

Flexible Scheduling of Transactional Memory on Trees¹

Costas Busch^a, Bogdan S. Chlebus^a, Maurice Herlihy^b, Miroslav Popovic^c, Pavan Poudel^d, Gokarna Sharma^e

^aAugusta University, Augusta, Georgia, USA

^bBrown University, Providence, Rhode Island, USA

^cUniversity of Novi Sad, Novi Sad, Serbia

^dATGWORK, Norcross, Georgia, USA

^eKent State University, Kent, Ohio, USA

Abstract

We study the efficiency of executing transactions in a distributed transactional memory system. The system is modeled as a **static** network with the topology of a tree. Contrary to previous approaches, we allow the flexibility for both transactions and their requested objects to move simultaneously among the nodes in the tree. Given a batch of transactions and shared objects, the goal is to produce a schedule of executing the transactions that minimizes the cost of moving the transactions and the objects in the tree. We consider both techniques for accessing a remote object with respect to a transaction movement. In the first technique, instead of moving, transactions send control messages to remote nodes where the requested objects are gathered. In the second technique, the transactions migrate to the remote nodes where the objects are gathered to access them. When all the transactions use a single object, we give an online algorithm that produces optimal schedules for both techniques. For the general case of multiple objects per transaction, in the first technique, we obtain a schedule with a constant-factor approximation of optimal. In the second technique, with transactions migrating, we give a k factor approximation where k is the maximum number of objects per transaction.

Keywords: Distributed system, transactional memory, shared object, network, communication cost

¹ A preliminary version of this article appears in the Proceedings of SSS'22 [1].

Corresponding author. Tel.: +1 800 551 7943

Email addresses: kbusch@augusta.edu (Costas Busch), bchlebus@augusta.edu (Bogdan S. Chlebus), herlihy@cs.brown.edu (Maurice Herlihy), miroslav.popovic@rt-rk.uns.ac.rs (Miroslav Popovic), poudelpavan@gmail.com (Pavan Poudel), gsharma2@kent.edu (Gokarna Sharma)

1. Introduction

Threads executed concurrently require synchronization to prevent inconsistencies while accessing shared objects. Traditional low-level thread synchronization mechanisms such as locks and barriers are prone to deadlock and priority inversion, among multiple vulnerabilities. The concept of transactional memory has emerged as a high-level abstraction of the functionality of distributed systems; see Herlihy and Moss [2] and Shavit and Touitou [3]. The idea is to designate blocks of program code as transactions to be executed atomically. Several commercial processors support transactional memory in hardware, for example, Haswell of Intel [4], Blue Gene/Q of IBM [5], zEnterprise EC12 of IBM [6], and Power8 of IBM [7].

Transactions are executed speculatively, in the sense that if a transaction aborts due to synchronization conflicts or failures then the transaction's execution is rolled back to be restarted later. A transaction commits if there are no conflicts or failures, and its effects become visible to all processes. If multiple transactions concurrently attempt to access the same object, then this creates a conflict for access and could trigger aborting some of the involved transactions. Scheduling transactions to minimize conflicts for access to shared objects improves the system's performance.

The processing units of a distributed transactional memory system are the nodes of a communication network, which is an integral part of the system. A transaction executing at a node may want to access shared memory objects residing in other nodes. This could be implemented such that the transaction coordinates access to the needed shared objects with the nodes hosting the objects. Such systems were studied by Herlihy and Sun [8], Sharma and Busch [9], and Siek and Wojciechowski [10]. The efficiency of executing a specific transaction may reflect the topology of the communication network that is part of a distributed system. For example, the amount of communication needed to execute a transaction interacting with some objects could be proportional to the distances in the network between all the nodes hosting the transaction and the objects.

To improve efficiency of processing transactions on shared objects, we may preemptively move objects and transactions among the nodes to schedule their presence at specific nodes at specific times. Moving transactions or program code among network nodes is currently used in several real-world applications. For example, Erlang Open Telecom Platform aids dynamic code upgrade by supporting transactional servers with hot code swapping whose call-back modules may be changed on the fly [11]. A job management system for a computer cluster may migrate a job to a different node, if the target node's load is below the migration threshold and the migration overhead is acceptable, in order to achieve better load balancing among the nodes, see Hwang et al. [12]. A related system that uses live virtual machine migration to support autonomic adaptation of virtual computation

1 environments is described by Ruth et al. [13].

2 Coordinating accessing objects to execute transactions may involve relocation
3 of objects or transactions. Efficiency of such coordination may depend on additional
4 model's specification which determines the very feasibility of moving transactions
5 and objects across the network. In the data-flow model, transactions are static and
6 objects move from one node to another to reach the nodes hosting transactions that
7 require interacting with them; see Tilevich and Smaragdakis [14] and Herlihy and
8 Sun [8]. In that model, a transaction initially requests the objects it needs, and
9 executes after assembling them. After a transaction commits, it releases its ob-
10 jects, possibly forwarding them to pending transactions. In the control-flow model,
11 objects are static and transactions move from one node to another to access the ob-
12 jects. Control-flow allows transactions to send control requests, in a manner similar
13 to remote procedure calls, to the nodes where the required objects are located; see
14 Arnold et al. [15] and Saad and Ravindran [16].

15 1.1. Contributions

16 We consider a flexible scheduling approach that combines the benefits of the
17 data-flow and control-flow models. We study the dual-flow model that allows for
18 both transactions and objects to move among the nodes to synchronize transactions
19 and objects. This model combines the functionality of the data-flow and control-
20 flow models. Assessing the cost of synchronizing transactions with objects in the
21 dual-flow model takes into account the communication cost of relocating objects,
22 and the communication cost of relocating transactions or the cost of sending control
23 messages from nodes hosting transactions to nodes hosting objects.

24 We consider distributed systems whose networks interpreted as graphs have
25 [static](#) tree topologies. This represents many real-world networks. For example,
26 the internet cloud consists of the cloud network, representing a root, the fog net-
27 work gateways and/or the edge network gateways, as internal nodes, and the IoT
28 devices as leaves, see Comer [17].

29 We study the efficiency of executing transactions by a distributed system repre-
30 sented as a tree in the dual-flow model. The efficiency is measured by the cost of
31 communication. Scheduling transactions is considered in a batch setting, in which
32 all the transactions are given at the outset, subject to the constraint that each node
33 is assigned at most one original transaction. The initial position of shared objects
34 are distributed arbitrarily among the nodes. We consider scheduling transactions in
35 the general case of arbitrarily many shared objects, and also in a special case of a
36 single shared object that needs to be accessed by all the transactions. Given a batch
37 of transactions and objects residing at nodes of the system, the goal is to produce a
38 schedule of executing transactions that minimizes the cost of moving transactions
39 and objects among the nodes and sending control messages to facilitate executing

1 the transactions. Such a schedule is computed by a centralized online algorithm 2
to be executed by the distributed system. We develop a centralized algorithm find-3
ing an optimal schedule in the case when all the transactions use a single object. 4
The general case of multiple objects is studied in two models that determine if ex-5
ecuting a transaction may involve sending control messages. For multiple shared 6
objects and with transactions sending control messages, we give a centralized algo-7
rithm that finds a schedule with a constant-factor approximation of communication 8
cost with respect to an optimal schedule. For multiple shared objects and with trans-9
actions migrating and not using control messages, we give a centralized algorithm 10
that finds a schedule approximating an optimal one by a factor k that equals the
11 maximum number of shared objects requested by a transaction.

12 1.2. Related work

13 We discuss related work on data-flow, control-flow, and the dual-flow models. 14
Attiya et al. [18], Busch et al. [19, 20, 21], and Sharma and Busch [9, 22] considered 15
transaction scheduling with provable performance bounds in the data-flow model. 16
Saad and Ravindran [16], Palmieri et al. [23], Siek and Wojciechowski [24, 10] 17
studied scheduling transactions in the control-flow model. Palmieri et al. [23] also 18
gave a comparative study of data-flow versus control-flow models for distributed 19
transactional memory. A prototype distributed transactional memory system de-20
scribed by Saad and Ravindran [25] supports experimentation for both data-flow 21
and control-flow models. Bocchino et al. [26] considered the dual-flow model by 22
allowing programmers to either bring the data to the code of computation (transac-23
tion) or send the code of computation to the data. Hendler et al. [27] studied a lease 24
based dual-flow model which dynamically determines whether to migrate transac-25
tions to the nodes that own the leases or to demand the acquisition of these leases 26 by
the node that originated the transaction.

27 Transaction scheduling in a distributed system with the goal of minimizing ex-28
ecution time was first considered by Zhang et al. [28]. Busch et al. [19] consid-29
ered minimizing both the execution time and communication cost simultaneously. 30
They showed that it is impossible to simultaneously minimize execution time and 31
communication cost for all the scheduling problem instances in arbitrary graphs 32
even in the online setting. Specifically, Busch et al. [19] demonstrated a tradeo 33
between minimizing execution time and communication cost and provided online 34
algorithms optimizing execution time and communication cost separately. Busch 35
et al. [21] considered transaction scheduling tailored to specific popular topologies 36
and provided online algorithms that minimize simultaneously execution time and 37
communication cost. In a follow-up work, Poudel and Sharma [29] provided an 38
evaluation framework for processing transactions in distributed systems. Busch et 39
al. [20] studied online algorithms to schedule transactions arriving continuously.

1 Distributed directory protocols have been designed by Herlihy and Sun [8], Sharma 2
and Busch [9], and Zhang et al. [28], with the goal to optimize communication cost 3 in
scheduling transactions. A distributed directory protocol has been designed by 4 Rai
et al. [30] in the data-flow model that reduces processing load of network nodes 5 in
addition to communication cost.

6 Alternative approaches to distributed transactional memory systems have been
7 proposed in the literature by way of replicating transactional memory on multiple
8 nodes and providing means to guarantee consistency of replicas. This includes work
9 by Couceiro et al. [31], Hirve et al. [32], Kobus et al. [33], Manassiev et al. [34],
10 Peluso et al. [35], and Peluso et al. [36]. In this work, we use a single copy of each
11 object. Replicas of objects help to improve reliability of the systems rather than
12 decrease the communication overhead. Other systems extend non-distributed trans-
13 actional memory with a communication layer, for example, the system presented
14 by Kotselidis et al. [37] extends the system described by Herlihy et al. [38] with
15 distributed coherence protocols.

16 Transaction scheduling has been widely-studied in shared memory multi-core 17
systems. Scheduling algorithms with provable upper bounds, along with lower 18
bounds and impossibility results were given by Attiya et al. [39], Dragojevic et 19
al. [40], Guerraoui et al. [41], Sharma and Busch [42, 43]. Several other schedul-20
ing algorithms were evaluated only experimentally, like the ones given in Yoo and 21 Lee
[44], Baldassin et al. [45], Manassiev et al. [46], and Kolli et al. [47].

