

Bibliography

UNIVERSITY OF CALGARY

Title of Thesis

Second Thesis Title Line

by

Student's name

A THESIS

SUBMITTED TO THE FACULTY OF GRADUATE STUDIES
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE
DEGREE OF NAME OF DEGREE IN FULL

DEPARTMENT OF NAME OF DEPARTMENT

CALGARY, ALBERTA

monthname, year

© Student's name year

Abstract

Acknowledgements

Table of Contents

Abstract	i
Acknowledgements	ii
Table of Contents	iii
List of Tables	iv
List of Figures	v
List of Symbols	vi
1 Introduction	1
1.1 Related Work	1
1.2 Statement of Results	1
2 Computational Model and Definitions	2
2.1 Computational Model	2
2.2 Definitions	3
3 Analysis	4
3.1 Correctness	4
3.1.1 Definitions and Terminology	4
4 The Adaptive To-Do Tree	9
4.0.2 Analysis and Proofs	11
4.1 Performance	15
4.1.1 DoTask Analysis	15
4.1.2 InsertTask Analysis	15

List of Tables

List of Figures and Illustrations

List of Symbols, Abbreviations and Nomenclature

Symbol

Definition

U of C

University of Calgary

Chapter 1

Introduction

In this chapter, we will introduce the task allocation problem as well as its dynamic version and survey the past work related to it.

Task allocation problem, also called do-all problem, is already a foundational topic in distributing computing. During the last two decades, significant research was dedicated to studying the task allocation problem in various models of computation, including message passing, shared memory, etc. under specific assumption about the asynchrony and failures.

In this thesis, we consider the dynamic version of the task allocation problem via randomized asynchronous shared memory. The dynamic version could be described as follows:

p processes cooperatively perform a set of tasks in the presence of adversity and the tasks are injected dynamically by the adversary during the execution.

1.1 Related Work

1.2 Statement of Results

Chapter 2

Computational Model and Definitions

In this chapter, we will formally describe the computational model and give the definitions.

2.1 Computational Model

The computational model we consider is the standard asynchronous shared memory model with n processes, each process has a unique identifier $i \in 1, \dots, n$. The set of processes is denoted $\{p_1, p_2, \dots, p_n\}$.

Each process executes its program by taking steps until it has reached a special state where it takes no further action. A step is defined to be the execution of shared memory operation on shared objects and receive corresponding response. Up to $n - 1$ processes may fail by crashing.

In this thesis, we consider the system that supports atomic compare-and-swap (CAS) object. A CAS object v is a shared memory object stores a value from some set and supports two atomic synchronization operations $v.read()$ and $v.CAS(x, y)$. Operation $v.read()$ returns the current value of v and leaves the value unchanged. A $v.CAS(x, y)$ operation writes new value into it only if it matches pre-supplied “expected” value x , i.e, if current value of v equals x , then operation $v.CAS(x, y)$ succeeds, and the value of v is changed to be y and *true* is returned. Otherwise, the operation fails, and the current value of v remains unchanged and *false* is returned.

A process can execute local coin flip operation that returns an integer value distributed uniformly at random from an arbitrary finite set of integers. In the following discussion, we use $random(s)$ to return a value which is distributed uniformly at random from set $\{0, 1, 2, \dots, s - 1\}$. We analyze our algorithm under the assumption of a strong adaptive

adversary. At any point of time, it can see the entire past history and know the states of all processes. Based on this information, it decides which process takes the next step.

2.2 Definitions

In the next section, we will first build the formal definition and terminology for linearizability, and then prove that our concurrent dynamic task allocation object is linearizable.

2.2.1 Definitions and Terminology

Definition 1. *History.* *A history (or execution) H , obtained by processes executing shared memory operations on concurrent objects, is a finite sequence of method invocation and response events.*

A history H may be related to several concurrent objects. When we consider the concurrency behavior of a specific object, we should focus on the subhistory of that object which is defined as follows:

Definition 2. *$H|obj$ of history H is the subsequence of all invocation and response events in H whose object names are all obj .*

If all invocation and response events in a history H have the same object name obj , then the $H|obj = H$. Thus, in the following discussion, when we discuss the concurrency behavior of a specific object obj , the history H and $H|obj$ are the same.

