

Sourcery VSIPL++

User's Guide

Version 3.1-11



Sourcery VSIPL++: User's Guide: Version 3.1-11

CodeSourcery, Inc.

Copyright © 2005-2011 CodeSourcery, Inc.

All rights reserved.

Table of Contents

I. Working with Sourcery VSIPL++	1
1. API overview	2
1.1. Library initialization	3
1.2. Views	3
1.3. Blocks	17
1.4. Matlab IO	23
2. Direct Data Access	29
2.1. Introduction	30
2.2. Basic usage	30
2.3. Restrictions on dda::Data objects	30
2.4. Non-dense blocks	30
2.5. Requesting a specific data layout	31
3. Using the Dispatch Framework	32
3.1. Introduction	33
3.2. Compile-time dispatch	33
3.3. Runtime dispatch	35
4. Custom Expression Evaluation	37
4.1. The problem	38
4.2. Expression templates	38
4.3. Expression templates in Sourcery VSIPL++	39
4.4. Creating custom expressions	39
4.5. Creating custom evaluators	42
5. Profiling	45
5.1. Enabling Profiling	46
5.2. Using the Profiler	47
5.3. Profiler Output	48
6. Benchmarking	51
6.1. Overview	52
6.2. Building the Benchmarks	52
6.3. Running Benchmarks	53
6.4. Benchmark Output	54
II. Example Application	57
7. Fast Convolution	59
7.1. Fast Convolution	60
7.2. Serial Optimization: Temporal Locality	64
7.3. Performing I/O with User-Specified Storage	66
7.4. Performing I/O with Direct Data Access	69
8. Parallel Fast Convolution	71
8.1. Parallel Fast Convolution	72
8.2. Improving Parallel Temporal Locality	75
8.3. Performing I/O	81
A. Command-line options supported by Sourcery VSIPL++	87
B. Examples	89
B.1. view	90
B.2. ustorage	90
B.3. solvers	90
B.4. signal	90
B.5. dispatch	90
B.6. eval	91
B.7. profile	91
B.8. cvsip	91

B.9. pi	91
Glossary	93

Part I. Working with Sourcery VSIPL++

Basic idioms and programming model

The chapters in this part provide information about fundamentals of Sourcery VSIPL++. You should also refer to the VSIPL++ API Specification and Sourcery VSIPL++ Reference Manual, both of which are available at <http://www.codesourcery.com/vsiplplusplus>.

Chapter 1

API overview

1.1. Library initialization

Before any Sourcery VSIPL++ object can be instantiated, the library itself needs to be initialized, by creating a `vsip::vsipl` object:

```
vsipl library;
```

Sourcery VSIPL++ allows a number of options, which may be passed from the command-line:

```
int main(int argc, char **argv)
{
    vsipl library(argc, argv);
    // ...
}
```

For example, to set the profiler mode of such an application to "accum" mode, the executable may be invoked as

```
./app --vsip-profile-mode=accum
```

A listing of all arguments that may be passed can be found in Appendix A, "Command-line options supported by Sourcery VSIPL++".

1.2. Views

Sourcery VSIPL++ defines a number of mathematical types for linear algebra: vectors, matrices, and (3D) tensors. They provide a high-level interface suitable for solving linear algebra equations. All these types give an intuitive access to their elements. They are collectively referred to as views, as they provide different ways to *look at* data.

You create a view type by instantiating one of the templates `Vector`, `Matrix` or `Tensor`. That is, your view implicitly encapsulates the number of dimensions. It *explicitly* encapsulates a reference to your data - see Section 1.3, "Blocks". It provides you with several related types, several methods to construct an object of the view type, and functions to access data elements, to access slices of data - see Section 1.2.1, "Domains" -, and to return information about the view and its underlying block.

To create a view type, you provide two compile-time parameters: the type of the data elements, and the type of the data block. To create an object of your view type, you provide run-time parameters to one of the constructors:

dimension sizes, initial element value, and block type	to create a new block of the desired shape, filled with the specified value.
a reference to an existing block	to look at the block in a different way.
a reference to an existing view	to look at its block in a different way.

Here are a few simple ways to use views.

```
// create an uninitialized vector of 10 elements
Vector<float> vector1(10);

// create a zero-initialized vector of 10 elements
Vector<float> vector2(10, 0.f);
```

```
// assign vector2 to vector1
vector1 = vector2;

// set the first element to 1.f
vector1(0) = 1.f;

// access the last element
float value = vector1(9);
```

Every view has an associated `Block`, which is responsible for storing or computing the data in the view. More than one view may be associated with the same block.

Depending on how a view is constructed it may allocate the block, or refer to a block from another view. All views created via copy-construction will share the blocks with the views they were constructed with.

When a view is passed by-value to a function, the copy-constructor performs what is known as a *shallow* copy. The function receives a new view that refers to the block of the original.

Assignment, for example, performs a *deep* copy - the contents of the block of the right-hand-side view is copied into the block of the left-hand-side.

```
// copy-construct a new vector from an existing one
Vector<float> vector3(vector1);

// modify the original vector
vector1.put(1, 1.f);

// the new vector reflects the new value
assert(vector3(1) == 1.f);
```

Here is the declaration of one of the view templates:

```
template <typename T      = VSIP_DEFAULT_VALUE_TYPE,
          typename Block = Dense<1, T> >
class Vector;
```

`T` and `Block` are template parameters. The templates for `Matrix` and `Tensor` are similar. The template name specifies its dimension: 1, 2 or 3. The parameter `T` specifies the type of each element in the view. Indexing is 0-based, similar to C.

You cannot use a class template as a type:

```
Vector vec; // ERROR: will not compile
```

You need to provide template arguments:

```
Vector<float, Dense<1, float>> vec; // GOOD
```

Conveniently there are defaults. They are noted above.

```
Vector<> vec; // GOOD: T == float (impl \
dependent)
```

All views provide several operations to the user.

- Accessors of view v

- **Properties**

`v.size()` Number of elements.

`v.block()` Reference to the underlying block of data.

- **Read Values**

`v.get(i)` Read element i .

`v(i)` Same using a more *natural* syntax.

`v.get(d)`,
`v(d)` Here d is a domain - see Section 1.2.1, “Domains”. The result of this kind of expression is a subview of v accessing the elements selected by the domain.

- **Write Values**

`v.put(i, t);` Read or write element i .

`v(i) = t;` Same using a more *natural* syntax.

`v.put(d, sv);`, d is a domain and sv is a subview.
`v(d) = sv;`

If you need more efficient access, Saurcery VSIPL++ provides methods to manipulate directly the storage block of a view. See Section 1.3.1.2, “User Storage”.

- Element-Wise Operations - see Section 1.2.2.1, “Elementwise Operations”

- Element-wise operations on one or more views
- Element-wise operations on views and scalars

- Overloaded Operators: “+”, “*”, ...

Saurcery VSIPL++ overloads the C++ built-in operators. For example, $v1+v2$ produces a view of the same size and shape as $v1$ and $v2$ filled with the element-wise addition of the two views. It is equivalent to `Add(v1, v2)`.

- Reduction Operations

A reduction operates on the values in a view and produces a scalar result. For example, `sumval(v)` adds together all the elements of v and returns the sum. Some reductions also return, via a reference parameter, the index of the scalar result.

```
Index<1> max_idx;  
max = maxval(v, max_idx);
```

The variable `max_idx` contains the index in v of the maximum value.

1.2.1. Domains

A domain represents a logical set of indices. Constructing a one-dimensional domain requires a start index, a stride, and a length. For convenience an additional constructor is provided that only takes a length argument, setting the starting index to 0 and the stride to 1.

```
// [0...9]
vsip::Domain<1> all(10);

// [0, 2, 4, 6, 8]
vsip::Domain<1> pair(0, 2, 5);

// [1, 3, 5, 7, 9]
vsip::Domain<1> impair(1, 2, 5);
```

Two- and three-dimensional domains are composed out of one-dimensional ones.

```
// [(0,0), (0,2), (0,4),..., (1,0),...]
vsip::Domain<2> dom(Domain<1>(10), Domain<1>(0, 2, 5));
```

Views provide convenient access to subviews in terms of subdomains. For example, to assign new values to every second element of a vector, simply write:

```
// assign 1.f to all elements in [0, 2, 4, 6, 8]
vector1(pair) = 1.f;
```

All complex views provide real and imaginary subviews:

```
// a function manipulating a float vector in-place
void filter(Vector<float>);

// create a complex vector
Vector<complex> vector(10);

// filter the real part of the vector
filter(vector.real());
```

1.2.2. View Operations

1.2.2.1. Elementwise Operations

Sourcery VSIPL++ provides elementwise functions and operations that are defined in terms of their scalar counterpart.

```
Vector<float> vector1(10, 1.f);

Vector<complex<float> > vector2(10, complex<float>(2.f, 1.f));

// apply operator+ elementwise
Vector<complex<float> > sum = vector1 + vector2;

// apply conj(complex<float>) elementwise
Vector<complex<float> > result = conj(sum);
```

For binary and ternary functions Sourcery VSIPL++ provides overloaded versions with mixed view / scalar parameter types:

```
// delegates to operator*=(complex<float>, complex<float>)
result *= complex<float>(2.f, 0.f);

// error: no operator*=(complex<float>, complex<double>)
result *= complex<double>(5., 0.);
```

Element-wise expressions allow us to express collective assignment:

```
LHS = RHS;
```

is equivalent to

```
(forall i) LHS(i) = RHS(i);
```

But what happens when LHS and RHS overlap?

```
A = A + B;
```

Ideally, VSIPL++ would evaluate RHS into temporary, then update LHS - but this would impose a performance overhead! If overlap is order-independent (i.e. LHS(i) only affects RHS(i)), then the answer is defined (that is, correct). If overlap is order-dependent, then the answer is undefined. If you have order-dependent overlap, you should introduce a temporary

```
T = A(Domain<1>(0, 1, 10)) + B; A(Domain<1>(1, 1, 10)) = T;
```

These element-wise functions are available:

arithmetic	neg, add, sub, mul, div, ...
trigonometry	cos, sin, ...
basic math	exp, log, pow, ...
comparison	gt, eq, ...
logical	lor, land, lnot, ...
rounding	floor, ceil, ...
complex	real, imag, mag, magsq, conj, ...

These scalar-view element-wise functions are available:

arithmetic	add, sub, mul, div, ...
fused	multiply-add (am), add-multiply (ma), multiply-subtract (msb), subtract-multiply (sbm).
exponential average	expoavg

1.2.2.2. Function Overloading

Many Sourcery VSIPL++ functions are overloaded. For example, `add()` is overloaded to work for different view types, of different value types. Overloading can be powerful. If you change your

algorithm data type from `float` to `double`, few (if any) function names will change. If you add a vector of `floats` to a vector of `complex`, the *right thing* happens.

There are some limits on function overloading. You can't add a vector and a matrix together, even if they're the same size.

1.2.2.3. Operator Overloading

Sourcery VSIPL++ overloads arithmetic operators. Using arithmetic operators

```
A = alpha*X + Y;
```

is more convenient than

```
A = add(mul(alpha, X), Y);
```

Operators are element-wise functions. $A*B$ corresponds to element-wise multiplication, not vector dot-product or matrix product. Terminology:

Multiply treat views as collection of elements

Product treat views as linear algebraic objects

1.2.2.4. Reduction Functions

Sourcery VSIPL++ has two flavors of reductions:

Reductions that return a result only	Logic	<code>alltrue</code> , <code>anytrue</code>
	Mean	<code>meanval</code> , <code>meansqval</code>
	Summation	<code>sumval</code> , <code>sumsqval</code>
Reductions that return a result and an element index	Maximum	<code>maxval</code> , <code>maxmgsqval</code> , <code>maxmgval</code>
	Minimum	<code>minval</code> , <code>minmgsqval</code> , <code>minmgval</code>

1.2.3. View Semantics

What are the semantics behind views? When do views copy data, versus reference the same data? When views are copied via assignment, they have value semantics:

```
A = B;
```

On assignment, values are copied from `B` into `A`. After assignment, `A` and `B` are fully distinct. When views are constructed from a view of the same type, the constructor has reference semantics:

```
Vector<T> B(10);
Vector<T> A(B); // A and B are aliased to same underlying data
```

On construction, `A` refers to the same data as `B`. After construction, changes to `A` will affect `B` (and visa-versa). When views of are constructed from a view of different type, the constructor has value semantics:

```
Vector<T1> B(10);
Vector<T2> A(B); // A and B are distinct
```

The following statements are equivalent:

```
Vector<float> A(B); Vector<float> A = B;
```

Both result in calling Vector's copy-constructor to construct A. The following are not equivalent:

```
Vector<float> B(size); // Create Vector B
Vector<float> A = B;   // A refer's to B's values.
```

And

```
Vector<float> B(size); // Create Vector B
Vector<float> A(size); // construct A,
A = B;                // copy B's values into A
```

Be careful when using '=' in a view declaration.

1.2.3.1. View Initialization Semantics

View construction from the same type has reference semantics.

```
Vector<float> B = ramp(0.f, 1.f, N_elem);
```

View construction from the same type has reference semantics.

```
Vector<float> A(B); // A refer's to B's block.
```

After construction, changes to A are reflected in B (and visa versa).

```
A.put(2, 3.14f);
B.get(2) => returns 3.14f
```

View construction from a different view type has value semantics.

```
Vector<float, BlockType1> B = ramp(0.f, 1.f, N_elem);
Vector<float, BlockType2> A(B);
```

Why would you want reference semantics?

- For creating subviews of existing data, row/column vector views of a matrix
- For passing data to functions, without copying data unnecessarily.

1.2.3.2. Subviews

Subviews allow you to create different views of data. For example

- Creating row and column vector subviews of a matrix.
- Creating a matrix subview of an existing matrix.
- Pulling a vector or matrix out of a tensor.

Subviews are *first-class* views. A row-vector from a matrix is a `Vector`. It can be used like any other vector.

1.2.3.3. Matrix Subviews

Several types of subviews are possible. Here are some vector subviews of a matrix:

```
Row-vector      matrix.row(index_type col);
Column-vector   matrix.col(index_type col);
Diagonal-vector matrix.diag(index_difference_t);
```

And some matrix subviews of a matrix:

```
Sub-matrix      matrix(Domain<2>( ... ));
Read-only       matrix.get(Domain<2>( ... ));
Transpose       matrix.transpose();
```

There are also real-part and imaginary-part subviews if data is complex.

Here is another way to create a row subview:

```
Matrix<float> M(4, 5);
Matrix<float>::row_type subview(M.row(1));
```

1.2.3.4. Subview Variables

`M.row(1)` accesses a row-vector of matrix `M`. You use it directly as a `Vector`:

```
sum = sumval(M.row(1));
```

However, it can be useful to hold a subview in a variable:

```
Vector<float> row = M.row(1); // (line 1)
sum = sumval(row);
```

As written, this works, but it is probably not what you want. Recall that view construction has value semantics if value types are the same and block types are the same. In (line 1) above, the value types are the same, but the `Vector`'s default block is a dense, whereas the subview's block is implementation-defined. How do you know what block type to give the row?

To help you define a vector with the proper block type, `Matrix` defines a subview typedef for `row()` called `row_type`. You can use it like so:

```
Matrix<float>::row_type row = M.row(1);
```

This guarantees that vector `row` will have the proper value type and block type so that references semantics are used. As you will see, `Matrix` defines other subview types as well: `col_type`, `subview_type`, `diag_type`, and so on.

```
Matrix<float>::col_type col = M.col(1); // Column vector
Matrix<float>::diag_type dia = M.diag(0); // Diagonal vector
```

You may also take a matrix subview of a matrix:

```
Matrix<float> M(4, 5);
Domain<2> dom(Domain<1>(1, 1, 2), Domain<1>(0, 2, 3));
Matrix<float>::subview_type subview(M(dom)); // 2 rows x 3 cols
```

It is possible to create a matrix transpose subview of a matrix:

```
Matrix<float> M(rows, cols);
Matrix<float>::transpose_type Mtrans = M.transpose();
```

Now, for all r and c : (where $0 < r < \text{rows}$ and $0 < c < \text{cols}$)

```
M.get(r, c) == Mtrans.get(c, r);
```

$Mtrans$ is a subview, no data is moved. Moreover, changes to $Mtrans$ affect M , and visa-versa. Assignment causes data to be moved:

```
Matrix<float> Mtrans_copy = M.transpose();
```

$Mtrans_copy$ is a transpose of M , but it is also a separate copy. The `trans()` function is synonymous:

```
Matrix<float> Mtrans_copy = trans(M);
```

For Tensors, subviews are created by fixing 1 or more dimensions. Vector subviews:

```
tensor(index_type, index_type, Domain<1>);
tensor(index_type, Domain<1>, index_type);
tensor(Domain<1>, index_type, index_type);
```

Matrix subviews:

```
tensor(index_type, Domain<1>, Domain<1>);
tensor(Domain<1>, index_type, Domain<1>);
tensor(Domain<1>, Domain<1>, index_type);
```

Tensor subviews:

```
tensor(Domain<3>);
```

You can also take a subview of a Vector: `vector(Domain<1>)`. For example, given a vector

```
Vector<float> vec(N_elem);
```

you can create a subview of the first 10 elements:

```
vec(Domain<1>(0, 1, 10));
```

and you can create a subview of every other element:

```
vec(Domain<1>(0, 2, N_elem/2));
```

For views of complex data, you take real and imaginary subviews. To extract the real subview, use `view.real()`; to extract an imaginary subview, use `view.imag()`. Both have the same dimensions as the original view. Examples:

```
Vector<complex<float> > view(...);
Vector<complex<float> >::realview_type real_view = view.real();
Vector<complex<float> >::imagview_type imag_view = view.imag();
```

1.2.3.5. Passing view parameters

Passing views as parameters should be done by reference. You don't want to be copying data! It should be possible to declare a parameter constant. Consider this example function:

```
float mysum1(Vector<float> vec)
{
    float sum = 0.f
    for (index_type i=0; i<vec.size(); ++i)
        sum += vec(i);
    return sum;
}
```

Function arguments have initialization semantics. Recall from the discussion of view semantics Section 1.2.3, “View Semantics” that initialization has reference semantics if both views have the same value type **and** the same block type. There are many situations in which a view's block type is not what you might expect. Suppose you call `mysum` with a row vector subview?

```
mysum(mat.row(0));
```

The argument is a vector but its block is unspecified, and it's safe to assume it is not `Dense`. The initialization of `vec` is equivalent to

```
Vector<float, Dense<1, T>> arg(mat.row(0));
```

This will copy data. Is there a way you can make sure that your parameter always has the same block type as the argument? Yes. Make the block a template parameter.

```
template <typename Block>
float mysum1(Vector<float, Block> vec)
{
    float sum = 0.f
    for (index_type i=0; i<vec.size(); ++i)
        sum += vec(i);
    return sum;
}
```

Template parameter deduction determines what `Block` is. This results in the initialization

```
Vector<float, [correct block]> vec(mat.row(0));
```

where `[correct block]` is the correct block. Thus, here the data is always passed by reference.