22 2. Technical Preliminaries

23 A distributed system consists of processing nodes with some pairs of nodes 24
connected by links. It is represented as a graph $G = (V; E)$. There are n vertices in 25
the set V , each representing a processing node. Edges in the set $E \subseteq V \times V$ represent 26
communication links between nodes. The function $w : E \rightarrow \mathbb{Z}^+$ assigns a weight to 27
each edge representing a communication delay. We let $\text{dist}(u; v)$ denote the shortest 28
path distance between two vertices u and v .

29 The initial configuration of the distributed system consists of a set of transac-30
tions and shared objects distributed among the nodes. Each node hosts at most one 31
transaction. During executing transactions, both shared objects and transactions 32
can move among the nodes of a network, which we call the dual-flow model. If a 33
transaction requests access to an object, that object may move to a different node, 34
possibly closer to the requesting transaction. At the same time, the transaction can 35
also migrate to the object's new location, or send a control message to that new 36
location to access the object. The combined cost of executing a transaction is mea-37
sured with relation to the distances traversed by the shared objects, transaction and 38
control messages.

1 In the dual-flow model, objects may not move to every transaction's node. In-2
 3 stead some transactions may need to access some objects at remote nodes. A trans-3
 4 action that performs multiple updates to a remote object iterates communication ex-4
 5 changes between the transaction's node and the object's node. In this case, migrat-5
 6 ing the transaction closer to the object's node decreases communication costs. To 6
 7 provide a formal framework for such situations, we consider the following two spe-7
 8 cializations of the dual-flow model for remote object access: (i) Control-message 8
 9 technique, where a transaction sends a control message to access the remote object. 9
 10 The control-message technique is motivated by a scenario in which each transaction
 11 performs a number of updates to an object bounded by a constant, with each update 11
 12 requiring a control message, for a total of a constant number of such messages. (ii) 12
 13 Transaction-migration technique, in which a transaction moves to the node where
 14 objects are located and no control messages are sent. This technique is motivated 14
 15 by the scenarios in which a transaction may issue a variable number of requests to 15
 16 an object, in which case it is advantageous to migrate the transaction to the object 16
 17 location to avoid potentially unbounded communication overhead.

17 We parameterize the costs of transmitting messages that carry transactions, ob-18
 19 jects, or control instructions. The cost of moving an object of size over a unit 19
 20 weight edge is denoted by α . We denote the cost of sending a control message over 20
 21 a unit weight edge by β . The cost of moving a transaction over a unit weight edge 21
 22 is denoted by γ .

22 A scheduling algorithm determines a schedule to execute transactions, including 23
 24 movements of objects and transactions. A centralized algorithm takes as input a 24
 25 configuration of transactions and objects in the system as arranged at the outset. 25
 26 We assume that each node has this input available so that it can execute it locally. 26
 27 A **schedule** E of executing transactions is a sequence of transactions $T_1; T_2; \dots$ that 27
 28 specifies for each transaction when to execute. Each transaction T_i is processed at 28
 29 its turn from the schedule E by accessing the required objects for T_i .

29 The communication cost of executing such a schedule is the sum of distances 30
 31 traversed by the shared objects, control messages, and transactions according to the 31
 32 schedule, weighted by the corresponding parameters α , β , and γ . The following 32
 33 example illustrates the cost of executing transactions in each of the three models 33
 34 and how dual-flow model minimizes total communication cost compared to the 34
 35 data-flow and control-flow models.

35 Consider a tree G with six nodes $v_1; v_2; \dots; v_6$ with the topology as shown in 36
 37 Figure 1. Let each edge of G has weight 1 and each node is assigned with one 37
 38 transaction requiring a single shared object o initially positioned at v_1 . Let the 38
 39 size of object o be 4 and size of a control message be 2, which means $\alpha = 4$ and 39
 40 $\beta = 2$. In the data-flow model, object o needs to visit each node of G to provide 40
 41 access to the transactions. The object may traverse the route $(v_1; v_2; v_3; v_4; v_5; v_6)$

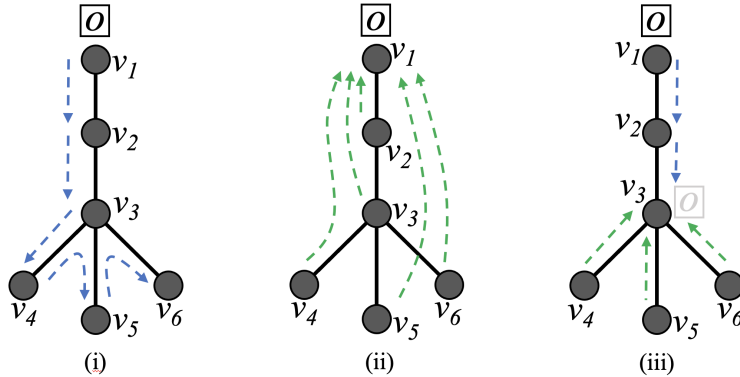


Figure 1: Illustration of communication overhead in (i) data-flow, (ii) control-flow, and (iii) dual-flow model.

1 to allow the transactions at each node to execute; see part (i) in Figure 1. The
 2 total communication cost becomes $4 + 4 + 8 + 8 + 4 = 28$, as the object moves
 3 twice in the edges between $v_3;v_4$ and $v_3;v_5$. In the control-flow model, the object
 4 resides at v_1 and all the transactions send control messages to v_1 to access o ; see
 5 part (ii) in Figure 1. The total communication cost for executing all six transactions
 6 becomes $0 + 2 + 4 + 6 + 6 + 6 = 24$. In the dual-flow model, object o may move up
 7 to node v_3 following the route $(v_1; v_2; v_3)$. The transactions at v_1 , v_2 and v_3 execute
 8 as soon as o arrives at each node. The remaining transactions at v_4 , v_5 and v_6
 9 send control messages to v_3 to access object o ; see part (iii) in Figure 1, where
 10 blue dashed arrows denote the movements of object and green dashed arrow denote
 11 the control messages sent by the transactions. The total communication cost for
 12 executing all six transactions in the dual-flow model becomes $4 + 4 + 2 + 2 + 2 = 14$
 13 which is less compared to that in both the data-flow and control-flow models.

14 3. A Single Object

15 We assume a single shared object o of size > 1 positioned at the root node of 16
 a tree G . Note that G is a static tree with a fixed set of nodes and links, and each 17
 node has a global view of the system. We develop an optimal scheduling algorithm 18
 denoted as Single-Object in the dual-flow model considering both techniques for 19
 accessing a remote object: control-message and transaction-migration. The algo-20
 rithm is called Single-Object, and its pseudocode is given as Algorithm 1.

21 3.1. Control-Message Technique

22 A general idea of the algorithm in the control-message technique is as follows. 23
 First we find a set of intermediate nodes in G to move the object o to. These nodes 24
 are referred to as supernodes. An intermediate node v becomes a supernode if the 25
 cost of moving o from v to one of its children is greater than the cost of sending

1 control messages from the transactions contained by the sub-tree of that child to v.
2 Each supernode contains a set of transactions in its sub-tree which send control
3 messages to that supernode to access object o. These transactions are added to
4 the local execution schedule of the supernode following an iterative pre-order tree
5 traversal in the sub-tree. We determine a subtree P containing paths in G that reach
6 the supernodes from the root of G. Starting from the root, object o travels all the
7 supernodes following the iterative pre-order tree traversal of P. Any transaction
8 that lies along the path is added to the execution schedule E as soon as o reaches the
9 respective node. When o reaches some supernode, the transactions from its local
10 execution schedule get added to E in the respective order. The execution ends when
11 all the transactions have been added to E. The algorithm can be modified as follows
12 if performed in the transaction-migration technique: (i) Determine supernodes with
13 respect to transaction migration cost rather than control messages cost; and (ii)
14 Migrate transactions to the corresponding supernodes instead of sending control
15 messages to access the object. These modifications result in creating an algorithm
16 of a comparable communication performance.

17 Remember that each node has a global view of the system which helps to deter-18
19 mine whether a given node is a supernode or not by comparing the possible cost of
20 control messages vs. object movement. During the computation of supernodes, no
21 object or transaction is moved or a control message is sent, thus there is no com-
22 munication cost involved. The actual communication cost will occur during the
23 execution of transactions following the respective algorithm. We elaborate on the
24 details of the algorithm next.

25 The cost of moving o over an edge of unit length is c . Let α represents the
26 control message cost for a transaction to access object o at one unit away and $\alpha > c$.
27 Let $T = \{T_1; T_2; \dots; T_n\}$ be the set of n transactions issued to the nodes of G, one
28 at each node. The first objectives are to determine the walk the object traverses and
29 to find transaction execution schedule. Intuitively, since it costs more to move the
30 object across a link than to send a control message through the link, we strive to
31 move the object minimally, only when this pays, and this approach is captured by
32 the concept of supernodes. The object o first travels from the root up to a supernode.
33 Transactions that lie along the path the object traverses execute as soon as the object
34 reaches the respective nodes. The remaining transactions beyond that supernode
35 and towards the leaves (i.e., in the sub-tree of the supernode) send control messages
36 to the supernode to access the object. Observe that the control messages are sent
37 only up to the supernode and the transactions will execute at that supernode after the
38 object reaches there. When all the transactions that have sent control messages to
39 the current supernode finish their executions, object o moves to the next supernode
40 and the remaining transactions get executed following a similar approach.

40 The communication cost of an execution of the algorithm is determined by the

1 location of supernodes. The set of supernodes is selected by referring to transaction 2
counts and transaction loads at all nodes, which are defined as follows. A transac-3
tion count at node v , denoted $\text{txnum}(v)$, is the total number of transactions contained 4
in the sub-tree of node v , including v (Line 4 of Algorithm 1). A transaction load of 5
a node v , denoted $\text{txload}(v)$, is the sum of distances from v to the positions of trans-6
actions contained in the sub-tree of v , including v (Lines 9–15 of Algorithm 1). 7
The transaction load of v represents the cost of sending control messages due to the 8
transactions contained in its sub-tree, assuming o is moved to v .

9 To identify supernodes, we start from the leaves of G and work through the an-10
cestors towards the root (Lines 16–24 of Algorithm 1). Let v_{cur} be a leaf node 11
and v_{next} be the parent of v_{cur} . During the computation of supernodes, we can 12
assume that the object is at the parent node v_{next} and check if it pays to move 13
the object down to v_{cur} , since object moves away from the root. Let $\text{txload}(v_{\text{cur}})$ 14
denote the control message cost incurred by the $\text{txnum}(v_{\text{cur}})$ number of transac-15
tions contained in the sub-tree of v_{cur} , including v_{cur} . If the object o moves to 16
 v_{cur} , the transactions contained in the sub-tree of v_{cur} can access o at v_{cur} and 17
the cost becomes $\text{txload}(v_{\text{cur}}) + \text{dist}(v_{\text{cur}}; v_{\text{next}})$. Here, $\text{dist}(v_{\text{cur}}; v_{\text{next}})$ is 18
the cost incurred by the movement of object o from v_{next} to v_{cur} . Otherwise,
19 these transactions send control messages to v_{next} to access o and the cost becomes 20
 $\text{txload}(v_{\text{cur}}) + \text{txnum}(v_{\text{cur}}) \text{dist}(v_{\text{cur}}; v_{\text{next}})$. Object o will move to v_{cur} from v_{next} 21
only if the control message cost from v_{cur} to v_{next} , due to the transactions contained 22
in the sub-tree of v_{cur} , is more than or equal to the object movement cost from v_{next} 23
to v_{cur} . This has been checked in Lines 19–23 of Algorithm 1.

