

Safe Open-Nested Transactions Through Ownership

Kunal Agrawal I-Ting Angelina Lee Jim Sukha

MIT Computer Science and Artificial Intelligence Laboratory

{kunal_ag, angelee, sukhaj}@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*.

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 the memory locations that each transaction accesses (using *read sets*

This research was supported in part by NSF Grants CNS-0615215 and CNS-0540248 and a grant from Intel corporation. Preliminary versions of this paper appeared as a poster at *PPoPP* 2008 and as a brief announcement at *SPAA* 2008.

and *write sets*), finding transaction 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.

When using TM, one issue that programmers must deal with is the semantics of *nested transactions*. If a transaction Y is *closed nested* [8] inside transaction X , 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, upon Y 's commit, 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 accesses to shared “low-level” memory done by nested transactions.

Researchers have proposed the methodology of *open-nested transactions* to increase concurrency in transactional programs by carefully 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 in [9] illustrates use of open nesting with an application that uses 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 ℓ_1 and ℓ_2 , and X does not care about changes made to ℓ_2 , but does care about ℓ_1 . The TM system can not 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 language 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, a change to memory location ℓ is committed globally if ℓ belongs to the module of T ; otherwise, the read or write to ℓ 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 module ever accesses data belonging to other modules. Thus, one corollary of our theorem is that open-nested transactions are serializable when modules 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 in [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 programs to 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 ℓ_1), and also writes some internal database metadata (location ℓ_2). After a vanilla open-nested commit of Y , the modifications to both ℓ_1 and ℓ_2 become visible globally. Assuming the `UserApp` does not care about the internal state of the database, committing the internal state of the `DB` (ℓ_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 ℓ_1 breaks the atomicity of transaction X . A transaction Z in parallel with transaction X can access this location ℓ_1 after Y commits, before the outer transaction X commits.¹ 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 modules. During program execution, a call to a method from `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.

¹Note 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.

In the ownership-aware TM system, every memory location is owned by exactly one Xmodule. If a memory location ℓ 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 ℓ ; in this case, we say that T is “responsible” for that access to ℓ . Otherwise, the read or write to ℓ 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 ℓ .

For ownership-aware TM to behave “nicely”, we must restrict interactions between Xmodules. For example, in the TM system, some transaction must be “responsible” for committing every memory access. Similarly, the TM system should guarantee some form of serializability. If Xmodules could arbitrarily call methods from or access memory owned by other Xmodules, then these two 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 ℓ directly only if ℓ is either owned by 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 can not 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 $\text{xid}(N) > \text{xid}(M)$ (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.²*

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., $\text{xid}(M) < \text{xid}(N)$), and thus N can not 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 the same Xmodule.

² 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 condition at the end of this section.

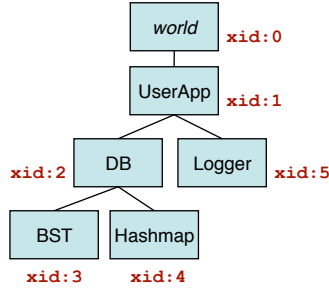


Figure 1. A module tree \mathcal{D} for the program described in Section 2. The `xid`'s are assigned according to a left-to-right depth-first tree walk, numbering Xmodules in increasing order, starting with `xid(world) = 0`.

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` 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` can not 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` can not 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 ℓ , and suppose Y writes value v_1 to location ℓ and then commits. Now a transaction Z in parallel with X can read this location ℓ , write value v_2 to ℓ , and commit, all before X commits. Therefore, X can now read this location ℓ and see

the value v_2 , which is neither the initial value v_0 (the value of ℓ 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 ℓ , X can not access ℓ , 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 ℓ , the ownership-aware commit of Y will not commit the changes to ℓ globally and ℓ will be propagated to X 's write set. Therefore, if Z tries to access ℓ before X commits, the TM system will detect a conflict. Thus X can not see an inconsistent value for ℓ .³

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 new 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. Therefore, 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 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.

³For 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.

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. *The type system in [3] 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, the type system of [3] 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. Within a class definition, 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).

It is shown in [3] 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 can not 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 in [3] exhibits some of the same properties as the module tree we described in Section 2; however, the type system and ownership scheme of [3] do 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 `Xmodule`s 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 the type system of [3]. 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 opted 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. In a class declaration, 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 $T2$ and $T1$ have the same type and all the formals match in name. In addition, if the indices are specified for the tags, then they must match.
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 $o_1.name \prec o_i.name$ or $o_1.name = o_i.name$ and $o_1.index < o_i.index$ (if both indices are known).⁴
4. The ownership tag `this` can only be used within the definition of a class that extends `Xmodule`.
5. `Xmodule` objects can not 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`.

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 $o_i.name \prec o_j.name$, 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.

