

# Safe Open-Nested Transactions Through Ownership

Kunal Agrawal I-Ting Angelina Lee Jim Sukha

MIT Computer Science and Artificial Intelligence Laboratory  
{kunal\_ag, angelee, sukha}@mit.edu

## ABSTRACT

Researchers in transactional memory (TM) have proposed open nesting as a methodology for increasing the concurrency of transactional programs. The idea is to ignore “low-level” memory operations of an open-nested transaction when detecting conflicts for its parent transaction, and instead perform abstract concurrency control for the “high-level” operation that the nested transaction represents. To support this methodology, TM systems use an open-nested commit mechanism that commits all changes performed by an open-nested transaction directly to memory, thereby avoiding low-level conflicts. Unfortunately, because the TM runtime is unaware of the different levels of memory, unconstrained use of open-nested commits can lead to anomalous program behavior.

We describe the framework of *ownership-aware transactional memory* which incorporates the notion of modules into the TM system and requires that transactions and data be associated with specific *transactional modules* or Xmodules. We propose a new *ownership-aware commit mechanism*, a hybrid between an open-nested and closed-nested commit which commits a piece of data differently depending on which Xmodule owns the data. Moreover, we provide a set of precise constraints on interactions and sharing of data among the Xmodules based on familiar notions of abstraction. The ownership-aware commit mechanism and these restrictions on Xmodules allow us to prove that ownership-aware TM has clean memory-level semantics. In particular, it guarantees *serializability by modules*, an adaptation of the definition of multilevel serializability from database systems. In addition, we describe how a programmer can specify Xmodules and ownership in a Java-like language. Our type system can enforce most of the constraints required by ownership-aware TM statically, and can enforce the remaining constraints dynamically. Finally, we prove that if transactions in the process of aborting obey restrictions on their memory footprint, then ownership-aware TM is free from *semantic deadlock*.

**Categories and Subject Descriptors** D.2.1 [Software Engineering]: Requirements/Specifications — Methodologies; D.3.3 [Pro-

---

This research was supported in part by NSF Grants CNS-0615215 and CNS-0540248 and a grant from Intel corporation.

A preliminary version of this paper appeared as a poster at *PPoPP* 2008 and as a brief announcement at *SPAA* 2008.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

*PPoPP'09*, February 14–18, 2009, Raleigh, North Carolina, USA.  
Copyright © 2009 ACM 978-1-60558-397-6/09/02...\$5.00.

*gramming Languages*]: Language Constructs and Features — Concurrent programming structures

**General Terms** Design, Languages

**Keywords** Abstract Serializability, Open-nested Transactions, Ownership-aware Transactions, Ownership Types, Safe Nesting, Semantic Deadlock, Serializability by Modules, Transactional Memory, Transactional Memory Semantics, XModules.

## 1. INTRODUCTION

Transactional memory (TM) [6] is meant to simplify concurrency control in parallel programming by providing a transactional interface for accessing memory; the programmer simply encloses the critical region inside an atomic block, and the TM system ensures that this section of code executes atomically. A TM system enforces atomicity by tracking memory locations accessed by transactions (using *read sets* and *write sets*), finding transactional conflicts, and aborting transactions that conflict. TM guarantees that transactions are *serializable* [12]; that is, transactions affect global memory as if they were executed one at a time in some order, even if in reality, several executed concurrently.

Transactions may be *nested*. If a transaction  $Y$  is *closed nested* [8] inside another transaction  $X$ , then for the purpose of detecting conflicts, the TM system considers any memory locations accessed by  $Y$  as conceptually also being accessed by its parent  $X$ . Thus, when  $Y$  commits, the TM system merges  $Y$ 's read and write sets into the read and write sets of  $X$ . TM with closed-nested transactions guarantees that transactions are serializable at the level of memory. Researchers have observed, however, that closed nesting might unnecessarily restrict concurrency in programs because it does not allow two “high-level” transactions to ignore conflicts due to “low-level” memory accessed by nested transactions.

Researchers have proposed the methodology of *open-nested transactions* to increase concurrency in transactional programs by breaking serializability at the memory level. The open-nesting methodology incorporates the *open-nested commit mechanism* [7, 10]. When an open-nested transaction  $Y$  (nested inside transaction  $X$ ) commits,  $Y$ 's changes are committed to memory and  $Y$ 's read and write sets are discarded. Thus, the TM system no longer detects conflicts with  $X$  due to memory accessed by  $Y$ . In this methodology, the programmer considers  $Y$ 's internal memory operations to be at a “lower level” than  $X$ ; thus  $X$  should not care about the memory accessed by  $Y$  when checking for conflicts. Instead,  $Y$  must acquire an *abstract lock* based on the high-level operation that  $Y$  represents and propagate this lock to  $X$ , so that the TM system can perform concurrency control at an abstract level. Also, if  $X$  aborts, it may need to execute *compensating actions* to undo the effect of its committed open-nested subtransaction  $Y$ . Moss [9] illustrates the use of open nesting with an application that employs a B-tree. Ni et al. [11] describe a software TM system that supports the open-nesting methodology.

An unconstrained use of the open-nested commit mechanism can lead to anomalous program behavior that can be tricky to reason about [2]. We believe that one reason for the apparent complexity of open nesting is that the mechanism and the methodology make different assumptions about memory. Consider a transaction  $Y$  open nested inside transaction  $X$ . The open-nesting methodology requires that  $X$  ignore the “lower-level” memory conflicts generated by  $Y$ , while the open-nested commit mechanism will ignore *all* the memory operations inside  $Y$ . Say  $Y$  accesses two memory locations  $\ell_1$  and  $\ell_2$ , and  $X$  does not care about changes made to  $\ell_2$ , but does care about  $\ell_1$ . The TM system cannot distinguish between these two accesses, and will commit both in an open-nested manner, leading to anomalous behavior.

Researchers *have* demonstrated specific examples [4, 11] that safely use an open-nested commit mechanism. These examples work, however, because the inner (open) transactions never write to any data that is accessed by the outer transactions. Moreover, since these examples require only two levels of nesting, it is not obvious how one can correctly use open-nested commits in a program with more than two levels of abstraction. The literature on TM offers relatively little in the way of formal programming guidelines which one can follow to have *provable* guarantees of safety when using open-nested commits.

### Contributions

In this paper, we bridge the gap between memory-level mechanisms for open nesting and the high-level view by explicitly integrating the notions of *transactional modules* (Xmodules) and *ownership* into the TM system. We believe the *ownership-aware TM system* allows the programmer to safely use the methodology of open nesting, because the runtime’s behavior more closely reflects the programmer’s intent. In addition, the structure imposed by ownership allows a compiler and runtime to enforce properties needed to provide provable guarantees of “safety” to the programmer. More specifically, the contributions of this paper are as follows:

1. We suggest a concrete set of guidelines for sharing of data and interactions between Xmodules.
2. We describe how the Xmodules and ownership can be specified in a Java-like language and propose a type system that enforces most of the above-mentioned guidelines in the programs written using this language extension.
3. We formally describe the operational model for ownership-aware TM, called the *OAT* model, which uses a new *ownership-aware commit mechanism*. The ownership-aware commit mechanism is a compromise between an open-nested and a closed-nested commit; when a transaction  $T$  commits, access to a memory location  $\ell$  is committed globally if  $\ell$  belongs to the same Xmodule as  $T$ ; otherwise, the access to  $\ell$  is propagated to  $T$ ’s parent transaction. Unlike an ordinary open-nested commit, the ownership-aware commit treats memory locations differently depending on which Xmodule owns the location. Note that the ownership-aware commit is still a mechanism; programmers must still use it in combination with abstract locks and compensating actions to implement the full methodology.
4. We prove that if a program follows the proposed guidelines for Xmodules, then the *OAT* model guarantees serializability by modules, which is a generalization of “serializability by levels” used in database transactions. Ownership-aware commit is the same as open-nested commit if no Xmodule ever accesses data belonging to other Xmodules. Thus, one corollary of our theorem is that open-nested transactions are serializable when Xmodules do not share data. This observation explains why researchers [4, 11] have found it natural to use open-nested

transactions in the absence of sharing, in spite of the apparent semantic pitfalls.

5. We prove that under certain restricted conditions, a computation executing under the *OAT* model can not enter a semantic deadlock.

In later sections, we distinguish between the variations of nested transactions as follows. We say that a transaction  $Y$  is *vanilla open nested* when referring to a TM system which performs the open-nested commit of  $Y$ . We say that  $Y$  is *safe nested* when referring to the ownership-aware TM system which performs the ownership-aware commit of  $Y$ . Finally, we say that a transaction  $Y$  is an open-nested transaction when we are referring to the abstract methodology, rather than a particular implementation with a specific commit mechanism.

### Outline

The paper is organized as follows. In Section 2 we present an overview of ownership-aware TM and highlight key features using an example application. Section 3 describes language constructs for specifying Xmodules and ownership. In Section 4, we review the transactional computation framework [2], and extend this framework to formally incorporate Xmodules and ownership. Section 5 describes the *OAT* model, and Section 6 gives a formal definition of serializability by modules, and a proof sketch that the *OAT* model guarantees this definition. Section 7 provides conditions under which the *OAT* model does not exhibit semantic deadlocks. Section 8 concludes with a discussion of some related work.

## 2. OWNERSHIP-AWARE TRANSACTIONS

In this section, we give an overview of ownership-aware TM. To motivate the need for the concept of ownership in TM, we first present an example application which might benefit from open nesting. We then introduce the notion of an Xmodule and informally explain the programming guidelines when using Xmodules. Finally, we highlight some of the key differences between ownership-aware TM and a TM with vanilla open nesting. In this section, we present the intuitive descriptions of the concepts in ownership-aware TM; we defer formal definitions until later sections.

### Example Application

We describe an example application for which one might use open-nested transactions. This example is similar to the one described by Moss [9], but it includes data sharing between nested transactions and their parents, and has more than two levels of nesting.

Since the open-nesting methodology is designed for programs that have multiple levels of abstraction, we choose a modular application. Consider a user application which concurrently accesses a database of many individuals’ book collections. The database stores records in a binary search tree, keyed by name. Each node in the binary search tree corresponds to a person, and stores a list of books in his/her collection. The database supports queries by name, as well as updates that add a new person or a new book to a person’s collection. The database also maintains a private hashmap, keyed by book title, to support a reverse query; given a book title, it returns a list of people who own the book. Finally, the user application wants the database to log changes on disk for recoverability. Whenever the database is updated, it inserts metadata into the buffer of a logger to record the change that just took place. Periodically, the user application is able to request a checkpoint operation which flushes the buffer to disk.

