



SERIALIZATION OF DISTRIBUTED THREADS IN JAVA

DANNY WEYNS, EDDY TRUYEN, PIERRE VERBAETEN*

Abstract. In this paper we present a mechanism for serializing the execution-state of a distributed Java application that is implemented on a conventional Object Request Broker (ORB) architecture such as Java Remote Method Invocation (RMI). To support serialization of distributed execution-state, we developed a byte code transformer and associated management subsystem that adds this functionality to a Java application by extracting execution-state from the application code. An important benefit of our mechanism is its portability. It can transparently be integrated into any legacy Java application. Furthermore, it does require no modifications to the Java Virtual Machine (JVM) or to the underlying ORB. Our serialization mechanism can serve many purposes such as migrating execution-state over the network or storing it on disk. In particular, we describe the implementation of a prototype for repartitioning distributed Java applications at run-time. Proper partitioning of distributed objects over the different machines is critical to the global performance of the distributed application. Methods for partitioning exist, and employ a graph-based model of the application being partitioned. Our mechanism enables then applying these methods at any point in an ongoing distributed computation. In the implementation of the management subsystem, we experienced the problem of losing logical thread identity when the distributed control flow crosses address space boundaries. We solved this well known problem by introducing the generic notion of distributed thread identity in Java programming. Propagation of a globally unique, distributed thread identity provides a uniform mechanism by which all the program's constituent objects involved in a distributed control flow can uniquely refer to that distributed thread as one and the same computational entity.

Key words. serialization of execution-state, distributed threads, Java

1. Introduction. In this paper we present a mechanism for serializing the execution-state of a distributed Java application. We describe this mechanism in the context of a system for run-time repartitioning of distributed Java applications. For distributed object-oriented applications, an important management aspect is the partitioning of objects such that workload is equally spread over the available machines and network communication is minimized. Traditional techniques for automatic partitioning of distributed object applications uses graph-based algorithms (e.g. [8]).

In a static approach an external monitor automatically determines the best possible partitioning of the application, based on observation of behavior of the application (i.e., the dispersal of costs) during a number of representative runs. This partitioning is fixed for the entire execution of the application. However in a dynamic environment the optimal object distribution may change during execution of the application. To cope with this, the external monitor may periodically check the workload at run-time on each separate machine. Whenever the workload on one or more machines crosses a certain threshold, e.g., following the low-water high-water workload model as described in [12], the monitor immediately triggers the repartitioning algorithm and relocates one or more objects to another machine.

The relocation of a running object involves the migration of its object code, data state and execution-state. Conventional Java-based Object Request Brokers (ORBs), such as the Voyager ORB [9], support passive object migration (migration of object code and data state, but no migration of execution-state). However, run-time repartitioning does not want to wait with object relocation until that object and eventually all objects involved in the execution of that object are passive. Instead it aims to handle the triggers for object repartitioning immediately. As a consequence existing methods for repartitioning must be adapted to be applied at any point in an ongoing distributed computation. As such, it is necessary to support object relocation with migration of execution-state. Migration of execution-state is in the literature often referred to as strong thread migration [7]. The fact that the Voyager ORB does not support strong thread migration is not just a missing feature, but the real problem is that migration of execution-state is simply not supported by current Java technology.

To solve this we developed a byte code transformer and associated management subsystem that enables an external control instance (such as the above load balancing monitor) to capture and reestablish the execution-state of a running distributed application. We call this serialization of distributed execution-state. The byte code transformer instruments the application code by inserting code blocks that extract the execution-state from the application code. The management subsystem, which is invoked by the inserted codes, is responsible for managing the captured execution-state efficiently. The management subsystem also provides operations by which an external control instance can initiate serialization of the distributed execution-state at its own will.

It is important to know that we solely focus on distributed applications that are developed using conventional

*Department of Computer Science, DistriNet, Katholieke Universiteit Leuven, Belgium
(email: danny-eddy-pv@cs.kuleuven.ac.be; web: www.cs.kuleuven.ac.be/~danny/DistributedBRAKES.html)

ORBs such as Java Remote Method Invocation (RMI) or Voyager. Programmers often use these middleware platforms because of their object-based Remote Procedure Call (RPC) programming model, which is very similar to the well-known object-oriented programming style.

1.1. Important aspects of our work. In this paper we first describe how we realized **serialization of a distributed execution-state**. Note that in the past, several algorithms have been proposed to capture the execution state of Java Virtual Machine (JVM) threads. Some require the modification of the JVM [3]. Others are based on the modification of source code [7]. Some algorithms rely on byte code rewrite schemes, e.g. [11, 17]. We too had already implemented such a byte code rewrite algorithm called Brakes [15]. However, most of these schemes are presented in the domain of mobile agents systems. Whenever a mobile agent wants to migrate, it initiates the capturing of its own execution-state. As soon as the execution-state is serialized the agent migrates with its serialized execution-state to the target host where execution is resumed. However, serialization of distributed execution-state of Java RMI-like applications introduces two aspects that are not covered in this existing thread capturing systems. First, contrary to how conventional Java threads are confined to a single address space, the computational entities in Java RMI applications execute as distributed flows of control that may cross physical JVM boundaries. In the remainder of this paper we call such distributed flows of control *distributed threads* [4]. Serializing the execution-state of a distributed application boils down to capturing all distributed threads that execute in that application. The second aspect is that mobile agents initiate the capturing/reestablishment of their execution-state themselves, whereas capturing/reestablishment of distributed execution-state must often be initiated by an external control instance. The Brakes thread serialization scheme is not designed for being initiated by such an external control instance. In this paper we describe how we have extended Brakes with (1) a mechanism for serialization of the execution-state of distributed threads that (2) can be initiated by an external control instance.

Subsequently, we show how we used this serialization mechanism to implement a prototype for **run-time repartitioning of distributed Java applications**. The idea is that a load balancing monitor, that plays the role of external control instance here, captures the execution-state of an application whenever it wants to repartitioning that application and reestablishes the execution-state after the repartitioning is finished. The advantage of having separate phases for migration of execution-state and object migration is that objects can migrate independently of their (suspended) activity. Requests for object migration can immediately be performed, without having to wait for the objects to become passive. This is possible because, by using the serialization mechanism, application objects can be turned passive on demand by the monitor, while their actual execution-state is safely stored in the management subsystem of the serialization mechanism. So when the actual repartitioning takes place, all application objects are a priori passive. As a result, we can still use conventional passive object migration support to implement the run-time repartitioning prototype. In this paper we assume that the application's underlying ORB supports passive object migration (e.g., the Voyager ORB), but existing work [6] has shown that support for passive object migration can also be added to the application by means of a byte code transformation.

Previous work [10][16] already offers support for run-time repartitioning, but this is implemented in the form of a new middleware platform with a dedicated execution model and programming model. A disadvantage of this approach is that Java RMI legacy applications, which have obviously not been developed with support for run-time repartitioning in mind, must partially be rewritten such that they become compatible with the programming model of the new middleware platform. Instead, our new approach is to develop a byte code transformer that transparently injects new functionality to an existing distributed Java application such that this application becomes automatically run-time repartitionable by the monitor. The motivation behind this approach taken is that programmers do not want to distort their applications to match the programming model of whatever new middleware platform.

Finally, an important benefit of our mechanism for capturing distributed threads is its **portability**: (1) byte code transformations integrate the required functionality transparently into existing Java applications. A custom class loader can automatically perform the byte code transformation at load-time. (2) Our serialization mechanism does require no modifications of the JVM. This makes the implementation portable on any system, as long as a standard JVM is installed on that system. (3) Our mechanism does require no modifications of the underlying ORB. Our mechanism works seamless on top of any ORB with an RPC-like programming model, provided that our byte code transformation is performed before stub code generation.

However, a limitation is that our serialization mechanism is only applicable on top of a dedicated cluster of machines where network latencies are low and faults are rare. This is not directly a dependability of our approach, but rather a dependability of the RPC-like programming model: performing blocking calls on remote objects is after all only feasible on a reliable, high-bandwidth and secure network. As such our serialization mechanism is not well suited for Internet or wireless applications.