⁴In 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 {
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
10 public class DB<db0, log0, data0>
11     extends Xmodule where (log0 < data0) {
12     private Logger<log0, data0> logger;
13     private BST<this[0], log0, data0> bst;
14     private Hashmap<this[1], log0, data0> hashmap;
15     public DB(Logger<log0, data0> logger) {
16         this.logger = logger;
17         ...
18     }
19 }

```

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

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.

We prove (by contradiction) that if least common ancestor of M and N in the module tree is $o_N.name$, and $xid(M) > xid(N)$, then M can not have a formal tag with value o_N . Therefore, it can not declare a type with owner tag o_N , and can not access N .

Let L be the least common ancestor of M . Since $L = o_N.name$, we know that L is N 's parent. 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 $\text{xid}(M) > \text{xid}(N)$, 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 .

Let $P_0 = Q$, and consider the chain of Xmodule instantiations where Xmodule P_i instantiated P_{i-1} . P_1 has to instantiate Q (which is the same as P_0) using its formal ownership tags $\langle f_a^1, f_b^1, \dots \rangle$, where f_a^1 has value o_Q and f_b^1 has value o_N . (We must have f_a^1 as the first formal, since o_Q is the owner of Q . Without loss of generality, we can have f_b^1 be the second formal since the type system does not care about the ordering of formal tags after the first one.)

Since $o_N.name = o_Q.name = L$, this chain of instantiations must lead back to L , since that is the only Xmodule that can create ownership tags with values o_N and o_Q in its class definition (using the keyword `this`).⁵ Let $P_k = L$. For the class declaration of each of the Xmodules P_i for $1 \leq i < k$, the following must be true.

- P_i must have formals f_a^i and f_b^i , with values o_Q and o_N , respectively, and P_i must pass these formals into the instantiation of P_{i-1} .
- In the type definition of P_i 's class, P_i must have the constraint $f_a^i < f_b^i$ on its formal tags (either because f_a^i is the owner tag, or through a `where` clause that enforces $f_a^i < f_b^i$).

The first condition must hold for us to be able to pass both o_N and o_Q down to $P_0 = Q$. The second condition is true for the Xmodules by induction. In the base case, P_1 must know that $f_a^1 < f_b^1$; otherwise, the type system will throw an error when it tries to instantiate $P_0 = Q$ with owner f_a^1 . Then, inductively, P_i must know $f_a^i < f_b^i$ to be able to instantiate P_{i-1} .

Finally, P_{k-1} is instantiated in the class file corresponding to $P_k = L$. In this declaration, the formal f_a^k with value o_Q is instantiated with `this[x]`. Similarly, f_b^k with value o_N is instantiated with `this[y]`. Since the class definition of P_k type checks, we must have $f_a^k < f_b^k$. This check contradicts our original assumption that $x > y$ however, since if $x > y$ our type check should fail. Therefore, we must have $o_Q.index < o_N.index$. □

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 can not 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.⁶ 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 allows a program with cyclic dependency between Xmodules to pass the type checks, which is not allowed by Definition 1.

⁵Note that L could be the `world` Xmodule, in which case both o_N and o_Q were created in the `main` function using the `world` keyword.

⁶Since all non-Xmodule objects are implicitly assigned higher indices than their Xmodule siblings, these non-Xmodule objects can not introduce cyclic dependencies between Xmodules.

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.⁷

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. First, we briefly review the transactional computation framework described in [2]. We then add Xmodules and ownership to this framework, and provide the formal statement of Definition 1.

Transactional Computations

In the framework from [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” Φ which characterizes the behavior of memory operations. A program execution is assumed to generate a *trace* (C, Φ) .

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 satisfies the *read predicate* $R(v, \ell)$ if v reads from location ℓ , while v satisfies the *write predicate* $W(v, \ell)$ if v writes to ℓ . 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 ℓ , we abstract away the values by using an *observer function* Φ . For a memory operation v that accesses a memory location ℓ , the node $\Phi(v)$ is defined to be the operation that wrote the value of ℓ 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 . Similarly, let $\text{desc}(X)$ denote the set of all X ’s descendants, including X itself. 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 $\text{xparent}(X)$, 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.

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 .

⁷It 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.

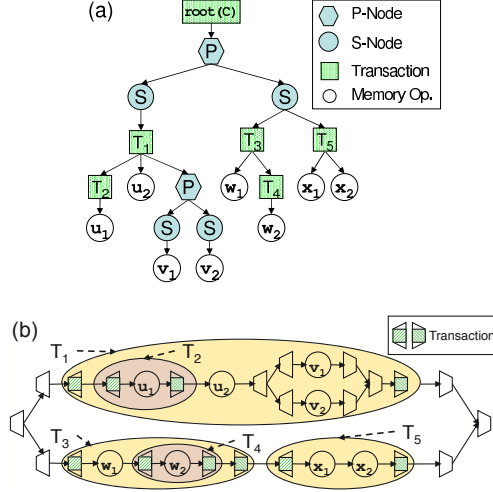


Figure 3. A sample (a) computation tree C and (b) its corresponding dag $G(C)$.

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{N} of Xmodules. Each Xmodule $M \in \mathcal{N}$ 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.⁸ 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 the module tree, denoted by \mathcal{D} . Each Xmodule is assigned an xid according to a left-to-right depth-first tree walk, with the root of \mathcal{D} being world with $\text{xid} = 0$. Denote the parent of Xmodule M in \mathcal{D} as $\text{modParent}(M)$, and the ancestors of M as $\text{modAnces}(M)$ (include M itself). Similarly, let $\text{modDesc}(M)$ be the set of M 's descendants. We say that $\text{xMod}(\text{root}(C)) = \text{world}$, i.e., the root of the computation tree is a transaction associated with the 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, Φ) which satisfy the following rules:*

1. *For any memory operation v which accesses a memory location ℓ , let $T = \text{xparent}(v)$. Then $\text{owner}(\ell) \in \text{modAnces}(\text{xMod}(T))$.*
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 $\text{modParent}(N) \in \text{modAnces}(M)$, and $\text{xid}(M) < \text{xid}(N)$.*

5. THE OAT MODEL

In this section, we describe 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

⁸As we explained in Section 2, callbacks are handled differently.

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, although these semantics can be used to guide an implementation.

Overview

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 and traverses a computation tree C as it executes instructions generated by the program. The *OAT* model maintains a set of **ready** nodes, denoted by $\text{ready}(C) \subseteq \text{nodes}(C)$, and at every step, the *OAT* model nondeterministically chooses one of these ready nodes $X \in \text{ready}(C)$ to issue the next instruction. The program then issues one of the following instructions (whose precondition is satisfied) on X 's behalf: `fork`, `join`, `xbegin`, `xend`, `xabort`, `read`, or `write`. For shorthand, we sometimes say that X issues an instruction.

The *OAT* model describes a sequential semantics, that is, we assume at every time step t , a program issues a single instruction. The parallelism in this model arises from the fact that at a particular time, several nodes can be ready, and the runtime nondeterministically chooses which node to issue an instruction.

In the rest of this section, we give a detailed description of the *OAT* model. First, we describe the state information maintained by the *OAT* model and define the notation we use to refer to this state. Second, we describe how the *OAT* model constructs and traverses the computation tree as instructions are issued. Then, we describe how the *OAT* model handles memory operations (i.e., `read` and `write`), conflict detection, and transaction commits, and transaction aborts.

5.1 State Information and Notation

As the *OAT* model executes instructions, it dynamically constructs the computation tree C . For each of the sets defined in Section 4 (e.g., $\text{nodes}(C)$, $\text{spNodes}(C)$, $\text{memOps}(C)$, $\text{xactions}(C)$, etc.), we define corresponding time-dependent versions of these sets by indexing them with an additional time argument. For example, we define the set $\text{nodes}(t, C)$ denotes the set of nodes in the computation tree after t time steps have passed. The generalized sets from Section 4 are monotonically increasing, i.e., once an element is added to the set, it is never removed at a later time t . Sometimes for shorthand, we omit the time argument when it is clear that we are referring to a particular fixed time t .

At any time t , each internal node $A \in \text{spNodes}(t, C)$ has a **status** field $\text{status}[A]$. These status fields change with time. If $A \in \text{xactions}(t, C)$, i.e., A is a transaction, then $\text{status}[A]$ can be one of `COMMITTED`, `ABORTED`, `PENDING`, or `PENDING_ABORT`. Otherwise, $A \in \text{spNodes}(t, C) - \text{xactions}(t, C)$ is either a P-node or a nontransactional S-node; in this case, $\text{status}[A]$ can either be `WORKING` or `SYNCHED`. We define several abstract sets for the tree based on this status field. The first 6 sets partition the $\text{spNodes}(t, C)$, the set of internal nodes of the computation tree. The last 4 sets categorize transactions and nodes as being either active or complete.

1. $\text{pending}(t, C) = \{X \in \text{xactions}(t, C) : \text{status}[X] = \text{PENDING}\}$ (Pending transactions).
2. $\text{pendingAbort}(t, C) = \{X \in \text{xactions}(t, C) : \text{status}[X] = \text{PENDING_ABORT}\}$ (Aborting transactions).
3. $\text{committed}(t, C) = \{X \in \text{xactions}(t, C) : \text{status}[X] = \text{COMMITTED}\}$ (Committed transactions).
4. $\text{aborted}(t, C) = \{X \in \text{xactions}(t, C) : \text{status}[X] = \text{ABORTED}\}$ (Aborted transactions).
5. $\text{working}(t, C) = \{Z \in \text{spNodes}(t, C) - \text{xactions}(t, C) : \text{status}[Z] = \text{WORKING}\}$ (Working nodes).
6. $\text{synched}(t, C) = \{Z \in \text{spNodes}(t, C) - \text{xactions}(t, C) : \text{status}[Z] = \text{SYNCHED}\}$ (Synched nodes).
7. $\text{activeX}(t, C) = \text{pending}(t, C) \cup \text{pendingAbort}(t, C)$ (Active transactions).
8. $\text{activeN}(t, C) = \text{activeX}(t, C) \cup \text{working}(t, C)$. (Active nodes).
9. $\text{doneX}(t, C) = \text{committed}(t, C) \cup \text{aborted}(t, C)$ (Complete transactions).

10. $\text{doneN}(t, C) = \text{doneX}(t, C) \cup \text{synched}(t, C)$ (Complete nodes).

The *OAT* model maintains a set of **ready** S-nodes, denoted as $\text{ready}(t, C)$. We discuss the properties of ready nodes later, in Section 5.2. Note that $\text{ready}(t, C)$, and the sets defined above which are subsets of $\text{activeN}(t, C)$ are not monotonic, because completing nodes removes elements from these sets.

For the purposes of detecting conflicts, at any time t , for any active transaction T , i.e., $T \in \text{activeX}(t, C)$, the *OAT* model maintains a **read set** $R(t, T)$ and a **write set** $W(t, T)$ for T . The read set $R(t, T)$ is a set of pairs (ℓ, v) , where $\ell \in \mathcal{L}$ is a memory location and $v \in \text{memOps}(t, C)$ is a memory operation that reads from ℓ . We define $W(t, T)$ similarly. We represent main memory as the read set/write set of $\text{root}(C)$. At time $t = 0$, we assume $R(0, \text{root}(C))$ and $W(0, \text{root}(C))$ initially contain a pair (ℓ, \perp) for all locations $\ell \in \mathcal{L}$.

