratory 545 Technology Square Cambridge, MA 02139				8. CONTRACT OR GRANT NUMBER(s) IRI-8616644 N00014-88-K-0487	
11. CONTROLLING OFFICE NAME AND ADDRESS Advanced Research Projects Agency 1400 Wilson Blvd. Arlington, VA 22209				10. PROGRAM ELEMENT PROJECT, TASK AREA & WORK UNIT NUMBERS	
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) Office of Naval Research Information Systems Arlington, VA 22217				12. REPORT DATE December 1989	
				13. NUMBER OF PAGES 50	
				15. SECURITY CLASS. (of this report) UNCLASSIFIED	
				15a. DECLASSIFICATION/DOWNGRADING SCHEDULE	
16. DISTRIBUTION STATEMENT (of this Report) Distribution is unlimited					
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)					
18. SUPPLEMENTARY NOTES None					
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) sequences, optimization, loops ; program transformations, functional programming.					
20. The benefits of programming in a functional style are well known. In particular, algorithms that are expressed as compositions of functions operating on sequences/vectors/streams of data elements are easier to understand and modify than equivalent algorithms expressed as loops. Unfortunately, this kind of expression is not used anywhere near as often as it could be, for at least three reasons: (1) Most programmers are less familiar with this kind of expression than with loops; (2) Most programming languages provide poor support for this kind of expression; and (3) When support is provided, it is seldom efficient. (cont.)					

AD-A219 961

DTIC ELECTED APR 2 1990 S B D

DD FORM 1 JAN 73 1473

EDITION OF 1 NOV 65 IS OBSOLETE S/N 0102-014-66011

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

DISTRIBUTION STATEMENT A Approved for public release; Distribution Unlimited

In any programming language, the second and third problems can be largely solved by introducing a data type called *series*, a comprehensive set of procedures operating on series, and a preprocessor (or compiler extension) that automatically converts most series expressions into efficient loops. A set of restrictions specifies which series expressions can be optimized. If programmers stay within the limits imposed, they are guaranteed of high efficiency at all times.

A Common Lisp macro package supporting series has been in use for some time. A prototype demonstrates that series can be straightforwardly supported in Pascal.

Keenan

1/19

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
ARTIFICIAL INTELLIGENCE LABORATORY

A.I. Memo No. 1083

December 1989

**Optimization of Series Expressions:
Part II: Overview of the Theory and Implementation**

by

Richard C. Waters

Abstract

The benefits of programming in a functional style are well known. In particular, algorithms that are expressed as compositions of functions operating on sequences/vectors/streams of data elements are easier to understand and modify than equivalent algorithms expressed as loops. Unfortunately, this kind of expression is not used anywhere near as often as it could be, for at least three reasons: (1) Most programmers are less familiar with this kind of expression than with loops; (2) Most programming languages provide poor support for this kind of expression; and (3) When support is provided, it is seldom efficient.

In any programming language, the second and third problems can be largely solved by introducing a data type called *series*, a comprehensive set of procedures operating on series, and a preprocessor (or compiler extension) that automatically converts most series expressions into efficient loops. A set of restrictions specifies which series expressions can be optimized. If programmers stay within the limits imposed, they are guaranteed of high efficiency at all times.

A Common Lisp macro package supporting series has been in use for some time. A prototype demonstrates that series can be straightforwardly supported in Pascal.

Submitted to ACM TOPLAS

Copyright © Massachusetts Institute of Technology, 1989

This report describes research done at the Artificial Intelligence Laboratory of the Massachusetts Institute of Technology. Support for the laboratory's artificial intelligence research has been provided in part by the National Science Foundation under grant IRI-8616644, in part by the IBM Corporation, in part by the NYNEX Corporation, and in part by the Advanced Research Projects Agency of the Department of Defense under Office of Naval Research contract N00014-88-K-0487

The views and conclusions contained in this document are those of the authors, and should not be interpreted as representing the policies, neither expressed nor implied, of the National Science Foundation, of the IBM Corporation, of the NYNEX Corporation, or of the Department of Defense.

90 02 29 079

Contents

1. Sequence Expressions	1
Getting Rid of Loops	3
2. Optimizable Sequence Expressions	5
The Static Analyzability Restriction	6
The Straight-Line Restriction	7
The Preorder Restriction	7
The Sequence Intermediate Value Restriction	8
The On-Line Cycle Restriction	9
Obeying the Restrictions	12
Other Approaches to Restrictions	13
3. Series Expressions	15
4. A Common Lisp Implementation	19
5. A Pascal Implementation	27
6. The Optimization Algorithms	32
Systems Based on Similar Algorithms	35
7. Comparisons	36
Other Support for Sequence Expressions	36
Looping Notations	41
8. Benefits	44
9. Bibliography	46

Acknowledgments

The concept of optimizable series expressions has benefited from the suggestions of a number of people. In particular, A. Meyer, C. Perdue, C. Rich, D. Wile, Y. Feldman, D. Chapman, Y. Tan, and P. Anagnostopoulos made suggestions that led to significant improvements.

1. Sequence Expressions

The mathematical term 'sequence' refers to a mapping from the non-negative integers (or some initial subset of them) to values. Under the name of sequences [5, 36], vectors [23, 31, 32, 36], lists [36], streams [6, 19, 25, 30], sets [35], generators [18, 48], and flows [34], data structures providing complete (or partial) support for mathematical sequences are ubiquitous in programming.

The most common use for sequence data structures is as mutable aggregate storage. Essentially every programming language provides operations for accessing and altering the elements of at least one such structure.

Sequences have another use that is potentially just as important and yet is supported by only a few languages. Most algorithms that can be expressed as loops can also be expressed as functional expressions manipulating sequences. For example, consider the problem of computing the sum of the squares of the odd numbers in a file `Data`. This can be done using a loop as shown in the following Pascal [24] program.

```

type FileOfReal = file of Real;
function FileSumLoop (Data: FileOfReal): Real;
  var Sum: Real;
begin
  Reset(Data);
  Sum := 0;
  while not eof(Data) do
    begin
      If Odd(Data↑) then Sum := Sum+Sqr(Data↑);
      Get(Data)
    end;
  FileSumLoop := Sum
end

```



Alternatively, the sum of the squares of the odd numbers in the file can be computed using the sequence expression shown below. This expression assumes that four subroutines have been previously defined: `CollectSum` computes the sum of the elements of a sequence; `MapFn` computes a sequence from a sequence by applying the indicated function to each element of the input; `ChooseIf` selects the elements of a sequence that satisfy a predicate; and `ScanFile` creates a sequence of the values in a file.

```

function FileSum (Data: FileOfReal): Real;
begin
  FileSum := CollectSum(MapFn(Sqr, ChooseIf(Odd, ScanFile(Data))))
end

```

For those who are not accustomed to functional programming, the greater familiarity of the program `FileSumLoop` may make it appear preferable. However, the program `FileSum` has two important advantages. First, the patterns of computation that are mixed together in the loop in `FileSumLoop` are pulled apart. Second, each of these subcomputations is distilled into a subroutine. For example, the pattern of initializing a variable to 0 and then repetitively accumulating a result by addition is distilled into `CollectSum`.

OF



7 Codes

and/or

Dist Special

A-1

Because the subcomputations are pulled apart, they can be understood in isolation. The action of the expression as a whole is the composition of the actions of the subcomputations. This makes `FileSum` more self-evidently correct than `FileSumLoop`. The separation of the subcomputations also means that they can be altered in isolation. This makes `FileSum` easier to modify. The distillation of the subcomputations into subroutines makes `FileSum` shorter and enhances the reusability of the subcomputations. It also enhances reliability in two ways. Since the subcomputations are being explicitly reused instead of regenerated by the programmer from memory, there is less chance of error. In addition, since each subroutine can be reused many times, it is practical to work very hard to ensure that the algorithm used in the subroutine is robust.

Unfortunately, there are two problems that inhibit most programmers from writing programs like `FileSum`. First, most programming languages provide very few predefined procedures that operate on sequences as aggregates, rather than merely operating on their individual elements. Second, even in languages such as APL and Common Lisp, where a wide range of sequence operations are available, sequence expressions are typically so inefficient (2 to 10 times slower than equivalent loops), that programmers are forced to use loops whenever efficiency matters.

The primary source of inefficiency when evaluating sequence expressions is the physical creation of intermediate sequence structures. This requires a significant amount of space overhead (for storing elements) and time overhead (for accessing elements and paging). The key to solving the efficiency problem is the realization that it is often possible to transform sequence expressions into a form where the creation of intermediate sequence structures is eliminated. For example, it is straightforward to transform the expression in `FileSum` into the loop in `FileSumLoop`.

A transformational approach to the efficient evaluation of sequence expressions has been used in a number of contexts. For example, it is used by optimizing APL compilers [11, 21], Wadler's Listless Transformer [38, 39] which can improve the efficiency of programs written in a Lisp-like language, and Bellegarde's transformation system [7, 8] which can improve the efficiency of programs written in the functional programming language FP [5]. In addition, Goldberg and Paige [17] have shown that the transformational approach can be used to improve the efficiency of data base queries.

Unfortunately, it is not possible to completely transform every sequence expression into an efficient loop. There are two basic ways to deal with this problem. First, one can hide the issue from the programmer and simply transform what can be transformed. Second, one can develop a set of restrictions defining what can be transformed and communicate with the programmer about the transformability of individual sequence expressions.

The hidden approach, which is followed by all the systems above, has the advantage that programmers can benefit from increased efficiency in some situations without having to think about efficiency in any situation. However, it makes it difficult for programmers to think about efficiency when they want to, because they have no way of knowing for sure whether a given sequence expression will be completely transformed. This is significant, because sequence expressions typically remain quite inefficient if any part of them fails to be transformed. In addition, quite simple changes in an algorithm often suffice to change an untransformable expression into a transformable one. As a result, it is not really a

favor to hide the issue of transformability from programmers.

The most important contribution of the research reported here is a set of restrictions that can serve as a basis for the communicative approach to the transformation of sequence expressions into loops. As discussed in Section 2, these restrictions identify a class of *optimizable* sequence expressions that can always be completely transformed. The restrictions are novel in two ways. First, they are explicit. While every system that optimizes sequence expressions implicitly embodies some set of restrictions, the restrictions used are not explicit except in the work of Wadler [40]. Second, the restrictions in Section 2 are less strict than most other sets of restrictions. In particular, they are much less strict than Wadler's restrictions.

Sections 4-6 show how the communicative approach can be used to add comprehensive and efficient support for sequence expressions into any given programming language. This is done by adding a new sequence data type called *series* and a preprocessor that can transform optimizable series expressions into loops. The support for series utilizes the optimizability restrictions in two ways. One of the key restrictions is enforced by selecting the set of series operations so that the restriction cannot be violated. The rest of the restrictions are explicitly checked by the preprocessor. Non-optimizable expressions are flagged with warning messages and left unoptimized. If users take the time to make each series expression optimizable, they can have complete confidence that every series expression is efficient. This is facilitated by the fact that simple series expressions that only use each series once can always be optimized.

Section 3 presents the series data type and a broad suite of associated functions. Currently, the most comprehensive support for series is in Common Lisp. This implementation [46, 47] is presented in Section 4, along with an extended Lisp example showing how series expressions can be used. A prototype implementation [28, 45] shows that series expressions can also be added into Pascal. This implementation is presented in Section 5, along with an extended Pascal example of how series expressions can be used. Readers are encouraged to focus on whichever of Sections 4 and 5 discusses the most familiar language.

Section 6 presents the algorithms used to transform optimizable series expressions into loops. Section 7 concludes by comparing series expressions with related concepts. The comparison includes both other implementations of sequences and other approaches to expressing loops in ways that are easy to understand and modify.

Getting Rid of Loops

To fully appreciate the practical impact of series expressions in general and optimizable ones in particular, one must return to the perspective of sequence expressions as a notational variant for loops. The program `FileSum` is an example based on the Pascal implementation of series. The series expression in it is optimizable and is transformed into a loop essentially identical to the one in `FileSumLoop`. As a result, it is not merely the case that `FileSumLoop` and `FileSum` compute the same result using the same abstract algorithm; the two programs denote exactly the same detailed computation. Using the expression in `FileSum`, one gains the advantages of functional form without paying any price in terms of efficiency or anything else, because there is no change in anything other

than the notation.

The value of optimizable series expressions as an alternate notation for loops is directly related to the percentage of loops that can be profitably replaced by them. Any loop can be expressed as an optimizable series expression by converting the subcomputations used in it into series operations and composing them together. (At worst, the entire loop becomes a single series operation.) The value of doing this depends on how many fragments the loop can be decomposed into and how many of these fragments correspond to familiar computations. In general, the change is advantageous as long as there is at least one familiar fragment, because at the least, there is value in separating the familiar from the unfamiliar.

An informal study [41] revealed that 80% of the loops programmers typically write are constructed solely by combining just a few dozen familiar looping fragments. (A somewhat similar study is reported in [15].) Experience with the Lisp implementation of series expressions indicates that at least 95% of loops contain some familiar computation. Given this, the practical benefit of optimizable series expressions can be summarized as follows:

**Optimizable series expressions are to loops
as structured control constructs are to gotos.**

Structured control constructs (*if...then...else*, *case*, *while...do*, *repeat...until*) are not capable of expressing anything that cannot be expressed as gotos. In addition, there are probably a few algorithms for which the use of gotos is preferable. Nevertheless, in almost every situation, structured control constructs are much better to use than gotos. They are better, not because they allow more algorithms to be expressed, but because they allow the same algorithms to be expressed in a way that is much easier to understand and modify.

Optimizable series expressions have the exact same advantage. They do not allow algorithms to be expressed that cannot be expressed as loops. However, they allow algorithms to be expressed in a much better way. The only place where the analogy with structured control constructs breaks down is that while one can argue that gotos are never needed, there are definitely some algorithms that can be expressed better as loops than as optimizable series expressions.

At the current time, most programs contain one or more loops and most of the interesting computation in these programs occurs in these loops. This is quite unfortunate, since loops are generally acknowledged to be one of the hardest things to understand in any program. If optimizable series expressions were used whenever possible, most programs would not contain any loops. This would be a major step forward in conciseness, readability, verifiability, and maintainability.

2. Optimizable Sequence Expressions

As noted above, the primary source of inefficiency when evaluating sequence expressions is the creation of intermediate sequence structures. There are two aspects to this. First, a subexpression may compute a sequence only part of which is used by the rest of the expression. Second, even when all the elements are used, a significant amount of space and time overhead is required to construct physical data structures containing them.

The first problem can be overcome by using lazy evaluation [16] to ensure that sequence elements are not computed until they are actually used. (This also makes it easy to support unbounded sequences.) However, lazy evaluation does little to assist with the second problem. In particular, in simple situations where the sequence elements are all used, lazy evaluation wastes time and does not save any space. Although the same elements are computed, time is wasted, because coordination overhead is required to decide when to compute the elements. The same space is used, because each element has to be stored after it is computed. (Otherwise, a later reuse of an element would require recomputation.)

Both of the problems above can be solved by *pipelining* the evaluation of sequence expressions. When this is possible, elements are computed on demand without coordination overhead and do not have to be stored.