1.2. Structure of the paper. This paper is structured as follows: first, in section 2 we discuss the problems with using Brakes for capturing distributed threads. In section 3 we introduce the notion of distributed thread identity as a concept to handle a distributed control flow as one computational entity. In section 4 we apply the concept of distributed thread identity to easier reuse Brakes for serialization of distributed threads. We further explain how serialization of execution-state can be initiated by an external control instance. In section 5 we introduce our prototype for run-time repartitioning and demonstrate how it works by means of a concrete example application. In section 6 we evaluate the performance overhead and byte code blowup that is generated by our approach. Section 7 discusses related work. Finally we conclude and look to future work in section 8.

2. Problem statement. In this section we first shortly describe the implementation of Brakes, our earlier developed mechanism for serialization of JVM threads. Then we discuss the problem we encountered when trying to reuse Brakes for capturing distributed execution-state.

2.1. Brakes for JVM serialization. In Brakes the execution-state of a JVM thread is extracted from the application code that is executing in that thread. For this, a byte code transformer inserts capture and reestablishing code blocks at specific positions in the application code. We will refer to this transformer as the Brakes transformer.

With each thread two flags (called `isSwitching` and `isRestoring`) are associated that represent the execution mode of that specific thread. When the `isSwitching` flag is on, the thread is in the process of capturing its state. Likewise, a thread is in the process of reestablishing its state when its `isRestoring` flag is on. When both flags are off, the thread is in normal execution. Each thread is associated with a separate Context object into which its execution-state is switched during capturing, and from which its execution-state is restored during reestablishing.

The process of capturing a thread's state (indicated by the empty-headed arrows in Fig. 2.1) is then implemented by tracking back the control flow, i.e., the sequence of nested method invocations that are on the stack of that thread. For this the byte code transformer inserts after every method invocation a code block that switches the stack frame of the current method into the context and returns control to the previous method on the stack, etc. This code block is only executed when the `isSwitching` flag is set. The process of reestablishing a

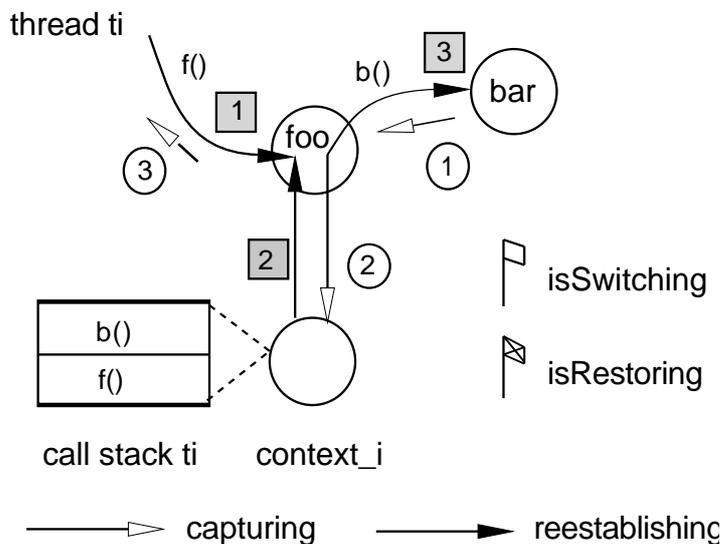


FIG. 2.1. Thread Capturing/Reestablishing in Brakes.

thread's state (indicated by the full-headed arrows in Fig. 2.1) is similar but restores the stack frames in reverse order on the stack. For this, the byte code transformer inserts in the beginning of each method definition a code block that restores stack frame data of the current method and subsequently creates a new stack frame for the next method that was on the stack, etc. This code block is only executed when the `isRestoring` flag is set.

A context manager per JVM manages both Context objects and flags. The inserted byte codes switch/restore the state of the current thread into/from its context via a context-manager-defined static interface. The context manager manages context objects on a per JVM thread basis. So every thread has its own Context object, exclusively used for switching the state of that thread. The context manager looks up the right context object with the thread identity as hashing key. For more information about Brakes, we refer the reader to [15].

2.2. Problem with Brakes to capture distributed execution-state. This section describes the problem that we encountered when trying to reuse Brakes for capturing distributed execution-state. In Brakes, execution-state is saved per local JVM thread. This works well for capturing local control flow but not for capturing a control flow that crosses system boundaries. Fig. 2.2 illustrates the problem. Once thread t_i in

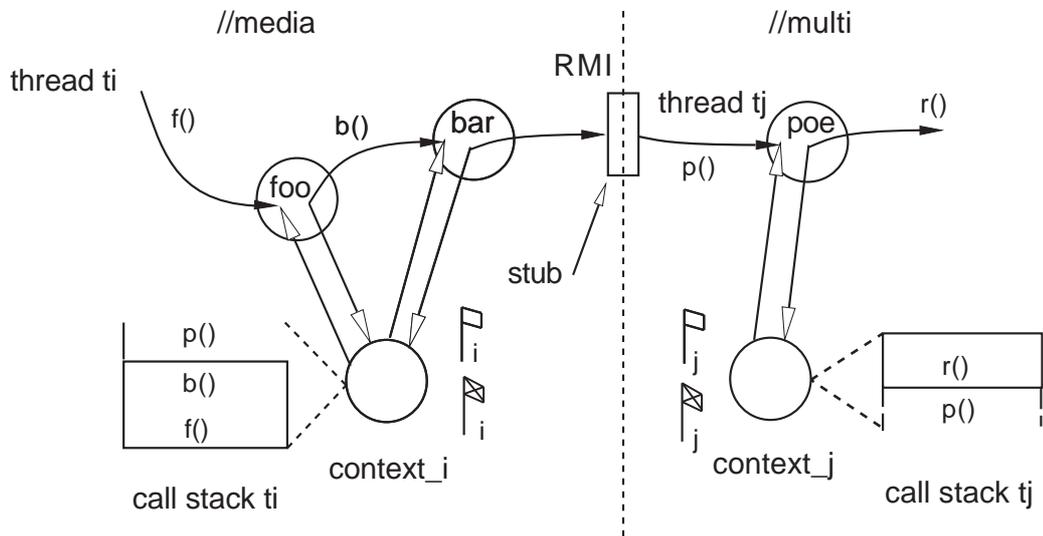


FIG. 2.2. Context per JVM Thread.

the example reaches method `b()` on object `bar`, the call `p()` on object `poe` is performed as a remote method invocation. This remote call implicitly starts a new thread t_j at host `media`. Physically, the threads t_i and t_j hold their own local subset of stack frames, but logically the total set of frames belongs to the same distributed control flow. The context manager however, is not aware of this logical connection between threads t_i and t_j . As a consequence Brakes will manage contexts and flags of these JVM threads as separate computational entities, although they should be logically connected. Without this logical connection, it becomes difficult to robustly capture and reestablish a distributed control flow as a whole entity. For example, it becomes quasi impossible for the context manager to determine the correct sequence of contexts that must be restored for reestablishment of a specific distributed control flow.

3. Distributed threads and distributed thread identity. The above discussed problem is a specific instance of a more general problem. The essence of this problem is that we can not rely on JVM thread identity to refer to a distributed thread as one and the same computational entity. We now discuss how we have extended Java programming to cope with this problem.