We use $inv_H(op)$ to denote the position of the invocation of operation op in history H . Similarly we use $rep_H(op)$ to denote the position of the corresponding response of op in H . Actually, for each operation op , it is not necessary to have a matching response. In this case, we call the operation op is pending and denote it as $rsp_H(op) = \infty$. Otherwise, we call the operation is complete in H .

Definition 3. Complete/Incomplete History. A history is complete if all its operations are complete. A completion of an incomplete history H is an extension H' of H , such that H' contains exactly the same operations as H but the responses are added for the pending operations in H .

Let H be a complete history. With each operation op in H we could associate a time interval $I_H(op) = [inv(op), rsp(op)]$. Similarly, for an uncomplete history, we denote the interval with respect to the pending operation op by $I_H(op) = [inv(op), \infty]$.

Definition 4. Sequential History. A history is sequential if the first event in the history is an invocation, and each invocation, except possibly the last one, is immediately followed by a corresponding response.

It is obvious that a sequential history defines a total order over all operations.

Definition 5. Sequential Specification. The sequential specification S_{obj} for a concurrent object obj is a set of sequential histories on obj .

We say a sequential history S satisfies the sequential specification S_{obj} of object obj , it means $S \in S_{obj}$. A sequential history S of object obj is valid (sometimes also called legal) if it satisfies the sequential specification of obj .

With the above definitions, we could now define linearization as follows:

Definition 6. Linearization. A history H linearizes to a sequential history S , if and only if S satisfies the following conditions: (1) S and the completion of H have the same operations, (2) sequential history S is valid, and (3) there is a mapping from each time interval $I_H(op)$ to a time point $t_H(op) \in I_H(op)$, such that the sequential history S could be obtained by sorting the operations in H by their $t_H(op)$ values.

A history is linearizable if and only if there exists a sequential history S that linearizes H . In this case, S is called linearization of H .

We choose one linearization of H and this linearization defines $t_H(op)$. For each operation op in history H , we call the point $t_H(op)$ linearization point of op .

Definition 7. *Linearizable Object.* *A concurrent object is linearizable if every history H of that object is linearizable.*

In the following subsection, we will prove that the concurrent dynamic task allocation object shown in Figure 1 is linearizable.

Every shared object has a *type* τ which can be defined mathematically ([?]) as follows,

$$\tau = (\mathcal{S}, s_{init}, \mathcal{O}, \mathcal{R}, \delta)$$

where \mathcal{S} is a set of states, $s_{init} \in \mathcal{S}$ is the initial state, \mathcal{O} is a set of operations, \mathcal{R} is the set of responses, and $\delta : \mathcal{S} \times \mathcal{O} \rightarrow \mathcal{S} \times \mathcal{R}$ is a state transition mapping.

Sequential specification of type τ describes, for each operation, the condition under which that operation can be invoked, and the effect produced after the operation was executed.

Our dynamic task allocation type (DTA) supports two operations, $DoTask()$ and $InsertTask(\ell)$, where ℓ is the task identifier that is unique for each task.

Now we formalize the notion of type DTA by specifying the above two operations. We assume that there exists an atomic operation $PutTask(M, \ell)$, and a process associates task ℓ with memory location M by calling $PutTask(M, \ell)$. It returns *success* if task ℓ is associated with location M , and returns *failure* if location M was already associated with another task. We say a task ℓ is inserted if it is associated to a shared memory location M .

Similarly, we assume there exists an atomic operation $TryTask(M)$, and the task ℓ associated with memory location M could be performed atomically by a process by calling $TryTask(M)$. Out of several processes calling $TryTask(M)$, only one receives *success* and the index ℓ of that task, while all the others receive *failure*. This assumption guarantees that only one process will actually perform task ℓ , i.e, each task ℓ can be performed at most once during the executing.

