



This is an electronic reprint of the original article. This reprint may differ from the original in pagination and typographic detail.

Development of Fault Tolerant Grid Applications Using Distributed B

Boström, Pontus; Walden, Marina

Published: 01/01/2005

Document Version Final published version

Link to publication

Please cite the original version: Boström, P., & Walden, M. (2005). *Development of Fault Tolerant Grid Applications Using Distributed B*. (TUCS Technical Reports; Vol. 706). Turku Centre for Computer Science (TUCS). https://urn.fi/URN:NBN:fife2022120168419

General rights Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.



Pontus Boström | Marina Waldén

Development of Fault Tolerant Grid Applications Using Distributed B

TURKU CENTRE for COMPUTER SCIENCE

TUCS Technical Report No 706, August 2005



Development of Fault Tolerant Grid Applications Using Distributed B

Pontus Boström
Åbo Akademi University, Department of Computer Science
Turku Centre for Computer Science (TUCS)
Lemminkäisenkatu 14 A, 20520 Turku, Finland
Pontus.Bostrom@abo.fi
Marina Waldén
Åbo Akademi University, Department of Computer Science
Turku Centre for Computer Science (TUCS)
Lemminkäisenkatu 14 A, 20520 Turku, Finland
Marina.Walden@abo.fi

TUCS Technical Report No 706, August 2005

Abstract

Computational grids have become popular for constructing large scale distributed systems. Grid applications typically run in a very heterogeneous environment and fault tolerance is therefore very important for their correctness. Since the construction of correct distributed systems is difficult with traditional development methods we propose the use of formal methods. We use Event B as our formal framework, which we extend with new constructs such as remote procedures and notifications for reasoning about grid systems. The extended language, called Distributed B, ensures that the application can handle both node and network failures. Furthermore, the new constructs in Distributed B enable straightforward implementation of the specifications, as well as automatic generation of the needed proof obligations.

Keywords: Event B, Grid computing, Fault tolerance, Domain specific languages, Language extensions, Stepwise development

TUCS Laboratory Distributed Systems Design Laboratory

1 Introduction

Computational grids have become a popular approach to handle vast amounts of available information and to manage computational resources. Examples of areas where grids have been successfully used for solving problems include biology, nuclear physics and engineering. Grid computing [8] is a distributed computing paradigm that differs from traditional distributed computing in that it is aimed toward large scale systems that even span organizational boundaries. Since grid applications run in a very heterogeneous computing environment fault tolerance is important in order to ensure their correct behaviour.

The development of correct grid applications is difficult with traditional software development methods. Hence, formal methods can be beneficial in order to ensure their correctness and structure their development from specification to implementation. The Action Systems formalism [4] is a formal method that is well suited for developing large distributed and concurrent systems, since it supports stepwise development. However, it lacks good tool support. The B Method [1], on the other hand, is a formal method provided with good tool support, but originally developed for construction of sequential programs. The B Method can be combined with Action Systems in order to formally reason about distributed systems as in the related methods B Action Systems [16] and Event B [2, 3]. B Action Systems models Action Systems in the B Method, while Event B also extends original B with new constructs. In this paper we use Event B as the basis for our work.

With generic formal languages like Event B, specifications are often unintentionally constructed in such a way that they cannot be implemented or are very difficult to implement efficiently. This can be due to synchronisation issues or the maintenance of atomicity of events. We have previously added extensions [5] to Event B to remedy this problem and to enable reasoning about grid applications using grid communication primitives. However, the extended Event B, referred to as Distributed B, did not consider fault tolerance. Here we will modify Distributed B to enable us to develop also fault tolerant grid applications.

The language Distributed B was designed for developing Grid applications using the Globus Toolkit [11] middleware. We here further develop this language to enable us to construct fault tolerant grid applications using only the fault tolerance mechanisms present in Globus toolkit. We add new constructs for handling exceptions raised during remote procedure calls and timeouts to handle lost notifications. These new features introduced to Event B force the developer to consider the grid environment and fault tolerance throughout the development. Furthermore, the constructs are introduced in such a manner that all needed proof obligations can be automatically generated and the system can be directly implemented.

In Section 2 we give an overview of the grid technology. Event B and the Distributed B extensions are presented in Section 3. In Section 4 we discuss the failure modes and the fault tolerance mechanisms used. Section 5 presents the extensions for developing fault tolerant grid applications. In Section 6 we present a case study using the new language. Implementation issues are discussed in Section 7 and in Section 8 we conclude.

2 Grid Systems

The purpose of grid systems is to share information and computing resources even over organizational boundaries. This requires security, scalability and protocols that are suited for Internet wide communication. The Open Grid Service Architecture (OGSA) [9] aims at providing a common standard to develop grid based applications. This standard defines what services a grid system should provide. A technical infrastructure specification defined by Open Grid Service Infrastructure (OGSI) [6] gives a precise definition of what a grid service is. The Globus Toolkit 3.x [11] is an implementation of the OGSI specification that has become de facto standard toolkit for implementing grid systems. We have chosen this toolkit as grid middleware for our extensions to Event B.

Grid systems usually have a client-server architecture, where the client initiates communication with the server that only responds to the client's request. A client may access several servers concurrently. In our grid applications a server is referred to as a grid service. Grid services as implemented in Globus Toolkit provide features such as remote procedures, notifications, services that contain state, transient services and service data. The main communication mechanism of grid services is remote procedure calls from client to grid service. If a call is unsuccessful the Globus toolkit will raise an exception for the programmer to handle. By using notifications a grid service can asynchronously notify clients about changes in its state. The state of grid services is preserved between calls and grid service instances can be dynamically created. Service data adds structured data to any grid service interface. This way not only remote procedures, but also variables are available to clients. Furthermore, Globus Toolkit contains an index service for managing information and keeping track of different types of services in the grid.

A grid application developed with Globus toolkit can be viewed as a collection of remote objects. A class defining a grid service can be used by first creating an instance of it with a specified grid service handle. These instances correspond to remote objects in Java RMI [14] or CORBA [14] and they can be used almost like normal local objects.

3 The Distributed B Extensions

Since constructing correct distributed applications is difficult with traditional development methods we here propose the use of formal methods to ensure their correctness. We have chosen to use Event B [2, 3], since it is a well supported formalism for modelling distributed systems. Event B is based

on the B Method by Abrial [1] and Action Systems by Back and Kurki-Suonio [4]. Grid applications can be specified in Event B. However, it is not straightforward to construct these specifications in such a manner that they can be efficiently implemented. Earlier we proposed extensions to Event B [5] for specifying and implementing grid applications. These extensions introduced grid features such as remote procedure calls and notifications to Event B. The semantics of the extensions is given by their translation to B. Note that we translate to B and not Event B, since the current tools for Event B also translate the specifications to B for verification.

We will here present the Distributed B extensions to Event B by a schematic example of a grid application. Furthermore, we show how these extensions are translated to B for verification. The abstract specification of an application is first written in Event B. Using a special grid refinement grid features are then introduced into the system in a stepwise manner. For a more formal treatment of the extensions the reader is referred to our previous paper [5].

3.1 Event B

The abstract specification of a grid application is given in a system-model written in Event B. A system-model contains constructs for defining sets, variables, an invariant, as well as events defining the behaviour of the component. The events consists of guarded substitutions. The semantics of the substitutions are given by the weakest precondition calculus introduced by Dijkstra [7]. When the guard of the event evaluates to true the event is said to be enabled. Enabled events are chosen non-deterministically for execution. When there are no enabled events left the event system terminates.

SYSTEM

$$C$$

VARIABLES
 v
INVARIANT
 $I(v)$
INITIALISATION
 $Init(v)$
EVENTS
 $E_1 \triangleq$
ANY y WHERE $G_1(v, y)$ THEN $S_1(v, y)$ END ;
 $E_2 \triangleq$
WHEN $G_2(v)$ THEN $S_2(v)$ END
END

The variables-clause defines a set of variables, v. The invariant I(v) defines the type of these variables as well as additional properties that should be preserved during the execution of the system. The initialisation Init(v) assign initial values to the variables. The *events*-clause contains the events, here E_1 and E_2 , that describe the behaviour of the system. In event E_1 the values of the local variables y are non-deterministically chosen. When the predicate $G_1(v, y)$ evaluates to true the event is enabled and the substitution $S_1(v, y)$ can be executed. The event E_2 has no local variables and when $G_2(v)$

evaluates to true, $S_2(v)$ can be executed. When both guards $G_1(v, y)$ and $G_2(v)$ evaluates to false the event system terminates.

3.2 The Grid Service Machine

The grid service machine [5] has been introduced for giving the abstract specification of grid services. It extends the system model of Event B with constructs for specifying remote procedures and notifications to send. The grid service machine can be viewed as a class of which clients can obtain remote objects.

The example grid service machine presented below has a set of variables x, the remote procedure *Proc* and the event E_A . Moreover, it can send notifications N_1 and N_2 . A client can call the remote procedure and thereby enable the event. The event is then executed concurrently with the events of the client. When all the events (here only one) in the grid service has become disabled a notification is sent to the client.

