

Dispel Language Tutorial

Paul William Martin

July 13, 2012

Contents

1	Getting started	2
1.1	Core concepts	2
1.2	The gateway	3
2	Building simple workflows	5
2.1	We'll skip 'Hello World'	5
2.2	Describing workflow input	7
2.3	Defining new types of processing element	8
2.4	Deriving new workflow elements	10
3	Constructing more sophisticated workflows	13
3.1	Dispel Language Types	13
3.2	Iterative workflow construction	16
3.3	Conditional workflow construction	19
4	Manipulating the flow of data	22
4.1	Dispel structural types	22
4.2	Type coercions	26
4.3	Connection modifiers	28
4.4	Dispel domain types	29
5	Case studies	31
5.1	The Sieve of Eratosthenes	31

Chapter 1

Getting started

Dispel is a strongly-typed imperative language for generating executable workflows for data-intensive distributed applications, particularly (but not exclusively) for use in computational sciences such as bioinformatics, astronomy and seismology — it has been designed to be a portable *lingua franca* by which researchers can interact with complex distributed research infrastructures *without* detailed knowledge of the underlying computational middleware, all in order to more easily conduct experiments in data integration, simulation and data-intensive modelling.

Dispel was created as part of the ADMIRE project¹, which sought to promote a model for advanced data mining and integration which insulates the computational scientist or domain expert from the specifics of how individual computational services are implemented or how data is moved between physical resources.

1.1 Core concepts

A *workflow* is simply a decomposition of a task into a number of sub-procedures which, when linked together, describe how that task can be carried out. Such workflows can be understood in terms of *data-flow*, wherein the output of one sub-procedure feeds directly into the next sub-procedure. If a task involves a continuous processing of data over some period of time, then each sub-procedure can begin enactment as soon as data produced by prerequisite sub-procedures starts to emerge. If these sub-procedures can be distributed amongst a number of different actors and a means to efficiently deliver data between actors can be provided, then the workflow can be effectively parallelised, and results can start being produced almost immediately regardless of the ultimate size of the base task.

A Dispel workflow is a decomposition of an application into a number of processing elements connected together via channels through which data can be streamed. Each processing element (PE) encapsulates some sub-procedure deemed pertinent to the task at hand and takes some number of input connections as well as some number of output connections. As a PE consumes data through its inputs, some operation is performed within the element which results in a number of new streams of data

¹EC Framework 7 ICT 215024 (<http://www.admire-project.eu>).

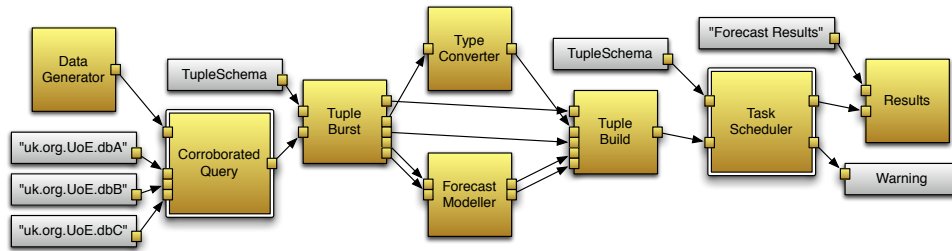


Figure 1.1: An example of a Dispel workflow.

produced through its outputs. The rate of data consumption and production depends on the behaviour of the PE and its neighbouring elements.

A Dispel script describes how to construct a workflow; scripts must be executed by an interpreter capable of mapping the resulting workflow onto a suitable *enactment platform*, being a collection of services and middleware distributed onto physical resources which can actually handle the execution of workflow sub-procedures. Such an interpreter is usually provided by a remote *gateway* which can execute any Dispel script submitted through it on an enactment platform provided by the (distributed) system to which it is attached. The gateway serves to conceal the vagaries of implementation and physical topology of resources from the casual user, instead presenting a library of PEs implementable by the enactment platform which can be enlisted in the construction of workflows. In order to provide a common library to Dispel users, the gateway usually defers to a central *registry* (though other arrangements are possible), whilst implementation code for registered PEs on a given enactment platform is kept in a nearby *repository*. Provided that a valid workflow has been submitted, a gateway will automatically map workflow components to services and processes implemented by the enactment platform, delegate their execution to suitable pre-configured computational resources and then choreograph the execution of the distributed ‘complete’ workflow.

1.2 The gateway

In order to use Dispel, there must be a service able to interpret Dispel scripts and map them onto an actual enactment platform. Whilst such a service can be configuration on a user’s own machine for simple tasks, the scenarios envisaged for Dispel generally involve submission of scripts to a remote gateway which acts as an interface onto some kind of distributed system — such a distributed system could be a federation of computational resources including large data archives and high-performance compute clusters.

An ADMIRE gateway maps Dispel requests onto OGSA-DAI² workflows — these workflows describe a concrete distributed implementation of the logical workflows described by a Dispel script by replacing individual PEs with compositions of OGSA-DAI activities. An ADMIRE gateway can be stand-alone, or part of a federation of gateways; a gateway can be configured to refer to an independently-hosted registry and repository, or have its own in-memory private registry for stand-alone use.

²<http://www.ogsadai.org.uk>

To execute the Dispel examples given in this tutorial, you must either use an existing gateway, or install a new one. The installation and deployment of a gateway is described at <http://sourceforge.net/apps/trac/admire/wiki/GatewayInstallation>; the ADMIRE gateway is a Java Servlet (JRE version 1.6+), typically hosted using Apache Tomcat (version 6+). It uses a RESTful HTTP interface for the submission of Dispel and the retrieval of workflow state.

Chapter 2

Building simple workflows

A basic *Dispel* script consists of a series of statements handling the importing of useful processing elements from the local gateway’s registry, the instantiation and configuration of those elements and either the registering of new composite processing elements or the submission of a workflow to be executed by the gateway (or both). This section shall introduce some of the core functionality of *Dispel*, demonstrating how to construct simple workflows for submission.

2.1 We’ll skip ‘Hello World’

Processing elements (PEs) describe the principal components which make up any workflow. An active gateway will provide a number of fundamental PEs which are commonly used in various data-intensive applications; more may be added by data analysis experts and software engineers. Such experts can also use *Dispel* itself to define *composite* PEs, building an increasingly sophisticated array of pre-fabricated components for users to exploit in their workflows.

Available PEs are described within the registry associated with the gateway to which a script is submitted. We can import a PE description from the local registry by invoking the `use` directive:

```
use dispel.db.SQLQuery;
```

The above command extracts a PE named `dispel.db.SQLQuery` from the registry; it also imports the identifier `SQLQuery` into the local namespace, which means that we can drop the prefix `dispel.db` when referring to this PE within this *Dispel* script. An `SQLQuery` converts queries written in SQL into responses returned by a selected database. Every useful PE has some combination of input and output connections — a few may act as sources (only outputs) and a few may act as sinks (only inputs), but the majority will have at least one input and one output. `SQLQuery` has two inputs and one output; `expression` (into which queries may be fed), `resource` (into which the location of the database to be queried must be fed in tandem with every query) and `data` (from which the results of any queries will emerge).

```

1 // Import PEs from the local registry.
2 use dispel.db.SQLiteQuery;
3 use dispel.core.Results;
4
5 // Create instances of PEs.
6 SQLiteQuery query = new SQLiteQuery;
7 Results results = new Results;
8
9 // Construct workflow and feed in data.
10 |-"SELECT * FROM littleblackbook WHERE id <= 10"-| => query.expression;
11 |-"uk.org.UoE.dbA"-| => query.resource;
12 |-"10 entries from the little black book"-| => results.name;
13 query.data => results.input;
14
15 // Submit the entire workflow.
16 submit results;

```

Figure 2.1: A simple Dispel script for submitting an SQL query.