Definition 1 (Pipelined) *The evaluation of a sequence expression E is pipelined if the evaluation proceeds in such a way that the following conditions hold for every sequence S computed by any part of E . First, each element of S is computed at most once. Second, when an element is computed, it is used wherever it needs to be used, and then discarded before any other element of S is computed.*

The primary implication of the definition above is that, while some of the procedures called by E may buffer sequence elements within themselves, no additional buffering is required when evaluating E . Each sequence is transmitted between the procedure that creates it and the procedures that use it, one element at a time.

When compile-time pipelining is possible, a sequence expression can be evaluated as efficiently as a loop. Unfortunately, pipelining is not always possible. Any system that does pipelining can only do so for a restricted class of sequence expressions. Whatever the system, it is valuable for these restrictions to be made explicit. If in addition, the programmer is given feedback about which expressions fail to meet the restrictions, two advantages can be obtained. First, the programmer is given a clear picture of which expressions are efficient and which are not. Second, the programmer has the opportunity to change inefficient expressions so that they can be pipelined.

It would be nice to have a set of necessary and sufficient restrictions that specifies exactly which sequence expressions can be pipelined. However, there are a number of reasons why a somewhat stricter set of restrictions are of greater pragmatic benefit. First, the restrictions must be associated with practical algorithms that can check whether the restrictions hold and can perform the pipelining. Second, the restrictions must be simple enough to understand that programmers can succeed in fixing expressions that violate them.

The primary contribution of the work presented here is a set of five restrictions that satisfy the subsidiary goals above without being excessively strict. These restrictions define a class of *optimizable sequence expressions*. It can be shown that every optimizable sequence expression can be pipelined at compile time by transforming it into a loop, using the algorithms in Section 6.

Definition 2 (Optimizable Sequence Expression) *An optimizable sequence expression is a sequence expression that satisfies the following restrictions:*

- 1) *Optimizable sequence expressions must be statically analyzable.*
- 2) *Optimizable sequence expressions must be straight-line computations.*
- 3) *Procedures called by optimizable sequence expressions must be preorder.*
- 4) *Intermediate values in optimizable sequence expressions must be sequences.*
- 5) *Every non-directed data flow cycle in an optimizable sequence expression must be on-line.*

The restrictions in Definition 2 can be divided into three groups. A restriction analogous to the first one is required for any optimization that is to be applied at compile time as opposed to run time. The second restriction greatly simplifies the algorithms in Section 6, but is undoubtedly stronger than necessary. It is hoped that it will be weakened in the future. The remaining three restrictions are the theoretical heart of the matter.

In this section, programs are discussed from the point of view of data and control flow graphs rather than program text. In this representation, procedure calls and data literals are represented by nodes. The nodes have labeled ports corresponding to data inputs and outputs. There are no limits on the numbers of inputs or outputs. Data flow is represented by directed arcs connecting output ports to input ports. A given output can be connected to several inputs. Control flow is modelled by additional arcs connecting special control inputs and outputs. In the context of these graphs, the term *sequence expression* is defined as follows.

Definition 3 (Sequence Expression) *A sequence procedure is a procedure that consumes or produces a sequence. A sequence data flow is a data flow arc that transmits a sequence. A sequence subexpression is a (not necessarily proper) subset of the nodes in a data and control flow graph that can be built up through repetitive application of the following two rules. Each node corresponding to a sequence procedure call or literal sequence is a sequence subexpression. If there is a sequence data flow from a node in a sequence subexpression X to a node in another sequence subexpression Y , then $X \cup Y$ is a sequence subexpression. A sequence expression is a sequence subexpression that is not a proper subset of any other sequence subexpression.*

The Static Analyzability Restriction

As with most other optimization processes, a sequence expression cannot be pipelined at compile time, unless it can be determined at compile time exactly what computation is being performed. To make sure that this will be possible, it is required that every use

of a sequence procedure in an optimizable sequence expression be an explicit call on a predefined sequence procedure. This ensures that it will always be clear exactly what sequence computation is being performed.

Similarly, it is required that every sequence value come directly from a sequence procedure call. This ensures that it will always be clear exactly how the sequence is being computed. Unfortunately, this also implies that a sequence cannot be stored in any kind of data structure. This is undoubtedly somewhat stronger than necessary.

Arbitrary storage of sequences in data structures is bound to block compile-time pipelining. However, certain limited cases could be allowed. For instance, one can sometimes determine how a sequence contained in another sequence is being computed. The practicality of this has been demonstrated by compilers for APL [11], Hibol [34], and Model [32].

The Straight-Line Restriction

In the interest of simplicity, optimizable sequence expressions are required to be straight-line computations, not subject to any conditional or looping control flow. That is to say, it must be the case that whenever a sequence expression is evaluated, every procedure call in it is evaluated exactly once.

While it is likely that looping control flow in sequence expressions must be prohibited, simple conditional control flow could probably be allowed. However, this would complicate pipelining in a number of ways and much of what can be done using conditional control flow can be done using operations like `ChooseIf` instead.

The fact that sequence expressions are straight-line computations means that they can be pipelined without worrying about control flow. As a result, control flow is not mentioned in the rest of this discussion.

The Preorder Restriction

Suppose that a procedure call \mathcal{F} uses a sequence computed by another procedure call \mathcal{G} . For the computation of these two calls to be pipelined, two conditions must be satisfied. First, it must be the case that the sequence elements are created and consumed one at a time. Second, the elements must be consumed in the same order they are created. A good way to ensure that this will always be the case is to pick some fixed order and require that every sequence procedure process every sequence one element at a time in that order. Given a desire to support unbounded sequences, a good order to pick is preorder.

Definition 4 (Preorder) *A sequence procedure is preorder if it processes the elements of each of its sequence inputs and outputs one at a time in ascending order starting with the first element.*

In this paper, the word 'function' is reserved for referring to mathematical functions, while the word 'procedure' is used to refer to the implementation of a function in a programming language. With that in mind, note that the definition above applies to procedures, not functions. It is a property of the way a computation is performed, not of the mathematical relationship between the input and the output. Any mathematical

sequence function can be implemented as a preorder procedure. At worst, one can simply read the input elements into a buffer in preorder, compute the result, and then write the elements of the result in preorder.

A property that does apply to mathematical functions is the amount of internal buffering that is required when implementing the function as a preorder procedure. Fortunately, most sequence functions can be implemented as preorder procedures that do not require internal buffering of input or output elements. For instance, the procedure `MapFn` operates as follows: it reads one element of each input sequence, applies the indicated function to them, writes the resulting output element, and then goes on to the next group of input elements.

The only common functions where internal buffering is required are ones that rearrange the input elements (e.g., reversal, rotation, and sorting). In some cases this buffering is solely due to the needs of preorder processing. For instance, while preorder reversal requires the entire input to be read before the first output element can be produced, some non-preorder implementations require no buffering. In other cases (e.g., sorting) the buffering is required no matter how the operation is implemented.

The goal of pipelining is the elimination of the external buffering of elements between procedure calls. The algorithms in Section 6 are applicable no matter how much internal buffering the individual procedures use. As a result, the restrictions in Definition 2 do not place any limits on internal buffering. (This issue is discussed further at the end of this section.)

Nevertheless, using large internal buffers clearly violates the spirit of what pipelining is trying to achieve. It is of no benefit to eliminate external buffering if this is replaced by internal buffering. Therefore, in the interest of overall efficiency, the functions operating on series (see Section 3) are limited to ones that can be implemented as preorder procedures without internal buffering of sequence elements.

The Sequence Intermediate Value Restriction

Even if a sequence expression satisfies the three restrictions above, it may still not be possible to pipeline its evaluation. The problem is that if a sequence is used in two places, the two uses may place incompatible constraints on the times at which sequence elements should be computed.

The following program shows an expression in which this problem arises. (This program, and the others below, are based on the Pascal implementation of series, see Section 5.) The sequence expression in `NormalizedMax` creates a sequence `X` of the numbers in a file `Data`. It then creates a normalized sequence by dividing each number by the sum of the numbers. Finally, the procedure returns the maximum of the normalized elements. (The procedure `Series` creates a sequence indefinitely repeating the value of its argument. The call on `MapFn` divides each element of `X` by this value.)

```
function NormalizedMax (Data: FileOfReal): Real;
  var X: series of Real;
begin
  X := ScanFile(Data);
  NormalizedMax := CollectMax(MapFn(/, X, Series(CollectSum(X))))
end
```

The two uses of X place contradictory constraints on the way pipelined evaluation must proceed. The procedure `CollectSum` requires all the elements of X to be produced before the sum can be returned and `Series` requires that its input be available before it can start producing its output. However, `MapFn` requires that the first element of X be available at the same time as the first element of the output of `Series`. For pipelining to work, this implies that the first element of the output of `Series` (and therefore the output of `CollectSum`) must be available before the second element of X is computed. Unfortunately, if X contains more than one element, this is impossible.

The essence of the inconsistency above is the cycle of constraints used in the argument. This in turn stems from a non-directed cycle in the data flow graph underlying the expression. Figure 2.1 shows the nodes in the sequence expression in `NormalizedMax` and the data flow between them. The nodes are represented by boxes and data flow is represented by arrows. Simple arrows indicate the flow of sequences and cross hatched arrows indicate the flow of other values.

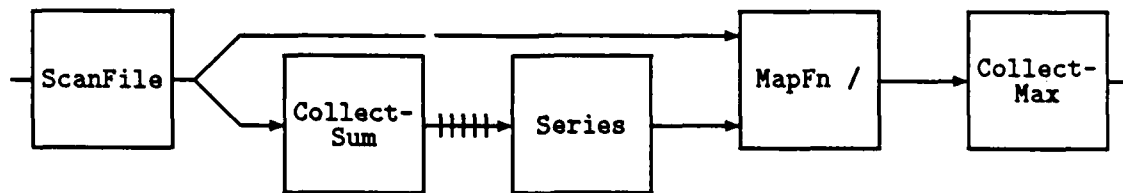


Figure 2.1: The sequence expression in `NormalizedMax`.

From the point of view of Figure 2.1, the problem in `NormalizedMax` can be summarized by noting that the non-directed data flow cycle has two conflicting parts. In the upper part, pipelining requires that each sequence element be used as soon as it is computed. In the lower part, the non-sequence data flow forces a delay—all of the elements at the left end of the lower part have to be available before any of the elements at the right end can be produced. If the upper part also contained a non-sequence data flow, then the delays on the two parts would balance. However, if there was a non-sequence data flow in the upper part, the expression would be broken into two separate sequence expressions: one on the left and one on the right.

Based on Definition 3, it can be shown that whenever a non-sequence value created by a node in a sequence expression is used by another node (either directly as in Figure 2.1 or via a chain of computation) the expression will be associated with a non-directed cycle like the one in Figure 2.1. To prevent this problem, it is required that intermediate values in optimizable sequence expressions must be sequences.

This restriction places significant limits on the qualitative character of optimizable sequence expressions. In particular, they all have the general form of creating some number of sequences, computing various intermediate sequences, and then computing one or more non-sequence results. A non-sequence value cannot be used in the intermediate computation unless it is the output of a disjoint expression.

The On-Line Cycle Restriction

The last situation that can block pipelining is illustrated in the program `OddSum` below. This program creates a sequence `X` of the numbers in a file. It then selects the odd elements of `X` and multiplies the i th odd element of `X` by the i th element of `X`. Finally, it sums the resulting products.

```
function OddSum (Data: FileOfReal): Real;
  var X: series of Real;
begin
  X := ScanFile(Data);
  OddSum := CollectSum(MapFn(*, X, ChooseIf(Odd, X)))
end
```

As in the program `NormalizedMax`, the two uses of `X` place contradictory constraints on the way pipelined evaluation must proceed. The key problem is that `ChooseIf` is inherently unsynchronized in the way it operates. In this example, it produces output elements only when it reads odd input elements. However, `MapFn` requires that the i th element of `X` be available at the same time as the i th element of the output of `ChooseIf`. For pipelining to work, this implies that the i th element of the output of `ChooseIf` must be available at the same time that `ChooseIf` reads the i th element of `X`. Unfortunately, if any of the first i elements of `X` are even, this is impossible.

As shown in Figure 2.2, the problem in `OddSum` is fundamentally much the same as the problem in `NormalizedMax`. In particular, it also stems from a non-directed cycle in the underlying data flow graph. Like the non-sequence data flow in Figure 2.1, the `ChooseIf` in Figure 2.2 introduces a delay that is not matched in the other part of the cycle. This delay depends on the input data, but may be arbitrarily large.

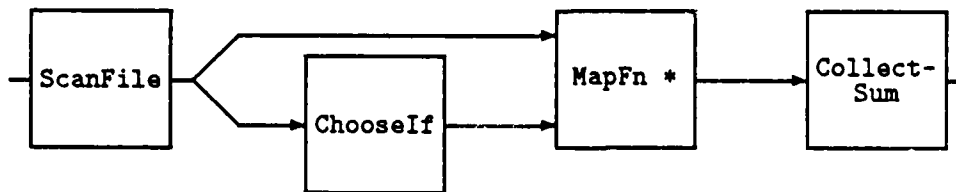


Figure 2.2: The sequence expression in `OddSum`.

In contrast to `OddSum`, consider the program `CosSum` below, which is identical except that `ChooseIf` is replaced by `MapFn` of `Cos`. A sum is computed of each element of `X` multiplied by its cosine.

```
function CosSum (Data: FileOfReal): Real;
  var X: series of Real;
begin
  X := ScanFile(Data);
  CosSum := CollectSum(MapFn(*, X, MapFn(Cos, X)))
end
```

Even though `CosSum` has exactly the same data flow graph as `OddSum` (with `ChooseIf` replaced by `MapFn` of `Cos`), `CosSum` can be pipelined without difficulty. The reason is that the `MapFn` of `Cos` produces an output element every time it reads an input element. Therefore, there is no problem synchronizing the arrival of the elements of the two sequence inputs of the `MapFn` of `*`.

The comparison of `OddSum` and `CosSum` shows that pipelining is blocked in `OddSum` by the delay introduced by `ChooseIf`, rather than by the mere existence of a non-directed data flow cycle. To develop a good restriction that rules out this problem, one has to develop a vocabulary for talking about where delays are introduced.

Definition 5 (On-Line and Off-Line) *An input or output port of a sequence procedure is on-line if and only if it operates in lock step with all the other on-line ports of the procedure as follows: The initial element of each on-line input is read, then the initial element of each on-line output is written, then the second element of each on-line input is read, then the second element of each on-line output is written, and so on. If all of the sequence ports of a procedure are on-line, the procedure as a whole is on-line. A non-directed cycle of data flow is on-line if every port it passes through is on-line. A non-directed cycle passes through an input or output port of a procedure call if and only if exactly one data flow arc in the cycle touches the port. (For example, the cycle in Figure 2.2 passes through every port it touches except for the output of `ScanFile`.) If a port, procedure, or data flow cycle is not on-line, it is off-line.*

Definition 5 extends the standard definition of the term 'on-line' [1, 22] so that it applies to individual ports as well as whole procedures. Like Definition 4, Definition 5 applies to procedures, rather than functions. While, some mathematical functions (e.g., choosing elements from a sequence) cannot be implemented in an on-line fashion, many can. For example, since the i th element of the result of mapping a function is computed solely from the i th elements of the inputs, it is easy for `MapFn` to be on-line.