In addition to the basic read and write sets, we also define **module read set** and **module write set** for all transactions $T \in \text{activeX}(t, C)$. Module read set is defined as

$$\text{modR}(t, T) = \{(\ell, v) \in R(t, T) : \text{owner}(\ell) = \text{xMod}(T)\}.$$

In other words, $\text{modR}(t, T)$ is the subset of $R(t, T)$ that accesses memory owned by T 's Xmodule $\text{xMod}(T)$. Similarly, we define the **module write set** as

$$\text{modW}(t, T) = \{(\ell, v) \in W(t, T) : \text{owner}(\ell) = \text{xMod}(T)\}.$$

The *OAT* model maintains two invariants on $R(t, T)$ and $W(t, T)$. First, $W(t, T) \subseteq R(t, T)$ for every transaction $T \in \text{xactions}(t, C)$, i.e., a write also counts as a read. Second, $R(t, T)$ and $W(t, T)$ each contain at most one pair (ℓ, v) for any location ℓ . Thus, we use the shorthand $\ell \in R(t, T)$ to mean that there exists a node u such that $(\ell, u) \in R(t, T)$, and similarly for $W(t, T)$. We also overload the union operator: at some time t , an operation $R(t+1, T) \leftarrow R(t, T) \cup \{(\ell, u)\}$ means we construct $R(t+1, T)$ by

$$R(t+1, T) = \{(\ell, u)\} \cup (R(t, T) - \{(\ell, u') \in R(t, T)\}).$$

In other words, we add (ℓ, u) to $R(t, T)$, replacing any $(\ell, u') \in R(t, T)$ that existed previously.

5.2 Constructing the Computation Tree

In the *OAT* model, the runtime constructs the computation tree in a straightforward fashion as instructions are issued. For completeness, however, we give a detailed description of this construction.

Initially, at time $t = 0$, we begin with only the root node in the tree, i.e., $\text{nodes}(0, C) = \text{xactions}(0, C) = \{\text{root}(C)\}$. This root node also begins as ready, i.e., $\text{ready}(0, C) = \{\text{root}(C)\}$. Throughout the computation, the status of the root node of the tree is always PENDING.

A new internal node is created if the *OAT* model picks ready node X and X issues a fork or `xbegin` instruction. If X issues a fork, then the runtime creates a P-node P as a child of X , and two S-nodes S_1 and S_2 as children of P , all with status WORKING. The fork also removes X from $\text{ready}(C)$ and adds S_1 and S_2 to $\text{ready}(C)$. If X issues an `xbegin`, then the runtime creates a new transaction $Y \in \text{xactions}(C)$ as a child of X , with $\text{status}[Y] = \text{PENDING}$, removes X from $\text{ready}(C)$, and adds Y to $\text{ready}(C)$.

The *OAT* model completes a nontransactional S-node $Z \in \text{ready}(t, C) - \text{xactions}(t, C)$ (which must have $\text{status}[Z] = \text{WORKING}$) by having Z issue a join instruction. The join instruction first changes $\text{status}[Z]$ to SYNCHED. In the tree, since $\text{parent}(Z)$ is always a P-node, Z has exactly one sibling. If Z is the first child of $\text{parent}(Z)$ to be SYNCHED, the *OAT* model removes Z from $\text{ready}(C)$. Otherwise, Z is the last child of $\text{parent}(Z)$ to be SYNCHED, and the *OAT* model removes Z and $\text{parent}(Z)$ from $\text{ready}(C)$, changes the status of both Z and $\text{parent}(Z)$ to SYNCHED, and adds $\text{parent}(\text{parent}(Z))$ to $\text{ready}(C)$.

The *OAT* model can complete a transaction $X \in \text{ready}(t, C)$ by having it issue either an `xend` or `xabort` instruction. If $\text{status}[X] = \text{PENDING}$, then X can issue an `xend` to change $\text{status}[X]$ to COMMITTED. Otherwise, $\text{status}[X] = \text{PENDING_ABORT}$, and X can issue an `xabort` to change its status to ABORTED. For

both xend and xabort , the *OAT* model removes X from $\text{ready}(C)$ and adds $\text{parent}(X)$ back into $\text{ready}(C)$. The xend instruction also performs an ownership-aware commit and changes read sets and write sets, which we describe later in Section 5.4.

Finally, a ready node X issues a read and write instruction, if the instruction does not generate a conflict, it adds a memory operation node v to $\text{memOps}(t, C)$, with v as a child of X . If the instruction would create a conflict, the runtime may change the status of one *PENDING* transaction T to *PENDING_ABORT* to make progress in resolving the conflict. For shorthand, we refer to the status change of a transaction T from *PENDING* to *PENDING_ABORT* as a *sigabort* of T .

This construction of the tree guarantees a few properties.

First, the sequence of instructions S generated by the *OAT* model is a valid topological sort of the computation dag $G(C)$. Second, the *OAT* model generates a tree of a canonical form, where the root node of the tree is a transaction, all transactions are S-nodes and every P-node has exactly two nontransactional S-node children. This canonical form is imposed for convenience of description; it is not important for any theoretical results. Finally, the *OAT* model maintains the invariant the active nodes form a tree, with the ready nodes at the leaves. This property is important for the correctness of the *OAT* model.

5.3 Memory Operations and Conflict Detection

The *OAT* model performs eager conflict detection; before performing a memory operation that would create a new $v \in \text{memOps}(C)$, the *OAT* model first checks whether creating v would cause a conflict, according to Definition 3.

DEFINITION 3. *Suppose at time t , the *OAT* model issues a read or write instruction that potentially creates a memory operation node v . We say that v generates a **memory conflict** if there exists a location $\ell \in \mathcal{L}$ and an active transaction $T_u \in \text{activeX}(t, C)$ such that*

1. $T_u \notin \text{xAncest}(v)$, and
2. either $R(v, \ell) \wedge ((\ell, u) \in W(t, T_u))$, or $W(v, \ell) \wedge ((\ell, u) \in R(t, T_u))$.

If a potential memory operation v would generate a conflict, then the memory operation v does not occur; instead, a *sigabort* of some transaction may occur. We describe the mechanism for aborts in Section 5.5. Otherwise, v does not generate a conflict and observes the value ℓ from $R(Y)$, where Y is the closest ancestor of v with ℓ in its readset (i.e., $(\ell, u) \in R(Y)$ and $\Phi(v) = u$). The read also adds v to X 's readset. A successful write operation v sets the observer function $\Phi(v)$ in the same way as a read. The write adds (ℓ, v) to both $R(X)$ and $W(X)$.

5.4 Ownership-Aware Transaction Commit

The *ownership-aware commit mechanism* employed by the *OAT* model 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 X-modules in a closed-nested fashion, merging those reads and writes into X 's read and write sets.

We can describe the *OAT* model's commit mechanism more formally in terms of module readsets and writesets. Suppose at time t , $Y \in \text{xactions}(t, C)$ with $\text{status}[Y] = \text{PENDING}$ issues an xend . This xend changes readsets and writesets as follows.

$$\begin{aligned}
 R(\text{root}(C)) &\leftarrow R(\text{root}(C)) \cup \text{modR}(Y) \\
 R(\text{xparent}(Y)) &\leftarrow R(\text{xparent}(Y)) \cup (R(Y) - \text{modR}(Y)) \\
 W(\text{root}(C)) &\leftarrow W(\text{root}(C)) \cup \text{modW}(Y) \\
 W(\text{xparent}(Y)) &\leftarrow W(\text{xparent}(Y)) \cup (W(Y) - \text{modW}(Y))
 \end{aligned}$$

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.

LEMMA 4. *Given a computation tree C , for any $T \in \text{xactions}(C)$, let $S_T = \{\text{xMod}(T') : T' \in \text{xAnces}(T)\}$. Then $\text{modAnces}(\text{xMod}(T)) \subseteq S_T$.*

PROOF. We prove this fact by induction on the nesting depth of transactions T in the computation tree. In the base case, the top-level transaction $T = \text{root}(C)$, and $\text{xMod}(\text{root}(C)) = \text{world}$. Thus, the fact holds trivially. For the inductive step, assume that $\text{modAnces}(\text{xMod}(T)) \subseteq S_T$ 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 $\text{modParent}(\text{xMod}(T^*)) \in \text{modAnces}(\text{xMod}(T))$. Thus, $\text{modAnces}(\text{xMod}(T^*)) \subseteq \text{modAnces}(\text{xMod}(T)) \cup \{\text{xMod}(T^*)\}$. By construction of the set S_T , we have $S_{T^*} = S_T \cup \{\text{xMod}(T^*)\}$. Therefore, using the inductive hypothesis, we have $\text{modAnces}(\text{xMod}(T^*)) \subseteq S_{T^*}$. \square

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

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

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

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 Definition 2, Rule 1, we know $\text{owner}(\ell) \in \text{modAnces}(\text{xMod}(T))$. We know $\text{modAnces}(\text{xMod}(T)) \subseteq S_T$ by Lemma 4. Thus, there exists some transaction $T^* \in \text{xAnces}(T)$ 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 xids strictly decrease walking up the tree from T . Thus, there can only be one ancestor transaction T^* of T with $\text{xid}(\text{xMod}(T^*)) = \text{xid}(\text{owner}(\ell))$.

To check the second condition of Theorem 5, consider any $X \in \text{pAnces}(T^*) \cap \text{xactions}(C)$. By Rule 1, X can access ℓ directly only if $\text{owner}(\ell) \in \text{modAnces}(\text{xMod}(X))$ implying that $\text{xid}(\text{owner}(\ell)) \leq \text{xid}(\text{xMod}(X))$. 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 ℓ 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 ℓ ; thus, it is “safe” for T^* to commit ℓ .

5.5 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`. Later, it will be useful to refer to the set of operations a transaction X issues while its status is `PENDING_ABORT`.

DEFINITION 4. *The set of operations issued by X or descendants of X after $\text{status}[X]$ changes to `PENDING_ABORT` are called X 's **abort actions**. This set is denoted by $\text{abortactions}(X)$.*

The `PENDING_ABORT` status 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]). First, we describe the definition of serializability in the transactional computation framework, as given in [2]. Next, we incorporate Xmodules into this definition and define serializability by modules. We then prove that the *OAT* model guarantees serializability by modules. Finally, we discuss the relationship between the definition of serializability by modules, and the notion of abstract serializability for the methodology of open nesting.

6.1 Transactional Computations and Serializability

In [2], serializability for a transactional computation with computation tree \mathcal{C} was defined in terms of topological sorts \mathcal{S} of the computation dag $G(\mathcal{C})$. Informally, a trace (\mathcal{C}, Φ) is serializable if there exists a topological sort order \mathcal{S} of $G(\mathcal{C})$ such that \mathcal{S} is “sequentially consistent with respect to Φ ”, and all transactions appear contiguous in the order \mathcal{S} . In this section, we give more precise, formal definitions of this concept.

Content Sets

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 \mathcal{C} that one visits when walking up the tree from u to T .

DEFINITION 5. 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, $\text{cContent}(T)$ is the set of operations that are neither committed in an “open” manner, nor aborted.

Sequential Consistency with Transactions

For computations with transactions, we can modify the classic notion of sequential consistency to account for transactions which abort. 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 6. 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 definition of serializability by modules requires that computations satisfy some notion of sequential consistency, generalized for the setting of TM.