A task is *done* if its index has been returned by a process after calling $DoTask()$.

A task is *available* at location M if it has been inserted to M and *success* is returned by a process after calling $InsertTask(\ell)$, but is not done yet. A task is *available*, if it is *available* at some memory location.

The aim of operation $DoTask()$ is to perform an available task on location M by calling $TryTask(M)$. Every $DoTask()$ may perform several $TryTask(M)$ operations. However, only one of them will succeed. Once one $TryTask(M)$ succeeds, there is no available task on M and the task index ℓ will be returned by $DoTask()$. Additionally, if there is no available task, then operation $DoTask()$ returns \perp .

The goal of $InsertTask(\ell)$ operation is to find a free memory location M and insert task ℓ by calling $PutTask(M, \ell)$. Because $PutTask(M, \ell)$ will fail if location M has been associated with another task, each $InsertTask(\ell)$ operation may perform several $PutTask(M, \ell)$ operations. But only one of them will succeed. Once one $TryTask(M)$ succeeds, then task ℓ is available on location M and *success* notification is returned by $InsertTask(\ell)$ operation.

Type DTA is required to satisfy: (*Validity*) If a $DoTask()$ operation returns ℓ , then before that, an $InsertTask(\ell)$ was executed and returned *success*. (*Uniqueness*) Each task is performed at most once, i.e, for each task ℓ , at most one $DoTask()$ operation returns ℓ .

In addition, the property that every inserted task is eventually done is also a desired progress property of the implementation of type DTA .

Chapter 3

The Adaptive To-Do Tree

Some description here.

Method 1 DoTask()

```
1: while (true) do
2:    $v \leftarrow root$ 
3:   if ( $v.surplus() \leq 0$ ) then
4:     return  $\perp$ 
5:   end if

   /* Descent */
6:   while  $v$  is not a leaf do
7:      $(x_L, y_L) \leftarrow v.left.read()$ 
8:      $(x_R, y_R) \leftarrow v.right.read()$ 
9:      $s_L \leftarrow \min(x_L - y_L, 2^{height(v)})$ 
10:     $s_R \leftarrow \min(x_R - y_R, 2^{height(v)})$ 
11:     $r \leftarrow random(0, 1)$ 
12:    if  $((s_L + s_R) = 0)$  then
13:      Mark-up( $v$ )
14:    else if  $(r < s_L / (s_L + s_R))$  then
15:       $v \leftarrow v.left$ 
16:    else
17:       $v \leftarrow v.right$ 
18:    end if
19:  end while

   /*  $v$  is a leaf */
20:   $(x, y) \leftarrow v.read()$ 
21:   $(flag, l) \leftarrow v.TryTask(task[y + 1])$ 
22:   $v.CAS((x, y), (x, y + 1))$  // Update Insertion Count
23:   $v \leftarrow v.parent$ 
24:  Mark-up( $v$ )
25:  if  $flag = success$  then
26:    return  $\ell$ 
27:  end if
28: end while
```

Method 2 Insert(task ℓ)

```
1: while (true) do
2:    $v \leftarrow root$ 

   /* Descent */
3:   while ( $v$  is not a leaf) do
4:      $(x_L, y_L) \leftarrow v.left.read()$ 
5:      $(x_R, y_R) \leftarrow v.right.read()$ 
6:      $s_L \leftarrow \min(x_L - y_L, 2^{height(v)})$ 
7:      $s_R \leftarrow \min(x_R - y_R, 2^{height(v)})$ 
8:      $r \leftarrow random(0, 1)$ 
9:     if  $((s_L + s_R) = 0)$  then
10:      Mark-up( $v$ )
11:     else if  $(r < s_L / (s_L + s_R))$  then
12:        $v \leftarrow v.left$ 
13:     else
14:        $v \leftarrow v.right$ 
15:     end if
16:   end while

   /*  $v$  is a leaf */
17:    $(x, y) \leftarrow v.read()$ 
18:    $flag \leftarrow v.PutTask(task[x + 1], \ell)$ 
19:    $v.CAS((x, y), (x + 1, y))$  // Update Insertion Count
20:    $v \leftarrow v.parent$ 
21:   Mark-up( $v$ )
22:   if ( $flag = success$ ) then
23:     return success
24:   end if
25: end while
```

