

Context Management and Its Applications to Distributed Transactions

George Samaras*, Ajay D. Kshemkalyani, and Andrew Citron

IBM Corporation
P. O. Box 12195, Research Triangle Park, NC 27709

Abstract

An emerging paradigm that handles multiple locii of control in a system allows multiple program threads to work on the same task, each thread to work on a different task, or a thread to work on multiple tasks for greater design flexibility or due to system constraints such as real-time demands and a high load on tasking. We use the definition of context to capture the notion of logical locus of control. The context of the work being currently executed must be identifiable uniquely by the application, the Resource Managers and the Transaction Manager because each context represents different work. In this paper, we define context management by defining a local Context Manager and its user interface. We then show why the notion of context is required to solve the problems that arise in local and distributed transaction processing due to the emerging paradigm. We present solutions to these problems in transaction processing using the proposed context management.

1. Introduction

Currently, operating systems that handle multiple applications provide a separate physical locus of execution for each individual application program. The physical locus of execution is the process for single-threaded processes provided by operating systems such as DOS and VM, and it is the thread¹ for multi-threaded processes provided by operating systems such as UNIX², OS/2³ [5, 12], Windows NT⁴ and Windows 95⁴ [8]. Operating systems support the client-server model of computing by dispatching a separate

server process or thread to handle a new request from the clients.

There are two recent trends which indicate that the existing support provided by systems is unsuitable for a range of application programs. First, as applications grow in number and get more distributed, the number of applications a server can support becomes limited by the operating system constraints such as the number of processes/threads allowed within the system. Second, as applications become more numerous and response times become critical for real-time systems, the servers cannot afford the overhead of process start-up and switching, or forking and dispatching, and the overhead of locking mechanisms for access to shared tables for each new application. A new processing paradigm is now evolving to overcome the above problems and to provide more flexibility to distribute tasks across processes/threads. The emerging paradigm is as follows: a server process or thread can concurrently support multiple applications, or an application can be distributed across multiple processes and/or threads. This paper proposes how this new paradigm can be supported, and discusses its interaction with transaction processing. We focus on distributed transaction processing as an application because it represents an important and growing class of applications, and it was our involvement in distributed transaction processing that triggered this work.

In the paradigm outlined above, each initiated transaction supported by a thread(s) is explicitly associated with a *context*. Thus, a context becomes a logical locus of control and represents a transaction and its associated resources. The above notion of context is similar to X/Open's notion of "thread-of-control" [16]. With the new processing paradigm, multiple transactions can be associated with a thread, representing multiple contexts per thread. However, at any instant, only one of these contexts will be active. The application and the system should be able to specify and determine (i) which context is currently being worked on by the thread, and (ii) all the resources associated with the processing of any context. There is an explicit need to coordinate the contexts within and across threads and processes,

*Currently with University of Cyprus, Nicosia, Cyprus.

¹For commit processing of a transaction, the locus of execution is still the process, not the thread. This is a drawback of existing transaction processing design for multithreaded systems.

²UNIX is a registered trademark in the United States and other countries licenced exclusively through X/Open Company Limited.

³OS/2 is a trademark of the IBM Corporation.

⁴Windows, Windows NT, and Windows 95 are trademarks of Microsoft Corporation.

and to coordinate access to resources by multiple contexts (within a thread, across threads, and across processes). This is achieved through a Context Manager mechanism and its associated user interface defined in this paper.

The contribution of this paper is that we define context, the *Context Manager (CM)* and its user interface, and show how context is used for local and distributed transaction processing. We highlight the role of context in a multithreaded, real-time, high tasking operating system environment, provide different practical styles of transaction management using context management, and show how to use context management to solve deadlocks, protocol violations and loopbacks. We have implemented a prototype of the Context Manager for use with the SNA LU6.2 Syncpoint Services [7]. The VM operating system Shared File System has been enhanced to provide a version of the context management support described here.

The paper is organized as follows: Section 2 describes the system model. Section 3 examines the requirements for a new notion of context, and defines and describes the operation and usage of context, along with a new Context Management user interface. Two examples of the use of context management in a multithreaded environment are given. Section 4 describes the problems that arise in distributed transaction processing when the new paradigm, viz., assigning multiple transactions to multiple threads, is used. It then shows how to use context management to solve the problems. Section 5 concludes.