This application is modular, with five natural modules — the user application (`UserApp`), the database (`DB`), the binary search

tree (BST), the hashtable (Hashtable), and the logger (Logger). The `UserApp` module calls methods from the `DB` module when it wants to insert into the database, or query the database. The database in turn maintains internal metadata and calls the `BST` module and the `Hashtable` module to answer queries and insert data. Both user application and the database may call methods from the `Logger` module.

If the modules use open-nested transactions, a TM system with vanilla open-nested commits can result in non-intuitive outcomes. Consider the example where a transactional method  $A$  from the `UserApp` module tries to insert a book  $b$  into the database, and the insert is an open-nested transaction. The method  $A$  (which corresponds to transaction  $X$ ) calls an insert method in the `DB` module and passes  $b$  (the `Book` object) to be inserted. This insert method generates an open-nested transaction  $Y$ . Suppose  $Y$  writes to some field of the book  $b$  (memory location  $\ell_1$ ), and also writes some internal database metadata (location  $\ell_2$ ). After a vanilla open-nested commit of  $Y$ , the modifications to both  $\ell_1$  and  $\ell_2$  become visible globally. Assuming the `UserApp` does not care about the internal state of the database, committing the internal state of the `DB` ( $\ell_2$ ) is a desirable effect of open nesting; this commit increases concurrency, because other transactions can potentially modify the database in parallel with  $X$  without generating a conflict. The `UserApp` does, however, care about changes to the book  $b$ ; thus, the commit of  $\ell_1$  breaks the atomicity of transaction  $X$ . A transaction  $Z$  in parallel with transaction  $X$  can access this location  $\ell_1$  after  $Y$  commits, before the outer transaction  $X$  commits.<sup>1</sup> To increase concurrency, we want the method from `DB` to commit changes to its own internal data; we do not, however, want it to commit the data that `UserApp` cares about.

To enforce this kind of restriction, we need some notion of **ownership of data**: if the TM system is aware of the fact that the book object “belongs” to the `UserApp`, then it can decide not to commit `DB`’s change to the book object globally. For this purpose, we introduce the notion of **transactional modules**, or **Xmodules**. When a programmer explicitly defines Xmodules and specifies the ownership of data, the TM system can make the correct judgement about which data to commit globally.

### *Xmodules and the Ownership-Aware Commit Mechanism*

The ownership-aware TM system requires that programs be organized into Xmodules. Intuitively, an Xmodule  $M$  is as a stand-alone entity that contains data and transactional methods; an Xmodule owns data that it privately manages, and uses its methods to provide public services to other Xmodules. During program execution, a call to a method from an Xmodule  $M$  generates a transaction instance (e.g.,  $X$ ). If this method in turn calls another method from an Xmodule  $N$ , an additional transaction  $Y$ , safe nested inside  $X$ , is created only if  $M \neq N$ . Therefore, defining an Xmodule automatically specifies safe-nested transactions.

In the ownership-aware TM system, every memory location is owned by exactly one Xmodule. If a memory location  $\ell$  is in a transaction  $T$ ’s read or write set, the ownership-aware commit of a transaction  $T$  commits this access globally only if  $T$  is generated by the same Xmodule that owns  $\ell$ ; in this case, we say that  $T$  is “responsible” for that access to  $\ell$ . Otherwise, the read or write to  $\ell$  is propagated up to the read or write set of  $T$ ’s parent transaction; that is, the TM system behaves as though  $T$  was a closed-nested transaction with respect to location  $\ell$ .

We wish to guarantee that ownership-aware TM behaves “nicely”. For example, in the TM system, some transaction must be “respon-

sible” for committing every memory access. Similarly, the TM system should guarantee some form of serializability. To guarantee these properties, we must restrict interactions between Xmodules; if Xmodules could arbitrarily call methods from or access memory owned by other Xmodules, then these properties might not be satisfied.

### *Rules for Xmodules*

Ownership-aware TM uses Xmodules to control both the structure of nested transactions, and the sharing of data between Xmodules (i.e., to limit which memory locations a transaction instance can access). In our system, Xmodules are arranged as a **module tree**, denoted as  $\mathcal{D}$ . In  $\mathcal{D}$ , an Xmodule  $N$  is a child of  $M$  if  $N$  is “encapsulated by”  $M$ . The root of  $\mathcal{D}$  is a special Xmodule called `world`. Each Xmodule is assigned an `xid` by visiting the nodes of  $\mathcal{D}$  in a left-to-right depth-first search order, and assigning ids in increasing order, starting with `xid(world) = 0`. Therefore `world` has the minimum `xid`, and “lower-level” Xmodules have larger `xid` numbers.

**DEFINITION 1.** *We impose two rules on Xmodules based on the module tree:*

1. **Rule 1:** *A method of an Xmodule  $M$  can access a memory location  $\ell$  directly only if  $\ell$  is owned by either  $M$  or an ancestor of  $M$  in the module tree. This rule means that an ancestor Xmodule  $N$  of  $M$  may pass data down to a method belonging to  $M$ , but a transaction from module  $M$  cannot directly access any “lower-level” memory.*
2. **Rule 2:** *A method from  $M$  can call a method from  $N$  only if  $N$  is the child of some ancestor of  $M$ , and `xid(N) > xid(M)` (i.e., if  $N$  is “to the right” of  $M$  in the module tree). This rule requires that an Xmodule can call methods of some (but not all) lower-level Xmodules.<sup>2</sup>*

The intuition behind these rules is as follows. Xmodules have methods to provide services to other higher-level Xmodules, and Xmodules maintain their own data in order to provide these services. Therefore, a higher-level Xmodule can pass its data to a lower-level Xmodule and ask for services. A higher-level Xmodule should not directly access the internal data belonging to a lower-level Xmodule.

If Xmodules satisfy Rules 1 and 2, TM can have a well-defined ownership-aware commit mechanism; some transaction is always “responsible” for every memory access (proved in Section 5). In addition, these rules and the ownership-aware commit mechanism guarantee that transactions satisfy the property of “serializability by modules” (proved in Section 6).

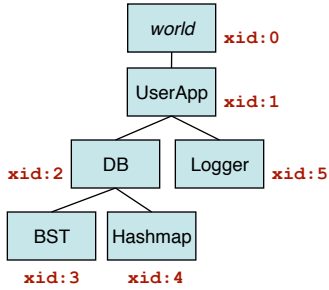
One potential limitation of ownership-aware TM is that some “cyclic dependencies” between Xmodules are prohibited. The ability to define one module as being at a lower level than another is fundamental to the open-nesting methodology. Thus, our formalism requires that Xmodules be partially ordered; if an Xmodule  $M$  can call Xmodule  $N$ , then conceptually  $M$  is at a higher level than  $N$  (i.e., `xid(M) < xid(N)`), and thus  $N$  cannot call  $M$ . If two components of the program call each other, then, conceptually, neither of these components is at a higher-level than the other, and we would require that these two components be combined into one Xmodule.

### *Xmodules in the Example Application*

Consider a Java implementation of the example application described earlier. It may have the following classes: `UserApp` as the top-level application that manages the book collections, `Person`

<sup>1</sup>Note that abstract locks [9] do not address this problem. Abstract locks are meant to disallow other transactions from noticing the fact that the book was inserted into the `DB`. They do not usually protect the individual fields of the book object itself.

<sup>2</sup>An Xmodule can, in fact, call methods within its own Xmodule or from its ancestor Xmodules, but we model these calls differently. We explain these cases at the end of this section.



**Figure 1.** A module tree  $\mathcal{D}$  for the program described in Section 2. The  $\text{xid}$ 's are assigned according to a left-to-right depth-first tree walk, numbering Xmodules in increasing order, starting with  $\text{xid}(\text{world}) = 0$ .

and `Book` as the abstractions representing book owners and books, `DB` for the database, `BST` and `HashMap` for the binary search tree and hashmap maintained by the database, and `Logger` for logging the metadata to disk. In addition, there are some other auxiliary classes: tree node `BSTNode` for the `BST`, `Bucket` in the `HashMap`, and `Buffer` used by the `Logger`.

For ownership-aware TM, not all of a program's classes are meant to be Xmodules; some classes only wrap data. In our example, we identified five Xmodules—`UserApp`, `DB`, `BST`, `HashMap`, and `Logger`; these classes are stand-alone entities which have encapsulated data and methods. Classes such as `Book` and `Person`, on the other hand, are data types used by `UserApp`. Similarly, classes like `BSTNode` and `Bucket` are data types used by `BST` and `HashMap` to maintain their internal state.

We organize the Xmodules of the application into the module tree shown in Figure 1. `UserApp` is encapsulated by `world`, `DB` and `Logger` are encapsulated under `UserApp`; `BST` and `HashMap` are encapsulated under `DB`. By dividing Xmodules this way, the ownership of data falls out naturally, i.e., an Xmodule owns certain pieces of data if the data is encapsulated under the Xmodule. For example, the instances of `Person` or `Book` are owned by `UserApp` because they should only be accessed by either `UserApp` or its descendants.

Let us consider the implications of Definition 1 for the example. Due to Rule 1, all of `DB`, `BST`, `HashMap`, and `Logger` can directly access data owned by `UserApp`, but the `UserApp` cannot directly access data owned by any of the other Xmodules. This rule corresponds to standard software-engineering rules for abstraction; the “high-level” Xmodule `UserApp` should be able to pass its data down, allowing lower-level Xmodules to access that data directly, but `UserApp` itself should not be able to directly access data owned by lower-level Xmodules. Due to Rule 2, the `UserApp` may invoke methods from `DB`, `DB` may invoke methods from `BST` and `HashMap`, and every other Xmodule may invoke methods from `Logger`. Thus, Rule 2 allows all the operations required by the example application. As expected, the `UserApp` can call the insert and search methods from the `DB` and can even pass its data to the `DB` for insertion. More importantly, notice the relationship between `BST` and `Logger`. The `BST` Xmodule can call methods from `Logger`, but the `BST` cannot pass data it owns directly into the `Logger`. It can, however, pass data owned by the `UserApp` to the `logger`, which is all this application requires.

### Advantage of Ownership-Aware Transactions

One of the major problems with vanilla open nesting is that some transactions can see inconsistent data. Say a transaction  $Y$  is open nested inside transaction  $X$ . Let  $v_0$  be the initial value of location

$\ell$ , and suppose  $Y$  writes value  $v_1$  to location  $\ell$  and then commits. Now a transaction  $Z$  in parallel with  $X$  can read this location  $\ell$ , write value  $v_2$  to  $\ell$ , and commit, all before  $X$  commits. Therefore,  $X$  can now read this location  $\ell$  and see the value  $v_2$ , which is neither the initial value  $v_0$  (the value of  $\ell$  when  $X$  started), nor  $v_1$  which was written by  $X$ 's inner transaction,  $Y$ . This behavior might seem counterintuitive.