In order to give our new constructs meaning and to use the tool support of B the grid service machine is translated to an ordinary B machine. Throughout this subsection the grid service machine is presented in the left column and its translation in the right one.

The B translation of the grid service machine generates an additional set, a constant and a variable. These are needed to model the instance management in B.

GRID SERVICE	MACHINE
\mathcal{A} –	\mathcal{A}_V
VARIABLES	SETS
x	A INSTS
INVARIANT	CONSTANTS
$x_i \in T \wedge I_A(x)$	A $null$
INITIALISATION	PROPERTIES
$Init_A(x)$	$A null \in A INSTS$
	VARIABLES
	x, A Insts
	INVARIANT
	$A Insts \subseteq A INSTS \land$
	A $null \notin A$ $\overline{I}nsts \land$
	$x_i \in A Insts \to T \land$
	$I_{AV}(x, \overline{A} \ Insts)$
	INITIALISATION
	$x_i := \emptyset \parallel$
	$A Insts := \emptyset$
	$A_1usus = \emptyset$

The deferred set A_INSTS models all possible instances of grid service machine \mathcal{A} . The constant A_null models the empty instance, while the variable A_Insts models the set of instances that are currently in use by the client. Note that a grid service machine instance can only be accessed from one client. The variables are translated to take into consideration the instances. Each variable becomes a function from instance to the type defined in the grid service machine. The variables in the translation are initialised to empty sets. The instance will be initialised to $Init_A(x(inst))$ according to the initialisation clause of the grid service machine in the constructor of the instance.

The remote procedures are translated to take the instance as an additional parameter. This is needed since a remote procedure call is always made to an instance of a grid service machine. Additionally, for each variable x_i in \mathcal{A} , a remote procedure, $GetX_i$, for enabling a client to read the variable value is automatically introduced. These procedures can then be used by a client as normal remote procedures. The events are then also translated to take into account every possible instance via an *any*-statement.

```
REMOTE PROCEDURES
                                                  OPERATIONS
Proc(p) =
                                                  Proc(inst, p) \stackrel{\circ}{=}
  PRE P(p)
                                                     PRE inst \in A Insts \wedge P(p)
  THEN S_p(x, p)
                                                     THEN S_p(x(inst), p)
  END;
                                                     END ;
                                                  rr \leftarrow GetX_1(inst) =
                                                     PRE inst \in A Insts
EVENTS
\begin{array}{l} E_A \stackrel{\circ}{=} \\ \textbf{WHEN} \ G_A(x) \end{array}
                                                     THEN rr := x_i(inst)
                                                     END
  THEN S_A(x)
  END;
                                                  E_A =
                                                     ANY inst WHERE
                                                       inst \in A\_Insts
                                                      THEN
                                                         WHEN G_A(x(inst))
                                                        THEN S_A(x(inst))
                                                         END
                                                     END
```

Notifications to be sent are defined in the *notifications*-clause.

```
\begin{array}{l} \textbf{NOTIFICATIONS} \\ N_1 \doteq \\ \textbf{GUARANTEES} \ Q_1(x) \\ \textbf{END} \ ; \\ N_2 \doteq \\ \textbf{GUARANTEES} \ Q_2(x) \\ \textbf{END} \ ; \\ \textbf{END} \ ; \\ \textbf{END} \end{array}
```

Each notification consists of a guarantees-statement. The guarantees-statement means that a certain condition $Q_i(x)$ holds in this grid service when the corresponding notification is sent. The notification handling is performed in the client. However, an invariant is added in the grid service machine to ensure that when the event system terminates the guarantees-statement of at least one notification evaluates to true, $\forall inst.(inst \in A_Insts \land \neg G_A(x(inst)) \Rightarrow (Q_1(x(inst)) \lor Q_2(x(inst)))).$

The translated B machine contains two additional automatically generated operations. These are the constructor and destructor of instances.

```
r \leftarrow A \quad GetNew \stackrel{}{=}
    ANY inst WHERE
        \begin{array}{l} A\_Insts \neq A\_INSTS - \{A\_null\} \land \\ inst \in A\_INSTS - A\_Insts \land \end{array}
        inst \neq A \quad null
    THEN
        A\_Insts := A \quad Insts \cup \{inst\} \parallel
        Init_A(x(inst))
        r := inst
    END :
A Destroy(inst) =
    PRE inst \in A Insts
    THEN
        A \quad Insts := A \quad Insts - \{inst\}
        x_i := \{inst\} \triangleleft \overline{x_i} \parallel
    END
END
```

The constructor adds the instance to the set of used instances and initialises it. In order to find instances to use, the addresses of all the instances are stored in the index service of Globus toolkit. The client can then locate a new instance by asking the index service to return the address to new free ones. The destructor of instances removes the instance from the set of used instances and modifies the variables x to reflect this change.

3.3 The Grid Refinement Machine

One of the benefits of Event B is that it supports refinement and thereby it enables stepwise development of systems. Refinement of a component preserves the behaviour of the abstract component while making it more concrete by adding variables and additional behaviour (events). An event is said to refine another event if the guard is strengthened and the behaviour of the substitutions is preserved.

In order to refine grid service machines and for introducing grid features to ordinary Event B specifications we introduce a grid refinement machine. The grid refinement machine extends the ordinary refinement of Event B with constructs for accessing grid service machine instances, refining remote procedures and handling notifications. The instances of grid service machines are modelled as ordinary variables. Remote procedure calls are modelled by ordinary operation calls. When all events in a grid service machine instance have become disabled a notification is received from it. The notification handlers in the grid refinement are then executed once for each notification. In order to verify that a grid refinement machine is a refinement of a more abstract specification it is translated to a refinement machine in B.

Grid refinement machines are here presented with the example C_1 that uses instances of grid service machine \mathcal{A} presented in Subsection 3.2. This grid refinement is a refinement of the abstract system \mathcal{C} presented in Subsection 3.1. The grid refinement C_1 contains the old variables v from \mathcal{C} , new variables w and a set of instances a_1, \ldots, a_n of grid service machine \mathcal{A} . The events E_1 and E_2 in the abstract specification are refined by more concrete events. The grid refinement machine also contains a new event F that only modifies the new variables w and a notification handler $N_1Handler$ for handling notifications of type N_1 from instances of \mathcal{A} . A notification handler is also considered to be a new event and, hence, it can only assign to new variables. The grid refinement machine C_1 is given below in the left column and its translation to B in the right.

In Distributed B a new structuring mechanism [5], *references*, is used in the grid refinement to handle instances of grid service machines. This construct is translated into an *includes* statement in B. The concurrent execution of the events in the grid service machine instances and grid refinement is modelled by promoting all the events in the grid service machine into the refinement in the B model.

GRID_REFINEMENT	REFINEMENT
\mathcal{C}_1	\mathcal{C}_{1V}
REFINES	REFINES
С	\mathcal{C}
REFERENCES	INCLUDES
\mathcal{A}	\mathcal{A}_V
VARIABLES	PROMOTES
v, w, a_1, \ldots, a_n	E_A
INVARIANT	VARIABLES
$a_1 \in \mathcal{A} \land$	$v, w, a_1, \ldots, a_n, A_notif$
	INVARIANT
$a_n \in \mathcal{A} \wedge$	$a_1 \in A_INSTS \land$
$J(v,w,a_1,\ldots,a_n)$	$a_1 \in A_Insts \cup \{A_null\} \land$
	$J_V(v,w,a_1,\ldots,a_n)\wedge$
	$A_notif \in A_Insts \to BOOL$

We introduce variables modelling instances of \mathcal{A} , a_i . The type of the variables a_i are given in the grid refinement as \mathcal{A} which is then translated to the corresponding representation of instances in the traditional B model. The invariant J is also modified by the translation to take into account the changed representation of instances.

The notification handler can only be executed once for each notification from \mathcal{A} . Hence, we introduce variable A_notif of type function from the instances in use to boolean values. If the value is *true* for an instance, notifications can be received from it. If the value is *false* the notification has already been handled.

The initialisation of the grid refinement is similar to its B translation. The automatically generated variable A_notif is initialised to the empty set, which corresponds to the value of A_Insts .