2. System Model

A distributed system consists of a set of computing nodes linked by a communications network. The nodes of the system cooperate with each other in order to process distributed computations. For the purpose of cooperation, the nodes communicate by exchanging messages via the communications network.

A multiprogramming/multiprocessing operating system runs at each node. A process which is an executing program has a single address space and a single thread of control for the program. The state information for the process consists of page tables, swap images, file descriptors, outstanding I/O requests, and saved register values. Multiple programs are handled by maintaining and switching between processes. If threads or lightweight processes are supported by the operating system, then the threads of a process concurrently execute within the same address space. Each thread uses a separate program counter, a stack of activation records, and a control block which contains information necessary for thread management. Most of the information that is part of a process is shared by all the threads executing in the same address space. This reduces the overhead in creating and maintaining information, and the information that has

to be saved when switching between threads of the same program.

A *distributed transaction* is a program of one or more statements that access data distributed on different nodes in the system. Each transaction has a unique identifier denoted TRANID. The execution of a distributed transaction requires a distributed *commit protocol* to ensure that the effects of the distributed transaction are *atomic*, i.e., either all the effects of the transaction persist or none persist, whether or not failures occur [4].

Once the computations of a transaction are completed, the application instructs the *transaction manager (TM)* of its node (site) to initiate and coordinate the commit protocol. The “logical locus of control” from which the application issues the commit command is the entity that identifies to the TM the transaction to be committed. At each node, the local *Resource Managers (RMs)*, such as database and file managers, and *Communication Resource Managers (CRMs)* participate in the commit protocol. The RMs/CRMs commit only those resources that are associated with the current “logical locus of control” and transaction. The CRM embodies the communication protocol and provides a local view of the remote processes and remote TMs. The TMs that participate in the commit processing include one *coordinator* and one or more *subordinates*. The coordinator coordinates the final outcome of the commit processing by issuing a COMMIT or ABORT, that is propagated to all subordinates. Subordinate TMs propagate the decision to their subordinate TMs or local RMs. The commit operation employs the well-known *two-phase commit (2PC)* protocol [4, 11, 15].

3. Context

Definition 1 *A context is a unique local logical locus of control shared by the application, TM and RMs to manage their resources, and relate their resources to the resources owned by other RMs within the system. A context represents a grouping of resources within the system, needed to perform a particular function in a logical locus of control, or to show the inter-relationship between diverse resources.*

3.1. Requirements

There are several emerging trends that require the notion of context.

1. Currently, a process or thread is associated with at most a single logical locus of control. This paradigm is proving inadequate for some applications because:
 - The client-server paradigm requires a server instance to accept multiple incoming requests. For example, the asynchronous RPC style of

distributed programming uses this model [1]. The following difficulties arise if a different thread/process is used for each request. (a) First, as applications grow and become more distributed, the number of transactions a server can support becomes limited by the operating system constraints such as the number of processes/threads allowed within the system. For example, OS/2 2.1 and Warp can support 4096 threads. (b) Second, as these transactions begin supporting larger volumes and response time becomes critical for real-time systems, the servers cannot afford the overhead of process start-up, switching, dispatching or forking, locking mechanisms for accessing shared tables, and extra storage for each new transaction. A thread should be able to support multiple transactions and temporarily suspend work on a long-running request to process work for another request. This minimizes demand on operating system resources, while allowing greater parallelism in servicing requests. The result is better response time and/or better throughput.

- A message routing program in a large database system acts as a router based on the content of the message. A database system typically uses long-lived programs that handle transactions from more than one end-user or transaction at a time and that can activate other conversations based on the database activity and the input. In both these cases, the same thread or process should accept the various incoming routing requests (locii of execution) rather than have separate threads handle the various routing requests for better efficiency.

A new paradigm that allows a server process/thread to support multiple transactions is required.

2. The notion of context provides useful functionality allowing process-oriented as well as thread-based systems the flexibility needed in today's complex and demanding environment. For operating systems that allow applications to spawn threads or fork processes, it is desirable to allow the server application to divide the incoming requests however the application chooses to. Some applications might want all threads to work on the same request. Other applications might want each thread to work on a different request. A context management service needs to allow each thread to identify each of the contexts it is associated with. The context must be independent from the operating system's task dispatching mechanism.