1.2.3.6. Returning view results

How do you return a view from a function? Simply declare it to return a view value:

```
template <typename Block>
Vector<float> my_sv_add(float arg1, Vector<float, Block> arg2)
{
    Vector<float> res(arg2.size());
    for (index_type i=0; i<arg2.size(); ++i)
```

```

    res.put(i, arg1 + arg2.get(i));
    return res;
}

```

You can also return a view result by side-effect. Modify the values in a view passed as a parameter. However, see the next section.

1.2.3.7. const view parameters

You can modify values in a view parameter. Since values are by reference, this has a side-effect - the caller sees changes too. You can use this to return values from a function. This is how `by_` reference functions work. But sometimes, you do not intend to modify the parameter. It would be nice to declare the parameter to be `const`. It finds bugs in your function at compile-time, and provides a guarantee to the caller. How do you declare that a `View` parameter is `const`? I.e. that the function will not change its values? Normally you would use the `const` keyword:

```
mysum(Vector<float, Block> const vec)
```

This works partially. You can't modify the view itself, but you **can** modify the data it references. This makes it is fairly easy to remove the `const`, by initializing a non-`const` view from it:

```
Vector<float> nonconstvec(vec);
```

This is what happens when you call another function. To address this, Sourcery VSIPL++ has `const_Views`. They are similar to regular views, but they cannot change their values and cannot be used to initialize a non-`const` view.

```

template <typename Block>
float mysum(const_Vector<float, Block> vec)

```

Now it is more difficult for `mysum` to mistakenly modify `vec`.

1.2.3.8. Views as class members

When using views as class members, it is necessary to use the constructor member initialization list:

```

class My_class {
    // member data
    Vector<float> data; // No arguments given
    ...
    My_class(length_type N) : data(N) { ... }
};

```

1.2.4. Vectors

Here is the declaration of the `Vector` template:

```

template <typename T                = VSIP_DEFAULT_VALUE_TYPE,
         typename Block             = Dense<1, T> >
class Vector;

```

The first template parameter specifies the value type of the view.

```

Vector<float> // vector of floats
Vector<complex<float> > // vector of complex

```

The simplest way to define a vector is to specify its length:

```
Vector<T>(length_type len);
```

For example:

```
Vector<float> vec(10);  
// vector has 10 elements (uninitialized)
```

You can also specify a length and an initial value:

```
Vector<T>(length_type len, T val);
```

This initializes each element to “val”, for example:

```
Vector<float> vec(10, 2.5f);  
// vector has 10 elements (initialized to 2.5f)
```

What else can you do with a `Vector`?

- Check its size:

```
size = v.size(); // returns a length_type
```

- Get the value stored at an index:

```
cur_value = v.get(idx); // returns a T
```

- Change the value stored at an index:

```
v.put(idx, new_value);
```

- Access a value by *view(index)* syntax:

```
cur_value = v(idx); v(idx) = new_value;
```

1.2.5. Matrices

Matrices provide a number of additional subviews.

```
Matrix<float> matrix(10, 10);  
//...  
  
// return the first column vector  
Matrix<float>::col_type column = matrix.col(0);  
  
// return the first row vector  
Matrix<float>::row_type row = matrix.row(0);  
  
// return the diagonal vector  
Matrix<float>::diag_type diag = matrix.diag();
```

```
// return the transpose of the matrix
Matrix<float>::transpose_type trans = matrix.trans();
```

The template declaration is similar to Vector:

```
template <typename T, typename Block>
  class Matrix;
// Declare some matrices.
Matrix<float> m(10, 15); // 10 row x 15 column matrix
Matrix<float> m(10, 15, -1f); // initialize data to -1f
```

Matrices have similar operations.

```
Matrix<float> m(n_rows, n_cols);
n_rows = m.size(0); // number of rows (or column length)
n_cols = m.size(1); // number of columns (or row length)
size = m.size(); // total size (n_rows x n_cols)
```

Element accessors take an additional index:

- Get the value stored at an index/element

```
cur_value = m.get(row_idx, col_idx);
```

- Change the value stored at an index/element

```
m.put(row_idx, col_idx, new_value);
```

- *view(index)* syntax

```
m(row, col) = m(row-1, col) + delta;
```

1.2.6. Tensors

Tensors are three-dimensional views. In addition to the types, methods, and operations defined for all view types, they provide additional methods to access specific subviews:

```
// a 5x6x3 cube initialized to 0.f
Tensor<float> tensor(5, 6, 3, 0.f);

// a subvector
Vector<float> vector1 = tensor(0, 0, whole_domain);
```

The template declaration is similar to Vector:

```
template <typename T, typename Block>
  class Tensor;
// Declare some tensors.
Tensor<float> m(64, 8, 256); // 64 x 8 x 256 tensor (or cube)
Tensor<float> m(2, 2, 2, 3.1f); // initialize data to 3.1f
```

Tensors have similar operations.

```
Tensor<float> m(d1, d2, d3);
n_d1 = m.size(0); // length of dimension 0
n_d2 = m.size(1); // length of dimension 1
n_d3 = m.size(2); // length of dimension 2
size = m.size(); // total size (n_d1 x n_d2 x n_d3)
```

Element accessors take a third index:

- Get the value stored at an index/element

```
cur_value = m.get(d1_idx, d2_idx, d3_idx);
```

- Change the value stored at an index/element

```
m.put(d1_idx, d2_idx, d3_idx, new_value);
```

- *view(index)* syntax

```
m(d1_idx, d2_idx, d3_idx) = m(d1_idx-1, d2_idx, d3_idx) + delta;
```

The symbolic constant `whole_domain` is used to indicate that the whole domain the target view holds in a particular dimension should be used. In the example above that not only provides a more compact syntax compared to explicitly writing `Domain<1>(6)` but it also enables better optimization opportunities.

```
// a submatrix
Matrix<float> plane = tensor(whole_domain, 0, whole_domain);

Tensor<float> upper_half = tensor(whole_domain, Domain<1>(3), \
whole_domain);
```

1.2.7. Miscellaneous

1.2.7.1. Promotion Rules

When combining values of different types,

- `float * int = ???`
- `float + double = ???`
- `float - complex<float> = ???`

what type should the result be? Saurcery VSIPL++ determines the result type through promotion rules:

```
Promotion<float, int>::type == float;
Promotion<float, double>::type == double;
Promotion<float, complex<float>>::type == complex<float>;
Promotion<double, complex<float>>::type == complex<double>;
```

1.3. Blocks

The data accessed and manipulated through the View API (see Section 1.2, “Views”) is actually stored in blocks. Blocks are reference-countable (see Section 1.3.2, “Reference Counting”), allowing multiple views to share a single block. However, blocks may themselves be proxies that access their data from other blocks (possibly computing the actual values only when these values are accessed). These blocks are thus not modifiable. They aren't allocated directly by users, but rather internally during the creation of subviews, for example.

Blocks allocated by users are most often Dense blocks. See Section 1.3.1, “Dense Blocks” below. To create a dense block type, you provide four compile-time parameters: dimension (1, 2 or 3), element type, dimension ordering and map type. See Section 1.3.1, “Dense Blocks” [18]. To create an object of your block type, provide run-time parameters to one of the constructors:

domain and map	to create a new, uninitialized block.
domain, fill value, map	to create a new block filled with the specified value.
domain, pointer, map	to create a new block that references data you allocated and initialized elsewhere.

In all the constructors, the map may be omitted. Its default is an instance of the map type passed as a parameter to the Dense template.

1.3.1. Dense Blocks

The default block type used by all views is Dense, meaning that `Vector<float>` is actually a shorthand notation for `Vector<float, Dense<1, float> >`. As such Dense is the most common block type directly used by users. Dense blocks are modifiable and allocatable. They explicitly store one value for each index in the supported domain:

```
// create uninitialized array of size 3
Dense<1, float> array1(Domain<1>(3));

// create array of size 3 with initial values 0.f
Dense<1, float> array2(Domain<1>(3), 0.f);

// assign array2 to array1
array1 = array2;

// access first item
float value = array1.get(0);

// modify first item
array1.put(0, 1.f);
```

A dense block get its size from the Domain argument passed to its constructor. You can also specify an initial value.

What is the difference between

```
Dense<1, float> blk1(Domain<1>(0, 1, 10));
```

and

```
Dense<1, float> blk2(Domain<1>(107, -3, 10));
```

Nothing. Only the domain's dimension sizes are used in creating a Dense block; its offset and stride are ignored.

You use the following class template to create a dense block type:

```
template <dimension_type Dim,
          typename T,
          typename OrderT,
          typename MapT>
class Dense;
```

It has the following parameters:

Dim the number of block dimensions
T the type of values stored in the block
OrderT the dimension ordering
MapT how data is distributed

Dim and *T* are used most commonly, but all parameters have default values.

Most often you will access data values and information with views. Using a block directly, you can check its size

```
size = blk.size();            // returns block's total size
size = blk.size(Dim, d);     // returns block size in dimension d
                              // useful if block is 1,x-dimensional.
```

and perform element-wise operations.

```
cur_value = blk.get(idx);    // return the value stored at index idx
blk.put(idx, new_value);    // change the value stored at index idx
```

You can also access user storage. See Section 1.3.1.2, “User Storage”. You can also attach a view to a block.

```
Dense<1, float>* blk = new Block(Domain<1>(n_elem));
Vector<float, Dense<1, float> > v(*blk);
```

1.3.1.1. Layout

Beside the two template parameters already discussed above, Dense provides an optional third parameter to specify its dimension ordering. Using this parameter you can explicitly control whether a 2-dimensional array should be stored in row-major or column-major format:

```
// array using row-major order
Dense<2, float, tuple<0, 1> > rm_array;

// array using column-major order
Dense<2, float, tuple<1, 0> > cm_array;
```

Row-major arrays store rows as contiguous chunks of memory. Iterating over its columns will thus access close-by memory regions, reducing cache misses and thus enhancing performance:

```
length_type size = rm_array.size(0);
for (index_type i = 0; i != size; ++i)
    rm_array.put(i, 1.f);
```

Layout can affect *locality* and is important for performance.

Sourcery VSIPL++ provides typedefs for tuple:

`rowD_type` indicates a D-dimensional row-major tuple.

`colD_type` indicates a D-dimensional column-major tuple.

To create a row-major 2-dimensional dense block:

```
Dense<2, float, row2_type> block(Domain<2>(3, 3));
```

To create a column-major 2-dimensional dense block:

```
Dense<2, float, col2_type> block(Domain<2>(3, 3));
```

It is good practice to use layout explicitly when declaring views. Instead of

```
Matrix<T> data(rows, cols);
```

use

```
typedef Dense<2, T, row2_type> block_type;
typedef Matrix<T, block_type> view_type;
view_type data(rows, cols);
```

This makes it easier to change dimension-ordering and add parallel mappings.

1.3.1.2. User Storage

They also allow user-storage to be provided, either at construction time, or later via a call to `rebind`:

```
float *storage = ...;

// create array operating on user storage
Dense<1, float> array3(Domain<1>(3), storage);

// create uninitialized array...
Dense<1, float> array4(Domain<1>(3));

// ...and rebind it to user-storage
array4.rebind(storage);
```

However, special care has to be taken in these cases to synchronize the user storage with the block using it. While the storage is being used via the block it was rebound to, it has to be *admitted*, and *released* in order to be accessed directly, i.e. outside the block.

```
// grant exclusive access to the block
array3.admit(true);

// modify it
array3.put(0, 1.f);
```

```
// force synchronization with storage
array3.release();

// access storage directly
assert(storage == 1.f);
```

`admit` and `release` take a boolean `update` flag:

- `true` indicates that consistency is required.
- `false` indicates that consistency is not required.

If the library does not re-arrange or copy data, the `update` may have no effect - but applications should not assume this. If library does re-arrange or copy data, `update` tells the library when consistency is required. If consistency is not required, the library may skip re-arrangement or copy.

Why would you set `update` to `false`? Consider a block used to bring external data into the system:

```
float ptr[10];
Dense<1, float> block(10, ptr);
for (...)
    fread(ptr, 10, sizeof(float), FILE); // read data from file
    block.admit(true); // admit data in block
    ... use data read in from file ... // require consistency
    block.release(false); // release block
    ... don't care about data ... // don't require consistency
```

After the `release`, the data referred to by `ptr` will not be used,

- it will just be overwritten;
- reorganizing block's data is unnecessary;
- if library reorganizes complex data (for example, converting from interleaved to split to facilitate SIMD), then setting `update` to `false` on `release` will avoid unnecessary reorganization.

You can check the `admit/release` state of a block with `admitted`. For example:

```
float buf[10];
Dense<1, float> block(10, buf);
assert(block.admitted() == false);
block.admit(true);
assert(block.admitted() == true);
block.release(true);
assert(block.admitted() == false);
```

A block without user-storage is admitted at all times. Calls to `admitted` will always return `true`.

1.3.2. Reference Counting

One important function of views is to reference count blocks.

- Each block maintains a reference count.
- Initially the reference count is one.

- When a new view is created on the block, the count is incremented.
- When a view is destroyed, the view's block's count is decremented.
- When a block's count goes to zero, it is freed.

For example:

```
for (int i=0; i<1000; ++i) {
Vector<float> vec(10); // allocate 10-element vector
... use vec ...
}
// vec goes out of scope,
// its block refcount is decremented
// if it is zero, block is freed.
```

For another example:

```
Dense<1, float>* blk = new Block(Domain<1>(n_elem)); // blk count \
is 1
Vector<float, Dense<1, float> > v(*blk);           // blk count \
is 2
```

Initializing a block's reference count to one prevents user created blocks from being freed inadvertently. Consider

```
Dense<1, float> blk(n_elem); // allocate blk on the stack.
Vector<float> v(blk);
```

If blk's reference count went to zero, it would attempt to destroy itself. This would cause undefined behavior.

1.3.3. Complex Data

There are three ways to specify complex storage:

Array Format Directly, with a pointer to complex values.

```
Dense<Dim, complex<float> >(length_type, \
complex<float>*);
```

Interleaved A scalar pointer to interleaved real/imaginary values.

```
Dense<Dim, complex<float> >(length_type, float*);
```

Split Two scalar pointers, one to real values, one to imaginary.

```
Dense<Dim, complex<float> >(length_type, float *, \
float *);
```

1.3.4. Storage Pointer

To find the pointer bound to a released user-storage block:

```
void Dense<Dim, T>::find(T*& pointer);
```

For example:

```
Dense<1, float> block(size, data_buffer);  
float *ptr;  
block.find(ptr);
```

find's behavior depends on the admit/release state:

- If block is released, find places the block's storage pointer in ptr.
- If block is admitted, find places NULL in ptr.
- Also, if block is not a user storage block, it places NULL in ptr.

For complex user-storage blocks:

- with interleaved format:

```
void Dense<1, complex<float> >::find(float*&);
```

- with split format:

```
void Dense<1, complex<float> >::find(float*&, float*&);
```

- with array format, the standard find() is used.

To change the storage pointer that a block is bound to:

```
void Dense<Dim, T>::rebind(T* pointer);
```

rebind only works for released blocks. Otherwise its behavior is undefined.

For complex user-storage blocks:

- with interleaved format:

```
void Dense<1, complex<float> >::rebind(float*);
```

- with split format:

```
void Dense<1, complex<float> >::rebind(float*, float*);
```

- with array format, the standard rebind is used.

Finally, the storage type of a block can be checked:

```
user_storage_t Dense<Dim, T>::user_storage();
```

user_storage_t is an enumeration whose values are described below.

If the block does have user-specified storage, user_storage returns:

array_format for array format, all scalar storage, and some complex storage.