A Java program is executed by means of a Java Virtual Machine thread (JVM thread). Such a thread is the unit of computation. It is a sequential flow of control within a single address space (i.e. JVM). However, distributed threads [4] execute as flows of control that may cross physical node boundaries. A distributed thread is a logical sequential flow of control that may span several address spaces (i.e. JVMs). As shown in Fig. 3.1 a distributed thread T is physically implemented as a concatenation of local (per JVM) threads $[t_1, \dots, t_4]$ sequentially performing remote method invocations when they transit JVM boundaries. In the context of

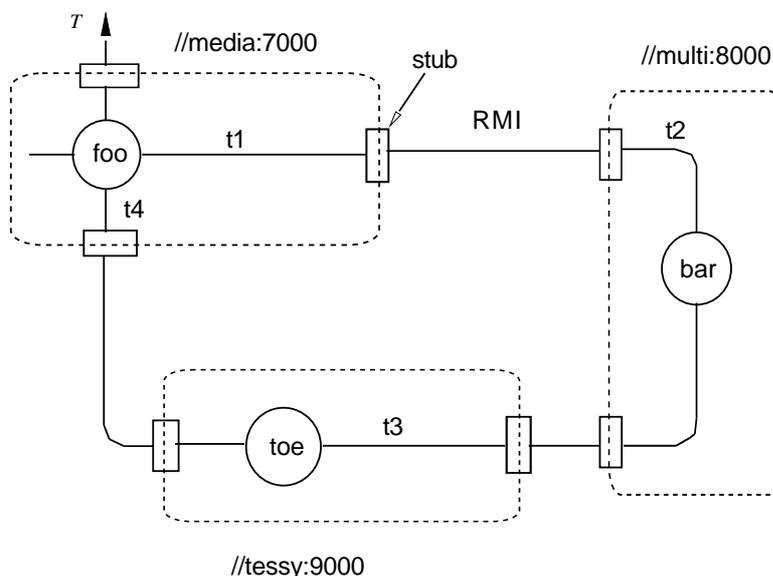


FIG. 3.1. A Distributed Thread.

a distributed control flow programming model, distributed threads offers a general concept for a distributed computational entity.

In a local execution environment, i.e., for centralized programs that run on one JVM, the JVM thread identifier offers a unique reference for a single computation entity. For a distributed application however distributed threads execute as flows of control that may cross physical node boundaries. Once the control flow crosses system boundaries a new JVM thread is used for continuing execution and so logical thread identity is lost. As a consequence JVM thread identifiers are too fine grained for identification of distributed threads. We extend Java programs with the notion of distributed thread identity. Propagation of a globally unique distributed thread identity provides a uniform mechanism by which all the program's constituent objects involved in a distributed thread can uniquely refer to that distributed thread as one and the same computational entity.

With respect to the subject of this paper, the implementation of distributed threads should enforce the following two rules:

1. A distributed thread identity must be created once, more specifically when the distributed thread is scheduled for the first time.
2. A distributed thread identity must not be modified after it is created. This is because it must provide physically dispersed objects with a unique and immutable reference to a distributed thread.

We used M. Dahm's tool for byte code transformation, BCEL [5], to develop a byte code transformer that extends Java programs with the notion of distributed threads, more specifically distributed thread identity. Hereafter we will refer to this transformer as the DTI transformer.

To achieve propagation of distributed thread identities, the DTI transformer extends the signature of each method with an additional argument of class `D_Thread_ID`. `D_Thread_ID` is a serializable class that implements an immutable, globally unique identifier. The signature of every method invoked in the body of the methods must be extended with the same `D_Thread_ID` argument type too. For example, a method `f()` of a class `C` is rewritten as:

```

//original method code      //transformed method code
f(int i, Bar bar) {        f(int i, Bar bar, D_Thread_ID id) {
    ...                      ...
    bar.b(i);                bar.b(i, id);
    ...                      ...
}                             }

```

When `f()` is called, the `D_Thread_ID` is passed as an actual parameter to `f()`. Inside the body of `f()`,

`b()` is invoked on `bar`. On its turn, the body of `f()` passes the `D_Thread_ID` it received as an extra argument to `b()`. This way the distributed thread identity is automatically propagated with the control flow along the method call graph.

3.1. Creation and modification of a distributed thread. The Java thread programming model offers the application programmer the possibility to start up a new thread from within the `run()` method of an object of a class that implements the `java.lang.Runnable` interface. The DTI transformer instrument this kind of classes such that they implement the `D_Runnable` interface instead. `D_Runnable` is defined as follows:

```
interface D_Runnable {
    void run(D_Thread_ID id);
}
```

The example class `Bar` that originally implements the `Runnable` interface illustrates this transformation:

```
//original class definition    //transformed class definition
class Bar implements          class Bar implements
    java.lang.Runnable {      D_Runnable {
    ...                          ...
    void run() {...}           void run(D_Thread_ID id) {...}
    }                            }
```

As stated in section 3, the identity of a distributed thread must be created at the moment the distributed thread is created. This behavior is encapsulated in the `D_Thread` class, which wraps each `D_Runnable` object in a `D_Thread` and as such serves as an abstraction for creating a new distributed thread. A new distributed thread can simply be started with:

```
Bar b = new Bar();
D_Thread dt = new D_Thread(b);
dt.start();
```

The `D_Thread` class itself is defined as:

```
class D_Thread implements java.lang.Runnable {
    public static D_Thread_ID getCurrentThreadID() {
    }
    public D_Thread(D_Runnable o) {
        object = o;
        id = new D_Thread_ID();
    }
    public void run() {
        object.run(id);
    }
    public void start() {
        new Thread(this).start();
    }
    private D_Runnable object;
    private D_Thread_ID id;
}
```

As stated in section 3, distributed thread identities may never be modified. As such `D_Thread` does not provide any method operations for this. Furthermore, `D_Thread_ID` objects are stored either as a private instance member of class `D_Thread` or as a local variable on a JVM stack. Nonetheless it remains possible for

a malicious person to modify distributed thread identities, e.g., by inserting malicious byte code that modifies the value of the local variable pointing to a `D_Thread_ID`. To prevent this, additional measures are necessary, which are subject of future work.

3.2. Inspection of distributed thread identity. Since distributed thread identities are propagated with method invocations as an additional last argument, it is possible to compute for every method definition which local variable on the JVM stack points to the corresponding `D_Thread_ID` object. This allows a management subsystem to inspect the value of the local `D_Thread_ID` variable at any point in the method code. This inspection requires two parts. First the management subsystem must define a static interface that includes the `D_Thread_ID` as an argument of its methods. Second the application code must be extended with byte codes at specific marking points in its program that invoke these methods. Then, each time a distributed thread passes such marking point, it notifies the management subsystem, identifying itself by means of the passed `D_Thread_ID` argument.

Our implementation of distributed thread identity also allows application objects to inspect `D_Thread_ID`, by means of a static operation `getCurrentThreadID()`. The implementation of this method is encapsulated in the DTI transformer, which transparently replaces any invocation of this method with the value of the `D_Thread_ID` variable.

3.3. Integration with the ORB architecture. Applying the byte code transformation to applications developed with an off-the-shelf ORB demands some attention. The programmer must be aware of generating the stub classes for the different remote interfaces only after byte code transformation has been applied. This to make sure that stubs and skeletons would propagate distributed thread identity appropriately.

4. Distributed thread serialization. Distributed thread identity provides an elegant solution to the problem of serializing distributed threads. More specifically if we augment the Brakes transformer with functionality for distributed thread identity inspection (i.e. `D_Thread_ID`) by the context manager via a static interface (see section 2.1), it is possible to implement a management subsystem that manages contexts on a per distributed thread basis (see Fig. 4.1). This update to the Brakes transformer entails only slight modification.

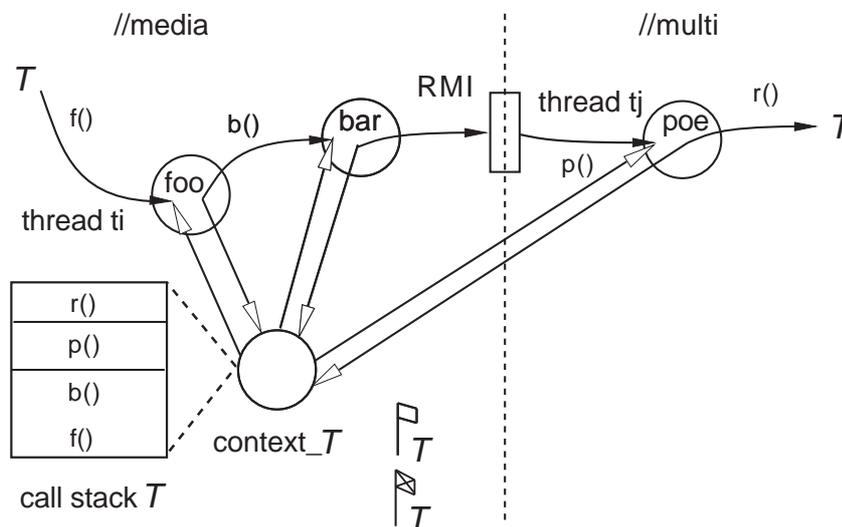


FIG. 4.1. Context per Distributed Thread.