In order to make use of any imported PE however, we must first create an instance of the PE. A new instance can be created by use of the `new` directive:

```
SQLiteQuery query = new SQLiteQuery;
```

Here we are defining a *PE instance* named `query`, which is of course an instance of the `SQLiteQuery` PE. To do anything useful, we need to connect each input of `query` to a suitable data stream. Such data streams typically come from the output of other PEs, but we can also feed in data directly from a script:

```

|-"SELECT * FROM littleblackbook WHERE id < 10"-| => query.expression;
|-"uk.org.UoE.dbA"-| => query.resource;

```

In this case, both inputs (`query.expression` and `query.resource`) read in text strings, one describing an SQL query and the other identifying a database (in this case via the identifier by which it is known within the distributed system associated with the gateway).

Theoretically, this is enough to describe a ‘useful’ workflow. We need only invoke the `submit` command, and the gateway will be able to provide an implementation of `SQLiteQuery` with which to query the referenced database with the specific query given:

```
submit query;
```

Of course, since we didn’t connect the output of `query` to anything in particular, we can only actually find out the result of our query if we are able to somehow inspect `query` directly as it executes, which we can assume to be unlikely and undesirable. Usually, what we want to do instead is channel the output of our workflows towards a sink PE which can report back its input, whether (for example) by saving it to a file in a known location or by directly visualising it using a portal widget or even in a terminal display. So instead of simply submitting the workflow as is, we direct the output of `query` towards an instance of the `Results` PE and submit that instead.

PE `dispel.core.Results` has two inputs — `input` which takes data to be recorded, and `name` which associates a name with the data recorded for ease of reference. In Figure 2.1, `query.data` is connected to `results.input`, storing the result of the given query somewhere which can be directly accessed after the workflow is enacted. The entire workflow is submitted by submitting `results` — submitting any PE instance connected to the rest of a workflow will serve to submit that workflow, but convention favours the ‘final’ process element.

2.2 Describing workflow input

In the previous section, we fed data into an instance of the `SQLQuery` PE directly, but only fed in a single set of inputs. We did this using what is referred to as a *stream literal*. A stream literal describes the content of a data stream explicitly as a sequence of data elements drawn from the Dispel script itself, rather than as the output of a PE instance:

```
|input1, input2, ...|
```

This sequence of elements can consist of any number of arbitrarily complex data structures built from a number of elementary data types (booleans, integers, character strings, *etc.*) as dictated by PE requirements and the whims of the workflow designer. For now however, we shall concentrate on simple inputs of text strings.

It is trivial to describe a list of homogeneous inputs as a stream literal — say we want to define not one but three queries to stream into an instance of the `SQLQuery` PE. We might write the following:

```
String query1 = "SELECT name FROM littleblackbook WHERE id <= 10";
String query2 = "SELECT name FROM littleblackbook"
    + "WHERE id > 10 AND id <= 20";
String query3 = "SELECT address FROM littleblackbook"
    + "WHERE name = 'David Hume'";

|-query1, query2, query3-| => query.expression;
```

This is not enough however. An instance of `SQLQuery` takes not one but two inputs, and consumes data from each input at the same rate. Thus, for each query made to a PE instance like `query`, a reference to a database to query must be given. If we wish to direct each query to a different database, then we could legitimately write the following:

```
String database1 = "uk.org.UoE.dbA";
String database2 = "uk.org.UoE.dbB";
String database3 = "uk.org.UoE.dbC";

|-database1, database2, database3-| => query.resource;
```

If we want to direct all queries to the *same* database however, then it seems rather inefficient to repeat the same string multiple times manually. Fortunately, we can use a `repeat` expression to write this more elegantly:


```

1 // Import PEs from the local registry.
2 use dispel.db.SQLiteQuery;
3 use dispel.tutorial.PrecociousChild;
4 use dispel.core.Results;
5
6 // Create instances of PEs.
7 PrecociousChild child = new PrecociousChild;
8 SQLiteQuery query = new SQLiteQuery;
9 Results results = new Results;
10
11 // Construct workflow and feed in data.
12 child.output => query.expression;
13 |-repeat enough of "uk.org.UoE.dbA"-| => query.response;
14 |-"Adult responses"-| => results.name;
15 query.data => results.input;
16
17 // Submit the entire workflow.
18 submit results;

```

Figure 2.2: PE instance `query` takes endless input from PE instance `child`, whilst its other input is locked to a single value.

```
|-repeat 3 of "uk.org.UoE.dbA"-| => query.resource;
```

In fact, if we know that we shall always be feeding the same input to a particular PE instance, then we can use the keyword `enough` to essentially lock the input to a given value no matter how much data flows through the PE instance's other inputs:

```
|-repeat enough of "uk.org.UoE.dbA"-| => query.resource;
```

This is especially useful if other inputs are connected to the output of 'black box' PEs which produce an unpredictable volume of data for processing — see Figure 2.2 for example.

It is worth noting that in Figure 2.2, the input to `results.name` is *not* repeated — this is because unlike instances of `SQLiteQuery`, instances of `Results` only read data through their `name` input once. If we want `query` to behave in the same way, we will need to adapt `SQLiteQuery` to better suit our purposes.

2.3 Defining new types of processing element

It might be felt that having to remember to ensure a continuous supply of identical inputs to a PE is undesirable, and that it should be possible to configure a PE which, for example, will always refer to the same database for a given stream of queries. It is possible to do such a thing in `Dispel` by adapting the functionality of existing PEs. These customised PEs can then be inserted into a workflow immediately, or registered in order to be used (and reused) later.

There are two ways to adapt a PE to serve our purpose. One is to modify the

behaviour of the connection interfaces associated with a PE, which we shall defer to later in this tutorial; the other is to wrap one or more PEs within a new PE. In order to produce a new PE, it is necessary to define its abstract behaviour. Since we consider PEs in `Dispel` essentially as black boxes taking in a set of inputs and producing a set of outputs, it is only to be expected that we classify PEs by their connections. Formally, a *connection* is a link between two *connection interfaces* — we have made several such connections already:

```
child.output => query.expression;
|-repeat enough of "uk.org.UoE.dbA"-| => query.resource;
query.data => results.input;
```

The connection operator `=>` connects the left interface (for example `child.output`) to the interface on the right (in this case `query.expression`). A stream literal is therefore an *ad-hoc* connection interface. Note that stream literals can only be found on the left-hand side of connections — we cannot feed the output of a PE into a stream literal, nor can we connect two stream literals together.

An *external* connection connects an output of one PE (or stream literal) to the input of another. A given output can be connected to many different inputs, though each input can only receive data from one source — if we wish to merge multiple outputs into a single input, then we must use a suitable PE which knows how to interpolate the outputs correctly, like `dispel.core.Combiner`. Meanwhile, an *internal* connection connects all inputs of a PE to all outputs of the same PE — in other words, we can describe a PE abstractly by the internal connection made between its interfaces with outside data streams. Thus, when we define a PE, we define it by describing such an internal connection in the following format:

```
PE( <Connection input1; Connection input2; ...> =>
    <Connection output1; Connection output2; ...> );
```

Consequently, we can define the PE `SQLQuery` as so:

```
PE( <Connection expression; Connection resource> =>
    <Connection data> );
```

Whilst for a source PE like `PrecociousChild` (used earlier in Figure 2.2):

```
PE( <> => <Connection output> );
```

Such definitions can quickly get cumbersome however, which is why we use aliases like `SQLQuery` and `PrecociousChild`. It is possible to define new aliases by using a `Type` declaration. Say that we wish to define a new PE which, like `SQLQuery`, reads queries from a data stream and produces a list of responses, but unlike `SQLQuery`, always refers to a specific database. We can define the abstract PE type as so:

```
Type SQLToTupleList is
    PE( <Connection expression> => <Connection data> );
```

From now on, whenever we want to create a new PE which takes in an expression and produces data, we can refer to it as a sub-type of `SQLToTupleList`. An `SQLToTupleList` is just an abstract PE however — to have a PE which we can instantiate and use, we

need PEs with internal architecture which can actually be implemented using concrete computational components. This is where the notion of a PE constructor comes in.