Returning to the issue of pipelining, it can be shown that if a non-directed data flow cycle in an optimizable sequence expression is on-line, the lock step processing of the ports involved guarantees that there will not be any conflicts between the constraints associated with the cycle. If every non-directed cycle in an expression is on-line, the evaluation of the expression can be pipelined using the following divide and conquer approach.

It can be shown that when a sequence expression in which every non-directed data flow cycle is on-line contains an off-line port, it is always possible to divide the expression into two non-overlapping subexpressions so that all data flow between the subexpressions originates (or terminates) on the off-line port in question. The expression as a whole can be pipelined by pipelining the evaluation of the two subexpressions separately and using a simplified form of lazy evaluation to interleave the evaluation of the subexpressions in a pipelined fashion. The lazy evaluation is simplified because the method for determining which subexpression to evaluate when is very simple. The first subexpression needs to be evaluated when the second subexpression needs to read a new value computed by it.

Once partitioning based on off-line ports has been applied as many times as possible, one is left with subexpressions where every data flow connects on-line ports. In such a

subexpression, pipelining can be achieved by simply evaluating every procedure call in lock step, one element at a time.

The limits imposed by the on-line cycle restriction are softened by the fact that a wide range of mathematical sequence functions can be implemented in an on-line way. For instance, any preorder procedure that has only one sequence port is trivially on-line. Beyond this, most of the functions in Section 3 can be implemented as on-line procedures.

Obeying the Restrictions

In the current implementations of series (see the discussion in the following sections), the preorder restriction is implicitly enforced by ensuring that every predefined series procedure is preorder and not providing any means for defining series procedures that are not preorder. (Note that the composition of two preorder procedures is preorder.)

The other four restrictions are explicitly checked. Whenever an expression satisfies these restrictions, the algorithms in Section 6 are used to transform the expression into an efficient loop. When they are not satisfied, a warning is issued. In the Pascal implementation, these warnings are fatal errors. However, in the Lisp implementation, the expressions are simply left as is and evaluated/compiled without optimization.

The best approach for programmers to take is to write expressions without worrying about the restrictions and then fix the expressions in the event that a problem is reported. The virtues of this approach are enhanced by the fact that simple expressions are very unlikely to violate any of the restrictions. In particular, it can be shown that if every sequence procedure in an expression has only one output and sequence outputs are not stored in variables, then the sequence intermediate value and on-line cycle restrictions cannot be violated.

Similarly, it can be shown that violations of the sequence intermediate value and on-line cycle restrictions can always be fixed using code copying to break the expression in two or to break the offending cycle. For instance, the program `NormalizedMax` can be brought into compliance with the sequence intermediate value restriction by duplicating the call on `ScanFile`. This breaks the sequence expression into two separate sequence expressions that each satisfy the sequence intermediate value restriction.

```
function NormalizedMaxA (Data: FileOfReal): Real;
  var Sum: Real;
begin
  Sum := CollectSum(ScanFile(Data));
  NormalizedMax := CollectMax(MapFn(/, ScanFile(Data), Series(Sum)))
end
```

It would be possible to automatically introduce code copying to resolve conflicts with the sequence intermediate value and on-line cycle restrictions. However, this can lead to significant inefficiencies. It is better to leave it up to the programmer to figure out how to fix conflicts. For example, the procedure `NormalizedMax` can be brought into compliance more efficiently, by realizing that the operations of computing the maximum and dividing by the sum commute.

```

function NormalizedMaxB (Data: FileOfReal): Real;
  var X: series of Real;
begin
  X := ScanFile(Data);
  NormalizedMax := CollectMax(X)/CollectSum(X)
end

```

This example brings up an important secondary goal underlying the restrictions presented here. This goal is to make it easy for programmers to reliably tell which expressions are efficient and which are not. It would be better to make every expression efficient. However, given that this is not possible, programmers need accurate information in order to decide what to do.

Other Approaches to Restrictions

The optimizability restrictions presented above are the result of research going back twelve years. The basic concept of representing common looping subcomputations as operations on sequences is descended from ideas developed in the Programmer's Apprentice project [33, 41]. The first attempt to state a formal set of restrictions appears in [42, 43]. A set of restrictions intermediate between the initial ones and the ones presented above appear in [44].

Optimizability restrictions are implicit in all the work on sequence expression optimizers. However, these restrictions are typically implicit in the way the optimizers work, rather than being explicitly stated. The only other research featuring explicit restrictions is that of Wadler [40].

The key difference between Wadler's restrictions and the ones presented here is that he assumes that procedures can only have one output and only considers the situation where two procedures are composed together. This can be straightforwardly generalized to expressions that are tree-like in form—i.e., ones where the output of each function is only used once. However, it is not applicable to more complex situations.

Wadler's implicit requirement that expressions be trees rules out non-directed data flow cycles. This obviates the need for anything like the sequence intermediate value or on-line cycle restrictions. However, it is unreasonably limiting. Since it is often possible to pipeline the evaluation of an expression even though it contains non-directed cycles of data flow, it is unreasonable (both from the point of view of readability and efficiency) to require that an intermediate value that is used in n places must always be computed n times.

The restriction that Wadler explicitly states (i.e., that procedures must be *preorder listless*) is basically equivalent to the preorder restriction stated here, except that it also limits the storage used, as discussed below. (Wadler's definition of the term preorder is different from the one used here.)

Another significant difference between Wadler's work and the work presented here is that he approaches the problem of efficiency from a different direction. Rather than considering pipelining directly, he focuses on the amount of storage required when evaluating an expression. He requires that preorder listless procedures evaluate using a bounded amount of storage (over and above the storage required for the inputs and outputs themselves) and shows that the composition of two such functions can be evaluated using a

bounded amount of storage (over and above the storage required for the net inputs and outputs of the composition).

It can be shown that when all the procedures called by an optimizable sequence expression use bounded storage, the pipelined evaluation of the expression will use bounded storage. However, it was decided not to introduce a restriction limiting sequence procedures to bounded storage for three reasons. First, such a restriction has nothing to do with pipeline *per se*. Pipelining can be applied just the same no matter how much storage the individual procedures use. The key thing is that pipelining allows the storage requirements of an expression to be the sum of the storage requirements of the procedures called by it. Second, when considering individual sequence procedures a requirement that they use bounded storage is too weak to be useful. As discussed in the section on the preorder restriction, one certainly does not want unbounded storage to be used to store input or output elements. (However, what if it is used for some other purpose?) In addition, large fixed storage needs are just as much of a practical problem as unbounded ones. What is really needed is for the implementor of the procedure to make a good faith effort to use as little storage as possible. Third, such a restriction cannot be usefully supported, because there is no practical way to determine whether a user-defined sequence procedure requires unbounded storage.

A final difference between Wadler's work and the work presented here is that Wadler chose to apply his restrictions to a pre-existing data type (lists). In the approach taken here, a new data type (series) is developed for three reasons. First, since lists cannot represent unbounded sequences, focusing on lists unduly limits the kind of procedures that can be expressed. Second, lists (and vectors) have evolved a style of use and a suite of associated operations that are appropriate for that use. These work well for their intended use, but are not as useful as they could be from the perspective of writing efficient sequence expressions. Developing a new data type makes it possible to create a new suite of operations that is more appropriate for this purpose. Third, an important part of the approach being advocated here is the error messages that report unoptimizable series expressions. These are important if programmers are to achieve efficiency. However, they would be irritating and counterproductive if they were constantly being reported for list expressions that were not intended to be optimizable. By adding a series data type, programmers can benefit from the restrictions when they choose to follow them, without being prevented from using lists, vectors, and streams in standard ways.

3. Series Expressions

Vectors, lists, and other such data structures differ in how closely they model the mathematical concept of a sequence and in the range of associated operations. The series data type embodies a set of design decisions that combine full support for the mathematical concept of a sequence, a wide range of operations, and high efficiency. As illustrated in the next two sections, support for series can be straightforwardly added to any programming language. High efficiency is obtained by using a preprocessor or compiler extension like the one presented in Section 6.

Informally speaking, series are like vectors except that they can be unbounded in length. Formally, series are defined by the way they can be operated on. The remainder of this section presents an illustrative selection of these operations as mathematical functions. (See [9] for an in depth discussion of the mathematical properties of many of these functions when applied to finite sequences.) The next two sections present procedures implementing these functions in specific programming languages.

In the following, lowercase letters (x, y) denote arbitrary values, uppercase letters (R, S) denote series, and calligraphic letters (\mathcal{F}, \mathcal{P}) denote functions and predicates. In addition, special notations are used to denote four of the most basic series functions. The construction function, which creates a series containing the indicated elements in the indicated order is denoted using angle brackets, i.e., $\langle x, y, \dots, z \rangle$. The concatenating function, which creates a series containing the elements of a series R followed by the elements of another series S is denoted using the operator “ \parallel ”, i.e., $R \parallel S$. The tail function, which creates a series containing the elements of a non-empty series S after the first one is denoted by drawing a line over the series, i.e., \overline{S} . The head function, which returns the first element of a non-empty series is denoted by subscripting the series with zero, i.e., S_0 . (For convenience, the series used as examples below contain numerical elements. However, any kind of object can be used as a series element.)

$$\begin{aligned} \langle 6, 7, 8 \rangle \parallel \langle 9, 10 \rangle &= \langle 6, 7, 8, 9, 10 \rangle \\ \overline{\langle 6, 7, 8 \rangle} &= \langle 7, 8 \rangle \\ \langle 6, 7, 8 \rangle_0 &= 6 \\ \langle 10, \langle 4, 5 \rangle_0, 20 \rangle &= \langle 10, 4, 20 \rangle \end{aligned}$$

Series functions can be divided into three categories: *scanners* produce series without consuming any, *collectors* compute non-series values from series, and *transducers* compute series from series. For instance, the head function is a collector and the concatenating function is a transducer.

There are two kinds of scanners. Some scanners create a series of the elements in an aggregate data structure, for instance, a series of the nodes in a tree. Other scanners create a series based on some formula, for instance, the successive powers of some number. The Lisp implementation of series supports 15 scanners. However, one of these (scanning, see below) can be used to define all the rest.

Scanning is a higher-order function—a function that takes functions as arguments. The first argument is an initial value z , which becomes the first element of the series created. The second argument is a stepping function \mathcal{F} , which is used to compute each

series element from the previous element. The third argument is a predicate \mathcal{P} , which is used to determine where the series should end. The series created contains the elements z , $\mathcal{F}(z)$, $\mathcal{F}(\mathcal{F}(z))$, and so on, up to but not including the first value satisfying \mathcal{P} . If no value satisfies \mathcal{P} , the series is unbounded in length. The example computes the powers of 2 that are less than 100.

$$\text{scanning}(z, \mathcal{F}, \mathcal{P}) = \begin{cases} \langle \rangle & \text{if } \mathcal{P}(z) \\ \langle z \rangle \parallel \text{scanning}(\mathcal{F}(z), \mathcal{F}, \mathcal{P}) & \text{otherwise} \end{cases}$$

$$\text{scanning}(1, \lambda x. x+x, \lambda x. x \geq 100) = \langle 1, 2, 4, 8, 16, 32, 64 \rangle$$

There are also two kinds of collectors. Some collectors create an aggregate data structure containing the elements of a series, for instance, a hash table containing the series elements. Other collectors create a summary value computed by some formula from the elements of a series, for instance, their sum. The Lisp implementation of series supports 18 scanners. However, one of these (the higher-order series function collection shown below) can be used to define all the rest.

Collection uses a binary function \mathcal{F} to combine the elements of a series S together. The combination process begins with an initial value z . Typically, z is chosen to be a left identity of \mathcal{F} . The example computes the sum of a series.

$$\text{collection}(z, \mathcal{F}, S) = \begin{cases} z & \text{if } S = \langle \rangle \\ \text{collection}(\mathcal{F}(z, S_0), \mathcal{F}, \overline{S}) & \text{otherwise} \end{cases}$$

$$\text{collection}(0, \lambda xy. x+y, \langle 1, 2, 3 \rangle) = 6$$

Transducers are more complex than scanners or collectors. In particular, there is no one transducer that serves as a basis for the rest. Nevertheless, four key higher-order transducers support wide classes of common transduction operations.

Collecting is the same as collection except that it returns a series of partial results, rather than just a final value. The length of the output is the same as the length of S . The example computes a series of partial sums.

$$\text{collecting}(z, \mathcal{F}, S) = \begin{cases} \langle \rangle & \text{if } S = \langle \rangle \\ \langle \mathcal{F}(z, S_0) \rangle \parallel \text{collecting}(\mathcal{F}(z, S_0), \mathcal{F}, \overline{S}) & \text{otherwise} \end{cases}$$

$$\text{collecting}(0, \lambda xy. x+y, \langle 1, 2, 3 \rangle) = \langle 1, 3, 6 \rangle$$

By far the most commonly used series function is mapping, which maps a function \mathcal{F} over some number of series producing a series of the results. Each element of the output is computed by applying \mathcal{F} to the corresponding elements of the inputs. The length of the output is the same as the length of the shortest input. The example adds the corresponding elements in two series.

$$\text{mapping}(\mathcal{F}, S^1, \dots, S^n) = \begin{cases} \langle \rangle & \text{if any } S^i = \langle \rangle \\ \langle \mathcal{F}(S_0^1, \dots, S_0^n) \rangle \parallel \text{mapping}(\mathcal{F}, \overline{S^1}, \dots, \overline{S^n}) & \text{otherwise} \end{cases}$$

$$\text{mapping}(\lambda xy. x+y, \langle 1, 2, 3 \rangle, \langle 4, 5, 6, 7 \rangle) = \langle 5, 7, 9 \rangle$$

Truncating cuts off a series by testing each element with a predicate and discarding all the remaining elements as soon as an element satisfying the predicate is encountered. The example truncates a series at the first negative element.

$$\text{truncating}(\mathcal{P}, S) = \begin{cases} \langle \rangle & \text{if } S = \langle \rangle \text{ or } \mathcal{P}(S_0) \\ \langle S_0 \rangle \parallel \text{truncating}(\mathcal{P}, \bar{S}) & \text{otherwise} \end{cases}$$

$$\text{truncating}(\lambda x. x < 0, \langle 0, 3, 2, -7, 1, -1 \rangle) = \langle 0, 3, 2 \rangle$$

Mingling combines two series into one under the control of a comparison predicate. The comparison is performed as indicated to ensure that the combination will be stable—if two elements are considered equal by the comparison predicate, the element from R will precede the element from S in the result. The example shows the combination of two sorted series into a sorted result.

$$\text{mingling}(R, S, \mathcal{P}) = \begin{cases} R & \text{if } S = \langle \rangle \\ S & \text{if } R = \langle \rangle \\ \langle R_0 \rangle \parallel \text{mingling}(\bar{R}, S, \mathcal{P}) & \text{if } \neg \mathcal{P}(S_0, R_0) \\ \langle S_0 \rangle \parallel \text{mingling}(R, \bar{S}, \mathcal{P}) & \text{otherwise} \end{cases}$$

$$\text{mingling}(\langle 1, 3, 7 \rangle, \langle 2, 4, 5 \rangle, \lambda xy. x < y) = \langle 1, 2, 3, 4, 5, 7 \rangle$$

Choosing selects the elements of a series that satisfy a predicate. The example picks out the negative elements of the input.