We now discuss these extensions.

4.1. New static interface for inspecting distributed thread identity. The Brakes byte code transformer must be modified to support inspection of distributed thread identity by the management subsystem. As stated in the description of Brakes, JVM thread stack frames are switched into the associated Context object via a static interface defined by the context manager. For example, switching an integer from the stack into the context is done by inserting in the method's byte code at the appropriate code position an invocation of the following static method:

```

static curPushInt(int i) {
    Context c = getContext(Thread.currentThread());
    c.pushInt(i);
}

```

Such push-methods are defined for all basic types as well as for object references. Complementary, there are pop-methods for restoring execution-state from the context. The appropriate target Context object is looked up with the current thread identity as hashing key. However, in order to allow the context manager to manage Context objects on a per distributed thread basis, this static interface must be changed such that the context manager can inspect the identity of the current distributed thread. Thus:

```

public static pushInt(int i, D_Thread_ID id) {
    Context c = getContext(id);
    c.pushInt(i);
}

```

Note that byte code instructions must be inserted for retrieving the `D_Thread_ID` local variable before calling the static methods.

4.2. Efficient management of serialized distributed execution-state. Extending the static interface of the Brakes context manager allow us to build an associated distributed management subsystem, illustrated in Fig. 4.2, that manages captured execution-state on a per distributed thread basis. The management subsystem

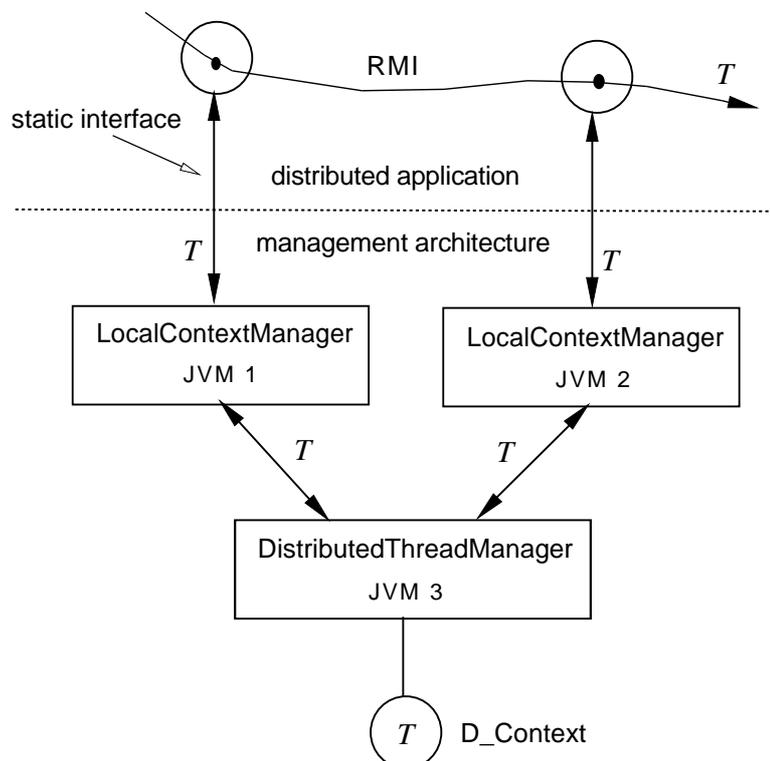


FIG. 4.2. *Distributed Architecture of the Context Manager.*

consists of a context manager for each JVM where the distributed application executes, to which we will refer as *local context manager*. Capturing and restoring code blocks still communicate with the static interface of the local context manager, but the captured execution-state is now managed per distributed thread by one central manager, the *distributed thread manager*.

To further deal with the problem that Brakes is not designed for capturing distributed execution-state, we also had to rearrange the management of the `isSwitching` and `isRestoring` flags. First, while in Brakes there was a separate `isSwitching` and `isRestoring` flag for each JVM thread, we now manage only one `isSwitching` and one `isRestoring` flag for the entire distributed execution-state of the application. Both flags are stored as static global variables on the distributed thread manager and are replicated with strong consistency on each local context manager. Furthermore we introduced a new flag, `isRunning`, associated with each individual distributed thread that marks the start of capturing and the end of reestablishing the execution-state of that distributed thread.

A positive side effect of the rearrangement of flags is that we drastically reduce the overhead during normal execution. When inspecting the global `isSwitching` and `isRestoring` flags during normal execution, no costly hashtable look up on distributed thread identity is performed anymore. Only during capturing and reestablishment the `isRunning` flag is looked up with `D_Thread_ID` as hashing key; however these look ups do not occur during normal execution. This choice may seem to be a trade-off between efficiency and flexibility. The rearrangement of flags results in a less flexible mechanism that can only capture execution-state at the level of the whole distributed execution state of the application. It is not possible to capture one distributed thread, without stopping other distributed threads. However, this coarse-grained scale is exactly what we want: it does not make sense to capture one thread, without stopping another thread when they are executing in the same application objects. In section 6.2 we discuss performance overhead more in detail.

4.3. External initiation of capturing and reestablishing. The Brakes JVM thread serialization mechanism is designed in the context of mobile agents and as such it is not designed for being initiated by an external control instance. To deal with this problem, we extended the Brakes transformer to insert extra byte codes at the beginning of each method body to verify whether there has been an external request for capturing the execution-state. We will refer to this code as `{external capturing request check}`. Furthermore, the distributed thread manager offers a public interface that enables an external control instance to initiate the capturing and reestablishing of distributed execution-state. Capturing of execution-state is started by calling the operation `captureState()` on the distributed thread manager. This method sets the `isSwitching` flag on all local context managers through broadcast. As soon as a distributed thread detects the `isSwitching` flag is set, (inside the first executed `{external capturing request check}` code) the distributed thread sets off its `isRunning` flag and starts switching itself into its context. Reestablishment of execution is initiated by calling the operation `resumeApplication()` on the distributed thread manager. This method sets the `isRestoring` flag on each local context manager and restarts the execution of all distributed threads. Each distributed thread detects immediately that the `isRestoring` flag is set, and thus restores itself from the context. Once the execution-state is reestablished the distributed thread sets on its `isRunning` flag (inside the `{external capturing request check}`) and resumes execution. When all distributed threads execute again, the distributed thread manager turns off the `isRestoring` flag on all local context managers through broadcast.

5. Run-time repartitioning at work. In this section we present our prototype for run-time repartitioning and demonstrate it for a simple text translator application. First we describe the process of run-time repartitioning. Then we give a sketch of the prototype. Next we illustrate the byte code transformations. Finally we explain the process of run-time repartitioning by means of an example.

5.1. The four phases of run-time repartitioning. run-time repartitioning aims to improve the global load balance or network communication overhead by repartitioning the object configuration of the application over the available physical nodes at run-time. We distinguish between 4 successive phases in the run-time repartitioning process. In the first phase, a load balancing monitor allows an administrator to monitor the application's execution and let him decide when to relocate the application objects over the available physical nodes. In the second phase, the monitor invokes the distributed thread manager to capture the distributed execution-state of the application. After this, the execution of all application objects is temporarily suspended and the corresponding distributed thread states are stored as serialized data in the context repository of the distributed thread manager. In the third phase, the monitor carries out the initial request for repartitioning by migrating the necessary objects over the network. In the final and fourth phase, the monitor invokes the distributed thread manager to reestablish the distributed execution-state of the application. As soon as the execution-state is reestablished, the application continues where it left off.