INITIALISATION	INITIALISATION
$a_i := A_null \parallel$	$a_1 := A_null \parallel$
$a_n := A_null \parallel$	$Init'(v,w) \parallel$
Init'(v,w)	$A_notif := \varnothing$

Events of the grid refinement and the B translation are again similar.

```
EVENTS
                                                   OPERATIONS
E_1 \stackrel{}{=}
                                                   E_1 \stackrel{}{=}
   ANY y WHERE
                                                      ANY y WHERE
     G_1'(v,y)
                                                         G_1'(v,y)
   THEN
                                                      THEN
      S'_1(v, y);
                                                         S'_{1}(v, y);
      a_i \leftarrow A\_GetNew
                                                         a_i \leftarrow A \quad GetNew;
   END;
                                                         A noti\overline{f}(a_i) := TRUE
                                                      END;
E_2 \stackrel{}{=}
   WHEN G'_2(v, w)
                                                   E_2 \stackrel{}{=}
   THEN S'_2(v, w)
                                                      WHEN G'_2(v, w)
                                                      THEN S'_2(v, w)
   END;
                                                      END;
F \stackrel{}{=}
   WHEN H(v, w, a_i)
                                                   F \stackrel{}{=}
   THEN
                                                      WHEN H(v, w, a_i)
                                                      THEN
      S_n(w) \parallel
      a_i.Proc(f(v,w))
                                                         S_n(w) \parallel
   END;
                                                         Proc(a_i, f(v, w));
                                                         A\_notif(a_i) := TRUE
                                                      END;
```

The remote procedure calls are translated in such a way that the instance is given as the first parameter to the procedure. Furthermore, the assignment $A_notif(inst) := TRUE$ is added after each call to a manually defined remote procedure and after a call to A_GetNew . This is done to enable the correct handling of notifications originating from that instance. However, it is not added after calls to the automatically generated remote procedures $GetX_i$, since they are read-only procedures that do not modify the state of the instance. Note that event above E_2 does perform any remote procedure calls and hence it does not contain assignment $A_notif(inst) := TRUE$.

Finally we study the notification handlers.

```
NOTIFICATION HANDLERS
                                              N_1Handler =
                                                 ANY inst WHERE
N_1Handler =
  NOTIFICATION N<sub>1</sub>
                                                   inst \in A Insts \land
                                                    A notif(inst) = TRUE \land
  SOURCE inst \in A
  THEN T_{N1}(inst, w, a_i, a_j)
                                                    \neg G_A(x(inst)) \land Q_1(x(inst))
  END
                                                 THEN
END
                                                   T_{N1}(inst, w, a_i, a_j);
                                                    A\_notif(inst) := FALSE
                                                 END
                                              END
```

Notifications of type N_1 from instance *inst* of \mathcal{A} are handled when all the events in the instance have become disabled, $\neg G_A(x(inst))$, and the condition in the guarantees-statement for that notification holds, $Q_1(x(inst))$. The variable A_notif is assigned false for this instance to denote that the notification has been handled.

3.3.1 **Proof obligations**

Since we develop grid applications in a stepwise manner, we need to show that each new specification is a refinement of the specification developed in the previous step. In order to show that an Event B or a Distributed B component is a refinement of another component the following properties must hold [2, 16]:

- 1. The initialisation of the concrete specification has to be a refinement of the initialisation of the abstract specification.
- 2. All events in the abstract specification have to be refined by corresponding events in the concrete specification.
- 3. New events that refine skip can be introduced in the concrete specification.
- 4. The new events have to terminate when executed in isolation.
- 5. The concrete system is not allowed to terminate more often than the abstract system.
- 6. The guard of a remote procedure cannot be strengthened.

The proof obligations for rules 1-5 [2, 16] can be automatically generated by the tools of Event B. The proof obligation for rule 6 [16] needs some additional constructs. Alternatively, the developer can be restricted to only use non-guarded statements in the remote procedures. The proof obligations can then be discharged by the automatic and interactive provers for B.

4 Fault Tolerance using Globus Toolkit

Grid applications run in a very heterogeneous computing environment. This means that fault tolerance is highly important for the correct behaviour of the application. In this paper we develop extensions to Event B for developing fault tolerant grid applications using the fault tolerance mechanisms in Globus toolkit. Currently the only fault tolerance mechanism in Globus toolkit is exceptions due to failed remote procedure calls. More advanced fault tolerance mechanism such as Replication [14] and Check pointing [14] are not yet supported by the toolkit. Even with support for advanced fault tolerance mechanisms in the middleware the application needs to consider faults, since these mechanisms in the middleware might not be sufficient to handle the faults transparently.

It is not feasible to construct a system that can tolerate all types of faults [14]. We will limit the fault tolerance of Distributed B to handle two types of faults, which we consider to be most important: Firstly, a grid service instance can stop (crash) and be restarted. Secondly, network connections can fail to deliver messages to desired destinations. We do not consider situations where grid services does not satisfy its specification or where an attacker can deliberately force a grid service instance to produce erroneous results since we have proved our application correct and Globus toolkit has

a very comprehensive security infrastructure. Furthermore, the grid middleware can use TLS or other secure protocols for communication and therefore we can assume that data has not been corrupted during transmission.

4.1 Faults and Fault Detection

We can identify five distinct faults that can occur in remote procedure calls from a client to a server grid service instance [14]:

- 1. The server grid service instance has crashed before the call.
- 2. The network connection fails when calling a remote procedure.
- 3. The server instance fails during the call.
- 4. The network connection fails when returning the result.
- 5. The client crashes during the remote procedure call.

In the client the faults 1-4 are not easily separated from each other and they are handled by not using the failed server grid service instance anymore. These faults are directly detected in the client by the Globus toolkit middleware that then raises an exception. However, the server grid service instance can crash and be restarted and it needs to detect the crash in order to respond correctly to the client. A restart can be detected if the instance creates a file on startup using its address as the filename. If the file creation fails due to the fact that the file already exists, it can be assumed that the grid service has been restarted. This holds, since file creation is atomic if the correct method is used.

In order to handle the fifth fault above the called instance needs to be able to identify when the calling client has failed, i.e., the instance has become an orphan [14]. To discover this we introduce an is-alive check in the client and an *is-alive timeout* in the server grid service instance. We first introduce a new remote procedure in every grid service instance called *isAlive.* This remote procedure returns *true*, if the instance is accessible. In the client the *is-alive check* for a given grid service instance consists of calling remote procedure *isAlive* of the instance periodically with a specified time interval. If the call fails, an exception is raised by the Globus toolkit and the client knows that there is a problem with the instance or with the network connection. The client then stops using the instance. If the isAliveprocedure was not called with the specified time interval in the server grid service instance, an *is-alive timeout* is triggered. This means that the client has either failed or stopped using the instance. The server grid service instance is then reset in order to enable other clients to use it. It is important to note that the times for the timeout mechanism should be chosen so that the client will get the exception for a failed call to *isAlive* before the server resets itself. This can be accomplished by choosing the timeout in the server grid service instance to be greater than the time period that the client calls

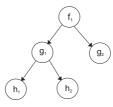


Figure 1: The tree formed by the grid service instances

isAlive with an appropriate amount of buffer time. This check for orphans is implemented in a layer on top of Globus toolkit that then provides an easy to use interface to Distributed B.

There is an additional fault, not directly related to remote procedure calls, that needs to be considered when allocating new grid service instances. There might not be any available grid service machine instances of the correct type when the client tries to obtain new instances to use. The problem can be due to broken network connections, a broken index service or exhaustion of the pool of available services. Note that the index service is here a single point of failure for the grid application. Hence, if it fails, the application might not function correctly. The index service could, however, be transformed into a replicated index service on several computers or we could use a peer-to-peer index service. The instance allocation is also implemented in a layer on top of Globus toolkit to present a more easy to use interface to Distributed B.

4.2 Use Cases Describing the Fault Tolerance Mechanisms

In order to clarify the behaviour of the system in the presence of faults we here present the use cases for the two most complicated failure modes. We consider a tree of grid service instances as shown in Figure 1. The grid service instance f_1 references instances g_1 and g_2 , where g_1 then references h_1 and h_2 .

4.2.1 The Middle Node Crashes

The first use case describes the actions taken when the grid service instance g_1 crashes. The instance g_1 is here a middle node that has client f_1 and has references to grid service instances h_1 and h_2 .

- 1. Instance g_1 crashes. There are now two cases to consider, steps 2 and 3.
- 2. If the instance g_1 is unreachable from its client, instance f_1 :
 - (a) An exception is raised in f_1 , due to a failed remote procedure call or *is-alive check*. (A notification timeout can also be triggered by this condition).

- (b) The instance g_1 is removed from the set of used instances in f_1 .
- 3. If the instance g_1 is reachable from f_1 but it has recovered from a failure and is in a restarted state:
 - (a) The instance g_1 waits for an *is-alive timeout*. It raises exceptions for the clients remote procedure calls and *is-alive checks*. That way g_1 can make sure that the client f_1 knows that g_1 has failed.
 - (b) Instance f_1 detects that instance g_1 has failed and it removes it from the set of used instances.
 - (c) The instance g_1 receives an *is-alive timeout* and resets itself.
- 4. The instances h_i notice that g_1 has failed when they receive *is-alive* timeouts and they are then reset.

4.2.2 The Network Connection Fails.

The second use case describes the behaviour of the system when a network connection failure between, the client f_1 and the middle node g_1 is noticed.

- 1. The network connection between f_1 and g_1 fails.
- 2. An exception is raised in f_1 due to failed remote procedure call or *is-alive check*.
- 3. The instance f_1 removes g_1 from its set of used instances.
- 4. The instance g_1 receives an *is-alive timeout*, since its services are not requested anymore. It then resets itself.
- 5. The instances h_i reset themselves when they receive *is-alive timeouts* due to the reset of g_1 .