The above requirements express the emergence of environments where:

- A server process or thread can accept requests from different end users, and the server is allowed to suspend work on one request to work on a different request, or
- Multiple server processes/threads are working on related work representing the same context of the application.

The application needs a way to inform the TM and the RMs which task the application is working on at a particular time, and the TMs and RMs need to coordinate to have a common understanding of which context is currently under execution, for the following reasons:

- To group together logically related work and separate logically unrelated work.
- Each request is likely to be a part of a different atomic transaction. The work a server process does on behalf of one transaction must commit or abort independently from other unrelated work that the server was handling.
- The security authorization of each request can be different. An application must make sure that the system's security manager and other RMs cooperate to ensure that the access granted at the particular instant is proper for the end user application that is currently being worked on.

For the above, a *thread_id* or *process_id* does not suffice to identify the logical locus of control; rather a *context_id* is required. The notion of context allows the management of multiple transaction program instances within a single process or single thread, as well as the management of a single transaction program instance across several processes/threads, implicitly externalizing the creation, coordination and deletion of transaction program instances via the Context Manager interface. It also allows RMs to associate their resources with a transaction program instance. This allows other subsystems, such as a TM to coordinate the usage of all resources related to a particular task (transaction). The commit issued by the application applies to the context from which it was issued; not to the thread that issued it.⁵ A server process has special responsibilities [3]. The server process must correctly indicate which end user context it is working on behalf of. When many threads are acting on behalf of the same context, the application must make sure that the work of all the threads is completed before kicking off 2PC processing. Similarly, when a thread does work on behalf of multiple transaction programs/applications, it must be ensured that all the work in the relevant context is completed before kicking off 2PC processing.

⁵Currently, for commit processing, the locus of control is confined to the process level and is not allowed to the thread level, even though multi-threaded environments are common. This is an existing anomaly.

3.2. Context Management

Context management is a local application-support mechanism that permits applications to manage logically separate pieces of work within a single physical locus of execution (e.g., a thread or a process) [2]. When an application program uses context management, the CM keeps track of the various contexts, allows the application program to create and set a context for work, and allows an application program to switch contexts when appropriate. The context manager shares the notion of contexts with the RMs and the application program.

A new incoming conversation is assigned a new context and the program's current context is set to the new context. When a new outgoing conversation is allocated to a partner program, the conversation is assigned the current context of the program.

While context management provides additional functions such as flexibility and avoidance of process switching to the application program, it also imposes extra burden on the application program. The application program has to keep track of the progress done in each context, and switch contexts to meaningfully exploit the features of context management. An application program can perform context management by exploiting the functions provided by the context manager using a suite of calls, described subsequently.

The notion of context provides a logical separation of work done by an application program; each logically separate piece of work is done in a separate context. Context is local to a system, and distributed work done for the same transaction is not part of the same context. When an application program is involved in multiple transactions at a time, each transaction is done in a separate context at the application program. The same context may be involved in multiple transactions sequentially. The application program assumes the responsibility of keeping track of its various contexts, coordinating the data spaces of the various contexts, and switching between contexts.

The CM administers the contexts independently of the operating system's task dispatching mechanism. At each node, the context is identified by a *context_id*; a CM maintains a context table that stores the following information per *context_id*:

- *context_id*, unique to a node
- (*thread_id*, *process_id*), unique to a node
- a boolean indicating whether this context is currently being worked on by this (*thread_id*, *process_id*)

The TM maintains the correspondence between the TRANID and the *context_id*, the RM maintains the correspondence between the resources it manages and the *context_id*; likewise for the managers of other subsystems. We define a Context Management interface that (i) implicitly

externalizes the creation, coordination and deletion of contexts, and (ii) allows the application and managers of various subsystems such as the TM, RM, and security manager to associate their resources with a context and coordinate the usage of all resources related to a particular context.

3.2.1 Context Management Calls

In an environment where there is no one-to-one association between threads and contexts (but rather a many-many association), explicit context management calls are required to control the association between threads and contexts. This is particularly important for transaction processing because the TM and the RMs must identify the TRANID on a context basis and not on a thread basis. The *thread_id*, TRANID, accounting and security information are all part of a context.

When a thread (more generally, a *locus of execution*) is created, it is assigned a context. The context can be a new one or an inherited one.

Transaction management in a multithreaded environment

We describe three transaction program styles that use threads [2] and suggest other CM functions to support the three styles. The proposals are commands to start and manage the threads within the context management framework.