2.4 Deriving new workflow elements

A PE constructor is a particular type of function which returns implementable descriptions of abstract PEs. Consider the constructor `lockSQLDataSource`:

```
PE<SQLToTupleList> lockSQLDataSource(String dataSource) {
    SQLQuery query = new SQLQuery;
    |-repeat enough of dataSource-| => query.resource;
    return PE( <Connection expression = query.expression> =>
              <Connection data      = query.data> );
}
```

Like any other function which might be found in a functional or imperative programming language, it has a function head and a function body. The function head consists of a return type (in this case `PE<SQLToTupleList>`), a function name (`lockSQLDataSource`) and an ordered set of parameters (here, only one — `String dataSource`) The function body is a set of statements ending with a `return` directive. A function which returns a description of a PE must return a PE internal connection which matches that of the abstract PE given in the function head (in this case, an input `expression` and an output `data`).

In the case of `lockSQLDataSource`, the function describes how to make implementable `SQLToTupleList` by taking a new `SQLQuery` instance named `query`, and then locking the `resource` connection interface within that PE instance to the string value passed to it. It would then simply return a version of `SQLToTupleList` wherein the `expression` and `data` interfaces are attached to those of `query`. Using this function, successive implementable versions `SQLToTupleList` can be defined by invoking `lockSQLDataSource` with a suitable instance of parameter `dataSource`:

```
PE<SQLToTupleList> TutorialQuery = lockSQLDataSource("uk.org.UoE.dbA");
```

Having created the new implementable PE `TutorialQuery`, it is now possible to create an instance of that PE and use it in a workflow:

```
TutorialQuery query = new TutorialQuery;
```

Thus a PE constructor describes how existing PEs can be used to produce a desired composite PE which implements a given abstract PE. In this case, `query` is a PE instance, an instance of `TutorialQuery`. Meanwhile `TutorialQuery` is a (rather trivial) composite PE, of type `PE<SQLToTupleList>` — in other words, `TutorialQuery` is an implementable version of the abstract PE `SQLToTupleList`. Thus `SQLToTupleList` exists principally to describe the kind of internal connection exhibited by PEs like `TutorialQuery`; the role of abstract PEs is to provide vessels into which to insert compositions of other PEs like `SQLQuery`. Composite PEs will be deconstructed on execution by the ADMIRE gateway into their constituent PEs, which will in turn be repeatedly deconstructed until only ‘primitive’ PEs remain for which there exists concrete implementations. Instances of these implementations may then be distributed

and executed on many different resources, but to the user, there is only one top-level PE instance to concern themselves with.

Having defined a new type of abstract PE, a constructor for making implementable versions of that abstract PE, and then an example of such an implementable PE, it would be helpful if those entities could be preserved for future workflows. That is the role of the `register` command:

```
register SQLToTupleList, lockSQLDataSource, TutorialQuery;
```

The command `register` directs the ADMIRE gateway to record the given entities within its local repository, registering their existence within that repository in its local registry. It is possible to register not only derivative PEs, but as demonstrated, abstract PEs and PE constructors (though not PE instances, which exist only for the lifetime of a given workflow).

However it is not generally a good idea to register new entities without placing them into some kind of intelligent hierarchy, so as to allow them to be easily located by other users, and to avoid overwriting other entities which happen to share the same name. This is done by placing entities into *packages*, which can themselves hold further sub-packages. We have already encountered packages before when importing PEs such as `SQLQuery`:

```
use dispel.db.SQLQuery;
```

The above directive states that PE `SQLQuery` resides in package `db`, which itself resides in `dispel` (as it happens, `dispel.db` is the main database processing package of `Dispel`).

When registering new entities, we want to put them within a particular package along with similar entities. To do this, we wrap entire `Dispel` scripts within `package` directives:

```
package tutorial.example {  
  ...  
}
```

Any invocations of `register` within a given `package` environment will automatically be registered within the package named (in the above example, `tutorial.example`). Thus, a full `Dispel` script introducing new workflow components to the ADMIRE framework will take on an appearance not unlike that of Figure 2.3. If we then wanted to make use of `TutorialQuery` in another script, we need only declare the following:

```
use tutorial.example.TutorialQuery;
```

In this manner can users simplify the construction of complex workflows, either by simplifying the use of certain standard PEs (by wrapping common configurations of complex PEs into simpler PEs with fewer inputs) or by encapsulating complex recurring tasks within a simple interface (by wrapping whole workflows into a single PE definition).

```

1 package tutorial.example {
2   // Import existing PE from the registry.
3   use dispel.db.SQLiteQuery;
4
5   // Define new PE type.
6   Type SQLiteToList is PE( <Connection expression> =>
7                           <Connection data> );
8
9   // Define new PE constructor.
10  PE<SQLiteToList> lockSQLiteDataSource(String dataSource) {
11    SQLiteQuery query = new SQLiteQuery;
12    |-repeat enough of dataSource-| => query.source;
13    return PE( <Connection expression = query.expression> =>
14              <Connection data = query.data> );
15  }
16
17  // Create new PEs.
18  PE<SQLiteToList> TutorialQuery = lockSQLiteDataSource("uk.org.UoE.dbA");
19  PE<SQLiteToList> MirrorQuery   = lockSQLiteDataSource("uk.org.UoE.dbB");
20
21  // Register new entities.
22  register TutorialQuery, MirrorQuery;
23 }

```

Figure 2.3: New PEs TutorialQuery and MirrorQuery are created by locking instances of SQLiteQuery to a given database.

```

1 // Import PEs from the local registry.
2 use tutorial.example.TutorialQuery;
3 use dispel.core.Results;
4
5 // Create instances of PEs (no import necessary).
6 TutorialQuery query = new TutorialQuery;
7 Results results = new Results;
8
9 // Connect PEI together to create workflow (no data source needed).
10 |-SELECT * FROM littleblackbook WHERE id <= 10-| => query.expression;
11 |-10 entries from the little black book-| => results.name;
12 query.data => results.input;
13
14 // Submit workflow.
15 submit results;

```

Figure 2.4: New PE TutorialQuery is used within a simple workflow.

Chapter 3

Constructing more sophisticated workflows

So far we have only considered very simple workflows and very simple composite PEs which are merely wrappers for only slightly more complex existing PEs. If we want to construct more complex workflows however, we need to be able to scale component composition to arbitrary degrees, and be able to exercise more control over the selection of components based on circumstances at execution time. In order to do that, we need to be able to refer to certain variable factors at execution time. The declaration and assignment of values to variables falls into the domain of the first of Dispel's three type systems — the *language* type system. Using language types, we can direct the iterative construction of workflows and impose conditions on certain elements, as well as configure functions such as the constructor functions used to build custom PEs.

3.1 Dispel Language Types

The Dispel language type system validates the variables, constants and functions used within Dispel scripts. A variable is simply a vessel for some value. Every variable has a language type, and its existence must be declared before use:

```
Integer number;
```

In this case, a variable `number` of language type `Integer` is declared. A variable name must begin with an alphabetic character and contain only alphanumeric characters or underscores (`_`). By convention, variable names use camel-case. Variables must be assigned an initial value before they can be used; afterwards, variables can be assigned new values as often as desired.

```
Integer numberOfSources;  
  
numberOfSources = 4;  
numberOfSources = -1;
```

Typically, variables are assigned an initial value upon declaration, as so:

```
Integer number = 4;
```

Variables can only be assigned to literals or expressions of the correct language type. Variables of language type `Integer` can only be assigned to integer values or expressions which evaluate to integer values — if another variable is provided, then the former variable is assigned the value of the latter variable at the point of evaluation:

```
Integer number = 0;

number = 3 + (-4);
number = factorial(7);

Integer square = number * number;
```

Dispel recognises five basic language types; `Boolean`, `Integer`, `Real`, `String` and `Stream`. Each of these language types has its own valid literal type:

`Boolean` variables can only be assigned to one of two values; `true` or `false`:

```
Boolean statement1 = false, statement2 = true;
```

`Integer` variables can be assigned any positive or negative integer value, as described above.