Now consider the same example for ownership-aware transactions. Say  $X$  is generated by a method of Xmodule  $M$  and  $Y$  is generated by a method of Xmodule  $N$ . If  $N$  owns  $\ell$ ,  $X$  cannot access  $\ell$ , since  $\text{xid}(M) < \text{xid}(N)$  (by Definition 1, Rule 2), and no transaction from a higher-level module can access data owned by a lower-level module (by Definition 1, Rule 1). Thus, the problem does not arise. If  $N$  does not own  $\ell$ , the ownership-aware commit of  $Y$  will not commit the changes to  $\ell$  globally and  $\ell$  will be propagated to  $X$ 's write set. Therefore, if  $Z$  tries to access  $\ell$  before  $X$  commits, the TM system will detect a conflict. Thus  $X$  cannot see an inconsistent value for  $\ell$ .<sup>3</sup>

### Callbacks

At first glance, the assumptions we have made regarding methods of Xmodules seem somewhat restrictive. In the description thus far, we prohibit an Xmodule  $M$  from calling another transactional method from  $M$  or a proper ancestor of  $M$ . In particular, it appears as though our model disallows callbacks. Our model, however, does permit both these cases; we simply model these calls differently.

If a method  $X$  from Xmodule  $M$  calls another method  $Y$  from an ancestor Xmodule  $N$ , this call does not generate a new safe-nested transaction instance. Instead,  $Y$  is subsumed in  $X$  using flat (or closed) nesting. Recall that Rule 1 in Definition 1 allows  $X$  to access data belonging to  $N$  or any of its ancestors directly. Thus, we can treat any data access by a flat (or closed) nested transaction  $Y$  as being accessed by  $X$  directly, provided that  $Y$  and its nested transactions access only memory belonging to  $N$  or  $N$ 's ancestors. We say that  $Y$  is a *proper callback* method for Xmodule  $N$  if its nested calls are all proper callback methods belonging to Xmodules which are ancestors of  $N$ . In our formal model in Section 4, we assume that we only have proper callbacks and model them as direct memory accesses, allowing us to ignore them in the formal definitions.

### Closed-Nested Transactions

In our model, every method call that crosses an Xmodule boundary automatically generates a safe-nested transaction. Ownership-aware TM can effectively provide closed-nested transactions, however, with appropriate specifications of ownership. If an Xmodule  $M$  owns no memory, but only operates on memory belonging to its proper ancestors, then transactions of  $M$  will effectively be closed nested. In the limit, if the programmer specifies that all memory is owned by the `world` Xmodule, then all changes in any transaction's read or write set are propagated upwards; thus all ownership-aware commits behave exactly as closed-nested commits.

## 3. OWNERSHIP TYPES FOR XMODULES

When using ownership-aware transactions, the Xmodules and data ownership in a program must be specified for two reasons. First, the ownership-aware commit mechanism depends on these concepts. Second, we can guarantee some notion of serializability only if a program has Xmodules which conform to the rules in Definition 1. In this section, we describe language constructs and a type system that can be used to specify Xmodules and ownership in a Java-like

<sup>3</sup> For simplicity, we have described the case where  $Y$  is directly nested inside  $X$ . The case where  $Y$  is more deeply open nested inside  $X$  behaves in a similar fashion.

language. Our type system — the *OAT type system* — statically enforces some of the restrictions described in Definition 1.

The *OAT* type system extends the ownership types of Boyapati et al. [3], which is described first in this section. We then describe extensions to this type system to enforce some of the restrictions in Definition 1. Next, we present code for parts of the example application described in Section 2. Finally, we discuss some restrictions required by Definition 1 which the *OAT* type system does not enforce statically. The type system’s annotations, however, enable dynamic checks for these restrictions.

### Boyapati et al.’s Parametric Ownership Type System

The type system of Boyapati et al. provides a mechanism for specifying ownership of objects. The type system enforces the properties stated in Lemma 1.

**LEMMA 1.** *Boyapati et al.’s type system enforces the following properties:*

1. Every object has a unique owner.
2. The owner can be either another object, or `world`.
3. The ownership relation forms an ownership tree (of objects) rooted at `world`.
4. The owner of an object does not change over time.
5. An object  $a$  can access another object  $b$  directly only if  $b$ ’s owner is either  $a$ , or one of  $a$ ’s proper ancestors in the ownership tree.

Boyapati et al.’s type system requires ownership annotations to class definitions and type declarations to guarantee Lemma 1. Every class type  $T1$  has a set of associated ownership tags, denoted  $T1\langle f_1, f_2, \dots, f_n \rangle$ . The first formal  $f_1$  denotes the owner of the current instance of the object (i.e., `this` object). The remaining formals  $f_2, f_3, \dots, f_n$  are additional tags which can be used to instantiate and declare other objects within the class definition. The formals get assigned with actual owners  $o_1, o_2, \dots, o_n$  when an object  $a$  of type  $T1$  is instantiated. By parameterizing class and method declarations with ownership tags, this type system permits owner polymorphism. Thus, one can define a class type (e.g. a generic hash table) once, but instantiate multiple instances of that class with different owners in different parts of the program.

The type system enforces the properties in Lemma 1 by performing the following checks:

1. Within the class definition of type  $T1$ , only the tags  $\{f_1, f_2, \dots, f_n\} \cup \{\text{this}, \text{world}\}$  are visible. The `this` ownership tag represents the object itself.
2. A variable  $c_2$  with type  $T2\langle f_2, \dots \rangle$  can be assigned to a variable  $c_1$  with type  $T1\langle f_1, \dots \rangle$  if and only if  $T2$  is a subtype of  $T1$  and  $f_1 = f_2$ .
3. If an object  $a$ ’s tags are instantiated to be  $o_1, o_2, \dots, o_n$  when  $a$  is created, then in the ownership tree,  $o_1$  must be a descendant of  $o_i, \forall i \in 2..n$ , (denoted by  $o_1 \preceq o_i$  henceforth).

Boyapati et al. show that these type checks guarantee the properties of Lemma 1.

In some cases, to enable the type system to perform check 3 locally, the programmer may need to specify a `where` clause in a class declaration. For example, suppose the class declaration of type  $T1$  has formal tags  $\langle f_1, f_2, f_3 \rangle$ , and inside  $T1$ ’s definition, some type  $T2$  object is instantiated with ownership tags  $\langle f_2, f_3 \rangle$ . The type system cannot determine whether or not  $f_2 \preceq f_3$ . To resolve this ambiguity, the programmer must specify `where (f2 <= f3)` at the class declaration of type  $T1$ . When an instance of type  $T2$  object is instantiated, the type system then checks that the `where` clause is satisfied.

### The OAT Type System

The ownership tree described by Boyapati et al. exhibits some of the same properties as the module tree we described in Section 2; however, this ownership scheme does not enforce two major requirements of our system.

- In [3], any object can own other objects. Our rules, however, require that only Xmodules own other objects.
- In [3], an object can call any of its ancestor’s siblings. Our rules (namely Definition 1), however, dictate that an Xmodule  $M$  can only call its ancestor’s siblings to the right.

With these requirements in mind, we extend Boyapati et al.’s type system to create the *OAT* type system.

The extensions to handle the first requirement are straightforward. The *OAT* type system explicitly distinguishes objects and Xmodules by requiring that Xmodules extend from a special Xmodule class. The *OAT* type system only allows classes that extend Xmodule to use `this` as an ownership tag. In the context of the Boyapati et al.’s ownership tree, this restriction creates a tree where all the internal nodes are Xmodules and all leaves are non-Xmodule objects. If we ignore any order imposed on the children of an Xmodule, for ownership-aware TM, the module tree (as described in Section 2) is essentially the ownership tree with all non-Xmodule objects removed.

The second requirement is more complicated to enforce. First, we extend each owner instance  $o$  to have two fields: *name*, represented by  $o.name$ ; and *index*, represented by  $o.index$ . The name field is conceptually the same as an ownership instance in Boyapati et al.’s type system. The index field is added to help the compiler to infer ordering between children of the same Xmodule in the module tree. The *OAT* type system allows the programmer to pass `this[i]` as the ownership tag (i.e., with an index  $i$ ) instead of `this`. Similarly, one can use `world[i]` as an ownership tag. Indices enable the type system to infer an ordering between two sibling Xmodules  $M$  and  $N$ ; for instance, if an Xmodule  $L$  instantiates  $M$  and  $N$  with owners `this[i]` and `this[i+1]`, respectively, then  $M$  appears to the left of  $N$  in the module tree.

Finally, for technical reasons, the *OAT* system prohibits all Xmodules  $M$  from declaring primitive fields. If  $M$  had primitive fields, then by Boyapati et al.’s type system, these fields are owned by the  $M$ ’s parent. Since this property seems counter-intuitive, we opt to disallow primitive fields for Xmodules.

In summary, the *OAT* type system performs these checks:

1. Within the class definition of type  $T1$ , only the tags  $\{f_1, f_2, \dots, f_n\} \cup \{\text{this}, \text{world}\}$  are visible.
2. A variable  $c_2$  with type  $T2\langle f_2, \dots \rangle$  can be assigned to a variable  $c_1$  with type  $T1\langle f_1, \dots \rangle$  if and only if  $T2 = T1$ , and all the formals are initialized to the same owners with the same indices, if indices are specified.
3. For a type  $T\langle o_1, o_2, \dots, o_n \rangle$ , we must have, for all  $i \in \{2, \dots, n\}$ , either  $o_1.name < o_i.name$  or  $o_1.name = o_i.name$  and  $o_1.index < o_i.index$  (if both indices are known).<sup>4</sup>
4. The ownership tag `this` can only be used within the definition of a class that extends Xmodule.
5. Xmodule objects cannot have primitive-type fields.

The first three checks are analogous to the checks in Boyapati et al.’s type system. The last two checks are added to enforce the additional requirements of Xmodules.

<sup>4</sup>In the ownership tree, for any Xmodule  $M$ , the *OAT* type system implicitly assigns non-Xmodule children of  $M$  higher indices than the Xmodule children of  $M$ , unless the user specifies otherwise.

```

1 public class UserApp<app0> extends Xmodule<app0> {
2     private Logger<this[1], this[2]> logger;
3     private DB<this[0], this[1], this[2]> db;
4     ...
5     public UserApp() {
6         logger = new Logger<this[1], this[2]>();
7         db = new DB<this[0], this[1], this[2]>(logger);
8     }
9
9 public class DB<db0, log0, data0>
10     extends Xmodule<db0> where (log0 < data0) {
11     private Logger<log0, data0> logger;
12     private BST<this[0], log0, data0> bst;
13     private Hashmap<this[1], log0, data0> hashmap;
14     public DB(Logger<log0, data0> logger) {
15         this.logger = logger;
16         ...
17     }

```

**Figure 2.** Specifying Xmodules and ownership for the example application described in Section 2.

The OAT type system supports where clauses of the form  $where (f_i < f_j)$ ; when  $f_i$  and  $f_j$  are instantiated with  $o_i$  and  $o_j$ , the type system ensures that either  $o_i.name < o_j.name$ , or  $o_i.name = o_j.name$  and  $o_i.index < o_j.index$ . The detailed type rules for the OAT type system are described in [1].