24 After reaching a supernode, object o may need to move back to the root or in-25
termediate nodes to visit other supernodes. To account for this and simplify the 26
argument, we assume that the object moves over each edge twice, but this assump-27
tion will be revisited when we optimize the algorithm (Line 25 of Algorithm 1). If 28
the following inequality holds:

$$\text{txload}(v_{\text{cur}}) + 2 \text{dist}(v_{\text{cur}}; v_{\text{next}}) \leq \text{txload}(v_{\text{cur}}) + \text{txnum}(v_{\text{cur}}) \text{dist}(v_{\text{cur}}; v_{\text{next}});$$

29 then we choose v_{cur} as a supernode. Otherwise, if v_{cur} is not the root, a new pair 30
of v_{cur} and v_{next} is checked such that current v_{next} becomes new v_{cur} and the parent 31
of current v_{next} becomes a new node v_{next} . If v_{cur} is the root, then it becomes a 32
supernode (Line 24 of Algorithm 1).

33 Let P denote the pruned tree, which contains only the supernodes and nodes that 34
need to be traversed on the way from the root to a supernode. Tree P is rooted at the 35
root of G . Figure 2 illustrates such a tree P . The object o is originally located at the 36
root, from which it moves to the supernodes in a pre-order traversal manner. The 37
transactions are executed along the way of the object's movement. Transactions at

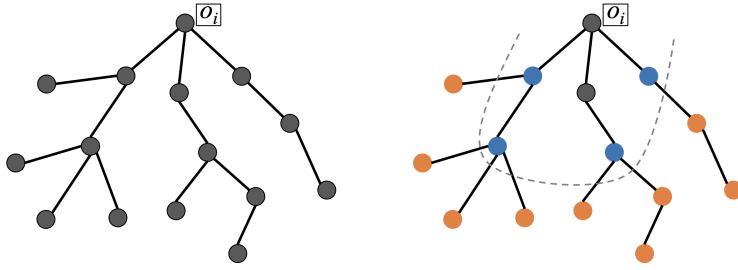


Figure 2: Identification of supernodes by algorithm Single-Object. The tree on the left is G. The tree on the right is the same G after determining the status of nodes. Supernodes are colored blue. Nodes on the path from the root to a blue node are colored black. The dashed line delineates P obtained from G by pruning G of vertices beyond the supernodes, which are colored orange.

1 the nodes beyond the pruned tree P, marked by color orange in Figure 2, either send 2
control messages or move to access o to their closest supernodes. When object o 3
reaches the respective supernode, these transactions are executed in order.

4 Lemma 1. If a node v does not belong to the pruned tree P, then the total number
5 of transactions contained in the sub-tree of v is less than 2.

6 Proof. This follows from the specification of the nodes that make P and how 7
supernodes are determined, while we assume that the object moves over each edge 8
twice.

9 After computing the set of supernodes, the object performs a pre-order tree 10
traversal starting from the root to visit all the supernodes. The transaction execution 11
schedule E is computed as follows. First add transaction at the root to E. During the 12
pre-order tree traversal to visit the supernodes, if E does not contain the transaction 13
at a visited node v , then add it to E. If the visited node v is a supernode, add to E the 14
transactions that sent control messages to v from the subtree rooted at v . Lines 27– 15 31
of Algorithm 1 represent how each transaction is added into the schedule E 16
following the pre-order tree traversal.

17 Next we show how to refine this approach, which is based on the assumption 18
that during the computation of supernodes if the object moves from some parent 19
node to the child node, then it will ultimately move back from that child node to the 20
parent. When the object reaches the last supernode, it does not move back because 21
there is no any other supernode remained to visit. A pseudocode to accomplish this 22
is given as Algorithm 2.

23 We define a one-way path to be such a path from v_{root} to the last supernode v_{last} , 24
all the edges of which the object traverses only once. This v_{last} must be chosen 25
in such a way that the total communication cost is minimized. A condition for 26
computing a supernode is:

$$2 \text{ dist}(v_{cur}; v_{next}) > \text{txnum}(v_{cur}) \text{ dist}(v_{cur}; v_{next}) \quad (1)$$

Algorithm 1: Single-Object

Input : Tree graph G of n nodes containing T transactions
 Output: Transaction execution schedule E

```

1  $v_{root}$  root of  $G$  at which  $o$  lies initially;
2  $L$  set of leaves of  $G$ ;
3 ; cost of moving  $o$  and a control message over a unit weight edge of  $G$ , respectively;
4  $txnum(v_i)$  transactions contained in the sub-tree of  $v_i$ ;
5  $txload(v_i)$  sum of distances from  $v_i$  to transactions in sub-tree of  $v_i$ ; // initialize 0
6  $S$  set of supernodes ; // initialize  $\emptyset$ 
7  $C(v_i)$  child nodes of  $v_i \in S$  towards which object  $o$  does not move further;
8  $D(v_i)$  set of candidates for the last supernode that are descendants of  $v_i \in S$ ;
/* Compute  $txload(v_i)$  iteratively */
9 for each transaction  $T_i \in T$  do
10  $v_i$  current node of  $G$  at which  $T_i$  is positioned;
11  $v_{cur} = v_i$ ;  $txnum(v_i)++$ ;
12 while  $v_{cur} \neq v_{root}$  do
13  $v_{next}$  parent node of  $v_{cur}$  in  $G$ ;  $txnum(v_{next})++$ ;
14  $txload(v_{next}) = txload(v_{next}) + dist(v_i, v_{next})$ ;
15  $v_{cur} = v_{next}$ ;
/* Compute set of supernodes */
16 for each node  $v \in L$  do
17  $v_{cur} = v$ ;  $v_{prev} = null$ ;  $v_{cand} = null$ ;
18  $v_{next}$  parent node of  $v_{cur}$  in  $G$ ;
19 while  $(v_{cur} \neq v_{root}) \wedge (txload(v_{cur}) + 2 \cdot dist(v_{cur}, v_{next}) >$ 
20  $txload(v_{cur}) + txnum(v_{cur}) \cdot dist(v_{cur}, v_{next}))$  do
21 if  $txload(v_{cur}) > txnum(v_{cur})$  then
22  $v_{cand} = v_{cur}$ ; add  $v_{prev}$  to  $C(v_{cur})$ ;
23  $v_{prev} = v_{cur}$ ;  $v_{cur} = v_{next}$ ;
24  $v_{next}$  parent node of  $v_{cur}$  in  $G$ ;
25 add  $v_{cur}$  to  $S$ ; add  $v_{prev}$  to  $C(v_{cur})$ ; add  $v_{cand}$  to  $D(v_{cur})$ 
/* Find last supernode to visit optimizing total cost */
26  $v_{last} = \text{FindLastSuperNode}(S; D; C)$ ;
27 reorder  $G$  so that each node on path from  $v_{root}$  to  $v_{last}$  becomes the last child of parent;
28 perform a pre-order tree traversal on  $G$  starting at  $v_{root}$ ;
29 if a visited node  $v_{cur}$  is a supernode then
30 move the object  $o$  to  $v_{cur}$ ;
31 add transactions in the sub-tree of each node  $v \in C(v_{cur})$  to  $E$ ;
32 add transaction at  $v_{cur}$  to  $E$ , unless added already;

```

1 Inequality 1 accounts for the object traversing each edge twice, which is not re-
 2 quired for v_{last} . The object can move further down until the following holds:

$$dist(v_{cur}, v_{next}) > txnum(v_{cur}) \cdot dist(v_{cur}, v_{next}) \quad (2)$$

3 We find the last supernode v_{last} and the one-way path as follows. Let S be the initial
 4 set of supernodes computed considering that the object moves twice on each edge

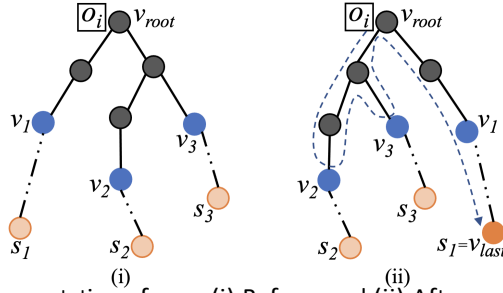


Figure 3: Illustration of computation of v_{last} . (i) Before; and (ii) After computing v_{last} and rearranging G .

1 up to the supernode. In a one-way path, the object may move further down towards
 2 the leaf node satisfying the condition in Inequality (2). For each node $v \in S$, if the
 3 sub-tree of v contains multiple branches, there could be a number of possible paths
 4 for the object to move. There will always be a unique one-way path that minimizes
 5 the total cost. In each sub-tree of $v \in S$, we find the set of nodes $D(v)$ that are
 6 candidates for v_{last} using the condition in Inequality (2) (Lines 20, 21 and 24 of
 7 Algorithm 1). Then the difference between the cost of selecting v as a supernode
 8 and $v_j \in D(v)$ as a supernode is computed (Lines 2–9 of Algorithm 2). Among
 9 these differences for every $v \in S$, the one with the highest difference is chosen as
 10 the last supernode v_{last} (Line 10 of Algorithm 2). Let $v_{ref} \in S$ and $v_k \in D(v_{ref})$ be
 11 the set of two nodes that provided the highest difference. Then v_k becomes v_{last} and
 12 is added to S . The path from v_{root} to v_{last} becomes the one-way-path and is visited
 13 at last following the pre-order tree traversal. Moreover, if a node between v_{ref} and
 14 v_{last} (including v_{ref}) in the one-way-path contains transactions in its sub-tree other
 15 than the one-way-path branch, it becomes a supernode to serve control requests to
 16 the transactions in those branches and is added to S (Lines 11–18 of Algorithm 2).
 17 Figure 3 illustrates the computation of v_{last} and rearrangement of G for the pre-order
 18 tree traversal.

19 Lemma 2. If v is a descendant of v_{last} , then the total number of transactions con-
 20 tained in the sub-tree of v is always less than .

21 Proof. This follows from the specification of the nodes that make P and how the
 22 last supernode v_{last} is determined. We assume that the object traverses only once on
 23 each link in the path from the root to v_{last} that is not needed for backtracking.

24 Lemma 3. For any transaction, the corresponding supernode for accessing the ob-
 25 ject always lies at or above its position along the path towards the root of G .