4.2.3 Final Notes

What applies to a middle node g_i also applies to a leaf node h_i . The only difference is that we do not have to consider any child nodes. If the root node f_1 crashes the program terminates and no result is obtained.

5 Fault Tolerance in Distributed B

Previously we have developed a language called Distributed B [5] based on Event B for constructing fault free grid applications. In this section we modify the language to also incorporate fault tolerance mechanisms as described in the previous section, Section 4. Distributed B provided a way to implement grid applications and the fault tolerance mechanisms will also be introduced in a manner that ensures that they can be implemented.

5.1 Specification of Instance Management

A challenging problem when introducing fault tolerance is the management of instances. In order to correctly reason about instances we need to have a model of the instance management. The model in the original version of Distributed B is described in Subsection 3.2. During the translation of grid service machine \mathcal{A} to original B we there added features such as a set A_INSTS modelling all instances of \mathcal{A} , a constant A_null modelling the empty instance, a variable A_Insts modelling the instances in use, as well as a constructor and destructor of instances. This model of the instance management is not changed when we introduce fault tolerance into Distributed B. However, we did not specify how an instance should be treated if a fault was encountered during a remote procedure call to it or if notification was never received from it. In order to also handle faults, an instance a_i of \mathcal{A} is immediately removed from the set of instances in use, A_Insts , if a problem with it is encountered.

The model of instance management described above serves as an abstract specification that need to be implemented. To increase the confidence in the system it can be developed in a stepwise manner using, e.g., Event B. The model is refined to consider the behaviour of grid service instances, the communication between them, as well as the faults described in Section 4.

We here give an overview of a refined, more concrete model of the behaviour of an instance of \mathcal{A} using the statechart diagram shown in figure 2. The users of fault tolerant Distributed B will not have to consider this refined model of instance management, since it is built into the language. It only serves as an illustration on how the instance management is implemented using Globus toolkit. The instance can be in three different states: *idle* when the instance is not reserved by a client; *busy* when a client has reserved the instance and *restarted* when the instance has detected that it has been restarted. The event *reserve* is generated when a client reserves the instance by a call to the operation A GetNew. The reservation is modelled by the transition from state *idle* to state *busy*. When the event *reserve* occurs the instance is initialised, modelled in Distributed B by executing the initialisation-clause of the grid service machine. The transition release from state *busy* to state *idle* models that the client releases the instance by a call to operation A Destroy. The event is timeout denotes that an is-alive timeout occurred and the instance is reset, i.e., it enters the state *idle*. The *restart* event models restart of the grid service instance by taking the instance to state *restarted*. When a grid service instance belongs to the set of instances in use, A Insts, it is in either state busy or restarted, since it has then been reserved by a client. Note that the instance will raise an exception when its remote procedures are called in state *restarted*.

5.2 Extensions to Distributed B

In order to handle faults as described in Section 4 and to handle instances as described in the previous subsection we add new constructs to the original

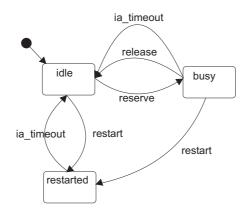


Figure 2: The behaviour of a grid service instance

version Distributed B. New constructs are needed for handling exceptions raised from remote procedures and timeouts for notifications. We also add constructs for handling orphans with the *is-alive check*. The mechanism for handling faults only affects the grid refinement machine and hence, the constructs in the grid service machine remains the same. For simplicity we will here consider only grid applications that do not have middle nodes. Hence, a refinement of a grid service machine cannot here reference other grid service machines.

As an example of a fault tolerant grid refinement machine we modify grid refinement C_1 in Subsection 3.3 to handle faults. The fault tolerant grid refinement C_{FT} has the same variables (v, w) and instances (a_1, \ldots, a_n) as C_1 . The remote procedure calls to instance a_i are modified to consider faults. The notification handling mechanism is also modified to consider timeouts due to lost notifications. Furthermore, we need an additional event to handle exceptions raised by *is-alive checks*. The fault tolerant Distributed B models are again translated to B for verification purposes. Below the fault tolerant Distributed B constructs are given in the left column and their translation to B is given to the right. The complete example can be found in Appendix A.

5.2.1 Remote Procedures

Remote procedure calls can fail in several ways as outlined in Subsection 4.1. We add a *call*-substitution for calling remote procedure in order to model the exception handling mechanism present in Java and Globus toolkit.

The *call*-part of this construct contains a remote procedure call. The *exception*-part then gives the non-guarded substitution that is executed if an

exception is raised during the call. The *call*-substitution is used whenever a remote procedure in a grid service instance is called as shown in event F. The *call*-substitution also needs to be used when calling A_GetNew to allocate a new instance since it can fail as described in Subsection 4.1.

```
F \stackrel{}{=}
   WHEN H(v, w, a_i)
                                                   FOk \stackrel{}{=}
   THEN
                                                      WHEN H(v, w, a_i)
      S_n(w) \parallel
                                                      THEN
      CALL a_i. Proc(f(v, w))
                                                         S_n(w) \parallel
      EXCEPTION T_f(w, a_i)
                                                         Proc(a_i, f(v, w));
      END
                                                         A\_notif(a_i) := TRUE
   END
                                                      END;
                                                   FFail =
                                                      WHEN H(v, w, a_i)
                                                      THEN
                                                         S_n(w) \parallel
                                                         A Destroy(a_i);
                                                         A notif := \{a_i\} \triangleleft A notif;
                                                         T_f(w, a_i)
                                                      END ;
```

The event F is translated to two separate operations in B, FOk and FFail, which both refine the abstract specification of F. One operation, FOk, models the successful execution of the remote procedure and the other, FFail, models the failed execution [2]. The operation FOk is translated as F previously in Distributed B, i.e., without the exception handler. In FFail the instance a_i is removed from the set of instances in use and the exception handler is executed. The guard is identical in both events modelling that any remote procedure call can fail.

5.2.2 Notifications

Failure of grid service instances can cause notifications to never be sent or a failure of network connections can cause notification messages to get lost. In order to detect this we introduce a timeout exception that occurs when a notification has not arrived within the desired time. We do not specify the time explicitly in B but we model the timeout with an event that is chosen non-deterministically for execution. The exact time is a detail that is considered later during the implementation phase.

The grammar of the notification handler substitution is modified in order to incorporate the timeout mechanism.

```
Notif_handler::="NOTIFICATION" Name
"SOURCE" Name ":" Name
"THEN" NG_Substitution
"TIMEOUT" NG_Substitution
"END"
```

The *notification*-part gives the name of the notification and the *source*-part gives the source of the notification as $<instance>:<grid\ service\ machine>$. When a notification is received the non-guarded substitution in the *then*-part

is executed. The *timeout*-handler denotes the substitution that is executed when a notification does not arrive within the desired time.

A notification handler for notification N_1 from instance *inst* of \mathcal{A} can be defined as follows.

```
NOTIFICATION HANDLERS
                                                  N_1HandlerOk =
N_1Handler =
                                                     ANY inst WHERE
   NOTIFICATION N<sub>1</sub>
                                                        inst \in A Insts \land
   SOURCE inst \in A
                                                        A notif(inst) = TRUE \land
   THEN T_{N1}(inst, w, a_i, a_j)
                                                         \neg \overline{G}_A(x(inst)) \land Q_1(x(inst))
   TIMEOUT T_t(inst, w, a_i, a_j)
                                                     THEN
   END
                                                        T_{N1}(inst, w, a_i, a_j) \parallel
                                                        A\_notif(inst) := FALSE
                                                     EN\overline{D};
                                                  N_1HandlerFail =
                                                      ANY inst WHERE
                                                        inst \in A Insts \land
                                                         A notif(inst) = TRUE \land
                                                      THEN
                                                         A Destroy(inst);
                                                         A \quad notif := \{inst\} \triangleleft A \_ notif;
                                                        T_t(inst, w, a_i, a_j)
                                                      END
```

As for remote procedure calls, the notification handler $N_1Handler$ is translated to two different operations in B, $N_1HandlerOk$ and $N_1HandlerFail$. The notification handler is a new event and hence, both translated operations have to refine *skip*. The operation, $N_1HandlerOk$, modelling successful notification handling is translated as in previous version of Distributed B presented in Subsection 3.3. Failed delivery of a notification is modelled by the execution of the timeout handler, $N_1HandlerFail$. Since it is not known if all events in the grid service machine instance *inst* of \mathcal{A} have become disabled before the failure, only $A_notif(inst) = TRUE$ is present in the guard. The instance *inst* is also here removed from the set of instances in use.

Note that T_{N1} can contain a remote procedure call, which means that the translation needs to be performed in two steps. First the the notification handler is split into two events as described above and then the event containing the remote procedure call is translated as in Subsubsection 5.2.1.