`interleaved_format` for interleaved complex storage.

`interleaved_format` for split complex storage.

If the block does not have user-specified storage, it returns `no_user_format`.

1.4. Matlab IO

Sourcery Sourcery VSIPL++'s `vsip_csl` library has routines that can read and write views from Matlab formatted text and binary files.

1.4.1. Matlab Text (.m) Files

The `Matlab_text_formatter` object writes a view to an output stream in Matlab text file format (Matlab text files commonly have an `.m` suffix).

The following example illustrates using the `Matlab_text_formatter` to write a matrix and a vector to the same file.

The first part of the example shows the necessary `includes` and declarations. In addition to any Sourcery VSIPL++ headers necessary for your program (the example includes headers for vectors, matrices, and generation functions), it is also necessary to include the `vsip_csl/matlab_text_formatter.hpp` header file. The example also uses the `vsip` and `vsip_csl` namespaces for convenience.

```
#include <iostream>
#include <fstream>

#include <vsip/initfin.hpp>
#include <vsip/vector.hpp>
#include <vsip/matrix.hpp>
#include <vsip/selgen.hpp>
#include <vsip/map.hpp>

#include <vsip_csl/matlab_text_formatter.hpp>

using namespace vsip;
using namespace vsip_csl;
```

The second part of the example shows writing a file. First matrix `m` and view `v` are created and filled with ramp data. Then an output file stream `out` is created. Finally, `Matlab_text_formatter` is used to write the views.

```
// Initialize matrix 'm'.
Matrix<float> m(3, 3);
for(index_type i=0;i<3;i++)
    m.row(i) = ramp<float>(3*i, 1, 3);

// Initialize vector 'v'.
Vector<float> v(3);
v = ramp<float>(0, 1, 3);

// Open output stream to file 'temp.m'.
std::ofstream out("text.m");
```

```
// Write 'm' and 'a' to output stream
out << Matlab_text_formatter<Matrix<float> >(m, "m");
out << Matlab_text_formatter<Vector<float> >(v, "v");
```

The output file `temp.m` contains the following

```
m =
[
  [ 0 1 2 ]
  [ 3 4 5 ]
  [ 6 7 8 ]
];
v =
[ 0 1 2 ];
```

This text file can be run inside of a Matlab console window to load matrix `a` and vector `v`.

1.4.2. Matlab Binary Files (.mat)

The `Matlab_bin_formatter` object can read and write views to a streams in Matlab binary file format (Matlab binary files commonly have a `.mat` suffix). For reading matlab binary files, the iterator interface described in the next section may be more convenient.

1.4.2.1. Writing a Matlab Binary Format File

Writing matlab binary format files is similar to writing text format files, except that a header must be written to the file with `Matlab_bin_header` before writing each view with `Matlab_binary_formatter`. The following example shows how to write a matrix and a vector to a `.mat` file.

The first part of the example shows the necessary `includes` and declarations. In addition to any Sourcery VSIPL++ headers necessary for your program (the example includes headers for vectors, matrices, and generation functions), it is also necessary to include the `vsip_csl/matlab_bin_formatter.hpp` header file. The example also uses the `vsip` and `vsip_csl` namespaces for convenience.

```
#include <iostream>
#include <fstream>

#include <vsip/initfin.hpp>
#include <vsip/vector.hpp>
#include <vsip/matrix.hpp>
#include <vsip/selgen.hpp>
#include <vsip/map.hpp>

#include <vsip_csl/matlab_bin_formatter.hpp>

using namespace vsip;
using namespace vsip_csl;
```

The second part of the example shows writing a file. First matrix `m` and view `v` are created and filled with ramp data. Then an output file stream `out` is created. Finally, `Matlab_text_formatter` is used to write the views.

```

// Initialize matrix 'm'.
Matrix<float> m(3, 3);
for(index_type i=0;i<3;i++)
    m.row(i) = ramp<float>(3*i, 1, 3);

// Initialize vector 'v'.
Vector<float> v(3);
v = ramp<float>(0, 1, 3);

// Open output stream to file 'sample.mat'.
std::ofstream out("sample.mat");

// Write matlab binary format header. This must be done once at the
// beginning of the file before any views can be written.
out << Matlab_bin_hdr("example");

// Write 'm' and 'v' to output stream
out << Matlab_bin_formatter<Matrix<float> >(m, "m");
out << Matlab_bin_formatter<Vector<float> >(v, "v");

```

This result file `sample.mat` can be read by Matlab, or other programs capable of reading matlab binary format files, such as Octave and Sourcery Sourcery VSIPL++ applications.

1.4.2.2. Reading a Matlab Binary Format File

Reading a matlab binary format file is similar to writing one. After the file is opened, it is necessary to read the file header. This header is used by each of the subsequent reads. The following example shows how to read the views back from the `sample.mat` file written in the previous example.

The same includes are used for this example.

First matrix `m` and view `v` are created. Their size must match the size of the views in the matlab binary file. Next an input stream is created to read the binary data. The header is read first into a `Matlab_bin_header`. Finally each view is read, using `Matlab_bin_formatter` objects.

```

// Create matrix and vector views of correct size.
Matrix<float> m(3, 3);
Vector<float> v(3);

// Open an input stream to read sample.mat.
std::ifstream in("sample.mat");

// Read matlab binary format file header. This must be done once
// after the file is opened before reading any views. The header \
is
// then used by Matlab_bin_formatter to determine global file
// parameters such as endianness, etc.
Matlab_bin_hdr h;
in >> h;

// Read the views.
in >> Matlab_bin_formatter<Matrix<float> >(m, "m", h);
in >> Matlab_bin_formatter<Vector<float> >(v, "v", h);

```

Note that when using `Matlab_bin_formatter` the size and types of the Sourcery VSIPL++ views `m` and `v` must match the size and type of the views stored in the binary file. If they do not match, an exception will be thrown. In situations where the size and type are not known in advance, it may be more convenient to use the Matlab iterator interface, described in the next section.

1.4.3. Matlab_file iterator interface

In situations where the size, type, and order of views written in a matlab file is not known in advance, the `Matlab_file` interface should be used to read the file. `Matlab_file` provides an iterator interface to step through each view in a file. The size, type, and name of each view can be queried before it is read. This allows an appropriate Sourcery VSIPL++ view to be constructed dynamically.

The `Matlab_file` object handles opening and reading a file. It provides a standard iterator interface with `begin` and `end` functions.

`Matlab_file::iterators` correspond to views in the file. Dereferencing the iterator returns a `Matlab_view_header` object. This contains information about the view, including its name, type, dimensionality, and size.

The `read_view` function reads the view referred to by the iterator.

The following example shows how to read the `m` matrix from the `sample.mat` binary file used in the previous examples.

The first part of the example shows the necessary `includes` and declarations. The `Matlab_file` interface is contained in the `vsip_csl/matlab_file.hpp` header file.

```
#include <iostream>>
#include <fstream>

#include <vsip/initfin.hpp>
#include <vsip/vector.hpp>
#include <vsip/matrix.hpp>
#include <vsip/selgen.hpp>
#include <vsip/map.hpp>

#include <vsip_csl/matlab_file.hpp>

using namespace vsip;
using namespace vsip_csl;
```

The second part of the example shows how to read the file. First the `Matlab_file` object `mf` is created. Then iterators `begin` and `end` are created to iterate over the views stored in the file. For each view, the name and size are checked to determine if it should be read. Finally, `read_view` is used to read the selected view.

```
// Create Matlab_file object for 'sample.mat' file.
Matlab_file mf("sample.mat");
Matlab_file::iterator cur = mf.begin();
Matlab_file::iterator end = mf.end();
Matlab_view_header* vhdr;

// Block pointer to hold the matrix. The block will be allocated
// once it's size is known.
Dense<2, float>* m_block = NULL;
```

```

// Iterate through views in file.
while (cur != end)
{
    vhdr = *cur;

    // Check if view is the one we're looking for.
    if(!strcmp(vhdr->array_name, "m") && vhdr->num_dims == 2)
    {
        // Check for multiple views named "m" in file.
        assert(m_block == NULL);

        // Create block and view.
        // At this point we can make the block size match size in \
the file.
        m_block = new Dense<2, float>(Domain<2>(vhdr->dims[0], \
vhdr->dims[1]));
        Matrix<float> tmp(*m_block);

        // Read view from file.
        mf.read_view(tmp, cur);
    }

    ++cur; // Move to next view stored in the file.
}

// Check that we found a view named "a" in file.
assert(m_block != NULL);

// Create a view to process "m".
Matrix<float> m(*m_block);

```

The handling of vectors in Matlab files requires special consideration. Matlab stores vectors as matrices with one dimensions of size 1. Sourcery VSIP++ can read matlab vectors as either Sourcery VSIP++ vectors or matrices.

The following example shows how to read the `v` vector from the `sample.mat` binary file used in the previous examples.

```

// Create Matlab_file object for 'sample.mat' file.
Matlab_file mf("sample.mat");
Matlab_file::iterator cur = mf.begin();
Matlab_file::iterator end = mf.end();
Matlab_view_header* vhdr;

// Block pointer to hold the vector. The block will be allocated
// once it's size is known.
Dense<1, float>* v_block = NULL;

// Iterate through views in file.
while (cur != end)
{
    vhdr = *cur;

```

```
// Check if view is the one we're looking for.
//
// Note: even though 'v' is a vector, it will be 2D because of \
how
//      matlab stores vectors.
if(!strcmp(vhdr->array_name, "v") && vhdr->num_dims == 2)
{
    // Check for multiple views named "m" in file.
    assert(v_block == NULL);

    // Determine the vector's size:
    length_type size = std::max(vhdr->dims[0], vhdr->dims[1]);

    // Create block and view.
    // At this point we can make the block size match size in \
the file.
    v_block = new Dense<1, float>(Domain<1>(size));
    Vector<float> tmp(*v_block);

    // Read view from file.
    mf.read_view(tmp, cur);
}

++cur; // Move to next view stored in the file.
}

// Check that we found a view named "a" in file.
assert(v_block != NULL);

// Create a view to process "m".
Vector<float> v(*v_block);
```

Chapter 2

Direct Data Access

2.1. Introduction

While views and blocks provide a powerful abstraction to handle data inside the Sourcery VSIPL++ API, it may be necessary to access the raw data, for example to pass them to third-party libraries or legacy code that requires data to be cast into different types.

To be able to cast Blocks into third-party types, without being required to copy the data, Sourcery VSIPL++ provides a *Direct Data Access* API.

2.2. Basic usage

Direct data access is provided by the `dda::Data<>` class. In its simplest form, it can be instantiated with a block, and provides access to the data as a raw pointer by means of the `ptr()` accessor.

```
void process(float *data, size_t size);
...
Vector<float> v(8);
dda::Data<Vector<float>>::block_type, dda::inout>
data(v.block());
process(data.ptr(), data.size());
```

This works, as long as the underlying block type is known to have unit stride. If this assumption can not be made, the access needs to take the stride into account.

2.3. Restrictions on `dda::Data` objects

In order to avoid data corruption, when a `dda::Data` object is created for a given block, certain other operations on the block are prohibited until the `dda::Data` object is destroyed. In particular, if a `dda::Data` object with `out` or `inout` access is created, its parent block may not be accessed by any other means (including additional `dda::Data` objects). If a `dda::Data` object with `in` access is created, the parent block may be read by other means (including the creation of additional `dda::Data` objects with `in` access), but it may not be modified, and no `dda::Data` objects with `out` or `inout` access referring to the same block may be created.

If a `dda::Data` object refers to a sub-block of another block, these restrictions apply to the parent block as well as the sub-block. Likewise, when these restrictions apply to a given block, they apply to any sub-blocks that may exist.

A pointer returned by the `ptr()` member function is invalid once the corresponding `dda::Data` object is destroyed.

2.4. Non-dense blocks

In the following example, we construct a subview of a dense view, aliasing every second value from a dense view, yielding a view with stride 2:

```
void process(float *data, ptrdiff_t stride, size_t size);
...
Vector<float> v(8);
Vector<float>::subview_type subview = view.get(Domain<1>(0, 2, 4));
dda::Data<Vector<float>>::block_type, dda::inout> data(v.block());
process(data.ptr(), data.stride(), data.size());
```

However, some functions may require unit-stride input. In that case, it is possible to force unit-stride access. The `dda::Data<>` object will copy the data into temporary storage, which the user then operates on, and synchronize back with the block it was constructed from. This synchronization may not always be necessary, and so it is possible to express whether to synchronize only from the block to the `dda::Data<>` object, the inverse, or both.

```
typedef Vector<float>::subview_type::block_type block_type;
typedef Layout<1, row1_type, unit_stride> layout_type;
dda::Data<block_type, dda::out, layout_type> data(subview.block());
ramp(data.ptr(), data.size());
```

As this access type may involve a performance penalty (temporary data allocation, as well as one or two copy operations), it is desirable to be able to query whether the direct data access comes with an extra cost. A user may decide to prefer unit-stride, as long as no copies are involved, but fall back to non-unit stride access otherwise:

```
typedef Vector<float>::subview_type::block_type block_type;
typedef Layout<1, row1_type, unit_stride> layout_type;
if (dda::Data<block_type, dda::in, layout_type>::ct_cost != 0)
{
    // If unit-stride access would require a copy,
    // choose non-unit stride access
    dda::Data<block_type, dda::out> data(subview.block());
    ramp(data.ptr(), data.stride(0), data.size());
}
else
{
    dda::Data<block_type, dda::out, layout_type> data(subview.block());
    ramp(data.ptr(), data.size());
}
```

Note that the `dda::Data<>::ct_cost` value is a compile-time constant, and thus can be used in compile-time expressions. Therefore, the above conditional may be done via compile-time decisions (e.g. template specializations).

In the above example we have considered unit- versus non-unit stride direct data access. However, there are other cases where data can never be accessed directly without copies, for example if the block represents an expression.

2.5. Requesting a specific data layout

The DDA API allows the user to express a very rich set of layout requirements. An aligned unit-stride (but not necessarily dense) row-major matrix can be requested like this:

```
typedef Layout<2, row2_type, aligned_32> layout_type;
dda::Data<block_type, dda::in, layout_type> data(view.block());
```

And a dense column-major matrix may be requested like this:

```
typedef Layout<2, col2_type, dense> layout_type;
dda::Data<block_type, dda::in, layout_type> data(view.block());
```

For details on the `Layout` class template, see Section 5.3, “The `Layout` template”

Chapter 3

Using the Dispatch Framework

3.1. Introduction

Writing High-Performance code for a wide range of hardware is very challenging. Typically, the software is targetted at particular hardware or optimized for a specific set of parameters. This article describes a mechanism to interface a set of functions covering the same functionality but for different hardware or types of input with a single API, using a mechanism to dispatch to the most appropriate backend.

Sourcery VSIPL++ is configurable to target a wide range of backend implementations for most of its functions. it achives portability by hiding these backends behind common interfaces, yet strives to minimize the calling overhead by doing as much as possible at compile-time.

When the user performs a particular operation (e.g., adding two vectors) the library must select an appropriate implementation. For example, if the vectors are single-precision floating-point types, then a special SIMD routine might be used to perform the addition efficiently. Or, if the vectors are distributed across processors, multi-processor communication might be required.

When determining how to implement a given operation, Sourcery VSIPL++ performs a two-step process. One step is performed at compile-time; the other at run-time. Conceptually, the process is as follows:

1. Sourcery VSIPL++ forms a list of all possible implementations of the operation.
2. At compile-time, those implementations which do not accept arguments of appropriate types, or which are otherwise inappropriate for reasons which can be determined statically, are eliminated.
3. At run-time, each implementation not yet eliminated at compile-time is queried to see whether it can perform the operation. The first implementation that is able to perform the operation is used.

Each implementation is provided as a (possibly partial) specialization of the Evaluator class template. The library checks the `ct_valid` static data member to determine compile-time suitability and calls the `rt_valid()` static member function at run-time to determine run-time suitability. The actual implementation of the operation is performed by the `exec()` static member function.

3.2. Compile-time dispatch

Let us assume we want to implement an operation `Operation` by means of three backend classes `A_impl`, `B_impl`, and `C_impl`, all modeling the same concept:

```
class A_impl { static void process(int*, size_t); };
template <typename T> class B_impl
{
    static void process(T*, size_t);
};
template <typename T> class C_impl
{
    static void process(T*, size_t);
};
```

For each backend, we define an Evaluator, that is, a meta-function that will be used to evaluate whether for a given type this backend is usable.

Then we expose these backends by declaring a list of backends for the operation in question, so the dispatcher can iterate over it to find the first match. (Note that typically this list will also depend on

configuration parameters, so the selection of the appropriate backend is in fact done in part at configure time and in part at compile time.

```
#include <vsip/opt/dispatch.hpp>
struct Operation;
struct A; struct B; struct C;
namespace vsip_csl
{
  namespace dispatcher
  {
    template <> struct Evaluator<Operation, A, int>
    { static
      { bool const ct_valid = true;
        typedef A_impl backend;
      };
    };
    template
      <typename T> struct Evaluator<Operation, B, T>
    { static bool const ct_valid =
      { has_feature<T>::value;
        typedef B_impl<T> backend;
      };
    };
    template
      <typename T> struct Evaluator<Operation, C, T>
    {
      static bool const ct_valid = true;
      typedef C_impl<T> backend;
    };
    template <> struct List<Operation>
    {
      typedef
        Make_type_list<A, B, C>::type type;
    };
  };
}
```

Here, the Evaluator declares A to be available for `int` types, B is defined in a way that allows some external meta-function `has_feature` to evaluate the availability of this backend for a given type, while Evaluator C declares C to be available for all types. (Typically, a catch-all backend is made available that is guaranteed to work for all types, but as it is slow, it is only available if no other backend matches.)