26 Proof. To compute a supernode in G , we start from a leaf node and proceed towards
 27 the root node following the shortest path.

Algorithm 2: FindLastSuperNode(S; D; C)

```

1 vlast ← null; vref ← null; di ← 0;
2 for each node v ∈ S do
3   ctrl(v) ← cost of control messages sent to v;
4   twowaycost ← 2 dist(vroot; v) + ctrl(v);
5   for each node vj ∈ D(v) do
6     onewaycost ← dist(vroot; vj) + ctrl(vj);
7     if twowaycost - onewaycost > di then
8       di ← twowaycost - onewaycost;
9       vlast ← vj; vref ← v;
10 add vlast to S; vcur ← vlast; set vnext to the parent of vlast;
11 while vcur ≠ vref do
12   if vnext has more children than vcur then
13     add vnext to S;
14     for each child node vk of vcur do
15       add vk to C(vnext);
16   if vnext == vref then
17     remove vcur from C(vnext);
18   vcur ← vnext; set vnext to the parent of vnext;
19 return vlast;

```

1 **Theorem 1.** Algorithm Single-Object schedules transactions with the optimal
2 communication cost.

3 **Proof.** Let S be the set of supernodes found for a tree G with respect to object o .
4 We will show that any other selection of supernodes gives strictly higher communi-
5 cation cost and hence, S provides optimal communication cost.

6 To simplify the problem, without loss of generality, we assume that each edge of
7 G has weight 1, $\alpha = 1$ and $\beta > 1$. Let P be the pruned tree containing nodes only up
8 to the supernodes starting from the root of G . Let $v_{last} \in S$ be the last supernode for
9 object o to visit. Let C be the total communication cost of Algorithm Single-Object.
10 Let $v \in S$ be a supernode in G , v_p be an ancestor of v with distance $\text{dist}(v_p; v) = 1$,
11 and v_q be a descendant of v with $\text{dist}(v; v_q) = 1$. Based on the positions of v and v_q ,
12 it can have one of the following three cases:

13 **Case (a):** $v = v_{last}$. Then, by Lemma 2, we have that

$$\text{txnum}(v_p) - \text{txnum}(v) > \text{txnum}(v_q) \quad (3)$$

14 **Case (b):** $v, v_{last}, v_q < P$, and the path from v to v_q contains no other supernode, in
15 that v is the bottommost supernode in the current branch. Then, by Lemma 1, we
16 have

$$\text{txnum}(v_p) - \text{txnum}(v) \geq 2 > \text{txnum}(v_q) \quad (4)$$

17 **Case (c):** Either $v_q \in P$ or $v_q < P$ and the path from v to v_q contains at least one

1 other supernode. Let $z \geq 1$ be the transactions that send control messages to v to access o .

3 We have following four subcases with respect to each supernode $v \in S$:

4 (i) Choosing an ancestor of v as a supernode instead of v increases communication:

5 Let S_p be the set of nodes contained between v and v_p (excluding both). Suppose
6 v_p be selected as a supernode instead of v . Then in Case (a) and Case (b), o moves
7 only up to v_p , and in addition to the transactions issued to the sub-tree of v , all
8 the transactions between v and v_p send control messages to v_p . But, in Case (c),
9 since the sub-tree of v (excluding v) still contains another supernode $v_k \in S$, o still
10 moves to v_k passing through v . When v was the supernode, $z \geq 1$ transactions could
11 access o at v . Now, since v_p is selected as the supernode instead of v , all those $z \geq 1$
12 transactions send control messages to v_p to access o . So, the total communication
13 cost C_{v_p} of selecting v_p as a supernode compared to that of selecting v in each case
14 becomes:

$$\begin{aligned}
 C_{v_p} &= \begin{cases} C + \sum_{v_k \in S_p} (\text{txnum}(v_k) - \text{txnum}(v)) \text{dist}(v_p; v) & \text{Case (a)} \\ C + 2 \sum_{v_k \in S_p} (\text{txnum}(v_k) - \text{txnum}(v)) \text{dist}(v_p; v) & \text{Case (b)} \\ C + z \text{dist}(v_p; v) & \text{Case (c)} \end{cases}
 \end{aligned}$$

15 In Case (a), from Inequality (3), since $\text{txnum}(v) \geq 1$, $C_{v_p} > C$. In Case (b), from
16 Inequality (4), since $\text{txnum}(v) \geq 2$, $C_{v_p} > C$. Also, in case (c), $C_{v_p} > C$.

17 (ii) Choosing a descendant of v as a supernode instead of v increases commu-
18 cation:

19 Now, we analyze the communication cost of selecting a descendant node v_q as
20 a supernode instead of $v \in S$. Let S_q be the set of nodes contained between v and v_q
21 (excluding both). As v_q is a new supernode, object moves up to it. So, in Case (a)
22 and Case (b), to get the change in total communication cost compared to C , we have
23 to add object movement cost of o from v to v_q and subtract the control message cost
24 for the transactions between v and v_q . Moreover, the transactions in the sub-tree of
25 v_q will also send control messages only up to v_q . Thus, the total communication
26 cost C_{v_q} of selecting node v_q as a supernode compared to C in Case (a) and Case (b)
27 becomes:

$$\begin{aligned}
 C_{v_q} &= \begin{cases} C + \sum_{v_k \in S_q} (\text{txnum}(v_k) - \text{txnum}(v_q)) \text{dist}(v; v_q) & \text{Case (a)} \\ C + 2 \sum_{v_k \in S_q} (\text{txnum}(v_k) - \text{txnum}(v_q)) \text{dist}(v; v_q) & \text{Case (b)} \end{cases}
 \end{aligned}$$

1 Let $\text{dist}(v; v_q) = k$ where $k \geq 1$. In Case (a), from Inequality (3), $\text{txnum}(v_q) < 2j$. Let
 2 $\text{txnum}(v_q) = 2j - 1$, $1 \leq j < \dots$. Following Lemma 2, the nodes between v and v_q (i.e., S_q)
 3 contain at most j number of transactions. The control message cost sent to v due to
 4 these transactions is: $\sum_{v_k \in S_q} (\text{txnum}(v_k) - \text{txnum}(v_q)) < j \cdot k$. Thus,

$$C_{v_q} > C + k \cdot (2j - 1) - j \cdot k > C :$$

5 In Case (b), $\text{txnum}(v_q) < 2$ by the Inequality (4). Let $\text{txnum}(v_q) = 2 - l$, for $1 \leq l < 2$. By Lemma 1, there are at most l transactions between v and v_q , and
 6 control message cost sent to v due to them is: $\sum_{v_k \in S_q} (\text{txnum}(v_k) - \text{txnum}(v_q)) < l \cdot k$.
 7 Thus

$$C_{v_q} > C + 2k - (2 - l)k - l \cdot k > C :$$

9 Now, we analyze Case (c). Based on the position of v_q , it has two sub-cases:
 10 Case (c.1): $v_q \geq P$. There is no extra movement of o and the $z - 1$ number
 11 of transactions that previously depend on v now send control messages to v_q to
 12 access o . So, the total communication cost C_{v_q} compared to C becomes: $C_{v_q} = C +$
 13 $z \cdot \text{dist}(v; v_q) > C$.

14 Case (c.2): $v_q < P$ but the path from v to v_q contains at least one other supernode
 15 in S . The node v_q lies below the bottommost supernode of current branch. Let
 16 $v_{\text{bot}} \in S$ be the bottommost supernode in the path between v and v_q . When v_q is
 17 selected as a supernode, there will be extra movement of object o from v_{bot} up to v_q .
 18 If $v_{\text{bot}} = v_{\text{last}}$, and o moves up to v_q . Otherwise, object o also needs to return back
 19 at v_{bot} . Let M represents the cost due to the movement of object o between v_{bot} and
 20 v_q , then, $M > \text{dist}(v_{\text{bot}}; v_q)$. Thus, the total communication cost C_{v_q} compared to
 21 C in this case becomes: $C_{v_q} = C + z \cdot \text{dist}(v; v_q) + M > C$.

22 (iii) Merging multiple supernodes at some ancestor node increases communication:
 23 Consider two supernodes $v_r, v_s \in S$ have a common ancestor v_y . Instead of v_r
 24 and v_s , let v_y be chosen as a supernode. Since v_y is ancestor of both v_r and v_s ,
 25 following argument (i), total communication cost C_{v_y} of selecting v_y as a supernode
 26 instead of v_r and v_s is more compared to C .

27 (iv) Splitting any supernode into multiple supernodes increases communication:
 28 Consider a supernode $v_j \in S$. Let v_x, v_z be two descendant nodes of v_j at two
 29 different sub-branches. Let v_x and v_z are chosen as two different supernodes instead
 30 of v_j . Since both v_x and v_z are descendants of v_j , following argument (ii), total cost
 31 C_{v_x, v_z} of selecting v_x, v_z as supernodes instead of v_j is more compared to C .

32 The set of supernodes S computed in algorithm Single-Object is unique. If any
 33 new node is added to S or any node in S is removed or replaced by another node,
 34 the total communication cost increases. It implies that scheduling by algorithm
 35 Single-Object provides the optimal communication cost.

1 Theorem 2. Algorithm Single-Object provides 2-approximation in communication
2 cost without optimization.

3 **Proof.** The algorithm in its optimized version has object o traverse each edge along
4 the path from v_{root} to v_{last} only once. The algorithm without optimization does not
5 identify and leverage the last supernode v_{last} for object o to visit and reorder G
6 accordingly. Without the optimization, object o still visits each edge along the path
7 from v_{root} to the supernode(s) at most twice.

8 3.2. Transaction-Migration Technique

9 Next we consider the transaction-migration technique. Let c be the cost of mov-
10 ing a transaction over a unit weight edge of G . Consider algorithm Single-Object
11 modified such that transactions are moved instead of sending control messages to
12 the supernodes and the cost of moving transaction replaces the cost of sending con-
13 trol messages, in that we use the parameter c instead of α . After these modifications
14 in algorithm Single-Object and its analysis, we obtain optimality similarly as stated
15 in Theorem 1.

16 4. Multiple Objects

17 We provide two scheduling algorithms for multiple shared objects, which ex-
18 tend the single object algorithm above. For the control-message technique, we
19 present the algorithm denoted as MultipleObjects-CtrlMsg, which provides an
20 $O(1)$ -approximation. For the transaction-migration technique, our algorithm is de-
21 noted as MultipleObjects-TxMigr, which provides $O(k)$ -approximation, where k is
22 the maximum number of shared objects accessed by a transaction.