5.2.3 Checking if an Instance is Alive

In order to handle orphans we introduce an *is-alive check* in the grid refinement machine. This check is performed in the grid refinement machine for all instances of all referenced grid service machines. If an instance is unavailable an exception is then raised. To handle these exceptions we introduce a new clause *is alive handlers*

"THEN" NG_Substitution "END"

The handler substitutions for *is-alive check* exceptions consists of two parts; a source given as $\langle instance \rangle :< grid service machine \rangle$ and a non-guarded substitution describing the behaviour of the system when an exception is raised. The handler for grid service machine \mathcal{A} is given below with its translation.

	$\begin{split} IAHandler &\doteq \\ \mathbf{ANY} \ inst \ \mathbf{WHERE} \\ inst \in A_Insts \\ \mathbf{THEN} \\ A_Destroy(inst); \\ A_notif := \{inst\} \preccurlyeq A_notif; \\ T_{ia}(inst, w, a_i, a_j) \\ \mathbf{END} \\ \mathbf{END} \end{split}$
--	--

The handler is translated to an *any*-substitution in B. The *any*-substitution models that an exception can be raised at any time for all instances in use. The handler for failed *is-alive checks* is, like the handlers for notifications, a new event that refines *skip*. As for the failed remote procedure calls and notifications we remove the failed instance from the set of used instances and remove all the variables of the instance.

6 Case study

To better illustrate the use of our fault tolerant extensions to Distributed B we provide a small case study. The application in the case study is abstract, but it illustrates how a problem can be decomposed into a distributed grid application. The application we develop computes the value of a function f of two values v_1 and v_2 . The function f could, for example, be a complicated matrix computation that takes two matrices as arguments and produces a third matrix as the result. We assume that the computation of $f(v_1, v_2)$ can split up into computations $f_1(v_1)$ and $f_1(v_2)$ and a computation that composes the final result $g(f_1(v_1), f_1(v_2))$. We thus have the invariant $\forall v_1, v_2.(f(v_1, v_2) = g(f_1(v_1), f_1(v_2)))$. First we model the computation abstractly in an Event B specification. We then decompose the problem and introduce grid features in the refinement.

6.1 The abstract specification

The first abstract specification of the application is written in Event B. The specification abstractly models the computation of f and the possible failures that can occur. First we need to introduce a machine GDEF that defines the set of possible values (matrices) V.

```
MACHINE

GDEF

SETS

V

END
```

The abstract Event B specification COMPF models the application that computes the function f of two values. The function f used in the computation is modelled as a constant total function. The values are given by the variable values, where values(1) is the first argument and values(2) is the second. The result of the computation is then stored in variable result. We introduce a variable state for modelling the progress of the computation. State *idle* models that the computation has not yet started and state working that the computation is in progress. When computation is successfully completed it terminates in state done, while a failed computation terminates in state failed.

```
SYSTEM
 COMPF
SEES
 GDEF
SETS
 STATE = \{idle, working, done, failed\}
CONSTANTS
PROPERTIES
   f \in (V \times V) \to V
VARIABLES
 state, values, result
INVARIANT
   state \in STATE \land
   values \in 1..2 \rightarrow V \wedge
   result \in V \land
   (state = done \Rightarrow result = f(values(1), values(2))))
INITIALISATION
   state := idle \parallel
ANY v_1, v_2 WHERE v_1 \in V \land v_2 \in V
   THEN values := \{1 \mapsto v_1, 2 \mapsto v_2\} END \parallel
   result :\in V
```

The abstract specification contains three events. The event *Start* assigns new values to the variables values(1) and values(2) and initialises the computation. The event *Finish* computes the function f(values(1), values(2)) and stores the result in variable *result*. The computation can fail while it is in progress as modelled by the event *Fail*.

```
EVENTS
Start =
  ANY v_1, v_2 WHERE
     v_1 \in V \land v_2 \in V \land
     state = idle
   THEN
     values := \{1 \mapsto v_1, 2 \mapsto v_2\} \parallel
     state := working
  END :
Finish =
  SELECT state = working
  THEN
     result := f(values(1), values(2)) \parallel
     state := done
  END;
Fail =
  SELECT state = working
  THEN state := failed
  END
END
```

6.2 The grid service machine

We like to use a distributed grid application to compute the function f and, hence we need to decompose the computation of f. The function f can be considered as the composition of functions g and f_1 . To enable decomposition of the computation of f into computation of f_1 and g we first introduce a grid service machine that computes the function f_1 of a value. The function f_1 is, as function f, modelled as a constant total function. The client can then use instances of this machine for the partial computation. The grid service machine is translated to a B machine for verification. Throughout this subsection, the grid service machine is presented in the left column and its translation is given to the right as in previous sections.

The grid service machine RF1 contains a set giving the possible states the service can be in. The translation also contains an automatically generated set and constant for modelling the grid service instances.

GRID SERVICE	MACHINE
RF1	$RF1_V$
SEES	SEES
GDEF	GDEF
SETS	SETS
RS STATE =	RS STATE =
$\{f1_init, f1_start, f1_done\}$	$\{\overline{f1}_init, f1_start, f1_done\};$
CONSTANTS	$RF1_INSTS$
f_1	CONSTANTS
PROPERTIES	$f_1, RF1_null$
$f_1 \in V \to V$	PROPERTIES
	$f_1 \in V \to V \land$
	$RF_null \in RF1_INSTS$

We also introduce variables $f1_state$, $f1_value$ and $f1_result$. The variable $f1_state$ gives the progress of the computation. The value to use in the computation of f_1 is given by the variable $f1_value$. The result is stored

in variable $f1_result$. The invariant states that the variable $f1_result$ contains the result of the computation when $f1_state$ has the value $f1_done$.

```
VARIABLES
  f1 \ state, f1 \ value, f1 \ result
                                               VARIABLES
INVARIANT
                                                 f1\_state, f1\_value, f1\_result
   f1 \ state \in RS \ STATE \land
                                                 RF1 Insts
   f1 value \in V \land
                                               INVARIANT
   f1\_result \in V \land
                                                  RF1\_Insts \subseteq RF1 INSTS \land
                                                  RF1\_null \notin RF1\_\overline{I}nsts \land
   (f1\_state = f1\_done \Rightarrow
   f1 \quad result = f_1(f1\_value))
                                                  f1\_state \in RF1\_Insts \rightarrow RS\_STATE \land
INITIALISATION
                                                  f1\_value \in RF1\_Insts \to V \land
   f1 \quad state := f1 \quad init \parallel
                                                  f1\_result \in RF1\_Insts \to V \land
   f1 value : \in V
                                                  \forall xx.(xx \in RF1\_Insts \Rightarrow
   f1 \quad result :\in V
                                                  (f1 \ state(xx) = f1 \ done \Rightarrow
                                                  f1\_result(xx) = f_1(f1\_value(xx))))
                                               INITIALISATION
                                                  f1\_state := \varnothing \parallel
                                                  RF1_Insts := Ø
```

The remote procedure StartF1 initiates the computation by assigning the variable $f1_value$ to the value supplied by the client and by setting the state, $f1_state$, to $f1_start$. The client can obtain the result of the computation with a call to the automatically generated read-only remote procedure $GetF1_result$.

REMOTE PROCEDURES $StartF1(v_1) \stackrel{\frown}{=}$ **OPERATIONS PRE** $v_1 \in V$ $StartF1(inst, v_1) =$ THEN **PRE** $inst \in RF1$ $Insts \land v_1 \in V$ $f1_value := v_1 ||$ THEN $f1_state := f1_start$ $f1_value(inst := v_1 \parallel$ END $f1_state(inst) := f1_start$ END ; $xx \leftarrow GetF1 \quad state(inst) \triangleq$ **PRE** $inst \in RF1$ Insts **THEN** xx := f1 state(inst) END; $xx \leftarrow GetF1_value(inst) \triangleq \ldots$ $xx \leftarrow GetF1 \quad result(inst) \doteq \ldots$

Grid service machine RF1 contains one event, Comp, that performs the computation of f_1 . The event is only executed once and it disables itself by setting the state to f1 done.