`Real` variables can be assigned any decimal value:

```
Real pi = 3.14, negative = -43.265;
```

`String` variables can be assigned to any character string; character strings must be enclosed within double quotes:

```
String text = "";
...
text = "Hello World!"
```

Special characters (such as tabs, carriage returns and double quote itself) are represented within strings by special escape characters preceded by a backslash (`\`). Longer strings can be split into segments, appended using the `+` operator:

```
text = "This is the first line...\n" +
      "This is the second line with text in \"quotes\".";
```

`Stream` variables can only be assigned to stream literals, as described in §2.2. Stream literals are enclosed in stream delimiters (`|-` and `-|`) and can contain either a comma-separated list of values (which can be literals, expressions or variables of any of the above language types *except* `Stream` itself) or a stream expression such as `repeat`:

```

Integer three      = 3;
Stream empty      = |--|, list = |-1, "2", three-|;
Stream repeating  = |-repeat 11 of "Eleven"-|;

```

As with strings, stream literal fragments can be appended together using the `+` operator:

```

Stream fragment = |-2-|;
Stream concat   = |-1-| + fragment + |-3-|;

```

Streams represent ordered sequences, with the left-most elements preceding elements to the right; concatenations of streams add to the end of the resulting stream.

Dispel also recognises an additional language type `Connection`, representing a connection interface. `Connection` is a ‘null’ type however, its variables never holding any value — instead connection interface ‘variables’ are simply handles for establishing connections between interfaces as well as streams and interface. As such, variables of type connection are simply declared and need never be assigned values, except to other connection interfaces:

```

Connection input;
Stream data = |-1, 2, 3, 4-|;
data => input;
Connection alias = input;

```

Arrays of variables can be created by first defining the type of the array’s constituent elements and the size of the array, and then assigning values to each individual element of the array in turn as if it was a new variable of the relevant type:

```

Boolean[] array = new Boolean[3];
array[0] = true;
array[1] = false;
array[2] = true;

```

Arbitrarily multi-dimensional arrays can be created by creating arrays of arrays.

```

Integer[][] matrix = new Integer[3][2];
matrix[0][0] = 0;
matrix[0][1] = 1;
matrix[1][0] = 34;
matrix[1][1] = -3;
matrix[2][0] = -245;
matrix[2][1] = 1111;

```

The length of an array `array` can be retrieved by referencing the `length` property of `array`:

```

Integer size = array.length;           // = 3
Integer outerSize = matrix.length;     // = 3
Integer innerSize = matrix[0].length;  // = 2

```


The length of an array is always an integer.

PEs are considered to be language types as well. As such, PE instances must be declared with a type (such as `SQLQuery`) and assigned a value — this will always be a new instantiation of the given PE type, made using the `new` directive:

```
SQLQuery query = new SQLQuery;
```

As already described in §2.3, new types of PEs can be created using a `Type` declaration. There are in fact two ways to create new PE types. The first is as shown in §2.3, using an internal connection signature:

```
Type SQLToTupleList is PE( <Connection expression> =>  
    <Connection data> );
```

This must then be followed by an invocation of a suitable PE constructor in order to provide an implementation of the abstract type, as shown in §2.4:

```
PE<SQLToTupleList> TutorialQuery = lockSQLDataSource("uk.org.UoE.dbA");
```

Another approach is to modify an existing PE type:

```
Type LockedSQLQuery is SQLQuery with initiator resource;
```

In this case, `SQLQuery` is modified such that its input interface `resource` is an `initiator`. Connection modifiers like `initiator` are used to specify restrictions on the use of connection interfaces, and can be used to refine how a given workflow element is implemented upon workflow submission — for now however, we shall defer further details until §4.3.

3.2 Iterative workflow construction

With the ability to define variables comes the ability to define iterators. Assume that we want to create an array of parallel `SQLQuery` PE instances where each instance queries a different data source, but otherwise all instances perform the same query operation, sending the results further along the workflow. To start with, we need to initialise the array of `SQLQuery` instances and provide a connection interface from which to extract query expressions.

```
Connection input;  
SQLQuery[] queries = new SQLQuery[numberOfSources];
```

We would also need an array of connection interfaces from which to acquire data source information and another array of connection interfaces to which to send the result (remember that we cannot simply connect multiple outputs to a single input, but will need to use a suitable PE to combine the outputs later):

```
Connection[] sources = new Connection[numberOfSources];  
Connection[] outputs = new Connection[numberOfSources];
```

The problem then lies with how to properly instantiate the constituent elements of array `queries` and connect each PE instance to the correct inputs and output for arbitrary values of `numberOfSources`. We need an iterator.

Iterators are used to repeatedly execute a block of statements whilst a given condition holds, and as such can be used to succinctly describe repetitive workflow patterns. Dispel supports two standard iteration constructs; `while` and `for`.

The `while` construct is the simplest type of iterator. At each cycle of the iterator, a condition is evaluated, and the statement block within the loop is then executed only if that condition evaluates as `true`; otherwise, execution proceeds beyond the loop. For example:

```
Integer i = 0;
while (i < numberOfSources) {
    queries[i] = new SQLQuery;
        input => queries[i].expression;
        sources[i] => queries[i].source;
    queries[i].data => outputs[i];
    i++;
}
```

Naturally if the loop is to terminate, the body of the iterator must do something which will eventually cause the evaluation of the condition to fail — in the above example, the statement `i++` increments variable `i`, ensuring that eventually, `i` will equal `numberOfSources`.

Since many iterators rely on a single control variable which is updated regularly during each cycle of the loop however, there exists a variant of the `while` construct known as a `for` loop. Each `for` loop consists of an initialisation part (where the control variable is initialised), a conditional part (which determines when the loop should terminate), and an update part (where the control variable is updated). For example:

```
for (Integer i = 0; i < numberOfSources; i++) {
    queries[i] = new SQLQuery;
        input => queries[i].expression;
        sources[i] => queries[i].source;
    queries[i].data => outputs[i];
}
```

In the above example, a new instance of `SQLQuery` is created within array `queries` and connected to surrounding interfaces a number of times equal to `numberOfSources`. First control variable `i` is initialised, which is incremented at the end of every iteration (as directed by the statement `i++`) as long as the condition `i < numberOfSources` holds.

An iterator is used in constructor `makeCorroboratedSQLQuery` in Figure 3.1 to implement abstract PE `MulticastQuery`. This particular implementation of `MulticastQuery` queries multiple data sources at once, and returns the intersection of results — in other words, it only returns results which can be corroborated by multiple sources. It does this using the `for` loop described above and an instance of PE `ListIntersect` to combine the outputs of the constituent `SQLQuery` instances.

Note the instantiation of `ListIntersect` PE instance `intersect` in Lines 18–19 specifies the required size of the array of inputs which it should merge:

```

1 package tutorial.example {
2   use dispel.db.SQLQuery;
3   use dispel.core.ListIntersect;
4
5   // A PE type which queries multiple sources.
6   Type MulticastQuery is
7     PE( <Connection expression; Connection[] sources> =>
8         <Connection data> );
9
10  // Use parallel SQLQuery instances to corroborate results.
11  PE<MulticastQuery>
12  makeCorroboratedSQLQuery(Integer numberOfSources) {
13    // Define aliases for workflow inputs in advance.
14    Connection expr;
15    Connection[] srcs = new Connection[numberOfSources];
16    // Create instances of internal PEs.
17    SQLQuery[] queries = new SQLQuery[sources];
18    ListIntersect intersect = new ListIntersect
19      with inputs.length = numberOfSources;
20
21    // Connect SQLQuery instances in parallel.
22    for (Integer = 0; i < sources; i++) {
23      queries[i] = new SQLQuery;
24      expr => queries[i].expression;
25      srcs[i] => queries[i].source;
26      queries[i].data => intersect.inputs[i]
27    }
28
29    // Return intersection of query results.
30    return PE( <Connection expression = expr;
31              Connection sources = srcs> =>
32              <Connection data = intersect.output> );
33  }
34
35  register MulticastQuery, makeCorroboratedSQLQuery;
36 }

```

Figure 3.1: A Dispel script which describes how to construct a composite PE which queries multiple data sources simultaneously and corroborates the results.

```

ListIntersect intersect = new ListIntersect
  with inputs.length = numberOfSources;

```

The use of `with` in this fashion can be used to specify configurable aspects of a PE instance's state and operation. We shall see other uses of `with` later in this tutorial.