5.2. Prototype. In the run-time repartitioning prototype, the serialization mechanism is integrated with a simple load balancing monitor. We demonstrate the repartitioning prototype for a simple text translator

system (see Fig. 5.1). The text translator is composed with a number of objects that can be distributed over some hosts. For each translation a new distributed thread is started. A client sends the text with a source and target language to a Server object. The Server forwards the job to a ParserBuilder, who sends each sentence for translation to a Translator. The Translator uses a Dictionary object for the translation of individual words. As soon as a sentence is translated the Translator returns it to the ParserBuilder. The ParserBuilder assembles the translated text. Finally the translated text is returned to the Server who sends it back to the client. Fig. 5.2 gives

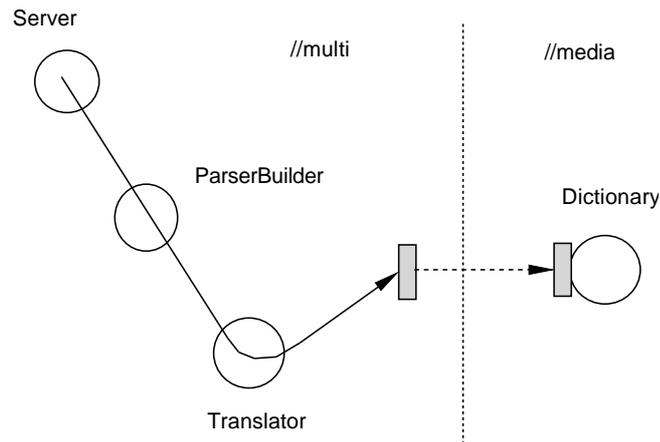


FIG. 5.1. *Prototype for Run-time Repartitioning.*

a snapshot of the load balancing monitor. The monitor offers a GUI that enables an administrator to do run-time repartitioning. The left and middle panels show the actual object distribution of the running application. The panels on the right show the captured context objects per distributed thread after a repartitioning request.

Since passive object migration, i.e., code and data migration, but no migration of run-time information like the program counter and the call stack, is necessary during the third phase of the run-time repartitioning process we used the mobile object system Voyager 2.0 [9] as distributed programming model in our prototype. To enable capturing and reestablishment of distributed execution-state during phases two and four, the im-

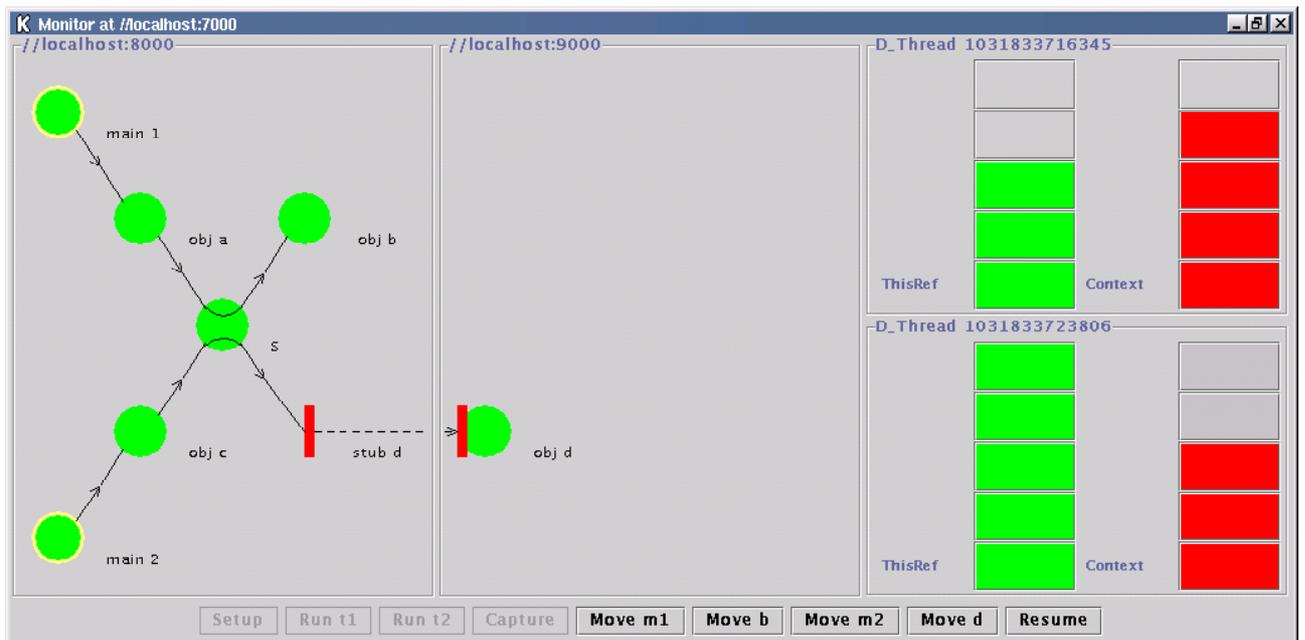


FIG. 5.2. *Snapshot of the Run-time Repartitioning Monitor.*

plementation code of the text translator system must be hauled through our byte code transformer. Finally, the distributed architecture of the associated management subsystem as shown in Fig. 4.2, defines an abstract framework that must be instantiated by a concrete distributed implementation. We used Voyager for this too, but another distribution platform like Java RMI was also possible. Thus in our prototype LocalContextManager and DistributedThreadManager are implemented as Voyager objects.

5.3. Byte code transformations. Before illustrating the process of run-time repartitioning we first give an overview of the transformation of the application code. We limit the extract to the principal code as in Fig. 5.3. The *Italic* marked code is inserted byte code. Each method signature as well as each method invocation is extended with an extra *D.Thread_ID* argument by the DTI transformer. The Brakes transformer has inserted code blocks for capturing the execution-state of a distributed thread (*switching code block*) and restoring it after-wards (*isRestoring code block*).

```

class Translator {
  Sentence analyze(Sentence sentence, D.Thread_ID, threadID) {
    {isRestoring code block}
    {external capturing request check}
    {switching code block}
    Sentence tSentence = new Sentence();
    Word word, tWord;
    sentence.resetCursor(threadID);
    {switching code block}
    while(!sentence.endOfSentence(threadID)){
      {switching code block}
      word = sentence.next(threadID);
      {switching code block}
      ** tWord = dictionary.translate(word, threadID);
      {switching code block}
      tSentence.add(tWord, threadID);
      {switching code block}
    }
    return tSentence;
  }
  ...
  private Dictionary dictionary;
}

class Dictionary {
  * public Word translate(Word word, D.Thread_ID threadID) {
    {isRestoring code block}
    {external capturing request check}
    {switching code block}
    ...
  }
  ...
}

```

FIG. 5.3. Application code after Byte code Transformation.

5.4. An example of run-time repartitioning. We now explain the process of run-time repartitioning starting from Fig. 5.1. Suppose the administrator decides to migrate the Dictionary object during the translation of a text. At a certain moment, lets say when the control flow enters the `translate()` method in Dictionary, the administrator pushes the capture button on the monitor. This point is marked in Fig. 5.3 with *. At that moment the execution-state of the distributed thread is scattered over two hosts as illustrated in Fig. 5.4. Pushing

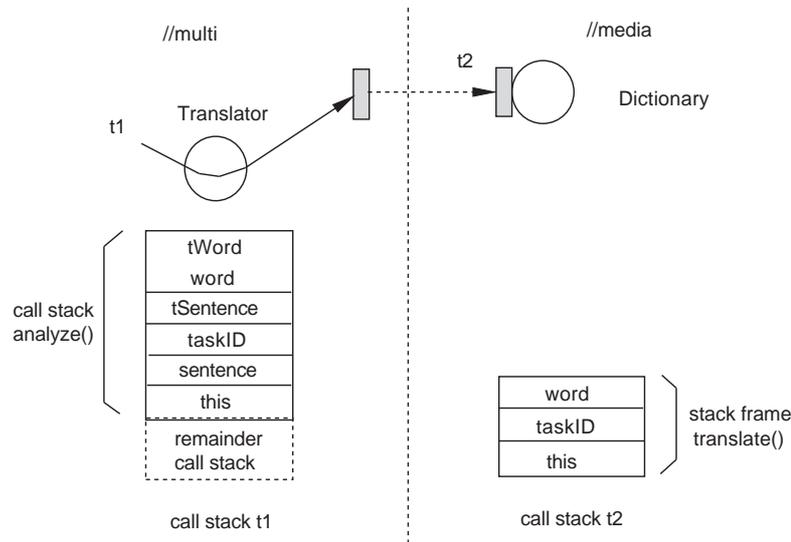


FIG. 5.4. *Distributed Execution—state before Capturing.*

the capture button invokes `captureState()` on the distributed thread manager that sets the `isSwitching` flag. The distributed thread detects this in the `{external capturing request check}` code. Immediately it set off its `isRunning` flag and saves the execution-state of the `translate()` method. This includes: (1) the stack frame for `translate()` (i.e. the only frame for thread `t2` on the call stack, see Fig. 5.4); (2) the index¹ of the last invoked method in the body of `translate` (i.e. zero for `{external capturing request check}`); (3) the object reference to the Dictionary object (i.e. the `this` reference). For reasons of efficiency, the execution-state is buffered per stack frame by the local context manager until completion of switching. Then the local manager forwards the serialized data to the distributed manager who stores it into the context repository of the distributed thread with the given `D_Thread_ID`.