Style 1

Server receives new work, and kicks off a thread to handle the new work. A new context is implicitly created whenever new work is accepted. This is a typical approach for an RPC server, or an OS/2 LU6.2 TP that issues RECEIVE_ALLOCATE and then waits for the next incoming work. The context management function to support this is:

- **START_THREAD_AND_HANDOFF_CONTEXT** : This function starts a new thread and disassociates the main (old) thread from the newly created context. The forked thread is associated with the new context.

Style 2

The server kicks off a thread but the forking thread continues to work on the same context. This is typical of an application that can take advantage of the parallelism that light-weight threads provide. The CM function associated with this is:

- **START_THREAD_AND_SHARE_CONTEXT** : This function starts a new operating system thread. Any one of the threads can initiate the commit operation. The commit operation affects all resources allocated to this context. It is up to the application's design to ensure that all threads are ready for commitment. If some threads are not ready, the commit call may

return a state-check, or it may backout, or accidentally commit work in-progress. The best approach is to have the main thread issue COMMIT, after all forked threads report that are ready (using an OS wait/post mechanism for example).

Style 3

The main thread receives a new work request, and then instead of forking a new thread, it hands the work to an existing thread. For performance reasons, it is better to avoid creating a new thread. So using prestarted threads is faster than creating a new one:

- **HANDOFF_CONTEXT** : This function gives exclusive ownership of a context to an existing thread and posts the thread to inform it a new context is available.
- **SHARE_CONTEXT** : This function permits shared ownership of a context with an existing thread and posts the thread to inform it a new context is available.
- **THREAD_DONE_WITH_CONTEXT** : This function allows a thread to disassociate itself with a context and get ready to be involved in a new context. If no other thread is associated with the context, an implicit commit is attempted.
- **GET_NEW_CONTEXT** : This function allows a thread to wait for a parent thread to issue HANDOFF or SHARE_CONTEXT. Blocking and non-blocking flavors are useful.

In addition, `Extract_Current_Context` and `Set_Context` are needed to support threaded applications and applications that are single-threaded, but process more than one context at a time.

- **EXTRACT_CURRENT_CONTEXT** : This function allows TM, RM, or the application to find out which context is currently active.
- **SET_CONTEXT** : This function allows an application that can run for more than one context (e.g., a Transaction program that processes many independent incoming requests) to inform the system which context the application is currently working on.

Along with these functions, a `START_NEW_CONTEXT` function is needed for the system scheduler and for applications that want to start new work that is independent of other work they are processing.

3.2.2 Examples of Styles

Example 1 - Thread handles multiple contexts

Table 1 gives Example 1 in which a single thread at the server switches context to handle requests from two clients.

Example 2-Threads Handle a Transaction Using Different/ Same Contexts

Example 2 given in Table 2 deals with one client and one server. The server has prestarted a number of threads. A database client/server application can use this design to optimize the performance of an application that opens more than one cursor, and fetches data from each cursor in an order that is not known in advance. A separate connection is used for each cursor to allow a fetch to be done on each open cursor independent of the data flowing on other connections. The server hands off work to the threads which are already initialized, and are waiting for incoming work. The waiting threads are not involved in work for any context. Once the thread is given the context and the connection associated with that context, the thread has exclusive use of the connection. In the example, multiple threads work on a single transaction using different contexts. There are two independent connections created at the server. By default, each connection starts a new context.

For performance reasons, it is better to have all the threads use the same context. Even though each connection has its own context, the server can understand, through application-specific logic, that the work is related work. So the server could choose to hand the same context to each of the threads. The example would change as follows: In step 1, the server would issue `SHARE_CONTEXT` instead of `HANDOFF_CONTEXT`. In step 2, the server would simply issue `SHARE_CONTEXT` instead of `START_NEW_CONTEXT` and `HANDOFF_CONTEXT`.

3.3. Interfaces between Context Manager, Transaction Manager and Resource Managers

The TM and protected RMs normally associate a thread with a `TRANID`. `Thread_id` and `TRANID`, along with accounting and security information are components of a "context". In an environment where an application can start a thread to take care of part or all of a transaction, the TM and the protected RMs need to share the current context of the thread that is executing on the protected resource. The mechanism that allows the TM and the RMs on a local system to share a common `context_id` is a matter of implementation. Two design choices we considered are given below.