3.3 Conditional workflow construction

In the previous section, we provided a constructor for abstract PE `MulticastQuery` which returned the intersection of results. Another reasonable implementation of `MulticastQuery` would be one which simply returned all results provided by all data sources. In question then is whether or not to remove duplicate results: if we simply want to see what results exist, we may prefer to remove duplicates; if we want to analyse how often certain results arise, we may prefer not to.

Let us assume that we have an iterator which provides an array of `SQLQuery` PE instances as described in the previous section. Assume also that we connect the outputs of these PE instances to an instance `merge` of PE `ListMerge`, a PE which reads a list from each of its inputs and combines them into one output list, happily duplicating results if they appear on more than one input stream. What we want then is to either return the output of `merge` unchanged, or pass that output through an instance `prune` of `DuplicatePrune`, a PE which will remove any duplicate entries in any list it consumes, first. We can decide this at execution-time based on a boolean variable:

```
Boolean removeDuplicatess;
```

If `removeDuplicatess` evaluates as `true`, then we want `prune` to be inserted into our workflow; otherwise we do not.

When statement blocks must be executed dependent upon certain conditions, we use the `if/else` or `switch/case` constructs. A basic `if` conditional executes a statement block only if a given expression evaluates as `true`:

```
if (removeDuplicatess) {
    DuplicatePrune prune = new DuplicatePrune;
    merge.output => prune.input;
    Connection output = prune.output;
}
```

Alternatively, an `if/else` conditional will execute one statement block if the condition evaluates as `true`, and another if it evaluates as `false`:

```
if (removeDuplicatess) {
    DuplicatePrune prune = new DuplicatePrune;
    merge.output => prune.input;
    Connection output = prune.output;
} else {
    Connection output = merge.output;
}
```

An `if/else` conditional is used in constructor `makeMassSQLQuery` in Figure 3.2. This constructor, which is also based on `MulticastQuery` like Figure 3.1, works similarly to `makeCorroboratedSQLQuery`, except that it returns all results provided by all data sources. In addition, this function takes `removeDuplicatess` as a parameter, a boolean variable which determines whether or not an instance of `DuplicatePrune` is used within the composite PE in order to prune out duplicate results.

It is possible to nest multiple `if/else` conditionals:

```

1 package tutorial.example {
2   use dispel.db.SQLiteQuery;
3   use dispel.core.ListMerge;
4   use dispel.core.DuplicatePrune;
5
6   // Use parallel SQLiteQuery instances to collect results.
7   PE<MulticastQuery>
8     makeMassSQLiteQuery(Integer numberOfSources,
9                         Boolean removeDuplicates) {
10    // Prepare components for connection.
11    Connection expr;
12    Connection[] srcs = new Connection[numberOfSources];
13    SQLiteQuery[] queries = new SQLiteQuery[numberOfSources];
14    ListMerge merge = new ListMerge;
15
16    // Connect SQLiteQuery instances in parallel.
17    for (Integer i = 0; i < numberOfSources; i++) {
18      queries[i] = new SQLiteQuery;
19      expr => queries[i].expression;
20      srcs[i] => queries[i].resource;
21      queries[i].data => merge.inputs[i];
22    }
23
24    if (removeDuplicates) {
25      // If removeDuplicates is true, prune duplicate responses.
26      DuplicatePrune prune = new DuplicatePrune;
27      merge.output => prune.input;
28      return PE( <Connection expression = expr;
29               Connection sources = srcs> =>
30               <Connection data = prune.output> );
31    } else {
32      // Otherwise, leave as is.
33      return PE( <Connection expression = expr;
34               Connection sources = srcs> =>
35               <Connection data = merge.output> );
36    }
37  }
38
39  register makeMassSQLiteQuery;
40 }

```

Figure 3.2: A Dispel script which describes how to construct a composite PE which queries multiple data sources simultaneously and collects the results.

```

if (dayOfTheWeek = "Monday") {
    colour = "gray";
} else if (dayOfTheWeek = "Tuesday") {
    colour = "yellow";
} else {
    colour = "green";
}

```

If the nesting of `if/else` conditionals becomes tedious, or when there are numerous choices for a given condition, the `switch/case` construct may be more useful:

```

switch (dayOfTheWeek) {
    case "Monday"    : colour = "gray";   break;
    case "Tuesday"   : colour = "yellow"; break;
    case "Wednesday" : colour = "red";    break;
    case "Thursday"  :
    case "Friday"    : colour = "blue";   break;
    default          : colour = "green";
}

```

We use the `break` keyword to exit from the `switch` construct — otherwise execution ‘falls through’ and executes all cases until the next `break` statement or until the end of the `switch` construct is reached (so in the above example case `"Thursday"` would execute `colour = "blue"`). The `default` keyword is used to mark the special case where none of the specified cases are satisfied.

The `break` keyword can be used to exit any statement block enclosed by braces (`{` and `}`). Thus `break` can be used to exit an iterator. When iterators are nested, the `break` keyword will only break the inner-most iterator, leaving the outer iterators to execute as normal:

```

for (Integer i = 0; i < 100; i++) {
    for (Integer j = 0; j < 100; j++) {
        if (j == 50) { break; }
        ... // Statement block A.
    }
    ... // Statement block B.
}

```

In the above example, the statement block A will be executed five thousand times whilst statement block B will be executed only one hundred times. Similarly, the `continue` can be used within an iterator to jump to the next iteration without breaking out of the iterator entirely:

```

for (Integer k = 0; k < 100; k++) {
    ... // Statement block A.
    if (j < 50) { continue; }
    ... // Statement block B.
}

```

In the above example, statement block A will be executed one hundred times, whereas statement block B will only be executed fifty times.

Chapter 4

Manipulating the flow of data

Thus far, we have constructed workflows with little concern as to the nature of the data being streamed between PE instances. However different PEs expect different inputs and produce output in accordance with their own specifications. Data may be consumed by inputs at different rates or require that data be consumed across all inputs synchronously, perhaps requiring data to be buffered whilst waiting for other parts of the workflow to catch up. Some PEs serve to push data through a workflow, others serve to pull data along. Many PEs are principally driven by one input, consuming data from other inputs only when needed — once the data-stream into that particular input is exhausted, any continuing input from other sources may be irrelevant. All of these factors play influence on data-intensive workflows.

In *Dispel*, many of these factors can be concealed from the casual user, the enactment platform hidden behind the gateway ensuring that data is adequately streamed and buffered across the length of an executing workflow, and ensuring that the data is of the correct form. The standard PEs provided by *Dispel* are designed to behave in the most intuitive fashion possible, so that most users will instinctively use them correctly.

Nonetheless, control over these factors can be exercised within *Dispel*. Expert users demand the ability to construct more sophisticated workflows, which require more careful attention to the flow of data; casual users rely on expert users to resolve data flow issues and conceal the detail behind the veil of composite PEs with simple interfaces. In this section we introduce the second of *Dispel*'s type systems, the *structural* type system, which concerns itself with the logical structure of data streamed through connections. We also introduce connection modifiers, which can be used to describe how a PE instance consumes and produces data. Finally, we look at how to move between the language and structural type systems using stream literals.