The last instruction of the `{switching code block}` is a return. This redirects the control to the previous frame on the call stack of the distributed thread, in our case the `analyze()` method. The control flow returns from host to host (i.e. from `media` to `multi`, see Fig. 5.4) which means that the execution-state of the JVM thread `t2` now is completely saved. In the code of Fig. 5.3, we then reach the point marked as `**`. Next the execution-state for `analyze()` is saved, i.e.: (1) the stack frame for `analyze()` (i.e. the top frame of JVM thread `t1` as in Fig. 5.4): (2) the index of the last invoked method in the body of `analyze` (i.e. 4, see Fig. 5.3); (3) the object reference to the Translator object. As soon as the buffered data is written to the context repository of the distributed thread manager, another return at the end of the `{switching code block}` redirects the control flow to the previous frame on the call stack. Subsequently, the execution-state for that method is saved. This process recursively continues until the JVM thread `t1` returns to the `run()` method of `D_Thread`. At that time the `DistributedContext` contains the complete distributed execution-state of the distributed thread.

Once the complete distributed execution-state is saved the Dictionary object can be migrated from `media` to `multi`. To this purpose the administrator pushes the corresponding migrate-button on the monitor. As explained in section 5.2, we used Voyager as distribution platform for our prototype. Voyager dynamically transfers the object references from remote to local and vice verse.

¹this index refers to the number the Brakes transformer associates with the subsequent `invoke` instructions in the body of each method, starting with 0 for external capturing request check code, 1 for the first `invoke` instruction and so on; the index of the last performed `invoke` instruction is saved in the context to remember which methods where on stack

Once the repartitioning is ready the administrator can push the resume button which invokes `resumeApplication()` on the distributed thread manager. That turns off the `isSwitching` flag and sets the `isRestoring` flag. Next a new JVM thread is created at multi to resume the translation. During the reestablishing process the relevant methods are called again in the order they have been on the stack when state capturing took place. The new thread takes the original `D_Thread_ID` with it. This is the key mechanism for reconstructing the original call stack. Each time inserted byte code reestablish the next stack frame the managers use the distributed thread identity to select the right context object for that particular distributed thread.

Fig.5.5 illustrates the reestablishment of the last two frames in our example. When `analyze()` on `Translator` is invoked, the `isRestoring` code block will be executed, based on the actual state of the flags (`isRestoring = on`, `isRunning = off`). The inserted byte code restores the values of the local variables of the `analyze()` frame one by one via the local context manager (see the left part of Fig. 5.5). At the end the index of the next method to invoke is picked up. For the `analyze()` frame an index 4 was earlier saved, so the fourth method, i.e. `translate()`, must be invoked on the `Dictionary` object. The managers pickup the reference to this object and `translate()` will be invoked. Again the `isRestoring` code restores the local execution-state of the `translate()` method (the right part of Fig. 5.5). This time the index for the next method is zero. This is the signal for the context manager to reset the `isRunning` flag of the current distributed thread and to resume its normal execution. At this point the expensive remote interaction between the translator and the dictionary objects is transferred into a cheap local invocation. Note that in the example the execution-state is captured with the

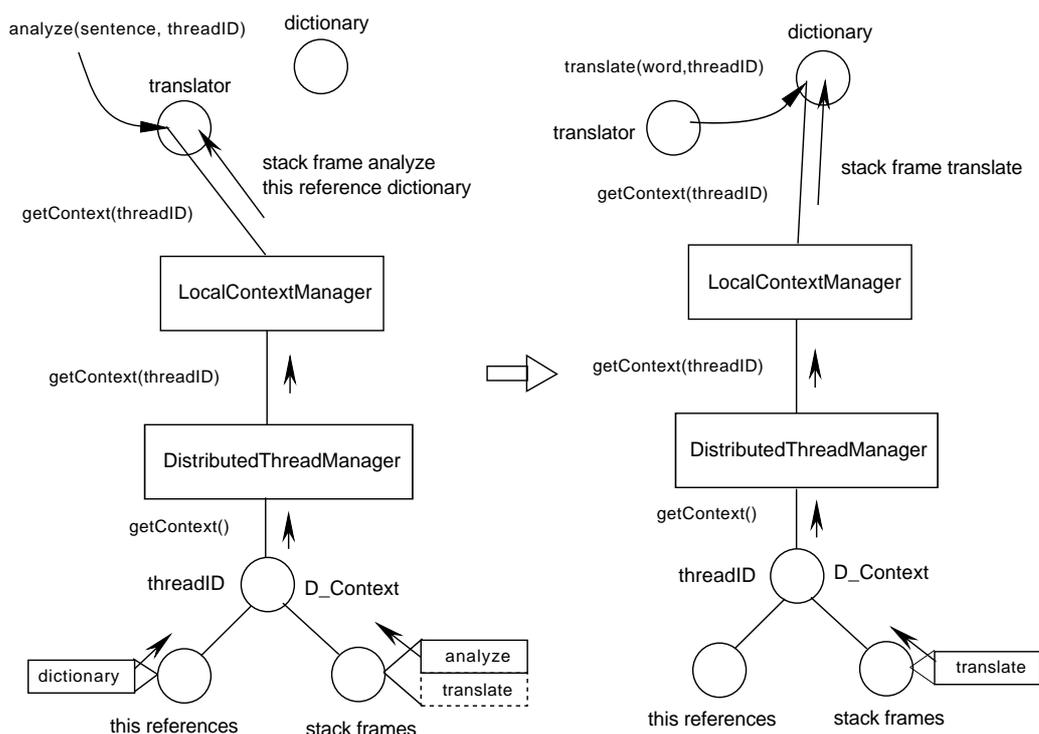


FIG. 5.5. Reestablishing a Distributed Execution-state.

translator object waiting on a pending reply of a remote method invocation (i.e. `dictionary.translate(word, threadID)`, see Fig. 5.3). This illustrates the advantage of phases two and four of our run-time repartitioning scheme in keeping execution and object migration completely orthogonal to each other.

6. Evaluation. In this section we evaluate the serialization mechanism for a distributed execution-state. Since inserting byte code introduces time and space overhead we look to the blowup of the class files and give results of performance measurements. To get a representative picture, we did tests on different types of applications. At the end we look at the limitations of our current implementation and outline the restrictions of our model.

6.1. Blowup of the byte code. The blowup of the byte code for a particular class highly depends on the number and kind of defined methods. Since the Brakes transformer inserts code for each invoke-instruction that occurs in the program, the space overhead is directly proportional to the total number of invoke-instructions that occur in the application code. Per invoke-instruction, the number of additional byte code instructions is a function of the number of local variables in the scope of that instruction, the number of values that are on the operand stack before executing the instruction and the number of arguments expected by the method to be invoked. The DTI transformer rewrites method and class signatures. This adds a space overhead proportional to the number of signature transformations. We measured the blowup for three kinds of applications:

1. Low degree of method invocation, i.e., the program has a structure `main{m1;}` thus the code is compacted in one method body;
2. Nested method invocations, i.e., the program has a structure `main{m1;}; m1{m2; m3;}; m3{m4;}` thus the code is scattered over a number of nested methods;
3. Sequential method invocations, i.e., the program has a structure `main{m1; m2; m3; m4;}` thus the code is scattered over a number of sequential non-nested methods.

Fig. 6.1 shows the results of our measurements. The dark part of the bars represents the size of the original code. The gray part represents the additional code for distributed thread functionality, while the white part represents the extra added code for full serialization functionality. We measured an average blowup for distributed thread identity of 27 % and 83 % for full serialization functionality. The expansion for Sequential is rather high, but its code is a severe test for blowup.