Method 3 Mark-up(v)

```
1: if ( $v$  is not Null) then
2:   for ( $i = 0; i < 2; i++$ ) do
3:      $(x, y) \leftarrow v.read()$ 
4:      $(x_L, y_L) \leftarrow v.left.read()$ 
5:      $(x_R, y_R) \leftarrow v.right.read()$ 
6:      $v.CAS((x, y), (max(x, x_L + x_R), max(y, y_L + y_R)))$ 
7:   end for
8: end if
```

3.1 Correctness

The standard correctness condition for shared memory algorithms is linearizability, which was introduced by Herlihy and Wing in 1990. The intuition of linearizability is that real-time behavior of method calls must be preserved, i.e, if one method call precedes another, then the earlier call must have taken effect before the later one. By contrast, if two method calls overlap, we are free to order them in any convenient way since the order is ambiguous. Informally, a concurrent object is linearizable if each method call appears to take effect instantaneously at some moment between its invocation and response.

3.1.1 Analysis and Proofs

By the definitions in Subsection 3.1.1, one way to show a concurrent object obj is linearizable is to prove every history H of obj is linearizable. Thus, we need to identify for each DoTask and InsertTask operation op (i.e, interval $I_H(op)$) in H a linearization point $t_H(op)$, and prove that the sequential history S obtained by sorting these operations according to their $t_H(op)$ satisfies the sequential specification S_{obj} of obj .

We notice that each complete DoTask or InsertTask operation can be associated with a unique task array slot based on the task it removed or inserted. Additionally, the removal and insertion count are both monotonically increasing. Thus, we could associate the node counts with operations which have been propagated to that node.

Now we define “an operation is counted at a node” recursively to formalize the operation propagation.

A DoTask operation is counted at leaf v when the removal count of v is updated with the index of the task array slot where the performed task is located. Symmetrically, an InsertTask operation is counted at v when the insertion count of v is updated with the index of the task array slot where the inserted task is located.

Now we only define DoTask operation is counted at an inner node v because counting an

InsertTask operation is symmetric as well.

Recall that the removal count of v is updated through CAS operation (line 6, method 3). Actually there could be more than one operations updating the count with the same value. We linearize all such CAS operations, which update the removal count of v with the same value y . We say for all these operations, only the first one in the linearization order counts the corresponding DoTask operation. In another word, a DoTask operation is counted at an inner node v as soon as the CAS updating operation that counts the DoTask is linearized. Based this definition, no operation will be counted twice at a node.

Please note that, the CAS operation counting the DoTask at node v is not necessary performed by the DoTask operation itself, i.e, suppose process p executes a DoTask operation and has successfully performed task ℓ at certain leaf. Then the CAS operation counting this $p.DoTask()$ at node v could be a different process q as long as q updates the removal count first in the linearization order.

Given the above concepts and properties, we could prove the following result:

Lemma 1. *Let v be a tree node,*

(1) If (x, y) is the return value of $v.read()$, then there exists a set of x InsertTask operations and a set of y DoTask operation that have been counted at node v by the end of the execution of $v.read()$.

(2) If there are x DoTask operations that have been counted at v before the execution of $v.read()$, then the removal count value returned by $v.read()$ is not less than x .