$$\text{choosing}(\mathcal{P}, S) = \begin{cases} \langle \rangle & \text{if } S = \langle \rangle \\ \langle S_0 \rangle \parallel \text{choosing}(\mathcal{P}, \bar{S}) & \text{if } \mathcal{P}(S_0) \\ \text{choosing}(\mathcal{P}, \bar{S}) & \text{otherwise} \end{cases}$$

$$\text{choosing}(\lambda x. x < 0, \langle 0, 3, 2, -7, 1, -1 \rangle) = \langle -7, -1 \rangle$$

In addition to the higher-order transducers above, some specific transducers are important as well. The function spreading is a quasi-inverse of choosing. Spreading takes a series of non-negative integers R and a series of values S , and creates a series containing the elements of S . In the output, the elements of S are spread out by interspersing them with copies of z . If the i th element of R is n , then the i th element of S is preceded by n copies of z . Taken together with the example above, the example below illustrates the relationship between choosing and spreading.

$$\text{spreading}(R, S, z) = \begin{cases} \langle \rangle & \text{if } R = \langle \rangle \text{ or } S = \langle \rangle \\ \langle S_0 \rangle \parallel \text{spreading}(\bar{R}, \bar{S}, z) & \text{if } R_0 = 0 \\ \langle z \rangle \parallel \text{spreading}(\langle R_0 - 1 \rangle \parallel \bar{R}, S, z) & \text{otherwise} \end{cases}$$

$$\text{spreading}(\langle 3, 1 \rangle, \langle -7, -1 \rangle, 0) = \langle 0, 0, 0, -7, 0, -1 \rangle$$

Subseries is a generalization of the tail function. It creates a series of the elements of its input from the n th up to but not including the m th. The first element in a series has

the index 0. If m is greater than or equal to the length of S , output stops as soon as the input runs out of elements. The example takes a chunk out of the middle of a series.

$$\text{subseries}(S, n, m) = \begin{cases} \langle \rangle & \text{if } m = 0 \text{ or } S = \langle \rangle \\ \text{subseries}(\overline{S}, n-1, m-1) & \text{if } n > 0 \\ \langle S_0 \rangle \parallel \text{subseries}(\overline{S}, 0, m-1) & \text{otherwise} \end{cases}$$

$$\text{subseries}(\langle 1, 1, 2, 2, 3, 3, 4, 4 \rangle, 2, 5) = \langle 2, 2, 3 \rangle$$

The function `chunk` is different from the ones above, because it can produce more than one output. It has the effect of breaking a series S into (possibly overlapping) chunks of width $m > 0$. Successive chunks are displaced $n > 0$ elements to the right, in the manner of a moving window. `chunk` produces m output series. The i th chunk is composed of the i th elements of the m outputs. Suppose that the length of S is l . The length of each output is $\lfloor 1 + (l-m)/n \rfloor$. By itself, `chunk` may appear somewhat unusual, however, it is quite useful in combination with other transducers. For instance, the example shows how `chunk` could be used as part of the computation of a moving average. (Programming languages differ in the mechanisms that could be used to channel the outputs of `chunk` to the inputs of `mapping`.)

$$\begin{aligned} \text{chunk}(m, n, S) &= R^1, \dots, R^m \text{ where} \\ R^k &= \text{chunk}(1, n, \text{subseries}(S, k-1, \infty)) \\ \text{chunk}(1, n, S) &= \begin{cases} \langle \rangle & \text{if } S = \langle \rangle \\ \langle S_0 \rangle \parallel \text{chunk}(1, n, \text{subseries}(S, n, \infty)) & \text{otherwise} \end{cases} \\ \text{chunk}(2, 1, \langle 1, 5, 3, 7 \rangle) &= \langle 1, 5, 3 \rangle, \langle 5, 3, 7 \rangle \\ \text{mapping}(\lambda xy. (x+y)/2, \langle 1, 5, 3 \rangle, \langle 5, 3, 7 \rangle) &= \langle 3, 4, 5 \rangle \end{aligned}$$

The list of functions above can be viewed as a recommendation for the kind of functions that can profitably be supported in conjunction with a sequence data type. As discussed in Section 7, some languages (e.g., APL) support most of these functions; others (e.g., Pascal) support almost none of them.

The list of functions is also interesting for what it does not contain. To start with, it does not contain functions for accessing arbitrary series elements or altering the value of series elements. This reflects the fact that, unlike vectors or lists, series are not intended to be used as mutable data storage.

In addition, the choice of functions is significantly influenced by the optimizability restrictions in the last section. The list only includes functions that can be implemented as preorder procedures using small fixed amounts of internal storage. (The only common functions ruled out by this criteria are ones like reversal, rotation, and sorting that rearrange the order of the elements of a sequence.) The list also favors functions that can be implemented as on-line procedures, because these are more useful in optimizable expressions. (The only functions in the list that require off-line implementation are concatenating, tail, mingling, choosing, spreading, subseries, and `chunk`.)

4. A Common Lisp Implementation

Series can be added into essentially any programming language by adding an implementation of the series data structure and defining a set of procedures supporting the series functions in Section 3. The optimization of series expressions can be supported by a preprocessor (see Section 6). It is in the nature of Lisp, that both of these things are easy to do using a macro package. Such a macro package has been in regular use for a number of years and is generally available (see [46, 47]).

Series. In the Lisp implementation, series are implemented lazily using closures. A series has a procedural part and a data part. The procedural part is a generator [18, 29] capable of computing the elements one by one. The data part records the elements computed so far.

The elements of a series are accessed using a second generator that enumerates the elements in the data part of the series and then uses the procedural part to compute more elements as needed. Each time the generator is called, it returns another element in the series. The generator takes a procedure argument that specifies what to do when the series runs out of elements.

The two-level generation scheme above ensures that: elements are not computed until needed, no element is computed twice, and each user of a series can access all

```
(setq first5          ; Implementation of <1,2,3,4,5>.
      (let ((x 0))
        (list #'(lambda (at-end)
                  (if (< x 5) (setq x (+ x 1)) (funcall at-end))))))

(defun generator (s)  ; Returns a generator for the elements of a series.
  (let ((g (car s)))
    #'(lambda (at-end)
        (when (null (cdr s))
          (setf (cdr s)
                (block nil
                  (list (funcall g #'(lambda () (return T)))))))
        (if (not (eq (cdr s) T))
            (car (setq s (cdr s)))
            (funcall at-end)))))

(defun choose-if (p s) ; Implementation of choosing( $\mathcal{P}$ ,  $S$ ).
  (let ((gen (generator s)))
    (list #'(lambda (end-action)
              (loop (let ((x (funcall gen end-action)))
                    (if (funcall p x) (return x)))))))

(defun collect-sum (s) ; Implementation of collection(0,  $\lambda xy.x+y$ ,  $S$ ).
  (let ((gen (generator s)))
    (sum 0)
    (loop (let ((x (funcall gen #'(lambda () (return sum)))))
          (setq sum (+ sum x)))))

(collect-sum (choose-if #'oddp first5))  $\Rightarrow$  9
```

Figure 4.1: Illustration of the Lisp implementation of unoptimized series.

the elements. For those familiar with Lisp, Figure 4.1 illustrates the implementation of series data structures. The same basic implementation approach is used in the language Seque [19].

The closure implementation of series is effective and straightforward; However, it is not very efficient. No effort has been expended on producing a more efficient implementation, because the focus of series expressions is on the situations where they can be optimized, eliminating the physical representation of series altogether. In situations where optimization is impossible, it is usually better to represent a sequence as a vector or list than as a series.

The protocol above for obtaining a generator for the elements of a series and thence the elements themselves is not an exported part of the series implementation. Users must manipulate series using the series procedures below. This is important in the interest of optimizability in general and static analyzability in particular.

Series procedures. The series functions described in Section 3 are all supported by Lisp procedures as shown in Figure 4.2. In addition, the # macro character syntax #Z(*x y ... z*) is provided for reading and printing literal series. The difference between `make-series` and #Z is that the arguments to #Z are implicitly quoted, while the arguments to `make-series` are evaluated one at a time as needed.

```
(catenate (make-series 1 (+ 2 2)) #Z(7 8)) => #Z(1 4 7 8)
(collect-first (choose-if #'oddp #Z(8 -7 6 -1))) => -7
(subseries (mingle #Z(1 5 9) #Z(2 6 8) #'<) 2 4) => #Z(5 6)
(multiple-value-bind (xs ys) (chunk 2 1 #Z(1 5 3 7))
  (map-fn T #'(lambda (x y) (/ (+ x y) 2)) xs ys)) => #Z(3 4 5)
```

The higher-order procedures implementing scanning, collecting, mapping, truncating, and collection are extended so that they can accept multiple series arguments and produce

Series Function	Lisp Implementation
$\langle x, y, \dots, z \rangle$	(<code>make-series</code> <i>x y ... z</i>)
S_0	(<code>collect-first</code> <i>S</i>)
\overline{S}	(<code>subseries</code> <i>S</i> 1)
$R \parallel S$	(<code>catenate</code> <i>R S</i>)
scanning(<i>z</i> , \mathcal{F} , \mathcal{P})	(<code>scan-fn</code> <i>type Z F P</i>)
collection(<i>z</i> , \mathcal{F} , <i>S</i>)	(<code>collect-fn</code> <i>type Z F S¹ ... Sⁿ</i>)
collecting(<i>z</i> , \mathcal{F} , <i>S</i>)	(<code>collecting-fn</code> <i>type Z F S¹ ... Sⁿ</i>)
mapping(\mathcal{F} , S^1, \dots, S^n)	(<code>map-fn</code> <i>type F S¹ ... Sⁿ</i>)
truncating(\mathcal{P} , <i>S</i>)	(<code>until-if</code> \mathcal{P} $S^1 \dots S^n$)
mingling(<i>R</i> , <i>S</i> , \mathcal{P})	(<code>mingle</code> <i>R S P</i>)
choosing(\mathcal{P} , <i>S</i>)	(<code>choose-if</code> \mathcal{P} <i>S</i>)
spreading(<i>R</i> , <i>S</i> , <i>z</i>)	(<code>spread</code> <i>R S z</i>)
subseries(<i>S</i> , <i>n</i> , <i>m</i>)	(<code>subseries</code> <i>S n m</i>)
chunk(<i>m</i> , <i>n</i> , <i>S</i>)	(<code>chunk</code> <i>m n S</i>)

Figure 4.2: Lisp support for series functions.

multiple values. (Series of tuples could be used to get the same effect in any given situation. However, using multiple series values is usually more convenient and almost always more efficient than using tuples.)

The examples below illustrate the Lisp procedure `scan-fn`, which supports scanning. In the second example, a two-valued stepping procedure is used and two series are returned (the unbounded series of natural numbers and a series of their partial sums). While scanning is in progress, two internal states are maintained. The stepping procedure must accept as many values as it returns. Each of these values is treated as a separate state variable.

```
(scan-fn 'list #'(lambda () '(a b c d)) #'cddr #'null)
  ⇒ #Z((a b c d) (c d))

(scan-fn '(values integer integer)
  #'(lambda () (values 1 1))
  #'(lambda (i sum)
      (setq i (+ i 1)) (values i (+ sum i))))
  ⇒ #Z(1 2 3 4 ...) and #Z(1 3 6 10 ...)
```

Three other features of `scan-fn` are worthy of note. First, a new first argument is introduced, which specifies the type (or types) of the values returned by the stepping procedure. Given the lack of typing information in Lisp, this argument is necessary to ensure that the number of arguments returned by the stepping procedure can be determined at compile time. Second, the initial value is replaced by a procedure that returns the initial values. This is convenient in situations where multiple initial values are needed. Third, the predicate argument is made optional. Omitting it is the same as supplying a predicate that is not true of any value. The first and second extensions are applied to `collect-fn`, `collecting-fn`, and `map-fn` as well as `scan-fn`.

As a convenience to the user, a number of specific scanners are provided in addition to `scan-fn`. These include: `series` which creates a series indefinitely repeating a given value, `scan` which enumerates the elements in a list, vector, or string, `scan-range` which enumerates the integers in a range, and `scan-plist` which creates a series of the indicators in a property list along with a second series containing the corresponding values. The first argument of `scan` specifies the type of object to be scanned. If omitted, the type defaults to `list`.

```
(series "test") ⇒ #Z("test" "test" "test" ...)
(scan '(a b c)) ⇒ #Z(a b c)
(scan 'vector '(a b c)) ⇒ #Z(a b c)
(scan 'string "Tuz") ⇒ #Z(#\T #\u #\z)
(scan-range :from 1 :upto 3) ⇒ #Z(1 2 3)
(scan-plist '(a 1 b 2)) ⇒ #Z(a b) and #Z(1 2)
```

Similarly, a number of specific collectors are provided including: `collect` which combines the elements of a series into a list, vector, or string, `collect-sum` which adds up the elements of a series, `collect-length` which returns the number of elements in a series, and `collect-last` which returns the last element of a series (or an optional default value if the series is empty). The first argument of `collect` specifies the type of object to be produced. If omitted, the type defaults to `list`.

```

(collect #Z(a b c)) ⇒ (a b c)
(collect 'simple-vector #Z(a b c)) ⇒ #(a b c)
(collect 'string #Z(#\T #\u #\z)) ⇒ "Tuz"
(collect-sum #Z(1 3 2)) ⇒ 6
(collect-length #Z("fee" "fi" "fo" "fum")) ⇒ 4
(collect-last #Z("fee" "fi" "fo" "fum")) ⇒ "fum"
(collect-last #Z() "none") ⇒ "none"

```

Finally, a number of additional transducers are provided including: `previous` (based on `map-fn`) which takes in a series and shifts it over one element by inserting the indicated value at the front and discarding the last element, `choose` (based on `choose-if`) which selects the elements of its second argument that correspond to non-null elements of its first argument, and `positions` (also based on `choose-if`) which returns the positions of the non-null elements in a series. If given only one argument, `choose` returns the non-null elements of this series.

```

(previous #Z("fee" "fi" "fo" "fum") " ") ⇒ #Z(" " "fee" "fi" "fo")
(choose #Z(T nil T) #Z(1 2 3)) ⇒ #Z(1 3)
(choose #Z(nil 3 4 nil)) ⇒ #Z(3 4)
(positions (map-fn #'oddp #Z(1 2 3 5 6 8))) ⇒ #Z(0 2 3)

```

Convenient support for mapping. In cognizance of the ubiquitous nature of mapping, the Lisp series implementation provides three mechanisms that make it easy to express particular kinds of mapping. The `#` macro character syntax `#M \mathcal{F}` converts a procedure \mathcal{F} into a transducer that maps \mathcal{F} .

```

(#Msqrt #Z(4 16)) ≡ (map-fn T #'sqrt #Z(4 16)) ⇒ #Z(2 4)

```

The form `mapping` can be used to specify the mapping of a complex computation over one or more series without having to write a literal lambda expression. It has the same basic syntax as `let`. For example,

```

(mapping ((x (scan '(2 -2 3))))
  (expt (abs x) 3)) ⇒ #Z(8 8 27)

```

is the same as

```

(map-fn T #'(lambda (x) (expt (abs x) 3))
  (scan '(2 -2 3))) ⇒ #Z(8 8 27)

```

The form `iterate` is the same as `mapping` except that the value `nil` is always returned.

```

(iterate ((x (scan '(2 -2 3))))
  (if (plusp x) (prin1 x))) ⇒ nil <after printing "23">

```

To a first approximation, `iterate` and `mapping` differ in the same way as `mapc` and `mapcar`. In particular, like `mapc`, `iterate` is intended to be used in situations where the body is being evaluated for side effect rather than for its result. However, due to the lazy evaluation nature of series, the difference between `iterate` and `mapping` is more than just a question of efficiency. If `mapping` is used in a situation where the output is not used, no computation is performed, because series elements are not computed until they are used.