Design 1: Each time an RM, including the security RM, is invoked, it can query the current context using the CM's `EXTRACT_CURRENT_CONTEXT`. The RM would then have to look up, in its own tables, and determine if it is already involved in work for this current context. If it is not already involved in work for the current context, the RM should add an entry to its internal tables. The entry should include the `context_id`, and whatever other information the particular RM needs. The TM needs to correlate information

Time	Client 1	Server	Client 2
1	BEGIN TRANSACTION request get reply	START_NEW_CONTEXT(C1) BEGIN TRANSACTION begin work for client 1; send reply	
2		START_NEW_CONTEXT(C2) BEGIN TRANSACTION begin work for client 2; send reply	BEGIN TRANSACTION request get reply
3	request get reply	SET_CONTEXT(C1) (server switch to Client 1's context C1) do some work for context 1; send reply	
4		SET_CONTEXT(C2) (server switch to Client 2's context C2) do some work for context 2; send reply	request get reply
5	COMMIT (client 1 commits)	SET_CONTEXT(C1) Change context to C1. Issue COMMIT to TM. TM commits all work associated with current context. (TM queries CM to get current context before processing commit. It will then associate current context with the TRANID and order all RMs to commit work associated with that TRANID and context.)	
6		SET_CONTEXT(C2) change context to C2. Issue COMMIT to TM. TM commits all work associated with current context. (TM interacts with CM as in previous step.)	COMMIT (client 2 commits)

Table 1. Example 1. Thread handles multiple contexts.

Time	Client 1	Server
0		Two threads are prestarted and are waiting to do work on behalf of a context by issuing to the main thread GET_NEW_CONTEXT
1	BEGIN TRANSACTION SQL request OPEN_CURSOR A (one protected connection with the server)	START_NEW_CONTEXT(C1) (to CM) BEGIN TRANSACTION (to TM) HANDOFF_CONTEXT(to thread 1) (This satisfies thread_one's GET_NEW_CONTEXT. Thread_one is now processing SQL request.)
2	SQL request OPEN_CURSOR B (one protected connection with the server) fetch from A fetch from B	START_NEW_CONTEXT(C2) (to CM) HANDOFF_CONTEXT(thread 2) (This satisfies thread_two's GET_NEW_CONTEXT. Thread_two is now processing SQL request).
3	COMMIT (If both connections are protected, the local TM will initiate commit processing on both connections.)	The two server threads will receive the commit message. Each thread issues COMMIT to TM. TM commits all work associated with current context of each thread. TM queries CM to get current context before processing commit. It then associates current context of each thread with TRANID and requests all RMs to commit work associated with that TRANID and context.

Table 2. Example 2. Threads handle a transaction using same/different contexts.

about the TRANID with the context_id. A security manager needs to associate the UserID, and possibly password or other security related tokens, with the context_id. A database RM would need to correlate table update information, lock information, and TRANID with the context_id.

Design 2: The CM provides a broadcast mechanism. With this approach, the CM broadcasts the current context_id to all “interested” RMs. The broadcast would occur each time an OS dispatchable unit (thread or process) changed the context it was working on. The context change would occur because the application has issued a SET_CONTEXT call to the CM. The RMs would do internal housekeeping when notified that the active context has switched from one context_id to another. The switch does not imply that the work has committed or aborted, but rather that a context has temporarily suspended execution. This is akin to X/OPEN’s xa_end(suspend). The mechanism for determining which RM is “interested” in the broadcast is also an implementation issue. There can be either dynamic registration where an RM calls the CM to request to be notified of all changes within a scope, or the CM can support automatic inclusion of RMs based on system definition.

The choice of query versus broadcast is a performance issue. In a particular environment queries might entail less overhead than broadcasts. In the prototype developed for OS/2, a query mechanism was used and the CM kept context_id on a process and thread basis. So each thread could be operating on a different context, or it could be operating on the same context as another thread depending on the application’s style.

The next section describes how context management can be used to solve complex issues in distributed transaction processing that arise due to the way multiple transactions are assigned across processes and threads.