4.1 *Dispel* structural types

Structural types are purely concerned with the data flowing through connections between and within PE instances. In that respect they differ from language types, which are principally used to guide the construction of workflows by providing variables which control iteration, selection and the behaviour of functions. Superficially, structural

types are very similar to language types — we have structural types like `Integer` and `String` — but we also have arbitrarily complex structures (involving lists and tuples of structural types) and partial descriptions (wherein we permit lists of undefined types or only define some of the elements in a tuple).

Conceptually, we consider a data-stream as being a sequence of data elements, each holding a single ‘unit’ of data. However what that ‘unit’ constitutes can vary under different circumstances — in one context, we might expect each unit to be a single integer value, in another we might expect a list of tuples, each tuple containing multiple elements. Generally, the interfaces of PEs are annotated with information about the structure of data units streamed through them:

```
Type ConvertIntegerToReal is PE( <Connection:Integer input> =>
                                <Connection:Real    output> );
```

It is possible to omit structural information when defining new types of PE as we have in the past, in which case the structural type expected by each connection is considered to be `Any`, meaning the data can be in any form. We can get away with this for simple workflows, but for more complex workflows we need to be more careful. Consider a scenario in which the following new PE type is registered:

```
Type Normalise is PE( <Connection input> => <Connection output> );
```

What is the expected input and output of `Normalise`? A user might be able to infer something from the package in which `Normalise` is registered, but essentially the input and output structural types can only be determined by experimentation or searching for additional documentation. Annotating new PE types with structural type information makes them more self-documenting and also permits the gateway through which workflows are enacted to perform type verification and validation of a submitted workflow prior to execution.

More complex structural types can be constructed in a number of ways. *Lists* can be of any length, but each element must have the same abstract structure. A list is defined by enclosing the structure type of the elements within the list within square brackets as so:

```
Type IntListToRealList is PE( <Connection:[Integer] input> =>
                                <Connection:[Real]    output> );
```

Note that connection interface type annotation only describes a single ‘unit’ of data carried within a data stream; it is already given that a description of the content of a stream over a set period of time is an ordered list of values. Thus if a stream outputs a simple sequence of real numbers one after another, then the structure type of the stream is `Real` not `[Real]`. On the other hand, if each element of the stream is itself a list of numbers (perhaps each of different length), then the structural type of the stream will be `[Real]` after all.

Conventional arrays also exist as structural types. Unlike lists, arrays can be multi-dimensional, but must be of fixed size. At the same time, we do not concern ourselves within `Dispel` with what that fixed size actually is. The reasoning behind this is that whilst it is important to know that the output of one PE is an array, or that another PE requires an array of certain dimensionality as input in order to validate a given workflow, we do not need to know the size of the array because we are always merely

piping data from one processing element to another. An array is represented just as for language type arrays:

```
Type MatrixToVectorList is PE( <Connection:Real [][] input> =>
                               <Connection:[Real[]] output> );
```

In this case, `MatrixToVectorList` takes as input a two-dimensional matrix of real numbers and outputs a list of real arrays (demonstrating the combination of list and array structural types).

The other important structural type is that of a tuple. A tuple is an unordered collection of elements of different types. A tuple is enclosed in angle brackets, within which must be found a sequence of typed keys to which values can be assigned:

```
Type GridLocToAngleLoc is
  PE( <Connection:<Integer x, y; String name> gridLoc> =>
      <Connection:<Real angle, magnitude; String label> angleLoc> );
```

Variable names must correspond to the keys used by processing elements themselves to identify the parts of the tuple; as demonstrated above, multiple keys of the same type can be defined together (so one can write `Integer x, y;` instead of `Integer x; Integer y;`).

Complex structural types can be given aliases in the same manner as PE types using an `Stype` declaration:

```
Stype GridLoc is <Integer x, y; String name>;
Stype AngleLoc is <Real angle, magnitude; String label>;

Type GridLocToAngleLoc is PE( <Connection:GridLoc gridLoc> =>
                              <Connection:AngleLoc angleLoc> );
```

What if we do not know (or care) about some or all of the structure of elements passing through a data stream however? In that case we have at our disposal the `Any` generic type and the `rest` identifier. An element of structural type `Any` can take any shape or form, and can be used anywhere, including within arrays, lists and tuples. Meanwhile, the `rest` identifier is used within tuples to encapsulate all tuple elements not referenced prior. To illustrate:

```
Stype GridLoc3D is <Integer x, y, z; String name>;
Stype AbsGridLoc is <Integer x, y, z; Any name>;
Stype GridLocXZ is <Integer x, z; rest>;
```

In the above example, each statement matches those prior to it (though not those after). Note that `rest` must be the last referenced element in a tuple, but can subsume any of the elements within that tuple (for example, it can happily subsume element `y` without subsuming element `z`).

Streams also have a structural type, determined by its content. This type will always be the least common sub-type of all the data elements described within the stream:

```
Stream first = |-1, 2, 3, 4-|;
Stream second = |-"one", "two", "three", "four"-|;
Stream third = |-1, "two", 3, 4.0-|;
Stream fourth = |-<key = 11; value = "eleven">,
                <key = 12; value = "twelve"; note = "2 * 6">-|;
```

In the above example, stream `first` has structural type `Integer`, whilst `second` has type `String`. Stream `third` has structural type `Any`, whilst `fourth` has type `<Integer key; String value; rest>`.

Earlier in Section 2.3 we defined a PE type `SQLToTupleList`. We can now annotate that type with the correct structural types:

```
Type SQLToTupleList is PE( <Connection:String expression =>
                          <Connection:[<rest>] data> );
```

The above declaration states that for each string of text representing an SQL query, a PE of type `SQLToTupleList` produces a list of tuples describing a response to that query (though we do not concern ourselves with the content of those tuples). As a consequence, when we define PE `TutorialQuery` using `lockSQLDataSource` and then register it, the ADMIRE gateway knows that `TutorialQuery` accepts only inputs of type `String` and outputs a sequence of tuple lists. This means that when validating a submitted Dispel workflow, the gateway can verify that any instance of `TutorialQuery` is being fed structurally correct input and output.

It is also possible to define (or re-define) the structural type of specific instances of PEs:

```
TutorialQuery query = new TutorialQuery
    with data as data: [<Integer key; String value>];
```

This can be useful if the user knows that the data being streamed into or out of a PE instance is limited to a particular subset of the input or output normally permitted, and that the PE instance is to be connected to other PE instances which *only* permit data limited to that particular subset:

```
Type DictionaryKeySort is
    PE( <Connection:[<Integer key; String value>] unsorted> =>
        <Connection:[<Integer key; String value>] sorted> );
```

```
DictionaryKeySort sort = new DictionaryKeySort;
query.data => sort.unsorted;
```

Ordinarily, this would be considered unsafe by the enactment gateway, being that superficially `TutorialQuery` instances can produce data which cannot be consumed by instances of `DictionaryKeySort`. Of course this approach only works if the data being produced really is limited to the specified structural type.

What happens if the data produced by one PE instance is incompatible with the expected input of another PE instance to which it has been connected? For example,

if the output of one PE is of structural type `Integer` whilst the expected input of another PE is of type `Real`:

```
Type Thermometer is PE( <> => <Connection:Integer measurement> );
```

```
Type SteamEngine is PE( <Connection:Real temperature> =>
                        <Connection:Real acceleration> );
```

```
Thermometer thermo = new Thermometer;
SteamEngine engine = new SteamEngine;
thermo.measurement => engine.temperature;
```

Implemented as is, `engine` cannot consume the data given. In this case, we need a *converter*. A converter is a PE, generally with one input and one output, which transforms data from one form into another without actually changing its semantic content — for example, PE `ConvertIntegerToReal` defined at the beginning of this section. In `Dispel`, we can explicitly introduce a converter into a workflow as just another PE, or we can rely on the gateway to which we submit our workflow to insert any necessary converters for us:

```
Thermometer      thermo = new Thermometer;
SteamEngine      engine = new SteamEngine;
ConvertIntegerToReal convert = new ConvertIntegerToReal;
thermo.measurement => convert.input;
convert.output => engine.temperature;
```

The ability of a gateway to perform automatic type conversions is useful, but not unlimited. In particular, the ability of a gateway to decompose complex structural types and perform conversions depends on there being a common abstract structure apparent in both the original and target structural type descriptions — one would not expect every possible output of `TutorialQuery` to be convertible into structural type `[<Integer key; String value>]` without user knowledge of the possible results which can be obtained from a given set of queries. In general, it is best not to rely too heavily on the intelligence of any given gateway without good cause, particularly within `Dispel` scripts which could be submitted to different gateways of possibly different type conversion capability.