EVENTS $Comp \stackrel{}{=}$ **ANY** inst WHERE Comp =**SELECT** $f1_state = f1_start$ $inst \in RF1$ Insts THEN THEN $f1_result := f_1(f1_value) \parallel$ **SELECT** $f1_state(inst) = f1_start$ $f1_state := f1_done$ THEN END $f1 \ result(inst) :=$ $f_1(f1 \ value(inst)) \parallel$ $f1 \ state(inst) := f1 \ done$ END END ;

Two different notifications can be sent by this grid service machine. The notification *InitNotification* is sent when the grid service instance has been initialised properly. The second notification *DoneNotification*, is sent when the computation has finished, i.e., when event *Comp* has become disabled.

```
NOTIFICATIONS

InitNotification \hat{=}

GUARANTEES

f1\_state = f1\_init

END;

DoneNotification \hat{=}

GUARANTEES

f1\_state = f1\_done \land

f1\_result = f1\_(f1\_value)

END

END
```

The constructor of instances $RF1_GetNew$ and destructor $RF1_Destroy$ are automatically generated during the translation.

```
xx \leftarrow RF1\_GetNew \mathrel{\hat{=}}
    ANY inst WHERE
       inst \in RF1\_INSTS - RF1\_Insts \land
       inst \neq RF1 null \land
       RF1 Insts \neq RF1 INSTS - \{RF1 \ null\}
    THEN
       RF1 \quad Insts := RF1 \quad Insts \cup \{inst\} \parallel
       f1 \quad s\overline{t}ate(inst):=f\overline{1}\_init \parallel
       ANY v_i WHERE v_i \in V
       THEN
           f1\_value(inst) := v_i \parallel
       f1\_result(inst) := v_i
END ||
       xx := inst
   END;
RF1 \quad Destroy(inst) \doteq
    PRE inst \in RF1 INSTS
    THEN
       RF1 \quad Insts := RF1 \quad Insts - \{inst\} \parallel
       f1 \quad state := \{inst\} \triangleleft f1\_state \parallel
       f1\_value := \{inst\} \triangleleft f1\_value \parallel
       f1 \quad result := \{inst\} \triangleleft f1 \quad result \parallel
    END
END
```

The constructor initialises the new instance to the values given in the *initialisations*-clause of the grid service machine and returns it to the client. The destructor removes the instance given as a parameter from the instances in use.

6.3 The grid refinement

The grid refinement machine $COMP_FT1$ is a refinement of the abstract specification COMPF presented in Subsection 6.1. The computation of f(value(1), value(2)) in one step is here refined by the distributed computation $g(f_1(value(1)), f_1(value(2)))$. The function g is modelled with a constant total function with the appropriate properties. Instances of the

grid service RF1 in Subsection 3.2 are used to compute partial results, $f_1(values(1))$ and $f_1(values(2))$, which is then composed into the final result, $g(f_1(values(1)), f_1(values(2)))$. The grid refinement is translated to a B refinement machine for verification. As with the grid service machine the grid refinement is given in the left column and its translation is given to the right throughout this subsection.

In this refinement step we introduce variables $f1_instances$, $f1_results$ and tries.

GRID_REFINEMENT	
$COMP_FT1$	REFINEMENT
REFINES	$COMP FT1_V$
COMPF	REFINES
SEES	COMPF
GDEF	SEES
REFERENCES	GDEF
RF1	INCLUDES
CONSTANTS	$RF1_V$
MaxTimes, g	PROMOTES
PROPERTIES	Comp
$g \in (V \times V) \to V \wedge$	CONSTANTS
$\forall (v_1, v_2). (v_1 \in V \land v_2 \in V \Rightarrow$	MaxTimes
$f(v_1, v_2) = g(f_1(v_1), f_1(v_2))) \land$	PROPERTIES
$MaxTimes \in \mathbb{N}$	$q \in (V \times V) \to V \wedge$
VARIABLES	$\forall (v_1, v_2). (v_1 \in V \land v_2 \in V \Rightarrow$
state, values, result,	$f(v_1, v_2) = g(f_1(v_1), f_1(v_2))) \wedge$
$f1_instances, f1_results, tries$	$MaxTimes \in \mathbb{N}$
INVARIANT	VARIABLES
$f1_instances \in 12 \rightarrow RF1 \land$	state, values, result,
$f1_results \in RF1 \rightarrow V \land$	f1 instances, $f1$ results, tries,
$tries \in \mathbb{N} \land$	RF1 notif
	INVARIANT
	$f1 \ instances \in 12 \rightarrow RF1 \ INSTS \land$
	$ran(f1 instances) \subset$
	$(RF1$ Insts $\cup \{RF1 \ null\}) \land$
	$f1 results \in RF1 INSTS \to V \land$
	$dom(f1 \ results) \subseteq$
	$(RF1 Insts \cup \{RF1 null\}) \land$
	$tries \in \mathbb{N} \land$
	$RF1 notif \in RF1 Insts \to BOOL$

The variable $f1_instances$ is an array containing the two instances that will compute the function values $f_1(values(1))$ and $f_1(values(2))$. The variable $f1_results$ is a partial function from instance to the result obtained from it. The number of times the application has unsuccessfully tried to use instances is modelled by the variable *tries*. The application can fail to use instances *MaxTimes* number of times before it enters state *failed* and aborts the computation. Initially, there has been no tries to use instances, tries = 0.

The instances, $f1_instances$, are initialised to empty instances, RF_null . Since no results have been obtained the variable $f1_results$ is initialised to the empty set.

INITIALISATION	INITIALISATION
$state := idle \parallel$	$state := idle \parallel$
ANY v_1, v_2 WHERE $v_1 \in V \land v_2 \in V$	ANY v_1, v_2 WHERE $v_1 \in V \land v_2 \in V$
THEN $values := \{1 \mapsto v_1, 2 \mapsto v_2\}$ END \parallel	THEN $values := \{1 \mapsto v_1, 2 \mapsto v_2\}$ END \parallel
$result :\in V \parallel$	$result :\in V \parallel$
$f1 instances := (12) \times \{RF1 null\} \parallel$	$f1 instances := (12) \times \{RF1 null\} \parallel$
$f1_results := \varnothing \parallel$	$f1$ results := $\emptyset \parallel$
tries := 0	$tries := 0 \parallel$
	$RF1_notif := \varnothing$

The event GetInstance performs an allocation to obtain a new instance of grid service machine RF1 to use. If a new instance could not be allocated the variable *tries* is incremented.

```
EVENTS
                                                       OPERATIONS
GetInstance =
                                                       GetInstanceOk \mathrel{\hat{=}}
   ANY xxWHERE
                                                          ANY xxWHERE
      xx \in 1..2 \wedge
                                                             xx \in 1..2 \land
      state = working \wedge
                                                              state = working \wedge
      f1 instances(xx) = RF1 null \land
                                                              f1 instances(xx) = RF1 null \land
      tries \leq MaxTimes
                                                             tries \leq MaxTimes
                                                          THEN
   THEN
     CALL f1 instances(xx) \leftarrow RF1 GetNew
                                                           f1 instances(xx) \leftarrow RF1 GetNew;
    EXCEPTION
                                                            RF1\_notif(f1\_instances(xx)) := TRUE
      f1 \quad instances(xx) := RF1\_null \parallel
                                                          END;
      tries := tries + 1
                                                       GetInstanceFail \stackrel{}{=}
    END
                                                          ANY xxWHERE
   END;
                                                             xx \in 1..2 \wedge
                                                              state = working \wedge
                                                              f1 instances(xx) = RF1 null \land
                                                              tries \leq MaxTimes
                                                          THEN
                                                            \begin{array}{l} RF1\_Destroy(f1\_instances(xx));\\ RF1\_notif:=\{f1\_instances(xx)\} \end{array}
                                                              \triangleleft R\overline{F}1 \quad notif;
                                                            f1 \quad instances(xx) := RF1 \quad null \parallel
                                                            tries := tries + 1
                                                          END;
```

Note that in the translation of the exception handler we actually remove $RF1_null$ (the value of $f1_instances(xx)$) from the set of instances in use, $RF1_Insts$. Since $RF1_null$ is not in the set, the destructor will have no effect.

Event Submit starts the computation of f_1 of values values(1) and values(2). Both the instances that will be used for the computation need to be different from $RF1_null$. When a failure of an instance is detected the corresponding variable is assigned the value $RF1_null$ since the instance can no longer be used.

SubmitOk =Submit $\hat{=}$ ANY xx WHERE ANY xx WHERE $xx:1..2\wedge$ $xx:1..2\wedge$ $state = working \wedge$ $state = working \wedge$ $f1_instances(1) \neq RF1 \quad null \land$ f1 instances(1) $\neq RF1$ null \wedge $\begin{array}{l} f1_instances(2) \neq RF1_null \land \\ f1_instances(xx) \notin dom(f1_results) \end{array}$ $\begin{array}{l} f1_instances(2) \neq RF1_null \land \\ f1_instances(xx) \notin dom(f1_results) \end{array}$ THEN THEN **CALL** f1 instances(xx).StartF1(values(xx)) StartF1(f1 instances(xx), values(xx)); EXCEPTION $RF1 \quad notif(f1 \quad instances(xx)) := TRUE$ $f1 \quad instances(xx) := RF1 \quad null \parallel$ END SubmitFail =tries := tries + 1END ANY xx WHERE END ; $xx:1..2\wedge$ $state = workina \wedge$ f1 $instances(1) \neq RF1$ $null \land$ $f1_instances(2) \neq RF1_null \land$ $f1_instances(xx) \notin dom(f1_results)$ THEN $RF1_Destroy(f1_instances(xx));$ $RF1 \quad notif := \{f1 \quad instances(xx)\}$ $\triangleleft R\overline{F}1 \quad notif;$ $f1 \quad instances(xx) := RF1 \quad null \parallel$ tries := tries + 1END;