4. Distributed Transaction Management

Each new incoming request accepted by a Transaction Program (TP) is handled by a new instance of the TP. Multiple such TP instances share the same TRANID but may be in the same thread/different threads/ different processes, in the proposed processing paradigm. Each TP instance is a different locus of control and context is necessary to identify it. We show that the TRANID, and thread_id or process_id are not enough to identify the TP instance. If context_id is not used, there is a possibility of deadlocks [9] or protocol violations during commit processing, due to the combination of the new processing paradigm and “loopback”, discussed next. The problems are more obvious in the peer-to-peer communication model where the commit can be initiated by any partner in the transaction tree and more than one transaction is in progress at the same time. The notion of context plays an important role in solving these problems

cleanly. We will discuss the problems in transaction processing without context, and the solutions offered by context management. The use of context management in reconciling communication protocol support between chained and unchained transactions has been presented in [14].

4.1. Loopback

Loopback is a system state in which a transaction reappears at a node that is already involved in the same transaction [7]. Multiple TPs in one commit tree present the same TRANID to the shared resource managers. The “re-infection” can be direct when a client invokes a server that happens to reside on the same node. Loopback can be indirect, when a server, say X, is invoked by a client on a different node, which in turn is a cascaded server for a client on the same node as the server X. Indirect loopback can also occur when two different servers on the same node are invoked as part of the same transaction. Figure 1 illustrates a loopback involving three partner programs X, Y, and Z.

Currently, when a loopback occurs, a process or thread is dispatched to handle the second occurrence of the transaction. The two locii of execution have the same TRANID but can be differentiated by using the process_id or thread_id. In 2PC, the TM needs to determine which resource needs to be sent Prepare [7]. (Note that the Prepare flows correspond to the TP-Prepare service of OSI TP [13], and in the X/Open model, they are triggered by xa_prepare and ax_prepare [16]). A particular resource associated with the locus of control can be identified by process_id (from which TRANID can be deduced) and connection_id (the ‘leg’ identifier, known in OSI TP as the branch_id and in LU6.2 as the conversation correlator). This is sufficient even if the application program has a conversation with another application on its own node whereby both branches of the connection/conversation have the same TRANID and connection_id. In the new processing paradigm, the TM can identify a particular resource associated with the locus of control by context_id, from which TRANID can be deduced, and connection_id. If TRANID or thread_id were used instead of context_id, two locii of control would satisfy (TRANID, connection_id) or (thread_id, connection_id) when the two locii of execution were allocated to the same thread and had a connection with each other. However, when the context_id is used with the connection_id, the context_id uniquely identifies the locus of execution, with respect to which the connection_id is used to determine which resource(s) should be sent Prepare.

The above occurrence of loopback arises in the new paradigm when a peer-to-peer communication model, such as modern SNA (APPN and APPC), is assumed. The peers in the transaction are considered to be ‘loosely coupled’ [16]. In the peer-to-peer model, the commit initiator can

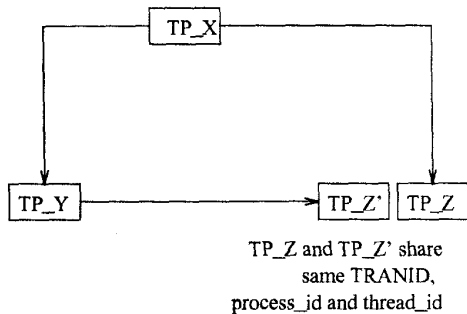


Figure 1. Example of Loopback

be different from the dialog initiator, and may reside on the same system. In this case, the (TRANID, connection_id) pair is not sufficient to identify the dialog on which the Prepare or Backout should be sent. But in the hierarchical communication model, (such as OSI TP [13] and RPC [1], whether blocking or nonblocking), only the dialog initiator can initiate commit processing or send Prepare; a TRANID, which can be deduced from process_id or thread_id, and connection_id are sufficient for the TM to determine the branch on which the CRM should send the Prepare. However, context is still needed for other requirements presented in Section 3.1.

4.2. Protocol Violations and Deadlocks

If the scope of the commit is the process and the TRANID and process_id are used to identify the different resources, protocol violations and deadlock [9] occur. If context management is used and the commit scope is the context (using the context_id), such problems do not occur. Consider the configuration shown in Figure 1. The distributed transaction involves three partner programs, X, Y and Z. An instance of program P is denoted by TP_P. TP_X has invoked servers TP_Y and TP_Z. The arrows on the lines indicate the direction of invocation of TPs. TP_Y has in turn invoked TP_Z (which is really instance TP_Z' that is distinct from TP_Z that was started before it). TP_Z and TP_Z' share the same TRANID, process_id, and thread_id, and hence the TM cannot distinguish between them. The commit tree is not a spanning tree if its nodes are defined by the TRANID, process_id, and thread_id.