### Example Application Using the OAT Type System

Figure 2 illustrates how one can specify Xmodules and ownership using ownership types. The programmer specifies an Xmodule by creating a class which extends from a special Xmodule class. The DB class has three formal owner tags – `db0` which is the owner of the DB Xmodule instance, `log0` which is the owner of the Logger Xmodule instance that the DB Xmodule will use, and `data0` which is the owner of the user data being stored in the database. When an instance of `UserApp` initializes Xmodules in lines 5–6, it declares itself as the owner of the Logger, the DB, and the user data being passed into DB. The indices on `this` are declaring the ordering of Xmodules in the module tree, i.e., the user data is lower-level than the Logger, and the Logger is lower level than the DB. lines 11–13 illustrate how the DB class can initialize its Xmodules and propagate the formal owner tags (i.e., `log0` and `data0`) down.

Note that in order for this code to type check, the DB class must declare `log0 < data0` using the `where` clause in line 10, otherwise the type check would fail at line 11, due to ambiguity of their relation in the module tree. The `where` clause in line 10 is checked whenever an instance of DB is created, i.e. at line 6.

### The OAT Type System’s Guarantees

The following lemma about the OAT type system can be proved in a reasonably straightforward manner using Lemma 1.

**LEMMA 2.** *The OAT type system guarantees the following properties.*

1. An Xmodule  $M$  can access a (non-Xmodule) object  $b$  with ownership tag  $o_b$  only if  $M \preceq o_b.name$ .
2. An Xmodule  $M$  can call a method in another Xmodule  $N$  with owner  $o_N$  only if one of the following is true:
  - (a)  $M = o_N.name$  (i.e.  $M$  owns  $N$ );
  - (b) The least common ancestor of  $M$  and  $N$  in the module tree is  $o_N.name$ .
  - (c)  $N \succeq M$  (i.e.  $N$  is an ancestor of  $M$ ).

Lemma 2 does not, however, guarantee all the properties we want from Xmodules (i.e., Definition 1). In particular, Lemma 2 does not consider any ordering of sibling Xmodules. The OAT type system can, however, provide stronger guarantees for a program which satisfies what we call the *unique owner indices* assumption: for all Xmodules  $M$ , all children of  $M$  in the module tree are instantiated with ownership tags with unique indices that can be statically determined. For such a program, the type system can order the children of every Xmodule  $M$  from smallest to largest index, and assign the `xid` to each Xmodule as described in Section 2. Then, the following result holds:

**THEOREM 3.** *For a program with unique owner indices, in addition to Lemma 2, the OAT type system guarantees that if the least common ancestor of Xmodules  $M$  and  $N$  in the module tree is  $o_N.name$ , then  $M$  can call a method in  $N$  only if  $xid(M) < xid(N)$ .*

### PROOF SKETCH.

We prove (by contradiction) that if the least common ancestor of  $M$  and  $N$  in the module tree is  $o_N.name$ , and  $xid(M) > xid(N)$ , then  $M$  cannot have a formal tag with value  $o_N$ . Therefore, it cannot declare a type with owner tag  $o_N$ , and cannot access  $N$ . We only sketch the proof here. For the details, please see [1].

Let  $L$  be the least common ancestor of  $M$  and  $N$ , let  $Q$  be the ancestor of  $M$  which is  $N$ ’s sibling, and let  $o_Q$  be  $Q$ ’s ownership tag (i.e., the tag with which  $Q$  is instantiated). Since  $N$  and  $Q$  have the same parent (i.e.  $L$ ) in the module tree, we have  $o_N.name = o_Q.name = L$ . Since  $xid(M) > xid(N)$ ,  $M$  is to the right of  $N$  in the ownership tree. Therefore,  $Q$ , which is an ancestor of  $M$ , is to the right of  $N$  in the ownership tree. Therefore, we have  $o_Q.index > o_N.index$ .

Assume for contradiction that  $M$  does have  $o_N$  as one of its tags. Using Lemma 1, one can show that the only way for  $M$  to receive tag  $o_N$  is if  $Q$  also has a formal tag with value  $o_N$ . Thus,  $Q$ ’s first formal owner tag has value  $o_Q$  and another one of its formals has value  $o_N$ . Therefore, the type system fails to type check, either at the point where  $Q$  is instantiated due to  $o_Q.index > o_N.index$  (check 3), or at some other place where a disambiguating where clause is used.  $\square$

Theorem 3 only modifies the Condition 2b of Lemma 2. Therefore, Lemma 2 along with Theorem 3 imposes restrictions on every Xmodule  $M$  which are only slightly weaker than the restrictions required by Definition 1. Condition 1 in Lemma 2 corresponds to Rule 1 of Definition 1. Conditions 2a and 2b are the cases permitted by Rule 2. Condition 2c, however, corresponds to the special case of callbacks or calling a method from the same Xmodule, which is not permitted by Definition 1. This case is modeled differently, as we explained in Section 2.

The OAT type system is a best-effort type system to check for the restrictions required by Definition 1. The OAT type system cannot fully guarantee, however, that a type-checked program does not violate Definition 1. Specifically, the OAT type system can not always detect the following violations statically. First, if the program does not have unique owner indices, then  $L$  may instantiate both  $M$  and  $N$  with the same index. Then, by Lemma 2,  $M$  and  $N$ , can call each other’s methods, and we can get cyclic dependencies between Xmodules.<sup>5</sup> Second, the program may perform improper callbacks. Say a method from  $M$  calls back to method  $B$  from  $L$ . An improper callback  $B$  can call a method of  $N$ , even though the type system knows that  $M$  is to the right of  $N$ . In both cases, the type system al-

<sup>5</sup>Since all non-Xmodule objects are implicitly assigned higher indices than their Xmodule siblings, these non-Xmodule objects cannot introduce cyclic dependencies between Xmodules.

lows a program with cyclic dependency between Xmodules to pass the type checks, which is not allowed by Definition 1.

To have an ownership-aware TM which guarantees exactly Definition 1, one needs to impose additional dynamic checks. The runtime system can use the ownership tags to build a module tree during runtime, and use this module tree to perform dynamic checks to verify that every Xmodule has unique owner indices and contains only proper callbacks. The runtime system can do this by dynamically inferring indices according to which Xmodule calls which other Xmodule, and reporting an error if there is any circular calling.<sup>6</sup>

#### 4. COMPUTATIONS WITH Xmodules

In this section, we formally define the structure of transactional programs with Xmodules. This section converts the informal explanation from Section 2 into a formal model that we later use to prove properties of ownership-aware TM. We briefly review the transactional computation framework [2] and add Xmodules and ownership to this framework, finally providing the formal statement of Definition 1.

##### Transactional Computations

In our framework [2], the execution of a program is modeled using a “computation tree”  $C$  that summarizes the information about both the control structure of a program and the nesting structure of transactions, and an “observer function”  $\Phi$  which characterizes the behavior of memory operations. A program execution is assumed to generate a *trace*  $(C, \Phi)$ .

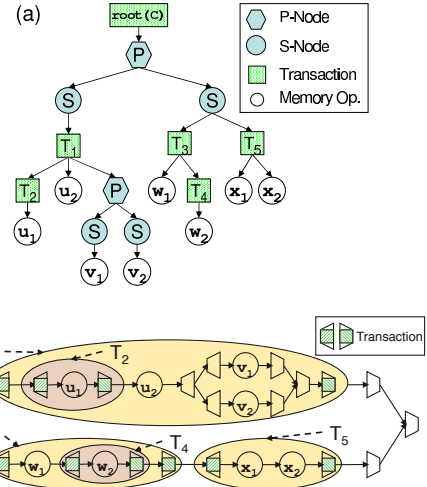
A computation tree  $C$  is defined as an ordered tree with two types of nodes: *memory-operation nodes*  $\text{memOps}(C)$  as leaves and *control nodes*  $\text{spNodes}(C)$  as internal nodes. A memory operation  $v$  either reads from or writes to a memory location. Control nodes are either  $S$  (series) or  $P$  (parallel) nodes. Conceptually, the children of an  $S$ -node must be executed serially, from left to right, while the children of  $P$  node can be executed in parallel. Some  $S$  nodes are labeled as transactions; define  $\text{xactions}(C)$  as the set of these nodes.

Instead of specifying the value that an operation reads or writes to a memory location  $\ell$ , we abstract away the values by using an *observer function*  $\Phi$ . For a memory operation  $v$  that accesses a memory location  $\ell$ , the node  $\Phi(v)$  is defined to be the operation that wrote the value of  $\ell$  that  $v$  sees.

We define several structural notations on the computation tree  $C$ . Denote the *root* of  $C$  as  $\text{root}(C)$ . For any tree node  $X$ , let  $\text{ances}(X)$  denote the set of all  $X$ ’s ancestors (including  $X$  itself) in  $C$ . Denote the set of proper ancestors of  $X$  by  $\text{pAnces}(X)$ . For any tree node  $X$ , we define the *transactional parent* of  $X$ , denoted by  $\text{xparent}(X)$ , as  $\text{parent}(X)$  if  $\text{parent}(X) \in \text{xactions}(C)$ , or  $\text{xparent}(\text{parent}(X))$  if  $\text{parent}(X) \notin \text{xactions}(C)$ . Define the *transactional ancestors* of  $X$  as  $\text{xAnces}(X) = \text{ances}(X) \cap \text{xactions}(C)$ . Denote the *least common ancestor* of two nodes  $X_1, X_2 \in C$  by  $\text{LCA}(X_1, X_2)$ . Define  $\text{xLCA}(X_1, X_2)$  as  $Z = \text{LCA}(X_1, X_2)$  if  $Z \in \text{xactions}(C)$ , and as  $\text{xparent}(Z)$  otherwise.

A computation can also be represented as a computation dag (directed acyclic graph). Given a tree  $C$ , the dag  $G(C) = (V(C), E(C))$  corresponding to the tree is constructed recursively. Every internal node  $X$  in the tree appears as two vertices in the dag. Between these two vertices, the children of  $X$  are connected in series if  $X$  is an  $S$  node, and are connected in parallel if  $X$  is a  $P$  node. Figure 3 show a computation tree and its corresponding computation dag.

<sup>6</sup>It is possible to statically check for unique owner indices by imposing additional restrictions on the program. We opted, however, to describe a more flexible programming model with weaker static guarantees.