With that, writing the wrapper function that will do the actual dispatch is very simple:

```
template <typename T>
void process(T *input, size_t size)
{ // Evaluate the dispatch (at compile-time): using
  vsip_csl::dispatcher::Dispatcher<Operation, T>;
  typedef
    Dispatcher<Operation, T>::backend backend_type;
  // Now use it.
  backend_type::process(input, size);
}
```

3.2.1. Implementation details

The `List<Operation>` type above creates a type-list of backend tags. The `Dispatcher<Operation>` class template constructs a type-list of Evaluators from that, which the compiler iterates over to select the first item for which the `ct_valid` member is true.

In other words, the `Dispatcher<Operation>` acts as a meta-function that takes a typename `T` as input, and returns a backend.

3.3. Runtime dispatch

Now let us make modifications to the above by stipulating that the function operates on an array of type `T` and size `size`, and backends have restrictions both on the type as well as the size of the array: Backend A only accepts ints, backend B accepts any type, but only buffers whose size is a power of 2, while backend C accepts any input:

```
namespace vsip_csl
{
    namespace dispatcher
    {
        template <>
        struct Evaluator<Operation, A, void(int*, size_t)>
        {
            static bool const ct_valid = true;
            static bool rt_valid(int*, size_t)
            { return true; }
            static void exec(int *input, size_t size)
            { A_impl::process(input, size); }
        };

        template <typename T>
        struct Evaluator<Operation, B, void(T*, size_t)>
        { static bool const ct_valid =
            has_feature<T>::value;
            static bool rt_valid(T*, size_t size)
            { return size^(size-1); }
            static void exec(T *input, size_t size)
            {
                B_impl<T>::process(input, size);
            }
        };

        template <typename T>
        struct Evaluator<Operation, C, void(T*, size_t)>
        {
            static bool const ct_valid = true;
            static bool rt_valid(T*, size_t)
            { return true; }
            static void exec(T *input, size_t size)
            { C_impl<T>::process(input, size); }
        };

        template <> struct List<Operation>
        {
            typedef
                Make_type_list<A, B, C>::type type;
        };
    };
};
```

```
    };  
  }  
}
```

The operation wrapper for this now simply becomes:

```
template <typename T>  
void process(T *input, size_t size)  
{  
    vsip_csl::dispatch<Operation, void(T*, size_t)>(input,  
        size);  
}
```

3.3.1. Implementation details

The runtime-dispatch works conceptually similar to the compile-time dispatch. However, in this case the dispatch is actually a two-phase process. The first phase is the same as in the compile-time dispatch. It reduces the type-list of Evaluators to those elements that match the given input type(s).

The second phase, then, is carried out at runtime, when this reduced type-list is traversed to evaluate all `rt_valid()` member functions until a match is found, based on runtime characteristics. In the case presented here this is a size parameter, though it could be anything else, as the signature of the operation to be carried out is a template parameter, too.

Chapter 4

Custom Expression Evaluation

4.1. The problem

The VSIP++ API allows operations to be expressed in a concise mathematical way. Using operator overloading, it is possible to write expressions involving vectors and other VSIP++ types, such as:

```
Vector<float> y(100), a(100), b(100), c(100), d(100);
y = (a+b)/(c-d);
```

However, this abstraction normally comes with a cost. Each binary operation is evaluated individually. This is equivalent to the following:

```
Vector<float> y(100), a(100), b(100), c(100), d(100);
Vector<float> tmp1 = a+b; Vector<float> tmp2 = c-d;
Vector<float> tmp3 = tmp1/tmp2; y = tmp3;
```

Thus, if the construction of those temporary objects is expensive, this adds considerable overhead. Further, for an element-wise definition of the above binary operations, each line performs a loop over all elements:

```
for (unsigned int i = 0; i != size; ++i)
    tmp1[i] = a[i] + b[i];
```

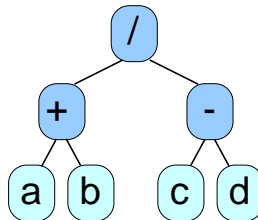
It would be desirable to avoid the temporary objects, and be able to fuse the different loops into one.

4.2. Expression templates

In Sourcery VSIP++, the above operators are defined such that they don't evaluate the operation directly. Instead, they return a representation of the operation that can then be evaluated later. Using this technique, an entire expression will map to a single *expression object*, which then can be evaluated in a single pass, during the assignment.

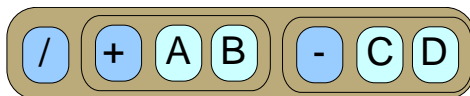
In other words, all such operators don't return a *value*, but a *parse tree* containing its operands, as well as operation

Figure 4.1. Parse tree representation of the expression $(a+b)/(c-d)$



With a suitable definition of those binary operators, this results into a result type like:

Figure 4.2. Expression type generated from the above expression $(a+b)/(c-d)$



If an assignment operator exists that takes such types as input, it may be able to perform a much faster assignment, without any temporaries.

4.3. Expression templates in Sourcery VSIPL++

4.3.1. Generating expression templates

Most functions and operators provided by Sourcery VSIPL++ are written to support such lazy evaluation. Thus, instead of calculating the result value directly, they return objects of such expression types. Moreover, they are written in such a way that their input may be expression type objects, too, such that expressions can be nested without being evaluated.

```
template <typename T,
         typename LeftBlock,
         typename RightBlock>
Vector<T,
      Binary_expr_block<1, Plus, LeftBlock, T, RightBlock, T> >
operator+(Vector<LeftBlock> left,
         Vector<RightBlock> right);
```

This version of `operator+` expects Vectors of arbitrary block types as input, and generates a Vector whose block type encodes the binary operation. (If the input arguments are themselves expressions, this will result in a composite expression tree.)

Sourcery VSIPL++ provides a range of types to represent such *non-terminal* expression nodes, as well as the means to traverse them, to make it possible to generate expression (parse) trees for a wide range of expressions.

4.3.2. Evaluating expression templates

Expressions are ideally evaluated only once all the relevant information has been gathered. That point is typically reached once the whole assignment-expression has been seen. That is the case during assignment-operator evaluation.

Here is a conventional implementation of a vector-assignment, doing element-wise assignment:

```
template <typename T,
         typename LeftBlock,
         typename RightBlock>
void assign(Vector<T, LeftBlock> left,
          Vector<T, RightBlock> right)
{
    for (length_type i = Vector<T, LeftBlock>::size(); i; --i)
        left.put(i - 1, right.get(i - 1));
}
```

If `RightBlock` is an expression-block with an elementwise operation, its implementation of `get` will perform the (elementwise) evaluation. Thus, in this case, at least one temporary has already been eliminated, and multiple loops have been fused into one.

For non-elementwise operations, this is not quite as simple.

4.4. Creating custom expressions

Sourcery VSIPL++ allows you to write custom functions that participate in expression template dispatch and evaluation. This optimizes handling of the functions return value and allows custom evaluators to recognize fused expressions containing the expression.

Let us work through an example, starting with the function `scale()`:

```
template <typename T,
         typename BlockType>
Vector<T> scale(Vector<T, BlockType> a, T value)
{
    Vector<T> r = a * value;
    return r;
}
```

This function takes a vector, scales it by a scalar value, and returns the result. As the result is returned by-value, it gets copied during assignment. In other words, the return value is a temporary object, which we may want to avoid.

To do that, we use a variant of a technique known as *return value optimization*. We rewrite `scale()` to return an expression type, which is only evaluated once the result object is available, so the computed value can be stored in-place. To do that, we capture the function logic into a functor, and rewrite the `scale()` function to return an *expression block* vector:

```
using vsip_csl::expr::Unary;
using vsip_csl::expr::Unary_functor;
// Scale implements a call operator that scales its input
// argument, and returns it by reference.
template
    <typename ArgumentBlockType> struct Scale :
        Unary_functor<ArgumentBlockType>
{
    Scale(ArgumentBlockType const &a,
         typename ArgumentBlockType::value_type s) :
        Unary_functor<ArgumentBlockType>(a), value(s)
    {}
    template <typename ResultBlockType>
    void apply(ResultBlockType &r) const
    {
        ArgumentBlockType const &a = this->arg();
        for (index_type i = 0; i != r.size(); ++i)
            r.put(i, a.get(i) * value);
    }
    typename
        ArgumentBlockType::value_type value;
};
// scale is a return-block optimised function returning an \
expression.
template <typename T,
         typename BlockType>
const_Vector<T,
            Unary<Scale, BlockType> const>
scale(const_Vector<T, BlockType> input, T value)
{
    Scale<BlockType> s(input.block(), value);
    Unary<Scale, BlockType> block(s);
    return const_Vector<T, Unary<Scale,
        BlockType> const>(block);
}
```

With that improvement, the `scale()` function in

```
Vector<> a(8);
Vector<> r = scale(a, 2.f);
```

is entirely evaluated during the assignment.

The `Unary_functor` above poses certain requirements on its function parameter. If they can't be met, we need to write a different functor. For example:

```
using vsip_csl::View_block_storage;
template <typename ArgumentBlockType>
struct Interpolator
{
public:
    typedef typename
        ArgumentBlockType::value_type value_type;
    typedef typename
        ArgumentBlockType::value_type result_type;
    typedef typename
        ArgumentBlockType::map_type map_type;
    static vsip::dimension_type const
        dim = ArgumentBlockType::dim;
    Interpolator(
        ArgumentBlockType const &a,
        Domain<ArgumentBlockType::dim> const &s)
        : argument_(a), size_(s) {}
    // Report the size of the new interpolated block
    length_type size() const
    { return size_.size(); }
    length_type size(dimension_type b,
                    dimension_type d) const
    {
        assert(b == ArgumentBlockType::dim);
        return size_[d].size();
    }
    map_type const &map() const
    {
        return argument_.map();
    }
    ArgumentBlockType const &arg() const
    {
        return argument_;
    }
    template <typename ResultBlockType> void
        apply(ResultBlockType &) const
    {
        std::cout << "apply interpolation !" << std::endl;
        // interpolate 'argument' into 'result'
    }
private:
    typename
        View_block_storage<ArgumentBlockType>::expr_type argument_;
```

```
Domain<ArgumentBlockType::dim> size_;
};
```

creates a new vector of different shape than the input Vector. To see the full requirements for the UnaryFunctor, see Section 6.7.3, “The Unary_functor class template”.

Here again, to write an interpolate() function that evaluates lazily, we need to return an *expression block* vector:

```
template <typename T, typename BlockType>
const_Vector<T, Unary<Interpolator, BlockType> const>
interpolate(const_Vector<T, BlockType> arg,
            Domain<1> const &size)
{
    typedef Unary<Interpolator, BlockType> expr_block_type;
    Interpolator<BlockType> interpolator(arg.block(), size);
    expr_block_type block(interpolator);
    return const_Vector<T, expr_block_type const>(block);
}
```

Now we can combine the above functions into a single expression:

```
Vector<float> a(8, 2.);
Vector<float> b = interpolate(scale(a, 2.f), 32);
```

The above demonstrates how to improve performance of an expression evaluation by using a technique that is a variant of the well-known *return value optimization*, where a copy operation (and a temporary object) may in certain cases be elided, if the result can be evaluated in-place.

4.5. Creating custom evaluators

In the previous section we have seen how to improve the expression evaluation by using the return-value optimization technique to avoid certain temporaries. However, there may be more that can be done to improve performance.

It may, for example, be possible to fuse multiple operations into one. Some platforms provide a fused "multiply-add" instruction that may be used, some algorithms are optimized for combined evaluation such as an FFT with a scalar multiplication, etc.

To be able to take advantage of those opportunities, we need to 'see' the whole expression at once, so we can dispatch the relevant sub-expression to such 'backends'.

For common cases, the library already performs this internally. However, sometimes users have their own optimized code that needs to be hooked into expression evaluation

In this section, we will develop an expression evaluator that matches the expression `interpolate(scale(a, 2.), 32)` from the last section.

Assignments are evaluated using the dispatch mechanism described in Chapter 3, “Using the Dispatch Framework”. To provide a custom evaluator for a particular expression assignment, it is thus necessary to specialize an Evaluator, using `op::assign<D>` as operation tag, and `be::user` as backend tag:

4.5.1. Specializing an evaluator for a particular expression type

To make Sourcery VSIPL++ use a custom evaluator, we need to specialize the `vsip_csl::dispatcher::Evaluator` template for the particular expression type we are interested in. Further, we need to model the evaluator concept.

The type of the expression can be discovered using `type_name()`:

```
std::cout
  << type_name(interpolate(scale(a, 2.),32))
  << std::endl;
```

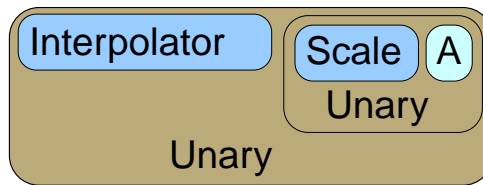
This yields (approximately):

```
vsip::Vector<
  float,
  vsip_csl::expr::Unary<
    example::Interpolator,
    vsip_csl::expr::Unary<
      example::Scale,
      vsip::Dense<1u, float, tuple<0u, 1u, 2u>,
        vsip::Local_map>,
      false>
    const>,
  false> const>
```

The block type is thus (with some details removed for clarity):

```
Unary<Interpolator,
  Unary<Scale, Dense<1>, false>
  const,
  false,
  const>
```

This can be visualized like this:



i.e., it is a `Unary` whose functor is an `Interpolator`. Its argument block, in turn, is a `Unary` whose functor is a `Scale`, and its argument block is a `Dense<1>`. This allows us to write a matching evaluator:

```
namespace vsip_csl {
namespace dispatcher {
  template <typename ResultBlockType,
            typename ArgumentBlockType>
  struct Evaluator<op::assign<1>,
                  be::user void(ResultBlockType &,
                                expr::Unary<Interpolator, \
expr::Unary<Scale, ArgumentBlockType>
const> const &)>
```

```

{
  typedef typename
    ArgumentBlockType::value_type value_type;
  typedef ResultBlockType LHS;
  typedef expr::Unary<Interpolator,
    expr::Unary<Scale, ArgumentBlockType> const>
    RHS;
  static bool const ct_valid = true;
  static bool rt_valid(LHS &, RHS const &)
  { return true;}
  static void exec(LHS &lhs, RHS const &rhs) {} // TBD
};

```

This evaluator will match the desired expression. As usual, the `ct_valid` and `rt_valid()` members can be used to refine the selection process.

4.5.2. Accessing the expression nodes.

The terminals in this expression are the block, the interpolator's target size, as well as the scale value. Once these three are available, the entire expression may be evaluated in a fused `scaled_interpolate()`, as shown here:

```

static void exec(LHS &lhs, RHS const &rhs)
{
  // rhs.arg() yields Unary<Scale, ArgumentBlockType>,
  // rhs.arg().arg() thus returns the terminal ArgumentBlockType \
  block...
  ArgumentBlockType &block = rhs.arg().arg();
  // ...and rhs.arg().operation() the Scale<ArgumentBlockType> \
  functor.
  value_type scale = rhs.functor().argument.functor().func.value;
  // rhs.operation() yields the Interpolator<Unary<Scale, ...> \
  functor.
  length_type new_size(rhs.operation().size(1, 0));
  // wrap terminal blocks in views for convenience, and evaluate.
  Vector<value_type, LHS> result(lhs);
  const_Vector<value_type, ArgumentBlockType const> argument(block);
  scaled_interpolate(result, argument, size, scale, new_size);
}

```

Chapter 5

Profiling

This reference explains how to compile a program with profiling statements enabled, how to use the profiling functions in order to investigate the execution timing of a program and finally how to interpret the resulting profile data.

5.1. Enabling Profiling

5.1.1. Configure and Compile Options

There are no configure options for profiling, instead it is enabled via compile-time options. However, to use profiling it is necessary to configure the library with a suitable high-resolution timer (refer to the Quickstart for details on this and other configuration options). For example,

```
--enable-timer=x86_64_tsc
```

Pre-built versions of the library enable a suitable timer for your system.

You may enable profiling per category by defining the corresponding macros for all the categories you want to enable. On many systems, these macros may be added to the `CXXFLAGS` variable in the project makefile, e.g.

```
CXXFLAGS=" -DVSIP_PROFILE_USER
          -DVSIP_PROFILE_MEMORY "
```

Table 5.1. Profiling Configuration Mask

Section	Description	Macro
memory	Memory management and block copying	VSIP_PROFILE_MEMORY
dispatch	Operation dispatch (incl. expression evaluation)	VSIP_PROFILE_DISPATCH
parallel	Parallel data I/O	VSIP_PROFILE_PARALLEL
func	Elementwise Functions	VSIP_PROFILE_FUNC
signal	Signal Processing	VSIP_PROFILE_SIGNAL
matvec	Matrix - Vector operations	VSIP_PROFILE_MATVEC
solver	Linear Algebra solvers	VSIP_PROFILE_SOLVER
user	User-defined Operations	VSIP_PROFILE_USER
all	Everything	VSIP_PROFILE_ALL

The above macros may be defined on the command line, or in a header, as long as that is included before any Sourcery VSIPL++ header.