23 The idea in the version MultipleObjects-CtrlMsg is as follows: first, we com-
24 pute separate sets of supernodes with respect to each object in O and the transactions
25 accessing that object. Each transaction will have a list of supernodes at which re-
26 quired objects can be accessed with the minimum cost. We select a random node of
27 G as the virtual root (v_{root}) of G . After that, transactions are added to the execution
28 schedule E by following the iterative pre-order tree traversal algorithm in G rooted
29 at v_{root} . Objects will then move to the respective supernodes following the iterative
30 pre-order tree traversal in G . Particularly, a transaction T_i is scheduled to execute at
31 time step $t(T_i)$ in such a way that the objects required by T_i are also reached at the
32 respective supernodes where T_i accesses them at time step $t(T_i)$. When an object
33 $o_i \in O$ reaches a supernode v_s , it stays there until all the transactions that access
34 object o at v_s finish their executions. Object o then moves to the next supernode in
35 order where other transactions are waiting for it. The algorithm ends when all the
36 transactions are added to the schedule E .

1 In algorithm MultipleObjects-TxMigr), if all the objects are positioned at the 2
 same node initially, we assume that node as the virtual root (v_{root}) of G . Otherwise, 3 if
 objects are at arbitrary nodes of G initially, we find the virtual root of G with 4
 respect to the initial positions of objects in O and the positions of transactions ac-5
 cessing those objects. We then move the objects to v_{root} . We compute separate sets 6 of
 supernodes with respect to each object $o_i \in O$ positioned at v_{root} of G . Then, 7 if a
 transaction requires multiple objects, it may have a set of different supernodes 8 to
 access each of them. Among them, we select a single supernode (called com-9 mon
 supernode) at which all the required objects for that transaction are gathered
 10 together and the transaction is also migrated during the execution. We find a com-11
 mon supernode for each transaction $T \in T$. At the common supernode, objects 12
 required by T (i.e., $objs(T)$) wait for each other like threads wait for each other at a 13
 barrier. Then, we find a pruned tree P containing the nodes only up to the common 14
 supernodes starting from the virtual root in G . We perform an iterative pre-order 15
 tree traversal in P to move objects from one common supernode to the next. At each 16
 common supernode $v_s \in P$, the transactions which have chosen v_s as the common 17
 supernode are scheduled for execution following the iterative pre-order tree traver-18
 sal in the sub-tree of v_s . As previously, at each common supernode v_s , object o 19
 stays until all the transactions requiring object o at v_s finish their executions. Ob-20
 ject o then moves to the next common supernode. The algorithm terminates once 21 all the
 transactions get scheduled.

22 We consider a set of shared objects $O = \{o_1, o_2, \dots, o_g\}$ initially positioned at
 23 arbitrary nodes of G . We assume that each object has size s . Each transaction in T 24
 accesses a subset of objects in O . Let $objs(T_i) \subseteq O$ be the set of objects accessed 25
 by transaction T_i . We assume that each object has a single copy and $home(o_i) \in V$ 26
 represents the home node at which object o_i is originally positioned. The owner-27
 ship of an object is also transferred with the movement of that object. Similarly, 28
 $home(T_i) \in V$ represents the node at which transaction T_i is positioned.

29 The idea in the algorithms is to provide synchronized accesses to the objects 30
 with minimum cost while executing the transactions in order. We achieve this ex-31
 tending the techniques used in algorithm Single-Object. In particular, we compute 32
 supernodes w.r.t. each object and the transactions requiring those objects. We then 33
 perform iterative pre-order tree traversal to move each object to the respective su-34
 pernodes and execute transactions in order.

35 For brevity, let T_i be a transaction that requires objects in $objs(T_i) = \{o_x, \dots, o_z\}$.
 36 Let $sv_i(o_x), \dots, sv_i(o_z)$ be the respective supernodes (computed using algorithm 37
 Single-Object w.r.t. each object) at which T_i can access o_x, \dots, o_z , respectively. 38
 Then, one way of providing synchronised access to the required objects by T_i is to 39
 bring each object in $objs(T_i)$ at the respective supernode (i.e., $sv_i(o_x), \dots, sv_i(o_z)$) 40
 at the same time so that T_i can access them by sending control messages. This

1 approach is used in the control-message technique. The other way is to gather all 2
the objects in $objs(T_i)$ at a single node $sv(T_i)$ (i.e., common supernode for T_i) and 3
access them at that node by migrating T_i . This approach is used in the transaction-4
migration technique.

5 We now describe how transactions are executed in order and the objects are 6
moved from one supernode to the next minimizing the communication cost. As in 7
algorithm Single-Object, this can be achieved using iterative pre-order tree traversal 8
algorithm in G , provided that there is a single reference point, i.e., root node. We 9
find a virtual root (v_{root}) of tree G as a single reference point.

10 In the control-message technique, any node of G can be selected as the virtual 11
root (v_{root}). In the transaction-migration technique, if all the objects are initially 12
positioned at the same node, that node is selected as the virtual root of G . If objects 13
are positioned at different nodes initially, we compute the virtual root with respect 14
to the initial positions (home nodes) of transactions and the objects they access. 15
The virtual root of tree G is the node in G from which the sum of distances to home 16
nodes of all the transactions and the objects they access is the minimum, that is,

$$v_{root}^c = v_i : W(v_i) = \min_{v \in V} W(v); \quad (5)$$

17 where,

$$W(v) = \sum_{j=1}^n \text{dist}(v; \text{home}(T_j)) + \sum_{o \in objs(T_j)} \text{dist}(v; \text{home}(o));$$

18 To compute the virtual root of G , we take into account both the initial positions 19
of transactions and the initial positions of objects. This is needed since we want to 20
minimize the distances from initial positions of objects to the virtual root as well as 21
the distances from virtual root to the transactions.

22 4.1. Multiple Objects with Control Messages

23 The algorithm for multiple objects in the control-message technique is named 24
MultipleObjects-CtrlMsg, and the pseudocode is given in Algorithm 3. The algo-25
rithm runs in two phases.

26 Phase 1: We compute sets of supernodes $S(o_i)$ with respect to each object $o_i \in O$ 27
individually following algorithm Single-Object without optimization. For each o_i , 28
 $\text{home}(o_i)$ is assumed as the root of G during the computation of respective supern-29
odes $S(o_i)$. If a transaction T_i requires an object o_j , T_i accesses o_j at supernode 30
 $sv(T_i(o_j)) \in S(o_j)$.

31 Phase 2: We find transaction execution schedule E and paths of movement for each
32 object $o_i \in O$ along their respective supernodes. For this, let a random node in G be

Algorithm 3: MultipleObjects-CtrlMsg

Input : Tree graph G of n nodes containing T transactions and O objects positioned at arbitrary nodes

Output: Transaction execution schedule E

- 1 v_{root}^0 virtual root of G ;
- /* Phase 1: Find set of supernodes for each object in O following algorithm Single-Object without optimization */
- 2 for each object $o_i \in O$ do
- 3 $S(o_i)$ set of supernodes with respect to o computed using algorithm Single-Object without optimization and assuming $home(o_i)$ as the root of G ;
- 4 $S = S \cup S(o_i)$;
- /* Phase 2: Pre-order tree traversal on G */
- 5 perform pre-order tree traversal on G rooted at v_{root}^0
- 6 v_{cur} current visited node on G ;
- 7 if v_{cur} is a supernode with respect to o then
- 8 move object o to v_{cur} ;
- 9 add transaction at v_{cur} to E , if not added previously;

1 selected as the virtual root v_{root}^0 of G . We perform an iterative pre-order tree traversal
 2 in G starting from v_{root}^0 . During the traversal, if there is a transaction T_j at current
 3 node v_{cur} , T_j is added to the schedule E . Each object o_k required by T_j (in notation, 4
 2 $objs(T_j)$) is scheduled to move to the respective supernode $sv_j(o_k)$. When the 5
 traversal of G completes, all the transactions get scheduled and the execution ends.

6 Lemma 4. An object o may traverse an edge along the path from $home(o_i)$ to v_{root}^0
 7 at most three times.

8 Proof. Let G be a tree and v_{root}^0 be the virtual root of G . Let the shared objects
 9 $O = \{o_1, \dots, o_n\}$ be positioned at arbitrary nodes of G . Then, with respect to each 10
 object $o_i \in O$, a set of supernodes $S(o_i)$ is computed using algorithm Single-Object. 11
 Let $S = \bigcup_{i=1}^n S(o_i)$ be the union of all supernodes and P be the pruned tree contain-12
 ing all the nodes in S starting from v_{root}^0 . Transactions are scheduled following an 13
 iterative pre-order tree traversal in G rooted at v_{root}^0 and objects are moved from one 14
 supernode to the next in S following the iterative pre-order tree traversal in P rooted 15
 at v_{root}^0 . We start from v_{root}^0 . For the transaction T assigned at v_{root}^0 , if v_{root}^0 is a supern-16
 ode with respect to some object $o_i \in objs(T)$, then object o moves from $home(o_i)$ 17
 to v_{root}^0 . Otherwise, $o_i \in objs(T)$ moves to the respective supernode $sv(T(o_i))$ below 18
 v_{root}^0 . By Lemma 3, the supernode for transaction T at v_{root}^0 for accessing object o 19
 is always at v_{root}^0 or below v_{root}^0 towards $home(o_i)$. After this, thanks to the proper-20
 ties of the pre-order tree traversal that each edge of P is visited no more than twice 21
 and hence each object $o_i \in O$ also traverses each edge along the path from v_{root}^0 22
 to $sv(T(o_i))$ in P at most two times. $home(o_i)$ itself can be one of the supernodes 23
 $sv(T(o_i))$ for T with respect to object o in P . Thus, in total, object o can traverse 24
 at most 3 times along the path from $home(o_i)$ to v_{root}^0 .

0

1 Theorem 3. Algorithm MultipleObjects-CtrlMsg provides a 3-approximation of
 2 communication cost.

3 Proof. Let $S(o_i)$ be the set of supernodes computed with respect to object o_i 2 O
 4 following algorithm Single-Object without optimization. Let P_i be the pruned tree
 5 containing nodes only up to the supernodes $S(o_i)$ starting from $\text{home}(o_i)$ in G . Let
 6 C_{obj} denotes the cost of moving object o_i at each edge inside P_i only once and C_{ctrl}
 7 denotes the communication cost incurred due to the control messages sent from
 8 transactions beyond P_i in G . By the analysis of algorithm Single-Object, o_i visits
 9 each edge of P_i at most twice during the execution. Theorem 1 shows that the set of
 10 supernodes computed in algorithm Single-Object provides the minimum commu-
 11 nication cost and Theorem 2 shows that algorithm Single-Object without optimiza-
 12 tion provides 2-approximation. Thus, if $C_{\text{OPT}}(o_i)$ be the optimal communication
 13 cost for accessing o_i by a set of transactions T , then,

$$C_{\text{obj}} + C_{\text{ctrl}} \leq C_{\text{OPT}}(o_i) \leq 2(C_{\text{obj}} + C_{\text{ctrl}}) \quad (6)$$

14 and $C_{\text{OPT}} = \sum_{o_i \in O} C_{\text{OPT}}(o_i)$.