DEFINITION 7. Consider a trace (C, Φ) and a topological sort S of $G(C)$. For all $v \in \text{memOps}(C)$ such that $R(v, \ell) \vee W(v, \ell)$, the **transactional last writer** of v according to S , denoted $X_S(v)$, is the unique $u \in \text{memOps}(C) \cup \{\perp\}$ that satisfies four conditions:

1. $W(u, \ell)$,
2. $u <_S v$,
3. $\neg(uHv)$, and
4. $\forall w (W(w, \ell) \wedge (u <_S w <_S v)) \Rightarrow wHv$.

DEFINITION 8. A trace (C, Φ) is **sequentially consistent** if there exists a topological sort S such that $\Phi = X_S$. We say that S is **sequentially consistent with respect to Φ** .

In other words, the transactional last writer of a memory operation u which accesses location ℓ , is the last write v to location ℓ in the order S , except we skip over writes w which are hidden from (i.e., aborted with respect to) u . Intuitively, Definition 8 requires that there exists an order S explaining all the memory operations of the computation.

Serializability

DEFINITION 9. A trace (C, Φ) is **serializable** if there exists a topological sort S that satisfies two conditions:

1. $\Phi = X_S$ (S is sequentially consistent with respect to Φ), and
2. $\forall T \in \text{xactions}(C)$ and $\forall v \in V(C)$, we have $\text{xbegin}(T) \leq_S v \leq_S \text{xend}(T)$ implies $v \in V(T)$.

Ordinary serializability can be thought of as a strengthening of sequential consistency which also requires that the order S both explains all memory operations, and also has all transactions appearing contiguous.

6.2 Defining Serializability by Modules

In [2], a trace (C, Φ) was said to be *serializable* if there exists a topological sort S of $G(C)$ such that S is sequentially consistent with respect to Φ , and all transactions appear contiguous in S . Serializability in this context can be thought of as a 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 can not 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, Φ) , for each Xmodule M , we transform the tree C into a new tree $\text{mTree}(C, M)$. The tree $\text{mTree}(C, M)$ 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 $(\text{mTree}(C, M), \Phi)$ is serializable. If the check holds for all Xmodules, then trace (C, Φ) is said to be serializable by modules.

Definition 10 formalizes the construction of $\text{mTree}(C, M)$.

DEFINITION 10. For any computation tree C , let $\text{mTree}(C, M)$ be the result of modifying C as follows:

1. For all memory operations $u \in \text{memOps}(C)$ with u accessing ℓ , if $\text{owner}(\ell) = N$ for some $\text{xid}(N) > \text{xid}(M)$, 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 10 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, Φ) is sequentially consistent, then $(\text{mTree}(C, M), \Phi)$ is a valid trace; an operation u that remains in the trace never attempts to observe a value from a $\Phi(u)$ which was turned into a nop due to Definition 10. In addition, the transformed trace is also sequentially consistent.

LEMMA 6. Let (C, Φ) be any sequentially consistent trace. Then for any Xmodule M , $(\text{mTree}(C, M), \Phi)$ is a valid trace. In other words, if $u \in \text{memOps}(\text{mTree}(C, M))$, then $\Phi(u) \in \text{memOps}(\text{mTree}(C, M))$. Furthermore, any S which is sequentially consistent with respect to Φ in (C, Φ) is also sequentially consistent with respect to Φ in $(\text{mTree}(C, M), \Phi)$.

PROOF. In the new tree $\text{mTree}(C, M)$, pick any $u \in \text{memOps}(\text{mTree}(C, M))$ 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 10, then we know because v accessed a memory location ℓ where $\text{xid}(\text{owner}(\ell)) > \text{xid}(M)$. Since u must access the same location ℓ , u must also be converted into a nop.

If v was turned into a nop in Step 2 of Definition 10, 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 $Y_0 = \text{xparent}(v)$, $\text{xparent}(Y_i) = Y_{i+1}$, 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, Φ) is sequentially consistent and $\Phi(v) = u$, by Definition 7, we know $\neg vHu$, leading to a contradiction.

(b) If Y_j is 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 Φ is the transactional last writer according to \mathcal{S} for (C, Φ) , 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 10 is not guaranteed to generate a valid trace.

Finally, we can define serializability by modules.

DEFINITION 11. A trace (C, Φ) is **serializable by modules** if it is sequentially consistent, and if for all Xmodules 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 Φ , 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, Φ) 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 $(\text{mtree}(C, M), \Phi)$ such that all transactions of M are contiguous in \mathcal{S} .

6.3 OAT Model Guarantees Serializability by Modules

In this section, we show that the OAT model described in Section 5 generates traces (C, Φ) that are serializable by modules, i.e., that satisfy Definition 11. The proof of this fact consists of three steps. First, we generalize the notion of “prefix race-freedom” described in [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.

Defining Prefix Race-Freedom

First, we define prefix races. These definitions are essentially the same as those in [2], except adapted for a system with an ownership-aware commit mechanism instead of an open-nested commit mechanism.

DEFINITION 12. For any execution order \mathcal{S} , for any transaction $T \in \text{xactions}(C)$, consider any $v \notin \text{memOps}(T)$ such that $\text{xbegin}(T) <_{\mathcal{S}} v <_{\mathcal{S}} \text{xend}(T)$. We say there exists a **prefix race** between T and v if there exists a memory operation $w \in \text{cContent}(T)$ s.t., $w <_{\mathcal{S}} v$, $\neg(vHw)$, v and w both access ℓ , and one of v, w writes to ℓ .

DEFINITION 13. A trace (C, Φ) is **prefix race-free** iff exists a topological sort \mathcal{S} of $G(C)$ satisfying two conditions:

1. $\Phi = \mathcal{X}_{\mathcal{S}}$ (\mathcal{S} is sequentially consistent with respect to Φ), and
2. $\forall v \in V(C)$ and $\forall T \in \text{xactions}(C)$ there is no prefix race between v and T .

\mathcal{S} is called a **prefix race-free sort** of the trace.

Properties of the OAT Model

Second, we prove several invariants that OAT model preserves, and then use these invariants to prove that the OAT model generates only traces (C, Φ) which are prefix race-free.

The sequence of instructions that the OAT model issues naturally generates a topological sort \mathcal{S} of the computation dag $G(C)$: the fork and xbegin instructions correspond to the begin nodes of a parallel or series blocks in the dag, the join, xend, and xabort instructions correspond to end nodes of parallel or series blocks, and the read or write instructions correspond to memory operation nodes $v \in \text{memOps}(C)$.

THEOREM 7. Suppose the OAT model generates a trace (C, Φ) and an execution order \mathcal{S} . Then, $\Phi = \mathcal{X}_{\mathcal{S}}$, i.e., \mathcal{S} is sequentially consistent with respect to Φ .

PROOF. This result is reasonably intuitive, but the proof is tedious and somewhat complicated. We defer the details of this proof to Appendix A. \square

Next, we describe an invariant on readsets and writesets that the *OAT* model maintains. Informally, Lemma 8 states that, if a memory operation u that reads (writes) location ℓ is in the $\text{cContent}(T)$ for some transaction T , then ℓ belongs to the read set (write set) of some active transaction under T 's subtree between the time when the memory operation is performed and the time when T ends.

LEMMA 8. *Suppose the OAT model generates a trace (C, Φ) with an execution order S . For any transaction T , consider a memory operation $u \in \text{cContent}(T)$ which accesses memory location ℓ at step t_0 . Let t_f be step when $\text{xend}(T)$ or $\text{xabort}(T)$ happens. At any time t such that $t_0 \leq t < t_f$ there exists some $T' \in \text{xDesc}(T) \cap \text{activeX}(t, C)$ (i.e., T' is an active transactional descendant of T) such that*

1. If $R(u, \ell)$, then $\ell \in \text{R}(t, T')$.
2. If $W(u, \ell)$, then $\ell \in \text{W}(t, T')$.

PROOF. Let X_1, X_2, \dots, X_k be the chain of transactions from $\text{xparent}(u)$ up to, but not including T , i.e., $X_1 = \text{xparent}(u)$, $X_j = \text{xparent}(X_{j-1})$, and $\text{xparent}(X_k) = T$. Since we assume that $u \in \text{cContent}(T)$ and since T completes at time t_f , for every j such that $1 \leq j < k$, there exists a unique time t_j (satisfying $t_0 \leq t_j < t_f$) when an xend changes $\text{status}[X_j]$ from PENDING to COMMITTED; otherwise, we would have $u \in \text{aContent}(T)$.

Also, by Theorem 5 and Definition 5, we know $\text{committer}(u) \in \text{xAnces}(T)$, i.e., none of the X_j 's will commit location ℓ in an open-nested fashion to the world; otherwise, we would have $u \in \text{oContent}(T)$.

First, suppose $R(u, \ell)$. At time t_i , when the memory operation u completes, (ℓ, u) is added to $\text{R}(X_1)$. In general, at time t_j , the ownership-aware commit mechanism, as described in Section 5.4, will propagate ℓ from $\text{R}(X_j)$ to $\text{R}(X_{j+1})$. Therefore, for any time t in the interval $[t_{j-1}, t_j)$, we know $\ell \in \text{R}(t, X_j)$, i.e., for Lemma 8, $T' = X_j$. Similarly, for any time t in the interval $[t_k, t_f)$, we have $\ell \in \text{R}(t, T)$, i.e., we choose $T' = T$.

The case where $W(u, \ell)$ is completely analogous to the case of $R(u, \ell)$, except we have both $\ell \in \text{R}(t, T')$ and $\ell \in \text{W}(t, T')$. \square

Angelina: Check this.

We use Theorem 7 and Lemma 8 to prove that the *OAT* model generates traces which are prefix race-free.

THEOREM 9. *Suppose the OAT model generates a trace (C, Φ) with an execution order S . Then S is an prefix race-free sort of (C, Φ) .*

PROOF.

For the first condition of Definition 13, we know by Theorem 7 that the *OAT* model generates an order S which is sequentially consistent with respect to Φ .

To check the second condition, assume for contradiction that we have an order S generated by the *OAT* model, but there exists a prefix race between a transaction T and a memory operation $v \notin \text{memOps}(T)$. Let w be the memory operation from Definition 12, i.e., $w \in \text{cContent}(T)$, $w <_S v <_S \text{xend}T$, $\neg(vHW)$, w and v access the same location ℓ , with one of the accesses being a write. Let t_w and t_v be the time steps in which operations w and v occurred, respectively, and let $t_{\text{end}T}$ be the time at which either $\text{xend}(T)$ or $\text{xabort}(T)$ occurs (i.e., either T commits or aborts). We argue that at time t_v , the memory operation v should not have succeeded because it generated a conflict.

There are three cases for v and w . First suppose $W(v, \ell)$ and $R(w, \ell)$. Since $t_w < t_v < t_{\text{end}T}$, by Lemma 8, at time t_v , ℓ is in the writeset of some active transaction $T' \in \text{desc}(T)$. Since $v \notin \text{memOps}(T)$, we know $T \notin \text{ances}(v)$. Thus, since T' is a descendant of T , we have $T' \notin \text{ances}(v)$. Since $T' \notin \text{ances}(v)$, by Definition 3, at time t_v , v generates a conflict with T' . The other two cases, where $R(v, \ell) \wedge W(w, \ell)$ or $W(v, \ell) \wedge W(w, \ell)$, are analogous. \square