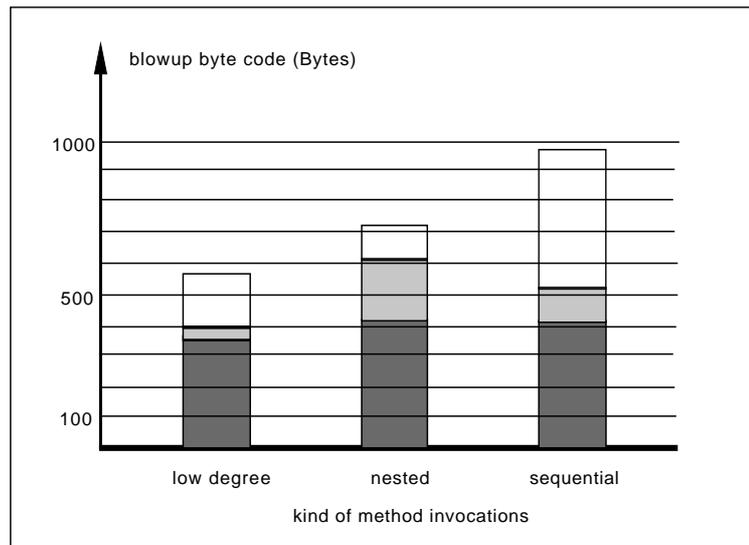


FIG. 6.1. *Byte code Blowup for Three kinds of Applications.*

6.2. Performance measurements. For performance measurements, we used a 500 MHz Pentium III machine with 128 MB RAM with Linux 2.2 and the SUN 2SDK, JIT enabled. We limited our tests to the overhead during normal execution. This overhead is a consequence of the execution of inserted byte code. Fig. 6.2 shows the results of our tests. The dark bars represent the execution speed of the original application code; the gray parts represent the overhead due to byte code transformation for distributed thread identity, while the white parts represents the additional overhead for distributed thread serialization. We measured an average overhead of 3 % for distributed thread identity. For full serialization functionality we get an average overhead of 17 %, a quite acceptable result. Note that “normal” applications typically are programmed in a nested invocation style. As such, the results for the Nested application are a good indication for blowup and performance overhead in practice. It is difficult to compare our measurement results with other systems, since to our knowledge, no related system truly covers functionality for serialization of *distributed* execution-state as our system does. In section 7 we discuss related work.

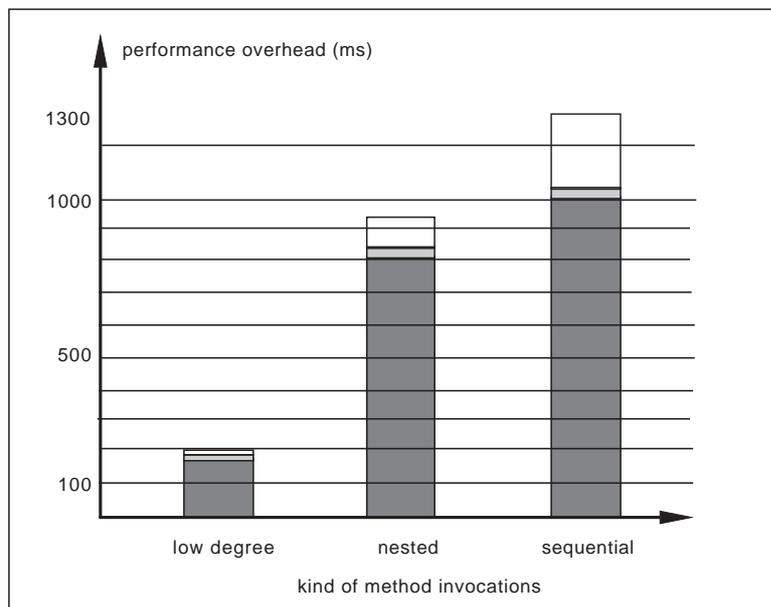


FIG. 6.2. Performance Overhead for Three kinds of Applications.

6.3. Limitations of the current implementation. Our byte code transformers have some limitations we intend to eliminate in the future. Although possible, we have not yet implemented state capturing during the execution of an exception handler. The major difficulty here is dealing with the `finally` statement of a `try` clause. Currently our byte code transformer throws away all debugging information associated with a Java class. This affects the ability to debug a transformed class with the source-code debugger. Furthermore, the DTI byte code transformer encapsulates each user defined JVM thread into a `D_Thread`, but currently ignores other JVM thread related code. Thus our current model doesn't support aspects such as thread locking, e.g. in synchronized code sections.

6.4. Restrictions of the model. Our model is intended and only applicable for applications running on top of a dedicated cluster of machines where network latencies are low and faults are rare. This is not directly a dependability of our approach, but rather a dependability of the RPC-like programming model: performing blocking calls on remote objects is after all only feasible on a reliable, high-bandwidth and secure network. Furthermore, in section 4.2 we already mention that the granularity of our repartitioning algorithm is at JVM level. As soon as the `isSwitching` flag is set all running distributed threads are suspended together irrespective of whatever application they belong. Thus applications that are transformed with our byte code transformer and that execute on the set of involved Java Virtual Machines will be suspended together.

Since we extract thread execution-state at byte code level we cannot handle a method call that causes a native method to be placed on the thread stack. Therefore programs that use reflection will not work properly with our repartitioning model.

In our prototype we transformed the application classes before deployment, but it is possible to defer this byte code transformation until run-time. In Java, this can easily be realized by implementing a custom class loader that automatically performs the transformation. In this regard, the overhead induced by the transformation process (which is not small for the current implementation) becomes a relevant performance factor.

7. Related work. The discussion on related work is organized according to the related fields that touches our work.

7.1. Distributed Threads. Existing work [4] in the domain of (real-time) distributed operating systems has already identified the notion of distributed thread as a powerful basis for solving distributed resource management problems. D. Jensen at CMU introduced the notion of distributed thread in the Alpha distributed real-time OS kernel. The main goal of distributed threads in the Alpha kernel was integrated end-to-end

resource management based on propagation of scheduling parameters such as priority and time constraints. In our project we adopted the notion of distributed thread (identity) at the application level. This allows the distributed management subsystem to refer to a distributed control flow as one and the same computational entity.

7.2. Strong Thread Migration. Several researchers developed mechanisms for strong thread migration in the context of Mobile Agents. Stefan Fünfroeken at TU Darmstadt [7] has implemented a transparent serialization mechanism for local JVM threads by processing the source code of the application. In this approach the Java exception mechanism is used to capture the state of an ongoing computation. A source code transformation requires the original Java files of the application. Besides, it is much easier to manipulate the control flow at byte code level than at source code level. As a result byte code transformation is much more efficient especially in terms of space overhead. Sakamoto et al. [11] developed a transparent migration algorithm for Java application by means of byte code transformation. They too used the exception mechanism of Java to organize the serialization of the execution-state. We have chosen not to use the exception mechanism, since entries on the operand stack are discarded when an exception is thrown, which means that their values cannot be captured from an exception handler. Sakamoto et al. solved this problem by copying all those values in extra local variables before method invocation, but this causes much more space penalty and performance overhead. In her dissertation [17], Wei Tao proposes another portable mechanism to support thread persistence and migration based on byte code rewriting and the Java exception mechanism. Contrary to the others, this mechanism also works for applications that use synchronization and locks.

7.3. Multi-Threading for Distributed Mobile Objects in FarGo. Abu and Ben-Shaul integrated a multi-threading model for distributed and mobile objects in the FarGo framework [1]. A FarGo application consists of a number of 'complets'. Complets are components similar to components in other frameworks, and in addition they are the unit of relocation in the model. The distributed mobile thread model of FarGo is based on a thread-partitioning scheme. Programmers must mark a migratable complet as thread-migratable (T-Migratable) by implementing the empty `T_Migratable` interface. The FarGo compiler uses this interface to generate proper thread partitioning code. Thread partitioning is integrated in the complet reference architecture, i.e. a stub and a chain of trackers. When a `T_Migratable` complet is referenced the invoking thread waits, and a new thread is started in the referenced complet. The migration itself is based on the source code transformation of Fünfroeken's migration scheme [7]. FarGo introduces a new distributed programming model to support migration, while our mechanism can be used for any existing Java RMI application.