5.1.2. Command Line Options

For programs that have been compiled with profiling enabled, the profiling mode and output file can be controlled from the command line. You may profile programs without modifying your source files using this method. Use this to choose the profiler mode:

```
--vsip-profile-mode=mode
```

where *mode* is either `accum` or `trace`.

Specify the path to the log file for profile output using:

```
--vsip-profile-output=/path/to/logfile
```

The second option defaults to the standard output on most systems, so it may be omitted if that is desirable.

The profiling command line options control profiling for the entire program execution. For finer grain control, such as enabling profiling during a specific portion of the program, or to mix different profiling modes, explicit profiling objects can be created.

5.2. Using the Profiler

5.2.1. Profiling Objects

The `Profile` object is used to enable profiling during the lifetime of the object. When created, it takes arguments to indicate the output file and the profiling mode (trace or accumulate). When destroyed (i.e. goes out of scope or is explicitly deleted), the profile data is written to the specified output file. For example:

```
profile::Profile profile("profile.txt", profile::pm_accum);
```

During the lifetime of the `Profile` object, timing data is stored through a simple interface provided by the `Scope` object. These objects are used to profile library operations for the different areas mentioned in Table 5.1, “Profiling Configuration Mask” above. Any `Scope` objects defined in user programs fall into the ‘user’ category.

The declaration of an instance of this object starts a timer and when it is destroyed, the timer is stopped. The timing data is subsequently reported when the `Profile` object is destroyed. For example:

```
profile::Scope<profile::user> scope("Scope  
Name", op_count);
```

The `profile::user` template argument indicates that this scope falls into the *user* category (and thus can be enabled with `-DVSIP_PROFILE_USER`). The first constructor argument is the name that will be used to display the scope's performance data in the log file (Section 5.3.2, “Scope names” describes the names used internally by the library.) The second parameter, `op_count`, is an optional unsigned integer specifying an estimate of the total number of operations (floating point or otherwise) performed. This is used by the profiler to compute the rate of computation. Without it, the profiler will still yield useful timing data, but the average rate of computation will be shown as zero in the log.

Creating a `Scope` object on the stack is the easiest way to control the region it will profile. For example, from within the body of a function (or as the entire function), use this to define a region of interest:

```
{  
    profile::Scope<profile::user> scope("Main computation:"); //  
    perform main computation // ... }  
}
```

The closing brace causes `scope` to go out of scope, logging the amount of time spent doing the computation.

5.2.2. Profiler Modes

In `trace` mode, the start and stop times where scopes begin and end are stored as profile data. The log will present these events in chronological order. This mode is preferred when a highly detailed view of program execution is desired.

In `accum` (accumulate) mode, the start and stop times are subtracted to compute the time spent in a scope and the cumulative sum of these durations are stored as profile data. The log will indicate the total amount of time spent in each scope. This mode is desirable when investigating a specific function's average performance.

5.3. Profiler Output

5.3.1. Log File Format

The profiler outputs a small header at the beginning of each log file which is the same accumulate and trace modes. The data that follows the header is different depending on the mode. The header describes the profiling mode used, the low-level timer used to measure clock ticks and the number of clock ticks per second.

5.3.1.1. Accumulate mode

```
# mode: pm_accum # timer: x86_64_tsc_time #
   clocks_per_sec: 3591375104 # # tag : total ticks : num \
calls : op
   count : mops
```

The respective columns that follow the header are:

tag	A descriptive name of the operation. This is either a name used internally or specified by the user.
total ticks	The time spent in the scope in processor ticks.
num calls	The number of times the scope was entered.
op count	The number of operations performed per scope.
mops	The calculated performance figure in millions of operations per second. $(\text{num_calls} \times \text{op_count} \times 10^{-6}) / (\text{total_ticks} / \text{clocks_per_sec})$

5.3.1.2. Trace mode

```
# mode: pm_trace # timer: x86_64_tsc_time # clocks_per_sec:
   3591375104 # # index : tag : ticks : open id : op count
```

The respective columns that follow the header are:

index	The entry number, beginning at one.
tag	A descriptive name of the operation. This is either a name used internally or specified by the user.
ticks	The current reading from the processor clock.

open id If zero, indicates the start of a scope. If non-zero, this indicates the end of an scope and refers to the index of corresponding start of the scope.

op count The number of operations performed per scope, or zero to indicate the end of an scope.

Note that the timings expressed in 'ticks' may be converted to seconds by dividing by the 'clocks_per_second' constant in the header.

5.3.2. Scope names

Sourcery VSIPL++ uses the following names for profiling objects and functions within the library. These names are readable text containing information that varies depending on the operation being profiled.

5.3.2.1. Signal Processing and Matrix Vector Operations

These operations follow this general format:

```
OPERATION
      [DIM] DATATYPE SIZE
```

OPERATION gives the object or function name, including direction for FFTs.

DIM is the number of dimensions (when needed).

DATATYPE describes the data types involved in the operation. FFTs have two listed, describing both the input type as well as the output type, which may be different. See Table 5.2, “Data Type Names” below.

SIZE is expressed by giving the number of elements in each dimension.

The specific operations profiled at this time are:

```
Convolution [1D|2D] T SIZE
Correlation [1D|2D] T SIZE
Fft 1D [Inv|Fwd] I-O [by_ref|by_val] SIZE
Fftm 2D [Inv|Fwd] I-O [by_ref|by_val] SIZE
Fir T SIZE
Iir T SIZE
dot T SIZE
cvjdot T SIZE
trans T SIZE
herm T SIZE
kron T SIZE_A SIZE_B
outer T SIZE
gemp T SIZE
gems T SIZE
cumsum T SIZE
modulate T SIZE
```

In all cases, data types T, I and O above are expressed using a notation similar to the BLAS/LAPACK convention as in the following table:

Table 5.2. Data Type Names

	Views	Scalars
single precision real	S	s
single precision complex	C	c
double precision real	D	d
double precision complex	Z	z

5.3.2.2. Elementwise Functions

Element-wise expression tags use a slightly different format:

```
EVALUATOR DIM EXPR SIZE
```

The EVALUATOR indicates which VSIPL++ evaluator was dispatched to compute the expression.

DIM indicates the dimensionality of the expression.

EXPR is mnemonic of the expression shown using prefix notation, i.e.

```
operator(operand, ...)
```

Each operand may be the result of another computation, so expressions are nested, the parenthesis determining the order of evaluation.

SIZE is expressed by giving the number of elements in each dimension.

At this time, these evaluators are profiled:

Expr_Loop - generic loop-fusion evaluator.

Expr_SIMD_Loop - SIMD loop-fusion evaluator.

Expr_Copy - optimized data-copy evaluator.

Expr_Trans - optimized matrix transpose evaluator.

Expr_Dense - evaluator for dense, multi-dimensional expressions. Converts them into corresponding 1-dim expressions that are re-dispatched.

Expr_SAL_* - evaluators for dispatch to the SAL vendor math library.

Expr_IPP_* - evaluators for dispatch to the IPP vendor math library.

Expr_SIMD_* - evaluators for dispatch to the builtin SIMD routines (with the exception of Expr_SIMD_Loop, see above).

For SAL, IPP and SIMD, the asterisk (*) denotes the specific function invoked.

Chapter 6

Benchmarking

Abstract

This chapter describes how to build and run the Saurcery VSIPL++ benchmark suite in order to determine how the library performs on a given platform.

This chapter explains how to build and run the performance benchmarks supplied with Sourcery VSIPL++.

6.1. Overview

The following tables describe the different benchmarks available currently. They are organized by type of operation, to allow more easy cross-referencing with the specification.

Table 6.1. Sourcery VSIPL++ Benchmark Descriptions

Operation	Source File
math, functions, elementwise - vector multiply	vmul.cpp
math, functions, elementwise - vector-matrix multiply	vmmul.cpp
math, functions, reductions - maximum value	maxval.cpp
math, functions, reductions - sum of values	sumval.cpp
math, matvec - matrix-matrix products	prod.cpp
math, matvec - matrix-matrix products, variations	prod_var.cpp
math, matvec - matrix copy, transpose	mcopy.cpp
math, matvec - vector dot product	dot.cpp
signal - convolution	conv.cpp
signal - correlation	corr.cpp
signal - fast convolution	fastconv.cpp
signal - Fast Fourier Transform	fft.cpp
signal - Fast Fourier Transform, multiple	fftm.cpp
signal - Finite Impulse Response filter	fir.cpp
view, vector, assign - memory write bandwidth	memwrite.cpp

All of the above source files are located in `share/benchmarks/` in the top-level install directory.

6.2. Building the Benchmarks

The benchmark sources may be set up for building as part of a Sourcery VSIPL++ Workspace. See Section 3.1, “Using a Sourcery VSIPL++ Example Workspace” for instructions on how to prepare that.

The benchmarks will be placed in the `benchmarks/` subdirectory of the newly created workspace. You may compile the benchmarks by simply invoking

```
> make
```

or, if you want to override the variant value,

```
> make variant=VARIANT
```

To rebuild a certain benchmark, specify its name. For example:

```
> make vmul
```

6.2.1. Optimization Settings

By default, the build process will use compilation and link flags provided by the Sourcery VSIPL++ installation. In order to experiment with different compiler flags, you may either override the default values, or add to them. Compilation flags are passed via the CXXFLAGS variable, link flags via the LDFLAGS variable. To redefine CXXFLAGS, edit the Makefile:

```
CXXFLAGS:= -your -flags
```

To add flags, use

```
CXXFLAGS:= $(CXXFLAGS) -your -flags
```

6.3. Running Benchmarks

Benchmarks are invoked as follows:

```
> benchmark -case-number [-option[ -option[...]]]
```

Case numbers are defined individually for each benchmark and represent various combinations of algorithms and parameters used for a given performance measurement. Valid test numbers begin at one. Most benchmarks utilize zero to display a list of valid case numbers. For example:

```
> vmul -0
vmul -- vector multiplication
single-precision:
  Vector-Vector:
    -1 -- Vector<          float > * Vector<          float>
    -2 -- Vector<complex<float> > * Vector<complex<float> >
double-precision:
    -101 -- Vector<          double > * Vector<          float>
    -102 -- Vector<complex<double> > * Vector<complex<float> >
```

While performance benchmarks may be run simply by specifying the case number, several options are provided making it possible to control a variety of useful parameters. These parameters affect either the way the test is run or the type of output provided in the results. A summary of the most useful options is provided below:

Controlling the range of problem sizes

`-start M` Sets the starting problem size to 2^M . Defaults to 2 (4 points).

`-stop M` Sets the stopping problem size to 2^M . Defaults to 21 (2097152 points).

Controlling the samples taken

`-samples S` Sets the number of samples taken to *S*. Defaults to one. When $S > 2$, the median value is reported.

`-ms time` Sets the goal time that each measurement should take, in hundredths-of-a-second. Defaults to 25 (250 ms).

Benchmark specific parameters

`-param value` Sets the user-specified parameter, where *value* is used in some instances to override a default value of a parameter, such as the number of rows or columns in an input matrix for example. The effect, if any, is dependent on each individual benchmark.

Reporting

`-pts` Report millions of points per second (MPT/s). This is the default.

`-ops` Reports millions of operations per second (MOP/s).

`-iob` Reports millions of input/output operations per second (MB/s) (by summing read and write `iob_per_point`).

`-all` Reports all three statistics: MPT/s, MOP/s, MB/s

6.4. Benchmark Output

The benchmark output depends on the command line options, but typically includes some meta information on the benchmark (name, ops/point, etc) and individual measurements for each problem size.

The header information, denoted by lines begging with “ # ”, contains three important factors that are used to convert timing data into other meaningful units. The number of floating point operations is shown as `ops_per_point` and the number of reads or writes to and from memory are shown as `riob_per_point` and `wiob_per_point` respectively.

Following the header information are performance results. Each line contains data for a certain problem size (number of points), which is given in the first column.

The second column contains the measured (or median) values calculated from the timing measurements. The default is in points-per-second as indicated in the header under “metric”. Alternatively, the values are in units as requested with the `-pts`, `-ops`, `-iob` option.

In other cases, three columns of measurements follow the size given in the first column. The values listed vary depending on the options specified, as outlined below:

`-all` Displays points per second, operations per second and the sum of the memory reads and writes per second (MPT/s, MOP/s, MB/s).

`-pts`, `-ops`, `-iob` Displays one of points per second, operations per second and the sum of the memory reads and writes per second, as requested.

`-samples S` With `-all`, three columns will be displayed, each containing the median value of the respective measurement. Without `-all`, the second column will contain the median value and columns three and four will contain the minimum and maximum value for the selected measurement. Note: *S* must be

greater than two in order to display the minimum and maximum values for `-pts`, `-ops` or `-iob`.

6.4.1. Examples

This example shows a very simple benchmark for vector-vector multiplication using complex values, defaulting to units of “points-per-second”:

```
> vmul -2
# what : t_vmul1
# ops_per_point(1) : 6
# riob_per_point(1): 16
# wiob_per_point(1): 8
# metric : pts_per_sec
# start_loop : 2981969
4 60.606903
8 123.195221
16 173.855408
32 207.837997
64 232.163071
...
```

The output is truncated, but continues on up until 2^{21} points per vector.

To measure operations per second instead, use:

```
> vmul -2 -ops
# what : t_vmul1
# ops_per_point(1) : 6
# riob_per_point(1): 16
# wiob_per_point(1): 8
# metric : ops_per_sec
# start_loop : 2973904
4 377.566650
8 765.744446
16 1055.679321
32 1261.269653
64 1425.231567
... \
```

To measure ops/sec, with the median of 3 samples of 0.5 seconds in duration each:

```
> vmul -2 -ops -samples 3 -ms 50
# what : t_vmul1
# ops_per_point(1) : 6
# riob_per_point(1): 16
# wiob_per_point(1): 8
# metric : ops_per_sec
# start_loop : 5934442
4 409.272583 398.208191 413.359711
8 854.137939 811.087402 854.964539
16 1132.262939 1087.544800 1137.489502
32 1317.902710 1297.148560 1342.707886
64 1483.941650 1453.872192 1501.823242 \
```

Note that this option is most often used when error bars are desired for plotting the performance data.

Part II. Example Application

Table of Contents

7. Fast Convolution	59
7.1. Fast Convolution	60
7.2. Serial Optimization: Temporal Locality	64
7.3. Performing I/O with User-Specified Storage	66
7.4. Performing I/O with Direct Data Access	69
8. Parallel Fast Convolution	71
8.1. Parallel Fast Convolution	72
8.2. Improving Parallel Temporal Locality	75
8.3. Performing I/O	81

Chapter 7

Fast Convolution

Abstract

This chapter describes how to create and run a serial VSIP++ program with Saurcery VSIP++ that performs fast convolution. You can modify this program to develop your own serial applications.

This chapter explains how to use Sourcery VSIPL++ to perform *fast convolution* (a common signal-processing kernel). First, you will see how to compute fast convolution using VSIPL++'s multiple FFT (Fftm) and vector-matrix multiply operations. Then, you will learn how to optimize the performance of the implementation.

7.1. Fast Convolution

Fast convolution is the technique of performing convolution in the frequency domain. In particular, the time-domain convolution $f * g$ can be computed as $F \cdot G$, where F and G are the frequency-domain representations of the signals f and g . A time-domain signal consisting of n samples can be converted to a frequency-domain signal in $O(n \log n)$ operations by using a Fast Fourier Transform (FFT). Substantially fewer operations are required to perform the frequency-domain operation $F \cdot G$ than are required to perform the time-domain operation $f * g$. Therefore, performing convolutions in the frequency domain can be substantially faster than performing the equivalent computations in the time domain, even taking into account the cost of converting from the time domain to the frequency domain.

One practical use of fast convolution is to perform the pulse compression step in radar signal processing. To increase the effective bandwidth of a system, radars will transmit a frequency modulated "chirp". By convolving the received signal with the time-inverse of the chirp (called the "replica"), the total energy returned from an object can be collapsed into a single range cell. Fast convolution is also useful in many other contexts including sonar processing and software radio.

In this section, you will construct a program that performs fast convolution on a set of time-domain signals stored in a matrix. Each row of the matrix corresponds to a single signal, or "pulse". The columns correspond to points in time. So, the entry at position (i, j) in the matrix indicates the amplitude and phase of the signal received at time j for the i th pulse.

The first step is to declare the data matrix, the vector that will contain the replica signal, and a temporary matrix that will hold the results of the computation:

```
// Parameters.
length_type npulse = 64; // number of pulses
length_type nrange = 256; // number of range cells
// Views.
typedef complex<float> value_type;
Vector<value_type> replica(nrange);
Matrix<value_type> data(npulse, nrange);
Matrix<value_type> tmp (npulse, nrange);
```

For now, it is most convenient to initialize the input data to zero. (In Section 7.3, "Performing I/O with User-Specified Storage", you will learn how to perform I/O operations so that you can populate the matrix with external data.)