Prefix Race-Freedom Implies Serializability by Modules

Finally, we show that a trace (C, Φ) which is prefix race-free is also serializable by modules.

THEOREM 10. *Any trace (C, Φ) which is prefix race-free is also serializable by modules.*

PROOF.

First, by Definition 10 and Lemma 6, it is easy to see that a prefix-race free sort \mathcal{S} of a trace (C, Φ) is also prefix-race free of the sort $(\text{mtree}(C, M), \Phi)$ for any Xmodule M . Now we shall argue that for any Xmodule M , we can transform \mathcal{S} into \mathcal{S}_M such that all transactions in $\text{xactions}(M)$ appear contiguous in \mathcal{S}_M .

Consider a prefix-race free sort \mathcal{S} of $(\text{mtree}(C, M), \Phi)$ which has k nodes v which violate the second condition of Definition 11. We show how to construct a new order \mathcal{S}' which is still a prefix race-free sort of $(\text{mtree}(C, M), \Phi)$, but which has only $k - 1$ violations.

We reduce the number of violations according to the following procedure:

1. Of all transactions $T \in \text{modXactions}(M)$ such that there exists an operation v such that $\text{xbegin}(T) \leq_{\mathcal{S}} v \leq_{\mathcal{S}} \text{xend}(T)$ and $v \notin V(T)$, choose the $T = T^*$ which has the latest $\text{xend}(T)$ in the order \mathcal{S} .
2. In T^* , pick the first $v \notin V(T^*)$ which causes a violation.
3. Create a new sort \mathcal{S}' by moving v to be immediately before $\text{xbegin}(T^*)$.

In order to argue that \mathcal{S}' is still a prefix race-free sort of $(\text{mtree}(C, M), \Phi)$, we need to show that moving v does not generate any new prefix races, and does not create a sort \mathcal{S}' which is no longer sequentially consistent with respect to Φ (i.e., that Φ is still the transactional last writer according to \mathcal{S}'). There are three cases: v can be a memory operation, an $\text{xbegin}(T')$, or an $\text{xend}(T')$.

1. Suppose v is a memory operation which accesses location ℓ . For all operations w such that $\text{xbegin}(T) <_{\mathcal{S}} w <_{\mathcal{S}} v$, we argue that w can not access the same location ℓ unless both w and v read from ℓ . Since we chose v to be the first memory operation such that $\text{xbegin}(T) <_{\mathcal{S}} v <_{\mathcal{S}} \text{xend}(T)$ such that $v \notin V(T)$, we know $w \in V(T)$. We know by construction of $\text{mtree}(C, M)$, that $w \in \text{cContent}(T)$ (if $w \in \text{oContent}(T)$ or $w \in \text{aContent}(T)$, then steps 1 or 2, respectively, in Definition 10 will turn w into a nop). Therefore, by Definition 12, unless w and v both read from ℓ , v has a prefix race with T , contradicting the fact that \mathcal{S} is a prefix race-free sort of the trace. Thus, moving v to be before $\text{xbegin}(T)$ can not generate any new prefix races or change the transactional last writer for any memory operation, and \mathcal{S}' is still a prefix race-free sort of the trace.
2. Next, suppose $v = \text{xbegin}(T')$. Moving $\text{xbegin}(T')$ can not generate any new prefix races with T' , because the only memory operations u which satisfy $\text{xbegin}(T) <_{\mathcal{S}} u <_{\mathcal{S}} \text{xbegin}(T')$ satisfy $u \notin \text{cContent}(T')$. Also, moving $\text{xbegin}(T')$ does not change the transactional last writer for any node v because the move preserves the relative order of all memory operations. Therefore, \mathcal{S}' is still a prefix race-free sort.
3. Finally, suppose $v = \text{xend}(T')$. By moving $\text{xend}(T')$ to be before $\text{xbegin}(T)$, we can only lose prefix races with T' that already existed in \mathcal{S} because we are moving nodes out of the interval $[\text{xbegin}(T'), \text{xend}(T')]$. Also, as with $\text{xbegin}(T')$, moving $\text{xend}(T')$ does not change any transaction last writers. Therefore, \mathcal{S}' is still a prefix race-free sort of the trace.

Since we can eliminate violations of the second condition of Definition 11 one at a time, we can construct a sort \mathcal{S}_M which satisfies serializability by modules by eliminating all violations. \square

Finally, we can prove the OAT model guarantees serializability by modules by putting the previous results together.

THEOREM 11. *Any trace (C, Φ) generated by the OAT model is serializable by modules.*

PROOF. By Theorem 9, the *OAT* model generates only trace (C, Φ) which are prefix race-free. By Theorem 6.3, any trace (C, Φ) which is prefix race-free is serializable by modules. \square

6.4 Abstract Serializability

By Theorem 11, 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 [13] assumes that the program is divided into levels, and that a transaction at level i can only call a transaction at level $i + 1$.⁹ 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 can not. 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 can not 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 ℓ .¹⁰ Locks associated with an Xmodule M are owned by `modParent(M)`. A module M is said to be **properly locked** if the following is true for all transactions T_1, T_2 with `xMod(T1) = xMod(T2) = 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 S in Definition 11 for Xmodule M with `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 operations before changing the status to `ABORTED`. Therefore, an aborting transaction can compensate for all its aborted subtransactions.

⁹We 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 (e.g. [13]), levels typically go in the opposite direction.

¹⁰More complicated locks can be modeled by generalizing the definition of conflict.

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 ℓ belonging to $\text{modAncest}(\text{xMod}(T))$ only if ℓ 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 the *OAT* model exhibits a semantic deadlock if it causes the TM system to enter a state in which it is impossible to “finish” a computation because of transaction conflicts. A computation might not finish for other reasons, such as an infinite loop or livelock. This section defines semantic deadlock precisely and distinguishes it from these other reasons for noncompletion.

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} can not complete the program no matter what it does. To eliminate programs which have infinite loops, we only consider ***bounded programs***.

DEFINITION 14. *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 run forever if it **livelocks**. We use the notion of a **schedule** to distinguish livelocks from semantic deadlocks.

DEFINITION 15. *A **schedule** Γ 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 16. *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 Γ 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 17. *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 ℓ and $\text{owner}(\ell) \in \text{modAnces}(\text{xMod}(T))$, 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 17 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 can not read from (write to) any location ℓ which is owned by any `Xmodules` which are ancestors of $\text{xMod}(T)$ (including $\text{xMod}(T)$ itself), unless ℓ 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 12, which states that a location ℓ 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 12. *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 ℓ 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 13. *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 ℓ 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 $M \in \text{modAnces}(\text{xMod}(T))$, or $\text{xid}(M) > \text{xid}(\text{xMod}(T))$. If $M \in \text{modAnces}(\text{xMod}(T))$, then by Definition 17, T already had ℓ in its read or write set. Therefore, using Definition 3, 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 ℓ in its read or write set. Therefore, from Corollary 12, $\text{xMod}(T') \in \text{modDesc}(M)$. Thus, we have $\text{xid}(\text{xMod}(T')) > \text{xid}(M) > \text{xid}(\text{xMod}(T))$. \square

THEOREM 14. *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 Γ 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 , $\text{status}[T] = \text{PENDING_ABORT}$ for all active transactions T . Otherwise, the first phase of the schedule Γ 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 Γ with k phases, numbered $k-1, k-2, \dots, 1, 0$. 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 xid larger than i . During phase i , we finish aborting all active transaction instances T with $\text{xid}(\text{xMod}(T)) = i$. By Lemma 13, 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 Γ 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 can not 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 ownership-aware transactions.

In [11], Ni, et al. 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 or 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 in [5] describe a technique of transactional boosting which allows transactions to call methods from a nontransactional 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 MIT-CSAIL-TR-2008-038, Laboratory of Computer Science and Artificial Intelligence, Massachusetts Institute of Technology, June 2008. Available online at <http://supertech.csail.mit.edu/~angelee/safeTech.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, 2003.
- [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.

A. THE OAT MODEL AND SEQUENTIAL CONSISTENCY

This appendix contains the details of the proof of Theorem 7: if the *OAT* model generates a trace (C, Φ) and a topological sort order S , then S satisfies Definition 8, i.e., S is sequentially consistent with respect to Φ .

In this appendix, we first define some useful notation for the proof. Next, we prove that the *OAT* model preserves several invariants about memory operations, read set, and write sets. Finally, we use these invariants to prove Theorem 7.

A.1 Notation

We define some notation that is useful later for stating operational invariants of the *OAT* model.

For any subset S of nodes in the computation tree C , i.e., $S \subseteq \text{nodes}(C)$, define

- $\text{low}(S) = \{X \in S : \text{pDesc}(X) \cap S = \emptyset\}$.
- $\text{high}(S) = \{X \in S : \text{pAnces}(X) \cap S = \emptyset\}$.

Intuitively, $\text{low}(S)$ represents the nodes in S closest to the leaves of the tree. Similarly, $\text{high}(S)$ represents the nodes in S closest to the root of the tree. In cases where the set S is guaranteed to fall along one root-to-leaf path in the tree, we define $\text{lowest}(S)$ as the only element $X \in \text{low}(S)$. Similarly, we define $\text{highest}(S)$ as the only element in $\text{high}(S)$.

We also define two time-dependent sets of transactions.

- The **reader set** $\text{readers}(t, \ell) = \{T \in \text{activeX}(t, C) : \ell \in \text{R}(t, T)\}$.
- The **writer set**, $\text{writers}(t, \ell) = \{T \in \text{activeX}(t, C) : \ell \in \text{W}(t, T)\}$.

Said differently, $\text{readers}(t, \ell)$ is the set of active transactions at time t which have location ℓ in their read set. Similarly, $\text{writers}(t, \ell)$ is the set of active transactions at time t with $\ell \in \text{W}(T)$.

Next, we generalize the content sets from Definition 5 and define a set of dynamic content sets.

DEFINITION 18. *At any time t , for any transaction $T \in \text{xactions}(t, C)$ and a memory operation $u \in \text{memOps}(t, C)$, define the sets $\text{cContent}(t, T)$, $\text{oContent}(t, T)$, $\text{aContent}(t, T)$, and $\text{vContent}(t, T)$ according to the $\text{ContentType}(t, u, T)$ procedure:*