15 The algorithm in MultipleObjects-CtrlMsg uses the same set of supernodes 16
 $S(o_i)$ computed in algorithm Single-Object without optimization and object o_i does 17
 not move beyond the pruned tree P_i . So, C_{ctrl} for MultipleObjects-CtrlMsg re-18
 mains the same. From Lemma 4, object o_i may traverse an edge inside P_i at most 19
 3 times. Thus, if $C_{\text{ALG}}(o_i)$ represents the total communication cost for accessing o_i 20
 by a set of transactions T , then,

$$C_{\text{ALG}}(o_i) \leq 3C_{\text{obj}} + C_{\text{ctrl}} \quad (7)$$

21 Inequalities (6) and (7) imply

$$C_{\text{ALG}}(o_i) \leq 3 C_{\text{OPT}}(o_i) \quad (8)$$

22 This gives the estimate

$$C_{\text{ALG}} = \sum_{o_i \in O} C_{\text{ALG}}(o_i) \leq \sum_{o_i \in O} (3 C_{\text{OPT}}(o_i)) \leq 3 C_{\text{OPT}} ;$$

23 where C_{ALG} represents the total communication cost in MultipleObjects-CtrlMsg 24
 for executing all the transactions accessing multiple objects and C_{OPT} represents 25
 that of any optimal algorithm.

26 4.2. Multiple Objects with Migration of Transactions

27 The algorithm for multiple objects implemented in the transaction-migration 28
 model is named MultipleObjects-TxMigr. First, we discuss the algorithm assum-29
 ing all the objects are initially positioned at the same node (i.e., $\text{home}(o_i)$) for each

Algorithm 4: MultipleObjects-TxMigr

Input: Tree graph G of n nodes containing T transactions and O objects
Output: Transaction execution schedule E

```
1  $v_{root}^0$  virtual root of  $G$ ;  
/* Phase 1: Compute supernodes */  
2 for each object  $o_i \in O$  do  
3   if  $home(o_i) = v_{root}^0$  then move  $o_i$  to  $v_{root}^0$ ;  $S(o_i)$  set of supernodes of  $o_i$ ;  
/* Phase 2: Find common supernode */  
4 for each transaction  $T \in T$  do  
5   for each object  $o_i \in objs(T)$  do  
6      $sv(T(o_i))$  node in  $S(o_i)$  where  $T$  accesses  $o_i$ ;  
7    $v$  the node in  $sv(T(o_j))$ ;  $o_j \in objs(T)$ ; which has the minimum  $dist(sv(T(o_j)), v_{root}^0)$ ;  
8    $sv(T) = v$ ;  $numtxs(sv(T))++$ ;  
9   for each object  $o_i \in objs(T)$  do  
10     $sv(T(o_i)) = v$ ;  $objs(sv(T)) = objs(sv(T)) \cup \{o_i\}$ ;  
11     $txs(sv(T)(o_i)) = txs(sv(T)(o_i)) \cup \{T\}$ ;  
/* Phase 3: Finalize supernodes */  
12  $L$  set of leaf nodes in  $G$ ;  
13 repeat  
14   for each node  $v \in L$  do  
15     while  $numtxs == 0$  do  $v = parent$  of  $v$  in  $G$ ;  
16     if  $numtxs(v) < 2$  then  
17       for each object  $o_i \in objs(v)$  do  
18         while  $v \neq v_{root}^0 \wedge |txs(v(o_i))| < 2$  do  
19            $p = parent$  node of  $v$  in  $G$ ;  
20            $txs(v(o_i)) = txs(v(o_i)) \cup txs(p(o_i))$ ;  $v = p$ ;  
21           for each transaction  $T$  in  $txs(v(o_i))$  do  $sv(T) = v$ ;  
22   FinalSV = fg;  
23   for each transaction  $T \in T$  do FinalSV = FinalSV  $\cup sv(T)$ ;  
24    $P$  pruned tree containing nodes up to FinalSV starting from  $v_{root}^0$  in  $G$ ;  
25    $L$  leaf nodes of  $P$ ;  
26   for each node  $v \in L$  do if  $|numtxs(v)| \geq 2$  then exit repeat; /*  
Phase 4: Find execution schedule */  
27 perform iterative pre-order tree traversal on  $P$ ;  
28   if current node  $v_{cur} \in FinalSV$  then  
29     move each object  $o \in objs(v_{cur})$  to  $v_{cur}$ ;  
30     add each  $T \in T$ , where  $sv(T) = v_{cur}$ , to  $E$ ;
```

¹ $o_i \in O$ is the same) which is the virtual root v_{root}^0 of G . Later, we relax the algorithm ² where objects can be positioned initially at arbitrary nodes in G . A pseudocode is ³ given as Algorithm 4.

⁴ The algorithm works in four phases. In Phase 1, we compute sets of supernodes ⁵ with respect to individual object $o_i \in O$. In Phase 2, we find a common supernode ⁶ for each transaction $T \in T$ where all the required objects for T can be gathered ⁷ together. In Phase 3, we finalize the set of common supernodes. Finally, in Phase 4, ⁸ we perform iterative pre-order tree traversal on G to create transaction execution

1 schedule and object movement paths along the respective common supernodes. We
 2 describe each phase in detail next.

3 Phase 1: In this phase, we compute supernodes with respect to each object $o_i \in O$
 4 using algorithm Single-Object without optimization where control message cost
 5 over an edge is replaced with the transaction migration cost. Let $S(o_i)$ be the set
 6 of supernodes with respect to object $o_i \in O$ and $sv(T(o_i)) \in S(o_i)$ represents the
 7 supernode for transaction T at which T accesses o_i . After this, each transaction
 8 $T_j \in T$ has a set of respective supernodes $sv(T_j(o_i))$ to access each required object
 9 $o_i \in \text{objs}(T_j)$. Since all the objects in $\text{objs}(T_j)$ need to gather at a single node, a
 10 common supernode $sv(T_j)$ for transaction T_j is selected out of all $sv(T_j(o_i))$ in the
 11 next phase.

12 **Observation 1.** For a transaction T , all the supernodes $sv(T(o_i))$, for $o_i \in \text{objs}(T)$;
 13 lie on the same path towards the root.

14 **Proof.** All the objects are initially positioned at the same node v_{root}^c , so it suces
 15 to refer to Lemma 3.

16 Phase 2: In this phase, we find a common supernode of objects $sv(T)$ for each
 17 transaction $T \in T$. The objective of selecting a common supernode for a transaction
 18 is to allow all the required objects for that transaction to gather together at the
 19 common supernode. After that, the transaction is also migrated at the common
 20 supernode and all the required objects are accessed locally. For a transaction T ,
 21 if all the supernodes $sv(T(o_i))$, $o_i \in \text{objs}(T)$; computed in Phase 1 are the same, it
 22 automatically becomes the common supernode for T . If they are dierent, then we
 23 select the one among $sv(T(o_i))$; $o_i \in \text{objs}(T)$; which is the closest from v_{root}^c .

24 **Theorem 4.** For a transaction $T \in T$, selecting the topmost node among $sv(T(o_i))$,
 25 $o_i \in \text{objs}(T)$, as a common supernode provides the minimum communication cost
 26 during the execution of T .

27 **Proof.** Consider a tree graph G rooted at v_{root}^c . Let T be a transaction requiring
 28 a set of objects $fo_1; o_2; \dots; o_k$. Let $sv_1; sv_2; \dots; sv_k$ be the supernodes for T with
 29 respect to objects $o_1; o_2; \dots; o_k$, respectively, computed in Phase 2. Among them,
 30 let sv_1 be the topmost, sv_k be the bottommost and the remaining nodes lie between
 31 sv_1 and sv_k . By Observation 1, all the supernodes $sv_1; sv_2; \dots; sv_k$ lie on the same
 32 path towards v_{root}^c . First, let sv_1 be selected as the common supernode for T . If
 33 C_{obj} be the cost of moving objects $o_1; o_2; \dots; o_k$ from v_{root}^c to sv_1 , then the cost of
 34 executing T at sv_1 becomes $C_{\text{obj}} + \text{dist}(\text{home}(T); sv_1)$. Now, instead of sv_1 , let
 35 sv_2 be selected as the common supernode for T . Then, all the objects $o_1; o_2; \dots; o_k$

1 move up to node sv_2 and T also moves to sv_2 to execute. Thus, total cost for exe-2
 cuting T becomes $C_{obj} + 2 |jobs(T)| \text{dist}(sv_1; sv_2) + \text{dist}(\text{home}(T); sv_2)$. On 3 the one
 hand, since o_1 moves below sv_1 , from the argument (ii) of Theorem 1, the 4
 communication cost increases. On the other hand, object movement cost increases 5 by
 $2|jobs(T)| \text{dist}(sv_1; sv_2)$ which is more than the decrease in transaction move-6
 ment (i.e., $\text{dist}(sv_1; sv_2)$). Thus, in any case, communication cost for executing T 7
 at sv_2 by selecting it as a common supernode is more compared to that by selecting 8
 sv_1 as a common supernode. Arguing similarly, selecting any descendant of sv_1 as 9
 a common supernode for executing T increases the communication cost than that
 10 with sv_1 . Let us also see what happens if we select an ancestor sv_1 as a common 11
 supernode. Note here that, during the merging of supernodes up to sv_1 , there is no 12
 extra cost due to the movement of transaction since sv_1 already accounts for the 13
 transaction movement cost up to it. If some ancestor of sv_1 is selected as a common 14
 super for executing T instead of sv_1 , T needs to move further upward up to that 15
 ancestor node, and from the argument (i) of Theorem 1, the communication cost 16
 increases. Hence, the theorem follows.

17 Optimization. By the end of Phase 2, each transaction T 2 T has a corresponding 18
 common supernode $sv(T)$ at which T can migrate for accessing all the required 19
 objects and execute. Since some of the initially computed supernodes with respect 20
 to individual objects are now merged together at some ancestor node, some common 21
 supernodes may contain fewer transactions and instead of moving objects to those 22
 common supernodes, moving the transactions upward may further decrease the total 23
 communication cost.

24 For example, let u and v be the two inner nodes of a tree rooted at v_{root} such 25
 that $u = \text{parent}(v)$. Let $n_u = 1$ and $n_v > 1$. Let sub-tree of node v contains 2 + 26
 1 transactions where $n_u + 1$ transactions access two objects $a; b$ and remaining 27
 transactions access only object a . Let the sub-tree of node u contains 2 + 1 number 28
 of transactions requiring both objects $a; b$, that is, extra transactions requiring 29
 the object b . From Phase 1, since the sub-tree of v contains 2 + 1 number of 30
 transactions accessing object a , it becomes a supernode with respect to a . Since the 31
 sub-tree of v contains only number of transactions requiring object b , it cannot 32
 be a supernode with respect to b . The sub-tree of node u contains 2 + 1 number 33
 of transactions requiring object b and number of transactions not included in the 34
 sub-tree of v requiring object a . Thus u becomes the supernode with respect to both 35
 a and b . In Phase 2, the 2 + 1 number of transactions requiring both objects $a; b$ 36
 select u as a common supernode and remaining number of transactions requiring 37
 only object a select v as a common supernode. When object a is moved from node u 38
 to v and again back to u , the object movement cost is 2. Instead, if we move those 39
 number of transactions from v to u and execute at u , transaction movement cost

1 will increase by reducing object movement cost by 2. That means it minimizes 2
the total cost by .