7.4. Byte Code Transformations for Distributed Execution of Java applications. The Doorastha system [6] allows implementing fine-grained optimization's for distributed applications just by means of code annotations. Doorastha is not an extension or a super-set of the Java language but instead is based on annotations to pure Java programs. By means of these annotations it is possible to dynamically select the required semantics for distributed execution. This allows to develop a program in a centralized (multi-threading) setting first and then prepare it for distributed execution by annotation. Byte code transformation will generate a distributed program whose execution conforms to the selected annotations. Researchers at the University of Tsukuba, Japan [14] developed a system named Addistant, which also enables the distributed execution of a Java program that originally was developed to run on a single JVM. We share the common point of view with the researchers that software developers often prefer to write software apart of non-functional aspects. It's only at deploy-time that those aspects have to be integrated, preferable in a transparent way. To avoid deadlock in the case of call back, Addistant guaranties that local synchronized method calls of one distributed control flow are always executed by one and the same JVM thread. Therefore it establishes a one-to-one communication channel for the threads that take part in such an invocation pattern. Such a communication channel is stored as thread local variable. In this approach it isn't necessary to pass distributed thread identity along the call graph of the distributed control flow. On the other hand for a run-time repartitioning system, thread identity must be propagated with every remote invocation anyway.

8. Conclusion and future work. In this paper we presented a mechanism for serialization of a distributed execution-state of a Java application that is developed by means of a distributed control-flow programming model such as Java RMI. This mechanism can serve many purposes such as migrating execution-state over the network or storing it on disk. An important benefit of our mechanism is its portability. It can be integrated into existing applications and requires no modifications of the JVM or the underlying ORB. However, because

of its dependability on the control-flow programming model, our mechanism is only applicable for distributed applications that execute on low latency networks where faults are rare. Our contribution consists of two parts. First we integrated Brakes, our existing serialization mechanism for JVM threads, in a broader byte code translation scheme to serialize the execution-state of a distributed control flow. Second we integrated a mechanism to initiate the serialization of the distributed execution-state from outside the application. We applied the serialization mechanism in a prototype for run-time repartitioning of distributed Java applications. Our repartitioning mechanism enables an administrator to relocate application objects at any point in an ongoing distributed computation.

Since for Java RMI-like applications, logical thread identity is lost when the control flows crosses JVM boundaries, we introduced the notion of distributed thread identity. Extending Java programming with distributed thread identity provides a uniform mechanism to refer to a distributed control flow as one and the same computational entity.

Often byte code transformation is criticized for blowup of the code and performance overhead due to the execution of inserted byte code. Based on a number of quantitative analyses we may conclude that the costs associated with our byte code translation algorithm are acceptable. Some limitations of our serialization mechanism for a distributed execution-state have to be solved. Finally it 's our intention to build a complete tool for run-time repartitioning for distributed Java RMI applications. Therefore we have to extend our current monitoring and management subsystem with several other features such as functionality for dynamic adaptation of object references after migration, support for different load balancing algorithms and an application adaptable monitor.

The latest implementation of the run-time repartitioning tool is available at:

<http://www.cs.kuleuven.ac.be/~danny/DistributedBrakes.html>

Acknowledgments. This research was supported by a grant from the Flemish Institute for the advancement of scientific-technological research in industry (IWT). We would like to thank Tim Coninx and Bart Vanhaute for their valuable contribution to this work. A word of appreciation also goes to Tom Holvoet and Frank Piessens for their usefully comments to improve this paper.

REFERENCES

- [1] M. ABU, I. BENSHAUL, *A Multi-Threading model for Distributed Mobile Objects and its Realization in FarGo*, in Proceedings of ICDCS 2001, The 21st International Conference on Distributed Computing Systems, April 16–19, 2001, Mesa, AZ, pp. 313–321.
- [2] M. BAKER, *Cluster Computing White Paper*, <http://www.dcs.port.ac.uk/~mab/tfcc/WhitePaper/> (2000).
- [3] S. BOUCHENAK, *Pickling threads state in the Java system*, ERSADS 1999, Third European Research Seminar on Advances in Distributed Systems, April 1999, Madeira Island, Portugal.
- [4] R. CLARK, D.E. JENSEN, AND F.D REYNOLDS, *An Architectural Overview of the Alpha Real-time Distributed Kernel*, In Proceedings of the USENIX Workshop, In 1993 Winter USENIX Conf., April 1993, pp. 127–146.
- [5] M. DAHM, *Byte Code Engineering*, in Proceedings of JIT'99, Clemens Cap ed., Java-Informationen-Tage 1999, September 20 - 21, 1999, Düsseldorf, Germany, pp. 267–277.
- [6] M. DAHM, *The Doorastha system*, Technical Report B-I-2000, Freie Universität Berlin, Germany (2001).
- [7] S. FUNFROCKEN, *Transparent Migration of Java-based Mobile Agents*, in Proceedings of the 2e International Workshop on Mobile Agents 1998, Lecture Notes in Computer Science, No. 1477, Springer-Verlag, Stuttgart, Germany, September 1998, pp. 26–37.
- [8] K.C. NWOSU, *On Complex Object Distribution Technique for Distributed Computing Systems*, in Proceedings of the 6th International Conference on Computing and Information (ICCI'94), Peterborough, Ontario, Canada, May 1994, (CD-ROM).
- [9] OBJECTSPACE INC., *VOYAGER, Core Technology 2.0*, <http://www.objectspace.com/products/voyager/> (1998).
- [10] B. ROBBEN, *Language Technology and Metalevel Architectures for Distributed Objects*, Ph.D thesis, Department of Computer Science, K.U.Leuven, Belgium, ISBN 90-5682-194-6, 1999.
- [11] T. SAKAMOTO, T. SEKIGUCHI, A. YONEZAWA, *Bytecode Transformation for Portable Thread Migration in Java*, in Proceedings of the Joint Symposium on Agent Systems and Applications / Mobile Agents (ASA/MA), Lecture Notes in Computer Science 1882, Springer-Verlag, ETH Zürich, Switzerland, September 13-15, 2000, pages 16–28.
- [12] W. SHU AND M. WU, *An Incremental Parallel Scheduling Approach to Solving Dynamic and Irregular Problems*, In Proceedings of the 24th International Conference on Parallel Processing, Oconomowoc, WI, 1995, pages II:143–150.
- [13] C. SZYPSKI, *Component Software: Beyond Object-Oriented Programming*, Addison-Wesley and ACM Press, ISBN 0-201-17888-5, 1998.
- [14] M. TATSUBORI, T. SASAKI, S. CHIBA AND K. ITANO, *A Bytecode Translator for Distributed Execution of Legacy Java Software*, in Proceedings of the 15th European Conference on Object Oriented Programming (ECOOP 2001), Lecture Notes in Computer Science 2072, Springer-Verlag, Budapest, Hungary, June 18-22, 2001, pp.236–255.
- [15] E. TRUYEN, B. ROBBEN, B. VANHAUTE, T. CONINX, W. JOOSEN AND P. VERBAETEN, *Portable Support for Transparent Thread Migration in Java*, in Proceedings of the Joint Symposium on Agent Systems and Applications / Mobile Agents

- (ASA/MA), Lecture Notes in Computer Science 1882, Springer-Verlag, ETH Zürich, Switzerland, September 13-15, 2000, pages 29–43.
- [16] E. TRUYEN, B. VANHAUTE, B. ROBBEN, F. MATTHIJS, E. VAN HOEYMISSEN, W. JOOSEN AND P. VERBAETEN, *Supporting Object Mobility, from thread migration to dynamic load balancing*, OOPSLA'99 Demonstration, November '99, Denver, USA, 1999.
- [17] W. TAO, *A Portable Mechanism for Thread Persistence and Migration*, PhD thesis, University of Utah, <http://www.cs.utah.edu/tao/research/index.html> 2001.

Edited by: Dan Grigoras, John P. Morrison, Marcin Paprzycki

Received: September 24, 2002

Accepted: December 15, 2002