```

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

```

The difference between Definition 18 and the previous statement in Definition 5 is that for dynamic content sets, if we encounter a PENDING or PENDING_ABORT transaction when walking up the tree from a memory operation u to a transaction T , we place u in the *active content* of T , i.e., $u \in \text{vContent}(t, T)$. If a transaction T completes at time t^* , it is not hard to see that the dynamic classification $\text{ContentType}(t, u, T)$ gives the same answer as the static classification $\text{ContentType}(u, T)$ for all times $t \geq t^*$.

A.2 OAT Model Invariants

Because the *OAT* model performs eager conflict detection according to Definition 3, it is not hard to prove the following invariant about the readers and writers to a particular memory location ℓ .

THEOREM 15. *At all times t , for all memory locations $\ell \in \mathcal{L}$, the OAT maintains the following invariants on the sets $\text{readers}(\ell)$ and $\text{writers}(\ell)$:*

1. For all $\ell \in \mathcal{L}$, $|\text{low}(\text{writers}(t, \ell))| = 1$, i.e., $\text{lowest}(\text{writers}(t, \ell))$ exists.
2. For any $T \in \text{readers}(t, \ell)$, either $\text{lowest}(\text{writers}(t, \ell)) \in \text{desc}(T)$ or $T \in \text{desc}(\text{lowest}(\text{writers}(t, \ell)))$.

PROOF. The proof is by induction on the instructions that the *OAT* model issues.

In the base case, for all locations $\ell \in \mathcal{L}$, we begin with $\text{readers}(0, \ell) = \text{writers}(0, \ell) = \{\text{root}(C)\}$, and no other nodes in the computation tree C except $\text{root}(C)$. Thus, Invariants 1 and 2 are satisfied.

In the inductive step, suppose at time $t - 1$, Invariants 1 and 2 are satisfied. A read or write instruction at time t can not break the invariants without causing a conflict according to Definition 3. Therefore, successful read and write operations preserve the invariant. An unsuccessful read or write operation can only trigger the `sigabort` of transactions, which does not affect either invariant.

An `xend` instruction that commits a transaction T can only add the transaction $\text{xparent}(T)$ to $\text{readers}(\ell)$ or $\text{writers}(\ell)$. Since $\text{xparent}(T)$ is an ancestor of T , it can not break either of the two invariants.

The remaining instructions preserve Invariants 1 and 2 trivially. A `fork` or `join` instruction at time t preserves the invariants because they do not change the set active transactions or any transaction read sets or write sets. An `xbegin` preserves the invariants because it creates new transactions T with empty read sets and write sets. The `xabort` instruction preserves the invariants because it can only remove transactions from $\text{readers}(t, \ell)$ or $\mathbb{W}(t, \ell)$. \square

The following invariant shows that, informally, the read sets of transactions act as caches for pairs (ℓ, u) stored in write sets.

LEMMA 16. *At any time t , for any $T \in \text{readers}(t, \ell)$, suppose $(\ell, u) \in \mathbb{R}(t, T)$. Let $T' = \text{lowest}(\text{xAnces}(T)) \cap \text{writers}(t, \ell)$. Then $(\ell, u) \in \mathbb{W}(t, T')$.*

PROOF. The proof is by induction on the instructions issued by the *OAT* model. In the base case, we know for all memory locations $\ell \in \mathcal{L}$, we start with $\text{readers}(0, \ell) = \text{writers}(0, \ell) = \{\text{root}(C)\}$ and $\mathbb{R}(\text{root}(C)) = \mathbb{W}(\text{root}(C))$. Since $T' = T = \text{root}(C)$, Lemma 16 is satisfied in the base case.

For the inductive step, assume the lemma is satisfied at time $t - 1$. We show after any *S*-node X issues an instruction at time t , the lemma is still satisfied.

For any $T \in \text{xactions}(t - 1, C)$, after a `fork`, `join`, or `xbegin` instruction in step t , we have $\mathbb{R}(t, T) = \mathbb{R}(t - 1, T)$ and $\mathbb{W}(t, T) = \mathbb{W}(t - 1, T)$. Thus, the lemma is satisfied after these instructions. An `xbegin` which creates a new transaction X at time step t starts with $\mathbb{R}(t, X) = \mathbb{W}(t, X) = \emptyset$; thus, the lemma is satisfied.

Next, consider an `xabort` issued by $X \in \text{xactions}(t - 1, C)$. Suppose, before the `xabort` of X there exists a transaction $T \in \text{readers}(t - 1, \ell)$ with $(\ell, u) \in \mathbb{R}(t - 1, T)$. Let $T' = \text{lowest}(\text{xAnces}(T)) \cap \text{writers}(t - 1, \ell)$. Then before the `xabort`, $(\ell, u) \in \mathbb{W}(t - 1, T')$. Assume for contradiction after the `xabort` of X , that there exists some transaction $T \in \text{xactions}(t, C)$ such that the invariant no longer holds for T , i.e., we no longer have $(\ell, u) \in \mathbb{W}(t, T')$. Since an `xabort` does not change the contents of any transaction's write set, but removes X from $\text{writers}(\ell)$, the only way to violate the invariant is if $X = T'$. Consider two cases: either $X = T' = T$, or $X = T' \neq T$. In the first case, we can not violate the invariant for T because T is aborted and removed from $\text{readers}(\ell)$. In the second case, we must have $T \in \text{pDesc}(X)$. But then, before the `xabort`, we have $T \in \text{pDesc}(X) \cap \text{activeN}(t - 1, C)$ and $X \in \text{ready}(t - 1, C)$, contradicting the property that the ready nodes are the leaves of tree of active nodes. Thus, the `xabort` must preserve the invariant.

A successful read operation v observes the value from the closest transactional ancestor X which has location ℓ in its read set. The only transaction whose read set changes is $\text{xparent}(v)$. The invariant is preserved because $\text{xAnces}(\text{xparent}(v)) \supseteq \text{xAnces}(X)$, and since the read does not change any write sets.

A successful write operation v only changes the write set of $\text{xparent}(v)$; this write can not break the invariant without generating a conflict.

Finally, suppose at time t , a ready node X issues an `xend`. Consider two cases:

1. $X \neq \text{owner}(\ell)$. The only transaction Y which has its read set or write set change after the `xend` (i.e., for which we could have $\mathbb{R}(t, Y) \neq \mathbb{R}(t - 1, Y)$ or $\mathbb{W}(t, Y) \neq \mathbb{W}(t - 1, Y)$) is $Y = \text{xparent}(X)$. The `xend` merges

X 's read and write sets into Y 's read and write sets, respectively; using Theorem 15, it is straightforward to show that the invariant is preserved for Y .

For all other transactions $T \in \text{readers}(t, \ell)$ with $T \neq Y$, since the read set or write set of T or any transaction in $\text{xAnces}(T)$ remains the same, the invariant is still preserved for T .

2. Suppose $X = \text{owner}(\ell)$. Then, the only transaction whose read set or write set can change is $Y = \text{root}(C)$. But the only way to break the invariant is if X commits a pair (ℓ, v) from $W(t-1, X)$ to $\text{root}(C)$, which corrupts the version $(\ell, u) \in R(t-1, T)$, for some transaction T parallel to X . But then, we would violate Theorem 15, and should have had a conflict earlier.

Since all possible choices for action $k+1$ preserve the invariant, the lemma holds by induction. \square

Theorem 17 characterizes when a transaction should have a location in its write set.

THEOREM 17. *At any time t , consider any transaction $T \in \text{activeX}(t, C)$ and any memory location ℓ such that $\text{xid}(\text{owner}(\ell)) \leq \text{xMod}(T)$. Let $S_\ell(t) = \{u \in \text{memOps}(t, C) : W(u, \ell)\}$. Exactly one of the following cases holds:*

1. $T = \text{root}(C)$, $(\ell, \perp) \in W(t, T)$, and two conditions are satisfied:
 - (a) $\text{cContent}(t, T) \cap S_\ell = \emptyset$.
 - (b) For all $v \in S_\ell(t)$, we have $v \in \text{aContent}(t, T) \cup \text{vContent}(t, T)$.
2. There exists an $(\ell, u) \in W(t, T)$ which happens at time t_u , and two conditions are satisfied:
 - (a) $u \in \text{cContent}(t, T) \cap S_\ell(t)$
 - (b) For any operation $v \in (S_\ell(t) - \{u\})$ which happens at time t_v , where $t_u < t_v \leq t$, we have $v \in \text{aContent}(t, T) \cup \text{vContent}(t, T)$.
3. We have $\ell \notin W(t, T)$, and $\text{cContent}(t, T) \cap S_\ell(t) = \emptyset$.

PROOF.

This theorem can be proved by a straightforward, albeit tedious, induction on time.

Note that because we assume $\text{xid}(\text{owner}(\ell)) \leq \text{xMod}(T)$, $S_\ell(t) \cap \text{oContent}(t, T) = \emptyset$, i.e., the theorem is only concerned with memory locations ℓ which belong to T 's open content. Because of the properties of ownership and Xmodules, any location ℓ with $\text{xid}(\text{owner}(\ell)) > \text{xMod}(T)$ can never propagate into T 's write set anyway. \square

The intuition for Theorem 17 lies mostly in Case 2; if at time t a pair (ℓ, u) is the write set of a transaction T , then u is the last write to ℓ in T 's subtree which is "committed with respect to" T . Any v which writes to ℓ after t_u (the time u occurs) must belong to T 's subtree; otherwise, there would have been a conflict. Furthermore, any v which happens after t_u must still be aborted or pending with respect to T (i.e., $v \in \text{aContent}(t, T) \cup \text{vContent}(t, T)$); otherwise, v should replace u in T 's write set.

Case 3 says the write set of T does not contain a location ℓ if no memory operation in T 's subtree commits ℓ to T . Case 1 of Theorem 17 handles the special case of the root.

A.3 Proof of Sequential Consistency

Finally, we can use the invariants from Lemma 16 and Theorem 17 to prove Theorem 7.

PROOF. [Theorem 7]

The first condition and second conditions are true by construction, since the *OAT* model can only set $\Phi(v) = u$ if $u <_S v$, $W(u, \ell)$ and $R(v, \ell) \vee W(v, \ell)$.

To check the third and fourth conditions, we require some setup. Suppose at time t_v , memory operation v happens and the *OAT* model sets $\Phi(v) = u$. Let $A = \text{lowest}(\text{readers}(t, \ell) \cap \text{ances}(v))$. Because the *OAT* model sets $\Phi(v) = u$, we must have $(\ell, u) \in R(t, A)$. Let $T = \text{lowest}(\text{xAnces}(A) \cap \text{writers}(t, \ell))$. By Lemma 16, we know $(\ell, u) \in W(t, T)$. By Theorem 17, since $(\ell, u) \in W(t, T)$, we know $u \in \text{cContent}(t, T)$. Let $X = \text{xLCA}(u, v)$. We must have $T \in \text{ances}(X)$; otherwise, we could not have $\{u, v\} \subseteq \text{memOps}(t, T)$.