The event modelling failure, Fail, takes the application to state *failed* when it has unsuccessfully used instances more than *MaxTimes* number of times. Successful execution of the application is modelled by the execution of the event *Finish*. This event is enabled when the result from both instances has been obtained via the execution of the notification handlers.

```
Fail =
                                                 Fail =
   \textbf{SELECT} \ state = working \wedge
                                                    \textbf{SELECT} \ state = working \wedge
    tries > MaxTimes
                                                     tries > MaxTimes
   THEN state := failed
                                                    THEN state := failed
  END ;
                                                    END;
Finish =
                                                 Finish =
  SELECT
                                                    SELECT
     state = working \land
                                                       state = working \land
     dom(f1\_results) = ran(f1\_instances)
                                                       dom(f1 \ results) = ran(f1 \ instances)
                                                    THEN
   THEN
     result := g(f1 \ results(f1 \ instances(1))),
                                                       result := g(f1 \ results(f1 \ instances(1))),
       f1\_results(f1\_instances(2))) \parallel
                                                        f1\_results(f1\_instances(2))) \parallel
                                                       state := done
      state := done
   END
                                                    END;
```

If a notification is received successfully the result in the corresponding instance is obtained with a call to the automatically generated remote procedure $GetF1_result$. When a timeout occurs the application stops using the instance. Since $GetF1_result$ is a remote procedure the translation is performed in two steps. The first step splits the notification handler into two events, one for successful handling the notification and the other for time-out handling. The event for successful handling of notifications is then split into two parts to consider both successful execution of the remote procedure $GetF1_result$ as well as the failed execution.

```
NOTIFICATION HANDLERS
F1DoneNotifHandler =
   NOTIFICATION DoneNotif
   SOURCE inst \in RF1
    THEN
       IF inst \in ran(f1 \ instances)
       THEN
         VAR xx, ok IN
           ok := TRUE; xx :: VV;
           \textbf{CALL} \ xx \gets inst.GetF1\_result
           EXCEPTION
               ok := FALSE
               f1 \quad instances := f1 \quad instances \Leftrightarrow
                 \{\overline{i}, v | i \in 1..2 \land
               \label{eq:constraint} \begin{array}{l} inst = f1\_instances(i) \land \\ v = RF1\_null\} \parallel \\ f1\_results := \{inst\} \end{array}
                 \triangleleft f1 results ||
               tries := tries + 1
           END :
           \mathbf{IF} \ ok = TRUE
           THEN f1 \ results(inst) := xx
           END
       END
   TIMEOUT
       f1\_instances := f1\_instances \Leftrightarrow
         \{\overline{i,v}|i\in 1..2 \land inst = f1\_instances(i) \land
         v = RF1\_null\} \parallel
       f1 \quad results := \{inst\}
         \triangleleft f1 \quad results \parallel
       tries := tries + 1
   END
```

 $F1DoneNotifHandlerOkOk \stackrel{}{=}$ ANY inst WHERE $\begin{array}{l} inst \in RF1_Insts \land \\ RF1_notif(inst) = TRUE \land \end{array}$ $\neg(f1_state(inst) = f1_start) \land$ $\begin{array}{l} (f1_state(inst) = f1_done \land \\ f1_result(inst) = f_1(f1_value(inst))) \end{array}$ THEN $RF1_notif(inst) := FALSE;$ **IF** $inst \in ran(f1_instances)$ THEN $\mathbf{VAR} \ xx, ok \ \mathbf{IN}$ ok := TRUE; xx :: V; $xx \leftarrow GetF1_result(inst);$ $\mathbf{IF} \ ok = TR\overline{U}E$ **THEN** $f1 \quad results(inst) := xx$ END END END END; F1DoneNotifHandlerOkFail =ANY inst WHERE $inst \in RF1 \quad Insts \land$ $RF1_notif(inst) = TRUE \land$ $\neg (f1_state(inst) = f1_start) \land (f1_state(inst) = f1_done \land$ $f1_result(inst) = f_1(f1_value(inst)))$ THEN $RF1 \quad notif(inst) := FALSE;$ **IF** $inst \in ran(f1_instances)$ THEN **VAR** xx, ok **IN** ok := TRUE; xx :: V; $RF1_Destroy(f1_instances(xx));$ $\begin{array}{l} RF1_notif := \{f1_instances(xx)\} \\ \triangleleft RF1_notif; \end{array}$ ok := FALSE; $f1 \ instances := f1 \ instances \Leftrightarrow$ $\{i,v|i\in 1..2 \wedge$ $inst = f1_instances(i) \land$ $v = RF1_null\} \parallel$ $f1 \quad results := \{inst\}$ $\triangleleft \overline{f}1_results \parallel$ tries := tries + 1END; **IF** ok = TRUE $\textbf{THEN} \ f1_results(inst) := xx$ END END END; F1DoneNotifHandlerFail =ANY inst WHERE $inst \in RF1 \ Insts \wedge$ $RF1_notif(inst) = TRUE$ THEN RF1 Destroy(inst); $RF1_notif := \{inst\} \triangleleft RF1_notif;$ $f1_instances := f1_instances \Leftrightarrow$ $\{\overline{i,v}|i\in 1..2 \land inst = f1_instances(i) \land$ $v = RF1_null\} \parallel$ $f1 \quad results := \{inst\}$ $\triangleleft f1_results \parallel$ tries := tries + 1END;

The variable giving the instance that failed is also here assigned the value $RF1_null$, the pending notification is removed and variable *tries* is incremented. Note that we do not introduce an assignment $RF1_notif(inst) := TRUE$ after the call to $GetF1_result$ since it is a automatically defined read-only remote procedure.

Successful *is-alive checks* are not modelled in Distributed B since they do not affect the behaviour of the application. The handling of failed *is-alive checks* is similar to handling notification timeouts. Like for notification timeouts the application stops using the instance that failed.

```
IS ALIVE HANDLERS
                                                              IAH and ler =
IA\overline{H}andler \stackrel{\frown}{=}
                                                                ANY inst WHERE
 SOURCE inst \in RF1
                                                                 inst \in RF1 Insts
  THEN
                                                                THEN
                                                                 RF1\_Destroy(inst);
   f1\_instances := f1\_instances \Leftrightarrow
     \{i, v | i \in 1..2 \land inst = f1 \ instances(i) \land
                                                                 RF1 \quad notif := \{inst\} \ll |RF1 \quad notif;
     v = RF1 \quad null \} \parallel
                                                                 f1 \quad instances := f1 \quad instances \Leftarrow
   f1 \quad results := \{inst\}
                                                                   \{i,v|i\in 1..2 \land inst = f1\_instances(i) \land
     \triangleleft f1 \quad results \parallel
                                                                   v = RF1 \quad null \} \parallel
   tries := tries + 1
                                                                  f1 \quad results := \{inst\}
 END
                                                                   {\triangleleft}f1\_results\parallel
END
                                                                 tries := tries + 1
                                                                END
                                                              END
```

We have now decomposed the computation of f(values(1), values(2))into a fault tolerant grid application. The grid service machine RF1 and the grid refinement $COMP_FT1$ can now be further refined using grid refinement until executable code can be generated. The different components in the development can be proved correct using the tool support for B.

7 Implementation in Java

The grid application development in fault tolerant Distributed B continues until all the non-determinism has been removed and all the used constructs can be implemented, i.e., they belong to the implementable subset of the B language, B0. When all substitutions of the system belong to the B0 language, they can be automatically translated to Java [15]. The event system composed of all events in the *events*-clause can be automatically translated to a *while*-loop in Java, if all the guards Gi and the substitutions Si belong to the B0 language.

	while (true) {
$E1 \stackrel{\circ}{=}$	if(G1) S1;
WHEN G1 THEN S1 END	
I HEIN SI EIND	else if(Gn) Sn;
$E_n \hat{=}$	· · · ,
2.0	else <stop loop="">;</stop>
WHEN Gn	}
THEN Sn END	

The place holder *<stop loop* > consists of statements that terminate the loop

in the main program and statements for sending notifications and pausing the execution in implementations of grid service machines.

In the client the grid service instances are translated to objects encapsulating the grid specific features. These features include instance allocation, remote procedure calls, notification handling, notification timeouts and management of *is-alive checks*. In the implementation of grid service machines an additional layer encapsulating grid features such as state management and *is-alive timeouts* is inserted between the Globus toolkit and the translated grid service machine. The code for the objects and the additional layer need to be manually created once, but can then be reused in all Distributed B applications.

Java code is generated for each concrete grid refinement machine that is a refinement of a grid service machine and for the root grid refinement machine that is a refinement of an Event B model. When the code has been generated we have implemented the grid application in a formal manner, where we can show that the implementation is a correct refinement of the abstract specification.