Example 1 of Protocol Violations / Deadlocks

(1) TP_X initiates the commit process. Both the conversation resources representing the connection to TP_Y and TP_Z belong to the current scope of TP_X (the scope is the process with the same TRANID) and the Prepare message [6] is sent

to both TP_Y and TP_Z. (2) In turn, the TM of TP_Y sends Prepare to TP_Z'. However, TP_Z and TP_Z' share the same TRANID, process_id, and thread_id, and hence cannot be distinguished by the TM handling them. The TM of TP_Z receives two Prepares for the same transaction, which is interpreted as a protocol violation and results in backing out the transaction.

Using context management, the two incoming requests from client TP_Y and client TP_X result in the creation of two separate contexts (with the same TRANID) for TP_Z and TP_Z'. The commit scope is now the context and the commit tree is now represented by a tree of contexts. This is now a spanning tree. When TP_X issues Commit, the Prepare for partner TP_Z is for the context representing the connection from TP_X to TP_Z. The other Prepare that arrives from TP_Y is for the context that represents the connection from TP_Y to TP_Z'. Based on context, the TM distinguishes between TP_Z and TP_Z' and there are no problems in the 2PC.

Example 2 of Protocol Violations / Deadlocks

(1) TP_X issues Prepare_For_Syncpt [6] on the branch to TP_Y. (2) TP_Y then issues the same call to TP_Z'. (3) The TM cannot distinguish between TP_Z and TP_Z' and in an effort to Prepare that branch of the tree identified by (TP_X, TP_Y) and (TP_Y, TP_Z'), the TM at Z propagates the Prepare along (TP_Z, TP_X). (4) The TM at TP_X will detect a protocol violation in the hierarchical communication model [16], or a deadlock will occur in the peer-to-peer communication model as follows: TP_X will not reply to the Prepare request of TP_Z; TP_Z' (which the TM cannot distinguish from TP_Z) will not reply to the Prepare request of TP_Y; TP_Y will not reply to the Prepare request of TP_X.

However, if context were used, TP_Z' would never have sent a Prepare to TP_X, and no problems would have occurred. Context is essential to keeping a spanning 2PC tree.

Example 3 of Protocol Violations / Deadlocks

The 2PC protocol has been extensively optimized by reducing the number of messages and force log writes [7, 15]. One such well-known optimization is the linear commit, otherwise known as the Last-Agent (LA) optimization. When a partner that initiates the 2PC protocol decides to use the LA optimization, it first chooses the agent that will act as the last agent, then prepares all the other subordinates and finally passes control to the LA (by sending it the YES vote).

(1) Let TP_X choose partner TP_Z as the last agent. When partner TP_X initiates the commit processing, the TM of TP_X will send Prepare only to TP_Y. (2) TP_Y will in turn send Prepare to TP_Z'. (3) The TM handling TP_Z' cannot distinguish between it and TP_Z if it does not use context. So it attempts to send Prepare to TP_X along (TP_Z, TP_X)

because it does not know TP_Z has been chosen as the LA of TP_X.

If the conversation between TP_X and TP_Z is half-duplex, TP_X has send control and TP_Z waits to receive send control in order to send Prepare to TP_X. This wait is indefinite because TP_X is blocked waiting for the response to the initial Prepare, TP_Y is blocked waiting for a response from TP_Z', and TP_Z, (which the TM at Z cannot distinguish from TP_Z') is blocked waiting for send control from TP_X. If the conversation between TP_X and TP_Z is full-duplex and a peer-to-peer model is used, TP_Z sends the Prepare to TP_X and waits for a reply from TP_X. But TP_X cannot receive the Prepare because its TM is blocked waiting for the response to the initial Prepare, TP_Y is blocked waiting for a response from TP_Z', and TP_Z, (which the TM at Z cannot distinguish from TP_Z') is blocked waiting for a response from TP_X. Thus, there is deadlock. If a hierarchical model is used and conversations are full-duplex, a protocol violation is detected by the TM of TP_X. None of these problems would arise if context were used because the TM of TP_Z' would not send a Prepare to TP_X.