Since $u \in \text{cContent}(t, T)$, we know $u \in \text{cContent}(t, X) \cup \text{oContent}(t, X)$. Therefore, we have $\neg(uHv)$, satisfying the third condition.

To check the fourth condition, assume for contradiction that there exists a w such that $W(w, \ell)$, and $u <_S w <_S v$. Let t_v be the time that v happens. Then, since $\Phi(v) = u$, we know $u \in W(t_v, T)$. Therefore, by Theorem 17 we know $w \in \text{memOps}(t_v, T)$, $w \in \text{aContent}(t_v, T) \cup \text{vContent}(t_v, T)$.

Let $Y = \text{xLCA}(w, v)$. Since $w \in \text{memOps}(t_v, T)$, we know $T \in \text{ances}(Y)$. Consider the two cases for w :

1. Suppose $w \in \text{aContent}(t_v, T)$. Since $T \in \text{ances}(Y)$, we know $w \in \text{cContent}(t_v, Y) \cup \text{aContent}(t_v, Y)$.

We can show by contradiction that we must have $w \in \text{aContent}(t_v, Y)$. If $Y = T$, then we already know $w \in \text{aContent}(t_v, Y)$. Otherwise, assume $T \in \text{pAnces}(Y)$. If we had $w \in \text{cContent}(t_v, Y)$, then by Theorem 17, we must have $(\ell, y) \in W(t_v, Y)$. This statement contradicts the fact that *OAT* model found (ℓ, u) from transaction T , since a closer transaction Y had ℓ in its read set.

But then, since $w \in \text{aContent}(t_v, Y)$, we have wHv .

2. Suppose $w \in \text{vContent}(t_v, T)$:

Then, we know $w \in \text{cContent}(t_v, Y) \cup \text{vContent}(t_v, Y)$. As in the previous case, we can show $w \notin \text{cContent}(t_v, Y)$.

If $w \in \text{vContent}(t_v, Y)$, then there exists some transaction $Z \in \text{activeX}(t_v, Y) - \{Y\}$ such that $\ell \in W(t_v, Z)$.

Since $w \in \text{memOps}(t_v, Z)$, we can strengthen this condition to $Z \in \text{activeX}(t_v, \text{LCA}(w, v)) - \{\text{LCA}(w, v)\}$.

This statement leads to a contradiction, however, because $w \in W(t_v, Z)$ must conflict with v .

More formally, by Invariant 2 of Theorem 15, any new read operation v at time t_v must satisfy $v \in \text{desc}(\text{low}(\text{writers}(t_v, \ell)))$ (i.e., v is a descendant of the base of the spine for ℓ). At time t_v , however, we must have $\text{low}(\text{writers}(t_v, \ell)) \in \text{desc}(Z)$.

□

B. RULES FOR TYPE CHECKING

This appendix contains the type rules for the *OAT* type system. The syntax for the type system is shown below. Within the syntax, we do not differentiate classes that are Xmodule types from the classes that are not. The additional restrictions that apply to the Xmodule types are specified as one of the type rules.

For simplicity, in our type system, we make the following assumptions. First, each class has only one constructor (specified by the term *init*), and that all fields are initialized properly after the call to the constructor. Second, all field names (whether inherited or declared) are distinct. Third, the call to *super* is explicit. Fourth, an index is always specified when the ownership tags *world* and *this* are used. Fifth, the class names *Object* and *Xmodule* are special and assumed to be properly defined by the system. Finally, the explicit use of *upcast* and *downcast* are not allowed, as specified in the abstract syntax.

$$\begin{aligned}
 P &= \text{defn}^*; e \\
 \text{defn} &= \text{class } cDecl \text{ extends } cDecl \text{ where } \text{constr}^* \{ \text{field}^*; \text{init}; \text{meth}^* \} \\
 cDecl &= \text{cn}\langle \text{formal}^+ \rangle \mid \text{Object}\langle \text{formal} \rangle \mid \text{Xmodule}\langle \text{formal} \rangle \\
 \text{constr} &= \text{formal} < \text{formal} \mid \text{formal} = \text{formal} \mid \text{formal} \neq \text{formal} \\
 \text{field} &= t \text{ fd} \\
 \text{init} &= \text{cn}\langle \text{formal}^+ \rangle(\text{param}^*) \{ \text{super}\langle \text{formal}^+ \rangle(e^*); \text{this.fld} = e;^* \} \\
 \text{meth} &= t \text{ mn}\langle \text{formal}^* \rangle(\text{param}^*) \text{ where } \text{constr}^* \{ e \} \\
 \text{param} &= t \ x \\
 \text{owner} &= \text{world}[i] \mid \text{formal} \mid \text{this}[i] \\
 \text{formal} &= f \\
 t &= \text{int} \mid \text{ct} \\
 \text{ct} &= \text{cn}\langle \text{owner}^+ \rangle \\
 e &= \text{new } \text{ct}\langle e^* \rangle \mid x \mid x = e \mid \text{let } (\text{param} = e) \text{ in } \{ e \} \mid x.\text{fld} \mid x.\text{fld} = e \mid x.\text{mn}\langle \text{owner}^+ \rangle(e^*) \\
 \\
 \text{cn} &= \text{a class name that is not Object nor Xmodule} \\
 \text{mn} &= \text{a method name that is not a constructor} \\
 \text{fd} &= \text{a field name} \\
 x, y &= \text{a variable name} \\
 f, g &= \text{an owner formal} \\
 i, j &= \text{an int literal}
 \end{aligned}$$

For the constraints on owners (*constr*), the notation $<$ is used as defined in Section 3: Assuming f_1 and f_2 are instantiated with o_1 and o_2 , $f_1 < f_2$ specifies that either $o_1.name \prec o_2.name$, or $o_1.name = o_2.name$ and $o_1.index < o_2.index$. Similarly, $f_1 = f_2$ specifies that $o_1.name = o_2.name$ and $o_1.index = o_2.index$. On the other hand, $f_1 \neq f_2$ specifies that either $o_1.name \neq o_2.name$, or $o_1.name = o_2.name$ and $o_1.index \neq o_2.index$.

In the type system, we use a few predicates. Before we define the predicates, we first define some notations:

Henceforth, for brevity, we use the notation \triangleleft in place of the keyword *extends* (i.e. *A extends B* is written as $A \triangleleft B$). We also use the notation \trianglelefteq between class names as the reflexive and transitive closure induced by

the \triangleleft relation. Note that the \triangleleft is not the same as subtyping (denoted as $<:$), because \triangleleft only considers the static relation defined by the `extends` keyword, and does not account for the ownership tags.

In addition, we define $field \in_d cn\langle \dots \rangle$ to mean that class $cn\langle \dots \rangle$ declares $field$, $field \in_i cn\langle \dots \rangle$ to mean that class $cn\langle \dots \rangle$ inherits $field$, and $field \in cn\langle \dots \rangle$ to mean that either $field \in_d cn\langle \dots \rangle$ or $field \in_i cn\langle \dots \rangle$. We use these notations for fd (field name), $meth$ (method), and mn (method name) similarly.

Now we define the predicates.

Predicate	Meaning
$ClassOnce(P)$	No class is declared twice in P $\forall cn, cn' \text{ in } P, cn \neq cn'$
$FieldsOnce(P)$	No class contains two fields with the same name $\forall ct \forall fd, fd' \in ct \text{ in } P, fd \neq fd'$
$MethodsOnce(P)$	No class declares two methods with the same name $\forall ct \forall mn, mn' \in_d ct \text{ in } P, mn \neq mn'$
$WFClasses(P)$	There are no cycles in the class hierarchy; i.e. the \trianglelefteq relation is antisymmetric $\forall cn, cn' \text{ in } P, cn \trianglelefteq cn' \wedge cn' \trianglelefteq cn \implies cn = cn'$

Our typing judgment has the form: $P; \Gamma \vdash e : t$, where P is the program being checked to provide information about class definitions; Γ is the typing environment, providing mappings from a variable name to its static type for the free variables in e ; finally, t is the static type of e .

The typing environment Γ is defined as

$$\Gamma ::= \emptyset \mid \Gamma, x : t \mid \Gamma, f : owner \mid \Gamma, constr$$

That is, the typing environment Γ contains the types of variables, the owner parameters and the constraints among owners. When checking for well-formedness of the typing environment, we assume the new entries are checking in the order listed, from left to right.

The typing system uses the following judgments.

Judgment	Meaning
$\vdash P : t$	program P yields type t
$P \vdash defn$	$defn$ is a well-formed class
$P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle g_{1..k} \rangle$	class $cn\langle f_{1..n} \rangle$ extends class $cn'\langle g_{1..k} \rangle$
$P \vdash cn \trianglelefteq cn'$	cn' is an ancestor of cn in the graph defined by the <code>extends</code> keyword
$P \vdash field \in_d cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ declares $field$
$P \vdash field \in_i cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ inherits $field$
$P \vdash field \in cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ declares / inherits $field$
$P \vdash init \in cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ declares $init$
$P \vdash meth \in_d cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ declares $meth$
$P \vdash meth \in_i cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ inherits $meth$
$P \vdash meth \in cn\langle \dots \rangle$	class $cn\langle \dots \rangle$ declares / inherits $meth$
$P; \Gamma \vdash field$	$field$ is a well-formed field
$P; \Gamma \vdash meth$	$meth$ is a well-formed method
$P; \Gamma \vdash wf$	typing environment Γ is well-formed
$P; \Gamma \vdash t$	t is a well-formed type
$P; \Gamma \vdash constr$	constraint $constr$ is satisfied
$P; \Gamma \vdash_{owner} o$	o is an owner
$P; \Gamma \vdash e : t$	expression e has type t
$P; \Gamma \vdash t <: t'$	t is a subtype of t'

In the type rules, we also use the following auxiliary rules:

The Extends Relation

$$\frac{P \vdash \text{class } cn\langle f_{1..n} \rangle \text{ extends } cn'\langle g_{1..m} \rangle \dots}{P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle g_{1..m} \rangle}$$

$$\frac{}{P \vdash cn \trianglelefteq cn} \quad \frac{P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle g_{1..m} \rangle}{P \vdash cn \trianglelefteq cn'} \quad \frac{P \vdash cn \trianglelefteq cn' \quad P \vdash cn' \trianglelefteq cn''}{P \vdash cn \trianglelefteq cn''}$$

Type Lookup