Nested loops. The equivalent of a nested loop is expressed by simply using a series expression in a procedure that is mapped over a series. This is typically done using `mapping`. In the example, a list of sums is computed based on a list of lists of numbers.

```
(let ((data '((1 2 3) (4 5 6) (7 8))))
  (collect
   (mapping ((number-list (scan data)))
            (collect-sum (scan number-list)))) ⇒ (6 15 15)
```

User-defined series procedures. As shown by the definitions of `collect-sum` and `mapping` below, the standard Lisp forms `defun` and `defmacro` can be used to define new series procedures. However, the Series macro package must be informed when a series procedure is being defined with `defun`. This is done by using the declaration `optimizable-series-function`. No special declaration is required when using `defmacro`.

```
(defun collect-sum (numbers)
  (declare (optimizable-series-function))
  (collect-fn 'number #'(lambda () 0) #' + numbers))

(defmacro mapping (var-value-pair-list &body body)
  (let* ((pairs (scan var-value-pair-list))
        (arg-list (collect (#Mcar pairs)))
        (value-list (collect (#Mcdr pairs))))
    `(map-fn T #'(lambda ,arg-list ,@ body) ,@ value-list)))
```

Example. The following example shows what it is like to use series expressions in a realistic programming context. The example consists of two parts: a pair of procedures that convert between sets represented as lists and sets represented as bits packed into an integer and a graph algorithm that uses the integer representation of sets.

Sets over a small universe can be represented very efficiently as binary integers where each 1 bit in the integer represents an element in the set. Here, sets represented as binary integers are referred to as *bit sets*.

Common Lisp provides a number of bitwise operations on integers, which can be used to manipulate bit sets. In particular, `logior` computes the union of two bit sets while `logand` computes their intersection.

The procedures in Figure 4.3 convert between sets represented as lists and bit sets. To perform this conversion, a mapping has to be established between bit positions and potential set elements. This mapping is specified by a *universe*. A universe is a list of elements. If a bit set integer b is associated with a universe u , then the i th element in u

```
(defun bset->list (bset universe)
  (collect (choose (#Mlogbitp (scan-range :from 0) (series bset))
                 (scan universe))))

(defun list->bset (items universe)
  (collect-fn 'integer #'(lambda () 0) #'logior
    (mapping ((item (scan items))
             (ash 1 (bit-position item universe))))))

(defun bit-position (item universe)
  (or (collect-first (positions (#Meq (series item) (scan universe))))
      (1- (length (nconc universe (list item))))))
```

Figure 4.3: Converting between lists and bit sets.

```

(defun collect-logior (bsets)
  (declare (optimizable-series-function))
  (collect-fn 'integer #'(lambda () 0) #'logior bsets))

(defun collect-logand (bsets)
  (declare (optimizable-series-function))
  (collect-fn 'integer #'(lambda () -1) #'logand bsets))

```

Figure 4.4: Operations on series of bit sets.

is in the set represented by b if and only if the i th bit in b is 1. For example, given the universe (a b c d e), the integer #b01011 represents the set {a,b,d}. (By Common Lisp convention, the 0th bit in an integer is the rightmost bit.)

Given a bit set and its associated universe, the procedure `bset->list` converts the bit set into a set represented as a list of its elements. It does this by scanning the elements in the universe along with their positions and constructing a list of the elements that correspond to 1s in the integer representing the bit set. (When no `:upto` argument is supplied, `scan-range` counts up forever.)

The procedure `list->bset` converts a set represented as a list of its elements into a bit set. Its second argument is the universe that is to be associated with the bit set created. For each element of the list, the procedure `bit-position` is called to determine which bit position should be set to 1. The procedure `ash` is used to create an integer with the correct bit set to 1. The procedure `collect-fn` is used to combine the integers corresponding to the individual elements together into a bit set corresponding to the list.

The procedure `bit-position` takes an item and a universe and returns the bit position corresponding to the item. The procedure operates in one of two ways depending on whether or not the item is in the universe. The first line of the procedure contains a series expression that determines the position of the item in the universe. If the item is not in the universe, the expression returns `nil`. (The procedure `collect-first` returns `nil` if it is passed an empty series.)

If the item is not in the universe, the second line of the procedure adds the item onto the end of the universe and returns its position. The extension of the universe is done by side effect so that it will be permanently recorded in the universe.

Figure 4.4 shows the definition of two collectors that operate on series of bit sets. The first procedure computes the union of a series of bit sets, while the second computes the intersection.

Live variable analysis. As an illustration of the way bit sets might be used, consider the following. Suppose that in a compiler, program code is being represented as blocks of straight-line code connected by possibly cyclic control flow. The top part of Figure 4.5 shows the data structure that represents a block of code. Each block B has several pieces of information associated with it. Two of these pieces of information are the blocks that can branch to B and the blocks B can branch to. A program is represented as a list of blocks that point to each other through these fields.

In addition to control flow information, each structure contains information about the way variables are accessed. In particular, it records the variables that are written by

```

(defstruct (block (:conc-name nil))
  predecessors ;Blocks that can branch to this one.
  successors   ;Blocks this one can branch to.
  written      ;Variables written in the block.
  used         ;Variables read before written in the block.
  live        ;Variables that must be available at exit.
  temp        ;Temporary storage location.

(defun determine-live (program-graph)
  (let ((universe (list nil)))
    (convert-to-bsets program-graph universe)
    (perform-relaxation program-graph)
    (convert-from-bsets program-graph universe))
  program-graph)

(defstruct (temp-bsets (:conc-name bset-))
  used written live)

(defun convert-to-bsets (program-graph universe)
  (iterate ((block (scan program-graph)))
    (setf (temp block)
          (make-temp-bsets
            :used (list->bset (used block) universe)
            :written (list->bset (written block) universe)
            :live 0))))

(defun perform-relaxation (program-graph)
  (let ((to-do program-graph))
    (loop
      (when (null to-do) (return (values)))
      (let* ((block (pop to-do))
             (estimate (live-estimate block)))
        (when (not (= estimate (bset-live (temp block))))
          (setf (bset-live (temp block)) estimate)
          (iterate ((prev (scan (predecessors block)))
                    (pushnew prev to-do)))))))

(defun live-estimate (block)
  (collect-logior
    (mapping ((next (scan (successors block)))
              (logior (bset-used (temp next))
                      (logandc2 (bset-live (temp next))
                                (bset-written (temp next)))))))

(defun convert-from-bsets (program-graph universe)
  (iterate ((block (scan program-graph)))
    (setf (live block)
          (bset->list (bset-live (temp block)) universe))
    (setf (temp block) nil)))

```

Figure 4.5: Live variable analysis.

the block and the variables that are used by the block (i.e., either read without being written or read before they are written). An additional field (computed by the procedure `determine-live` discussed below) records the variables that are *live* at the end of the block. (A variable is live if it has to be saved, because it can potentially be used by a following block.) Finally, there is a temporary data field, which is used by procedures (such as `determine-live`) that perform computations involved with the blocks.

The remainder of Figure 4.5 shows the procedure `determine-live`, which given a program represented as a list of blocks, determines the variables that are live in each block. To perform this computation efficiently, the procedure uses bit sets. The procedure operates in three steps. The first step (`convert-to-bsets`) looks at each block and sets up an auxiliary data structure containing bit set representations for the written variables, the used variables, and an initial guess that there are no live variables. This auxiliary structure is defined by the third form in Figure 4.5 and is stored in the `temp` field of the block. The integer 0 represents an empty bit set.

The second step (`perform-relaxation`) determines which variables are live. This is done by relaxation. The initial guess that there are no live variables in any block is successively improved until the correct answer is obtained.

The third step (`convert-from-bsets`) operates in the reverse of the first step. Each block is inspected and the bit set representation of the live variables is converted into a list, which is stored in the `live` field of the block.

On each cycle of the loop in `perform-relaxation`, a block is examined to determine whether its live set has to be changed. To do this (see the procedure `live-estimate`), the successors of the block are inspected. Each successor needs to have available to it the variables it uses, plus the variables that are supposed to be live after it, minus the variables it writes. (The procedure `logandc2` takes the difference of two bit sets.) A new estimate of the total set of variables needed by the successors as a group is computed by using `collect-logior`.

If this new estimate is different from the current estimate of what variables are live, then the estimate is changed. In addition, if the estimate is changed, `perform-relaxation` has to make sure that all the predecessors of the current block will be examined to see whether the new estimate for the current block requires their live estimates to be changed. This is done by adding each predecessor onto the list `to-do` unless it is already there. As soon as the estimates of liveness stop changing, the computation stops.

Summary. Figure 4.5 is a particularly good example of the way series expressions are intended to be used in three ways. First, all the series expressions are optimizable. Second, series expressions are used in a number of places to express computations that would otherwise be expressed less clearly as loops or less efficiently using operations on lists or vectors. Third, the main relaxation algorithm in `perform-relaxation` is expressed as a loop. This is done, because the data flow in this algorithm prevents it from being decomposed into two or more fragments. This highlights the fact that optimizable series expressions are not intended to render iterative programs entirely obsolete, but rather to provide a greatly improved method for expressing the vast majority of loops.

5. A Pascal Implementation

Series can be added to Pascal in much the same way as they are added to Lisp. A prototype system has been constructed that demonstrates this [28, 45]. However, the prototype is written in Lisp rather than Pascal and only supports optimizable series expressions. A fatal error is issued whenever optimization is blocked. Although less complete than the approach of the Lisp implementation, this still allows loops to be replaced by optimizable series expressions.

Series. The Pascal series preprocessor supports the declaration of series in analogy with array declarations as shown below. (This is the only syntactic extension required to support series. Further, the output of the preprocessor is standard Pascal without any references to series.) In line with the general philosophy of Pascal, it is required that all the elements of a series have the same type. However, the length of a series is not part of its type. This is important to facilitate the definition of series procedures operating on series of arbitrary length.

```
type Integers = series of Integer;
var InputData: series of Real;
```

Series procedures. The series functions described in Section 3 are all supported by Pascal procedures as shown in Figure 5.1. In general, these have the same names as the corresponding Lisp procedures with any hyphens removed (e.g., scan-fn becomes ScanFn). Since Pascal does not support the concept of a function procedure that returns multiple values, the outputs of chunk are turned into arguments. (The examples in [28, 45] show an obsolete set of names linked to an earlier Lisp implementation of series.) The Pascal implementation does not provide a syntax for series literals. (The mathematical syntax is used in the examples below.)

Series Function	Pascal Implementation
$\langle x, y, \dots, z \rangle$	MakeSeries(x, y, \dots, z)
S_0	CollectFirst(S)
\overline{S}	Subseries($S, 1$)
$R \parallel S$	Catenate(R, S)
scanning($z, \mathcal{F}, \mathcal{P}$)	ScanFn($z, \mathcal{F}, \mathcal{P}$)
collection(z, \mathcal{F}, S)	CollectFn(z, \mathcal{F}, S)
collecting(z, \mathcal{F}, S)	CollectingFn(z, \mathcal{F}, S)
mapping($\mathcal{F}, S^1, \dots, S^n$)	MapFn($\mathcal{F}, S^1, \dots, S^n$)
truncating(\mathcal{P}, S)	TruncateIf(\mathcal{P}, S)
mingling(R, S, \mathcal{P})	Mingle(R, S, \mathcal{P})
choosing(\mathcal{P}, S)	ChooseIf(\mathcal{P}, S)
spreading(R, S, z)	Spread(R, S, z)
subseries(S, n, m)	Subseries(S, n, m)
chunk(m, n, S)	Chunk(m, n, S, R^1, \dots, R^m)

Figure 5.1: Pascal support for series functions.

```

Catenate(MakeSeries(1, 2+2), <7,8>) ⇒ <1,4,7,8>
CollectFirst(ChooseIf(odd, <8,-7,6,-1>)) ⇒ -7
Subseries(Mingle(<1,5,9>, <2,6,8>, <>) 2 4) ⇒ <5 6>
Chunk(2, 1, <1,5,3,7>, Xs, Ys) ⇒ Xs := <1,5,3> and Ys := <5,3,7>
MapFn(Average, Xs, Ys) ⇒ <3,4,5>
function Average (x,y: Integer): Integer;
begin
  Average := (x+y)/2
end;

```

The Pascal implementation does not extend the higher-order procedures over their specifications in Section 3 for two reasons. Given the strong typing in Pascal, the pre-processor can obtain type information without needing type arguments. Since, Pascal does not support the concept of multiple return values, some other method needs to be employed to avoid the need for tuples.

The procedures in Figure 5.1 do not follow the usual Pascal restrictions on the parameters of procedures. Some of the procedures allow the number of arguments they receive to vary and they all allow considerable flexibility in the types of their arguments. This is important because the series procedures are inherently generic in character. For instance, `MapFn` is naturally applicable to any number and any type of series as long as the element types are compatible with the procedure being mapped.

Due to their generic nature, the procedures in Figure 5.1 could not be implemented as user-defined procedures in Pascal. However, as an extension to the language, they do not violate the spirit of Pascal. In particular, the predeclared Pascal procedures are generic in exactly the same way. Several (e.g., `Read` and `Write`) allow variable numbers of arguments and most of them are applicable to more than one type of object. Using a more flexible language such as Ada [50], it would be possible to implement (at least most of) the higher-order series functions as user-defined procedures.

All of the specific scanners, collectors, and transducers from the Lisp implementation that are applicable to Pascal are supported by the Pascal implementation as well. Given the strong typing in Pascal, `Scan` and `Collect` do not need type arguments. Since Pascal has sets, but not lists, these functions apply to sets and not lists. In keeping with the general style of Pascal, `Collect` takes the destination vector/string/set as its first argument rather than returning an aggregate value.

```

Series('test') ⇒ <'test','test','test', ...>
Scan('Tuz') ⇒ <'T','u','z'>
Scan([Mon,Wed,Fri]) ⇒ <Mon,Wed,Fri>
ScanRange(1, 3) ⇒ <1,2,3>

Collect(X, <'T','u','z'>) { Places 'Tuz' in X. }
CollectSum(<1,3,2>) ⇒ 6
CollectLength(<'fee','fi','fo','fum'>) ⇒ 4
CollectLast(<'fee','fi','fo','fum'>) ⇒ 'fum'
CollectLast(<>, 'none') ⇒ 'none'

Previous(<'fee' 'fi' 'fo' 'fum'>, ' ') ⇒ <' ', 'fee', 'fi', 'fo'>
Choose(<true,false,true>, <1,2,3>) ⇒ <1,3>
Positions(MapFn(Odd, <1,2,3,5,6,8>)) ⇒ #Z(0 2 3)

```

Implicit mapping. To avoid making syntactic extensions, the Pascal implementation does not support constructs analogous to the Lisp forms `mapping` and `iterate`. However, it supports a related concept that is in many ways even more useful. Whenever a non-series procedure is applied to a series, it is automatically mapped over the elements of the series. For example, in the expression below, `Sqr` is automatically mapped over the series of numbers created by scanning the set.

```
CollectSum(Sqr(Scan([2,4])))
≡ CollectSum(MapFn(Sqr, Scan([2,4]))) ⇒ 20
```

The key virtue of implicit mapping is that it reduces the number of helping procedures that have to be defined. For instance, in the example of a moving average above, you can write the following instead of defining a procedure `Average` and explicitly mapping it.

```
Chunk(2, 1, (1,5,3,7), Xs, Ys); (Xs+Ys)/2 ⇒ (3,4,5)
```

The concept of implicit mapping is completely separate from the other concepts associated with series expressions. As such, it could easily be dispensed with. However, as shown by experience with APL and the other languages that support it, implicit mapping is extremely useful. (The lack of reliable compile-time type information makes it impractical to support implicit mapping in Lisp.)