4.2 Type coercions

To recap, the language type system is concerned with control flow within `Dispel` scripts, whilst the structural type system is concerned with data flow in workflows. It is vitally important that there is no confusion between the two type systems — even in circumstances where language types and structural types appear identical (`Integer` / `Integer` for example), there may be implementational differences between the types used by the `Dispel` parser for validating and executing scripts, and the types used within the actual implementation of the workflow described by a `Dispel` script.

If we want to step between the realm of language types and structural types, then we need the means to convert data between type systems. To this end, we have at our disposal both *implicit* and *explicit* type coercion. Implicit type coercion occurs

between primitive language and structural types, and occurs within stream literals. For example, when we write:

```
Integer x = 3, y = 64, z = -4;

|-x, y, z-| => counter.input;
```

What we are doing is making a conversion of variables `x`, `y` and `z` from the `Integer` language type to the `Integer` structural type. This type of coercion also occurs whenever variables defined within the `Dispel` script are wrapped within tuples, arrays or lists and inserted into stream literals. Note that stream literals are the only way to directly inject information from a `Dispel` script into the data stream described by a PE workflow.

The implicit type coercion described above is adequate for quickly constructing simple data streams to inject into a workflow, but is less suitable for more complex or dynamically constructed streams. For example, how do we feed an array of strings into a stream as a sequence of elements? Assume that we have the stated the following:

```
TutorialQuery query = new TutorialQuery;
String[] queries = new String[3];
queries[0] = "SELECT name FROM littleblackbook";
queries[1] = "SELECT address FROM littleblackbook";
queries[2] = "SELECT number FROM littleblackbook";
```

How can we inject `queries` into `query`? We can insert it bodily into a data stream, but then we have a single element containing an array of strings, which we now know is the incorrect structure for data going into an instance of `TutorialQuery`, being of type `SQLToTupleList`:

```
|-queries-|:String[] => query.expression; // Incorrect structure.
```

We can inject it manually element by element, but this requires that we know the size of the array, and is particularly tedious for large arrays:

```
|-queries[0], queries[1], queries[2]-|:String => query.expression;
```

A more promising approach is to use a loop to build the stream based on the length of the array at the time of computation. We do not want to have to define the loop anew every time however, so we use a *stream function*. Stream functions are simply functions which return elements of type `Stream`, which can then be connected directly to a connection interface or concatenated with other streams. For example, to feed an array element by element into a stream, we only need a function like `stringArrayToStream`:

```
Stream stringArrayToStream(String[] array) {
    Stream stream = |--|;
    for (Integer i = 0; i < array.length; i++) {
        stream += array[i];
    }
    return stream;
}
```

Note that we are still using implicit type coercion to add elements of `array` to `stream`. Having registered this function (or in this case, imported it from package `dispel.stream`), we can then invoke it when injecting the array with the desired content into our workflow:

```
stringArrayToStream(queries) => query.expression;
```

Stream functions are thus the preferred means by which to move between the language and structural type systems for constructed types.

4.3 Connection modifiers

Connection modifiers are a class of modifier attached to connection interfaces which describe how that interface produces / consumes data and how it does that in relation to other interfaces within the same PE. For example, the `initiator` modifier, when applied to an input interface, asserts that the modified interface consumes data before all other input interfaces within a given PE; the other interfaces will only begin to consume data once the initiator terminates.

```
Type LockedSQLQuery is PE( <Connection:String expression;  
                           Connection:String initiator source> =>  
                           <Connection:[<rest>] data> );
```

In this case, an instance of `LockedSQLQuery` handles queries in a similar fashion to an instance of `SQLQuery`; however unlike `SQLQuery`, this PE only reads the data source input once, prior to any queries, and returns results only from that data source. As alluded to earlier, there is a simpler way to define PEs which are basically minor modifications on existing ones:

```
Type LockedSQLQuery is SQLQuery with initiator source;
```

If we are only concerned with a one-off use of a modified PE, we can apply connection modifiers upon instantiation:

```
SQLQuery query = new SQLQuery with initiator source;
```

Connection modifiers are both descriptive and prescriptive; they inform users about the behaviour of certain PEs by their type descriptions, and they permit users to make ad-hoc modifications of existing PEs. Consider the two PEs `Combiner` and `SynchronisedCombiner`:

```
Type Combiner is PE( <Connection[] inputs> => <Connection output> );  
Type SynchronisedCombiner is Combiner with roundrobin inputs;
```

The basic behaviour of `Combiner` is simply to unite its input streams into one stream — no commitment is made as to how data from each input is ordered in the output stream. `SynchronisedCombiner` however does make such a commitment — the effect of the `roundrobin` modifier, when applied to an array of input connection interfaces, is to assert that each interface in the array will consume a single data element (as defined by the highest-level abstract structure of the data given) in order, further

consumption restricted into a full cycle has been performed. Thus, we know that the output of an instance of `SynchronisedCombiner` will consist of all first data elements of each input in order, followed by all second elements in order, and so on until all inputs terminate.

The ability of users to apply *ad-hoc* modifications to existing PEs is limited to the ability of the enactment platform to implement the modified PE; in most cases however, this is a simple matter, since most modifications in `Dispel` can be implemented over existing code via interim interfaces. The full list of available connection modifiers can be found in the `Dispel` reference manual.

4.4 Dispel domain types

The final type system used in `Dispel` is the domain type system. As with structural types, domain type information is appended onto connection interface declarations — in this case using double-colons (`::`). Unlike structural types however, domain types describe not the structure of streamed data, but what a given data flow *is* from the perspective of a domain expert. For example:

```
Type SQLToTupleList is
  PE( <Connection:String::"db:SQLQuery"      expression> =>
      <Connection:[<rest>]::"db:TupleRowSet" data> );
```

From this we can infer that the text strings consumed by interface `expression` represent queries in the SQL language. We can also infer that the lists of tuples produced by interface `data` are sets of entries (rows) drawn from a database. Domain types permit ontological information to be embedded into `Dispel` workflows, which can be used for mapping equivalences between ontological terms, or simply to assist the user in understanding how to use a given PE type correctly.

Domain type descriptors are enclosed in quotation marks and prefixed with a *namespace* identifier. The namespace identifier serves as an alias for the ontology from which a domain type is drawn. Namespace identifiers must be introduced within a `Dispel` script using a `namespace` declaration:

```
namespace db "http://www.admire-project.eu/ontology/db#";
```

In this case, the namespace `db` is defined, drawing upon an ontology for databases (`db` is in fact one of the core `Dispel` namespaces used internally for PE types, along with `dispel`, which need not be explicitly defined in scripts).

Like with structural types, domain type mismatches can be dealt with using converters, albeit in this case ones that often maintain the same structure of data flowing into and out of the converter, but rescale the data based on the perceived and required domain types:

```
Type CelsiusToKelvin is
  PE( <Connection:Real::"measure:Celsius" celsius> =>
      <Connection:Real::"measure:Kelvin"> );
```

For the `CelsiusToKelvin` converter, both input and output are of structural type `Real`, but all celsius values going in are offset to match the equivalent kelvin value. Again, implicit type conversions performed automatically by the gateway to which a `Dispatch` script is submitted may be available, but should not be relied upon too heavily.