Proof. TBD □

Lemma 2. *Consider a history H and an arbitrary operation op in H , let $t_H(op)$ be the point when op is counted at the root, then $t_H(op)$ is between $inv_H(op)$ and $rsp_H(op)$.*

Proof. Without loss of generality, suppose process p executes DoTask operation op which has performed task ℓ successfully at the leaf. When it reaches the *root* and executes $Markup(root)$, there will be two cases:

Case 1: It increments the removal count of *root* successfully via CAS (line 6, method 3) at point τ . If it is the first one in the linearization order of all CAS operation updating the removal count with the same value, then $t_H(op) = \tau$, therefore $inv_H(op) < t_H(op) < rsp_H(op)$. Otherwise, we could let $t_H(op) = \tau'$, where τ' is the time when the first CAS operation in the linearization order updated the removal count. Thus $t_H(op) < \tau < rsp_H(op)$, Because only if the task has been performed then the removal count of root could be updated. so $t_H(op) > inv_H(op)$. Therefore, in this case, $inv_H(op) < t_H(op) < rsp_H(op)$ holds.

Case 2: It fails to increment the removal count. The CAS operation of p fails if and only if the value of the counts were updated by another process at a point before τ , suppose it is τ' . Please note that, we say the counts were updated, it means the removal count or the insertion count was updated because the two counts are stored in one memory location. Thus, there are two subcases.

Subcase 2.1: If it is the removal count that was incremented at τ' , it means the DoTask operation has already been counted by another process. Thus, $t_H(op) \leq \tau' < rsp_H(op)$.

Subcase 2.2: If it is not the removal count but the insertion count that was updated at τ' . We notice that the CAS operation (line 6, method 3) will be repeated by p , during the second iteration, if the CAS of p succeed at τ'' , then we could deduce that $t_H(op) \leq \tau''$, therefore $t_H(op) < rsp_H(op)$. If it fails again at point τ'' , then there must be another process updated the counts of root again. This time, the counts of the children will be noticed and the DoTask operation will propagate to the root. We could deduce $t_H(op) < \tau''$. Thus, $t_H(op) < rsp_H(op)$ as well. \square

Under the above definitions and properties, we claim the point $t_H(op)$ when op is counted at the root is the linearization point of operation op .

Lemma 3. *The dynamic task allocation object in the above figure is linearizable.*

Proof. Consider an arbitrary history H containing DoTask and InsertTask operations. We should prove for any execution of our algorithm, the total order given by the linearization

point is verifies the uniqueness and validity. If H is not complete, then we let all processes that have not finished their operations continue to take steps in an arbitrary order until all operations are completed. Every operation will finally be done is ensured by our computation model and the randomness of our algorithm. This way we obtain a completion H' of H and it suffices to prove H' is linearizable. Thus, to prove this lemma, we should prove the total order obtained by sorting the operations by their $t_H(op)$ values is valid.

The uniqueness is obvious. When multiple processes are calling $\text{TryTask}(\ell)$ at the memory location, only one of them will receive success and the index ℓ of the task, all the other competitor processes will get failure. Task ℓ is performed successfully as long as the value of corresponding memory location is turned to 1. The following process will never repeatedly turn it be to 0 and turn 0 to 1 which is guaranteed by the our semantics of task insertion and removing.

Now we prove the validity, i.e. each task that is performed successfully must have been inserted before. We should prove the insertion operation of a task is always counted at the root before the removal operation. To prove this result holds for the root, we now prove it by induction from the leaf.

At the leaf, this holds because if and only if the insertion count of newly inserted task ℓ has been incremented (line 19, method 2) then the following removal operation could read that (line 20, method 1) to know the available task ℓ at the leaf and then try to perform it. In another word, suppose the task ℓ is inserted but the insertion count is not incremented (i.e. insertion has not been counted yet), then the following removal operation has no way to know the available task ℓ , perform it and increment the removal count. Thus, the removal will not be counted.

For an arbitrary inner node v , we suppose, by the induction step and lemma 1, the result holds for children $v.left$ and $v.right$. We could notice that any process updates the insertion count and removal count as an atomic operation (line 6, method 3). If some remove operation

has been counted at node v , then the corresponding insert operation for that task must have been counted at v simultaneously by the double compare-and-swap operation. Apply this to the root, then validity condition holds. \square

3.2 Performance

3.2.1 DoTask Analysis

3.2.2 InsertTask Analysis

Bibliography