User-defined series procedures. As shown in the examples below, series procedures in Pascal are simply procedures that either have series inputs or return series values. As with series in general, all such definitions are handled directly by the preprocessor. There is no need for any special kind of declaration. Pascal does not support the concept of macros.

Example. The following example illustrates how series expressions can best be used in Pascal. As in the last section, all of the expressions are optimizable. The example revolves around a job queue data abstraction that might be used in an operating system. The basic type definition is shown below. A `JobQ` is a pointer to a chain of entries that point to records describing jobs. These records have a number of fields including a numerical priority.

```
type JobQ = ↑JobQentry;
type JobQentry = record; Job: JobInfo; Rest: JobQ end;
type JobInfo = ↑JobRecord;
type JobRecord = record Priority: Real; ... end;
```

There are a number of procedures defined that operate on job queues. These procedures include putting a new job onto a queue (shown below) and removing a job from a queue (discussed near the end of this section). To add a job onto a queue, one merely needs to allocate a new queue entry and attach it to the front of the queue.

```
procedure AddToJobQ (J: JobInfo; var Q: JobQ);
  var E: ↑JobQentry;
begin
  new(E);
  E↑.Job := J;
  E↑.Rest := Q;
  Q := E
end
```

In addition to ordinary procedures that operate on job queues, it is useful to define a number of series procedures that operate on job queues. In particular, as with any aggregate data structure, it is useful to have procedures `ScanJobQ` and `CollectJobQ` that convert job queues to series of jobs and vice versa. It also turns out to be useful to have a procedure `ScanJobQtails` that enumerates all of the tails of a queue (i.e., $\langle Q, Q\uparrow.\text{Rest}, Q\uparrow.\text{Rest}\uparrow.\text{Rest}, \dots \rangle$). As shown below, `ScanJobQtails` can be implemented using the higher-order series procedure `ScanFn` and two special-purpose procedures operating on job queues.

```
function ScanJobQtails (Q: JobQ): series of JobQ;
  function JobQrest (Q: JobQ): JobQ;
    begin JobQrest := Q↑.Rest end;
  function JobQnull (Q: JobQ): Boolean;
    begin JobQnull := Q=nil end;
begin
  ScanJobQtails := ScanFn(Q, JobQrest, JobQnull)
end
```

Among other things, `ScanJobQtails` can be used to implement `ScanJobQ` as shown below. The expression `Qs↑.Job` causes the operations of following a pointer and selecting the job field of a `JobQentry` to be implicitly mapped over the pointers returned by `ScanJobQtails`.

```
function ScanJobQ (Q: JobQ): series of JobInfo;
  var Qs: series of JobQ;
begin
  Qs := ScanJobQtails(Q);
  ScanJobQ := Qs↑.Job
end
```

The procedure `RemoveFromJobQ` removes a job from the end of a queue. It can be implemented using `ScanJobQtails` as shown below. To start with, `RemoveFromJobQ` enumerates the tails of the queue and uses `CollectLast` and `Previous` to obtain pointers to the last and next to last entries in the queue. The job field of the last queue entry is returned as the result of `RemoveFromJobQ`. (It is assumed that there must be at least one job in the queue.) The rest pointer in the next to last entry is set to nil, in order to remove the last entry from the queue. (If there is no next to last entry, the queue variable itself is set to nil.) The storage associated with the last entry is then freed.

```
function RemoveFromJobQ (var Q: JobQ): JobInfo;
  var Qs: series of JobQ;
      NextToLast, Last: JobQ;
begin
  Qs := ScanJobQtails(Q);
  Last := CollectLast(Qs, nil);
  NextToLast := CollectLast(Previous(Qs, nil), nil);
  RemoveFromJobQ := Last↑.Job;
  if NextToLast=nil
    then Q := nil
    else NextToLast.Rest := nil;
  dispose(Last)
end
```

The first three lines of the body of `RemoveFromJobQ` are a series expression, while the remainder are non-series expressions. From an efficiency standpoint, it should be noted that since there is only one instance of `ScanJobQtails`, the series expression is converted into a loop that only traverses the queue once.

As a final example of the use of optimizable series expressions, consider the procedure `SuperJob` below. This procedure inspects a job queue and returns the last (i.e., longest queued) job in the queue whose priority is more than two standard deviations larger than the average priority of the jobs in the queue. If there is no such job, `nil` is returned. The first five statements in the procedure compute the mean and deviation of the priorities. The next two statements select the jobs that have sufficiently large priorities. The final line selects the last of these jobs, if any.

```
function SuperJob (Q: JobQ): JobInfo;
  var Jobs, SuperJobs: series of JobInfo;
      N: Integer;
      Mean, SecondMoment, Deviation, Limit: Real;
begin
  Jobs := ScanJobQ(Q);
  N := CollectLength(Jobs);
  Mean := CollectSum(Jobs.Priority)/N;
  SecondMoment := CollectSum(Sqr(Jobs.Priority))/N;
  Deviation := Sqrt(SecondMoment-Sqr(Mean));
  Limit := Mean+2*Deviation;
  SuperJobs := Choose(Jobs.Priority>Limit, ScanJobQ(Q));
  SuperJob := CollectLast(SuperJobs, nil)
end
```

The programs above are a good example of the way series expressions are intended to be used. To start with, all of the programs are straightforward in nature. This reflects the fact that the primary goal of series expressions is to convert the vast majority of programs that are in fact straightforward programs into dirt simple programs. When a program is straightforward, it is usually easy to write it in a loop-free form without having to use anything other than very simple series expressions.

6. The Optimization Algorithms

A preprocessor or compiler extension that transforms optimizable series expressions into loops can be implemented in three stages: *parsing* which locates optimizable expressions and converts them into equivalent data flow graphs, *pipelining* which converts the expressions into loops, and *unparsing* which converts the resulting loops into appropriate program code and inserts this code in place of the original expressions.

Below, the Pascal implementation of series is used as a concrete illustration. The Common Lisp implementation works in the same way, except that the characteristics of Lisp simplify the parsing and unparsing stages.

Parsing. When the preprocessor is applied to a program, it begins by parsing the program. Series expressions are located by inspecting the types of the procedures called by the program. While this is being done, the static analyzability and straight-line computation restrictions are checked and any violations reported.

In a language like Pascal where complete compile-time type information is available, implicit mapping can be supported by noting places where non-series functions are applied to series. Each such application is replaced by an appropriate use of `MapFn`.

The final action of the parsing stage is to create a data flow graph corresponding to each optimizable series expression located. Since each of these expressions is a straight-line computation, this is easy to do. Each procedure call becomes a node in the graph and the data flow between the nodes is derived from the way procedure calls are nested and the way variables are used.

Pipelining. The operation of the pipelining stage is illustrated in Figure 6.1. The series expression in the procedure `SumSqr5` (which computes the sum of the squares of the odd elements of a vector) is transformed into the loop shown in the procedure `SumSqr5Pipelined`. The readability of the loop code is reduced by the fact that it contains a number of internally generated variables. However, the code is quite efficient. The only significant problem is that the pipeliner sometimes uses more variables than strictly necessary (e.g., `Result5`). However, this need not lead to inefficiency during execution as long as a compiler capable of simple optimizations is available.

The pipelining process operates in several steps. In the first step, the divide and conquer strategy discussed in Section 2 is used to partition the data flow graph for a series expression into subexpressions where all the data flow connects on-line ports. While doing this, the pipeliner checks that the expression obeys the sequence intermediate value and on-line cycle restrictions.

Once partitioning is complete, the procedures in each subexpression are combined into a single procedure. The resulting procedures are then combined based on the data flow between the subexpressions. To support the combination process, each series procedure is represented as a loop fragment with one or more of the following parts:

```

function SumSqrS (V: array [1..N] of Integer): Integer;
begin
  SumSqrS := CollectSum(Sqr(ChooseIf(Odd, Scan(V))))
end

```

⇓

```

function SumSqrSPipelined (V: array [1..N] of Integer): Integer;
  label 0,1;
  Var Element12, Index15, Result5, Sum2: Integer;
begin
[1]   Index15 := 0;
[4]   Sum2 := 0;
[1] 1: Index15 := 1+Index15;
[1]   if Index15>N then goto 0;
[1]   Element12 := V[Index15];
[2]   if Odd(Element12) then goto 1;
[3]   Result5 := Sqr(Element12);
[4]   Sum2 := Sum2+Result5;
      goto 1;
0: SumSqrSPipelined := Sum2
end

```

```

[1] -- Scan of a vector -----
      inputs- Vector: array [K..L] of ElementType;
      outputs- Element: Series of ElementType;
      vars- Index: Integer;
      prolog- Index: 1-K;
      body- Index := 1+Index;
           if Index>L then goto 0;
           Element := Vector[Index];

[2] -- ChooseIf -----
      inputs- function P(X: ElementType): Boolean;
           Item: Series of ElementType;
      outputs- Item: Series of ElementType;
      labels- 2;
      body- 2: NextIn(Item); if P(Item) then goto 2;

[3] -- Implicit mapping of Sqr -----
      inputs- Item: Series of ElementType;
      outputs- Result: Series of ElementType;
      body- Result := Sqr(Item);

[4] -- CollectSum -----
      inputs- Number: Series of ElementType;
      outputs- Sum: ElementType;
      prolog- Sum := 0;
      body- Sum := Sum+Number;

```

Figure 6.1: Transforming optimizable series expressions into loops.

- inputs-** Input variables.
- outputs-** Output variables.
- vars-** Auxiliary variables used by the computation.
- labels-** Labels used by the computation.
- prolog-** Statements that are executed before the computation starts.
- body-** Statements that are repetitively executed.
- epilog-** Statements that are executed after the loop terminates.

The bottom part of Figure 6.1 shows the fragments that represent the procedures called by the series expression in `SumSqrs`. These fragments are combined to create the loop in `SumSqrsPipelined`. The numbers in the left hand margin indicate which fragment each line of the loop comes from. Two different combination algorithms are used: one corresponding to data flow between on-line ports and one corresponding to data flow touching off-line ports.

When two procedures are connected by data flow connecting on-line ports (e.g, the data flow from the output of `Sqrs` to the input of `CollectSum`), the procedures are combined by simply concatenating the various parts of the corresponding fragments together. In addition, the variables and labels in the fragments are renamed so that there will be no possibility of conflicts. The data flow between the procedures is implemented by renaming the input variable of the destination so that it is the same as the output variable of the source. (The process above is much the same as an application of the standard compiler optimization technique of loop fusion [3].)

When two procedures are connected by series data flow terminating on an off-line input (e.g., the data flow from the output of `Scan` to the series input of `ChooseIf` in the figure), the fragment representing the destination procedure contains an instance of the form `NextIn`, which specifies when elements of the input should be computed. The two fragments are combined exactly as in the on-line combination algorithm except that the body of the source fragment is substituted in place of the call on `NextIn`, rather than being concatenated with the body of the destination fragment. (This process essentially compiles in support for a simple case of lazy evaluation [16].)

Unparsing. The result of pipelining is a single loop fragment that corresponds to the series expression as a whole. In the unparsing stage, this fragment is converted into a loop as indicated below. The combination process eliminates the inputs. The other parts of the fragment appear directly in the loop except for the outputs.

```

label 0,1,labels;
  var vars;
begin
  prolog;
1: body;
  goto 1;
0: epilog;

```

The outputs are connected up to the surrounding code when the loop is substituted into the program in place of the original series expression. Once each series expression has been replaced by a loop, the resulting code can be passed to a standard Pascal compiler.

Side-effects. The correctness preserving nature of the transformations above has been shown under the assumption that there are no side-effects involved. It is believed that if unoptimized series are implemented as shown in the beginning of Section 4, the transformations are also correctness preserving in the presence of side-effects. The reason for this is that the transformed code exactly mimics the lazy evaluation of the untransformed expression. For instance, the off-line port combination algorithm involves code motion; however, this motion simply moves the generation of the series elements to the place where lazy evaluation requires the elements of the series to be first computed.

The above can be made more formal from the point of view of path analysis [10]. Path analysis seeks to determine at compile time where in a program each lazy value will be first used and where it will be reused. This information can be used to optimize lazy evaluation in two ways. If there is an identifiable place of first use of a given value then ordinary evaluation can be used instead of lazy evaluation for that value. If there is an identifiable last use for a value, the value does not have to be stored beyond that time.

The restrictions in Section 2 guarantee that for each series, there is an identifiable place where each element of the series is first used and that, for each element, the last use precedes the computation of the next element. The transformations above merely position the computation of the elements at their place of first use and omit their long term storage.

The above notwithstanding, one should realize that side-effects are still problematical, because lazy evaluation makes it difficult for programmers to figure out what the net results of side-effects will be. Some situations can be readily understood. For example, one can depend on the fact that mapping will apply the mapped function \mathcal{F} first to the first element of the input, then to the second, and so on in strict temporal order. Thus, if \mathcal{F} interacts with itself or the environment outside of the containing series expression X via side-effects (e.g., by doing input or output), but does not interact with anything else in X , the result is easy enough to understand. More complex uses of side-effects should be avoided.

Systems Based on Similar Algorithms

The algorithms above have evolved into their current form over the past twelve years. The first generally available implementation was a Lisp macro package called LetS [42, 43]. The current Lisp implementation [46] is available in Portable Common Lisp.

The same basic approach to representing and combining sequence procedures was independently developed by Wile [48]. However, he does not explicitly address the question of restrictions and his approach does not guarantee that every intermediate sequence can be eliminated. Much the same can be said about APL compilers [11].

A quite similar approach is also used internally by the Loop macro [12]. However, as discussed in the next section, the Loop macro is externally very different from series expressions. In particular, it uses an idiosyncratic English-like syntax rather than representing computations as compositions of procedures operating on series.

7. Comparisons

There are two primary vantage points from which to compare series expressions with related concepts. The most obvious comparison is with other support for sequence expressions. From this perspective, the key feature of series expressions is that they support most of the operations supported by the vector operations of APL [31], Common Lisp sequence operations [36], and the stream operations of Seque [19] (along with a few additional ones) while being more efficient.

Another way to view series expressions is that they are a logical continuation of the trend in programming language design toward supporting the reuse of loop fragments. From this point of view, series expressions extend the approach taken by iterators in CLU [27] and the Lisp Loop macro [12]. The key feature of series expressions in this context is that they support the reuse of a wider variety of fragments and are easier to understand and modify, without being any less efficient.