5. Conclusions

Operating systems that support threads within a process need the notion of context to efficiently support the paradigm where a single thread can be concurrently associated with several transactions or where several threads work on the same transaction. The first contribution of this paper was that it defined context management by defining a Context Manager and the primitives in its associated user interface. Systems that do not support threads, but support server processes can take advantage of the context management services. The context management services permit a transaction processing application to specify which work is to be handled at any instant. The application using these services can then divide the work within and among the threads or processes, and be assured the resource managers will know which transaction the work belongs to. While context management fits naturally in the peer-to-peer transactional paradigm, it also allows legacy, process-oriented systems to increase the transactional throughput by allowing multiplexing of transactions within one process.

A second contribution of this paper is that it showed how context management is necessary to solve problems such as deadlocks and protocol violations, and to handle loopback situations that arise in distributed transaction processing when using the paradigm in which a thread could be associated with multiple transactions. The paper then gave solutions to these problems in transaction processing. We expect context management will become increasingly important to efficiently handle higher tasking demands and real-time demands on the operating system. We have im-

plemented a prototype of the Context Manager and its user interface for use with transaction processing. The VM operating system Shared File System has been enhanced to provide several features of context management support described here.

References

- [1] Ananda, A. L., Tay, B. H., Koh, E. K., A Survey of Asynchronous RPC, *ACM Operating Systems Review*, 1992.
- [2] Citron, A., Context Manager, *Proc. 4th Int. Workshop on High Performance Transaction Systems*, Asilomar, September 1991.
- [3] Comer, D., *Internetworking with TCP/IP: Principles, Protocols and Architecture*, Prentice Hall, Englewood Cliffs, N.J. ISBN 0-13-470154-2, 1988.
- [4] Gray, J.N., Notes on Data Base Operating Systems, In *Operating Systems - An Advanced Course*, R. Bayer, R. Graham, and G. Seegmuller (Eds.), LNCS, Vol. 60, Springer-Verlag, 1978.
- [5] OS/2 2.0 Technical Library, Control Point Programmers' Reference, IBM, 1986, 1991.
- [6] Systems Network Architecture : Transaction Programmers' Reference Manual for LU Type 6.2, Document Number SC30-3084-5, IBM, June 1993.
- [7] Systems Network Architecture Sync Point Services Architecture Reference, Document Number SC31-8134-0, IBM, Sept. 1994.
- [8] King A., *Inside Windows 1995*, Microsoft Press, ISBN 1-55615-626-x, 1994.
- [9] Kshemkalyani, A.D., Singhal, M., On Characterization and Correctness of Distributed Deadlock Detection, *Jour. of Parallel and Distrib. Computing*, 22(1), 44-59, July 1994.
- [10] Lampson, B.W., Atomic Transactions, In "Distributed Systems: Architecture and Implementation - An Advanced Course", B.W. Lampson (Ed.), LNCS, Vol. 105, Springer-Verlag, 246-265, 1981.
- [11] Mohan, C., Lindsay, B., Obermarck, R., Transaction Management in the R* Distributed Data Base Management System, *ACM Trans. on Database Syst.*, 11(4), Dec. 1986.
- [12] Moskowitz, D., Kerr, D., et al., OS/2 2.1 Unleashed, Sam's Publishing (Prentice-Hall), ISBN 0-672-30240-3, 1993.
- [13] Information Technology - Open Systems Interconnection - Distributed Transaction Processing - Part 1: OSI TP Model; Part 2: OSI TP Service, ISO/IEC JTC 1/SC 21 N, April 1992.
- [14] Samaras, G., Citron A., Kshemkalyani, A. D., Unchained Transactions and SNA's LU6.2, *Proc. 5th Intern. Workshop on High-Performance Transaction Systems*, 28.1-28.19, Asilomar, Sept. 1993.
- [15] Samaras, G., K. Britton, Citron A., C. Mohan, Two-Phase Commit Optimizations in a Commercial Distributed Environment, *Jour. of Distrib. and Par. Databases*, 3(4), 325-360, Oct. 1995.
- [16] Distributed TP: a) The TX Specification P209, b) The XA Specification C193 6/91, c) The XA+ Specification S201, X/Open Consortium, Nov. 1992, Feb. 1992, Apr. 1993.