In C++, you can use the constructor syntax $T()$ to perform "default initialization" of a type $T()$. The default value for any numeric type (including complex numbers) is zero. Therefore, the expression $value_type()$ indicates the complex number with zero as both its real and imaginary components. In the VSIPL++ API, when you assign a scalar value to a view (a vector, matrix, or tensor), all elements of the view are assigned the scalar value. So, the code below sets the contents of both the data matrix and replica vector to zero:

```
data = value_type();
replica = value_type();
```

The next step is to define the FFTs that will be performed. Typically (as in this example) an application performs multiple FFTs on inputs with the same size. Since performing an FFT requires that some set-up be performed before performing the actual FFT computation, it is more efficient to set up the FFT just once. Therefore, in the VSIPL++ API, FFTs are objects, rather than operators. Constructing the FFT performs the necessary set-up operations.

Because VSIPL++ supports a variety of different kinds of FFT, FFTs are themselves template classes. The parameters to the template allow you to indicate whether to perform a forward (time-domain to frequency-domain) or inverse (frequency-domain to time-domain) FFT, the type of the input and output data (i.e., whether complex or real data is in use), and so forth. Then, when constructing the FFT objects, you indicate the size of the FFT. In this case, you will need both an ordinary FFT (to convert the replica data from the time domain to the frequency domain) and a "multiple FFT" to perform the FFTs on the rows of the matrix. (A multiple FFT performs the same FFT on each row or column of a matrix.) So, the FFTs required are:

```
// A forward Fft for computing the frequency-domain version of
// the replica.
typedef Fft<const_Vector, value_type,
           value_type, fft_fwd, by_reference>
    for_fft_type;
for_fft_type for_fft(Domain<1>(nrange), 1.0);
// A forward Fftm for converting the time-domain data matrix to the
// frequency domain.
typedef Fftm<value_type, value_type, row, fft_fwd, by_reference>
    for_fftm_type;
for_fftm_type for_fftm(Domain<2>(npulse, nrange), 1.0);
// An inverse Fftm for converting the frequency-domain data back to
// the time-domain.
typedef Fftm<value_type, value_type, row, fft_inv,
            by_reference>
    inv_fftm_type;
inv_fftm_type
    inv_fftm(Domain<2>(npulse, nrange), 1.0/(nrange));
```

Before performing the actual convolution, you must convert the replica to the frequency domain using the FFT created above. Because the replica data is a property of the chirp, we only need to do this once; even if the radar system runs for a long time, the converted replica will always be the same. VSIPL++ FFT objects behave like functions, so you can just "call" the FFT object:

```
for_fft(replica);
```

Now, you are ready to perform the actual fast convolution operation! You will use the forward and inverse multiple-FFT objects you've already created to go into and out of the frequency domain. While in the frequency domain, you will use the `vmmul` operator to perform a vector-matrix multiply. This operator multiplies each row (dimension zero) of the frequency-domain matrix by the replica. The `vmmul` operator is a template taking a single parameter which indicates whether the multiplication should be performed on rows or on columns. So, the heart of the fast convolution algorithm is just:

```
// Convert to the frequency domain.
for_fftm(data, tmp);
// Perform element-wise multiply for each pulse.
tmp = vmmul<0>(replica, tmp);
// Convert back to the time domain.
inv_fftm(tmp, data);
```

A complete program listing is show below. You can copy this program directly into your editor and build it. (You may notice that there are a few things in the complete listing not discussed above, including in particular, initialization of the library.)

```

\
/*****
  Included Files
*****/

#include <vsip/initfin.hpp>
#include <vsip/support.hpp>
#include <vsip/signal.hpp>
#include <vsip/math.hpp>

using namespace vsip;

/*****
  Main Program
*****/

int
main(int argc, char** argv)
{
    // Initialize the library.
    vsipl vpp(argc, argv);

    typedef complex<float> value_type;

    // Parameters.
    length_type npulse = 64; // number of pulses
    length_type nrange = 256; // number of range cells

    // Views.
    Vector<value_type> replica(nrange);
    Matrix<value_type> data(npulse, nrange);
    Matrix<value_type> tmp(npulse, nrange);

    // A forward Fft for computing the frequency-domain version of
    // the replica.
    typedef Fft<const_Vector, value_type, value_type, fft_fwd, \
by_reference>
    for_fft_type;
    for_fft_type for_fft (Domain<1>(nrange), 1.0);

    // A forward Fftm for converting the time-domain data matrix to the
    // frequency domain.
    typedef Fftm<value_type, value_type, row, fft_fwd, by_reference>
    for_fftm_type;
    for_fftm_type for_fftm(Domain<2>(npulse, nrange), 1.0);

    // An inverse Fftm for converting the frequency-domain data back to

```

```
// the time-domain.
typedef Fftm<value_type, value_type, row, fft_inv, by_reference>
    inv_fftm_type;
inv_fftm_type inv_fftm(Domain<2>(npulse, nrange), 1.0/(nrange));

// Initialize data to zero.
data      = value_type();
replica   = value_type();

// Before fast convolution, convert the replica to the the
// frequency domain
for_fft(replica);

// Perform fast convolution.

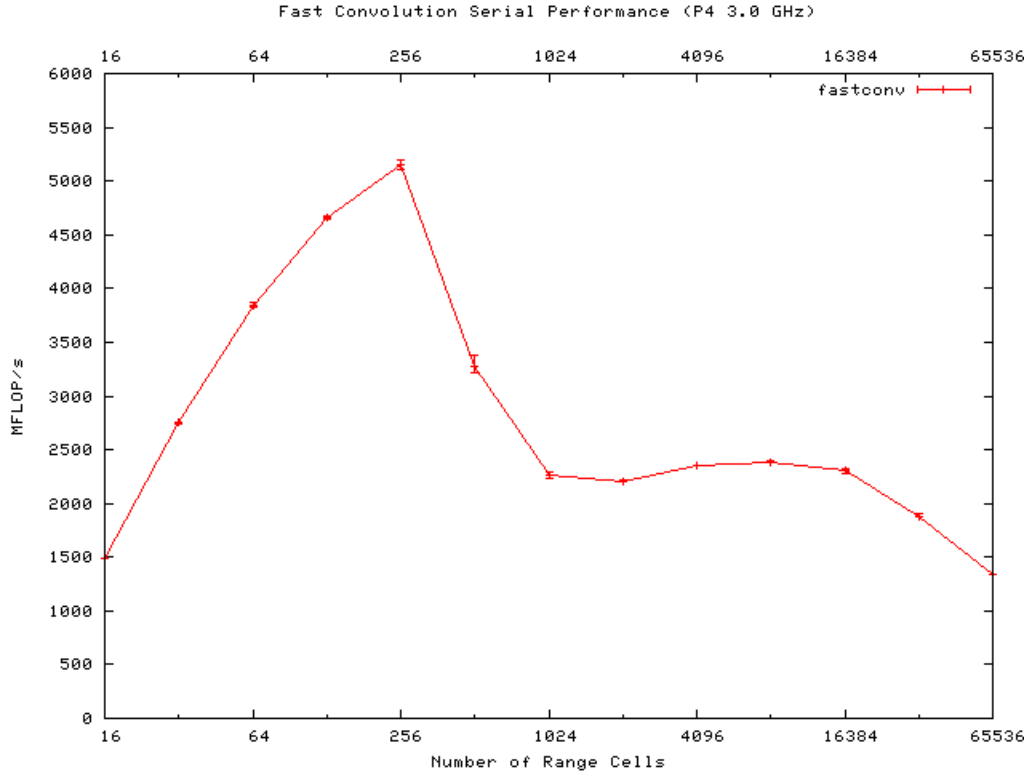
// Convert to the frequency domain.
for_fftm(data, tmp);

// Perform element-wise multiply for each pulse.
tmp = vmmul<0>(replica, tmp);

// Convert back to the time domain.
inv_fftm(tmp, data);
}

\
```

The following figure shows the performance in MFLOP/s of fast convolution on a 3.06 GHz Pentium Xeon processor as the number of range cells varies from 16 to 65536.



7.2. Serial Optimization: Temporal Locality

In this section, you will learn how to improve the performance of fast convolution by improving *temporal locality*, i.e., by making accesses to the same memory locations occur near the same time.

The code in Section 7.1, “Fast Convolution” performs a FFT on each row of the matrix. Then, after all the rows have been processed, it multiplies each row of the matrix by the `replica`. Suppose that there are a large number of rows, so that data is too large to fit in cache. In that case, while the results of the first FFT will be in cache immediately after the FFT is complete, that data will likely have been purged from the cache by the time the vector-matrix multiply needs the data.

Explicitly iterating over the rows of the matrix (performing a forward FFT, elementwise multiplication, and an inverse FFT on each row before going on to the next one) will improve temporal locality. You can use this approach by using an explicit loop, rather than the implicit parallelism of `Fftm` and `vmmul`, to take better advantage of the cache.

You must make a few changes to the application in order to implement this approach. Because the application will be operating on only a single row at a time, `Fftm` must be replaced with the simpler `Fft`. Similarly, `vmmul` must be replaced with `*`, which performs element-wise multiplication of its operands. Finally, `tmp` can now be a vector, rather than a matrix. (As a consequence, in addition to being faster, this new version of the application will require less memory.) Here is the revised program:

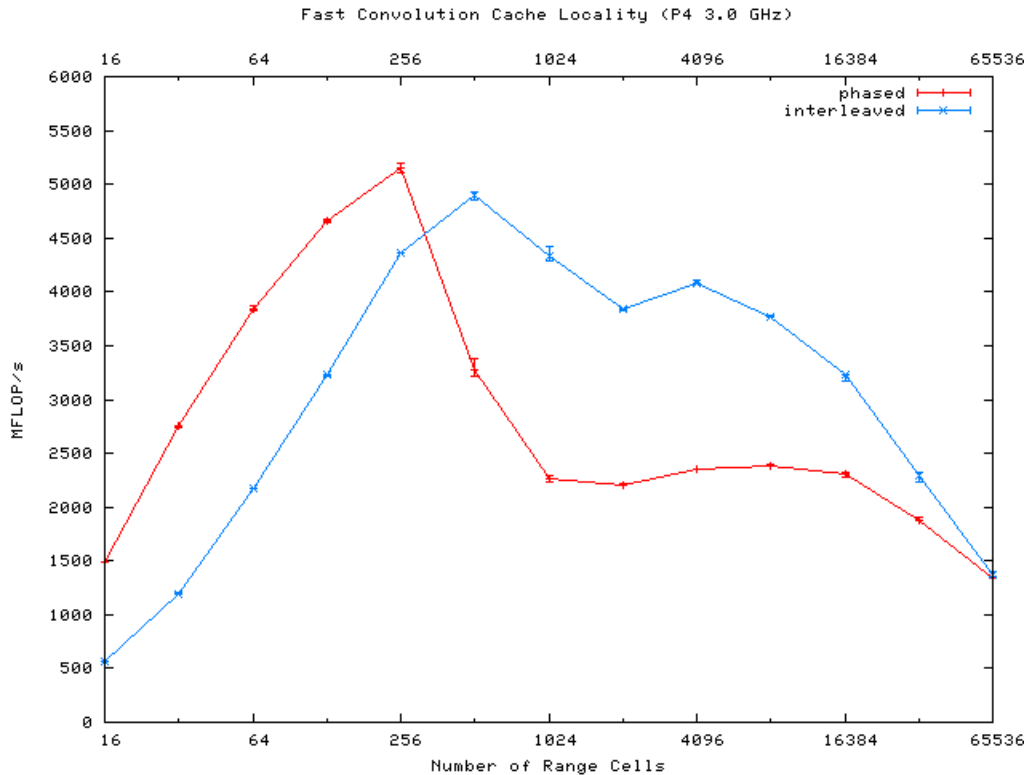
```
// Create the data cube.
Matrix<value_type> data(npulse, nrange);
Vector<value_type> tmp(nrange);
// tmp is now a vector
// Create the pulse replica
Vector<value_type> replica(nrange);
```

```

// Define the FFT typedefs.
typedef Fft<const_Vector,
          value_type, value_type, fft_fwd, by_reference>
  for_fft_type;
typedef Fft<const_Vector, value_type, value_type, fft_inv, \
by_reference>
  inv_fft_type;
// Create the FFT objects.
for_fft_type for_fft(Domain<1>(nrange), 1.0);
inv_fft_type inv_fft(Domain<1>(nrange), 1.0/(nrange));
// Initialize data to zero
data = value_type();
replica = value_type();
// Before fast convolution, convert the replica into the
// frequency domain
for_fft(replica);
// Perform fast convolution:
for (index_type r=0; r < nrange; ++r)
{
  for_fft(data.row(r), tmp);
  tmp *= replica;
  inv_fft(tmp, data.row(r));
}

```

The following graph shows that the new "interleaved" formulation is faster than the original "phased" approach for large data sets. For smaller data sets (where all of the data fits in the cache anyhow), the original method is faster because performing all of the FFTs at once is faster than performing them one by one.



7.3. Performing I/O with User-Specified Storage

The previous sections have ignored the acquisition of actual sensor data by setting the input data to zero. This section shows how to initialize data before performing the fast convolution.

To perform I/O with external routines (such as the POSIX `read` and `write` functions) it is necessary to obtain a pointer to the raw data used by Sourcery VSIPL++. Sourcery VSIPL++ provides two ways to do this: you may use either *user-defined storage* or *external data access*. In this section you will use user-defined storage to perform I/O. Later, in Section 7.4, “Performing I/O with Direct Data Access” you will see how to use external data access for I/O.

VSIPL++ allows you to create a block with user-specified storage by giving VSIPL++ a pointer to previously allocated data when the block is created. This block is just like a normal block, except that it now has two states: “admitted” and “released”. When the block is admitted, the data is owned by VSIPL++ and the block can be used with any VSIPL++ functions. When the block is released, the data is owned by you allowing you to perform operations directly on the data. The states allow VSIPL++ to potentially reorganize data for higher performance while it is admitted. (Attempting to use the pointer while the block is admitted, or use the block while it is released will result in unspecified behavior!)

The first step is to allocate the data manually.

```
std::vector<value_type>
    buffer(npulse*nrange);
```

Next, you create a VSIPL++ Dense block, providing it with the pointer.

```
Dense<2, value_type>
    block(Domain<2>(nrange, npulse), &buffer.front());
```

Since the pointer to data does not encode the data dimensions, it is necessary to create the block with explicit dimensions.

Finally, you create a VSIPL++ view that uses this block.

```
Matrix<value_type> data(block);
```

The view determines its size from the block, so there is no need to specify the dimensions again.

Now you're ready to perform I/O. When a user-specified storage block is first created, it is released.

```
... setup IO
...
read(..., &buffer.front(),
      sizeof(value_type)*nrange*npulse);
... check for errors (of course!)
...
```

Finally, you need to admit the block so that it and the view can be used by VSIPL++.

```
data.block().admit(true);
```

The `true` argument indicates that the data values should be preserved by the admit. In cases where the values do not need to be preserved (such as admitting a block after output I/O has been performed and before the block will be overwritten by new values in VSIPL++) you can use `false` instead.

After admitting the block, you can use `data` as before to perform fast convolution. Here is the complete program, including I/O to output the result after the computation.

```

\
/*****
  Included Files
*****/

#include <vsip/initfin.hpp>
#include <vsip/support.hpp>
#include <vsip/signal.hpp>
#include <vsip/math.hpp>
#include <vector>

using namespace vsip;

/*****
  Main Program
*****/

int
main(int argc, char** argv)
{
    // Initialize the library.
    vsipl vpp(argc, argv);

    typedef complex<float> value_type;

    // Parameters.
    length_type npulse = 64; // number of pulses
    length_type nrange = 256; // number of range cells

    // Allocate data.
    std::vector<value_type> ptr(npulse*nrange);

    // Blocks.
    Dense<2, value_type> block(Domain<2>(npulse, nrange), \
&data.front());

    // Views.
    Vector<value_type> replica(nrange);
    Matrix<value_type> ptr(block);
    Matrix<value_type> tmp(npulse, nrange);

    // A forward Fft for computing the frequency-domain version of
    // the replica.
    typedef Fft<const_Vector, value_type, value_type, fft_fwd, \
by_reference>
    for_fft_type;
    for_fft_type for_fft (Domain<1>(nrange), 1.0);

```

```

// A forward Fftm for converting the time-domain data matrix to the
// frequency domain.
typedef Fftm<value_type, value_type, row, fft_fwd, by_reference>
    for_fftm_type;
for_fftm_type for_fftm(Domain<2>(npulse, nrange), 1.0);

// An inverse Fftm for converting the frequency-domain data back to
// the time-domain.
typedef Fftm<value_type, value_type, row, fft_inv, by_reference>
    inv_fftm_type;
inv_fftm_type inv_fftm(Domain<2>(npulse, nrange), 1.0/(nrange));

// Initialize data to zero.
data    = value_type();
replica = value_type();

// Before fast convolution, convert the replica to the the
// frequency domain
for_fft(replica);

// Read input.
view.block().release(false);
size_t size = read(0, &data.front(), \
sizeof(value_type)*nrange*npulse);
assert(size == sizeof(value_type)*nrange*npulse);
view.block().admit(true);

// Perform fast convolution.

// Convert to the frequency domain.
for_fftm(data, tmp);

// Perform element-wise multiply for each pulse.
tmp = vmmul<0>(replica, tmp);

// Convert back to the time domain.
inv_fftm(tmp, data);

// Write output.
view.block().release(true);
size_t size = write(0, &data.front(), \
sizeof(value_type)*nrange*npulse);
assert(size == sizeof(value_type)*nrange*npulse);
view.block().admit(false);
}

\

```

The program also includes extra `release()` and `admit()` calls before and after the input and output I/O sections. For this example, they are not strictly necessary. However they are good practice because they make it clear in the program where the block is admitted and released. They also make it easier to modify the program to process data repeatedly in a loop, and to use separate buffers for input and output data. Because the extra calls have a `false` update argument, they incur no overhead.