8 Conclusions

We have earlier developed a language Distributed B [5] that extends Event B for designing and implementing correct grid applications. In that paper we introduced two new types of machines, grid service machine and grid refinement machine, for handling grid specific issues in Event B. In this paper we have modified Distributed B for developing fault tolerant and robust grid applications. We introduced exception handling to take care of exceptions raised due to failed remote procedure calls, as well as timeouts to discover lost notifications. In order to handle orphan grid service instances we introduced a special checking mechanism and facilities to handle exceptions raised by it. These new features force the developer to consider the fault tolerance of grid applications throughout the development and enables formal proofs of correctness. Hence, we have proposed a method for implementing fault tolerant grid applications where the implementation can be proved correct with respect to its specification.

There are several middlewares comparable to Globus toolkit that support advanced fault tolerance mechanisms. For example, fault tolerance in CORBA [10, 14] is based on replication. In CORBA replicas of an object form an object group. Object groups then provide replication transparency and failure transparency. Replication is complementary to the mechanisms presented here and both can be used together for developing very reliable grid applications. Fault detection is an important part of a fault tolerant application. A service for detecting faulty components in a grid environment using *unreliable fault detectors* is presented in a paper by Stelling et al. [13]. However, the service is no longer part of the Globus toolkit due to a number of deficiencies, such as excessive resource usage and difficulties in supporting From a performance perspective our fault tolerant language for developing grid applications could be improved. One of the most challenging problems when constructing fault tolerant grid systems is to handle the orphan grid service instances. Our approach uses timers and extra remote procedure calls, which might not be optimal. The handling of orphans could be improved as discussed in detail by Panzieri and Shrivastava [12] and by Tanenbaum and Steen [14]. Furthermore, the fault tolerance mechanism will immediately delete references to grid service instances that experiences problems. This can lead to the deletion of an entire subtree of instances for one failure. Hence, more efficient mechanisms that would maintain the partial results should be investigated.

The language we proposed in this paper provides a convenient formal development process for fault tolerant grid applications. The applications will by construction have an architecture that is fault tolerant and implementable. Furthermore, the applications are modelled in terms of grid primitives with a precise meaning and the specifications of them will therefore be clear to understand. Our approach to adapt Event B to the Globus Toolkit middleware is not limited to that specific middleware, but it can be applied to other middlewares for distributed systems as well.

References

- J. R. Abrial. The B-Book: Assigning Programs to Meanings. Cambridge University Press, 1996.
- [2] J. R. Abrial, D. Cansell and D. Méry. Refinement and Reachability in Event B. In H. Treharne et al, editors, proceedings of the 4th international conference of Z and B users: ZB2005. LNCS 3455, Guildford, UK, pp. 144-163, Springer-Verlag, 2005.
- [3] J. R. Abrial and L. Mussat. Event B Reference Manual, 2001. http://www.atelierb.societe.com/ressources/evt2b/ eventb_reference_manual.pdf. (accessed 10.08.2005)
- [4] R. J. R. Back and R. Kurki-Suonio. Decentralization of process nets with centralized control. In *Proceedings of the 2nd ACM SIGACT-SIGOPS* Symposium of Principles of Distributed Computing, pp. 131-142, 1983.
- [5] P. Boström and M. Walden. An extension of Event B for developing grid systems. In H. Treharne et al, editors, proceedings of the 4th international conference of Z and B users: ZB2005. LNCS 3455, Guildford, UK, pp. 144-163, Springer-Verlag, 2005.
- [6] K. Czajkowski, et. al. Open Grid Services Infrastructure, 2003. http://www-unix.globus.org/toolkit/ draft-ggf-ogsi-gridservice-33_2003-06-27.pdf. (accessed 10.08.2005)

it.

- [7] E. W. Dijkstra. A Discipline of Programming. Prentice-Hall International, 1976.
- [8] I. Foster, C. Kesselman and S. Tuecke. The Anatomy of the Grid: Enabling Scalable Virtual Organizations. The International Journal of Supercomputer Applications, 15(3), 2001.
- [9] I. Foster, C. Kesselman, J. Nick and S. Tuecke. The Physiology of the Grid: An Open Grid Services Architecture for Distributed Systems Integration. Open Grid Service Infrastructure WG, Global Grid Forum, 2002. http://www.globus.org/alliance/publications/papers/ogsa.pdf. (accessed 10.08.2005)
- [10] Object Management Group. Fault tolerant CORBA, 2001, http://www.omg.org/docs/formal/01-09-29.pdf. (accessed 10.08.2005)
- [11] The Globus Alliance. Globus Toolkit. 2005. http://www.globus.org/. (accessed 10.08.2005)
- [12] F. Panzieri and S. K. Shrivastava. Rajdoot: A Remote Procedure Call Mechanism Supporting Orphan Detection and Killing. *IEEE Transac*tions on Software Engineering, 14(1), pp 30-37, 1988.
- [13] P. Stelling, C. DeMatteis, I. Foster, C. Kesselman, C. Lee, G. von Laszewski. A Fault Detection Service for Wide Area Distributed Computations, *Cluster Computing*, 2, pp. 117-128, 1999
- [14] A. S. Tanenbaum and M. Van Steen. Distributed systems principles and paradigms. Prentice Hall. 2002
- [15] J. C. Voisinet, B. Tatibouet and A. Hammand. JBTools: An experimental platform for the formal B method. In Proceedings of the inaugural conference on the Principles and Practice of programming and Proceedings of the second workshop on Intermediate representation engineering for virtual machines. National University of Ireland, 2002
- [16] M. Waldén and K. Sere. Reasoning About Action Systems Using the B-Method. Formal Methods in Systems Design, 13:5-35, 1998.

A A fault tolerant grid refinement machine

GRID REFINEMENT C_{FT} REFINES \mathcal{C} REFERENCES \mathcal{A} VARIABLES v, w, a_1, \ldots, a_n INVARIANT $a_1 \in \mathcal{A} \wedge$ $\ldots a_n \in \mathcal{A} \land$ $J(v, w, a_1, \ldots, a_n)$ INITIALISATION $a_1 := A_null \parallel$ Init'(v, w)EVENTS $E_1 =$ ANY y WHERE $G_1'(v,y)$ THEN $S'_{1}(v, y);$ **CALL** $a_i \leftarrow A_GetNew$ **EXCEPTION** T_e END \mathbf{END} ; $E_2 =$ WHEN $G'_2(v, w)$ THEN $S'_2(v, w)$ END; $F \stackrel{}{=}$ **WHEN** $H(v, w, a_i)$ THEN $S_n(w) \parallel$ **CALL** a_i . Proc(f(v, w))**EXCEPTION** $T_f(w, a_i)$ END END NOTIFICATION HANDLERS $N_1Handler =$ **NOTIFICATION** N_1 **SOURCE** $inst \in A$ **THEN** $T_{N1}(inst, w, a_i, a_j)$ **TIMEOUT** $T_t(inst, w, a_i, a_j)$ END IS ALIVE_HANDLERS $IAHandler \stackrel{\frown}{=}$ **SOURCE** *inst* $\in A$ **THEN** $T_{ia}(inst, w, a_i, a_j)$ END

REFINEMENT C_{FTV} REFINES \mathcal{C} INCLUDES \mathcal{A}_V PROMOTES E_A VARIABLES $v, w, a_1, \ldots, a_n, A_notif$ INVARIANT $a_1 \in A_INSTS \land$ $a_1 \in A_Insts \cup \{A_null\} \land$ $J_V(v,w,a_1,\ldots,a_n)\wedge$ $A_notif \in A_Insts \to BOOL$ INITIALISATION $a_1 := A \quad null \parallel$ $Init'(v,w)\parallel$ $A \quad notif := \emptyset$ OPERATIONS $E_1Ok =$ ANY y WHERE $G_1'(v,y)$ THEN $S_1'(v,y);$ $a_i \leftarrow A$ GetNew; $A_notif(a_i) := TRUE$ END; $E_1Fail = \ldots$ $E_2 \stackrel{.}{=} \ldots$ $FOk \stackrel{}{=}$ **WHEN** $H(v, w, a_i)$ THEN $S_n(w) \parallel$ $Proc(a_i, f(v, w));$ A $notif(a_i) := TRUE$ \mathbf{END} ; FFail = ... $N_1HandlerOk =$ **ANY** inst **WHERE** $inst \in A$ $Insts \land$ A $notif(inst) = TRUE \land$ $\neg G_A(x(inst)) \land Q_1(x(inst))$ THEN $T_{N1}(inst,w,a_i,a_j)\parallel$ $A \quad notif(inst) := FALSE$ END; $N_1HandlerFail = \ldots$ $IAH and ler \stackrel{}{=}$ **ANY** inst **WHERE** $inst \in A_Insts$ THEN $A_Destroy(inst);$ $A_notif := \{inst\} << |A_notif;$ $T_{ia}(inst, w, a_i, a_j)$ END END

END





University of Turku

- Department of Information Technology
- Department of Mathematics



Åbo Akademi University

- Department of Computer Science
- Institute for Advanced Management Systems Research



Turku School of Economics and Business AdministrationInstitute of Information Systems Sciences

ISBN 952-12-1594-1 ISSN 1239-1891