3 We provide such optimization in Phase 3 of the algorithm.

4 Phase 3: In this phase, we compute the final set of supernodes FinalSV in G where 5
respective transactions and the required objects are gathered together. From Phase 6 2,
we have a set of common supernodes $sv(T)$ for each transaction $T \in T$. For each 7
common supernode $v \in sv()$, the following information is maintained separately:

- 8 – $numtxs(v)$: total number of transactions that have selected v as a common
9 supernode.
- 10 – $objs(v)$: set of objects with respect to which the node v is a supernode.
- 11 – $txs(v(o_i))$; $o_i \in objs(v)$: set of transactions requiring object o_i that have se-
12 lected v as the common supernode.

13 Let P be the pruned tree containing the nodes of G only up to the common supern-14
odes moving down from v_{root}^0 . Starting from every leaf node of P towards v_{root}^0 , we 15
check at each node how many transactions have selected it as a common supernode. 16
Particularly, if $v \in P$ is a leaf node in P and is selected as a common supernode with 17
respect to the set of objects $objs(v)$, then, we check if $numtxs(v) \geq |objs(v)|$. If 18 the
condition is satisfied, v belongs to FinalSV with respect to all objects in $objs(v)$. 19
Otherwise, for each object $o_i \in objs(v)$, we check how many transactions requiring 20
the object o_i have selected v as the common supernode in Phase 2. Let $txs(v(o_i))$ be 21 the
set of transactions requiring object o_i that have selected v as a common supern-22 ode.
If $|txs(v(o_i))| \geq 2$, v belongs to FinalSV. But if $|txs(v(o_i))| < 2$, we 23 visit its parent
node $parent(v)$, find the set of transactions $txs(parent(v)(o_i))$ requir-24 ing object o_i that
have selected $parent(v)$ as the common supernode. At the parent 25 node $parent(v)$, we
again check if $(|txs(v(o_i))| + |txs(parent(v)(o_i))|) \geq 2$. If the 26 condition is met, $parent(v)$
belongs to FinalSV and all the transactions in $txs(v(o_i))$ 27 that previously selected node
 v as the common supernode now select $parent(v)$ as 28 the common supernode.
Otherwise, if the condition is not met, we repeat the same 29 procedure by selecting the
parent of $parent(v)$ and so on until the inequality

$$(|txs(v(o_i))| + |txs(parent(v)(o_i))| + \dots) \geq 2$$

30 is satisfied or reach at v_{root}^c . We apply this approach recursively until at each leaf 31
node $v \in P$, $numtxs(v) \geq 2$ where P is the pruned tree containing nodes only 32 up to
final set of common supernodes FinalSV starting from v_{root}^c .

33 Phase 4: In this phase, we find the transaction execution schedule E and the paths
34 of movement for each object $o_i \in O$ along their respective supernodes. We find

1 the pruned tree P containing the nodes up to the common supernodes in FinalSV
2 starting from v_{root}^c . Then we perform iterative pre-order traversal on P starting from
3 v_{root}^c . At each current visited node v , if $v \in \text{FinalSV}$, then all the transactions which
4 have selected v as their common supernode (i.e., $sv(T) = v$) are added to the execu-
5 tion schedule E . Additionally, the objects in O for which v is a common supernode
6 (i.e., $\text{objs}(v)$) are scheduled to move at v . An object $o_k \in \text{objs}(v)$ remains at v until
7 all the transactions that require o_k finish their executions. After all the transactions
8 that require object $o_k \in \text{objs}(v)$ finish their executions, o_k can move to the next com-
9 mon supernode in the order where other transactions are waiting for it. When the
10 traversal of P completes, all the transactions get scheduled and the algorithm ends.

11 **Theorem 5.** Algorithm `MultipleObjects-TxMigr` provides k -approximation in
12 communication cost, where k is the maximum number of objects accessed by a
13 transaction.

14 **Proof.** After computing the final set of common supernodes FinalSV , at the bottom-
15 most common super $v \in \text{FinalSV}$ in each branch of G , the number of transactions
16 that require object o are at least 2. These 2 number of transactions in the sub-
17 tree of v may require k number of objects in O . Thus node v can be a common
18 supernode for all those 2 transactions with respect to k objects. During the ex-
19 ecution, these k objects are moved from v_{root} to v and the cost is $k \cdot 2 \cdot \text{dist}(v_{\text{root}}; v)$.
20 Instead, if we move those 2 transactions up towards some closest common su-
21 pernode v_j that contains at least $k \cdot 2$ number of transactions, then the cost due to
22 transaction migration increases by $2 \cdot \text{dist}(v_j; v)$ reducing the object movement cost
23 by $k \cdot 2 \cdot \text{dist}(v_j; v)$. That means the total cost may increase by at most a k factor
24 from optimal.

25 **Arbitrary Initial Positions of Objects.** Here, we discuss algorithm
26 `MultipleObjects-TxMigr` with the relaxed setting where objects are positioned at
27 arbitrary nodes of G initially. In this case, before Phase 1, we compute the virtual
28 root v_{root}^0 of G using Equation 5. All the objects in O are then moved to v_{root}^0 . After
29 this, algorithm continues with Phase 1 to Phase 4 as it is. There is an extra cost
30 incurred before Phase 1 due to the movements of objects from their home nodes to
31 the virtual root. Let C_{extra} represents this extra cost due to the movements of objects
32 from their home nodes to v_{root}^0 which is:

$$C_{\text{extra}} = \sum_{o_i \in O} \text{dist}(\text{home}(o_i); v_{\text{root}}^0) \quad (9)$$

33 Let FinalSV be the finalized set of common supernodes computed in Phase 3 of
34 algorithm `MultipleObjects-TxMigr` after moving all objects in O to v_{root}^c . Let C_{mov}

1 be the total cost due to the movements of objects from v_{root}^c to their respective com-2
mon supernodes in FinalSV following the iterative pre-order tree traversal. Now, 3 let
4 $S(o_i)$ be the sets of supernodes computed with respect to each object $o_i \in O$
5 positioned at the respective home node and using algorithm Single-Object without
6 optimization. Let C_{opt_mov} denotes the total cost due to the movements of objects in
7 their respective supernodes in $S(o)$ following iterative pre-order tree traversal. By
8 Theorem 2, we have that C_{opt_mov} is asymptotically optimal with respect to the ob-
9 jects movement cost. If $C_{extra} + C_{mov} \leq k C_{opt_mov}$ then algorithm MultipleObjects-
TxMigr has performance as in Theorem 5 in the relaxed setting as well.

10 On the other hand, consider a case where objects are already on the respec-11
12 tive common supernodes during the initial configuration. Then since there are
13 total objects, from Equation 5, algorithm MultipleObjects-TxMigr provides
14 $O(D)$ -approximation in the relaxed setting where D is the diameter of tree G
and $dist(home(o); v_{root}) \leq D$.^c

15 5. Conclusion

16 In this paper, we studied transaction scheduling problem on trees to minimize 17
18 communication cost in the dual-flow model where both data and transactions are
19 mobile. When transactions access the same single shared object, we provided an
20 optimal schedule for executing the transactions. Extending it for multiple shared
21 objects, we provided algorithms that are within k factor away from the optimal,
where k is the maximum number of shared objects requested by a transaction.

22 In the future work, it will be interesting to study the dynamic online setting of 23
24 the scheduling problem in the dual-flow model. It will also be interesting to con-
25 sider other kinds of graph or extend the algorithms towards the general networks.
26 Consideration of other performance metrics such as execution time and congestion
is also a probable future direction.

27 References

28 [1] C. Busch, B. S. Chlebus, M. Herlihy, M. Popovic, P. Poudel, G. Sharma, Flexible 29
30 scheduling of transactional memory on trees, in: Stabilization, Safety, and Security of
31 Distributed Systems - 24th International Symposium, SSS 2022, Clermont-Ferrand,
France, November 15-17, 2022, Vol. 13751 of LNCS, Springer, 2022, pp. 146–163.

32 [2] M. Herlihy, J. E. B. Moss, Transactional memory: Architectural support for lock-33
34 free data structures, in: Proceedings of the 20th Annual International Symposium on
Computer Architecture, ACM, 1993, pp. 289–300.

35 [3] N. Shavit, D. Touitou, Software transactional memory, Distributed Computing 10 (2)
36 (1997) 99–116.

- ¹ [4] P. Hammarlund, A. J. Martinez, A. A. Bajwa, D. L. Hill, E. G. Hallnor, H. Jiang, ²
M. G. Dixon, M. Derr, M. Hunsaker, R. Kumar, R. B. Osborne, R. Rajwar, R. Singhal, ³
R. D'Sa, R. Chappell, S. Kaushik, S. Chennupati, S. Jourdan, S. Gunther, T. Piazza, ⁴
T. Burton, Haswell: The fourth-generation Intel core processor, *IEEE Micro* ⁵ 34 (2) ⁵
(2014) 6–20.
- ⁶ [5] R. A. Haring, M. Ohmacht, T. W. Fox, M. Gschwind, D. L. Satterfield, K. Sugavanam, ⁷
P. Coteus, P. Heidelberger, M. A. Blumrich, R. W. Wisniewski, A. Gara, G. L. Chiu, ⁸ P.
A. Boyle, N. H. Christ, C. Kim, The IBM blue gene/q compute chip, *IEEE Micro* ⁹ 32 (2)
(2012) 48–60.
- ¹⁰ [6] T. Nakaike, R. Odaira, M. Gaudet, M. M. Michael, H. Tomari, Quantitative com-¹¹
parison of hardware transactional memory for Blue Gene/Q, zEnterprise EC12, Intel ¹²
Core, and POWER8, in: *Proceedings of the 42nd Annual International Symposium ¹³ on*
Computer Architecture, ACM, 2015, pp. 144–157.
- ¹⁴ [7] H. W. Cain, M. M. Michael, B. Frey, C. May, D. Williams, H. Q. Le, Robust archi-¹⁵
tectural support for transactional memory in the power architecture, in: *Proceedings ¹⁶ of*
the 40th Annual International Symposium on Computer Architecture (ISCA 2013), ¹⁷
ACM, 2013, pp. 225–236.
- ¹⁸ [8] M. Herlihy, Y. Sun, Distributed transactional memory for metric-space networks, *Dis-¹⁹*
tributed Computing 20 (3) (2007) 195–208.
- ²⁰ [9] G. Sharma, C. Busch, Distributed transactional memory for general networks, *Dis-²¹*
tributed Computing 27 (5) (2014) 329–362.
- ²² [10] K. Siek, P. T. Wojciechowski, Atomic RMI: A distributed transactional memory
²³ framework, *International Journal of Parallel Programming* 44 (3) (2016) 598–619.
- ²⁴ [11] J. Armstrong, *Programming Erlang: Software for a Concurrent World*, Pragmatic
²⁵ Bookshelf, 2007.
- ²⁶ [12] K. Hwang, J. Dongarra, G. C. Fox, *Distributed and Cloud Computing: From Parallel
²⁷ Processing to the Internet of Things*, Morgan Kaufmann Publishers, 2011.
- ²⁸ [13] P. Ruth, J. Rhee, D. Xu, R. Kennell, S. Goasguen, Autonomic live adaptation of vir-²⁹
tual computational environments in a multi-domain infrastructure, in: *Proceedings ³⁰*
of the 3rd International Conference on Autonomic Computing (ICAC 2006), IEEE ³¹
Computer Society, 2006, pp. 5–14.
- ³² [14] E. Tilevich, Y. Smaragdakis, J-Orchestra: Automatic java application partitioning, ³³
in: *Proceedings of the 16th European Conference on Object-Oriented Programming ³⁴*
(ECOOP 2002), Vol. 2374 of LNCS, Springer, 2002, pp. 178–204.