7.4. Performing I/O with Direct Data Access

In this section, you will use *Direct Data Access* to get a pointer to a block's data. You can use this method with any block, even if the block does not use user-specified storage. The external data access method is useful in contexts where you cannot control how the block is allocated. For example, in this section, you will create a utility routine for I/O that works with any matrix or vector, even if it was not created with user-defined storage.

To access a block's data with external data access, you create a `dda::Data` object.

```
dda::Data<block_type, dda::inout, layout_type> data(block);
```

`dda::Data` is a class template that takes template parameters to indicate the block type `block_type` and the requested layout `layout_type`. The constructor takes two parameters: the block being accessed, and the type of synchronization necessary.

The `layout_type` parameter is a specialized `Layout` class template that determines the layout of data that `dda::Data` provides. If no type is given, the natural layout of the block is used. However, in some cases you may wish to specify row-major or column-major layout.

The `Layout` class template takes 4 parameters to indicate dimensionality, dimension-ordering, packing format, and complex storage format (if complex). In the example below you will use the `layout_type` to request the data access to be dense, row-major, with interleaved real and imaginary values. This layout corresponds to a common storage format used for binary files storing complex data.

The synchronization type is analogous to the update flags for `admit()` and `release()`. `dda::in` indicates that the block and pointer should be synchronized when the `dda::Data` object is created (like `admit(true)`) `dda::out` indicates that the block and pointer should be synchronized when the `dda::Data` object is destroyed (like `release(true)`) `dda::inout` indicates that the block and pointer should be synchronized at both points.

Once the object has been created, the pointer can be accessed with the `ptr` method.

```
value_type* ptr = data.ptr();
```

The pointer provided is valid only during the life of the `dda::Data` object. Moreover, the block referred to by the `dda::Data` object must not be used during this period.

Using these capabilities together, you can create a routine to perform I/O into a block. This routine will take two arguments: a filename to read, and a view in which to store the data. The amount of data read from the file will be determined by the view's size.

```
template <typename ViewT>
void
read_file(ViewT view, char const* filename)
{
    using namespace vsip;
    using vsip::impl::Row_major;

    dimension_type const dim = ViewT::dim;
    typedef typename ViewT::block_type block_type;
    typedef typename ViewT::value_type value_type;
```

```
typedef Layout<dim, typename Row_major<dim>::type, dense, \
interleaved_complex>
    layout_type;

dda::Data<block_type, dda::out, layout_type> data(view.block());

std::ifstream ifs(filename);

ifs.read(reinterpret_cast<char*>(data.ptr()), view.size() * \
sizeof(value_type));
}

\
```

Chapter 8

Parallel Fast Convolution

Abstract

This chapter describes how to create and run parallel VSIPL++ programs with Saurcery VSIPL++. You can modify the programs to develop your own parallel applications.

This chapter explains how to use Sourcery VSIPL++ to perform parallel computations. You will see how to transform the fast convolution program from the previous chapter to run in parallel. First you will convert the `Fftm` based version. Then you will convert the improved cache locality version. Finally, you will learn how to handle input and output when working in parallel.

8.1. Parallel Fast Convolution

The first fast convolution program in the previous chapter makes use of two implicitly parallel operators: `Fftm` and `vmmul`. These operators are implicitly parallel in the sense that they process each row of the matrix independently. If you had enough processors, you could put each row on a separate processor and then perform the entire computation in parallel.

In the VSIPL++ API, you have explicit control of the number of processors used for a computation. Since the default is to use just a single processor, the program in Section 7.1, “Fast Convolution” will not run in parallel, even on a multi-processor system. This section will show you how to use *maps* to take advantage of multiple processors. Using a map tells Sourcery VSIPL++ to distribute a single block of data across multiple processors. Then, Sourcery VSIPL++ will automatically move data between processors as necessary.

The VSIPL++ API uses the Single-Program Multiple-Data (SPMD) model for parallelism. In this model, every processor runs the same program, but operates on different sets of data. For example, in the fast convolution example, multiple processors perform FFTs at the same time, but each processor handles different rows in the matrix.

Every map has both compile-time and run-time properties. At compile-time, you specify the *distribution* that will be applied to each dimension. In this example, you will use a *block distribution* to distribute the rows of the matrix. A block distribution divides a view into contiguous chunks. For example, suppose that you have a 4-processor system. Since there are 64 rows in the matrix `data`, there will be 16 rows on each processor. The block distribution will place the first 16 rows (rows 0 through 15) on processor 0, the next 16 rows (rows 16 through 31) on processor 1, and so forth. You do not want to distribute the columns of the matrix at all, so you will use a *whole distribution* for the columns.

Although the distributions are selected at compile-time, the number of processors to use in each dimension is not specified until run-time. By specifying the number of processors at run-time, you can adapt your program to the configuration of the machine on which your application is running. The VSIPL++ API provides a `num_processors` function to tell you the total number of processors available. Of course, since each row should be kept on a single processor, the number of processors used in the column dimension is just one. So, here is the code required to create the map:

```
typedef Map<Block_dist, Whole_dist> map_type;
map_type map = map_type(/*rows=*/num_processors(),
                       /*columns=*/1);
```

Next, you have to tell Sourcery VSIPL++ to use this map for the relevant views. Every view has an underlying *block*. The block indicates how the view's data is stored. Until this point, you have been using the default `Dense` block, which stores data in a contiguous array on a single processor. Now, you want to use a contiguous array on *multiple* processors, so you must explicitly distribute the block. Then, when declaring views, you must explicitly indicate that the view should use the distributed block:

```
typedef Dense<2, value_type, row2_type, map_type>
    block_type;
typedef Matrix<value_type, block_type> view_type;
```

```
view_type data(npulse, nrange, map);
view_type tmp(npulse, nrange, map);
```

Performing the vector-matrix multiply requires a complete copy of `replica` on each processor. An ordinary map divides data among processors, but, here, the goal is to copy the same data to multiple processors. Sourcery VSIP++ provides a special `Replicated_map` class to use in this situation. So, you should declare `replica` as follows:

```
Replicated_map<1> replica_map;
typedef Dense<1,
    value_type, row1_type, Replicated_map<1> >
    replica_block_type;
typedef Vector<value_type, replica_block_type>
    replica_view_type;
replica_view_type replica(nrange, replica_map);
```

Because the application already uses implicitly parallel operators, no further changes are required. The entire algorithm (i.e., the part of the code that performs FFTs and vector-matrix multiplication) remains unchanged.

The complete parallel program is:

```

\
/*****
  Included Files
*****/

#include <vsip/initfin.hpp>
#include <vsip/support.hpp>
#include <vsip/signal.hpp>
#include <vsip/math.hpp>
#include <vsip/map.hpp>

using namespace vsip;

/*****
  Main Program
*****/

int
main(int argc, char** argv)
{
    // Initialize the library.
    vsipl vpp(argc, argv);

    typedef complex<float> value_type;

    typedef Map<Block_dist, Whole_dist>          map_type;
    typedef Dense<2, value_type, row2_type, map_type> block_type;
    typedef Matrix<value_type, block_type>       view_type;

```

```

typedef Dense<1, value_type, row1_type, Replicated_map<1> >
\
replica_block_type;
typedef Vector<value_type, replica_block_type> \
replica_view_type;

// Parameters.
length_type npulse = 64; // number of pulses
length_type nrange = 256; // number of range cells

// Maps.
map_type      map = map_type(num_processors(), 1);
Replicated_map<1> replica_map;

// Views.
replica_view_type replica(nrange, replica_map);
view_type      data(npulse, nrange, map);
view_type      tmp (npulse, nrange, map);

// A forward Fft for computing the frequency-domain version of
// the replica.
typedef Fft<const_Vector, value_type, value_type, fft_fwd, \
by_reference>
for_fft_type;
for_fft_type for_fft (Domain<1>(nrange), 1.0);

// A forward Fftm for converting the time-domain data matrix to the
// frequency domain.
typedef Fftm<value_type, value_type, row, fft_fwd, by_reference>
for_fftm_type;
for_fftm_type for_fftm(Domain<2>(npulse, nrange), 1.0);

// An inverse Fftm for converting the frequency-domain data back to
// the time-domain.
typedef Fftm<value_type, value_type, row, fft_inv, by_reference>
inv_fftm_type;
inv_fftm_type inv_fftm(Domain<2>(npulse, nrange), 1.0/(nrange));

// Initialize data to zero.
data      = value_type();
replica   = value_type();

// Before fast convolution, convert the replica to the the
// frequency domain
for_fft(replica);

// Perform fast convolution:

// Convert to the frequency domain.
for_fftm(data, tmp);

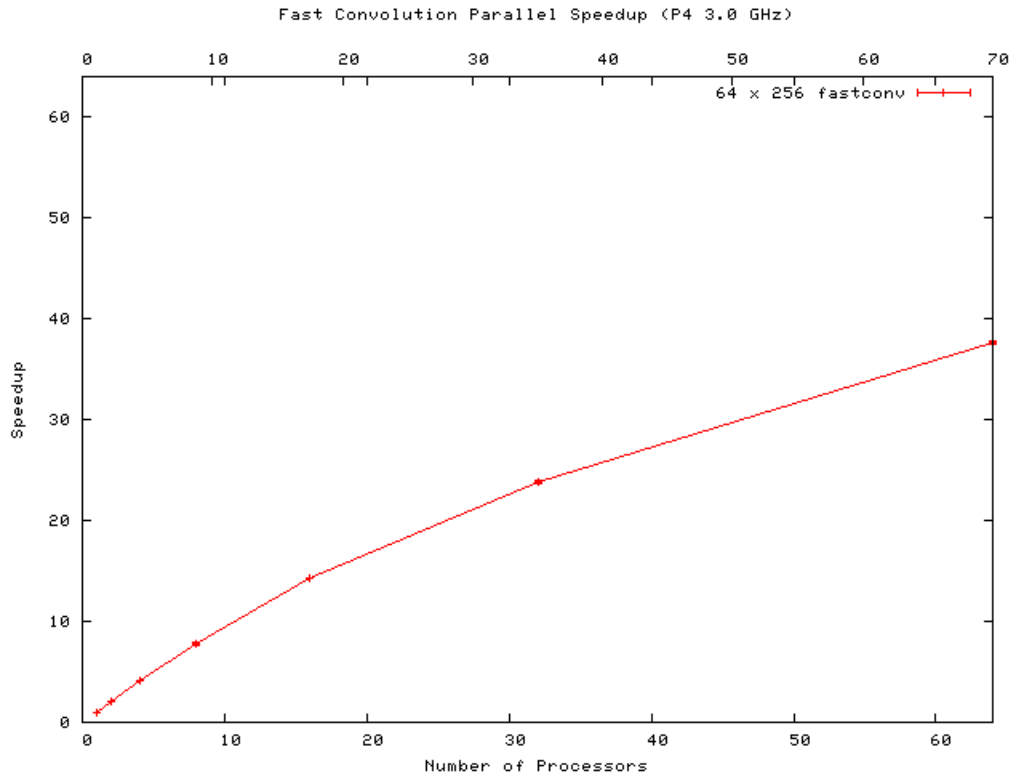
// Perform element-wise multiply for each pulse.
tmp = vmmul<0>(replica, tmp);

```

```
// Convert back to the time domain.
inv_fftm(tmp, data);
}

\
```

The following graph shows the parallel speedup of the fast convolution program from 1 to 32 processors using a 3.0 GHz Pentium cluster system. As you can see, increasing the number of processors also increases the performance of the program.



8.2. Improving Parallel Temporal Locality

In the previous chapter, you improved the performance of the fast convolution program by exploiting temporary cache locality to process data while it was "hot" in the cache. In this section, you will convert that program to run efficiently in parallel.

If we apply `maps` (as in Section 8.1, "Parallel Fast Convolution"), but do not adjust the algorithm in use, the code in Section 7.2, "Serial Optimization: Temporal Locality" will not run faster when deployed on multiple processors. In particular, every processor will want to update `tmp` for every row. Therefore, all processors will perform the forward FFT and vector-multiply for each row of the matrix.

VSIP++ provides *local subviews* to solve this problem. For a given processor and view, the local subview is that portion of the view located on the processor. You can obtain the local subview of any view by invoking its `local` member function:

```
view_type::local_type l_data =
    data.local();)
```

Every view class defines a type (`local_type`) which is the type of a local subview. The `local_type` is the same kind of view as the view containing it, so, in this case, `l_data` is a matrix. There is virtually no overhead in creating a local subview like `l_data`. In particular, no data is copied; instead, `l_data` just refers to the local portion of data. We can now use the same cache-friendly algorithm from Section 7.2, “Serial Optimization: Temporal Locality” on the local subview:

```
rep_view_type::local_type l_replica = replica.local();
for (index_type l_r=0; l_r < l_data.size(0); ++l_r)
{
    for_fft(l_data.row(l_r), tmp);
    tmp *= l_replica;
    inv_fft(tmp, l_data.row(l_r));
}
```

Because each processor now iterates over only the rows of the matrix that are local, there is no longer any duplicated effort. Applying maps, as in Section 8.1, “Parallel Fast Convolution” above, results in the following complete program:

```

\
/*****
Included Files
*****/

#include <vsip/initfin.hpp>
#include <vsip/support.hpp>
#include <vsip/signal.hpp>
#include <vsip/math.hpp>
#include <vsip/map.hpp>

using namespace vsip;

/*****
Main Program
*****/

int
main(int argc, char** argv)
{
    // Initialize the library.
    vsipl vpp(argc, argv);

    typedef complex<float> value_type;

    typedef Map<Block_dist, Whole_dist>          map_type;
    typedef Dense<2, value_type, row2_type, map_type> block_type;
    typedef Matrix<value_type, block_type>       view_type;

    typedef Dense<1, value_type, row1_type, Replicated_map<1> >
replica_block_type;
    typedef Vector<value_type, replica_block_type>

```

```

replica_view_type;

// Parameters.
length_type npulse = 64; // number of pulses
length_type nrange = 256; // number of range cells

// Maps.
map_type      map = map_type(num_processors(), 1);
Replicated_map<1> replica_map;

// Views.
replica_view_type replica(nrange, replica_map);
view_type      data(npulse, nrange, map);
Vector<value_type> tmp(nrange);

// A forward Fft for converting the time-domain data to the
// frequency domain.
typedef Fft<const_Vector, value_type, value_type, fft_fwd, \
by_reference>
for_fft_type;
for_fft_type  for_fft(Domain<1>(nrange), 1.0);

// An inverse Fft for converting the frequency-domain data back to
// the time-domain.
typedef Fft<const_Vector, value_type, value_type, fft_inv, \
by_reference>
inv_fft_type;
inv_fft_type  inv_fft(Domain<1>(nrange), 1.0/nrange);

// Initialize data to zero.
data      = value_type();
replica  = value_type();

// Before fast convolution, convert the replica into the
// frequency domain
for_fft(replica.local());

view_type::local_type      l_data      = data.local();
replica_view_type::local_type l_replica = replica.local();

for (index_type l_r=0; l_r < l_data.size(0); ++l_r)
{
    for_fft(l_data.row(l_r), tmp);
    tmp *= l_replica;
    inv_fft(tmp, l_data.row(l_r));
}
}
\

```

8.2.1. Implicit Parallelism: Parallel Foreach

You may feel that the original formulation using implicitly parallel operators was simpler and more intuitive than the more-efficient variant using explicit loops. Sourcery VSIPL++ provides an extension

to the VSIPL++ API that allows you to retain the elegance of that formulation while still obtaining good temporal locality.

In particular, Sourcing VSIPL++ provides a "parallel foreach" operator. This operator applies an arbitrary user-defined function (or an object that behaves like a function) to each of the rows or columns of a matrix. In this section, you will see how to use this approach.