**Figure 3.** A sample (a) computation tree  $C$  and (b) its corresponding dag  $G(C)$ .

Classical theories on serializability refer to a particular execution order for a program as a *history* [12]. In our framework, a history corresponds to a topological sort  $S$  of the computation dag  $G(C)$ . We define our models of TM using these sorts. Reordering a history to produce a serial history is equivalent to choosing a different topological sort  $S'$  of  $G(C)$  which has all transactions appearing contiguously, but which is still “consistent” with the observer function associated with  $S$ .

##### Xmodules and Computation Tree

As mentioned in Section 2, in this paper, we consider programs that contain Xmodules. In our theoretical framework, we consider traces generated by a program which is organized into a set  $\mathcal{X}$  of Xmodules. Each Xmodule  $M \in \mathcal{X}$  has some number of methods and a set of memory locations associated with it.

We partition the set of all memory locations  $\mathcal{L}$  into sets of memory owned by each Xmodule. Let  $\text{modMemory}(M) \subseteq \mathcal{L}$  denote the set of memory locations owned by  $M$ . For a location  $\ell \in \text{modMemory}(M)$ , we say that  $\text{owner}(\ell) = M$ . When a method of Xmodule  $M$  is called by a method from a different Xmodule, a safe-nested transaction  $T$  is generated.<sup>7</sup> We use the notation  $\text{xMod}(T) = M$  to associate the instance  $T$  with the Xmodule  $M$ . We also define the instances associated with  $M$  as

$$\text{modXactions}(M) = \{T \in \text{xactions}(C) : \text{xMod}(T) = M\}.$$

As mentioned in Section 2, Xmodules of a program are arranged as a module tree, denoted by  $\mathcal{D}$ . Each Xmodule is assigned an  $\text{xid}$  according to a left-to-right depth-first tree walk, with the root of  $\mathcal{D}$  being  $\text{world}$  with  $\text{xid} = 0$ . Denote the parent of Xmodule  $M$  in  $\mathcal{D}$  as  $\text{modParent}(M)$ , the ancestors of  $M$  as  $\text{modAnces}(M)$ , and the descendants of  $M$  as  $\text{modDesc}(M)$ . We say that  $\text{xMod}(\text{root}(C)) = \text{world}$ , i.e., the root of the computation tree is a transaction associated with the  $\text{world}$  Xmodule.

We use the module tree  $\mathcal{D}$  to restrict the sharing of data between Xmodules and to limit the visibility of Xmodule methods according to the rules given in Definition 2.

**DEFINITION 2** (Formal Restatement of Definition 1). *A program with a module tree  $\mathcal{D}$  should generate only traces  $(C, \Phi)$  which satisfy the following rules:*

<sup>7</sup>As we explained in Section 2, callbacks are handled differently.

1. For any memory operation  $v$  which accesses a memory location  $\ell$ , let  $T = \text{xparent}(v)$ . Then  $\text{owner}(\ell) \in \text{modAncest}(\text{xMod}(T))$ .
2. Let  $X, Y \in \text{xactions}(C)$  be transaction instances such that  $\text{xMod}(X) = M$  and  $\text{xMod}(Y) = N$ . We can have  $X = \text{xparent}(Y)$  only if  $\text{modParent}(N) \in \text{modAncest}(M)$ , and  $\text{xid}(M) < \text{xid}(N)$ .

## 5. THE OAT MODEL

In this section, we informally sketch the *OAT* model, an abstract execution model for TM with ownership and Xmodules. The novel feature of the *OAT* model is that it uses the structure of Xmodules to provide a commit mechanism which can be viewed as a hybrid of closed- and open-nested commits. The *OAT* model presents an operational semantics for TM, and is not intended to describe an actual implementation. For the full formal description of the model, see [1].

### Basic Operation

The TM system is modeled as a nondeterministic state machine with two components: a *program* and a *runtime system*. The runtime system, which we call the *OAT* model, dynamically constructs a computation tree  $C$  as it executes instructions generated by the program. This sequence of instructions is a valid topological sort  $S$  of  $G(C)$ . During execution, each transaction  $T$  in the tree maintains a *status* field, which can be one of COMMITTED, ABORTED, PENDING, or PENDING\_ABORT. The *OAT* model maintains a set of *ready* nodes, denoted by  $\text{ready}(C)$ , and at every step, the *OAT* model nondeterministically chooses one of these ready nodes to issue the next instruction. The program then issues an instruction on this node's behalf.

To detect conflicts, the *OAT* model maintains a read set  $R(T)$  and a write set  $W(T)$  for all  $T \in \text{xactions}(C)$ . The read set  $R(T)$  is a set of pairs  $(\ell, v)$ , where  $\ell \in \mathcal{L}$  is a memory location, and  $v \in \text{memOps}(C)$  is a memory operation that reads from  $\ell$ . We define  $W(T)$  similarly. We also assume that a write is implicitly a read as well; thus,  $W(T) \subseteq R(T)$ .

The *OAT* model performs eager conflict detection; whenever a memory operation  $v$  accesses a location  $\ell$ , the *OAT* model checks to see if  $v$  creates any conflicts. Informally, a  $v$  which is a read (write) generates a conflict if there is another active transaction  $T \notin \text{xAncest}(v)$  ( $T$  is active if its status is PENDING or PENDING\_ABORT) which has  $\ell$  in its write (read) set. If  $v$  generates a conflict, then some transaction must be aborted, using the mechanism explained at the end of this section.

If  $v$  does not generate a conflict, then  $v$  succeeds and observes the value  $\ell$  from  $R(Y)$ , where  $Y$  is the closest transactional ancestor of  $v$  with  $\ell$  in its read set (i.e.,  $(\ell, u) \in R(Y)$ ). Let  $X = \text{xparent}(v)$ . Then, if  $v$  is a read,  $(\ell, v)$  is added to  $R(X)$ . If  $v$  is a write,  $(\ell, v)$  is added to both  $R(X)$  and  $W(X)$ .

### Ownership-Aware Commit

The *OAT* model implements an *ownership-aware commit mechanism* for nested transactions which contains elements of both closed-nested and open-nested commits. A PENDING transaction  $Y$  issues an `xend` instruction to commit  $Y$  into  $X = \text{xparent}(Y)$ . This `xend` commits locations from its read and write sets which are owned by  $\text{xMod}(Y)$  in an open-nested fashion to the root of the tree, while it commits locations owned by other Xmodules in a closed-nested fashion, by merging those reads and writes into  $X$ 's read and write sets.

### Unique Committer Property

Definition 2 guarantees certain properties of the computation tree which are essential to the ownership-aware commit mechanism. Theorem 5 proves that every memory operation has one and only

one transaction that is responsible for committing the memory operation. The proof of the theorem requires the following lemma which we prove by induction on the nesting depth of transactions.

**LEMMA 4.** *Given a computation tree  $C$ , consider any transaction  $T \in \text{xactions}(C)$ . Let  $S_T = \{\text{xMod}(T') : T' \in \text{xAncest}(T)\}$ . Then we have  $\text{modAncest}(\text{xMod}(T)) \subseteq S_T$ .*

**PROOF.** We prove this fact by induction on the nesting depth of transactions  $T$  in the computation tree. In the base case, the lemma holds trivially, since the top-level transaction  $T = \text{root}(C)$ , and  $\text{xMod}(\text{root}(C)) = \text{world}$ . For the inductive step, assume that  $\text{modAncest}(\text{xMod}(T)) \subseteq S_T$  holds for any transaction  $T$  at depth  $d$ . We show that the fact holds for any  $T^* \in \text{xactions}(C)$  at depth  $d+1$ .

For any such  $T^*$ , we know  $T = \text{xparent}(T^*)$  is at depth  $d$ . Then, by Rule 2 of Definition 2, we have  $\text{modParent}(\text{xMod}(T^*)) \in \text{modAncest}(\text{xMod}(T))$ . Thus, we know that  $\text{modAncest}(\text{xMod}(T^*)) \subseteq \text{modAncest}(\text{xMod}(T)) \cup \{\text{xMod}(T^*)\}$ . By construction of the set  $S_T$ , we have  $S_{T^*} = S_T \cup \{\text{xMod}(T^*)\}$ . Therefore, using the inductive hypothesis, we have  $\text{modAncest}(\text{xMod}(T^*)) \subseteq S_{T^*}$ .  $\square$

**THEOREM 5.** *If a memory operation  $v$  accesses a memory location  $\ell$ , then there exists a unique transaction  $T^* \in \text{xAncest}(v)$ , such that*

1.  $\text{owner}(\ell) = \text{xMod}(T^*)$ , and
2. For all transactions  $X \in \text{pAncest}(T^*) \cap \text{xactions}(C)$ ,  $X$  cannot directly access location  $\ell$ .

*This transaction  $T^*$  is the **committer** of memory operation  $v$ , denoted by  $\text{committer}(v)$ .*

**PROOF.** This result follows from the properties of the module tree and computation tree stated in Definition 2.

Let  $T = \text{xparent}(v)$ . First, by Rule 1, we know  $\text{owner}(\ell) \in \text{modAncest}(\text{xMod}(T))$ . We know  $\text{modAncest}(\text{xMod}(T)) \subseteq S_T$  by Lemma 4. Thus, there exists some transaction  $T^* \in \text{xAncest}(T)$  such that  $\text{owner}(\ell) = \text{xMod}(T^*)$ . We can use Rule 2 to show that the  $T^*$  is unique. Let  $X_i$  be the chain of ancestor transactions of  $T$ , i.e., let  $X_0 = T$ , and let  $X_i = \text{xparent}(X_{i-1})$ , up until  $X_k = \text{root}(C)$ . By Rule 2, we know  $\text{xid}(\text{xMod}(X_i)) < \text{xid}(\text{xMod}(X_{i-1}))$ , that is, the  $\text{xids}$  strictly decrease walking up the tree from  $T$ . Thus, there can only be one ancestor transaction  $T^*$  of  $T$  with  $\text{xid}(\text{xMod}(T^*)) = \text{xid}(\text{owner}(\ell))$ .

For any  $X \in \text{pAncest}(T^*) \cap \text{xactions}(C)$ , we can check the second condition. By Rule 1,  $X$  can access  $\ell$  directly only if  $\text{owner}(\ell) \in \text{modAncest}(\text{xMod}(X))$ ; thus, we have  $\text{xid}(\text{owner}(\ell)) \leq \text{xid}(\text{xMod}(X))$ . But we know that  $\text{owner}(\ell) = \text{xMod}(T^*)$  and  $\text{xid}(\text{xMod}(T^*)) > \text{xid}(\text{xMod}(X))$ .  $\square$

Intuitively,  $T^* = \text{committer}(v)$  is the transaction which “belongs” to the same Xmodule as the location  $\ell$  which  $v$  accesses, and is “responsible” for committing  $v$  to memory and making it visible to the world. The second condition of Theorem 5 states that no ancestor transaction of  $T^*$  in the call stack can ever directly access  $\ell$ ; thus, it is “safe” for  $T^*$  to commit  $\ell$ .