Chapter 5

Case studies

5.1 The Sieve of Eratosthenes

The *Sieve of Eratosthenes* is a simple algorithm for finding prime numbers. The algorithm works by counting natural numbers and filtering out numbers which are composite (non-prime). We start with the integer 2 and discard every integer greater than 2 that is divisible by 2. Then, we take the smallest of all the remaining integers, which is definitely a prime, and discard every integer greater than that prime (in this case 3). We continue this process with the next integer and so on, until the desired number of primes have been discovered.

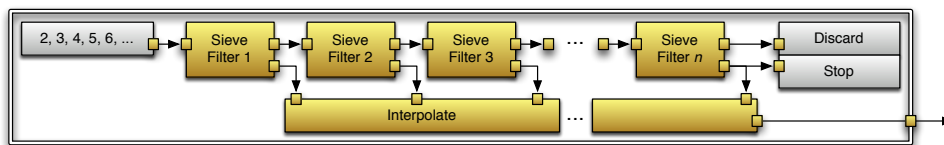


Figure 5.1: The internal composition of the Sieve of Eratosthenes.

The Sieve of Eratosthenes, whilst ultimately a toy application, serves as a useful device to demonstrate such Dispel concepts as connection modifiers, stream comprehensions and cascading termination. The Sieve can be implemented by a pipeline pattern described by a PE constructor. Using such a constructor, it is possible to implement the pipeline for arbitrary numbers of primes. This pipeline pattern will take the form shown in Figure 5.1.

The principal component of a Sieve of Eratosthenes is the filtering element used to determine whether or not a given integer is divisible by the last encountered prime. We define an abstract filter as so:

```
Type AbstractFilter is
  PE( Stype Element is Any;
    <Connection:Element input> =>
    <Connection:Element filtered; Connection:Element unfiltered> );
```


A filter is expected to take a stream of inputs and split it into two streams (the ‘filtered’ stream and the remainder). A filter is not supposed to be transformative, so the output streams should be of the same type as the input stream. Note that an implementation of `AbstractFilter` is not actually expected to discard either filtered or unfiltered elements, since it cannot be predicted which elements will be of most interest in a particular use-case; if a workflow designer has no use of a given output, that output can be redirected to `discard`.

Our filtering element must filter the first integer encountered on its input stream as a prime (assuming the correct construction of the sieve as a whole), discard all successive integers divisible by that prime, and pass onto the next filter all remaining input. We can use a `ProgrammableIntegerFilter` to do the heavy lifting:

```
Type ProgrammableIntegerFilter is
  PE( <Connection:Integer terminator input;
      Connection:String initiator expression;
      Connection[]:Integer lockstep parameters> =>
      <Connection:Integer filtered; Connection:Integer unfiltered> );
```

We need only specify the filtering behaviour via input `expression` (filter all x where x is divisible by some set integer y) and bind any free variables within the filter specification via `parameters` (in this case y , provided by the first integer to enter our filter). To split off the first input, we use `HeadFilter`:

```
Type HeadFilter is AbstractFilter
  with filtered as head, unfiltered as tail,
  @description = "Diverts the head of a stream.";
```

We can then create a constructor `makeSieveFilter` as defined in Figure 5.2; because the constructor has no variable parameters, we immediately construct `SieveFilter` and export it.

We can now define a constructor for the Sieve. The Sieve consists of an array of filters, which sequentially redirect primes to an `Interpolate` PE:

```
Type Interpolate is
  PE( Stype Element is Any;
      <Connection[]:Element inputs> => <Connection:Element output> );
```

In order to ensure that the primes are output in order of discovery, we modify the interpolator’s inputs to be `roundrobin`. We connect each filter’s unfiltered output (being the sequence of numbers not divisible by the first prime encountered) to the next filter, except for the last, which discards all such values (having found the last prime of interest). We use a stream comprehension to generate all integers in sequence from 2 onwards indefinitely, and connect that to the first filter.

That is enough to implement the Sieve; however we also want the Sieve to shutdown once all required primes have been found rather than pour integers into the ether indefinitely. So we specify each filter’s unfiltered output steam as `terminator`, except for the last (which is discarded) — we instead create a connection from that filter’s prime output to `stop` and declare that as being `terminator`. The effect of this is to create a backwards termination cascade once the last prime is generated, which will

```

1 package tutorial.example {
2     // Import filters.
3     use dispel.filter.AbstractFilter;
4     use dispel.filter.HeadFilter;
5     use dispel.filter.ProgrammableIntegerFilter;
6
7     // Define sieve element constructor.
8     PE<AbstractFilter> makeSieveElement() {
9         // Create reference to input connection.
10        Connection:Integer input;
11        // Instantiate internal components.
12        HeadFilter split = new HeadFilter;
13        ProgrammableIntegerFilter divide =
14            new ProgrammableIntegerFilter with parameters.length = 1;
15
16        // Construct internal workflow.
17        |-"x if (x % $0) == 0"-| => divide.expression;
18        input                    => split.input;
19        split.head                => divide.parameter[0];
20        split.tail                => divide.input;
21        divide.unfiltered         => discard;
22
23        // Output first integer received and all indivisible integers.
24        return PE( <Connection input = input> =>
25            <Connection filtered = split.head;
26                Connection unfiltered = divide.filtered> );
27    }
28
29    // Create the sieve element PE.
30    PE<AbstractFilter> SieveElement = makeSieveElement();
31
32    // Register sieve element.
33    register SieveElement;
34 }

```

Figure 5.2: Construction of a filter for the Sieve of Eratosthenes.

```

1  package tutorial.example {
2      // Import sieve components and abstract type.
3      use dispel.core.Interpolate;
4      use tutorial.example.SieveFilter;
5      use dispel.math.PrimeGenerator;
6
7      // Define sieve constructor.
8      PE<PrimeGenerator> makeSieveOfEratosthenes(Integer count) {
9          // Instantiate internal components.
10         SieveFilter filter      = new SieveFilter[count];
11         Interpolate interpolate =
12             new Interpolate with roundrobin inputs, input.length = count;
13
14         // Initialise sieve elements.
15         for (Integer i = 0; i < count - 1; i++)
16             filter[i] = new SieveFilter with terminator output;
17         filter[count - 1] = new SieveFilter with terminator prime;
18
19         // Construct internal workflow.
20         |-x for x in 2..-| => filter[0].input;
21         for (Integer i = 0; i < count - 1; i++) {
22             filter[i].unfiltered => filter[i + 1].input;
23             filter[i].filtered  => interpolate.inputs[i];
24         }
25         filter[count - 1].unfiltered => discard;
26         filter[count - 1].filtered  => interpolate.input[count - 1];
27         filter[count - 1].filtered  => stop;
28
29         // Return all primes generated.
30         return PE( <Connection input = numbers> =>
31                 <Connection primes = interpolate.output> );
32     }
33
34     // Register constructor.
35     register makeSieveOfEratosthenes;
36 }

```

Figure 5.3: The Sieve of Erathosthenes, as a workflow pattern encapsulated within a PE constructor.

```

1 package tutorial.example {
2   // Import abstract PE type and constructor.
3   use dispel.math.PrimeGenerator;
4   use tutorial.example.makeSieveOfEratosthenes;
5
6   // Construct the sieve.
7   PE <PrimeGenerator> SoE100 = makeSieveOfEratosthenes(100);
8   SoE100 sieve100 = new SoE100;
9   Results results = new Results;
10
11  // Construct the top-level workflow.
12  |-"100 prime numbers"-| => results.name;
13  sieve100.primes      => results.input;
14
15  // Submit the workflow.
16  submit results;
17 }

```

Figure 5.4: An execution script for generating the first 100 primes in the Sieve of Erathosthenes.

cascade back through all filters and end the infinite integer stream as well as close the interpolator — thus ensuring the efficient shutdown of the entire Sieve.

The Sieve of Eratosthenes for one hundred prime numbers can now be executed as shown in Figure 5.4.