$$\begin{aligned} \text{type}() &= () \\ \text{type}(t \ x) &= t \\ \text{type}(t \ fd) &= t \\ \text{type}(t_1 \ x_1, t_2, x_2, \dots) &= t_1, t_2, \dots \\ \text{type}(t \ mn\langle g_{1..k} \rangle(\text{param}^*)\{ \dots \}) &= t \rightarrow \langle g_{1..k} \rangle \rightarrow \text{type}(\text{param}^*) \end{aligned}$$

Field Lookup

$$\frac{P \vdash \text{class } cn\langle f_{1..n} \rangle \dots \{ \dots \text{field} \dots \}}{P \vdash \text{field} \in_d cn\langle f_{1..n} \rangle} \quad \frac{P \vdash \text{field} \in cn'\langle g_{1..m} \rangle \quad P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle o_{1..m} \rangle}{P \vdash \text{field} [o_1/g_1]..[o_m/g_m] \in_i cn\langle f_{1..n} \rangle}$$

$$\frac{P \vdash \text{field} \in_d cn\langle f_{1..n} \rangle \vee P \vdash \text{field} \in_i cn\langle f_{1..n} \rangle}{P \vdash \text{field} \in cn\langle f_{1..n} \rangle}$$

Init Lookup

$$\frac{P \vdash \text{class } cn\langle f_{1..n} \rangle \dots \{ \dots \text{init} \dots \}}{P \vdash \text{init} \in cn\langle f_{1..n} \rangle}$$

Method Lookup

$$\frac{P \vdash \text{class } cn\langle f_{1..n} \rangle \dots \{ \dots \text{meth} \dots \}}{P \vdash \text{meth} \in_d cn\langle f_{1..n} \rangle} \quad \frac{P \vdash \text{meth} \in cn'\langle g_{1..m} \rangle \quad P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle o_{1..m} \rangle}{P \vdash \text{meth} [o_1/g_1]..[o_m/g_m] \in_i cn\langle f_{1..n} \rangle}$$

$$\frac{P \vdash \text{meth} \in_d cn\langle f_{1..n} \rangle \vee P \vdash \text{meth} \in_i cn\langle f_{1..n} \rangle}{P \vdash \text{meth} \in cn\langle f_{1..n} \rangle}$$

Override Ok

$$\frac{\begin{aligned} P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle o_{1..m} \rangle \\ P \vdash \text{meth} \in_d cn\langle f_{1..n} \rangle \\ P \vdash \text{meth} \notin cn'\langle g_{1..m} \rangle \end{aligned}}{\text{OverrideOk}(cn\langle f_{1..n} \rangle, cn'\langle o_{1..m} \rangle, \text{meth})} \quad \frac{\begin{aligned} P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle o_{1..m} \rangle \\ P \vdash \text{meth} \in_d cn\langle f_{1..n} \rangle \\ P \vdash \text{meth}' \in cn'\langle g_{1..m} \rangle \\ \text{type}(\text{meth}) = \text{type}(\text{meth}') [o_1/g_1]..[o_m/g_m] \end{aligned}}{\text{OverrideOk}(cn\langle f_{1..n} \rangle, cn'\langle o_{1..m} \rangle, \text{meth})}$$

We present the type rules next.

$\boxed{\vdash P : t}$

[PROG]

$$\frac{WFClasses(P) \quad ClassOnce(P) \quad FieldsOnce(P) \quad MethodsOnce(P)}{P = defn_{1..n}; e \quad P \vdash defn_i \quad P; \emptyset \vdash e : t} \vdash P : t$$

$\boxed{P \vdash defn}$

[CLASS]

$$\frac{\begin{array}{c} P \vdash cn \not\leq Xmodule \\ \Gamma = f_{1..n} : owner, f_1 < f_i, constr^*, this : cn\langle f_{1..n} \rangle \\ P; \Gamma \vdash wf \quad P; \Gamma \vdash cn'\langle f_1, o^* \rangle \quad P; \Gamma \vdash field_i \quad P; \Gamma \vdash init \quad P; \Gamma \vdash meth_i \\ OverrideOk(cn\langle f_{1..n} \rangle, cn'\langle f_1, o^* \rangle, meth_i) \end{array}}{P \vdash class cn\langle f_{1..n} \rangle \text{ extends } cn'\langle f_1, o^* \rangle \text{ where } constr^* \{ field^*; init; meth^* \}}$$

[XMODULE CLASS]

$$\frac{\begin{array}{c} P \vdash cn \leq Xmodule \\ \Gamma = f_{1..n} : owner, f_1 < f_i, constr^*, this : cn\langle f_{1..n} \rangle, this : owner, this[i] < f_1 \\ P; \Gamma \vdash wf \quad P; \Gamma \vdash cn'\langle f_1, o^* \rangle \quad P; \Gamma \vdash field_i \quad P; \Gamma \vdash init \quad P; \Gamma \vdash meth_i \\ type(field_i) \neq int \quad OverrideOk(cn\langle f_{1..n} \rangle, cn'\langle f_1, o^* \rangle, meth_i) \end{array}}{P \vdash class cn\langle f_{1..n} \rangle \text{ extends } cn'\langle f_1, o^* \rangle \text{ where } constr^* \{ field^*; init; meth^* \}}$$

$\boxed{P; \Gamma \vdash field}$

[FIELD]

$\boxed{P; \Gamma \vdash init}$

[INIT]

$$\frac{\begin{array}{c} P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle f_1, o_{2..m} \rangle \\ \Gamma' = \Gamma, param^* \quad P; \Gamma' \vdash wf \quad P; \Gamma' \vdash this.fd_j = e_j \\ P; \Gamma \vdash t \quad P \vdash init\langle g_{1..m} \rangle(t_i \ x_i^{i \in 1..k}) \in cn'\langle g_{1..m} \rangle \quad P; \Gamma' \vdash e_i : t_i [f_1/g_1][o_2/g_2]..[o_m/g_m] \end{array}}{P; \Gamma \vdash t \quad P; \Gamma \vdash cn\langle f_{1..n} \rangle(param^*) \{ super\langle f_1, o_{2..m} \rangle(e_i^{i \in 1..k}); \quad this.fd = e;^* \}}$$

$\boxed{P; \Gamma \vdash meth}$

[METHOD]

$$\frac{\begin{array}{c} \Gamma^c = f_{1..n} : owner, constr^* \quad P; \Gamma^c \vdash wf \\ \Gamma' = \Gamma, \Gamma^c, param^* \quad P; \Gamma' \vdash wf \quad P; \Gamma' \vdash e : t \end{array}}{P; \Gamma \vdash t \ mn\langle f_{1..n} \rangle(param^*) \text{ where } constr^* \{ e \}}$$

$$\boxed{P; \Gamma \vdash wf}$$
[ENV \emptyset]

[ENV X]

[ENV OWNER]

$$\frac{}{P; \emptyset \vdash wf} \quad \frac{P; \Gamma \vdash t \quad x \notin \text{Dom}(\Gamma) \quad P; \Gamma \vdash wf}{P; \Gamma, x : t \vdash wf} \quad \frac{f \notin \text{Dom}(\Gamma) \quad P; \Gamma \vdash wf}{P; \Gamma, f : \text{owner} \vdash wf}$$

[ENV CONSTR]

$$\frac{\begin{array}{c} \text{constr} = (o < o') \vee (o = o') \vee (o \neq o') \\ P; \Gamma \vdash wf \quad P; \Gamma \vdash_{\text{owner}} o, o' \quad \Gamma' = \Gamma, \text{constr} \\ \exists_{x,y} (P; \Gamma' \vdash x < y) \wedge (P; \Gamma' \vdash y < x) \\ \exists_{x,y} (P; \Gamma' \vdash x < y) \wedge (P; \Gamma' \vdash x = y) \quad \exists_{x,y} (P; \Gamma' \vdash x = y) \wedge (P; \Gamma' \vdash x \neq y) \end{array}}{P; \Gamma, \text{constr} \vdash wf}$$

$$\boxed{P; \Gamma \vdash t}$$

[TYPE INT]

[TYPE OBJECT]

[TYPE XMODULE]

$$\frac{}{P; \Gamma \vdash \text{int}} \quad \frac{P; \Gamma \vdash_{\text{owner}} o}{P; \Gamma \vdash \text{Object}\langle o \rangle} \quad \frac{P; \Gamma \vdash_{\text{owner}} o}{P; \Gamma \vdash \text{Xmodule}\langle o \rangle}$$

[TYPE CT]

$$\frac{P \vdash \text{class } cn\langle f_{1..n} \rangle \dots \text{ where } \text{constr}^* \dots \quad P; \Gamma \vdash_{\text{owner}} o_i \quad P; \Gamma \vdash o_1 < o_i \quad P; \Gamma \vdash \text{constr} [o_1/f_1]..[o_n/f_n]}{P; \Gamma \vdash cn\langle o_{1..n} \rangle}$$

$$\boxed{P; \Gamma \vdash \text{constr}}$$

[CONSTR ENV]

[< WORLD I]

[< WORLD II]

[< THIS]

$$\frac{\Gamma = \Gamma', \text{constr}, \Gamma''}{P; \Gamma \vdash \text{constr}} \quad \frac{P; \Gamma \vdash_{\text{owner}} o \quad P; \Gamma \vdash o \neq \text{world}}{P; \Gamma \vdash o < \text{world}[i]} \quad \frac{i < j}{P; \Gamma \vdash \text{world}[i] < \text{world}[j]} \quad \frac{i < j \quad P; \Gamma \vdash_{\text{owner}} \text{this}}{P; \Gamma \vdash \text{this}[i] < \text{this}[j]}$$

[< TRANS]

[= WORLD]

[= THIS]

[= TRANS]

$$\frac{P; \Gamma \vdash o_1 < o_2 \quad P; \Gamma \vdash o_2 < o_3}{P; \Gamma \vdash o_1 < o_3} \quad \frac{i = j}{P; \Gamma \vdash \text{world}[i] = \text{world}[j]} \quad \frac{i = j \quad P; \Gamma \vdash_{\text{owner}} \text{this}}{P; \Gamma \vdash \text{this}[i] = \text{this}[j]} \quad \frac{P; \Gamma \vdash o_1 = o_2 \quad P; \Gamma \vdash o_2 = o_3}{P; \Gamma \vdash o_1 = o_3}$$

[= REFL]

$$\frac{P; \Gamma \vdash_{owner} o \quad P; \Gamma \vdash o \neq world \quad P; \Gamma \vdash o \neq this}{P; \Gamma \vdash o = o}$$

[≠ WORLD]

$$\frac{i \neq j}{P; \Gamma \vdash world[i] \neq world[j]}$$

[≠ THIS]

$$\frac{i \neq j \quad P; \Gamma \vdash_{owner} this}{P; \Gamma \vdash this[i] \neq this[j]}$$

[≠ WORLD]

$$\frac{P; \Gamma \vdash_{owner} this[i]}{P; \Gamma \vdash this[i] \neq world}$$

[SUBSTITUTION]

$$\frac{P; \Gamma \vdash o_1 = o_2 \quad P; \Gamma \vdash constr}{P; \Gamma \vdash constr [o_1/o_2]}$$

[RELATION]

$$\frac{P; \Gamma \vdash o_1 < o_2}{P; \Gamma \vdash o_1 \neq o_2}$$

 $P; E \vdash_{owner} o$

[OWNER WORLD]

[OWNER FORMAL]

[OWNER THIS]

$$\frac{}{P; \Gamma \vdash_{owner} world[i]}$$

$$\frac{\Gamma = \Gamma', f : owner, \Gamma''}{P; \Gamma \vdash_{owner} f}$$

$$\frac{\Gamma = \Gamma', this : owner, \Gamma''}{P; \Gamma \vdash_{owner} this[i]}$$

 $P; E \vdash e : t$

[EXP TYPE]

[EXP SUB]

[EXP NEW]

[EXP VAR]

$$\frac{P; E \vdash t \quad P; E \vdash e : t' \quad P; E \vdash t' <: t}{P; E \vdash e : t}$$

$$\frac{P; \Gamma \vdash cn\langle o_{1..n} \rangle \quad init = cn\langle f_{1..n} \rangle(t_i \ x_i^{i \in 1..k}) \{ \dots \} \quad P \vdash init \in cn\langle f_{1..n} \rangle \quad P; \Gamma \vdash e_i : t_i [o_1/f_1]..[o_n/f_n]}{P; \Gamma \vdash new \ cn\langle o_{1..n} \rangle(e_i^{i \in 1..k}) : cn\langle o_{1..n} \rangle}$$

$$\frac{\Gamma = \Gamma', x : t, \Gamma''}{P; \Gamma \vdash x : t}$$

[EXP VAR ASSIGN]

[EXP LET]

[EXP REF]

$$\frac{P; \Gamma \vdash x : t \quad P; \Gamma \vdash e : t}{P; \Gamma \vdash x = e : t}$$

$$\frac{param = t' \ x \quad P; \Gamma \vdash e' : t' \quad P; \Gamma, param \vdash wf \quad P; \Gamma, param \vdash e : t}{P; E \vdash let (param = e') in \{ e \} : t}$$

$$\frac{P; \Gamma \vdash x : cn\langle o_{1..n} \rangle \quad P \vdash (t \ fd) \in cn\langle f_{1..n} \rangle}{P; \Gamma \vdash x.f d : t [o_1/f_1]..[o_n/f_n]}$$

[EXP REF ASSIGN]

[EXP INVOKE]

$$\frac{P; \Gamma \vdash x : cn\langle o_{1..n} \rangle \quad P \vdash (t \ fd) \in cn\langle f_{1..n} \rangle \quad P; \Gamma \vdash e : t [o_1/f_1]..[o_n/f_n]}{P; \Gamma \vdash x.f d = e : t [o_1/f_1]..[o_n/f_n]}$$

$$\frac{P \vdash t \ mn\langle f_{(k+1)..n} \rangle(t_i \ y_i^{i \in 1..h}) \text{ where } constr^* \dots \in cn\langle f_{1..k} \rangle \quad P; \Gamma \vdash x : cn\langle o_{1..k} \rangle \quad P; \Gamma \vdash e_i : t_i [o_1/f_1]..[o_n/f_n] \quad P; \Gamma \vdash constr [o_{k+1}/f_{k+1}]..[o_n/f_n]}{P; \Gamma \vdash x.mn\langle o_{(k+1)..n} \rangle(e_{1..h}) : t [o_1/f_1]..[o_n/f_n]}$$

$$\boxed{P; \Gamma \vdash t <: t'}$$

[SUBTYPE]

[SUBTYPE TRANS]

[SUBTYPE REFL]

$$\frac{\begin{array}{l} P; \Gamma \vdash cn\langle o_{1..n} \rangle \\ P \vdash cn\langle f_{1..n} \rangle \triangleleft cn'\langle f^+ \rangle \end{array}}{P; \Gamma \vdash cn\langle o_{1..n} \rangle <: cn'\langle f^+ \rangle[o_1/f_1]..[o_n/f_n]}$$

$$\frac{\begin{array}{l} P; \Gamma \vdash t <: t' \\ P; \Gamma \vdash t' <: t'' \end{array}}{P; \Gamma \vdash t <: t''}$$

$$\frac{P; \Gamma \vdash t}{P; \Gamma \vdash t <: t}$$