To lend depth to the comments above, the remainder of this section presents detailed comparisons between the Common Lisp implementation of series expressions and five other systems. Each of these comparisons features the example below. This example shows the definition of a procedure that computes the sum of the positive elements of a vector. It also illustrates how a new series procedure can be defined. (The example assumes the existence of a procedure `Positive` that tests whether an integer is greater than zero.)

```
function SumPositive (V: array [1..N] of Integer): Integer;
begin
  SumPositive := CollectSum(ChooseIf(Positive, Scan(V)))
end

function CollectSum (S: series of Integer): Integer;
begin
  CollectSum := CollectFn(0, +, S)
end
```

Other Support for Sequence Expressions

There are many programming languages that provide support for sequence expressions, e.g., [5, 6, 19, 25, 32, 34, 35, 36]. So many, that it would not be practical to make detailed comparisons between each one of these languages and series expressions. Three representative languages are discussed below.

APL. One of the oldest and most used languages that takes a functional approach is APL [23, 31]. A style of writing APL has evolved where vector expressions are used instead of loops. The correspondence between the series functions discussed in Section 3 and the APL vector operators is summarized in Figure 7.1. The APL concept of the extension of scalar operations to vectors corresponds to implicit mapping.

As illustrated below, both the vector summation algorithm and user-defined sequence procedures can be very compactly represented in APL. Since each sequence is directly represented as a vector, there is no need for an operation analogous to `Scan`.

```

[1] ▽ SUM←SUMPOSITIVEAPL V
    SUM←COLLECTSUM((V>0)/V)
    ▽
[1] ▽ SUM←COLLECTSUM NUMBERS
    SUM←+/NUMBERS
    ▽

```

The key differences between APL and series expressions are that APL vectors cannot represent unbounded sequences, the set of APL operations is somewhat different, and users are not given any feedback about what is efficient and what is not.

Although all the series operations could have been supported in APL, there is no direct built-in support for scanning, mingling, or chunk. In addition, APL does not support higher-order operations as well as it might initially appear. For instance, the reduction operator appears to be a higher-order operation. However, at least in standard APL, the operation to be reduced must be one of the predefined scalar operations—user-defined operations cannot be used. As a result, the reduction operator is actually just part of a naming scheme for a small set of specific collectors. (This has the collateral benefit of allowing the initial identity value to be implicit.)

APL supports four operations (reversal, membership, grade up, and grade down) that are not supported by series expressions because they cannot be implemented in a preorder fashion using bounded storage. APL also allows the modification of the elements of a sequence. When using series expressions, one has to rely on other constructs in the host language when performing any of these operations. For instance, to sort a series in the Lisp implementation of series, one must first collect the series into a list and then sort the list. Explicitly creating a list makes the expensive nature of sorting more obvious,

Series Function	APL Operation	name
$\langle x, y, \dots, z \rangle$	x, y, \dots, z	scalar catenation
S_0	$S[1]$	indexing first element
\overline{S}	$1 \downarrow S$	dropping first element
$R \parallel S$	R, S	catenation
scanning($z, \mathcal{F}, \mathcal{P}$)	missing	
collection(z, \mathcal{F}, S)	\mathcal{F} / S	reduction
collecting(z, \mathcal{F}, S)	$\mathcal{F} \setminus S$	scanning
mapping($\mathcal{F}, S^1, \dots, S^n$)	$S^1 \mathcal{F} S^2$	extension of scalar operations
truncating(\mathcal{P}, S)	$((\mathcal{P} S) \iota 1) - 1 \uparrow S$	take
mingling(R, S, \mathcal{P})	missing	
choosing(\mathcal{P}, S)	$(\mathcal{P} S) / S$	compression
spreading(R, S, z)	idiom based on	expansion
subseries(S, n, m)	$(m-n) \uparrow n \downarrow S$	take and drop
chunk(m, n, S)	missing	
simple idioms	index generation, membership, inner product, etc.	
missing	reversal, rotation, grade up/down, modifying elements	

Figure 7.1: The correspondence between series functions and APL operations.

however, it does not make sorting more expensive, because sorting requires that some physical representation of the sequence be created.

A few APL compilers [11, 21] are capable of producing efficient code in most of the situations where series expressions can be optimized; however, most are not. As a result, optimizable series expressions are typically much more efficient. Further, even when the compiler supports optimization, programmers are not given any feedback about whether optimization will occur. Rather, programmers are (at least implicitly) encouraged not to think about such issues.

An area where APL is fundamentally more powerful than series expressions is that the standard intermediate structure in APL is the array. APL has a number of powerful array operators (not shown in Figure 7.1) and a few APL compilers can optimize some array expressions. In contrast, while it is possible to have series of series, there are no special series operations for operating on them, and they are never optimized.

Finally, a superficial but striking difference between series expressions and APL is that series expressions use standard subroutine calling notation while APL uses a special set of concise, but cryptic, operators.

Common Lisp Sequence Operations. While many (if not most) Lisp programmers use loops extensively, a style of writing Lisp has evolved where expressions computing intermediate lists and vectors are used instead of loops. Unfortunately, until recently, Lisp supported an impoverished set of predefined sequence operations—it supported `mapcar`, but not much else. When Common Lisp was designed [36], this defect was rectified by introducing a relatively comprehensive suite of sequence operations.

In Common Lisp, the term 'sequence' is used to refer to either a list or a vector. How-

Series Function	Sequence Operation
$\langle x, y, \dots, z \rangle$	(list $x y \dots z$)
S_0	(elt 0 S)
\overline{S}	(subseq S 1)
$R \parallel S$	(concatenate type $R S$)
scanning($z, \mathcal{F}, \mathcal{P}$)	missing
collection(z, \mathcal{F}, S)	(reduce $\mathcal{F} S$)
collecting(z, \mathcal{F}, S)	missing
mapping($\mathcal{F}, S^1, \dots, S^n$)	(map type $\mathcal{F} S^1 \dots S^n$)
truncating(\mathcal{P}, S)	(subseq S 0 (position-if $\mathcal{P} S$))
mingling(R, S, \mathcal{P})	(merge type $R S \mathcal{P}$)
choosing(\mathcal{P}, S)	(remove-if-not $\mathcal{P} S$)
spreading(R, S, z)	missing
subseries(S, n, m)	(subseq S $n m$)
chunk(m, n, S)	missing
simple idioms	elt, length, count, find, some, etc.
missing	reverse, sort, modifying elements

Figure 7.2: The correspondence between series functions and sequence operations.

ever, since both of these structures are limited to representing bounded sequences, Lisp sequences are not a complete implementation of mathematical sequences. The correspondence between the basic series functions and the Lisp sequence operations is summarized in Figure 7.2. Lisp does not support implicit mapping.

The example below shows how the Common Lisp sequence operations can be used to express the vector summation algorithm and a user-defined sequence procedure. Since a vector is a Lisp sequence, there is no need for a procedure analogous to `scan`.

```
(defun sum-positive-sequence (v)
  (collect-sum (remove-if-not #'plusp v)))
(defun collect-sum (numbers)
  (reduce #'+ numbers))
```

Except for the fact that Lisp sequences cannot represent unbounded sequences, there is no reason why all the series functions could not be supported by sequence operations. However, there is no direct built-in support for scanning, collecting, spreading, or chunk. In addition, the identity value to use for `reduce` (collection) is specified in an odd way. The procedure argument must be implemented in such a way that when called with zero arguments it returns the identity value.

Current Lisp compilers do not optimize sequence expressions. As a result, optimizable series expressions are much more efficient. In light of the lack of optimization, it is not surprising that Lisp provides no feedback about optimizability. As in APL, there is no bias toward preorder functions and modification of sequence elements is allowed. It is also common to have sequences of sequences, however, Lisp does not provide any special operations for manipulating them.

An interesting aspect of the Lisp sequence operations is that they typically support a number of keyword arguments that modify their behaviors. For example, consider the sequence operation `count-if`, which takes a predicate and a sequence, and returns a count of the number of elements in the sequence that satisfy the predicate.

```
(count-if #'plusp '(1 -2 3 4 -5)) ⇒ 3
```

The Lisp operation `count-if` takes two keyword arguments `:start` and `:end`, which can be used to specify a subsequence of the input in which counting is to occur. In addition, a keyword argument `:key` can be used to specify an access procedure that will be used to fetch the part of each sequence element that should be tested by the predicate. Finally, an operation `count-if-not` exists, which is the same as `count-if` except that it automatically negates the values returned by the predicate.

As illustrated by the example below, none of these options is strictly necessary. The `:start` and `:end` keywords can be dispensed with by using `subseq`. The `:key` keyword and `count-if-not` can be dispensed with by specifying complex predicates.

```
(count-if-not #'plusp '((1) (-2) (3) (4) (-5))
  :start 0 :end 3 :key #'first)
≡ (count-if #'(lambda (element) (not (plusp (car element))))
  (subseq '((1) (-2) (3) (4) (-5)) 0 3)) ⇒ 1
```

Nevertheless, the various options described above are important for two reasons. First, they promote efficiency. (Using `subseq` instead of the `:start` and `:end` keywords is inefficient, because it creates an intermediate sequence.) Second, they increase the probability that predefined operations can be used as procedure arguments instead of `lambda` expressions. This makes uses of `count-if` more concise and easier to read.

Using series expressions, neither of these issues comes up. In the Lisp series expression below, the use of `subseries` does not lead to inefficiency, since pipelining eliminates the physical creation of its output series. Convenient support for mapping makes it possible to avoid the need for an explicit `lambda` expression. The desired test and key is simply mapped over the series in question (again without inefficiency). Finally, `count-if` itself can be dispensed with by using a combination of `choose` and `collect-length`. The approach taken by series expressions allows the individual procedures to be simpler and makes things more functional in appearance.

```
(let ((elements (subseries (scan '((1) (-2) (3) (4) (-5))) 0 3)))
  (collect-length (choose (#Mnot (#Mplusp (#Mcar elements)))))) ⇒ 1
```

Seque. Under the name of streams, sequences are the central data type of the language Seque [19]. Using the same basic lazy evaluation technique discussed in the beginning of Section 4, Seque supports both bounded and unbounded sequences. The correspondence between the series functions discussed in Section 3 and the stream operations provided by Seque is summarized in Figure 7.3. As in APL, many of these operations are provided by means of special syntax. In addition, implicit mapping is supported for many non-stream operations when they are applied to streams.

Series Function	Seque Operation	name
$\langle x, y, \dots, z \rangle$	$\{x, y, \dots, z\}$	sequence of expressions
S_0	$S ! 1$	referencing first element
\overline{S}	$S \% 1$	pre-truncation
$R \parallel S$	$R \rightarrow S$	concatenation
$\text{scanning}(z, \mathcal{F}, \mathcal{P})$	<i>idioms based on</i>	generators
$\text{collection}(z, \mathcal{F}, S)$	$\text{Red}(S, \mathcal{F}) ! \text{Length}(S)$	reduction
$\text{collecting}(z, \mathcal{F}, S)$	$\text{Red}(S, \mathcal{F})$	reduction
$\text{mapping}(\mathcal{F}, S^1, \dots, S^n)$	$[\mathcal{F}(S^1 ! i, \dots, S^n ! i)]$	derived stream
$\text{truncating}(\mathcal{P}, S)$	$S \setminus \setminus [\text{if } \mathcal{P}(S ! i) \text{ then } i-1] ! 1$	post-truncation
$\text{mingling}(R, S, \mathcal{P})$	<i>missing</i>	
$\text{choosing}(\mathcal{P}, S)$	$[\text{if } \mathcal{P}(S ! i) \text{ then } S ! i]$	filtering
$\text{spreading}(R, S, z)$	<i>missing</i>	
$\text{subseries}(S, n, m)$	$S\{n-1, m-2\}$	sectioning
$\text{chunk}(m, n, S)$	<i>missing</i>	
<i>simple idioms</i>	Length , reference, operations over streams, etc.	
<i>missing</i>	modifying elements	

Figure 7.3: The correspondence between series functions and Seque operations.

As illustrated below, both the vector summation algorithm and user-defined sequence procedures can be easily represented in Seque. Since streams are a distinct data type from vectors, an operation ! equivalent to Scan is required.

```

procedure SumPositive (V)
  return CollectSum([!V: lambda(e) if e>0 then e])
end

procedure CollectSum (S)
  L := Length(S)
  return if L=0 then 0 else Red(S, "+")!L
end

```

Since unbounded sequences are supported, it would be easy to completely support all the series functions in Seque. However, there is no direct support for mingling, spreading, or chunk. In addition, collection is only indirectly supported by collecting, this leads to awkwardness when collection is applied to an empty sequence, because there is no specification of the correct default value to return. There is also no direct support for the higher-order function scanning. However, there is an impressive array of facilities for defining scanning functions, both in Seque and in the language Icon [18], which Seque is based on. It is interesting to note that like the series operations, all of Seque's stream operations are preorder.

Seque does not attempt to optimize the evaluation of stream expressions by eliminating the computation of unnecessary intermediate streams. As a result, series expressions are never less efficient and often much more efficient. Seque programmer's are encouraged to think in terms of streams of streams and to make use of assigning to the elements of streams without any regard for the consequences on efficiency.

Summary. In comparison with the languages above, series expressions have three principal advantages. They support a wider range of operations than any one of the languages. Except in comparison with the best of APL compilers, they are much more efficient. They give clear feedback about what is efficient and what is not. As part of this, they make the use of non-preorder procedures more awkward, by forcing the programmer to use the facilities of the host language.

Looping Notations

The fundamental virtues of series expressions in comparison with looping constructs are illustrated by the discussion in the beginning of Section 1. However, this discussion is colored by the fact that it illustrates the use of only the most basic kind of looping construct. More complex looping constructs support several of the features of series expressions. In particular, they allow the equivalent of scanners and collectors (but not transducers) to be expressed as localized forms rather than as a statements dispersed in a loop.

Most programming languages contain a **for** construct, which makes it easy to express loops that are based on enumerating a range of integers (i.e., they support a standard loop fragment analogous to **scan-range**). Some languages go beyond this by providing special looping forms corresponding to a few additional scanners. For example, Common Lisp provides a form **dolist** that makes it easy to implement a loop based on scanning

the elements of a list. A couple of languages go further still by supporting a relatively wide range of standard looping fragments. Some of the oldest and most comprehensive support for this is in Lisp.

The Lisp Loop macro. The Lisp Loop macro [12] (which is based on the iterative statements in the InterLisp Clisp facility [37]) introduces two concepts into Lisp. First, it supports a looping construct analogous to `for` that uses an Algol-like keyword syntax. Second, it goes way beyond most `for` constructs by supporting a wide range of looping fragments analogous to scanners and collectors. Because of its non-Lisp syntax, the Loop macro has always been controversial. However, because of the utility of its predefined looping fragments, it has gained wide use.