First, declare a `Fast_convolution` template class. The template parameter `T` is used to indicate the value type of the fast convolution computation (such as `complex<float>`):

```
template <typename T> class Fast_convolution
{
```

This class will perform the forward FFT and inverse FFTs on each row, so you must declare the FFTs:

```
typedef Fft<const_Vector, T, T, fft_fwd,
           by_reference>
for_fft_type;
typedef Fft<const_Vector, T, T,
           fft_inv, by_reference>
inv_fft_type;
Vector<T> replica_;
Vector<T> tmp_;
for_fft_type for_fft_;
inv_fft_type inv_fft_;
```

Next, define a constructor for `Fast_convolution`. The constructor stores a copy of the replica, and also uses the replica to determine the number of elements required for the FFTs and temporary vector.

```
template <typename Block>
Fast_convolution(Vector<T, Block> replica)
: replica_(replica.size()),
  tmp_(replica.size()),
  for_fft_(Domain<1>(replica.size()), 1.0),
  inv_fft_(Domain<1>(replica.size()), 1.0/replica.size())
{
  replica_ = replica;
}
```

The most important part of the `Fast_convolution` class is the `operator()` function. This function performs a fast convolution for a single row of the matrix:

```
template <typename Block1,
         typename Block2,
         dimension_type Dim>
void operator()( Vector<T, Block1> in,
                Vector<T, Block2> out,
                Index<Dim> /*idx*/)
{
  for_fft_(in, tmp_);
  tmp_ *= replica_;
  inv_fft_(tmp_, out);
}
```

The `foreach_vector` template will apply the new class you have just defined to the rows of the matrix:

```
Fast_convolution<value_type>
  fconv(replica.local());
foreach_vector<tuple<0, 1> >(fconv, data);
```

The resulting program contains no explicit loops, but still has good temporal locality. Here is the complete program, using the parallel foreach operator:

```

\
/*****
Included Files
*****/

#include <vsip/initfin.hpp>
#include <vsip/support.hpp>
#include <vsip/signal.hpp>
#include <vsip/math.hpp>
#include <vsip/map.hpp>
#include <vsip/parallel.hpp>

using namespace vsip;

/*****
Main Program
*****/

template <typename T>
class Fast_convolution
{
  typedef Fft<const_Vector, T, T, fft_fwd, by_reference> \
for_fft_type;
  typedef Fft<const_Vector, T, T, fft_inv, by_reference> \
inv_fft_type;

public:
  template <typename Block>
  Fast_convolution(
    Vector<T, Block> replica)
    : replica_(replica.size()),
      tmp_      (replica.size()),
      for_fft_(Domain<1>(replica.size()), 1.0),
      inv_fft_(Domain<1>(replica.size()), 1.0/replica.size())
  {
    replica_ = replica;
  }

  template <typename          Block1,
            typename          Block2,
            dimension_type Dim>

```

```

void operator()(
    Vector<T, Block1> in,
    Vector<T, Block2> out,
    Index<Dim>          /*idx*/)
{
    for_fft_(in, tmp_);
    tmp_ *= replica_;
    inv_fft_(tmp_, out);
}

// Member data.
private:
    Vector<T>    replica_;
    Vector<T>    tmp_;
    for_fft_type for_fft_;
    inv_fft_type inv_fft_;
};

int
main(int argc, char** argv)
{
    // Initialize the library.
    vsipl vpp(argc, argv);

    typedef complex<float> value_type;

    typedef Map<Block_dist, Whole_dist>          map_type;
    typedef Dense<2, value_type, row2_type, map_type> block_type;
    typedef Matrix<value_type, block_type>        view_type;

    typedef Dense<1, value_type, row1_type, Replicated_map<1> >
    replica_block_type;
    typedef Vector<value_type, replica_block_type> \
    replica_view_type;

    // Parameters.
    length_type npulse = 64; // number of pulses
    length_type nrange = 256; // number of range cells

    // Maps.
    map_type          map = map_type(num_processors(), 1);
    Replicated_map<1> replica_map;

    // Views.
    replica_view_type replica(nrange, replica_map);
    view_type          data(npulse, nrange, map);
    view_type          tmp (npulse, nrange, map);

    // A forward Fft for computing the frequency-domain version of
    // the replica.
    typedef Fft<const_Vector, value_type, value_type, fft_fwd, \

```

```

by_reference>
  for_fft_type;
  for_fft_type for_fft (Domain<1>(nrange), 1.0);

  Fast_convolution<value_type> fconv(replica.local());

  // Initialize data to zero.
  data      = value_type();
  replica  = value_type();

  // Before fast convolution, convert the replica into the
  // frequency domain
  for_fft(replica.local());

  // Perform fast convolution.
  foreach_vector<tuple<0, 1> >(fconv, data);
}
\

```

8.3. Performing I/O

The previous sections have ignored the acquisition of actual sensor data by setting the input data to zero. This section shows how to extend the I/O techniques introduced in the previous chapter to initialize data before performing the fast convolution.

Let's assume that all of the input data arrives at a single processor via DMA. This data must be distributed to the other processors to perform the fast convolution. So, the input processor is special, and is not involved in the computation proper.

To describe this situation in Sourcery VSIPL++, you need two maps: one for the input processor (`map_in`), and one for the compute processors (`map`). These two maps will be used to define views that can be used to move the data from the input processor to the compute processors. Let's assume that the input processor is processor zero. Then, create `map_in` as follows, mapping all data to the single input processor:

```

typedef Map<> map_type;
Vector<processor_type> pvec_in(1);
pvec_in(0) = 0;
map_type map_in(pvec_in, 1, 1);

```

In contrast, `map` distributes rows across all of the compute processors:

```

// Distribute computation across all processors:
map_type
  map (num_processors(), 1);

```

Because the data will be arriving via DMA, you must explicitly manage the memory used by Sourcery VSIPL++. Because VSIPL++ uses the SPMD model, each processor must allocate the memory for its local portion the input block, even though all processors except the actual input processor will allocate zero bytes. The code required to set up the views is:

```

block_type data_in_block(npulse, nrange, 0, map_in);
view_type data_in(data_in_block);

```

```
view_type data (npulse, nrange, map);
size_t size = subblock_domain(data_in).size();
std::vector<value_type> buffer(size);
data_in.block()->rebind(&buffer.front());
```

Now, you can perform the actual I/O. The I/O (including any calls to low-level DMA routines) should only be performed on the input processor. The `subblock` function is used to ensure that I/O is only performed on the appropriate processors:

```
if (subblock(data_in) != no_subblock)
{
    data_in.block().release(false);
    // ... perform IO into data_in ...
    data_in.block().admit(true);
}
```

Once the I/O completes, you can move the data from `data_in` to `data` for processing. In the VSIP++ API, ordinary assignment (using the `=` operator) will perform all communication necessary to distribute the data. So, performing the "scatter" operation is just:

```
data = data_in;
```

The complete program is:

```

\
/*****
  Included Files
*****/

#include <vsip/initfin.hpp>
#include <vsip/support.hpp>
#include <vsip/signal.hpp>
#include <vsip/math.hpp>
#include <vsip/map.hpp>
#include <vsip/parallel.hpp>

using namespace vsip;

/*****
  Main Program
*****/

template <typename          ViewT,
          dimension_type Dim>
ViewT
create_view_wstorage(
    Domain<Dim> const&          dom,
    typename ViewT::block_type const& map)
{
    typedef typename ViewT::block_type block_type;
    typedef typename ViewT::value_type value_type;

```

```

    block_type* block = new block_type(dom, (value_type*)0, map);
    ViewT view(*block);
    block->decrement_count();

    if (subblock(view) != no_subblock)
    {
        size_t size = subblock_domain(view).size();
        value_type* buffer = vsip::impl::alloc_align<value_type>(128, \
size);
        block->rebind(buffer);
    }

    block->admit(false);

    return view;
}

template <typename ViewT>
void
cleanup_view_wstorage(ViewT view)
{
    typedef typename ViewT::value_type value_type;
    value_type* ptr;

    view.block().release(false, ptr);
    view.block().rebind((value_type*)0);

    if (ptr) vsip::impl::free_align((void*)ptr);
}

template <typename ViewT>
ViewT
create_view_wstorage(
    length_type          rows,
    length_type          cols,
    typename ViewT::block_type::map_type const& map)
{
    return create_view_wstorage<ViewT>(Domain<2>(rows, cols), map);
}

template <typename ViewT>
ViewT
create_view_wstorage(
    length_type          size,
    typename ViewT::block_type::map_type const& map)
{
    return create_view_wstorage<ViewT>(Domain<1>(size), map);
}

```

```

}

int
main(int argc, char** argv)
{
    // Initialize the library.
    vsipl vpp(argc, argv);

    typedef complex<float> value_type;

    typedef Map<Block_dist, Block_dist>          map_type;
    typedef Dense<2, value_type, row2_type, map_type> block_type;
    typedef Matrix<value_type, block_type>       view_type;

    typedef Dense<1, value_type, row1_type, Replicated_map<1> >
replica_block_type;
    typedef Vector<value_type, replica_block_type> \
replica_view_type;

    typedef Dense<1, value_type, row1_type, Map<> > \
replica_io_block_type;
    typedef Vector<value_type, replica_io_block_type> \
replica_io_view_type;

    // Parameters.
    length_type npulse = 64; // number of pulses
    length_type nrange = 256; // number of range cells

    length_type np = num_processors();

    // Processor sets.
    Vector<processor_type> pvec_in(1); pvec_in(0) = 0;
    Vector<processor_type> pvec_out(1); pvec_out(0) = np-1;

    // Maps.
    map_type          map_in (pvec_in,  1, 1);
    map_type          map_out(pvec_out, 1, 1);
    map_type          map_row(np, 1);
    Replicated_map<1> replica_map;

    // Views.
    view_type data(npulse, nrange, map_row);
    view_type tmp (npulse, nrange, map_row);
    view_type data_in (create_view_wstorage<view_type>(npulse, nrange, \
map_in));
    view_type data_out(create_view_wstorage<view_type>(npulse, nrange, \
map_out));
    replica_view_type  replica(nrange);
    replica_io_view_type replica_in(
        create_view_wstorage<replica_io_view_type>(nrange, map_in));

```

```

// A forward Fft for computing the frequency-domain version of
// the replica.
typedef Fft<const_Vector, value_type, value_type, fft_fwd, \
by_reference>
for_fft_type;
for_fft_type for_fft (Domain<1>(nrange), 1.0);

// A forward Fftm for converting the time-domain data matrix to the
// frequency domain.
typedef Fftm<value_type, value_type, row, fft_fwd, by_reference>
for_fftm_type;
for_fftm_type for_fftm(Domain<2>(npulse, nrange), 1.0);

// An inverse Fftm for converting the frequency-domain data back to
// the time-domain.
typedef Fftm<value_type, value_type, row, fft_inv, by_reference>
inv_fftm_type;
inv_fftm_type inv_fftm(Domain<2>(npulse, nrange), 1.0/(nrange));

// Perform input IO
if (subblock(data_in) != no_subblock)
{
    data_in.block().release(false);
    // ... perform IO ...
    data_in.block().admit(true);

    replica_in.block().release(false);
    // ... perform IO ...
    replica_in.block().admit(true);

    data_in = value_type();
    replica_in = value_type();

    // Before fast convolution, convert the replica into the
    // frequency domain
    for_fft(replica_in.local());
}

// Scatter data
data = data_in;
replica = replica_in;

// Perform fast convolution.
for_fftm(data, tmp); // Convert to the frequency domain.
tmp = vmmul<0>(replica, tmp); // Perform element-wise multiply.
inv_fftm(tmp, data); // Convert back to the time domain.

// Gather data
data_out = data;

// Perform output IO
if (subblock(data_out) != no_subblock)
{

```

```
    data_out.block().release(true);  
    // ... perform IO ...  
    data_out.block().admit(false);  
}  
  
// Cleanup  
cleanup_view_wstorage(data_in);  
cleanup_view_wstorage(data_out);  
cleanup_view_wstorage(replica_in);  
}  
  
\
```

The technique demonstrated in this section extends easily to the situation in which the sensor data is arriving at multiple processors simultaneously. To distribute the I/O across multiple processors, just add them to `map_in`'s processor set `pvec_in`:

```
Vector<processor_type> pvec_in(num_io_proc);  
pvec_in(0) = 0;  
...  
pvec_in(num_io_proc-1) = ...;
```

Appendix A

Command-line options supported by Sourcery VSIPL++

<code>--vsip-profile-mode</code>	Either "accum" or "trace".
<code>--vsip-profile-output</code>	The name of the output file to which to write profile output.
<code>--vsip-num-spes</code>	The number of SPEs to use. (This option is available on Cell binaries only.)
<code>--vsip-cuda-device</code>	The id of the CUDA device to run on. Available devices can be inspected using the <code>vsip-cuda-device-info</code> tool.

Appendix B

Examples

B.1. view

These examples demonstrate the basic use of Sourcery VSIPL++ views and blocks.

tensor_subview.cpp Illustrate how to access tensor subviews.

B.2. ustorage

These examples demonstrate the use of user-defined storage with Sourcery VSIPL++ views and blocks.

dense.cpp Demonstrate the use of user storage with Dense blocks.

B.3. solvers

These examples demonstrate the use of Sourcery VSIPL++ Linear Algebra Solvers.

lu.cpp LU solver example.

qr.cpp QR solver example.

covsol.cpp Covariance solver example.

llsqsol.cpp Least Squares solver example.

B.4. signal

These examples demonstrate the use of Sourcery VSIPL++ Signal processing types.

fft.cpp Simple FFT example.

fft2d.cpp 2-Dimensional FFT example. Separate a two-dimensional dataset into spectral components and find the strongest of two mixed-frequency signals.

iir.cpp IIR filter example. Implement a 2nd order low pass filter and compute the step response.

B.5. dispatch

These examples demonstrate the use of the Dispatcher harness. See Chapter 3, “Using the Dispatch Framework” for details.

ct_custom.cpp This example demonstrates the use of a compile-time dispatch for a custom operation.

rt_custom.cpp This example demonstrates the use of a runtime dispatch for a custom operation.

mprod.cpp This example illustrates the process of creating a custom evaluator for a VSIPL++ standard function -- in this case, a fixed-size 4x4 matrix product that would be amenable to SIMD optimizations.

mprod-dda.cpp This example builds on the ``mprod.cpp`` example, illustrating the process of creating a custom evaluator that references its data directly through data

pointers, as might be necessary for interfacing to external libraries or using SIMD vectors.

fft.cpp This example illustrates the process of creating a custom FFT evaluator. FFT evaluators are more complex than simple function evaluators, because they must create a VSIPL++ Fft functor object rather than simply processing data.

B.6. eval

These examples illustrate how to use Sourcery VSIPL++ expression templates, and how to write custom evaluators for them. See Chapter 4, “Custom Expression Evaluation” for details.

scale.cpp Example scale function using return block optimization.

interpolate.cpp Example interpolate function using return block optimization.

evaluation.cpp Custom expression evaluator for a fused scale & interpolate operation.

fft_expression.cpp Demonstrate how to disambiguate different expressions when specializing evaluators.

B.7. profile

These examples illustrate how to profile Sourcery VSIPL++ applications to fine-tune performance. See Chapter 5, “Profiling” for details.

basic.cpp This example illustrates the basics of using user-defined profiling to trace a program's execution and time operations. Run the resulting program with `--vsip-profile-mode=accum` or `--vsip-profile-mode=trace` to see the profile output.

memory.cpp Trace memory allocation and library-internal block copying.

eval.cpp Profile expression evaluation.

diagnostics.cpp Print diagnostic information about operation dispatching, as well as expression evaluation.

B.8. cvsip

These examples demonstrate the use of Sourcery VSIPL.

fft.c This example performs FFTs using the VSIPL API.

dda.c This example shows how to access VSIPL blocks via a direct data access extension API.

B.9. pi

These examples illustrate the use of Parallel Iterators.

sobel.cpp

stencil.cpp Use of a stencil expression.

<code>sumsqval.cpp</code>	Use of Parallel Iterators to evaluate a row-wise reduction.
<code>iterator_map.cpp</code>	Use of iterator map to reorder a vector.

Glossary

Block	<p>A block is an interface to a logically contiguous array of data. Blocks provide a means to organize the access to the data. They may store the data themselves, or access the data through other blocks. This abstraction provides important latitude for optimizations such as expression templates, or parallelism.</p> <p>Block types have to fulfill the requirements outlined in table 6.1 of the specification.</p>
Dense Block	<p>Dense blocks are modifiable, allocatable blocks that explicitly store one value for each index in its domain. The data layout is specified in terms of a template parameter, allowing storage to be optimized for particular operations (see dimension ordering).</p> <p>Dense blocks allow users to supply data storage, either at construction time, or later, in which case the block is 'rebound' to an alternate user storage.</p>
Dimension Ordering	<p>Dimension ordering refers to the layout of data in a multi-dimensional block, such as row-major or column-major. Dimension ordering has an impact on performance in operations involving loops over the data, as adjacent reads / writes may require a new cache-line to be fetched first.</p>
Domain	<p>A domain represents a logical set of indices for which views provide data. It may be a contiguous set of indices for dense matrices, or a non-contiguous set of indices for subviews.</p>
Expression Block	<p>Expression blocks are used to store mathematical expressions, allowing optimized evaluation. Conventionally, in an equation 'View A = B + C * D' the computation of A would require at least two temporaries, representing the results of the two binary operations. Additionally, the evaluation of each of these subexpressions implies a loop, resulting in suboptimal performance.</p>

With expression blocks, the above expression will generate a block representing $B + C * D$, which is evaluated when assigned to 'A'. Specializations of expression blocks may use highly optimized functions to be called, depending on the specific types and subexpressions involved.

Map

A map specifies how a block can be divided into subblocks for the purpose of parallel execution. It defines how subblocks are to be assigned to processors.

Map types have to fulfill the requirements outlined in table 3.1 of the parallel specification.

View

A view represents the base for mathematical linear algebra operations, such as vectors, matrices, tensors. It has a dimension, a `value_type`, and a number of accessors to access and manipulate its values. The actual data are stored in blocks, to which views hold references internally.

Multiple views may share the same data, making copy operations for those views an inexpensive operation. All views are parameterized for two types: the view's `value_type`, as well as the underlying block type.

View types have to fulfill the requirements outlined in table 6.3 of the specification.