### Transaction Abort

When the *OAT* model detects a conflict, it aborts one of the conflicting transactions by changing its status from PENDING to PENDING\_ABORT. In the *OAT* model, a transaction  $X$  might not abort immediately; instead, it might continue to issue more instructions after its status has changed to PENDING\_ABORT. The set of operations issued by  $X$  or its descendants after  $X$ 's status changes to PENDING\_ABORT are called  $X$ 's *abort actions*, denoted by  $\text{abortactions}(X)$ . This condition allows  $X$  to compensate for the safe-nested transactions that may have committed; if transaction  $Y$  is nested inside  $X$ , then the abort actions of  $X$  contain



the compensating action of  $Y$ . Eventually a `PENDING_ABORT` transaction issues an `xend` instruction, which changes its status from `PENDING_ABORT` to `ABORTED`.

If a potential memory operation  $v$  generates a conflict with  $T_u$ , and  $T_u$ 's status is `PENDING`, then the *OAT* model can nondeterministically choose to abort either `xparent`( $v$ ), or  $T_u$ . In the latter case,  $v$  waits for  $T_u$  to finish aborting (i.e., change its status to `ABORTED`) before continuing. If  $T_u$ 's status is `PENDING_ABORT`, then  $v$  just waits for  $T_u$  to finish aborting before trying to issue `read` or `write` again.

This operational model uses the same conflict detection algorithm as *TM* with ordinary closed-nested transactions does; the only subtleties are that  $v$  can generate a conflict with a `PENDING_ABORT` transaction  $T_u$ , and that transactions no longer abort instantaneously because they have abort actions. Some restrictions on the abort actions of a transaction may be necessary to avoid deadlock, as we describe later in Section 7.

## 6. SERIALIZABILITY BY MODULES

In this section, we define *serializability by modules*, a definition inspired by the database notion of multilevel serializability (e.g., as described in [13]). We then provide a proof sketch that the *OAT* model from Section 5 guarantees serializability by modules. For more details about the proof, see [1].

### Notation and Definitions

We first describe some notation needed to formally describe serializability by modules. All definitions in this section are *a posteriori*, i.e., they are defined on the computation tree after the program has finished executing.

We define “content” sets for every transaction  $T$  by partitioning `memOps`( $T$ ) (all the memory operations enclosed inside  $T$  including those belonging to its nested transactions) into three sets: `cContent`( $T$ ), `oContent`( $T$ ) and `aContent`( $T$ ). For any  $u \in \text{memOps}(T)$ , we define the content sets based on the final status of transactions in  $C$  that one visits when walking up the tree from  $u$  to  $T$ .

**DEFINITION 3.** For any transaction  $T$  and memory operation  $u$ , define the sets `cContent`( $T$ ), `oContent`( $T$ ), and `aContent`( $T$ ) according the `ContentType`( $u, T$ ) procedure:

```

ContentType( $u, T$ )  ▷ For any  $u \in \text{memOps}(T)$ 
1   $X \leftarrow \text{xparent}(u)$ 
2  while ( $X \neq T$ )
3    if ( $X$  is ABORTED)    return  $u \in \text{aContent}(T)$ 
4    if ( $X = \text{committer}(u)$ ) return  $u \in \text{oContent}(T)$ 
5     $X \leftarrow \text{xparent}(X)$ 
6  return  $u \in \text{cContent}(T)$ 

```

Recall that in the *OAT* model, the safe-nested commit of  $T$  commits some memory operations in an open-nested fashion, to `root`( $C$ ), and some operations in a closed-nested fashion, to `xparent`( $T$ ). Informally, `oContent`( $T$ ) is the set of memory operations that are committed in an “open” manner by  $T$ 's subtransactions. Similarly, `aContent`( $T$ ) is the set of operations that are discarded due to the abort of some subtransaction in  $T$ 's subtree. Finally, `cContent`( $T$ ) is the set of operations that are neither committed in an “open” manner, nor aborted.

Transactional semantics dictate that memory operations belonging to an aborted transaction  $T$  should not be observed by (i.e., are *hidden* from) memory operations outside of  $T$ .

**DEFINITION 4.** For  $u \in \text{memOps}(C)$ ,  $v \in V(C)$ , let  $X = \text{xLCA}(u, v)$ . We say that  $u$  is *hidden* from  $v$  if  $u \in \text{aContent}(X)$ .

Our notion of serializability requires *sequential consistency*. Without transactions, a trace  $(C, \Phi)$  is said to be sequentially con-

sistent if there exists a topological sort  $S$  of the computation dag  $G(C)$  in which a memory operation  $u$  that accesses  $\ell$  observes the value written by the last writer to  $\ell$  in  $S$ ; that is, the observer function  $\Phi$  is the same as the *last writer* function. For transactional sequential consistency, we define the *transactional last writer* of memory operation  $u$  as a memory operation  $v$  that is the last write in the order  $S$  before  $u$ , skipping over nodes  $w$  which are hidden from (i.e., aborted with respect to)  $u$ . Henceforth, we say that a sort order  $S$  is *sequentially consistent with respect to*  $\Phi$  if  $\Phi$  is the transactional last writer.

### Defining Serializability by Modules

In [2], a trace  $(C, \Phi)$  was said to be *serializable* if there exists a topological sort  $S$  of  $G(C)$  such that  $S$  is sequentially consistent with respect to  $\Phi$ , and all transactions appear contiguous in  $S$ . Serializability in this context can be thought of as sequential consistency plus the requirement that transactions are atomic. This definition of serializability is the “correct definition” for flat or closed-nested transactions. This definition of serializability is too strong, however, for ownership-aware transactions. A *TM* system that enforces this definition of serializability cannot ignore lower-level memory accesses when detecting conflicts for higher-level transactions.

Instead, we describe a definition of serializability by modules which checks for correctness of one *Xmodule* at a time. Given a trace  $(C, \Phi)$ , for each *Xmodule*  $M$ , we transform the tree  $C$  into a new tree `mTree`( $C, M$ ). The tree `mTree`( $C, M$ ) is constructed in such a way as to ignore memory operations of *Xmodules* which are lower-level than  $M$ , and also to ignore all operations which are hidden from transactions of  $M$ . For each *Xmodule*  $M$ , we check that the transactions of  $M$  in the trace `(mTree`( $C, M$ ),  $\Phi$ ) is serializable. If the check holds for all *Xmodules*, then trace  $(C, \Phi)$  is said to be serializable by modules.

Definition 5 formalizes the construction of `mTree`( $C, M$ ).

**DEFINITION 5.** For any computation tree  $C$ , let `mTree`( $C, M$ ) be the result of modifying  $C$  as follows:

1. For all memory operations  $u \in \text{memOps}(C)$  with  $u$  accessing  $\ell$ , if `owner`( $\ell$ ) =  $N$  for some `xid`( $N$ ) > `xid`( $M$ ), convert  $u$  into a *nop*.
2. For all transactions  $T \in \text{modXactions}(M)$ , convert all  $u \in \text{aContent}(T)$  into *nops*.

The intuition behind Condition 1 of Definition 5 is the following. When looking at *Xmodule*  $M$ , we throw away memory operations belonging to a lower-level *Xmodule*  $N$ , since by Theorem 5, transactions of  $M$  can never directly access the same memory as those operations anyway. In Condition 2, we ignore the content of any aborted transactions nested inside transactions of  $M$ ; those transactions might access the same memory locations as operations which we did not turn into *nops*, but those operations are aborted with respect to transactions of  $M$ .

Lemma 6 argues that if a trace  $(C, \Phi)$  is sequentially consistent, then `(mTree`( $C, M$ ),  $\Phi$ ) is a valid trace; an operation  $u$  that remains in the trace never attempts to observe a value from a memory operation  $v = \Phi(u)$  which was turned into a *nop* due to Definition 5. In addition, the transformed trace is also sequentially consistent.

**LEMMA 6.** Let  $(C, \Phi)$  be any sequentially consistent trace. Then for any *Xmodule*  $M$ , `(mTree`( $C, M$ ),  $\Phi$ ) is a valid trace. If  $u \in \text{memOps}(\text{mTree}(C, M))$ , then we have  $\Phi(u) \in \text{memOps}(\text{mTree}(C, M))$ . Furthermore, any  $S$  which is sequentially consistent for  $\Phi$  in  $(C, \Phi)$  is also sequentially consistent for  $\Phi$  in `(mTree`( $C, M$ ),  $\Phi$ ).

PROOF. In the new tree  $\text{mtree}(C, M)$ , pick any memory operation  $u \in \text{memOps}(\text{mtree}(C, M))$  which remains. Assume for contradiction that  $v = \Phi(u)$  was turned into a nop in one of Steps 1 and 2.

If  $v$  was turned into a nop in Step 1 of Definition 5, then we know that  $v$  accessed a memory location  $\ell$  where  $\text{xid}(\text{owner}(\ell)) > \text{xid}(M)$ . Since  $u$  must access the same location  $\ell$ ,  $u$  must also be converted into a nop. If  $v$  was turned into a nop in Step 2 of Definition 5, then  $v \in \text{aContent}(T)$  for some  $\text{xMod}(T) = M$ . Then we can show that either  $vHu$ , or  $u$  should have also been turned into a nop. Let  $X = \text{xLCA}(v, u)$ . Since  $X$  and  $T$  are both ancestors of  $v$ , either  $X$  is an ancestor of  $T$  or  $T$  is a proper ancestor of  $X$ .

1. First, suppose  $T$  is a proper ancestor of  $X$ . Consider the path of transactions  $Y_0, Y_1, \dots, Y_k$ , where  $Y_0 = \text{xparent}(v)$ , for each  $0 < i < k$ , we have  $\text{xparent}(Y_i) = Y_{i+1}$ , and  $\text{xparent}(Y_k) = T$ . Since  $v \in \text{aContent}(T)$ , for some  $Y_j$  for  $0 \leq j \leq k$  must have  $\text{status}[Y_j] = \text{ABORTED}$ . Since  $T$  is a proper ancestor of  $X$ ,  $X = Y_x$  for some  $x$  satisfying  $0 \leq x \leq k$ .
  - (a) If  $\text{status}[Y_j] = \text{ABORTED}$  for any  $j$  satisfying  $0 \leq j < x$ , then we know  $v \in \text{aContent}(X)$ , and thus  $vHu$ . Since we assumed  $(C, \Phi)$  is sequentially consistent and  $\Phi(u) = v$ , we know  $\neg vHu$ , leading to a contradiction.
  - (b) If  $Y_j$  is  $\text{ABORTED}$  for any  $j$  satisfying  $x \leq j \leq k$ , then  $\text{status}[Y_j] = \text{ABORTED}$  implies that  $u \in \text{aContent}(X)$ , and thus,  $u$  should have been turned into a nop, contradicting the original setup of the statement.
2. Next, consider the case where  $X$  is an ancestor of  $T$ . Since  $v \in \text{aContent}(T)$ , we have  $v \in \text{aContent}(X)$ . Therefore, this case is analogous to Case 1a above.