- 1 [15] K. Arnold, R. Scheifler, J. Waldo, B. O’Sullivan, A. Wollrath, *Jini Specification*,
2 Addison-Wesley Longman Publishing, 1999.
- 3 [16] M. M. Saad, B. Ravindran, Snake: Control flow distributed software transactional
4 memory, in: *Proceedings of the 13th International Symposium on Stabilization, 5*
Safety, and Security of Distributed Systems (SSS 2011), Vol. 6976 of LNCS, Springer, 6
2011, pp. 238–252.
- 7 [17] D. E. Comer, *The Cloud Computing Book: The Future of Computing Explained*,
8 Chapman and Hall/CRC, 2021.
- 9 [18] H. Attiya, V. Gramoli, A. Milani, Directory protocols for distributed transactional
10 memory, in: *Transactional Memory. Foundations, Algorithms, Tools, and Applica-11*
tions, Vol. 8913 of LNCS, Springer, 2015, pp. 367–391.
- 12 [19] C. Busch, M. Herlihy, M. Popovic, G. Sharma, Time-communication impossibility
13 results for distributed transactional memory, *Distributed Computing* 31 (6) (2018) 14
471–487.
- 15 [20] C. Busch, M. Herlihy, M. Popovic, G. Sharma, Dynamic scheduling in distributed
16 transactional memory, in: *Proceedings of the 2020 IEEE International Parallel and 17*
Distributed Processing Symposium (IPDPS 2020), IEEE, 2020, pp. 874–883.
- 18 [21] C. Busch, M. Herlihy, M. Popovic, G. Sharma, Fast scheduling in distributed transac-
19 tional memory, *Theory of Computing Systems* 65 (2) (2021) 296–322.
- 20 [22] G. Sharma, C. Busch, A load balanced directory for distributed shared memory ob-
21 jects, *Journal of Parallel and Distributed Computing* 78 (2015) 6–24.
- 22 [23] R. Palmieri, S. Peluso, B. Ravindran, Transaction execution models in partially repli-23
cated transactional memory: The case for data-flow and control-flow, in: *Trans-24*
actional Memory. Foundations, Algorithms, Tools, and Applications, Vol. 8913 of 25
LNCS, Springer, 2015, pp. 341–366.
- 26 [24] K. Siek, P. T. Wojciechowski, Atomic RMI 2: Highly parallel pessimistic distributed
27 transactional memory, *CoRR* abs/1606.03928.
- 28 [25] M. M. Saad, B. Ravindran, HyFlow: a high performance distributed software trans-29
actional memory framework, in: *Proceedings of the 20th ACM International Sympto-30*
sium on High Performance Distributed Computing (HPDC 2011), ACM, 2011, pp. 31
265–266.
- 32 [26] R. L. Bocchino Jr., V. S. Adve, B. L. Chamberlain, Software transactional memory
33 for large scale clusters, in: *Proceedings of the 13th ACM SIGPLAN Symposium 34*
on Principles and Practice of Parallel Programming (PPOPP 2008), ACM, 2008, pp. 35
247–258.

- 1 [27] D. Hendler, A. Naiman, S. Peluso, F. Quaglia, P. Romano, A. Suissa, Exploiting local-2
ity in lease-based replicated transactional memory via task migration, in: Proceedings 3
of the 27th International Symposium on Distributed Computing (DISC 2013), Vol. 4
8205 of LNCS, Springer, 2013, pp. 121–133.
- 5 [28] B. Zhang, B. Ravindran, R. Palmieri, Distributed transactional contention manage-6
ment as the traveling salesman problem, in: Proceedings of the 21st International 7
Colloquium on Structural Information and Communication Complexity (SIROCCO 8
2014), Vol. 8576 of LNCS, Springer, 2014, pp. 54–67.
- 9 [29] P. Poudel, G. Sharma, GraphTM: An efficient framework for supporting transactional 10
memory in a distributed environment, in: Proceedings of the 21st International Con-11
ference on Distributed Computing and Networking (ICDCN 2020), ACM, 2020, pp. 12
11:1–11:10.
- 13 [30] S. Rai, G. Sharma, C. Busch, M. Herlihy, Load balanced distributed directories, in: 14
Proceedings of the 20th International Symposium on Stabilization, Safety, and Secu-15
rity of Distributed Systems (SSS 2018), LNCS, Springer, 2018, pp. 221–238.
- 16 [31] M. Couceiro, P. Romano, N. Carvalho, L. E. T. Rodrigues, D2STM: dependable dis-17
tributed software transactional memory, in: Proceedings of the 15th IEEE Pacific Rim 18
International Symposium on Dependable Computing (PRDC 2009), IEEE Computer 19
Society, 2009, pp. 307–313.
- 20 [32] S. Hirve, R. Palmieri, B. Ravindran, Hipertm: High performance, fault-tolerant trans-
21 actional memory, *Theoretical Computer Science* 688 (2017) 86–102.
- 22 [33] T. Kobus, M. Kokocinski, P. T. Wojciechowski, Hybrid replication: State-machine-23
based and deferred-update replication schemes combined, in: Proceedings of the IEEE 24
33rd International Conference on Distributed Computing Systems (ICDCS 2013), 25
IEEE Computer Society, 2013, pp. 286–296.
- 26 [34] K. Manassiev, M. Mihailescu, C. Amza, Exploiting distributed version concurrency in 27
a transactional memory cluster, in: Proceedings of the ACM SIGPLAN Symposium 28
on Principles and Practice of Parallel Programming (PPOPP 2006), ACM, 2006, pp. 29
198–208.
- 30 [35] S. Peluso, P. Romano, F. Quaglia, SCORE: A scalable one-copy serializable partial 31
replication protocol, in: Proceedings of ACM/IFIP/USENIX 13th International Mid-32
dleware Conference (Middleware 2012), Vol. 7662 of LNCS, Springer, 2012, pp. 456–33
475.
- 34 [36] S. Peluso, P. Ruivo, P. Romano, F. Quaglia, L. E. T. Rodrigues, When scalability 35
meets consistency: Genuine multiversion update-serializable partial data replication, 36
in: Proceedings of the 2012 IEEE 32nd International Conference on Distributed Com-37
puting Systems, IEEE Computer Society, 2012, pp. 455–465.

- 1 [37] C. Kotselidis, M. Ansari, K. Jarvis, M. Luján, C. C. Kirkham, I. Watson, DiSTM: A 2
software transactional memory framework for clusters, in: Proceedings of the 2008 3
International Conference on Parallel Processing (ICPP 2008), IEEE Computer Soci-4
ety, 2008, pp. 51–58.
- 5 [38] M. Herlihy, V. Luchangco, M. Moir, A flexible framework for implementing software 6
transactional memory, in: Proceedings of the 21th Annual ACM SIGPLAN Confer-7
ence on Object-Oriented Programming, Systems, Languages, and Applications (OOP-8
SLA 2006), ACM, 2006, pp. 253–262.
- 9 [39] H. Attiya, L. Epstein, H. Shachnai, T. Tamir, Transactional contention management
10 as a non-clairvoyant scheduling problem, *Algorithmica* 57 (1) (2010) 44–61.
- 11 [40] A. Dragojevic, R. Guerraoui, A. V. Singh, V. Singh, Preventing versus curing: avoid-12
ing conflicts in transactional memories, in: Proceedings of the 28th Annual ACM 13
Symposium on Principles of Distributed Computing, (PODC 2009), ACM, 2009, pp. 14
7–16.
- 15 [41] R. Guerraoui, M. Herlihy, B. Pochon, Toward a theory of transactional contention 16
managers, in: Proceedings of the Twenty-Fourth Annual ACM Symposium on Prin-17
ciples of Distributed Computing (PODC 2005), ACM, 2005, pp. 258–264.
- 18 [42] G. Sharma, C. Busch, A competitive analysis for balanced transactional memory
19 workloads, *Algorithmica* 63 (1-2) (2012) 296–322.
- 20 [43] G. Sharma, C. Busch, Window-based greedy contention management for transactional
21 memory: theory and practice, *Distributed Computing* 25 (3) (2012) 225–248.
- 22 [44] R. M. Yoo, H. S. Lee, Adaptive transaction scheduling for transactional memory sys-23
tems, in: Proceedings of the 20th Annual ACM Symposium on Parallelism in Algo-24
rithms and Architectures (SPAA 2008), ACM, 2008, pp. 169–178.
- 25 [45] A. Baldassin, R. Murari, J. P. L. de Carvalho, G. Araujo, D. Castro, J. Barreto, P. RO-26
mano, NV-PhTM: an ecient phase-based transactional system for non-volatile mem-27
ory, in: Proceedings of the 26th International Conference on Parallel and Distributed 28
Computing (Euro-Par 2020), Vol. 12247 of LNCS, Springer, 2020, pp. 477–492.
- 29 [46] M. Liu, M. Zhang, K. Chen, X. Qian, Y. Wu, W. Zheng, J. Ren, DudeTM: building
30 durable transactions with decoupling for persistent memory, in: Proceedings of the 31
Twenty-Second International Conference on Architectural Support for Programming 32
Languages and Operating Systems (ASPLOS 2017), ACM, 2017, pp. 329–343.
- 33 [47] A. Kolli, S. Pelley, A. G. Saidi, P. M. Chen, T. F. Wenisch, High-performance transac-34
tions for persistent memories, in: Proceedings of the Twenty-First International Con-35
ference on Architectural Support for Programming Languages and Operating Systems 36
(ASPLOS 2016), ACM, 2016, pp. 399–411.