The example below shows a program that uses the Loop macro to implement the vector summation algorithm. It also shows how a keyword `vector-element` (which corresponds to scanning a vector) could be defined. (The Loop macro does not support the definition of new collector-like fragments.) The code produced by the Loop macro is more or less identical to the code produced when optimizable series expressions are pipelined. As a result, the two approaches are equally efficient.

```
(defun sum-positive-loop-macro (v)
  (loop for item being each vector-element of v
        when (plusp item)
          sum item))

(define-loop-path vector-element scan-vector (of))

(defun scan-vector (path-name variable data-type prep-phrases
                  inclusive? allowed-prepositions data)
  (declare (ignore path-name data-type inclusive?
                  allowed-prepositions data))
  (let ((vector (gensym))
        (i (gensym))
        (end (gensym)))
    '(((,vector) (,i 0) (,end) (,variable))
      ((setq ,vector ,(cadar prep-phrases))
       (setq ,end (- (length ,vector) 1)))
      (> ,i ,end)
      (,variable (aref ,vector ,i))
      nil
      (,i (+ ,i 1))))))
```

Loop supports a keyword `when` that is similar to the transducer choosing (see the example). A call on the Loop macro can also contain an arbitrary body that is mapped over the values computed by the scanner-like fragments (this is not shown in the example). However, the Loop macro does not support any other transducer-like looping fragments.

A subtle difficulty with the Loop macro is that there are no restrictions on the computation that can be in the body and there is no attempt to prevent the body from interfering with the computation specified by the looping fragments. As a result, programmers cannot depend on the fact that these fragments will necessarily do what they are intended to do.

Another problem with the Loop macro is that the facilities provided for defining the equivalent of new scanners are quite cumbersome. The user has to define a procedure that

deals with parsing parts of the Loop syntax and that returns a list of six parts analogous to the parts of the loop fragments discussed in Section 6. (Recently, the Common Lisp standardization committee decided to adopt most of the Loop macro as part of Common Lisp. However, on the grounds that it is too complex, they decided not to include the scanner-defining form.)

The concept of Generators and Gatherers presented in [29], provides essentially the same capabilities as the Loop Macro, but with a more functional syntax and simpler defining forms.

Iterators in CLU. Among Algol-like languages, some of the most powerful support for the use of looping fragments is provided by CLU [27]. In CLU, scanner-like fragments called *iterators* can be used in `for` loops to generate a series of elements that are processed by the body of the `for`. CLU provides a number of predefined scanners including one corresponding to scanning a vector and users can define new ones. (Alphard [49] supports a construct called a generator that is essentially identical to a CLU iterator.)

As an illustration, the example below shows how the vector summation algorithm can be expressed in CLU. It also shows the definition of a user-defined scanner. This is done by writing a coroutine that yields the scanned elements one at a time.

```
SUM_POSITIVE_CLU = proc(V: ARRAY[INT]) returns(INT)
  SUM: INT := 0
  for ITEM: INT in SCAN_VECTOR(V) do
    if ITEM>0 then SUM := SUM+ITEM end
  end
  return(SUM)
end SUM_POSITIVE_CLU

SCAN_VECTOR = iter(V: ARRAY[INT]) yields(INT)
  I: INT := ARRAY[INT]$LOW(V)
  END: INT := ARRAY[INT]$HIGH(V)
  while I<=END do
    yield(V[I])
    I := I+1
  end
end SCAN_VECTOR
```

Taken together, CLU iterators and the `for` statement are essentially the same as the Loop macro except for three things. Nothing besides mapping and scanning is supported. (In the example above, the operations of choosing and summing are both represented in non-local ways in the body of the loop.) Each `for` can only contain one iterator instead of many. CLU's method for defining iterators as coroutines is significantly easier to use than the Loop macro's scanner-defining form.

A method for supporting multiple iterators in a CLU `for` statement is described in [14]. Going beyond this, [13] describes how one could support collectors (again restricted to only one in each loop). While both of these papers merely present proposals rather than describing actual implementations, there is no doubt that everything supported by the Lisp Loop macro could be straightforwardly supported in an Algol-like language.

Summary. The key difference between the looping constructs above and series expressions is that while the looping constructs support looping fragments corresponding to

(potentially unbounded) scanners and collectors and are highly efficient, they do not support the concept of a sequence data structure nor the idea of treating loop fragments as procedures. This preserves the iterative feel of the constructs, however, it is significantly limiting in several ways.

The lack of a sequence data type prevents the constructs from supporting anything other than a few simple transducers. (It is not clear how one could support transducers in general without having some kind of object that they can act upon.)

The fact that the loop fragments are not procedures means that the way the fragments can be used is intimately tied up with the syntax of the constructs. One has to learn a new language of combination rather than simply using standard functional composition. In addition, this new language of combination is much more restricted than functional composition. For instance, the only thing that can be done with a scanner-like fragment is to use it in a call on `loop` or `for` and map some computation over the values scanned.

An interesting thing to note about the looping constructs above is the way they avoid getting involved with a discussion of the restrictions in Section 2. By not allowing a sequence data structure to be stored in a variable, only supporting preorder fragments, largely ignoring transducers (particularly off-line ones), and limiting the way fragments can be combined, the constructs implicitly enforce these restrictions without having to talk about them. Unfortunately, the total restrictions they embody are much stronger than the ones in Section 2. This unnecessarily limits what can be expressed.

8. Benefits

There are three principal perspectives from which to view series expressions. To start with, series expressions can be looked at as embodying relatively complete support for sequence expressions in such a way that they can be included in any programming language without: removing any preexisting features of the language, requiring the use of unusual syntax, or causing inefficiency. This support includes most of the operations provided by languages such as APL, Lisp, and Seque along with a few additional ones.

An alternate perspective focuses on the fact that programmers are given clear and immediate feedback about the efficiency of the series expressions they write. Series expressions that do not violate the restrictions in Section 2 are guaranteed to be as efficient as they look. By means of error messages, programmers are encouraged to think of efficient methods for computing the results they want. In particular, unlike APL, Lisp, or Seque, programmers are never tempted to think that all sequence expressions are equally efficient.

A final perspective is summarized by the statement that "optimizable series expressions are to loops as structured control constructs are to `gotos`." By using optimizable series expressions, it is possible to banish loops from most programs. Given that expressions are much easier to understand and modify than loops, this has the potential for being a step forward at least as important as banishing `gotos`.

While the idea has not yet been explored, optimizable series expressions might also be helpful in the context of parallelism. Even though it is optimized for sequential machines, the pipelining applied to optimizable series expressions is very much the same as the 'software pipelining' of loops for execution on very large instruction word machines [26]

and for execution by the processors of a systolic array [2, 20]. If programs were written using series expressions, the process of analyzing the programs to determine a good schedule for the pipelined computation might be simplified. In addition, the restrictions in Section 2 appear relevant, because buffering of elements also causes inefficiency in a parallel context.

The application of optimizable series expressions to non-pipelined parallelism is less clear. The emphasis in such situations is on locating opportunities for evaluating sub-computations completely in parallel with no data flow between them. This is appropriate for mapping, but not for most of the other series operations. Nevertheless, using optimizable series expressions might make it easier to detect where such parallelism exists. For example, this might make it easier to 'vectorize' [4] programs.

9. Bibliography

- [1] A. Aho, J. Hopcraft, and J. Ullman, *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading MA, 1974.
- [2] A. Aiken and A. Nicolau, "Optimal Loop Parallelization", *Proc. of the SIGPLAN '88 Conf. on Programming Language Design and Implementation*, :308-317, Atlanta GA, June 1988.
- [3] F. Allen and J. Cocke, "A Catalogue of Optimizing Transformations" in *Design and Optimization of Compilers*, R. Rustin (ed.), Prentice Hall, 1971.
- [4] R. Allen and K. Kennedy, "Automatic Translation of Fortran Programs to Vector Form", *ACM Transactions on Programming Languages and Systems*, 9(4):491-542, October 1987.
- [5] J. Backus, "Can Programming be Liberated from the Von Neuman Style? A Functional Style and Its Algebra of Programs", *Comm. of the ACM*, 21(8), Aug. 1978.
- [6] D. Barstow, "Automatic Programming for Streams", *Proc. of the 9th Int. Joint Conf. on Artificial Intelligence*, :232-237, Aug. 1985.
- [7] F. Bellegarde, "Rewriting Systems on FP Expressions That Reduce the Number of Sequences They Yield", *Proc. ACM Symp. on Lisp and Functional Programming*, :63-73, Aug. 1984.
- [8] F. Bellegarde, "Convergent Term Rewriting Systems Can Be Used for Program Transformation", *Proc. Workshop on Programs as Data Objects*, LNCS 217, :24-41, H. Ganziger and N. Jones (eds), Springer-Verlag, 1985
- [9] R. Bird, "An Introduction to the Theory of Lists", :5-42, in *Logic of Programming and Calculi of Discrete Design*, M. Broy (ed.), NATO ASI series, Springer Verlag, 1986.
- [10] A. Bloss, P. Hudak, and J. Young, "Code Optimizations for Lazy Evaluation", *Lisp and Symbolic Computation*, 1(2)147-164, September 1988.
- [11] T. Budd, *An APL Compiler*, Springer-Verlag, New York NY, 1988.
- [12] G. Burke and D. Moon, *Loop Iteration Macro*, MIT/LCS/TM-169, July 1980.
- [13] R. Cameron, "Efficient High-Level Iteration with Accumulators", *ACM Trans. on Programming Languages and Systems*, 11(2)194-211, February 1989.
- [14] J. Eckart, "Iteration and Abstract Data Types", *ACM SIGPLAN Notices*, 22(4):103-110, April 1987.
- [15] J. Emery. "Small-Scale Software Components", *ACM SIGSOFT Software Engineering Notes*, 4(4):18-21, Oct. 1979.
- [16] D. Friedman and D. Wise, *CONS Should Not Evaluate Its Arguments*, Indiana Tech. Rep. 44, Nov. 1975.
- [17] A. Goldberg and R. Paige, *Stream Processing*, Rutgers report LCSR-TR-46, Aug. 1983.
- [18] R. Griswold and M. Griswold, *The Icon Programming Language*, Prentice-Hall, Englewood Cliffs NJ, 1983.
- [19] R. Griswold and J. O'Bagy, "Seque: A Programming Language for Manipulating Sequences", *Computer Languages*, 13(1)13-22, January 1988.
- [20] T. Gross and A. Sussman, "Mapping a Single-Assignment Language onto the Warp Systolic Array", in *Functional Programming Languages and Computer Architecture*, LNCS 274, :347-363, G. Goos and J. Hartmanis (eds.), Springer-Verlag, 1987.
- [21] L. Guibas and D. Wyatt, "Compilation and Delayed Evaluation in APL", *Proc. 1978 ACM Conf. on the Principles of Programming Languages*, Sept. 1978.
- [22] J. Hartmanis, P. Lewis, and R. Stearns, "Classification of computations by time and memory requirements", *proc. IFIP Congress 65*, :31-35, Spartan Books, Washington DC,

- 1965.
- [23] K. Iverson, "Operators", *ACM Trans. on Programming Languages and Systems*, 1(2):161-176, Oct. 1979.
 - [24] K. Jensen and N. Wirth, *Pascal User Manual and Report*, Springer-verlag, New York, 1985.
 - [25] G. Kahn and D. MacQueen, "Coroutines and Networks of Parallel Processes", *Proc. 1977 IFIP congress*, North-Holland, Amsterdam, 1977.
 - [26] M. Lam, "Software Pipelining: An Effective Scheduling Technique for VLIW Machines", *Proc. SIGPLAN '88 Conf. on Programming Language Design and Implementation*, :318-328, Atlanta GA, June 1988.
 - [27] B. Liskov, et. al., *CLU Reference Manual*, Lecture Notes in Computer Science, 114, G. Goos and J. Hartmanis eds., Springer-Verlag, New York, 1981.
 - [28] J. Orwant, *Support of Obviously Synchronizable Series Expressions in Pascal*, Massachusetts Institute of Technology, technical report AI-WP-312, September 1988.
 - [29] C. Perdue and R. Waters, "Generators and Gatherers", in *Common Lisp: the Language*, Second Edition, G. Steele Jr., Digital Press, Maynard MA, (in press).
 - [30] K. Pingali and Arvind, "Efficient Demand-Driven Evaluation. Part 1", *ACM Transactions on Programming Languages and Systems*, 7(2):311-333, April 1985.
 - [31] R. Polivka and S. Pakin, *APL: The Language and Its Usage*, Prentice-Hall, Englewood Cliffs NJ, 1975.
 - [32] N. Prywes, A. Pnueli, and S. Shastri, "Use of a Non-Procedural Specification Language and Associated Program Generator in Software Development", *ACM Trans. on Programming Languages and Systems*, 1(2), Oct. 1979.
 - [33] C. Rich and R. Waters, "The Programmer's Apprentice Project: A Research Overview", *IEEE Computer*, 21(11):10-25, November 1988.
 - [34] G. Ruth, S. Alter, and W. Martin, *A Very High Level Language for Business Data Processing*, MIT/LCS/TR-254, 1981.
 - [35] J. Schwartz et al, *Programming with Sets: An Introduction to SETL*. Springer-Verlag, 1986.
 - [36] G. Steele Jr., *Common Lisp: the Language*, Digital Press, Maynard MA, 1984.
 - [37] W. Teitelman, *InterLisp Reference Manual*, Xerox PARC, 1978.
 - [38] P. Wadler, "Applicative Languages, Program Transformation, and List Operators", *Proc. ACM Conf. on Functional Programming Languages and Computer Architecture*, :25-32, Oct. 1981.
 - [39] P. Wadler, "Listlessness is Better than Laziness; Lazy Evaluation and Garbage Collection at Compile-Time", *Proc. ACM Symp. on Lisp and Functional Programming*, :45-52, Aug. 1984.
 - [40] P. Wadler, "Listlessness is Better than Laziness II: Composing Listless Functions", *Proc. workshop on Programs as Data Objects*, LNCS 217, H. Ganziger and N. Jones (eds). Springer-Verlag 1985
 - [41] R. Waters, "A Method for Analyzing Loop Programs", *IEEE Trans. on Software Engineering*, 5(3):237-247, May 1979.
 - [42] R. Waters, *LetS: an Expressional Loop Notation*, MIT/AIM-680a, Oct. 1982.
 - [43] R. Waters, "Expressional Loops", *Proc. 1984 ACM Conf. on the Principles of Programming Languages*, :1-10, Jan. 1984.
 - [44] R. Waters, "Efficient Interpretation of Synchronizable Series Expressions" *Proc. ACM SIGPLAN '87 Symposium on Interpreters and Interpretive Techniques*, *ACM SIGPLAN Notices*, 22(7):74-85, July 1987.

- [45] R. Waters, "Using Obviously Synchronizable Series Expressions Instead of Loops" *Proc. 1988 International Conference on Computer Languages*, 338-346, Miami FL, IEEE Computer Society press, October 1988.
- [46] R. Waters, *Optimization of Series Expressions: Part I: A User's Manual for the Series Macro Package*, Massachusetts Institute of Technology technical report AIM-1082, January 1989.
- [47] Waters R.C., "Series", in *Common Lisp: the Language*, Second Edition, G. Steele Jr., Digital Press, Maynard MA, (in press).
- [48] D. Wile, *Generator Expressions*, USC Information Sciences Institute Technical Report ISI/RR-83-116, 1983.
- [49] W. Wulf, R. London, and M. Shaw, "An Introduction to the Construction and Verification of Alphard Programs", *IEEE Trans. on Software Eng.*, 2(4):253-265, December 1976.
- [50] *Military Standard Ada Programming Language*, ANSI/MIL-STD-1815A-1983, U.S. Government Printing Office, February 1983.