Finally, if  $\Phi$  is the transactional last writer according to  $\mathcal{S}$  for  $(C, \Phi)$ , it is still the transactional last writer for  $(\text{mtree}(C, M), \Phi)$  because the memory operations which are not turned into nops remain in the same relative order. Thus, the last condition is satisfied.  $\square$

Note that Lemma 6 depends on the restrictions on Xmodules described in Definition 2. Without this structure of modules and ownership, the construction of Definition 5 is not guaranteed to generate a valid trace.

Finally, we can define serializability by modules.

**DEFINITION 6.** A trace  $(C, \Phi)$  is *serializable by modules* if it is sequentially consistent, and if for all Xmodule  $M$  in  $\mathcal{D}$ , there exists a topological sort  $\mathcal{S}$  of  $C_M = \text{mtree}(C, M)$  such that:

1.  $\mathcal{S}$  is sequentially consistent with respect to  $\Phi$ , and
2. For the tree  $C_M$ ,  $\forall T \in \text{modXactions}(M)$  and  $\forall v \in V(C_M)$ , if we have  $\text{xbegin}(T) \leq_{\mathcal{S}} v \leq_{\mathcal{S}} \text{xend}(T)$ , then  $v \in V(T)$ .

Informally, a trace  $(C, \Phi)$  is serializable by modules if it is sequentially consistent, and if for every Xmodule  $M$ , there exists a sequentially consistent order  $\mathcal{S}$  for the trace  $(\text{mtree}(C, M), \Phi)$  which has all transactions of  $M$  contiguous.

### OAT Model Guarantees Serializability by Modules

We can show that the OAT model described in Section 5 generates traces that satisfy Definition 6.

**THEOREM 7.** Any trace  $(C, \Phi)$  generated by the OAT model is serializable by modules.

PROOF SKETCH. The proof consists of three steps. First, we generalize the notion of “prefix race freedom” [2] to computations with Xmodules. Second, we prove that the OAT model guarantees that a program execution is prefix race free. Finally, we argue that any trace which is prefix race free is also serializable by modules. See [1] for details.  $\square$

### Abstract Serializability

By Theorem 7, the OAT model guarantees serializability by modules. We now relate this definition to the notion of *abstract serializability* used in multilevel database systems [13]. As we mentioned in Section 1, the ownership-aware commit mechanism is a part of a methodology which includes abstract locks and compensating actions. In this section we argue that OAT model provides enough flexibility to accommodate abstract locks and compensating actions. In addition, if a program is “properly locked and compensated,” then serializability by modules guarantees abstract serializability.

The definition of abstract serializability in databases [13] assumes that the program is divided into levels, and that a transaction at level  $i$  can only call a transaction at level  $i + 1$ .<sup>8</sup> In addition, transactions at a particular level have predefined commutativity rules, i.e., some transactions of the same Xmodule can commute with each other and some cannot. The transactions at the lowest level (say  $k$ ) are naturally serializable; call this schedule  $Z_k$ . Given a serializable schedule  $Z_{i+1}$  of level- $i + 1$  transactions, the schedule is said to be serializable at level  $i$  if all transactions in  $Z_{i+1}$  can be reordered, obeying all commutativity rules, to obtain a serializable order  $Z_i$  for level- $i$  transactions. The original schedule is said to be abstractly serializable if it is serializable for all levels.

These commutativity rules might be specified using abstract locks [11]: if two transactions cannot commute, then they grab the same abstract lock in a conflicting manner. In the application described in Section 2, for instance, transactions calling insert and remove on the BST using the same key do not commute and should grab the same write lock. Although abstract locks are not explicitly modeled in the OAT model, we can model transactions acquiring the same abstract lock as transactions writing to a common memory location  $\ell$ .<sup>9</sup> Locks associated with an Xmodule  $M$  are owned by  $\text{modParent}(M)$ . A module  $M$  is said to be *properly locked* if the following is true for all transactions  $T_1, T_2$  with  $\text{xMod}(T_1) = \text{xMod}(T_2) = M$ : if  $T_1$  and  $T_2$  do not commute, then they access some  $\ell \in \text{modMemory}(\text{modParent}(M))$  in a conflicting manner.

If all transactions are properly locked, then serializability by modules implies abstract serializability (as defined above) in the special case when the module tree is a chain (i.e., each non-leaf module has exactly one child). Let  $S_i$  be the sort  $\mathcal{S}$  in Definition 6 for Xmodule  $M$  with  $\text{xid}(M) = i$ . This  $S_i$  corresponds to  $Z_i$  in the definition of abstract serializability.

In the general case for ownership-aware TM, however, by Rule 2 of Definition 1, we know a transaction at level  $i$  might call transactions from multiple levels  $x > i$ , not just  $x = i + 1$ . Thus, we must change the definition of abstract serializability slightly; instead of reordering just  $Z_{i+1}$  while serializing transactions at level- $i$ , we have to potentially reorder  $Z_x$  for all  $x$  where transactions at level  $i$  can call transactions at level  $x$ . Even in this case, if every module is properly locked (by the same definition as above), one can show serializability by modules guarantees abstract serializability.

The methodology of open nesting often requires the notion of compensating actions or inverse actions. For instance, in a BST, the inverse of insert is remove with the same key. When a transaction  $T$  aborts, all the changes made by its subtransactions must be inverted. Again, although the OAT model does not explicitly model compensating actions, it allows an aborting transaction with status `PENDING_ABORT` to perform an arbitrary but finite number of opera-

<sup>8</sup> We assume level number increases as you go from a higher level to a lower-level to be consistent with our numbering of `xid`. In the literature, levels typically go in the opposite direction.

<sup>9</sup> More complicated locks can be modeled by generalizing the definition of conflict.

tions before changing the status to ABORTED. Therefore, an aborting transaction can compensate for all its aborted subtransactions.

## 7. DEADLOCK FREEDOM

In this section, we argue that the OAT model described in Section 5 can never enter a “semantic deadlock” if we impose suitable restrictions on the memory accessed by a transaction’s abort actions. In particular, an abort action generated by transaction  $T$  from  $\text{xMod}(T)$  should read (write) from a memory location  $\ell$  belonging to  $\text{modAncest}(\text{xMod}(T))$  only if  $\ell$  is already in  $\text{R}(T)$  ( $\text{W}(T)$ ). Under these conditions, we show that the OAT model can always “finish” reasonable computations.

An ordinary TM without open nesting and with eager conflict detection never enters a semantic deadlock because it is always possible to finish aborting a transaction  $T$  without generating additional conflicts; a scheduler in the TM runtime can abort all transactions, and then complete the computation by running the remaining transactions serially. Using the OAT model, however, a TM system can enter a semantic deadlock because it can enter a state in which it is impossible to finish aborting two parallel transactions  $T_1$  and  $T_2$  which both have status PENDING\_ABORT. If  $T_1$ ’s abort action generates a memory operation  $u$  which conflicts with  $T_2$ , then  $u$  will wait for  $T_2$  to finish aborting (i.e., when the status of  $T_2$  becomes ABORTED). Similarly,  $T_2$ ’s abort action can generate an operation  $v$  which conflicts with  $T_1$  and waits for  $T_1$  to finish aborting. Thus  $T_1$  and  $T_2$  can both wait on each other, and neither transaction will ever finish aborting.

### Defining Semantic Deadlock

Intuitively, we want to say that a TM system exhibits a semantic deadlock if it might enter a state from which it is impossible to “finish” a computation because of transaction conflicts. This section defines semantic deadlock precisely and distinguishes it from these other reasons for noncompletion, such as livelock or infinite loop.

Recall that our abstract model has two entities: the program, and a generic operational model  $\mathcal{F}$  representing the runtime system. At any time  $t$ , given a ready node  $X \in \text{ready}(C)$ , the program chooses an instruction and has  $X$  issue the instruction. If the program issues an infinite number of instructions, then  $\mathcal{F}$  cannot complete the program no matter what it does. To eliminate programs which have infinite loops, we only consider **bounded programs**.

**DEFINITION 7.** We say that a program is **bounded** for an operational model  $\mathcal{F}$  if any computation tree that  $\mathcal{F}$  generates for that program is of a finite depth, and there exists a finite number  $K$  such that at any time  $t$ , every node  $B \in \text{nodes}^{(t)}(C)$  has at most  $K$  children with status PENDING or COMMITTED.

Even if the program is bounded, it might still run forever if it **livelocks**. We use the notion of a **schedule** to distinguish livelocks from semantic deadlocks.

**DEFINITION 8.** A **schedule**  $\Gamma$  on some time interval  $[t_0, t_1]$  is the sequence of nondeterministic choices made by an operational model in the interval.

An operational model  $\mathcal{F}$  makes two types of nondeterministic choices. First, at any time  $t$ ,  $\mathcal{F}$  nondeterministically chooses which ready node  $X \in \text{ready}(C)$  executes an instruction. This choice models nondeterminism in the program due to interleaving of the parallel executions. Second, while performing a memory operation  $u$  which generates a conflict with transaction  $T$ ,  $\mathcal{F}$  nondeterministically chooses to abort either  $\text{xparent}(u)$  or  $T$ . This nondeterministic choice models the contention manager of the TM runtime. A program may livelock if  $\mathcal{F}$  repeatedly makes “bad” scheduling choices.

Intuitively, an operational model deadlocks if it allows a **bounded computation** to reach a state where **no schedule** can complete the computation after this point.

**DEFINITION 9.** Consider an  $\mathcal{F}$  executing a bounded computation. We say that  $\mathcal{F}$  does not exhibit a **semantic deadlock** if for all finite sequences of  $t_0$  instructions that  $\mathcal{F}$  can issue that generates some intermediate computation tree  $C_0$ , there exists a finite schedule  $\Gamma$  on  $[t_0, t_1]$  such that  $\mathcal{F}$  brings the computation tree to a rest state  $C_1$ , i.e.,  $\text{ready}(C_1) = \{\text{root}(C_1)\}$ .

This definition is sufficient, since once the computation tree is at the rest state, and only the root node is ready,  $\mathcal{F}$  can execute each transaction serially and complete the computation.

### Restrictions to Avoid Semantic Deadlock

The general OAT model described in Section 5 exhibits semantic deadlock because it may enter a state where two parallel aborting transactions  $T_1$  and  $T_2$  keep each other from completing their aborts. For a restricted set of programs, where a PENDING\_ABORT transaction  $T$  never accesses new memory belonging to Xmodules at  $\text{xMod}(T)$ ’s level or higher, however, we can show the OAT model is free of semantic deadlock.

More formally, for all transactions  $T$ , we restrict the memory footprint of  $\text{abortactions}(T)$ .

**DEFINITION 10.** An execution (represented by a computation tree  $C$ ) has **abort actions with limited footprint** if the following condition is true for all transactions  $T \in \text{aborted}(C)$ . At time  $t$ , if a memory operation  $v \in \text{abortactions}(T)$  accesses location  $\ell$  and  $\text{owner}(\ell) \in \text{modAncest}(\text{xMod}(T))$ , then (1) if  $v$  is a read, then  $\ell \in \text{R}(T)$ , and (2) if  $v$  is a write then  $\ell \in \text{W}(T)$ .

Intuitively, Definition 10 requires that once a transaction  $T$ ’s status becomes PENDING\_ABORT, any memory operation  $v$  which  $T$  or a nested transaction inside  $T$  performs to finish aborting  $T$  cannot read from (write to) any location  $\ell$  which is owned by any Xmodules which are ancestors of  $\text{xMod}(T)$  (including  $\text{xMod}(T)$  itself), unless  $\ell$  is already in the read (or write set) of  $T$ .

First, we show that the properties of Xmodules from Theorem 5 in combination with the ownership-aware commit mechanism imply that transaction read sets and write sets exhibit nice properties. In particular, we have Corollary 8, which states that a location  $\ell$  can appear in the read set of a transaction  $T$  only if  $T$ ’s Xmodule is a descendant of  $\text{owner}(\ell)$  in the module tree  $\mathcal{D}$ .

**COROLLARY 8.** For any transaction  $T$  if  $\ell \in \text{R}(T)$ , then  $\text{xMod}(T) \in \text{modDesc}(\text{owner}(\ell))$ .

**PROOF.** Follows from Definition 1 and Theorem 5, and induction on how a location  $\ell$  can propagate into readsets and writsets using the ownership-aware commit mechanism.  $\square$

If all abort actions have a limited footprint, we can show that operations of an abort action of an Xmodule  $M$  can only generate conflicts with a “lower-level” Xmodule.

**LEMMA 9.** Suppose the OAT model generates an execution where abort actions have limited footprint. For any transaction  $T$ , consider a potential memory operation  $v \in \text{abortactions}(T)$ . If  $v$  conflicts with transaction  $T'$ , then  $\text{xid}(\text{xMod}(T')) > \text{xid}(\text{xMod}(T))$ .

**PROOF.** Suppose  $v \in \text{abortactions}(T)$  accesses a memory location  $\ell$  with  $\text{owner}(\ell) = M$ . Since  $\text{abortactions}(T) \subseteq \text{memOps}(T)$ , by the properties of Xmodules given in Definition 2, we know that either  $M \in \text{modAncest}(\text{xMod}(T))$ , or  $\text{xid}(M) > \text{xid}(\text{xMod}(T))$ . If  $M \in \text{modAncest}(\text{xMod}(T))$ , then by Definition 10,  $T$  already had  $\ell$  in its read or write set. Therefore,  $v$  can not generate a conflict with  $T'$  because then  $T$  would already have had a conflict with  $T'$  before

$v$  occurred, contradicting the eager conflict detection of the *OAT* model.

Thus, we have  $\text{xid}(M) > \text{xid}(\text{xMod}(T))$ . If  $v$  conflicts with some other transaction  $T'$ , then  $T'$  has  $\ell$  in its read or write set. Therefore, from Corollary 8,  $\text{xMod}(T')$  is a descendant of  $M$ . Thus, we have  $\text{xid}(\text{xMod}(T')) > \text{xid}(M) > \text{xid}(\text{xMod}(T))$ .  $\square$

**THEOREM 10.** *In the case where aborted actions have limited footprint, the OAT model is free from semantic deadlock.*

PROOF. Let  $C_0$  be the computation tree after any finite sequence of  $t_0$  instructions. We describe a schedule  $\Gamma$  which finishes aborting all transactions in the computation by executing abort actions and transactions serially.

Without loss of generality, assume that at time  $t_0$ , all active transactions  $T$  have  $\text{status}[T] = \text{PENDING\_ABORT}$ . Otherwise, the first phase of the schedule  $\Gamma$  is to make this status change for all active transactions  $T$ .

For a module tree  $\mathcal{D}$  with  $k = |\mathcal{D}|$  Xmodules (including the world), we construct a schedule  $\Gamma$  with  $k$  phases,  $k-1, k-2, \dots, 1, 0$ . The invariant we maintain is that immediately before phase  $i$ , we bring the computation tree into a state  $C^{(i)}$  which has no active transaction instances  $T$  with  $\text{xid}(\text{xMod}(T)) > i$ , i.e., no instances  $T$  from Xmodules with  $\text{xid}$  larger than  $i$ . During phase  $i$ , we finish aborting all active transaction instances  $T$  with  $\text{xid}(\text{xMod}(T)) = i$ . By Lemma 9, any abort action for a  $T$ , where  $\text{xid}(\text{xMod}(T)) = i$ , can only conflict with a transaction instance  $T'$  from a lower-level Xmodule, where  $\text{xid}(\text{xMod}(T')) > i$ . Since the schedule  $\Gamma$  executes serially, and since by the inductive hypothesis we have already finished all active transaction instances from lower levels, phase  $i$  can finish without generating any conflicts.  $\square$

### Restrictions on compensating actions

If transactions  $Y_1, Y_2, \dots, Y_j$  are nested inside transaction  $X$  and  $X$  aborts, typically abort actions of  $X$  simply consists of compensating actions for  $Y_1, Y_2, \dots, Y_j$ . Thus, restrictions on abort actions translate in a straightforward manner to restrictions on compensating actions: a compensating action for a transaction  $Y_i$  (which is part of the abort action of  $X$ ), should not read (write) any memory owned by  $\text{xMod}(X)$  or its ancestor Xmodules unless the memory location is already in  $X$ 's read (write) set. Assuming locks are modeled as accesses to memory locations, the same restriction applies, meaning, a compensating action cannot acquire new locks that were not already acquired by the transaction it is compensating for.

## 8. CONCLUSIONS

In this paper, we describe ownership-aware transactions, which provide a disciplined methodology for open nesting while guaranteeing abstract serializability. In this section, we describe two other approaches for improving open-nested transactions, and distinguish them from our work.

Ni, et al. [11] propose using an `open_atomic` class to specify open-nested transactions in a Java-like language with transactions. Since the private fields of an object with an `open_atomic` class type can not be directly accessed outside of that class, one can think of the `open_atomic` class as defining an Xmodule. This mapping is not exact, however, because neither the language nor TM system restrict exactly what memory can be passed into a method of an `open_atomic` class, and the TM system performs a vanilla open-nested commit for a nested transaction, not a safe-nested commit. Thus, it is unclear what exact guarantees are provided with respect to serializability and/or deadlock freedom.

Herlihy and Koskinen [5] describe a technique of transactional boosting which allows transactions to call methods from a non-transactional module  $M$ . Roughly, as long as  $M$  is linearizable and

its methods have well-defined inverses, the authors show that the execution appears to be “abstractly serializable.” Boosting does not, however, address the cases when the lower-level module  $M$  writes to memory owned by the enclosing higher-level module, or when programs have more than two levels of modules.

### Acknowledgements

We thank James Noble of Victoria University of Wellington, Derek Rayside, Martin Rinard, Amy Williams, and Charles Leiserson and other members of the Supercomputing Technologies Group at MIT CSAIL for helpful discussions and comments on the paper. We also thank all the reviewers of this and prior versions of the paper for their comments. In particular, we are grateful to Bill Scherer for his help in improving the paper.

## REFERENCES

- [1] K. Agrawal, I.-T. A. Lee, and J. Sukha. Safe open-nested transactions through ownership (Technical report). Technical report, Laboratory of Computer Science and Artificial Intelligence, Massachusetts Institute of Technology, June 2008. Available at: <http://supertech.csail.mit.edu/papers/safe-tech.pdf>.
- [2] K. Agrawal, C. E. Leiserson, and J. Sukha. Memory models for open-nested transactions. In *Proceedings of the ACM SIGPLAN Workshop on Memory Systems Performance and Correctness (MSPC)*, October 2006. In conjunction ASPLOS.
- [3] C. Boyapati, B. Liskov, and L. Shrira. Ownership types for object encapsulation. In *Proceedings of the ACM Symposium on Principles of Programming Languages (POPL)*, New Orleans, Louisiana, Jan. 2003.
- [4] B. D. Carlstrom, A. McDonald, M. Carbin, C. Kozyrakis, and K. Olukotun. Transactional collection classes. In *Proceedings of the ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming (PPoPP)*, pages 56–67, New York, NY, USA, 2007. ACM Press.
- [5] M. Herlihy and E. Koskinen. Transactional boosting: a methodology for highly-concurrent transactional objects. In *Proceedings of ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming (PPoPP)*, pages 207–216, New York, NY, USA, Feb 2008. ACM.
- [6] M. Herlihy and J. E. B. Moss. Transactional memory: Architectural support for lock-free data structures. In *Proceedings of the International Symposium on Computer Architecture (ISCA)*, pages 289–300, 1993.
- [7] A. McDonald, J. Chung, B. D. Carlstrom, C. Cao Minh, H. Chafi, C. Kozyrakis, and K. Olukotun. Architectural semantics for practical transactional memory. In *Proceedings of the International Symposium on Computer Architecture (ISCA)*, June 2006.
- [8] J. E. B. Moss. *Nested Transactions: An Approach to Reliable Distributed Computing*. MIT Press, Cambridge, MA, USA, 1985.
- [9] J. E. B. Moss. Open nested transactions : Semantics and support. In *Proceedings of the Workshop on Memory Performance Issues (WMPI)*, Austin, Texas, Feb 2006.
- [10] J. E. B. Moss and A. L. Hosking. Nested transactional memory: Model and architecture sketches. In *Science of Computer Programming*, volume 63, pages 186–201. Elsevier, Dec 2006.
- [11] Y. Ni, V. Menon, A. Adl-Tabatabai, A. L. Hosking, R. L. Hudson, J. E. B. Moss, B. Saha, and T. Shpeisman. Open nesting in software transactional memory. In *Proceedings of ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming (PPoPP)*, Mar. 2007.
- [12] C. H. Papadimitriou. The serializability of concurrent database updates. *Journal of the ACM*, 26(4):631–653, 1979.
- [13] G. Weikum. A theoretical foundation of multi-level concurrency control. In *Proceedings of the ACM SIGACT-SIGMOD symposium on Principles of database systems (PODS)*, pages 31–43, New York, NY, USA, 1986. ACM Press.